

A Multi-City COVID-19 Categorical Forecasting Model Utilizing Wastewater-Based Epidemiology

Naomi Rankin^{a,*}, Samee Saiyed^a, Hongru Du^a, Lauren Gardner^{a,b}

^aDepartment of Civil and Systems Engineering, Johns Hopkins University, Baltimore, MD, USA

^bDepartment of Epidemiology, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD, USA

Abstract

The COVID-19 pandemic highlighted shortcomings in forecasting models, such as unreliable inputs/outputs and poor performance at critical points. As COVID-19 remains a threat, it is imperative to improve current forecasting approaches by incorporating reliable data and alternative forecasting targets to better inform decision-makers.

Wastewater-based epidemiology (WBE) has emerged as a viable method to track COVID-19 transmission, offering a more reliable metric than reported cases for forecasting critical outcomes like hospitalizations. Recognizing the natural alignment of wastewater systems with city structures, ideal for leveraging WBE data, this study introduces a multi-city, wastewater-based forecasting model to categorically predict COVID-19 hospitalizations.

Using hospitalization and COVID-19 wastewater data for six US cities, accompanied by other epidemiological variables, we develop a Generalized Additive Model (GAM) to generate two categorization types. The Hospitalization Capacity Risk Categorization (HCR) predicts the burden on the healthcare system based on the number of available hospital beds in a city. The Hospitalization Rate Trend (HRT) Categorization predicts the trajectory of this burden based on the growth rate of COVID-19 hospitalizations. Using these categorical thresholds, we create probabilistic forecasts to retrospectively predict the risk and trend category of six cities over a 20-month period for 1, 2, and 3 week forecasting windows.

We also propose a new methodology to measure forecasting model performance at change points, or time periods where sudden changes in outbreak dynamics occurred. We also explore the influence of wastewater as a predictor for hospitalizations, showing its inclusion positively impacts the model's performance. With this categorical forecasting study, we are able to predict hospital capacity risk and disease trends in a novel and useful way, giving city decision-makers a new tool to predict COVID-19 hospitalizations.

Keywords: COVID-19, Wastewater-Based Epidemiology, City-level, Categorical Forecasting

1. Introduction

The COVID-19 pandemic had an unprecedented impact on public health and healthcare systems in the United States, resulting in more than 6.7 million hospitalizations and 1.1 million deaths as of February 3, 2024 (Centers for Disease Control and Prevention, 2024a). Various surveillance methods have been used to track the spread of the disease, including individual testing and surveys. However, these traditional approaches face challenges such as inconsistent sample sizes, which fluctuate over time, biases in individuals' willingness to get tested or participate in surveys, and the underreporting of at-home test results (Nixon et al., 2022; Li et al., 2023; Reese et al., 2021; Rubin, 2021). In order to address these challenges, wastewater-based epidemiology (WBE) has emerged as a viable method to track virus spread in a community by measuring SARS-CoV-2 concentration trends in wastewater (Feng et al., 2021). WBE offers a non-invasive, aggregated, low-cost approach for monitoring COVID-19 infection trends and potential incidence in a given catchment area (Shah et al., 2022; Bibby et al., 2021). Several studies have found correlations between SARS-CoV-2 levels in wastewater and lagged COVID-19 cases, hospitalizations, and deaths (Shah et al., 2022; Galani et al., 2022; Kanchan et al., 2024). These findings suggest that WBE is a viable representation of community disease prevalence and has the potential to serve as a valuable input for COVID-19 forecasting models.

Recent studies have incorporated SARS-CoV-2 load in wastewater as a key variable to help forecast COVID-19 hospitalizations. A recent study forecasted weekly COVID-19 hospitalizations for 159 counties in the U.S.

*Corresponding author.

Email address: nrankin2@jh.edu (Naomi Rankin)

42 using wastewater data, demonstrating that forecasting models that incorporate WBE can be effective at predicting
43 county-level COVID-19 hospital admissions 1 to 4 weeks in advance (Li et al., 2023). Hill et al., incorporated
44 WBE into a generalized linear mixed model to predict COVID-19 hospitalizations at 4 different geographic scales
45 (sewershed level, county level, regional level, and state level) for 56 counties in New York State to demonstrate the
46 higher predictive power provided by incorporating wastewater into forecasting models (Hill et al., 2023). Many
47 studies have also incorporated WBE into different modeling frameworks to predict COVID-19 hospitalizations
48 outside of the United States. For example, Schenk et al., incorporated WBE in multivariate regression models
49 to infer national hospitalization bed occupancy in Austria by aggregating SARS-CoV-2 load from local wastew-
50 ater treatment plants to a national level. The study determined that, at a national level, increasing the number
51 of monitored wastewater plants provides more accurate forecasts. (Schenk et al., 2023). Zamarreño et al., de-
52 veloped a dynamic artificial neural network to predict the number of hospitalized patients in Valladolid, Spain,
53 utilizing wastewater-based epidemiology. The accuracy of the model was improved by forecasting a categorical
54 risk level which aligns with warning levels established by the Regional Health Administration in Castile and León
55 (Zamarreño et al., 2024). Finally, Vaughan et al. incorporated WBE into a random forest algorithm to forecast
56 COVID-19 hospitalizations in Scotland, Catalonia, Ohio, the Netherlands, and Switzerland. To improve model
57 performance, the authors advocate for inputting data with a high sample frequency to ensure adequate training set
58 size (Vaughan et al., 2023).

59 While these studies demonstrate the promise of WBE for forecasting COVID-19 hospitalizations, certain
60 areas offer potential for further exploration and improvement. Firstly, most existing studies forecast continuous
61 COVID-19 health outcomes as targets, such as the raw number of COVID-19 cases and hospitalizations. However,
62 these continuous forecasts struggle to effectively communicate uncertainty to decision-makers, which can result
63 in misinterpretation of model results and can be challenging to translate into actionable outcomes (Nixon et al.,
64 2022). There is also a need to compare the efficacy of various non-continuous forecasts, such as those based on
65 resource availability or previous week trends, to ensure that the result is accurate and clear. Secondly, few models
66 are built at the city level, which aligns with wastewater testing catchment areas due to the presence of centralized
67 wastewater treatment facilities that offer the best support for local-level decision making (Sen et al., 2023; Hillary
68 et al., 2020). Thirdly, the potential of wastewater data to enhance model performance during critical periods,
69 particularly those characterized by rapid or sudden trends, remains uncertain. Therefore, model evaluations should
70 focus specifically on these periods to better understand and evaluate the model’s utility.

71 In this paper, we expand on the current literature by presenting an interpretable categorical short-term forecast-
72 ing model using Generalized Additive Models (GAMs) that incorporates WBE data to predict weekly COVID-19
73 hospitalizations for 6 cities in the United States from January 2021 to November 2022. Our multi-city model gen-
74 erates two actionable categorical outputs for city-level decision-makers: Hospitalization Capacity Risk (HCR),
75 which predicts the strain on hospital resources based on available beds, and Hospitalization Rate Trend (HRT),
76 which forecasts the future trajectory of disease transmission based on hospitalization trends. This model generates
77 1-week, 2-week, and 3-week ahead forecasts by incorporating SARS-CoV-2 wastewater load, recent hospitaliza-
78 tion rates, vaccination rates, prior infection data, and static variables like the COVID-19 Community Vulnerability
79 Index. We assessed model performance for both categorical outputs across six cities for over 80 weeks, utiliz-
80 ing five distinct error metrics. To enhance our understanding of the model’s accuracy and utility, we examined
81 performance at change points, which are defined as any week where the true hospitalization category is different
82 than the previous week. The GAM framework allows for the evaluation of variable contribution, highlighting
83 the value of using WBE data for forecasting COVID-19 health outcomes at the city level. The proposed model
84 offers a valuable tool for decision-making support due to its simple design, reliance on readily available data, and
85 production of easily interpretable and reliable forecasts.

86 **2. Material and Methods**

87 In this section, we present the design of our study, beginning with an overview of the data collection and
88 preprocessing procedures in Section 2.1, as well as the design of the target variables in Section 2.2. A summary
89 table of all variables is provided in Table 3. We then introduce the GAMs utilized in this study in Section 2.3 and
90 discuss the error metrics used to evaluate the model’s performance in Section 2.4.1. Furthermore, we propose a
91 novel approach for measuring forecasting capabilities at critical points in Section 2.4.2.

92 *2.1. Input Variables*

93 This study utilizes SARS-CoV-2 wastewater surveillance data alongside other epidemiological metrics such as
94 hospitalizations, previous infections, vaccination coverage, and local vulnerability indicators to forecast weekly

95 COVID-19 hospitalizations for six U.S. cities. City boundaries were delineated using USPS ZIP codes identified
 96 through the USPS ZIP code lookup tool (United States Postal Service, 2024). Detailed descriptions of each input
 97 variable are provided below.

98 **Hospitalizations (H):** In this study, we utilize COVID-19 hospitalization data as both an input to our model
 99 and to define our categorical targets, the latter of which is presented in Section 2.2. We obtained raw hospitalization
 100 data from the U.S. Department of Health and Human Services (HHS) COVID-19 Reported Patient Impact and
 101 Hospital Capacity by Facility dataset (United States Government Department of Health & Human Services, 2024).
 102 For each city i , the weekly reported hospitalizations at week t (H_i^t) were summed across all facilities f within the
 103 zip codes belonging to the city as follows:

$$H_i^t = \sum_{f \in F} H_f^t, \quad (1)$$

104 where f is a healthcare facility within the set of all facilities F in the zip codes of a the city at week t . The
 105 H_i^t were then logged to generate the input variable for the model to normalize for city-specific hospitalization
 106 scales. The detailed data cleaning procedures of the H_i^t are described in **Supplementary A**. Previous COVID-19
 107 hospitalizations are considered indicators of future trends, (Hill et al., 2023) thus at each week t , we include the
 108 hospitalization data from the previous three weeks as separate input variables ($H_i^{t-1}, H_i^{t-2}, H_i^{t-3}$). This approach
 109 enhances the model’s predictive capabilities by incorporating recent hospitalization dynamics.

110 **SARS-CoV-2 Wastewater Viral Load (WW):** Given that SARS-CoV-2 concentrations in wastewater are
 111 leading indicators of COVID-19 hospitalizations, we included city-level SARS-CoV-2 wastewater viral loads as
 112 an input variable to forecast city-level hospitalizations. The SARS-CoV-2 load in wastewater was obtained from
 113 state and city dashboards (see detailed reference in Table 1). For each city, we include SARS-CoV-2 concentrations
 114 for all sewersheds partially or completely within the city boundaries.

Table 1: Cities of interest in the study and their wastewater collection systems and dates.

City	Population	Number of Treatment Plants	Sampling Frequency	Source
Charlotte, NC	942,437	3	Every 3 days	(North Carolina Department of Health and Human Services, 2024)
Denver, CO	1,399,707	1	Daily	(Colorado Department of Public Health and Environment, 2024)
Houston, TX	3,206,416	39	Daily	(Kinder Institute Urban Data Platform, 2024)
New York, NY	8,570,761	14	Thrice Weekly	(New York City Department of Health and Mental Hygiene, 2024)
San Diego, CA	1,387,376	1	Every 2 days	(Lab, 2024)
San Francisco, CA	865,933	2	Daily	(California State Water Resources Control Board, 2024)

115 The SARS-CoV-2 wastewater viral loads were reported as viral gene copies/L at varying frequencies across
 116 cities. All viral load data was aggregated to a one-week reporting frequency and logged to normalize the input
 117 variable across cities. We evaluated the preceding relationship between wastewater viral loads and COVID-19
 118 hospitalization rates by identifying the leading weeks where there was a Pearson correlation coefficient of at least
 119 0.75 between the two. A detailed description of these correlations can be found in Figure S2 in **Supplementary**
 120 **A**. Accordingly, each week, we utilize the wastewater viral loads from the previous three weeks as separate input
 121 variables.

122 **Past Infections (PI):** Natural immunity has been shown to confer significant protection against COVID-19
 123 reinfection and severe outcomes (Pooley et al., 2023). To capture the impact of natural immunity, we adapted the
 124 past infection metric (PI) from a prior study as a proxy for population natural immunity (Du et al., 2024a). This
 125 variable quantifies the number of reported COVID-19 infections within the past three months, using daily reported
 126 county level case data from the Johns Hopkins COVID-19 Dashboard (Dong et al., 2020) aggregated to a weekly
 127 level. Due to the lack of consolidated reported case data at the city level, we use the data from the county that
 128 encompasses the majority of the city as a proxy for city-level case counts. The PI formulation is as follows:

$$PI_i^t = \frac{\sum_{j=t-16}^{t-4} C_i^j}{p_i}, \quad (2)$$

129 where PI_i^t denotes the total reported infections in county i during the previous 12 weeks, C_i^j is the number of
 130 reported cases for city i at week j , and p_i is the population for city i .

131 **Full Vaccination Coverage (VC):** In addition to the incorporation of PI as a proxy for natural immunity,
 132 we included the cumulative percentage of the population that is fully vaccinated to account for the impact of
 133 vaccine-induced immunity on COVID-19 hospitalizations. We obtained weekly COVID-19 vaccination data from

134 the Immunization Information System at the Centers for Disease Control and Prevention (Centers for Disease
 135 Control and Prevention, 2023). Similar to the reported case data, this national vaccination dataset tracks vaccine
 136 uptake at the county level. Therefore, we mapped city vaccination coverage data from the most representative
 137 county. Fully vaccinated is defined as the total number of individuals who have completed a primary vaccine
 138 series, either receiving the second dose of a two-dose series or one dose of a single-dose series. The percentage of
 139 fully vaccinated is calculated by dividing the fully vaccinated population by the total county population.

140 **COVID-19 Community Vulnerability Index (CCVI):** In order to account for population-level susceptibility
 141 to adverse disease outcomes, we utilize the COVID-19 Community Vulnerability Index (CCVI), which was de-
 142 veloped by the Surgo Foundation (Ventures, 2021). This index is an adaptation of the CDC Social Vulnerability
 143 Index with a focus on the specific risk factors of COVID-19 for all zip codes in the United States. The CCVI
 144 covers seven themes (Socioeconomic Status, Minority Status and Language, Housing Type, and Transport, Epi-
 145 demiological Factors, Healthcare System, High Risk Environments, and Population Density), giving each zip code
 146 a numerical value on a scale from 0 to 1, where a 1 indicates a community has a high vulnerability in that theme.
 147 For each of the cities in the study, we quantify the vulnerability themes by aggregating the zip-code level CCVI
 148 values as follows:

$$CCVI_i^k = \frac{\sum_{z=1}^N CCVI_z^k \times p_z}{p_i}, \quad (3)$$

149 where N is the set of all zip codes z within a city i , and p is the population. The zip code-level CCVI values were
 150 aggregated to the city level by weighting each by its zip code population and then normalizing by the total city
 151 population to ensure the variable ranges between 0 and 1. We repeat this aggregation for each of the seven themes
 152 k and use them as static inputs for our model to represent local vulnerability to adverse outcomes of COVID-19
 153 throughout the study period.

154 2.2. Target Design: Hospitalization Categorization

155 This study aims to provide city-level decision-makers with interpretable and actionable forecasts for COVID-
 156 19 hospitalizations by enhancing robustness against reporting issues. Rather than relying on traditional numerical
 157 predictions of hospitalization counts, we forecast risk categories that are derived from the rate of COVID-19
 158 hospitalizations per 100,000 people. We propose two distinct categorization models, each emphasizing different
 159 aspects of the future impact of COVID-19 on a city healthcare system: the Hospitalization Capacity Risk (HCR)
 160 Categorization, which predicts the burden on the healthcare system and the Hospitalization Rate Trend (HRT)
 161 Categorization, which predicts the trajectory of COVID-19 hospitalization trends.

162 **Hospitalization Capacity Risk (HCR) Categorization:** Hospital demand can be used to measure times of
 163 overburdened healthcare systems, indicating when personal risk of infection is at its highest. We develop the
 164 HCR as a 5-tier categorization model based on hospital demand, defined as the static ratio between observed
 165 hospital admission rate and average number of available hospital beds for COVID-19 forecasted in a city. The
 166 HCR categorization is defined as follows:

$$HCR_i^t = \begin{cases} \text{Very High Risk} & \text{if } HR_i^t > 10\% \text{ of } B_i \\ \text{High Risk} & \text{if } 10\% \text{ of } B_i > HR_i^t > 7.5\% \text{ of } B_i \\ \text{Moderate Risk} & \text{if } 7.4\% \text{ of } B_i > HR_i^t > 5\% \text{ of } B_i \\ \text{Low Risk} & \text{if } 4.9\% \text{ of } B_i > HR_i^t > 2\% \text{ of } B_i \\ \text{Very Low Risk} & \text{if } 2\% \text{ of } B_i < HR_i^t \end{cases} \quad (4)$$

167 Where HCR_i^t is the Hospitalization Capacity Risk category for city i at week t , HR_i^t is the forecasted COVID-19
 168 hospitalization rate per 100,000 people for city i at week t , and B_i is the average number of available hospital
 169 beds per 100,000 people for city i across the study period. The thresholds of this categorization reflect the static
 170 capacities of each specific city's healthcare system.

171 The thresholds are designed after the City of Austin stage alert system (Yang et al., 2021). They created a static
 172 4-tier alert system based on the percent of filled ICU beds and the 7-day average of COVID-19 hospitalizations
 173 at a city level from March 2020 to September 2021. We calculate the percentage of utilized COVID-19 beds for
 174 Austin equivalent to the threshold number of hospitalizations determining each stage threshold. Additionally, we
 175 introduce a fifth category of highest risk to further specify the intensity of burden during times of peak hospitaliza-
 176 tions. This categorization provides a more detailed and city-specific understanding of risk than the CDC 3-tier risk
 177 metric that identifies high, medium, and low hospital admission rates per 100,000 population based on a single
 178 threshold for all US locations (Centers for Disease Control and Prevention, 2024b). By increasing the number of

179 high-risk categories, we expand on the CDC 3-tier community risk system and the Austin 4-tier system while still
 180 maintaining thresholds that reflect predetermined foundations for risk levels.

181 Figure 1 illustrates the five risk categories of the HCR for the six cities, where the black line indicates the
 182 observed hospitalization rate per 100,000 population, and the colored bars indicate the risk category. Across
 183 cities, there are common times of higher risk, such as January 2022 during the peak of the Omicron wave, but the
 184 city-specific dynamics vary by risk level at other times. Table 2 presents the distribution of categories for the HCR
 185 across the study period and all locations. The majority of weeks are assigned as Low and Very Low risk, indicating
 186 that times of very high community risk are less prevalent across multiple locations over the study period.

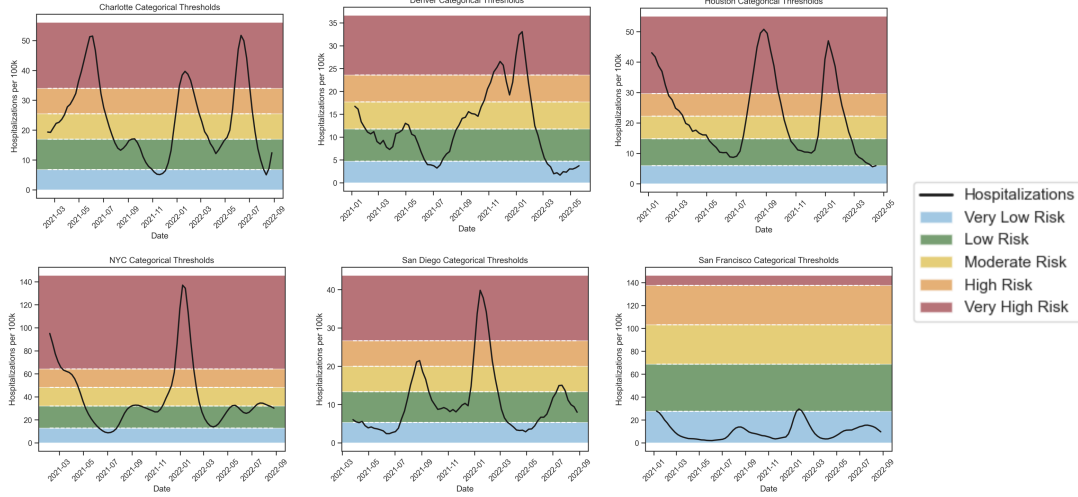


Figure 1: **Hospitalization Capacity Risk Categorization Visualization.** HCR category thresholds by color in Charlotte, NC; Denver, CO; Houston, TX; New York, NY; San Diego, CA; and San Francisco, CA. The background color indicates the risk category. The dotted white lines indicate the thresholds, and the black line is the true hospitalization rate.

187 **Hospitalization Rate Trend (HRT) Categorization:** Observing trends in the rate of change of COVID-19
 188 hospitalizations provides public health decision makers with a dynamic understanding of the change in burden on
 189 a city's healthcare system. We develop the HRT categorization to capture city-level COVID-19 hospitalization
 190 trends to support broader city-level decision-making to describe how quickly the number of hospitalizations is
 191 changing between weeks, providing a complementary metric to the static categorization of the HCR. This categorization
 192 is also derived from the weekly average COVID-19 hospitalization rate per 100,000 people (HR_i^t). This
 193 metric is useful to forecast periods of large changes for public planning purposes and is inspired by the FluSight
 194 forecasting target (FluSight-forecast-hub, 2024). We define the rate trend RT_i^t as the change in COVID-19 hospitalizations
 195 for a city i in a given week t relative to the average change of the prior three weeks. The rate change is
 196 calculated as:

$$RT_i^t = 100\% \times \frac{HR_i^t - \frac{1}{3} \sum_{j=1}^3 HR_i^{t-j}}{\frac{1}{3} \sum_{j=1}^3 HR_i^{t-j}}, \quad (5)$$

197 where HR_i^t is the hospitalizations per 100,000 at week t for city i , and the summation is the mean hospitalization
 198 rate for the previous 3 weeks in the same city. Using this growth rate, we define the HRT categorization as:

$$HRT_i^t = \begin{cases} \text{Large Increase} & \text{if } RT_i^t > 20\% \\ \text{Increase} & \text{if } 19.9\% > RT_i^t > 10\% \\ \text{Stable} & \text{if } 10\% > RT_i^t > -10\% \\ \text{Decrease} & \text{if } -10\% > RT_i^t > -19.9\% \\ \text{Large Decrease} & \text{if } -20\% > RT_i^t \end{cases}, \quad (6)$$

199 The thresholds of this categorization reflect the changing dynamics of the COVID-19 hospitalization trends.

200 Figure 2 illustrates the rate trend categorical assignment for all six cities. The black line indicates the change
 201 in hospitalizations compared to previous weeks. The rate trend differs between cities based on their individual
 202 dynamics, not an aggregated national trend, providing more specific and actionable information. Table 2 presents
 203 the distribution of categories for the HRT across the study period for all locations. Over our study period, most

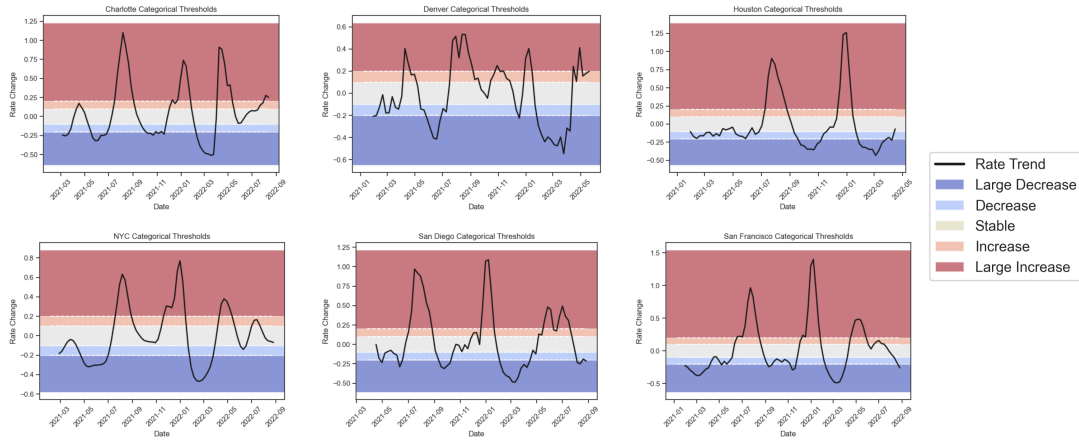


Figure 2: **Hospitalization Rate Trend Categorization Visualization.** HRT category thresholds by color in Charlotte, NC; Denver, CO; Houston, TX; New York, NY; San Diego, CA; and San Francisco, CA. The dotted white lines indicate the thresholds, the black line represents the observed rate change, and the background color indicates the category.

	Category	Threshold	Percent of Study Period
HCR	Very High Risk	>10%	14.2%
	High Risk	7.5 to 10%	8.8%
	Moderate Risk	5 to 7.4%	16.6%
	Low Risk	2 to 4.9%	31.9%
	Very Low Risk	<2%	28.5%
HRT	Large Increase	>20%	22.5%
	Increase	10 to 19.9%	12.1%
	Stable	-10 to 10%	19.8%
	Decrease	-10 to -19.9%	17.1%
	Large Decrease	<-20%	28.5%

Table 2: Thresholds for Hospitalization Rate Trend Categorization and distribution of categories

204 cities for most weeks are in the Large Increase or Large Decrease category, indicating the volatile trends of
 205 COVID-19 that have been historically difficult to predict.

206 2.3. Model

207 We develop 1, 2, and 3-week out forecasts of weekly COVID-19 hospitalizations for the 82-week period
 208 from January 2021 to September 2022 using wastewater viral load, past infections, the full vaccination coverage,
 209 previous hospitalizations, and 7 social vulnerability indices, summarized in Table 3. These forecasts are created
 210 via Generalized Additive Models (GAMs), an additive model that can learn non-linear predictor variables by
 211 modeling the outcome as a sum of spline functions of each predictor. We fit a distinct GAM for each forecasting
 212 window (1-week, 2-week, and 3-week), each utilizing the same model formulation:

$$\begin{aligned}
 & H_i^{t+l} \sim \text{Gaussian}(\mu^l) \\
 & \log(\mu^l) = \alpha + f_1(WW_i^{t-1}) + f_2(WW_i^{t-2}) + f_3(WW_i^{t-3}) + \\
 & \quad f_3(HR_i^{t-1}) + f_4(HR_i^{t-2}) + f_5(HR_i^{t-3}) + f_6(PI_i^t) + \\
 & \quad f_7(VC_i^t) + s_1(CCVI_1) + s_2(CCVI_2) + s_3(CCVI_3) + s_4(CCVI_4) + \\
 & \quad s_5(CCVI_5) + s_6(CCVI_6) + s_7(CCVI_7).
 \end{aligned}$$

218 Where l is the forecasting window, taking values of 1, 2, or 3. Each f_i coefficient indicates the spline smooth
 219 function of the dynamic variables, α indicates the intercept and μ^l is the mean of the hospitalization distribution
 220 l weeks ahead. Each s_i represents the parametric coefficient of the static variables. In order to mimic real-time
 221 forecasting practices, we used an expanding training window. This method involves initially training the model
 222 on a small set of data and then gradually expanding the training set as new data becomes available. In this study,
 223 the model is first trained on data before February 26, 2021. As each new week of data is collected, it is added to
 224 the training set, thereby expanding the window of historical data used for model training. This analysis was done
 225 using the pyGAM package in Python 3.9.13.

Variable	Calculation/Collection	Units	Type	Usage
SARS-CoV-2 Load in Wastewater	State/City Dashboards (3 preceding weeks)	log(N gene copies/L)	Dynamic	Input
Past Infections	HHS	Percent	Dynamic	Input
Percentage Vaccinated	IIS	Percent	Dynamic	Input
CCVI	Surgo Ventures	Ratio	Static	Input
Log Hospitalization Rate	HHS (3 preceding weeks)	log(hospitalizations/100k people)	Dynamic	Input
Hospitalization Capacity Risk Categorization	HCR Methodology	Risk Categories from 1 to 5	Static	Target
Hospitalization Rate Trend Categorization	HRT Methodology	Trend Categories from -2 to 2	Dynamic	Target

Table 3: Inputs and Targets of the model

226 The continuous predictions are then converted into discrete categories to generate categorical predictions for
227 each forecasting window. To create a probabilistic distribution over these categories, we performed prediction
228 interval sampling to generate multiple continuous outcomes. We then compute the frequency with which these
229 interval sample falls into each category. Specifically, for each week and city, we utilize 100 H_i^t values representing
230 the upper and lower bounds of 1% to 99% prediction intervals. Each interval value is then converted to a cate-
231 gorical label based on the procedure described in Section 2.2. The final forecast for each week and city is the
232 most frequent label among 100 probabilistic samples, selecting the categorical prediction with the greatest overall
233 agreement.

234 2.4. Evaluation

235 This section details the various ways in which we measure forecasting performance. We define and justify
236 the error metrics used to determine categorical forecast performance, as well as describe a novel procedure to
237 determine model accuracy at the critical points where the disease burden is increasing or decreasing.

238 2.4.1. Error Metrics

239 We evaluate our model using five error metrics: 1) Accuracy, 2) Mean Square Error (MSE), 3) Weighted Mean
240 Square Error (WMSE), 4) Brier Score, and 5) Brier Skill Score, each of which is established metrics for measuring
241 categorical forecasting error (Bradley et al., 2008; Du et al., 2024b).

242 Accuracy measures the percentage of predicted labels that match the true labels. Although widely used as
243 a baseline for evaluating model performance, accuracy alone does not convey how far off the predictions are
244 from the true values. To address this, we use the MSE, which reflects the magnitude of error, with smaller
245 values indicating better performance. However, MSE does not account for the uncertainty in the predictions. To
246 incorporate uncertainty, we also use the WMSE, where the weights are the probabilities of each category. This
247 metric penalizes the confidence in inaccurate forecasts, where a smaller value indicates a better performance. We
248 also measure our inaccuracies with the Brier Score, an error metric designed specifically to measure the accuracy
249 of probabilistic forecasts. It is calculated as the mean square difference between the predicted probability of each
250 category to the actual probability of each category. A BS value lies between 0 and 1, where a 0 reflects perfect
251 accuracy. From the Brier Score, we derive the Brier Skill Score, which compares our model performance to a
252 baseline of random guessing with uniform probabilities for each category. A BSS less than 0 indicates that our
253 predictions perform worse than the baseline, equal to 0 indicates that the prediction is equivalent to the baseline,
254 and greater than 0 indicates that our predictions perform better than the baseline. The equations for each of these
255 error metrics can be found in **Supplementary A**. Employing all five of these error metrics enables a comprehensive
256 evaluation of our model performance for both precision and confidence, specific to categorical predictions.

257 2.4.2. Model Performance at Change Points

258 A significant issue with hospitalization forecasting models during the peak of the COVID-19 pandemic was
259 the failure to accurately predict time periods where sudden changes in outbreak dynamics occurred (Lopez et al.,
260 2024). The ability to forecast such periods of rapid fluctuation sudden increases or decreases in reported hospi-
261 talization rates is critical for effective decision-making. To address this, we developed a novel method to evaluate
262 model performance during these periods of rapid change, referred to as “change points”. We define a change point
263 for both the HCR and HRT categorization as any week (t) where the true hospitalization category is different than
264 the category of the previous week (t_{-1}). We further classify change points into two categories: “upward shift”,

265 where the category is higher than that of the previous week, such as going from a Low Risk to a Moderate Risk
 266 for the HCR, or a "downward shift", where the category is lower than the previous week, such as going from a
 267 Large Increase to Stable for HRT.

268 We then evaluate each model's performance during periods of rapid changes by assessing how accurately each
 269 model captures the correct category over all change points, as well as at upward and downward shifts specifically.
 270 We evaluated performance using MSE, WMSE, Brier Score, and Brier Skill Score for all change points in the
 271 model, the detailed definitions of which are in Section 2.4.1.

272 3. Results

273 In this section, we describe the model performance using the error metrics described in Section 2.4.1. We then
 274 provide an in-depth analysis of the HCR and HRT categorization model performances for the 2-week forecasting
 275 window in Section 3.1. The equivalent results for the 1 and 3-week forecasting windows are provided in Figures
 276 S1-S4 in **Supplementary B**. In Section 3.2, we present results for the change point evaluation. Finally, in Section
 277 3.3, we present the results regarding variable importance.

278 3.1. Model Performance

279 Using the error metrics described in Section 2.4.1, we demonstrate the performance of the HCR and HRT
 280 categorization models for the 1, 2, and 3-week forecasting windows in Table 4. The Accuracy and BSS are error
 281 metrics where higher values indicate better model performance, whereas lower MSE, WMSE, and BS values
 282 indicate better model performance. The model performance reveals several key findings:

		Accuracy ↑	MSE ↓	WMSE ↓	BS ↓	BSS ↑
HCR	1-week	88.8%	0.114	0.106	0.157	0.885
	2-week	82.3%	0.205	0.181	0.262	0.822
	3-week	74.8%	0.382	0.269	0.341	0.748
HRT	1-week	69.6%	0.494	0.458	0.406	0.493
	2-week	56.4%	1.24	1.044	0.566	0.292
	3-week	45.2%	2.610	1.704	0.715	0.107

Table 4: A summary of model performance across various error metrics. ↑ / ↓ denotes if a higher/lower metric value signifies better performance.

283 Both categorization models significantly outperform random guesses. A random guess model would yield an
 284 accuracy of 20% and a BSS of 0, whereas the HCR accuracy spans from 75 to 89% and the HRT accuracy spans
 285 from 46 to 70% across all forecasting windows. Notably, all BSS values remain above 0, further underscoring the
 286 predictive power of the models.

287 For both the HRT and HCR categorization models, performance worsens as the forecasting window increases,
 288 demonstrated by the decrease in accuracy and BSS and the increase in MSE, BS, and WMSE. The decrease in
 289 model performance over increasing forecasting windows mirrors the decrease in correlations between wastewater
 290 and future hospitalizations over time. This diminishing relationship reflects rapidly changing COVID-19 dynamics
 291 in a community which makes the further future more difficult to predict.

292 Across all error metrics, model performance when predicting capacity (HCR) is consistently better than when
 293 predicting trends (HRT). The fluctuating shifts in COVID-19 hospitalization trends tend to be more volatile than
 294 changes in capacity, making the HRT more difficult to accurately predict. We detail these specific categorization
 295 model performances in the following sections.

296 3.1.1. Hospitalization Capacity Risk Categorization Results

297 In this section, we provide a detailed description of the performance of the HCR categorization model for the 2-
 298 week forecasting window. The 1 and 3-week forecast results are provided in Figures S1 and S2 of **Supplementary**
 299 **B**. In Figure 3, we demonstrate the overall categorical assignment accuracy, as well as the city specific categorical
 300 assignments over the study period.

301 The HCR model performance can be described with a confusion matrix, which visualizes the accuracy of
 302 the HCR model across all cities and the entire study period for the 2-week forecast. In the matrix, each row
 303 represents the true categorical label of a week, and each column is the predicted label of that week. Each matrix
 304 entry contains the number of weeks where the true label is categorized as that predicted label, with incorrect
 305 categorical assignments are indicated by dissimilar row and column labels. In the confusion matrix in Figure

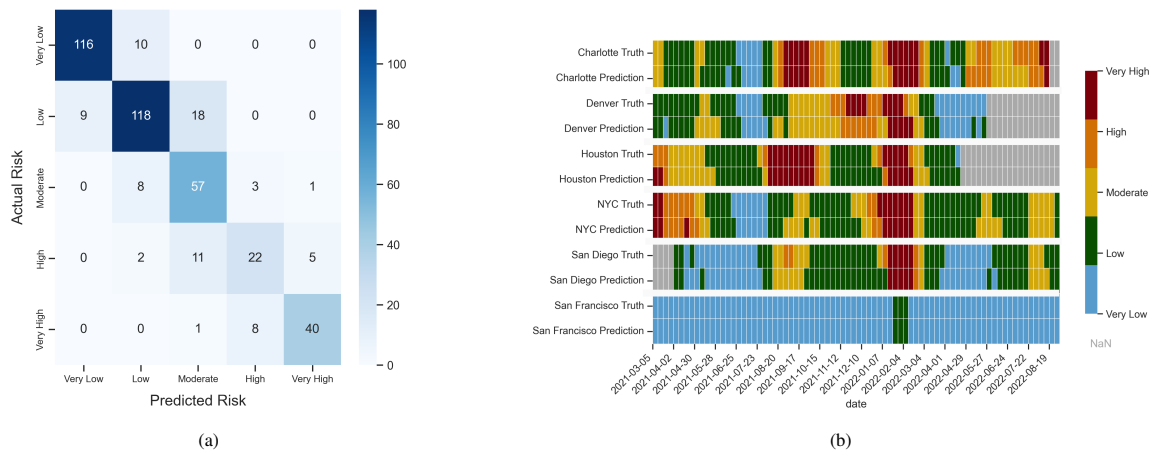


Figure 3: **Summary of Model performance of Hospitalization Capacity Risk (HCR) Categorization, January 2021-September 2022.** (a) Confusion Matrix detailing the accuracy of categorization for the HCR. Darker colors indicate that more assigned labels match the true labels for that category. (b) Hospital Capacity Risk true and predicted label for 2-week forecasts by city for January 2021 to September 2022. The color indicates the predicted category. Grey squares indicate that a prediction was not made for that week due to a lack of available data.

306 3a, the concentration of values along the diagonal indicates that the model accurately labels most forecasts. Off-
 307 diagonal values are primarily adjacent to the correct category, indicating that even when the model is incorrect,
 308 the magnitude of inaccuracy is relatively small.

309 Figure 3b compares the weekly HCR predictions to the true HCR categorizations for each city at a 2-week
 310 forecasting window, where the first row of the city illustrates the true weekly category and the second row illus-
 311 trates the predicted category. The color indicates the categorical assignment, with grey indicating that predictions
 312 were not made due to a lack of available data over the study period from January 2021 to September 2022. Al-
 313 though the HCR is dependent on specific city dynamics, there are similar trends across all cities. There is a
 314 similarly high risk for all cities in January 2022, during the Omicron wave. San Francisco maintains a relatively
 315 low risk throughout the study period due to the low rate of hospitalizations per 100,000 people, but the subsequent
 316 increase at this time is still present. When the model predictions do not match the true category, it tends to predict
 317 the correct category with a minor time lag.

318 3.1.2. Hospitalization Rate Trend Categorization Results

319 In this section, we focus on the performance of the HRT categorization model at the 2-week forecasting
 320 window. The 1 and 3-week forecast results are provided in Figures S3 and S4 in **Supplementary B**. In Figure 4,
 321 we demonstrate the overall categorical assignment accuracy with a confusion matrix, as well as the city specific
 322 categorical assignments with a time series plot equivalent to those shown in Figure 3.

323 The performance of the HRT categorization model for the 2-week forecasting window is illustrated in the
 324 confusion matrix in Figure 4a. The concentration of values at the top left and bottom right corners indicate that
 325 most COVID-19 hospitalization rate changes are large changes, and that the model can accurately predict Large
 326 Increase and Large Decrease growth rates. The concentration of values in the center of the confusion matrix
 327 indicates that weeks categorized as Stable, of which there are many, are also accurately predicted. The HRT
 328 model performance is worst during weeks categorized as Decrease and Increase, illustrated in their respective
 329 rows in the confusion matrix. In these rows, we can see that errors are nearly evenly split between overestimations
 330 and underestimations.

331 Figure 4b compares the weekly HRT predictions to the true HRT categorizations for each city at a 2-week
 332 forecasting window from January 2021 to September 2022 equivalently to Figure 3b. The top rows for each city
 333 illustrate the true HRT trends, revealing that periods of large increases are quickly followed by periods of large
 334 decreases, indicating the turbulent changes in COVID-19 hospitalizations over the study period. In the bottom
 335 row of each city we present the HRT predictions. Although the exact category is not always perfectly predicted,
 336 true increases generally yield forecasted increases and true decreases typically yield forecasted decreases. This
 337 phenomenon is expanded upon in Section 3.2.

338 3.2. Change Points Performance

339 In addition to building and evaluating interpretable city-level forecasting models, we also aim to develop a
 340 novel method to evaluate categorical forecasting model performance during periods of increases or decreases

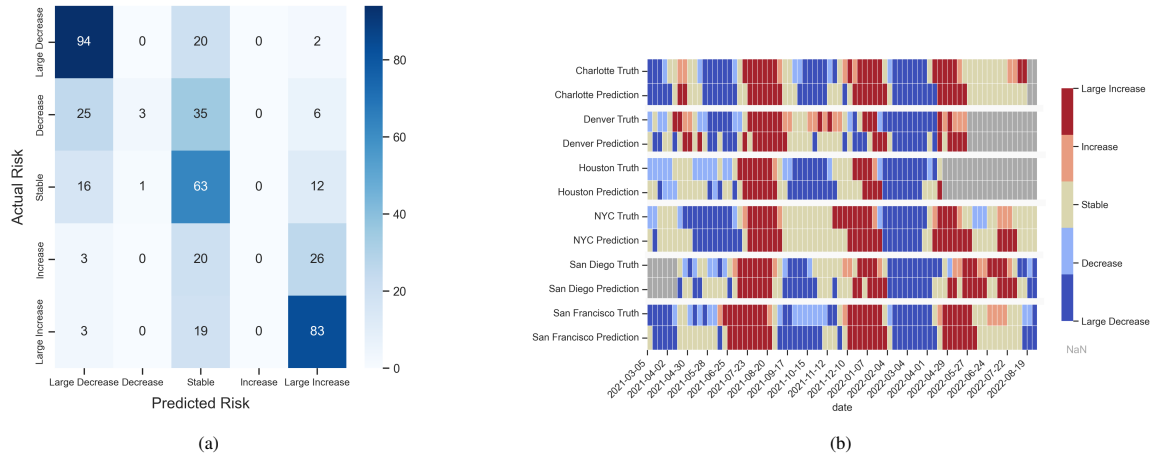


Figure 4: **Summary of Model performance of Hospitalization Rate Trend (HRT) Categorization, January 2021-September 2022.** (a) Confusion matrix detailing the accuracy of the categorization of the HRT. Darker colors indicate that more assigned labels match the true labels for that category. (b) City-specific Rate Trend target and 2-week prediction over time from January 2021 to September 2022. Grey squares indicate that a prediction was not made for that week due to a lack of available data.

		Upward Shift	Downward Shift	All Change Points	Overall Model
HCR	1 Week	0.224	0.262	0.244	0.106
	2 Week	0.288	0.359	0.325	0.181
	3 Week	0.446	0.462	0.454	0.269
HRT	1 Week	0.651	0.607	0.628	0.458
	2 Week	1.183	1.335	1.262	1.044
	3 Week	1.864	1.741	1.801	1.704

Table 5: A summary of model performance using WMSE for change points for 1. all change points in the model, 2. from a lower to higher category (upward shift), 3. from a higher to lower category (downward shift). The final column 4. is model WMSE over all predictions (same as the WMSE column in Table 4). ↓ denotes that a lower error demonstrates better performance.

341 in reported hospitalization rates, denoted as change points. As outlined in the section 2.4.2, we measured model
 342 performance by determining model accuracy at all change points, as well as at upward and downward shifts specifi-
 343 cally. We present the resulting WMSE scores in Table 5, and provide the MSE, Brier Score, and Brier Skill Score
 344 error metrics in Table T1 and T2 in **Supplementary B**. These results indicate that the HCR model consistently
 345 outperforms the HRT model, with the HCR predictions usually within a one category margin. Additionally, the
 346 error metrics generally worsen for farther ahead forecasts. Finally, we observe that both models perform worse at
 347 change points than in their overall evaluations.

348 3.3. Variable Importance Analysis

349 This work aims to demonstrate the value of wastewater based epidemiology for improving disease surveillance
 350 capabilities. The following analysis illustrates the importance of each variable in improving forecast performance.
 351 We measure variable importance as the difference in mean average percent error (MAPE) of the continuous hospi-
 352 talization forecasts when each variable is removed. This variable importance is measured for the 1, 2, and 3-week
 353 forecasts. This approach allows us to evaluate variable contribution independently of any categorical threshold
 354 selection biases. In Figure 5, we demonstrate the changes in MAPE after the removal of each variable, namely
 355 wastewater viral load, full vaccination coverage, past infections, and CCVI. For the 1 and 2-week forecasting
 356 window, the removal of wastewater leads to the highest increase in MAPE, indicating its critical role in improv-
 357 ing forecasting performance. For the 3-week forecast, the removal of the CCVI leads to the largest increase in
 358 MAPE, indicating that in addition to wastewater, the community vulnerability indices serve an important role in
 359 understanding community risk and creating useful predictions. Removing the full vaccination rate from the model
 360 does not generate significant differences in model performance. Additionally, the importance of past infections
 361 diminishes as the forecasting window extends.

362 4. Discussion

363 Integrating Wastewater Surveillance Data Enhances COVID-19 Hospitalization Forecasting

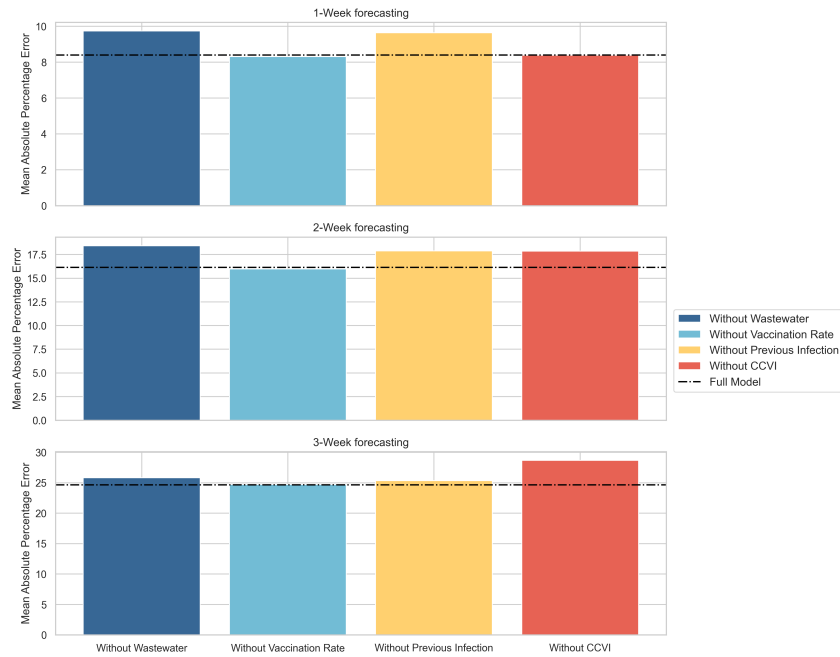


Figure 5: **Variable contribution in model performance for 1, 2, and 3 week out continuous hospitalization forecasts.** The black dotted line indicates the full model performance, and the colored bars indicate the MAPE with that variable removed. A larger increase in MAPE indicates a worse performance and a higher importance in model performance.

364 As traditional epidemiological datasets, such as reported cases, become less available and more subject to
 365 reporting issues over time, alternative data streams must be considered to inform short-term forecasting models.
 366 In this study, we built a comprehensive forecasting model to predict near term COVID-19 hospitalization risks in
 367 six US cities. This model incorporates novel data streams, including wastewater surveillance data and community
 368 vulnerability indices, alongside traditional epidemiological inputs such as previous infections and vaccinations.
 369 Our model performance confirms the effectiveness of this approach, particularly emphasizing the importance of
 370 wastewater surveillance data, as explored in Section 3.3.

371 Our variable importance analysis identified wastewater surveillance data as a critical factor in our models,
 372 particularly for forecasting city-level COVID-19 hospitalizations 1 and 2 weeks into the future (see Figure 5).
 373 Therefore, publicly available, reliable wastewater data is crucial for improving our capabilities to predict future
 374 hospitalization risks and inform both the public and public health practitioners. Wastewater disease surveillance
 375 systems, such as the CDC National Wastewater Surveillance System (NWSS), represent a centralized database
 376 that has the potential to inform disease risk assessment and bolster modeling capabilities. Thus, we encourage
 377 continued investment in and expansion of this critical surveillance tool, along with the centralized sharing of data
 378 in standardized formats.

379 Beyond the immediate utility of wastewater surveillance, our findings underscore the importance of local
 380 vulnerability indices on model performance. We determined that the CCVI is the most influential variable for our
 381 3-week continuous forecasting model. This finding aligns with previous studies demonstrating the critical role
 382 of vulnerability data in predicting weekly COVID-19 hospital admissions at less-immediate forecasting windows,
 383 such as 3 and 4 weeks in the future (Li et al., 2023). The importance of CCVI in our analysis suggests that more
 384 severe health outcomes from COVID-19 infection may be amplified in vulnerable communities. As the sole static
 385 variable in the model, the community vulnerability index avoids any information decay that may arise with the
 386 temporally dependent variables, such as wastewater surveillance data. This enduring relevance underscores the
 387 need to incorporate data on both current disease dynamics, such as wastewater surveillance data, and persistent
 388 vulnerabilities, such as the CCVI, when forecasting near future infectious disease risk.

389 **Categorical Forecasts Support Public Health Decision Making**

390 The well-documented challenges of forecasting traditional continuous disease outcomes during the COVID-19
 391 pandemic (Nixon et al., 2022) motivated this study's focus on categorical targets for hospitalization forecasting.
 392 This approach, exemplified by the Hospitalization Capacity Risk (HCR) and Hospitalization Rate Trend (HRT)
 393 categorizations defined in this study, enhances model interpretability, actionability, and robustness.

394 The HCR model, derived from hospitalization rates, provides hospital administrators with an intuitive under-

standing of capacity risk, enabling proactive resource allocation during times of high risk. Our models achieve a 75% accuracy in predicting HCR three weeks in advance, proving consistently reliable prediction across cities and forecasting windows. Exceptionally, HCR forecasts for San Francisco resulted in a uniquely high accuracy due to the city's uniquely low hospitalization rates. While the HRT model, based on the trends of hospitalization growth rates, is a less predictable outcome, the model successfully captures 64% of weeks classified as Large Increases and 67% of weeks classified as Large Decreases for the 2-week forecasting window (Figure 4a). The focus of the HRT on directional trends, when considered in conjunction with HCR, offers decision-makers a more comprehensive understanding of potential outbreak scenarios, facilitating timely and targeted interventions.

The inherent stability of the two categorical target variables defined in this study, which are less sensitive to minor data fluctuations than continuous outputs, further reinforces the reliability of our approach and provides a robust foundation for data-driven decision-making in the face of evolving public health challenges.

Evaluating Forecasting Performance at Change Points Illuminates Model Capabilities

A notable issue with COVID-19 forecasting during the pandemic was the inability to accurately predict inflection points in case and hospitalization trends (Lopez et al., 2024; Cramer et al., 2021). To assess the effectiveness of our categorical models in addressing this challenge, we evaluated performance at change points, examining the models' ability to accurately account for shifts in the COVID-19 hospitalization category. As outlined in section 3.2, we evaluated model performance by determining accuracy at all change points, as well as at upward and downward shifts specifically.

The results outlined in section 3.2 presented several important insights that merit further discussion. First, we note that our HCR model generally outperforms the HRT model for all change point evaluation methods. This aligns with the overall model performance in section 3.1, where the HRT model's weaker performance could be attributed to its prediction of highly variable hospitalization growth rates. Second, we highlight that error metrics at change points generally worsen as forecasting window increases. This is also consistent with the results in section 3.1, which suggests that this decline is likely due to the decreased correlations between wastewater and future hospitalizations over time. Third, we recognize that the performance at change points is weaker for both models compared to overall model accuracy, as per Table 5 the difference in WMSE is less than 0.3 across all models and forecast horizons. While there is still room for improvement, this indicates that the model is reasonably effective at predicting the correct category at change points. Nevertheless, this discrepancy in model performance at change points further highlights the need for a greater focus on improving forecasting accuracy during periods of rapid changes.

Our change point analysis methodology evaluates the utility of categorical forecasting models at critical points. Sudden shifts in population level disease dynamics are some of the most crucial trends to understand in order to better inform decision makers at times of great uncertainty. Our change point analysis enhances model utility by evaluating how well each categorization captures hospitalization dynamics during periods of rapid change. We believe that this straightforward yet comprehensive approach to assessing categorical model accuracy at change points is an effective and transferable method for evaluating model performance and developing timely policies to mitigate the effects of a virus.

City-Level COVID-19 Hospitalization Forecasting is Effective and Useful

We propose that achieving a balance between forecast accuracy and actionable spatial scales, such as cities, is crucial for providing key insights that can be translated into effective mitigation strategies for COVID-19 and other infectious diseases. COVID-19 forecasting using wastewater surveillance has been implemented at several different granularities, such as country, state, county, city, and sewershed level. (Schenk et al., 2023; Hill et al., 2023; Li et al., 2023; Galani et al., 2022). We contend that forecasting at a city-level, which has received the least attention by modelers and is well-aligned with the spatial boundaries of wastewater data, is particularly advantageous for decision-makers. Densely populated urban areas are often the epicenter of COVID-19 infections and have government structures that allow for more targeted public health interventions. Thus, accurate city-level forecasts can more effectively help inform local government policies to prevent harm to the population and minimize economic burden at these disease epicenters (Lopez et al., 2024). Despite this potential, city-level forecasting is an under-utilized tool, with most forecasts targeting larger, aggregated populations such as county-level, state-level, and national-level (Sen et al., 2023). This is often due to the lack of reliable city-level disease surveillance, a challenge addressed by the natural alignment between wastewater treatment plants and city boundaries as demonstrated in Figure S1 in **Supplementary A**. Additionally, many cities have multiple wastewater treatment facilities which can be aggregated to provide a smoother prediction (Medina et al., 2022). Multiple wastewater treatment facilities also enable higher modeling granularities within cities. Overall, WBE provides an opportunity to advance higher resolution infectious disease forecasting and permits local governments to implement mitigation strategies that best suits their communities.

We note that the quality of wastewater data significantly impacts model performance across cities. Houston

and New York City consistently demonstrate superior performance due to their high-quality wastewater data, characterized by consistent reporting and a large number of sewersheds. Conversely, as shown in Figures 3b and 4b, Charlotte and Denver exhibit poorer model performance for both HCR and HRT due to lesser data quality, stemming from recorded lapses in reporting or changes in reporting styles (Keshaviah et al., 2023; Ghanbari et al., 2024). This finding underscores the importance of consistent and transparent WBE data for both understanding local trends and enhancing forecasting capabilities.

Limitations and Future Work.

Wastewater disease surveillance for COVID-19 is a rapidly evolving field. Different cities and wastewater treatment plants have varying timelines for reporting and normalization techniques that make retrospective data difficult to compare to currently reported information. Despite the expansion of the NWSS to expand wastewater disease surveillance systems, a lack of data accessibility and inconsistencies in data collection, standardization and preprocessing can still impact data quality across locations. Consequently, the observed variation in model performance across cities can be attributed to these inherent data quality differences.

The models in this study are better equipped to capture downward shifts, which may be reflective of limited input data streams. Therefore, we propose including alternative disease prevalence information, such as circulating variants, which may improve model performance, as demonstrated in (Du et al., 2023). Wastewater-based epidemiology is a vehicle to understand community disease prevalence, and variant information can be easily recovered from it to strengthen model performance, something considered for future extensions of this work.

In this study, we focus on aggregating epidemiological data to a city level to create a more granular forecasting scale than the traditional state or county-level forecasts. By utilizing wastewater disease surveillance, we are able to provide an even more specific geographical scale at which to understand disease dynamics. In future work, we aim to expand our modeling capabilities to understand wastewater signals at the sewershed level as an early warning signal for smaller spatial scales.

5. Conclusions

In this study, we explored the use of wastewater disease surveillance data for categorical COVID-19 forecasting at the city level with two models: the Hospitalization Capacity Risk (HCR) categorization and the Hospitalization Rate Trend (HRT) categorization. The HCR model provides an understanding of general risk based on bed availability to inform resource allocation strategies. The HRT Categorization forecasts the growth rate of hospitalizations to understand the predicted trajectory of burden on a healthcare system. We also proposed a methodology to determine the circumstances under which models are useful with our change point analysis. Utilizing these models, we demonstrate that wastewater disease surveillance data is a crucial input for COVID-19 hospitalization forecasting. Due to the natural alignment of city borders and sewersheds, and the ability to generate reliable forecasts, city-level forecasting utilizing wastewater disease surveillance data is an advantageous approach. As the community dynamics of COVID-19 change and data availability shifts, we must update our forecasting capabilities to use more impactful inputs and design more beneficial and interpretable targets.

Acknowledgements

Funding sources

This work was supported by NSF RAPID Award ID 2333435.

References

- Bibby, K., Bivins, A., Wu, Z., North, D., 2021. Making waves: plausible lead time for wastewater based epidemiology as an early warning system for covid-19. *Water Research* 202, 117438.
- Bradley, A.A., Schwartz, S.S., Hashino, T., 2008. Sampling uncertainty and confidence intervals for the brier score and brier skill score. *Weather and Forecasting* 23, 992–1006.
- California State Water Resources Control Board, 2024. Sewershed surveillance for covid-19 updates. URL: <https://www.waterboards.ca.gov/resources/>. accessed: 2024-07-10.
- Centers for Disease Control and Prevention, 2023. Vaccines for COVID-19. <https://www.cdc.gov/coronavirus/2019-ncov/vaccines/>.
- Centers for Disease Control and Prevention, 2024a. Cdc covid data tracker. URL: <https://covid.cdc.gov/covid-data-tracker/>. accessed: 2024-07-10.
- Centers for Disease Control and Prevention, 2024b. Indicators for monitoring covid-19 community levels. URL: <https://archive.cdc.gov/>. accessed: 2024-07-10.
- Colorado Department of Public Health and Environment, 2024. CDPHE COVID19 Wastewater Dashboard Data . <https://data-cdphe.opendata.arcgis.com/datasets>.

504 Cramer, E.Y., Huang, Y., Wang, Y., Ray, E.L., Cornell, M., Bracher, J., Brennen, A., Castro Rivadeneira, A.J., Gerding,
505 A., House, K., Jayawardena, D., Kanji, A.H., Khandelwal, A., Le, K., Niemi, J., Stark, A., Shah, A., Wattanachit, N.,
506 Zorn, M.W., Reich, N.G., Consortium, U.C.F.H., 2021. The united states covid-19 forecast hub dataset. medRxiv URL:
507 <https://www.medrxiv.org/content/10.1101/2021.11.04.21265886v1>, doi:10.1101/2021.11.04.21265886.

508 Dong, E., Du, H., Gardner, L., 2020. An interactive web-based dashboard to track covid-19 in real time. *The Lancet infectious diseases* 20,
509 533–534.

510 Du, H., Dong, E., Badr, H.S., Petrone, M.E., Grubaugh, N.D., Gardner, L.M., 2023. Incorporating variant frequencies data into short-term
511 forecasting for covid-19 cases and deaths in the usa: A deep learning approach. *Ebiomedicine* 89.

512 Du, H., Saiyed, S., Gardner, L.M., 2024a. Association between vaccination rates and covid-19 health outcomes in the united states: a
513 population-level statistical analysis. *BMC Public Health* 24, 220.

514 Du, H., Zhao, J., Zhao, Y., Xu, S., Lin, X., Chen, Y., Gardner, L.M., Yang, H.F., 2024b. Advancing real-time pandemic forecasting using large
515 language models: A covid-19 case study [arXiv:2404.06962](https://arxiv.org/abs/2404.06962).

516 Feng, S., Roguet, A., McClary-Gutierrez, J.S., Newton, R.J., Kloczko, N., Meiman, J.G., McLellan, S.L., 2021. Evaluation of sampling,
517 analysis, and normalization methods for sars-cov-2 concentrations in wastewater to assess covid-19 burdens in wisconsin communities.
518 *ACS ES&T Water* 1, 1955–1965.

519 FluSight-forecast-hub, 2024. FluSight 2023-2024. <https://github.com/cdcepi/FluSight-forecast-hub/tree/main>.

520 Galani, A., Aalizadeh, R., Kostakis, M., Markou, A., Alygizakis, N., Lytras, T., Adamopoulos, P.G., Peccia, J., Thompson, D.C., Kontou, A.,
521 et al., 2022. Sars-cov-2 wastewater surveillance data can predict hospitalizations and icu admissions. *Science of The Total Environment*
522 804, 150151.

523 Ghanbari, M., Huang, J., Luc, A., Arabi, M., Goldman, J.E., Byrne-Nash, R., Kane, S.J., Ferrell, R., Fielder, T., De Long, S.K., Wilusz, C.J.,
524 2024. View of an evolving pandemic: Changes in the relationship between clinical cases and levels of sars-cov-2 rna in colorado wastew-
525 ater. *ACS ES&T Water* 4, 2018–2030. URL: <https://doi.org/10.1021/acsestwater.3c00615>, doi:10.1021/acsestwater.3c00615,
526 [arXiv:https://doi.org/10.1021/acsestwater.3c00615](https://doi.org/10.1021/acsestwater.3c00615).

527 Hill, D.T., Alazawi, M.A., Moran, E.J., Bennett, L.J., Bradley, I., Collins, M.B., Gobler, C.J., Green, H., Insaf, T.Z., Kmush, B., et al., 2023.
528 Wastewater surveillance provides 10-days forecasting of covid-19 hospitalizations superior to cases and test positivity: A prediction study.
529 *Infectious disease modelling* 8, 1138–1150.

530 Hillary, L.S., Malham, S.K., McDonald, J.E., Jones, D.L., 2020. Wastewater and public health: the potential of wastewater surveillance for
531 monitoring covid-19. *Current Opinion in Environmental Science & Health* 17, 14–20.

532 Kanchan, S., Ogden, E., Kesheri, M., Skinner, A., Miliken, E., Lyman, D., Armstrong, J., Sciglitano, L., Hampikian, G., 2024. Covid-19
533 hospitalizations and deaths predicted by sars-cov-2 levels in boise, idaho wastewater. *Science of The Total Environment* 907, 167742.

534 Keshaviah, A., Diamond, M.B., Wade, M.J., Scarpino, S.V., Ahmed, W., Amman, F., Aruna, O., Badilla-Aguilar, A., Bar-Or, I., Bergthaler,
535 A., et al., 2023. Wastewater monitoring can anchor global disease surveillance systems. *The Lancet Global Health* 11, e976–e981.

536 Kinder Institute Urban Data Platform, 2024. Dataset catalog. URL: <https://www.kinderudp.org>. accessed: 2024-07-10.

537 Lab, A., 2024. Sars-cov-2 wastewater san diego. URL: <https://github.com/andersen-lab>. accessed: 2024-07-10.

538 Li, X., Liu, H., Gao, L., Sherchan, S.P., Zhou, T., Khan, S.J., Van Loosdrecht, M.C., Wang, Q., 2023. Wastewater-based epidemiology predicts
539 covid-19-induced weekly new hospital admissions in over 150 usa counties. *Nature Communications* 14, 4548.

540 Lopez, V.K., Cramer, E.Y., Pagano, R., Drake, J.M., O’Dea, E.B., Adee, M., Ayer, T., Chhatwal, J., Dalgic, O.O., Ladd, M.A., Linas,
541 B.P., Mueller, P.P., Xiao, J., Bracher, J., Castro Rivadeneira, A.J., Gerding, A., Gneiting, T., Huang, Y., Jayawardena, D., Kanji, A.H.,
542 Le, K., Mühlemann, A., Niemi, J., Ray, E.L., Stark, A., Wang, Y., Wattanachit, N., Zorn, M.W., Pei, S., Shaman, J., Yamana, T.K.,
543 Tarasewicz, S.R., Wilson, D.J., Baccam, S., Gurung, H., Stage, S., Suchoski, B., Gao, L., Gu, Z., Kim, M., Li, X., Wang, G., Wang, L.,
544 Wang, Y., Yu, S., Gardner, L., Jindal, S., Marshall, M., Nixon, K., Dent, J., Hill, A.L., Kaminsky, J., Lee, E.C., Lemaitre, J.C., Lessler,
545 J., Smith, C.P., Truelove, S., Kinsey, M., Mullany, L.C., Rainwater-Lovett, K., Shin, L., Tallaksen, K., Wilson, S., Karlen, D., Castro,
546 L., Fairchild, G., Michaud, I., Osthus, D., Bian, J., Cao, W., Gao, Z., Lavista Ferres, J., Li, C., Liu, T.Y., Xie, X., Zhang, S., Zheng,
547 S., Chinazzi, M., Davis, J.T., Mu, K., Pastore y Piontti, A., Vespignani, A., Xiong, X., Walraven, R., Chen, J., Gu, Q., Wang, L., Xu,
548 P., Zhang, W., Zou, D., Gibson, G.C., Sheldon, D., Srivastava, A., Adiga, A., Hurt, B., Kaur, G., Lewis, B., Marathe, M., Peddireddy,
549 A.S., Porebski, P., Venkatramanan, S., Wang, L., Prasad, P.V., Walker, J.W., Webber, A.E., Slayton, R.B., Biggerstaff, M., Reich, N.G.,
550 Johansson, M.A., 2024. Challenges of covid-19 case forecasting in the us, 2020–2021. *PLOS Computational Biology* 20, 1–25. URL:
551 <https://doi.org/10.1371/journal.pcbi.1011200>, doi:10.1371/journal.pcbi.1011200.

552 Medina, C.Y., Kadonsky, K.F., Roman Jr, F.A., Tariqi, A.Q., Sinclair, R.G., D’Aoust, P.M., Delatolla, R., Bischel, H.N., Naughton, C.C., 2022.
553 The need of an environmental justice approach for wastewater based epidemiology for rural and disadvantaged communities: a review in
554 california. *Current Opinion in Environmental Science & Health* 27, 100348.

555 New York City Department of Health and Mental Hygiene, 2024. Sars-cov-2 concentrations measured in nyc wastewater. URL:
556 <https://data.cityofnewyork.us/Health/SARS-CoV-2-concentrations-measured-in-NYC-Wastewat>. accessed: 2024-07-
557 10.

558 Nixon, K., Jindal, S., Parker, F., Marshall, M., Reich, N.G., Ghobadi, K., Lee, E.C., Truelove, S., Gardner, L., 2022. Real-time covid-19
559 forecasting: challenges and opportunities of model performance and translation. *The Lancet Digital Health* 4, e699–e701.

560 North Carolina Department of Health and Human Services, 2024. NC COVID-19 Dashboard Data .
561 <https://covid19.ncdhhs.gov/dashboard/data-behind-dashboards>.

562 Pooley, N., Abdool Karim, S.S., Combadière, B., Ooi, E.E., Harris, R.C., El Guerche Seblain, C., Kisomi, M., Shaikh, N., 2023. Durability of
563 vaccine-induced and natural immunity against covid-19: a narrative review. *Infectious Diseases and Therapy* 12, 367–387.

564 Reese, H., Iuliano, A.D., Patel, N.N., Garg, S., Kim, L., Silk, B.J., Hall, A.J., Fry, A., Reed, C., 2021. Estimated incidence of coronavirus
565 disease 2019 (covid-19) illness and hospitalization—united states, february–september 2020. *Clinical Infectious Diseases* 72, e1010–e1017.

566 Rubin, R., 2021. Covid-19 testing moves out of the clinic and into the home. *Jama* 326, 1362–1364.

567 Schenk, H., Heidinger, P., Insam, H., Kreuzinger, N., Markt, R., Nägele, F., Oberacher, H., Scheffknecht, C., Steinlechner, M., Vogl, G., et al.,
568 2023. Prediction of hospitalisations based on wastewater-based sars-cov-2 epidemiology. *Science of The Total Environment* 873, 162149.

569 Sen, A., Stevens, N.T., Tran, N.K., Agarwal, R.R., Zhang, Q., Dubin, J.A., 2023. Forecasting daily covid-19 cases with gradient boosted
570 regression trees and other methods: evidence from us cities. *Frontiers in Public Health* 11, 1259410.

571 Shah, S., Gwee, S.X.W., Ng, J.Q.X., Lau, N., Koh, J., Pang, J., 2022. Wastewater surveillance to infer covid-19 transmission: A systematic
572 review. *Science of The Total Environment* 804, 150060.

573 United States Government Department of Health & Human Services, 2024. COVID-19 Reported Patient Impact and Hospital Capacity by
574 Facility. <https://healthdata.gov/d/t7zc-4t6g>.

575 United States Postal Service, 2024. Zip code lookup. URL: <https://tools.usps.com/zip-code-lookup.htm>. accessed: 2024-07-10.
576 Vaughan, L., Zhang, M., Gu, H., Rose, J.B., Naughton, C.C., Medema, G., Allan, V., Roiko, A., Blackall, L., Zamyadi, A., 2023. An
577 exploration of challenges associated with machine learning for time series forecasting of covid-19 community spread using wastewater-
578 based epidemiological data. *Science of The Total Environment* 858, 159748.
579 Ventures, S., 2021. Vulnerable communities and covid-19: the damage done, and the way forward. Washington, DC: Surgo Ventures .
580 Yang, H., Sürer, Ö., Duque, D., Morton, D.P., Singh, B., Fox, S.J., Pasco, R., Pierce, K., Rathouz, P., Valencia, V., et al., 2021. Design of
581 covid-19 staged alert systems to ensure healthcare capacity with minimal closures. *Nature communications* 12, 3767.
582 Zamarreño, J.M., Torres-Franco, A.F., Gonçalves, J., Muñoz, R., Rodríguez, E., Eiros, J.M., García-Encina, P., 2024. Wastewater-based
583 epidemiology for covid-19 using dynamic artificial neural networks. *Science of The Total Environment* 917, 170367.

1. Supplementary A: Methods

1.1. Data Collection

Sewersheds and City Boundaries The sewershed maps are provided by city services for each location, the references for which can be found in the image caption. The city boundary and ZIP code maps were created using the Python package GEOPANDAS. The city boundaries are generated using data obtained from the U.S. Department of Health and Human Services (US HHS, 2023). We delineated these maps with the USPS ZIP codes, as identified by the USPS ZIP code lookup tool (United States Postal Service, 2024). To properly represent the USPS ZIP codes on the map, we overlaid the USPS ZIP codes with the ZIP Code Census Tract Areas, obtained from the United States Census Bureau (Census, 2020). Therefore, any ZIP code that did not correspond with a physical polygon (i.e., ZIP code that referred to a single building) was removed from these maps.

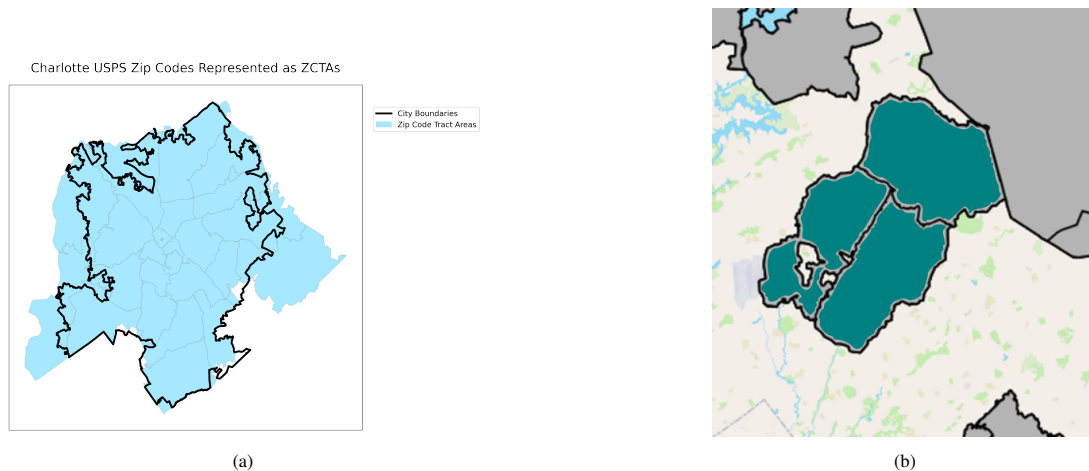


Figure 1: A side by side comparison of A) Charlotte city boundaries delineated by USPS Zip Code and B) the sewershed map of Charlotte. The sewersheds included in the study are highlighted in blue. (North Carolina Department of Health and Human Services, 2024)

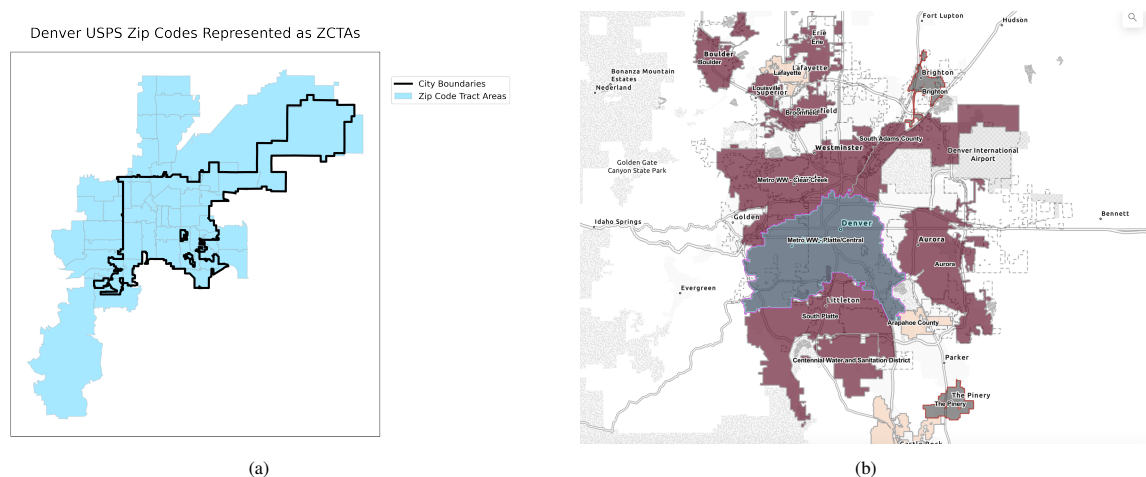
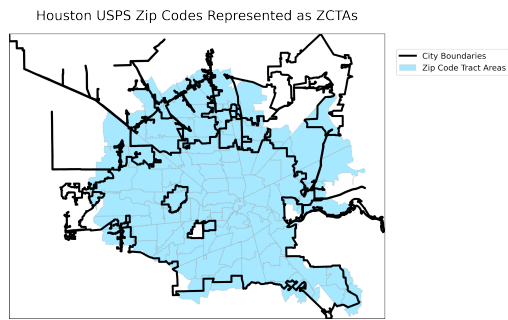
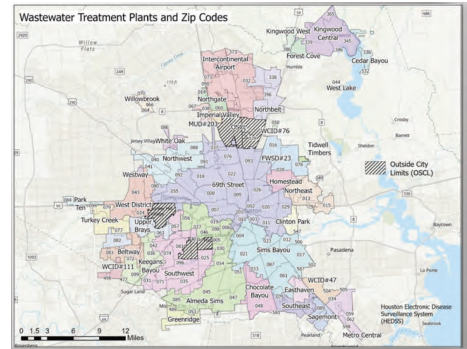


Figure 2: A side by side comparison of A) Denver city boundaries delineated by USPS Zip Code and B) the sewershed map of Colorado. The included sewershed is highlighted in blue. (Colorado Department of Public Health and Environment, 2024)

Cross Correlations Between Wastewater and Hospitalizations In order to determine the temporal relationship between wastewater SARS-CoV-2 rates and COVID-19 hospitalizations, we examine the cross correlation between these two time series. As shown in Figure 7, the weeks when the lag is greater than 0.75 indicate that wastewater SARS-CoV-2 rates are highly correlated with hospitalizations that many weeks in advance. For example, in Houston, TX, the wastewater SARS-CoV-2 rates are highly correlated for the observed week of hospitalizations, as well as the preceding 35 weeks. We included the SARS-CoV-2 rates from the previous three weeks in the model as predictor variables to align with the most common correlations observed across all cities.

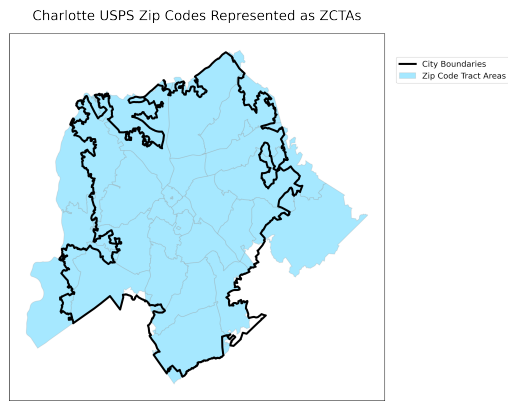


(a)

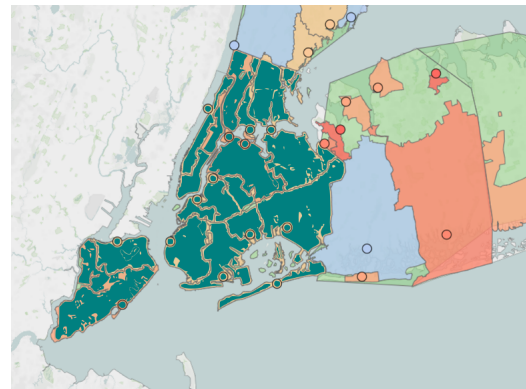


(b)

Figure 3: A side by side comparison of A) Houston city boundaries delineated by USPS Zip Code and B) the sewershed map of Houston (Hopkins et al., 2022)

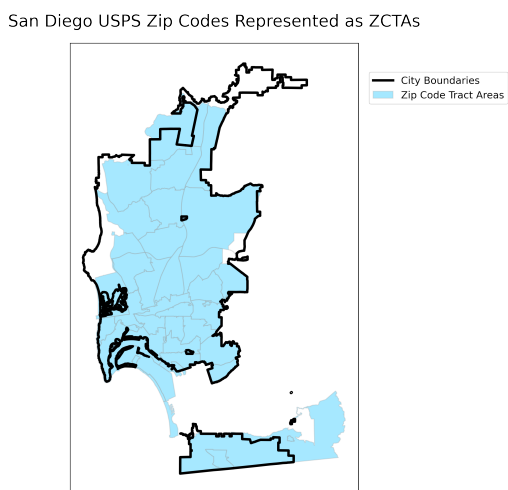


(a)

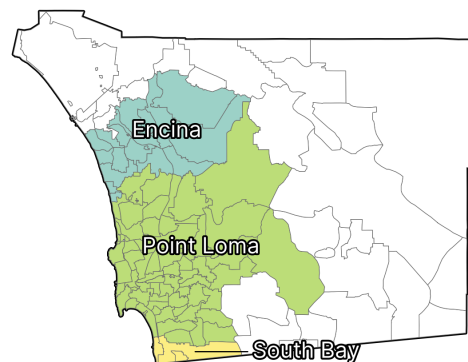


(b)

Figure 4: A side by side comparison of A) New York City boundaries delineated by USPS Zip Code and B) the sewershed map of New York City CITATION



(a)



(b)

Figure 5: A side by side comparison of A) San Diego city boundaries delineated by USPS Zip Code and B) the sewershed map of San Diego: (Lab, 2024)

San Francisco USPS Zip Codes Represented as ZCTAs

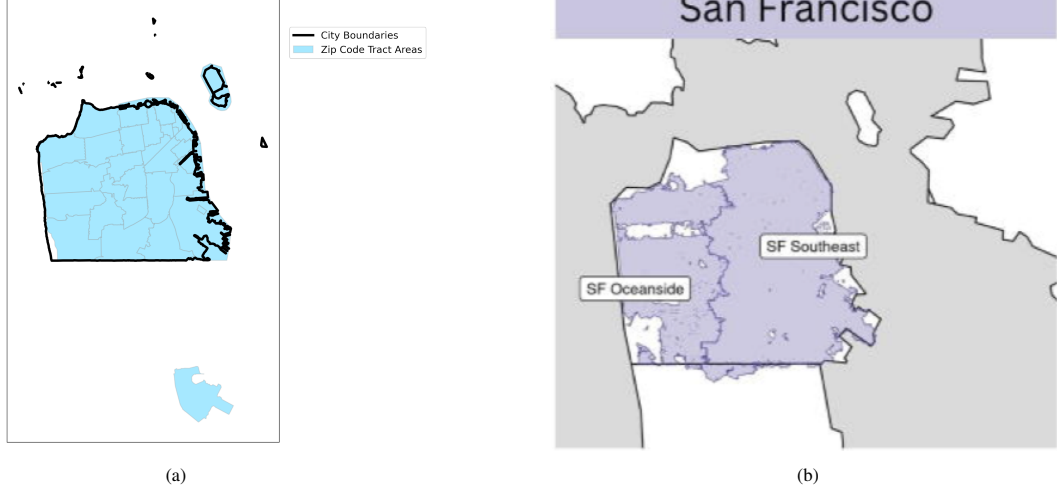


Figure 6: A side by side comparison of A) San Francisco city boundaries delineated by USPS Zip Code and B) the sewershed map of San Francisco: (Ravuri et al., 2024)

1.2. Input Variables

1.3. Target Design

The HR_i^t is obtained from the U.S. Department of Health and Human Services (HHS) COVID-19 Reported Patient Impact and Hospital Capacity by Facility dataset at HealthData.gov (United States Government Department of Health & Human Services, 2024). This dataset has weekly aggregated facility-level data for hospital utilization, derived from reports from HHS TeleTracking and reports made directly to HHS Protect from state and territorial health departments. The number of hospitalizations for each facility can be derived from the reported 7-day sum of confirmed and suspected adult COVID-19 cases, which details the total number of patients counted at that facility that week, and the 7-day coverage, which details the number of days that the number of patients were recorded at each facility that week. In the dataset, if the reported 7-day sum was less than 4, the cell was replaced with a -999999 to protect patient confidentiality. For the indicated low values, we interpolate the number of hospitalizations that week as the mean value of the surrounding weeks. If there were multiple missing values in a row, we replace the missing sum with a 2, as the average value between the 1 and 3 patients that the missing value represents. Using these reported and interpolated values, we calculate the average number of patients seen that week as the $\frac{7\text{-day-sum}}{7\text{-day-coverage}}$ for all facilities.

1.4. Error Metrics

In this section we provide the equations that define our five error metrics: 1) Accuracy, 2) Mean Square Error (MSE), 3) Weighted Mean Square Error (WMSE), 4) Brier Score, and 5) Brier Skill Score. To calculate the error metrics, we convert our forecasted categories to numerical values; For the HRT, the labels of [Large Decrease, Decrease, Stable, Increase, and Large Increase] map to [-2, -1, 0, 1, 2]. For the HCR, the labels of [Very Low, Low, Moderate, High, Very High] map to [1, 2, 3, 4, 5].

Accuracy measures the percent of predicted labels that match the true labels, and is defined as

$$\text{Accuracy} = \frac{\eta_{TP} + \eta_{TN}}{\eta_{TP} + \eta_{FP} + \eta_{TN} + \eta_{FN}}, \quad (1)$$

where $\eta_{TP}, \eta_{FP}, \eta_{TN}, \eta_{FN}$ indicate true positives, false positives, true negatives and false negatives respectively.

MSE demonstrates the magnitude of error and is defined as

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2, \quad (2)$$

where N is the size of the test set, y_i indicates the actual category label, and \hat{y}_i is the predicted category.

The WMSE weights the magnitude of the error, where the weights are the probabilities of each category.

$$\text{WMSE} = \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K f_k^{(i)} (k - \hat{y}_i)^2, \quad (3)$$

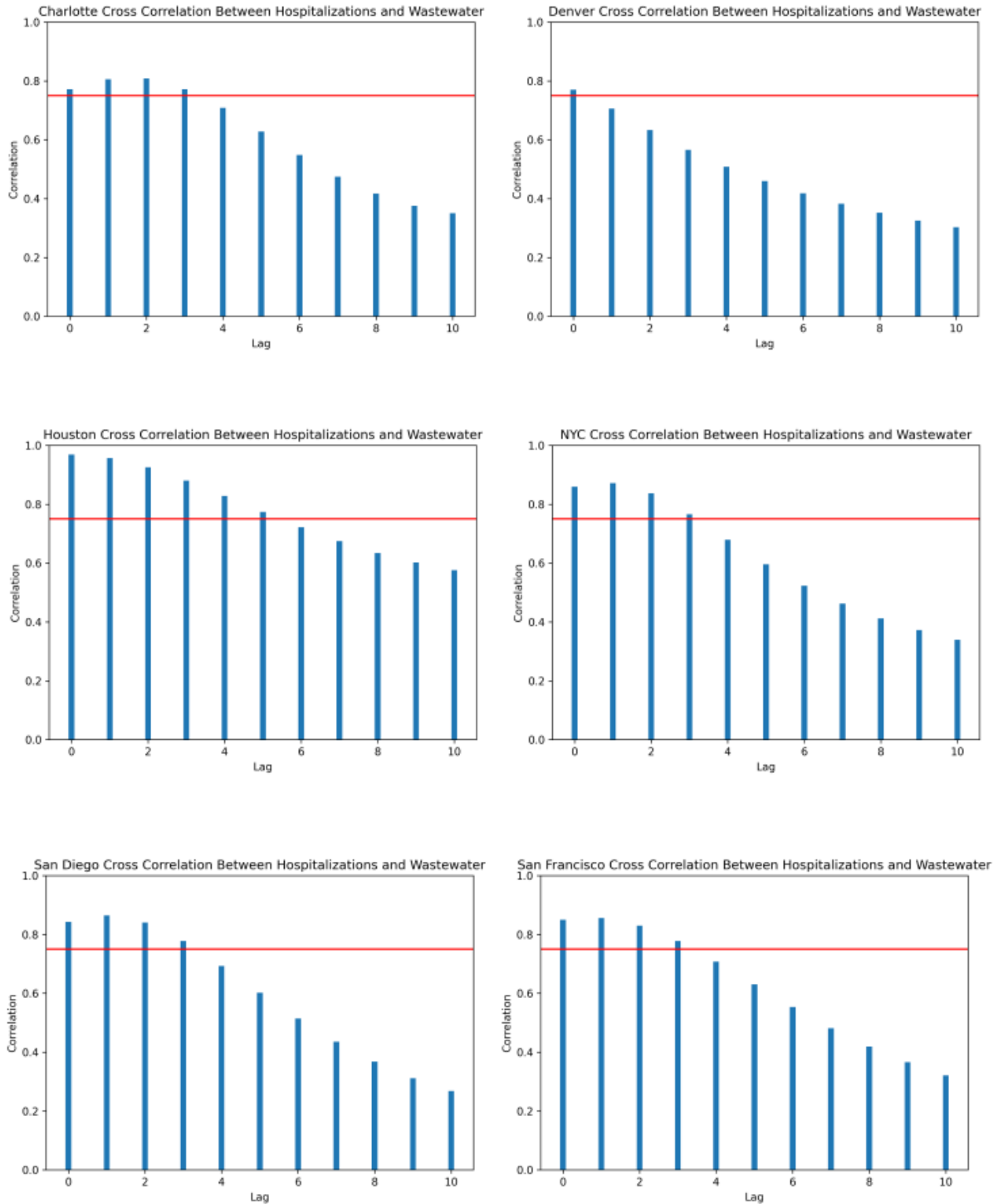


Figure 7: **Figure S2:** Cross correlations between wastewater SARS-CoV-2 rates and hospitalizations per 100,000 people in each city included in the study from January 2021 to September 2022. The red line indicates the critical threshold of 0.75 we consider for inclusion as a predictor variable.

where K denotes the set of categories, $f_k^{(i)}$ indicates the probability of k -th category being forecasted for the i -th item in the test set.

The Brier Score measures the accuracy of probabilistic forecasts. It is calculated as the mean square difference between the predicted probability of each category to the actual probability of each category. The Brier Score is defined as

$$BS = \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K (f_k^{(i)} - o_k^{(i)})^2, \quad (4)$$

where $f_k^{(i)}$ indicates the probability of k -th category being forecasted for the i -th item in the test set and $o_k^{(i)}$ is a binary variable which is equal to 1 for the true category k_i and a 0 for all other categories for the i -th item in the test set.

The Brier Skill Score, compares model performance that of a random guess, where there is a uniform probability for each category, defined as

$$BSS = \frac{BS_{ref} - BS}{BS_{ref}} \quad (5)$$

Where BS is the calculated Brier Score from our test set, and BS_{ref} is the Brier Score from of a random guess on the test set. A BSS less than 0 indicates that our predictions perform worse than the baseline, equal to 0 indicates that the score is equivalent to the baseline, and greater than one indicates that the our predictions perform better than the baseline.

References

- Census , 2020. Census zip code tabulation areas. URL: <https://www.census.gov/cgi-bin/geo/shapefiles>. accessed: 2024-07-30.
- Colorado Department of Public Health and Environment, 2024. CDPHE COVID19 Wastewater Dashboard Data . <https://data-cdphe.opendata.arcgis.com/datasets>.
- Hopkins, L., Ensor, K., Stadler, L.B., 2022. Best practices for wastewater-based epidemiology URL: <https://www.hou-wastewater-epi.org/sites/g/files/bxs4786/files/2022-06/COH-Wastewater-Epi-Best-Practices.pdf>. accessed: 21 June 2024.
- Lab, A., 2024. Sars-cov-2 wastewater san diego. URL: <https://github.com/andersen-lab>. accessed: 2024-07-10.
- North Carolina Department of Health and Human Services, 2024. NC COVID-19 Dashboard Data . <https://covid19.ncdhhs.gov/dashboard/data-behind-dashboards>.
- Ravuri, S., Burnor, E., Routledge, I., Linton, N., Thakur, M., Boehm, A., Wolfe, M., Bischel, H.N., Naughton, C.C., Yu, A.T., White, L.A., León, T.M., 2024. "real-time county-aggregated wastewater-based estimates for sars-cov-2 effective reproduction numbers". medRxiv URL: <https://www.medrxiv.org/content/early/2024/05/03/2024.05.02.24306456>, doi:10.1101/2024.05.02.24306456, arXiv:<https://www.medrxiv.org/content/early/2024/05/03/2024.05.02.24306456.full.pdf>.
- United States Government Department of Health & Human Services, 2024. COVID-19 Reported Patient Impact and Hospital Capacity by Facility. <https://healthdata.gov/d/t7zc-4t6g>.
- United States Postal Service, 2024. Zip code lookup. URL: <https://tools.usps.com/zip-code-lookup.htm>. accessed: 2024-07-10.
- US HHS, 2023. 500 cities: City boundaries. URL: <https://catalog.data.gov/dataset/500-cities-city-boundaries>. accessed: 2024-07-30.

1 Supplementary B: Results

1.1 Model Performance

This section describes the model performance for both the HCR and HRT over the 1- and 3-week forecasting window.

1.1.1 HCR

HCR 1-week forecast

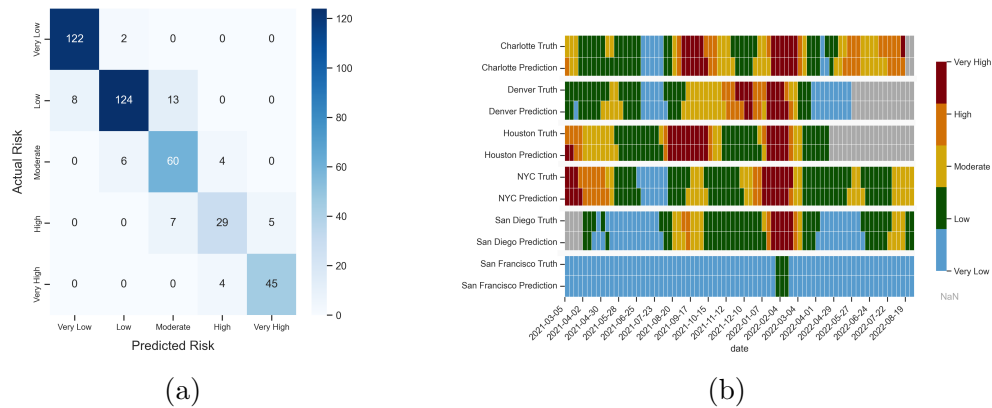


Figure 1: **Figure S1: Model performance of Hospitalization Capacity Risk (HCR) Categorization for the 1 Week forecast.** (a) Confusion Matrix detailing the accuracy of categorization for the HCR at a 1 week forecasting window. Darker colors indicate more assignment-true label combos lie in that category. (b) Hospital Capacity Risk true and predicted label for 1-week forecasts by city for January 2021 to September 2022. Color indicates predicted category. Grey squares indicate that a prediction was not made for that week due to a lack of available data.

HCR 3-week forecast

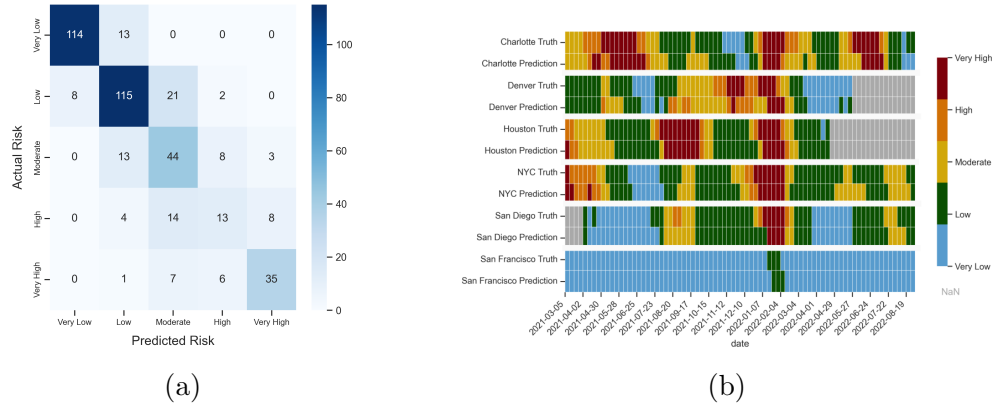


Figure 2: **Figure S2: Model performance of Hospitalization Capacity Risk (HCR) Categorization for the 3-week forecast.** (a) Confusion Matrix detailing the accuracy of categorization for the HCR over a 3 week forecasting window. Darker colors indicate more assignment-true label combos lie in that category. (b) Hospital Capacity Risk true and predicted label for 3-week forecasts by city for January 2021 to September 2022. Color indicates predicted category. Grey squares indicate that a prediction was not made for that week due to a lack of available data.

1.1.2 HRT

HRT 1-week forecast

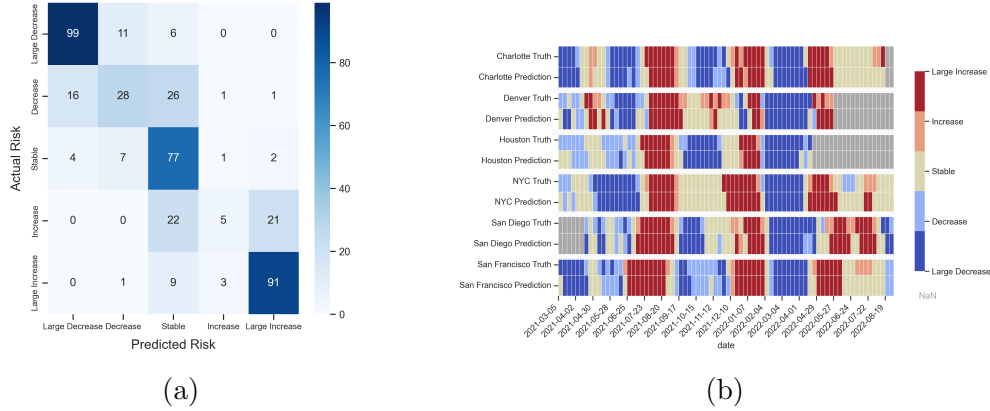


Figure 3: **Figure S3: Model performance of Hospitalization Rate Trend (HRT) Categorization for the 1 week forecast** (a) Confusion matrix detailing the accuracy of the categorization of the HRT at a 1 week forecasting window. Darker colors indicate more assignment-true label combos lie in that category. (b) City-specific Rate Trend target and 1-week prediction over time from January 2021 to September 2022. Grey squares indicate that a prediction was not made for that week due to a lack of available data.

HRT 3-week forecast

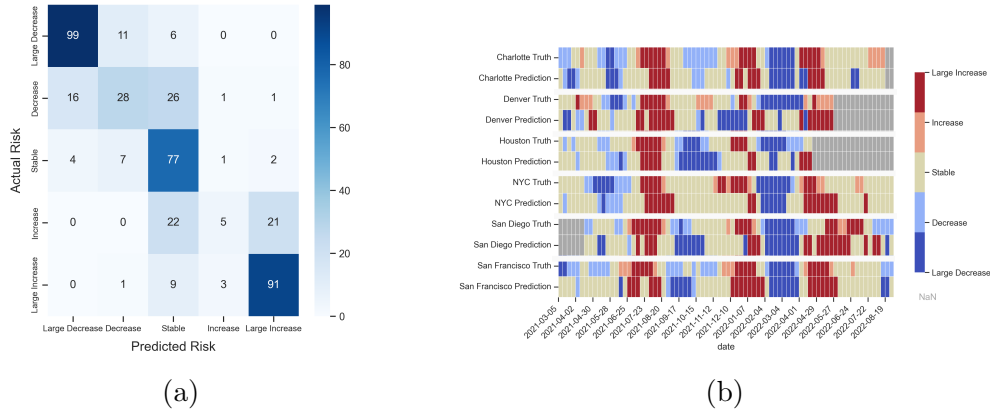


Figure 4: **Figure S4: Model performance of Hospitalization Rate Trend (HRT) Categorization.** (a) Confusion matrix detailing the accuracy of the categorization of the HRT at a 3 week forecasting window. Darker colors indicate more assignment-true label combos lie in that category. (b) City-specific Rate Trend target and 1-week prediction over time from January 2021 to September 2022. Grey squares indicate that a prediction was not made for that week due to a lack of available data.

1.2 Change Point Performance

		All Change Points			
		MSE ↓	WMSE ↓	BS ↓	BSS
HCR	1 Week	0.345	0.244	0.464	0.420
	2 Week	0.614	0.325	0.703	0.121
	3 Week	1.034	0.454	0.821	-0.027
HRT	1 Week	0.944	0.628	0.707	0.117
	2 Week	2.304	1.262	0.881	-0.101
	3 Week	2.076	1.100	0.909	-0.136

Table 1: A summary of model performance for all change points. \uparrow / \downarrow denotes if an upward or downward shift change point demonstrates better performance.

		Upward Shift				Downward Shift			
		MSE ↓	WMSE ↓	BS ↓	BSS ↑	MSE ↓	WMSE ↓	BS ↓	BSS ↑
HCR	1 Week	0.405	0.224	0.554	0.308	0.289	0.262	0.380	0.525
	2 Week	0.786	0.288	0.845	-0.056	0.457	0.359	0.574	0.282
	3 Week	1.419	0.446	0.903	-0.128	0.674	0.462	0.746	0.068
HRT	1 Week	1.195	0.651	0.694	0.133	0.711	0.607	0.719	0.102
	2 Week	2.571	1.183	0.907	-0.134	2.060	1.335	0.857	-0.071
	3 Week	2.181	1.131	0.846	-0.057	1.973	1.069	0.971	-0.213

Table 2: A summary of model performance for forecasting change points from a lower to higher category (upward shift) and a higher to lower category (downward shift). \uparrow / \downarrow denotes if an upward or downward shift change point demonstrates better performance.