# E. coli Testing: Safeguarding Public Health*

Tianrui Fu

September 27, 2024

Swimming in beaches with high E. coli levels can lead to gastrointestinal issues, rashes, and infections, particularly affecting young children, the elderly, and those with weakened immune systems. In 1998, the implementation of the Ministry of Health Beach Management Protocol project led the Toronto department to start collecting beach water quality data in 2007, aiming to reduce the incidence of waterborne diseases in the population. This paper prove that the mean E.coli is decrease by time and site location. Higher than 50% time and site location satisfy the E.coli standard.

## Table of contents

---

*Code and data are available at: https://open.toronto.ca/dataset/toronto-beaches-water-quality/.

# 1 Introduction

The increase in E.coli concentration poses a serious risk to human health. Elevated levels render beach water unsafe, impeding recreational use. While some strains of E.coli are harmless, pathogenic microorganisms may still be present. In cases of heightened concentration, healthy individuals can often recover on their own; however, vulnerable populations—such as children, the elderly, and those with compromised immune systems—are at greater risk of severe illnesses, including gastrointestinal discomfort, rashes, ear and eye infections, and, in extreme cases, bloody diarrhea, kidney failure, and even death. The rise in E.coli concentration is attributed to various factors, including runoff from animal feces carrying pathogens, aging sewage pipelines, and issues with wastewater infrastructure. Given the rapid growth of E.coli and the associated public health risks, real-time monitoring of beach water quality is essential.

In accordance with the Ministry of Health's Beach Management Protocol, the Beach Water Sampling Program for the City of Toronto was implemented to reduce the incidence of waterborne diseases. Toronto's beaches are certified under the Blue Flag program. From June to September each year, relevant departments collect water samples daily from all regulated beaches in the city to test for E.coli concentrations.

E.coli levels in freshwater are determined by counting the yellow and yellow-brown colonies that grow after applying a 0.45-micron filter on m-TEC medium and incubating it at 35.0ºC for 22-24 hours Alliance (n.d.). The water quality standard for E. coli in Ontario and federally is 200 E.coli per 100 milliliters of water, while Toronto's beach water quality standard is 100 E.coli per 100 milliliters of water *About Beach Water Quality* (n.d.). This article utilizes 22,000 data points collected by government departments from June to September between 2007 and 2024 to assess whether E.coli levels exceed Toronto's water quality standard. Additionally, it aims to determine if E.coli concentrations are influenced by changes over time and variations in site location.

# 2 Data

## 2.1 Raw Data

The data used in this paper is from Open Data Toronto and download by the Gelfand (2022). The Toronto Public Health (2024) is used to analyze whether the E.coli value in the beach during June and September is satisfy for people to swim. All the data analysis is through the R Core Team (2023), the templete from Alexander (2023) and the completed by following packages Wickham et al. (2019), Wickham et al. (2023), Firke (2023), Wickham (2016), Müller (2020), Zhu (2024) and Xie (2023). The dataset is published by Toronto Public Health and is part of The Beach Water Sampling Program, which collects E. coli values from two different beaches. The data is updated daily, and the data used for this analysis spans from June 3, 2007, to September 8, 2024. The raw dataset contains 21,882 water quality samples testing for E.coli concentrations. The dataset also includes beach itendifiers, names, sample site names, collection dates for each sample, and geographic coordinates (latitude and longitude).

## 2.2 Cleaned Data

The data cleaning is processed by R Core Team (2023). For detailed steps, please refer to the Section 5.

Table 1: E.coli quality summary by date

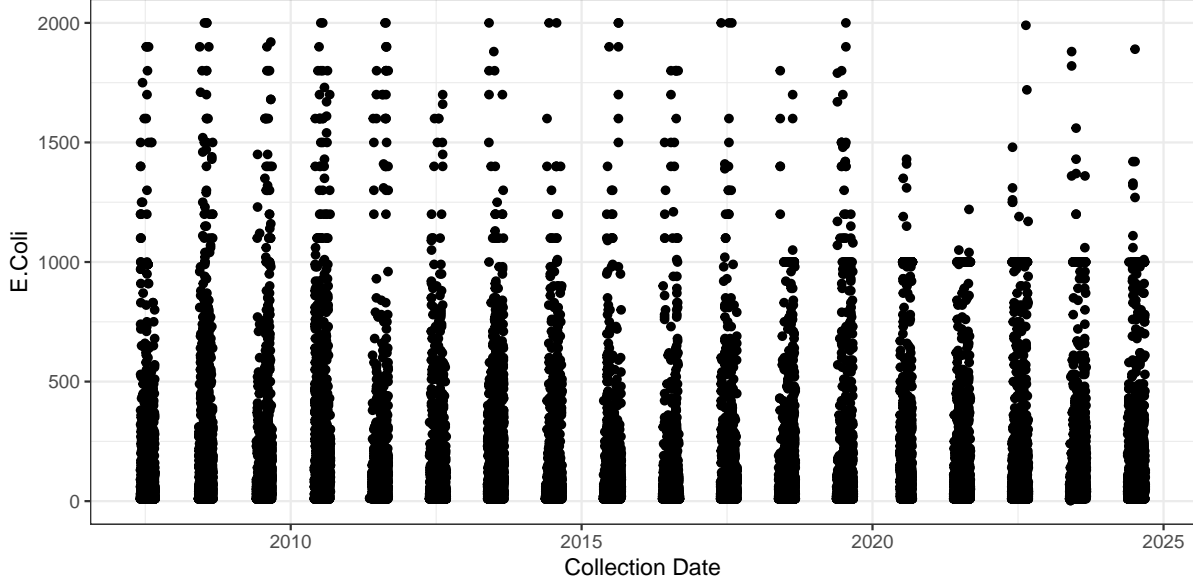| Beach Name | Date | E.coli | Site Location |
|---|---|---|---|
| Marie Curtis Park East Beach | 2024-09-01 | 230 | 29 |
| Marie Curtis Park East Beach | 2024-09-01 | 220 | 33 |
| Marie Curtis Park East Beach | 2024-09-01 | 240 | 32 |
| Marie Curtis Park East Beach | 2024-09-01 | 200 | 31 |
| Marie Curtis Park East Beach | 2024-09-01 | 200 | 30 |
| Sunnyside Beach | 2024-09-01 | 910 | 18 |
| Sunnyside Beach | 2024-09-01 | 210 | 17 |
| Sunnyside Beach | 2024-09-01 | 170 | 20 |
| Sunnyside Beach | 2024-09-01 | 80 | 21 |
| Sunnyside Beach | 2024-09-01 | 400 | 19 |
| Sunnyside Beach | 2024-09-01 | 360 | 22 |

Figure 1: Relationship between wing length and width

## 2.3 Basic Summary of cleaned Dataset

Figure 1 shows the attributes of the data, with most E.coli values being below 500, meaning there are 500 E.coli per 100 ml of water. However, the scatter plot also reveals that the data from 2020 stands out compared to other years, but a more precise graphical output is needed for a better understanding of the dataset. Therefore, in Figure 2, a bar chart was used to show the total amount of valid data collected each year. It is clear that due to the COVID-19 pandemic, the data collected by Toronto Public Health in 2020 is significantly less compared to other years. Table 2 also provides a calculation of the distribution of collected E.coli data. Based on Canada's water quality standard of 200 E.coli per 100 ml, about 78% of the data falls within this standard, meaning that approximately 78% of the days from 2007 to the present were safe for beach activities. However, based on Toronto stricter standard of 100 E.coli per 100 ml, only around 60% of the data meets this standard, meaning that about 60% of the days were safe for beachgoers. However, considering the additional variables of location and date, it is not yet possible to draw a reasonable conclusion about the overall water quality. A more in-depth analysis will be conducted in Section 3.
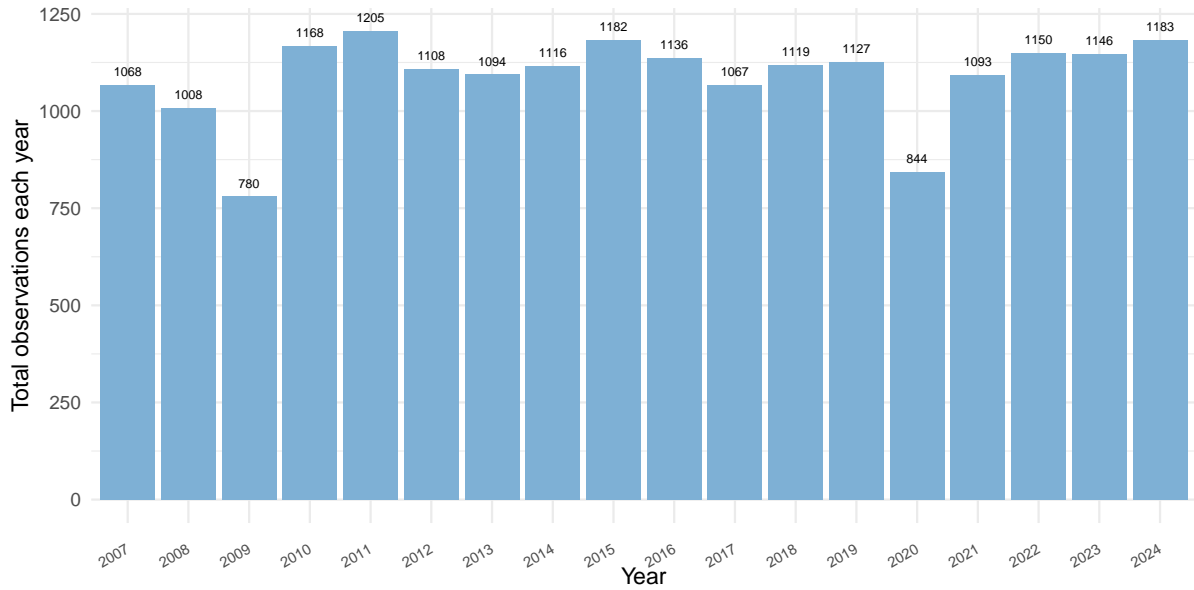
Figure 2: Number of observations by year

Table 2: Proportion of E.coli levels in different ranges

| E.Coli Levels | Sample Count | Proportion |
|---|---|---|
| 0-50 | 8588 | 43.82974 |
| 50-100 | 3609 | 18.41890 |
| 100-200 | 3009 | 15.35674 |
| >200 | 4388 | 22.39461 |

## 2.4 Dataset justification

The reason for choosing this dataset is that, during this summer, there were news reports about several incidents of people defecating on Toronto beaches, with even photos of the feces circulating. This caused some panic among people in Toronto who were planning to visit the beach. To verify whether the beach closures were indeed caused by these incidents and to reduce potential bias, as well as out of personal interest, this paper uses data published by Toronto Public Health. The dataset spans 18 years, from 2007 to the present, containing E.coli test results for beach water quality. It effectively reveals the trends in water quality over the years and helps to project future trends, reducing doubts about the validity of the analysis due to a small data sample.

# 3 Result

## 3.1 Studying the relationship between time and E. coli concentration

Though Section 2.3 finds that about 60% of the beaches meet the Toronto Beach water quality standard, we need to further explore the changes in the average E.coli concentration over the years and whether the average values truly reflect that the water quality meets the standard more than half of the time. Figure 3 shows the yearly average E.coli values. We can observe that none of the years had an average E.coli concentration that meets the stricter Toronto Beach standard, but only four years exceeded the Canadian water quality standard. The average E.coli concentration was at its lowest in 2016 which is about 107 E.coli/100ml, approaching 100 E.coli/100ml (Toronto Beach standard), and peaked in 2008, about 240 E.coli/100ml. Over time, while there was a slight increase in the average E.coli concentration from 2016 to 2020, the overall trend shows a decline, stabilizing around 150 E.coli/100ml. This trend suggests that while water quality has improved over the years, it still fluctuates and does not consistently meet the Toronto Beach water quality standard.
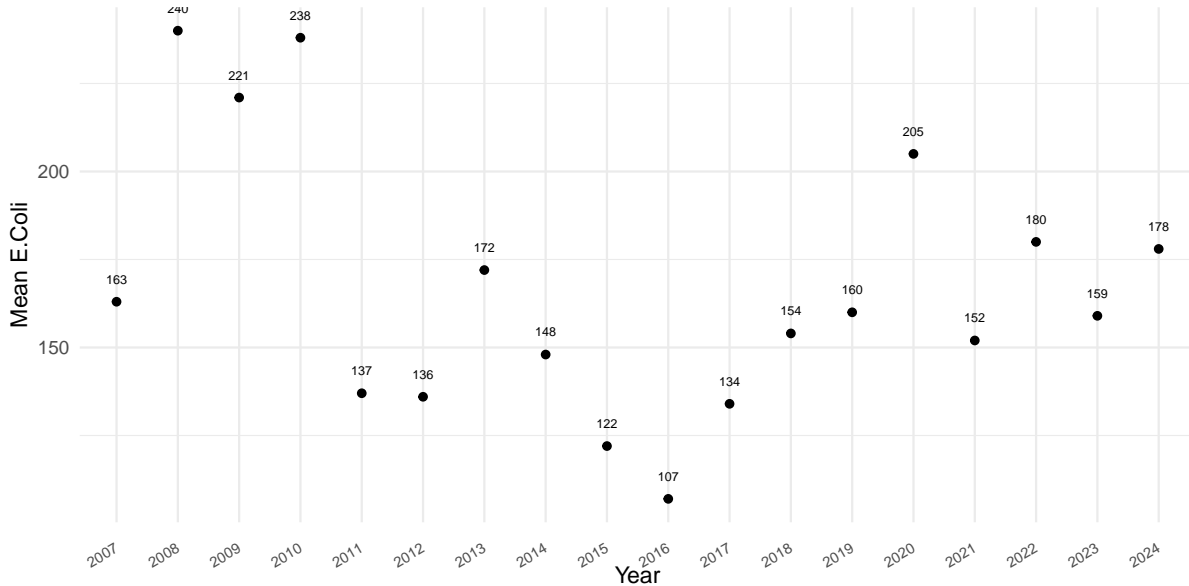


Figure 3: Mean E.coli levels by year

## 3.2 Examing the Proportion of E.coli higher than the standard level by year

To further understand the data, we calculated the proportion of E.coli levels exceeding the Toronto Beach standard each year, providing a clearer view of water quality trends at Toronto beaches. From Figure 3, we can see that the highest and lowest average E.coli concentrations, observed in 2008 and 2016 respectively, correspond to the highest and lowest proportions of

exceedance. In 2008, over 50% of the samples exceeded the Toronto standard, while in 2016, this proportion dropped to as low as 21%. Additionally, the overall proportion of exceedances has been declining over time and has stabilized in recent years.
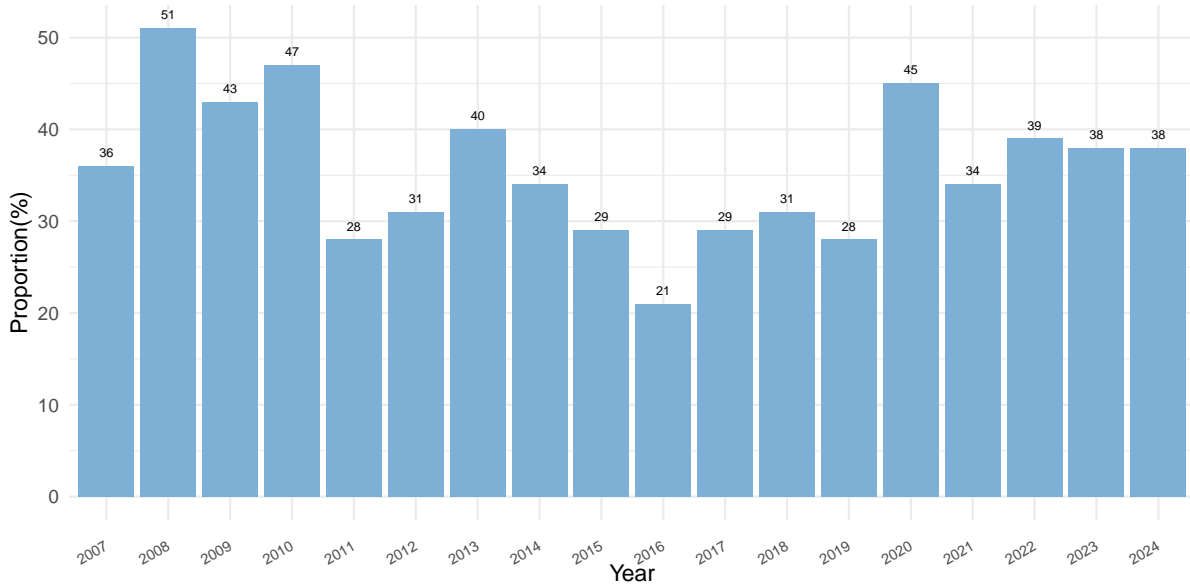


Figure 4: Proportion of higher than safe E.Coli level by year

## 3.3 Investigating the relationship with E.coli and site location

By observing the trend of average E.coli levels across testing sites, as shown in Figure 5, we can gain further insights into the data. The E.coli exceedance percentages for each site are presented in Table 3. Sites 17W-23W correspond to Sunnyside Beach, while sites 29W-33W are located at Marie Curtis Park East Beach. Although all of these sites exceed Toronto's beach standards, the E.coli levels at sites 17W-23W are generally lower than those at sites 29W-33W. The lowest E.coli levels were recorded at the 20W testing site. Table 3 so shows that the data from site 23W is noticeably lower than the other sites, making it less significant for reference. The exceedance percentage at site 20W is the lowest, at 35.09%, and the exceedance rates for most sites range between 30% and 40%.
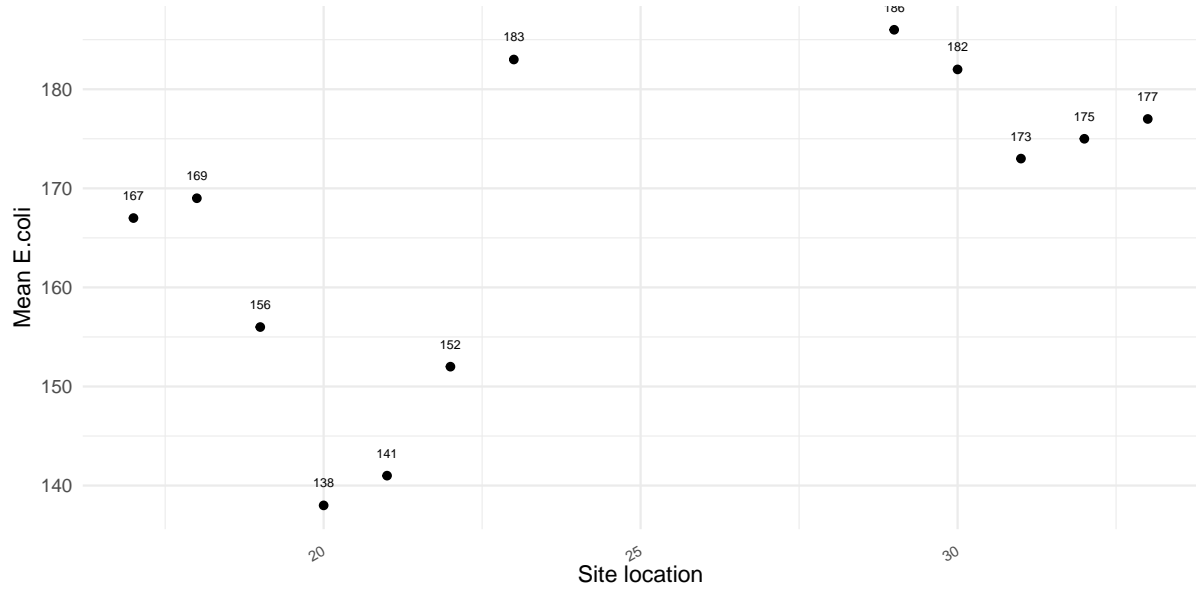
Figure 5: Mean E.Coli levels by each site

Table 3: Proportion of E.Coli levels above safe level in different site

| Site name | Total Samples | Above 100 count | Proportion(%) |
|---|---|---|---|
| 17 | 1778 | 690 | 38.80765 |
| 18 | 1774 | 681 | 38.38782 |
| 19 | 1778 | 628 | 35.32058 |
| 20 | 1782 | 554 | 31.08866 |
| 21 | 1783 | 564 | 31.63208 |
| 22 | 1770 | 625 | 35.31073 |
| 23 | 460 | 197 | 42.82609 |
| 29 | 1691 | 610 | 36.07333 |
| 30 | 1693 | 609 | 35.97165 |
| 31 | 1697 | 577 | 34.00118 |
| 32 | 1696 | 590 | 34.78774 |
| 33 | 1692 | 596 | 35.22459 |

# 4 Discussion

## 4.1 The time influence E.coli

From a temporal perspective, the data presented in the previous section's graphs and tables indicate significant changes in the water quality of Toronto's beaches from 2007 to 2024. Notably, in 2008, E.coli levels reached a peak, with average concentrations approaching 300 E.coli per 100 ml. Subsequently, E.coli levels exhibited a downward trend and became more stable. This improvement may be attributed to several factors mentioned in the introduction, including animal waste, human activities, and the aging stormwater drainage system.

The trend illustrated in Figure 3 shows a clear negative correlation between average E.coli levels and time. However, since the average value alone cannot serve as a definitive indicator of this trend, we can also consider the exceedance rates displayed in Figure 4. In 2008, the exceedance rate of E. coli surpassed 50%, indicating poor water quality from June to September of that year, which likely curtailed recreational activities and posed potential health risks. In contrast, 2016 marked the lowest levels of E.coli, aligning closely with the Toronto beach water quality standard of 100 E.coli per 100 ml, with an exceedance rate of only about 20%. This improvement may be linked to government initiatives aimed at enhancing the beach environment, as well as favorable weather conditions during the summer months.

Data collection in 2020 was limited due to restrictions imposed by the COVID-19 pandemic; however, there was no significant increase in E.coli levels compared to other years, suggesting that overall water quality remained stable. Over time, both graphs demonstrate that after experiencing two peaks, E.coli levels have generally declined and stabilized.

## 4.2 Whether the Site location influence E.coli

In addition to temporal factors, evaluating whether testing locations impact E.coli levels can further highlight differences in water quality across various measurement sites. As illustrated in Figure 5, average E.coli levels vary significantly based on the test location. Sunnyside Beach generally exhibits lower E.coli levels, indicating relatively better water quality, particularly at the 20W test site, where the exceedance rate is only 35.09%, according to Table 3. In contrast, Marie Curtis Park East Beach shows relatively higher E. coli levels, especially at the 29W and 33W test sites, where exceedance rates approach 40%.

Despite these differences among test sites, average E.coli levels at all locations still exceed Toronto's water quality standards, although they remain below Canadian water quality standards. Since average concentration alone cannot fully determine whether an area's water quality is compliant, the exceedance rates calculated in Table 3 provide additional clarity, revealing that only a small proportion of water quality measurements fall outside compliance. The elevated exceedance values contribute to the higher average levels. By calculating

9

exceedance rates by test location, the analysis becomes more insightful and evidence-based, yielding more compelling results.

## 4.3 Conclusion

Overall, the analysis of E.coli data from Toronto beaches between 2007 and 2024 has revealed trends in concentration changes over time and across various locations. While there has been a notable improvement in overall water quality, certain years and measurement sites still exhibit significant exceedances in E.coli concentrations. Moreover, the influence of time on E.coli levels cannot be fully confirmed due to the variability of factors beyond the data and the uncertainties surrounding governmental interventions.

## 4.4 Limitations and Next Step

Although the visualizations from the graphs and tables provide an overview of E.coli concentrations in Toronto beaches over the past 18 years, there remain limitations and the need for more detailed data to address these issues. The article outlines the overall concentration trend, but Toronto Public Health did not mention whether interventions were conducted in the middle years, which may explain the overall downward trend in concentration. To avoid the impact of outliers on average value plots, this study excluded large outliers at the start of the analysis. However, TPH did not specify whether all data were accurate, so the analysis after removing some data may lack reference value.

Additionally, while there were many measurement points, the beach locations were limited to only two, narrowing the scope of dataset. As a result, the data may not fully represent the trends in Toronto beach water quality, leading to potential inaccuracies in observations and analysis. Future analyses could increase the number of beaches to provide more comprehensive data and improve credibility. Moreover, by controlling for other variables and confirming outliers, tracking E.coli levels from the same source over time or by location would provide deeper insights into beach water quality, helping the government take effective intervention measures.

# 5 Appendix

In the raw data provided by Toronto Public Health, there were missing values (NA). During the data cleaning process, rows containing these NA values were completely removed to ensure that the NA values would not affect the analysis output and to simplify the analysis. Since the raw data also contained some particularly large outliers, those were also removed. We also removed the "W" in the site location name which is easy for us to simulate, test and analysis dataset. The cleaned data only includes the necessary columns for analysis, such as collection date, E.coli, site name, and beach name. Figure 1 shows the cleaned data samples, listing the results of tests conducted at different locations on the same day.

# References

*About Beach Water Quality.* n.d. https://www.toronto.ca/community-people/health-wellness-care/health-inspections-monitoring/swimsafe/beach-water-quality/about-beach-water-quality/.

Alexander, Rohan. 2023. *Telling Stories with Data: With Applications in r.* Chapman; Hall/CRC.

Alliance, Clean Lakes. n.d. *Bacteria and Our Beaches.* https://www.cleanlakesalliance.org/e-coli/.

Firke, Sam. 2023. *Janitor: Simple Tools for Examining and Cleaning Dirty Data.* https://CRAN.R-project.org/package=janitor.

Gelfand, Sharla. 2022. *Opendatatoronto: Access the City of Toronto Open Data Portal.* https://CRAN.R-project.org/package=opendatatoronto.

Müller, Kirill. 2020. *Here: A Simpler Way to Find Your Files.* https://CRAN.R-project.org/package=here.

R Core Team. 2023. *R: A Language and Environment for Statistical Computing.* Vienna, Austria: R Foundation for Statistical Computing. https://www.R-project.org/.

Toronto Public Health. 2024. *Toronto Beaches Water Quality.* https://open.toronto.ca/dataset/toronto-beaches-water-quality/.

Wickham, Hadley. 2016. *Ggplot2: Elegant Graphics for Data Analysis.* Springer-Verlag New York. https://ggplot2.tidyverse.org.

Wickham, Hadley, Mara Averick, Jennifer Bryan, Winston Chang, Lucy D'Agostino McGowan, Romain François, Garrett Grolemund, et al. 2019. "Welcome to the tidyverse." *Journal of Open Source Software* 4 (43): 1686. https://doi.org/10.21105/joss.01686.

Wickham, Hadley, Romain François, Lionel Henry, Kirill Müller, and Davis Vaughan. 2023. *Dplyr: A Grammar of Data Manipulation.* https://cran.r-project.org/package=dplyr.

Xie, Yihui. 2023. *Knitr: A General-Purpose Package for Dynamic Report Generation in r.* ttps://yihui.org/knitr/.

Zhu, Hao. 2024. *kableExtra: Construct Complex Table with 'Kable' and 'Knitr' Packages.* https://CRAN.R-project.org/package=kableExtra.