# ST117 Individual DRAFT Written Report - Task A

My WARWICK ID 5645242 (Report Pod 041)

2025-04-08

## Setting up data frames before for phase 2 part A

```r
#Import the dataframes from phase 1
df_stream_wide <- readRDS("/Users/danielguo/Desktop/University/Year 1/ST117/5645242_041_WR/df_stream_wide.rds")
df_precipitation_wide <- readRDS("/Users/danielguo/Desktop/University/Year 1/ST117/5645242_041_WR/df_precipitati
df_soil_wide <- readRDS("/Users/danielguo/Desktop/University/Year 1/ST117/5645242_041_WR/df_soil_wide.rds")

#Averaging the repeated sampling by the variable LCODE
df_stream_avg <- df_stream_wide %>%
  group_by(SDATE, SITECODE, LCODE) %>%
  summarise(across(where(is.numeric), mean, na.rm = TRUE), .groups = "drop")
```

```
## Warning: There was 1 warning in `summarise()`.
## i In argument: `across(where(is.numeric), mean, na.rm = TRUE)`.
## i In group 1: `SDATE = 1992-10-06`, `SITECODE = "T04"`, `LCODE = 1`.
## Caused by warning:
## ! The `...` argument of `across()` is deprecated as of dplyr 1.1.0.
## Supply arguments directly to `.fns` through an anonymous function instead.
##
##   # Previously
##   across(a:b, mean, na.rm = TRUE)
##
##   # Now
##   across(a:b, \(x) mean(x, na.rm = TRUE))
```

```r
#weighted averaging considering the sampling volume (VOLUME) by the variable RID
df_soil_avg <- df_soil_wide %>%
  group_by(SDATE, SITECODE, RID) %>%
  summarise(across(where(is.numeric), ~ weighted.mean(., w = VOLUME, na.rm = TRUE)), .groups = "drop")

#STAGE, VACUUM, and VOLUME are only contained in some of the datasets and may be dropped unless needed.
df_stream_avg <- df_stream_avg %>% select(-STAGE)
df_precipitation_wide <- df_precipitation_wide %>% select(-VOLUME)
df_soil_avg <- df_soil_avg %>% select(-VACUUM, -VOLUME)

#filter out the assigned data for our pod: 2002-2008: T02, T04, T06
df_stream_Afiltered <- df_stream_avg %>%
  filter(format(SDATE, "%Y") %in% 2002:2008,
         SITECODE %in% c("T02", "T04", "T06"))
df_precipitation_Afiltered <- df_precipitation_wide %>%
  filter(format(SDATE, "%Y") %in% 2002:2008,
         SITECODE %in% c("T02", "T04", "T06"))
df_soil_Afiltered <- df_soil_avg %>%
  filter(format(SDATE, "%Y") %in% 2002:2008,
         SITECODE %in% c("T02", "T04", "T06"))
```
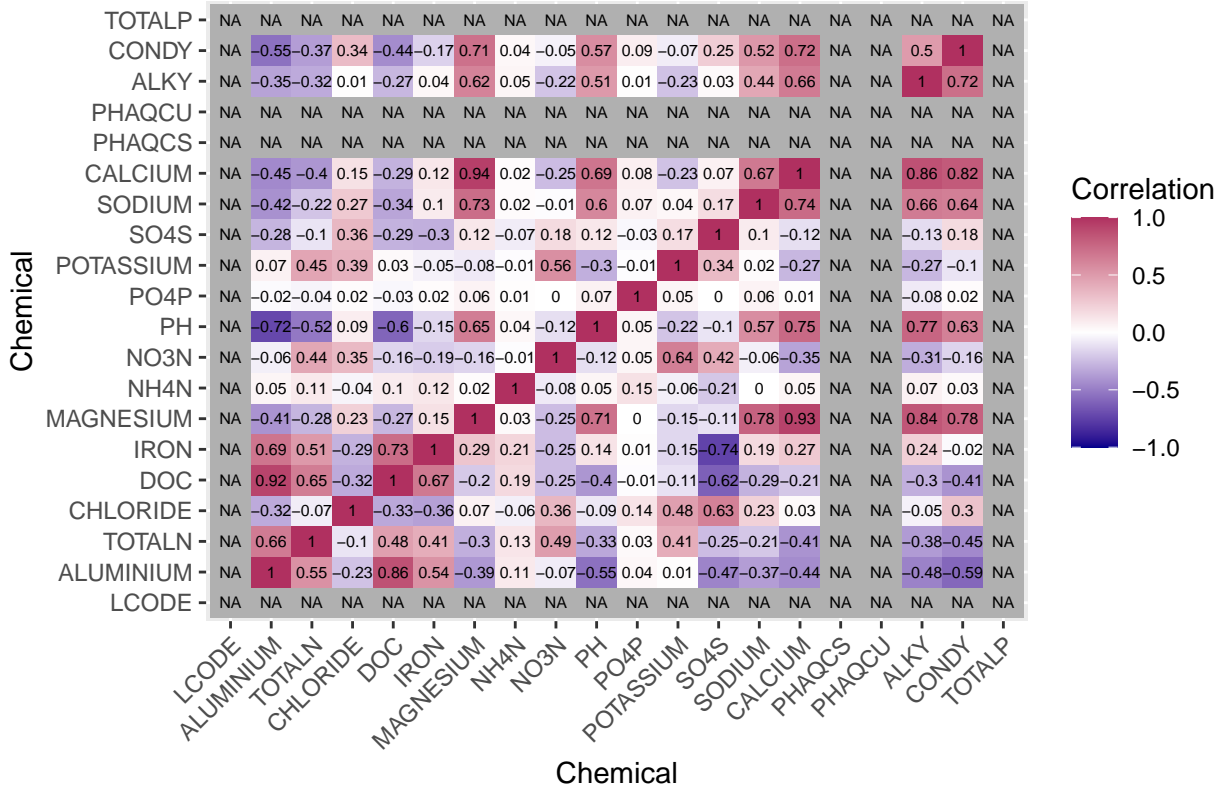
# Question 1

# Question 2

```r
plot_and_find_differences <- function(sitecode, df, dataname = "data") {
  # Filter for the selected site
  df_site <- df %>%
    filter(SITECODE == sitecode) %>%
    ungroup() %>%
    select(where(is.numeric))

  #Pearson and Spearman correlation matrices
  correlation_pearson <- suppressWarnings(round(cor(df_site, method = "pearson", use = "pairwise.complete.obs"),
  correlation_spearman <- suppressWarnings(round(cor(df_site, method = "spearman", use = "pairwise.complete.obs"

  # Combine: Pearson upper triangle, Spearman lower triangle
  correlation_combined <- correlation_pearson
  correlation_combined[lower.tri(correlation_combined)] <- correlation_spearman[lower.tri(correlation_spearman)]

  # Long format
  correlation_long <- as.data.frame(as.table(correlation_combined)) %>%
    rename(Variables1 = Var1, Variables2 = Var2, Correlation = Freq)

  # plot
  heatmap_plot <- ggplot(correlation_long, aes(Variables1, Variables2, fill = Correlation)) +
    geom_tile() +
    scale_fill_gradient2(low = "darkblue", high = "maroon",
                         midpoint = 0, limit = c(-1, 1), na.value = "grey69") +
    geom_text(aes(label = ifelse(is.na(Correlation), "NA", round(Correlation, 2))), size = 2) + #adjusting to th
    theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
    labs(
      title = paste("Correlation Heatmap for", sitecode, "in", dataname, "dataset (Pearson above, Spearman below
      x = "Chemical", y = "Chemical"
    )
   # identify large absolute differences bigger than 0,5
  abs_diff <- abs(abs(correlation_pearson) - abs(correlation_spearman))
  diff_table <- as.data.frame(as.table(abs_diff)) %>%
    rename(Variables1 = Var1, Variables2 = Var2, AbsDiff = Freq) %>%
    filter(AbsDiff > 0.5) %>%
    mutate(Site = sitecode) %>%
    arrange(desc(AbsDiff))

  return(list(
    plot = heatmap_plot,
    differences_table = diff_table
  ))
}

result_T02_stream <- plot_and_find_differences("T02", df_stream_Afiltered, "Stream Water")
result_T04_stream <- plot_and_find_differences("T04", df_stream_Afiltered, "Stream Water")
result_T06_stream <- plot_and_find_differences("T06", df_stream_Afiltered, "Stream Water")
result_T02_precipitation <- plot_and_find_differences("T02", df_precipitation_Afiltered, "Precipitation")
result_T04_precipitation <- plot_and_find_differences("T04", df_precipitation_Afiltered, "Precipitation")
result_T06_precipitation <- plot_and_find_differences("T06", df_precipitation_Afiltered, "Precipitation")
result_T02_soil <- plot_and_find_differences("T02", df_soil_Afiltered, "Soil Solution")
result_T04_soil <- plot_and_find_differences("T04", df_soil_Afiltered, "Soil Solution")
result_T06_soil <- plot_and_find_differences("T06", df_soil_Afiltered, "Soil Solution")


# View the plot
```
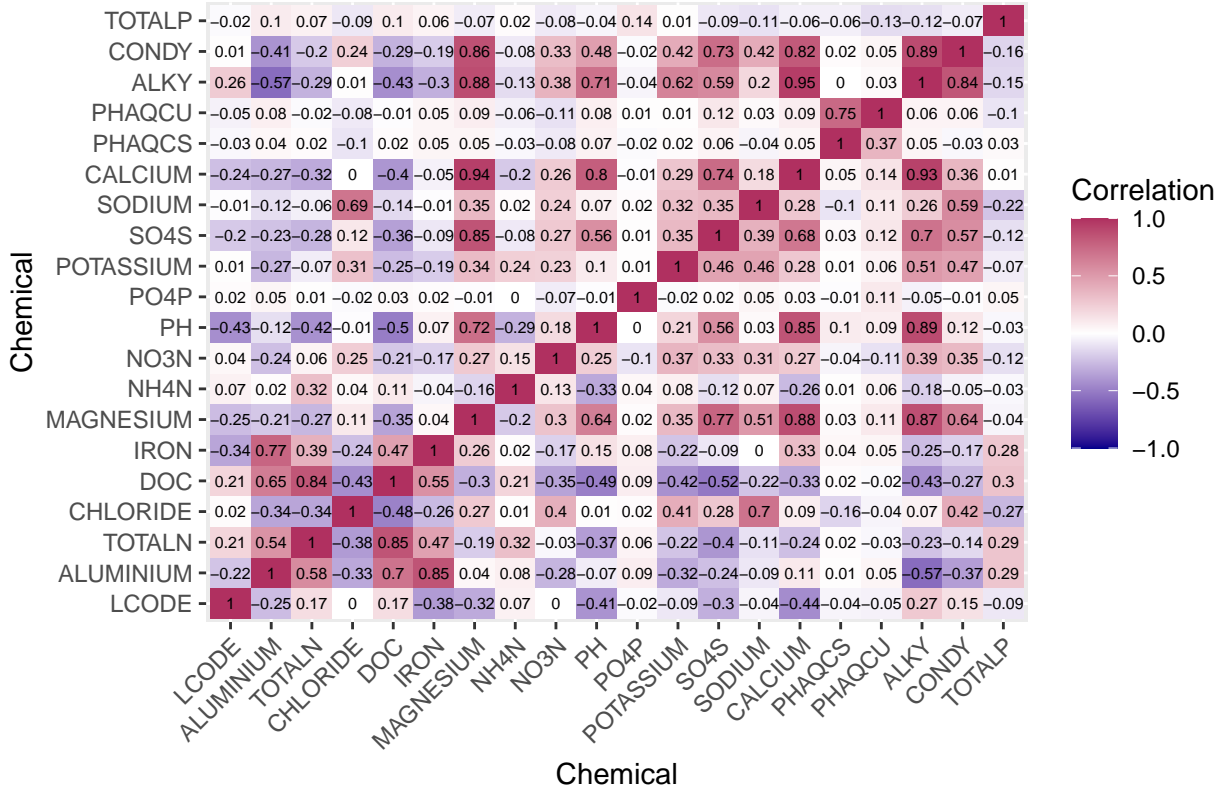
```
print(result_T02_stream$plot)
```

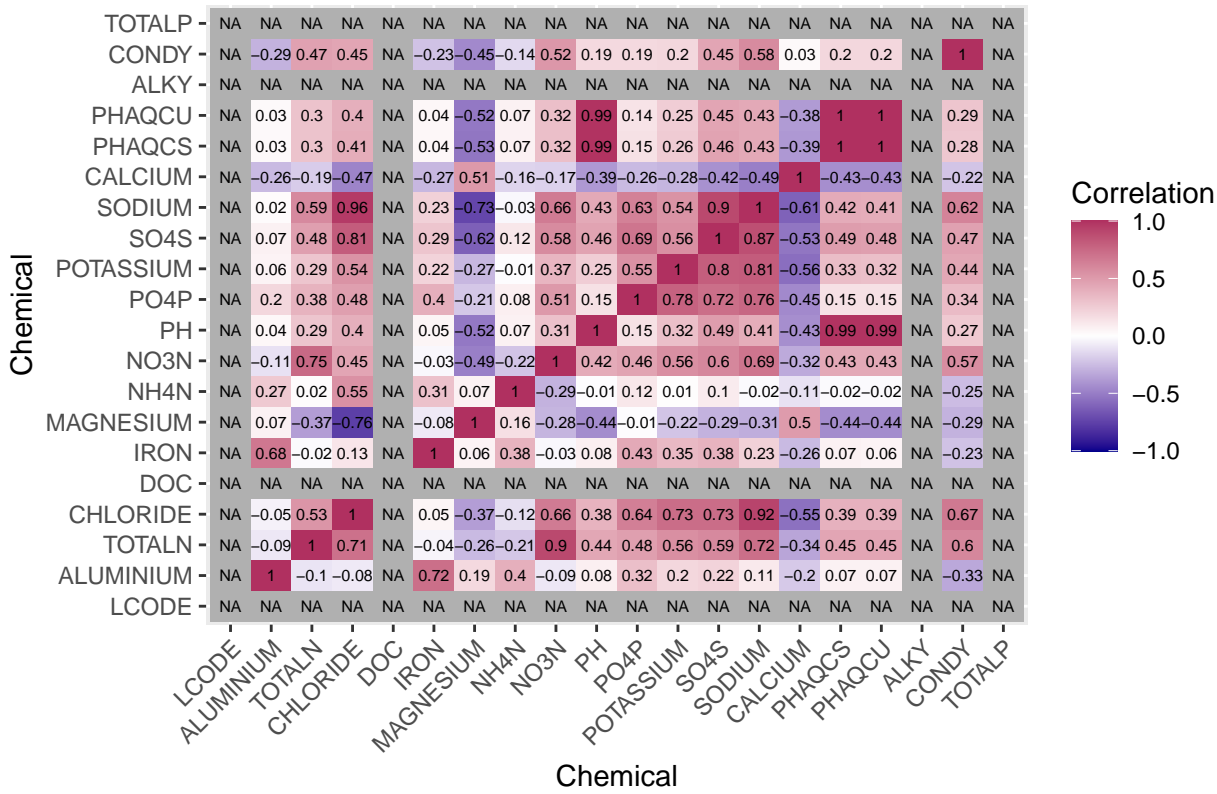## Correlation Heatmap for T02 in Stream Water dataset (Pearson abov



```
print(result_T04_stream$plot)
```

## Correlation Heatmap for T04 in Stream Water dataset (Pearson abov
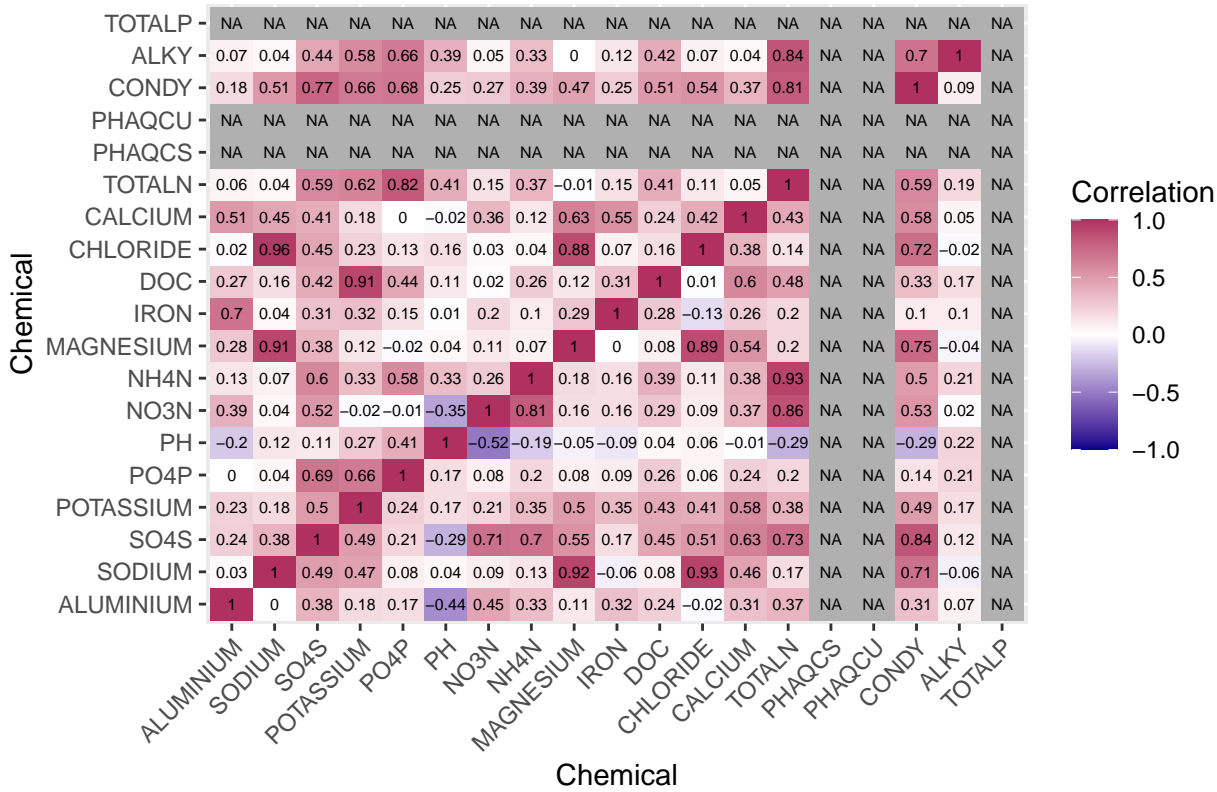


```
print(result_T06_stream$plot)
```

# Correlation Heatmap for T06 in Stream Water dataset (Pearson abov



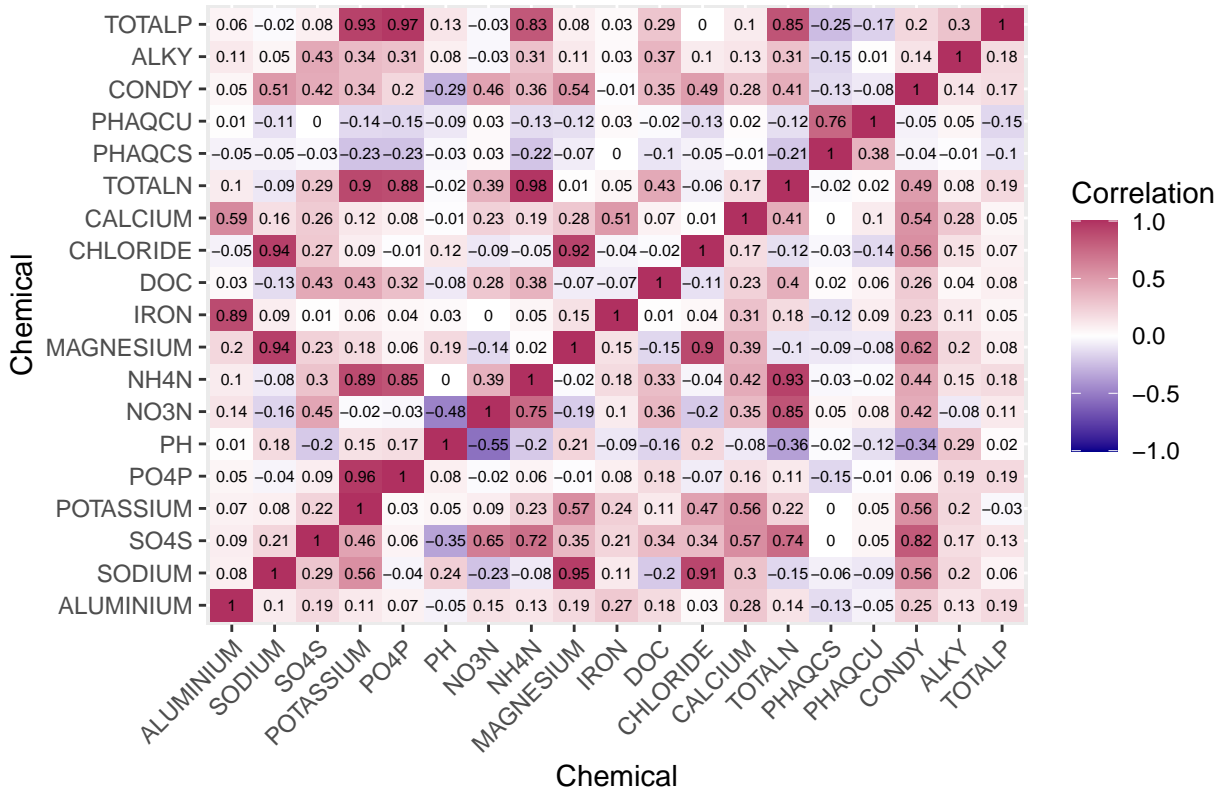```
print(result_T02_precipitation$plot)
```

# Correlation Heatmap for T02 in Precipitation dataset (Pearson above



```
print(result_T04_precipitation$plot)
```

4

## Correlation Heatmap for T04 in Precipitation dataset (Pearson above
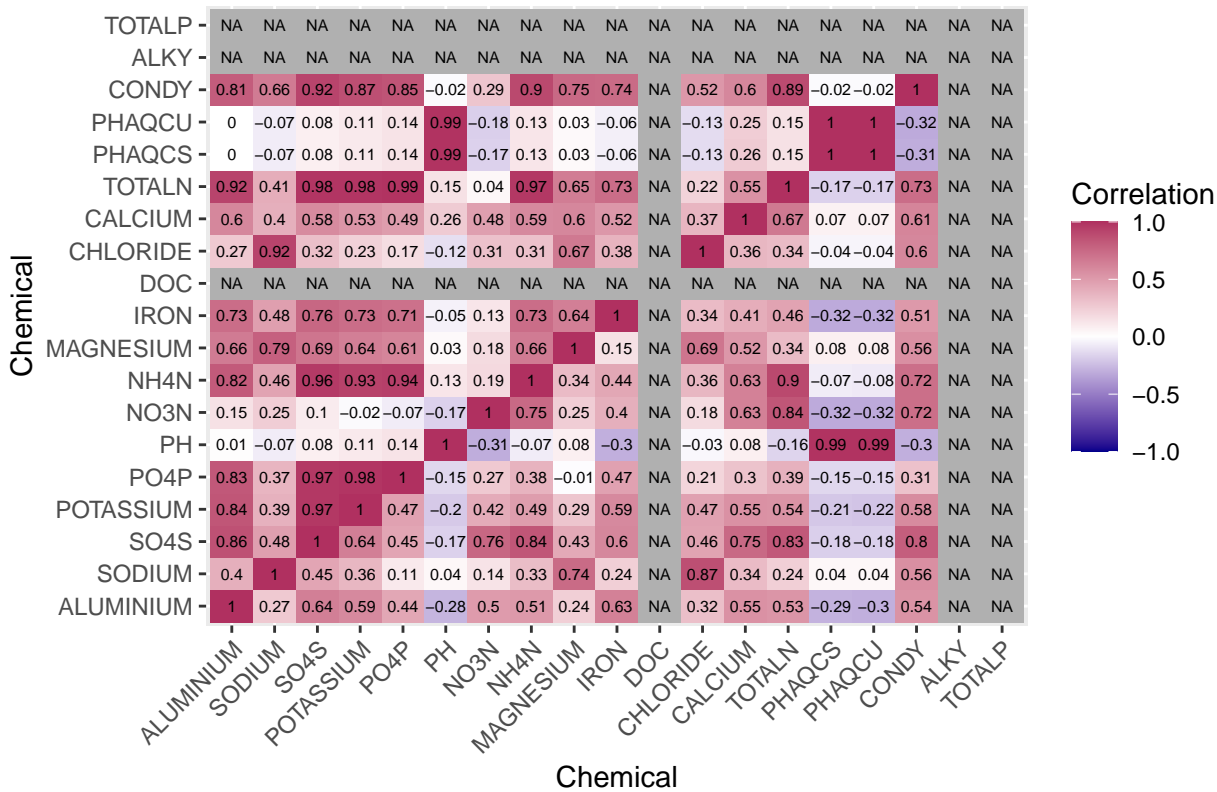


```
print(result_T06_precipitation$plot)
```
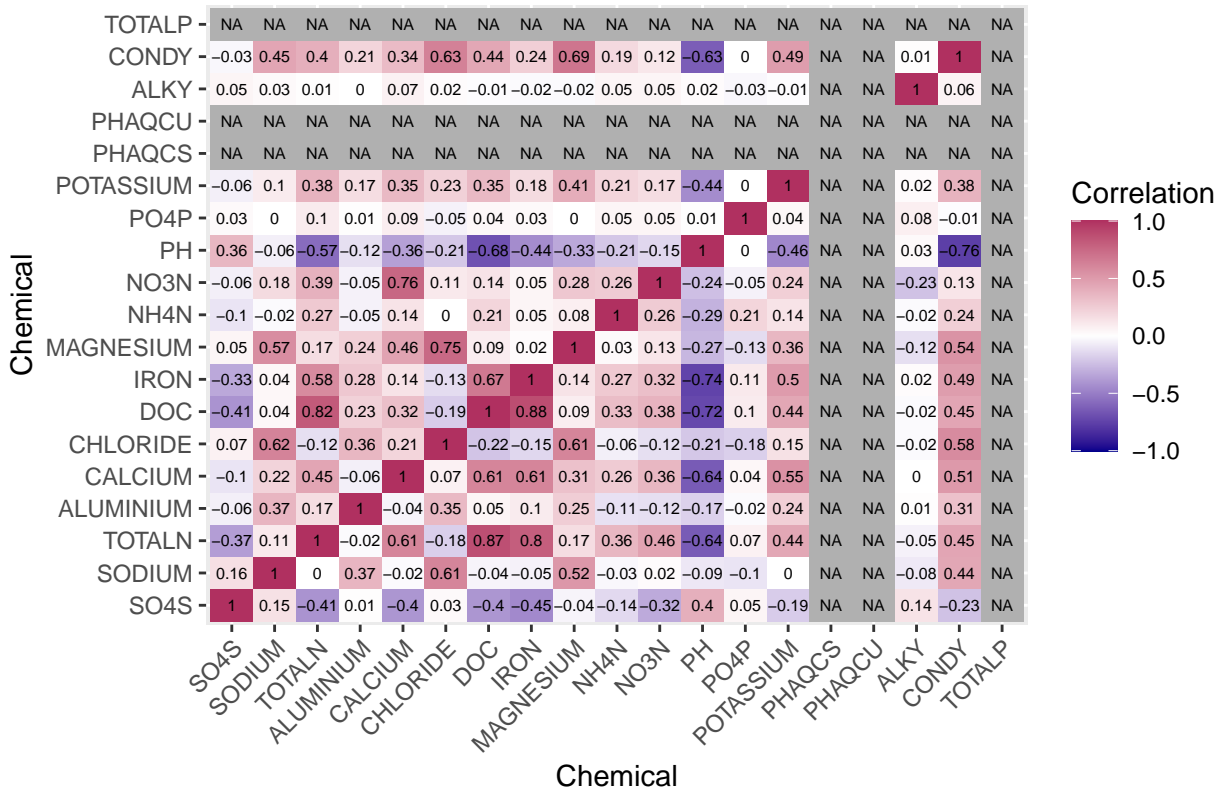
## Correlation Heatmap for T06 in Precipitation dataset (Pearson above



```
print(result_T02_soil$plot)
```
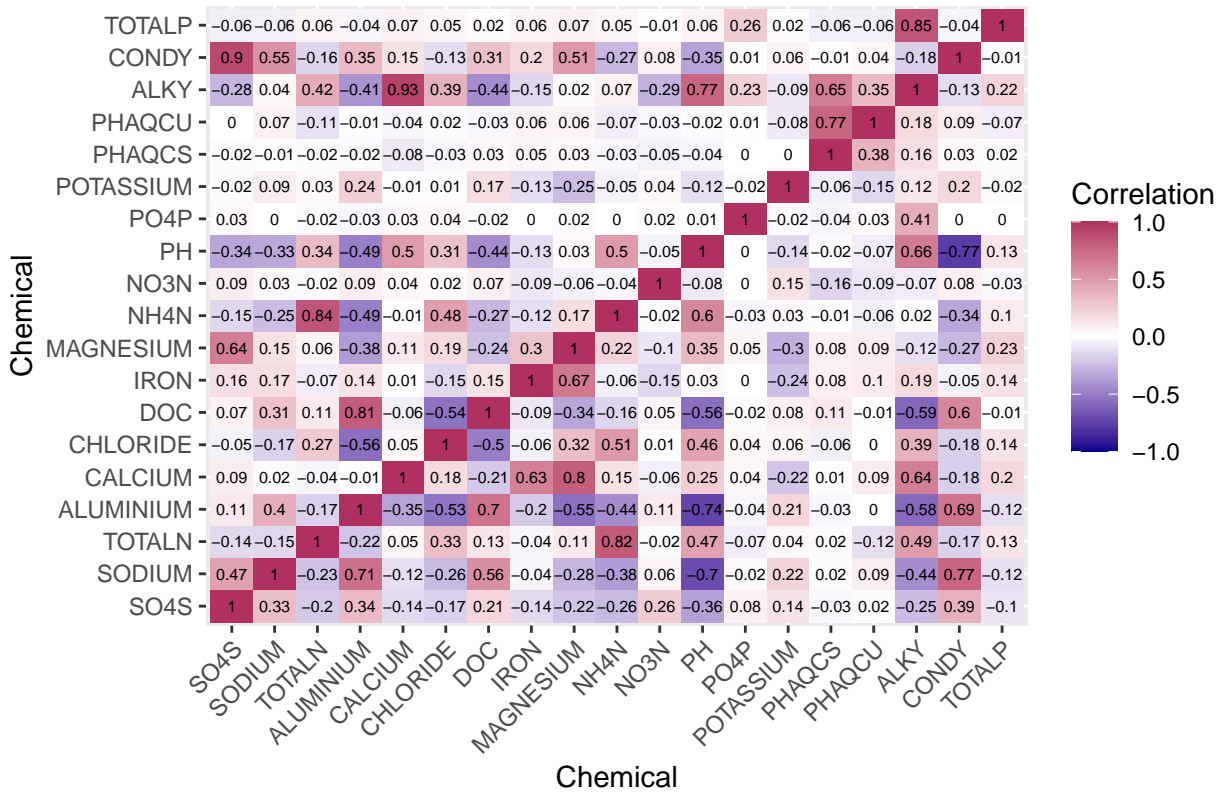
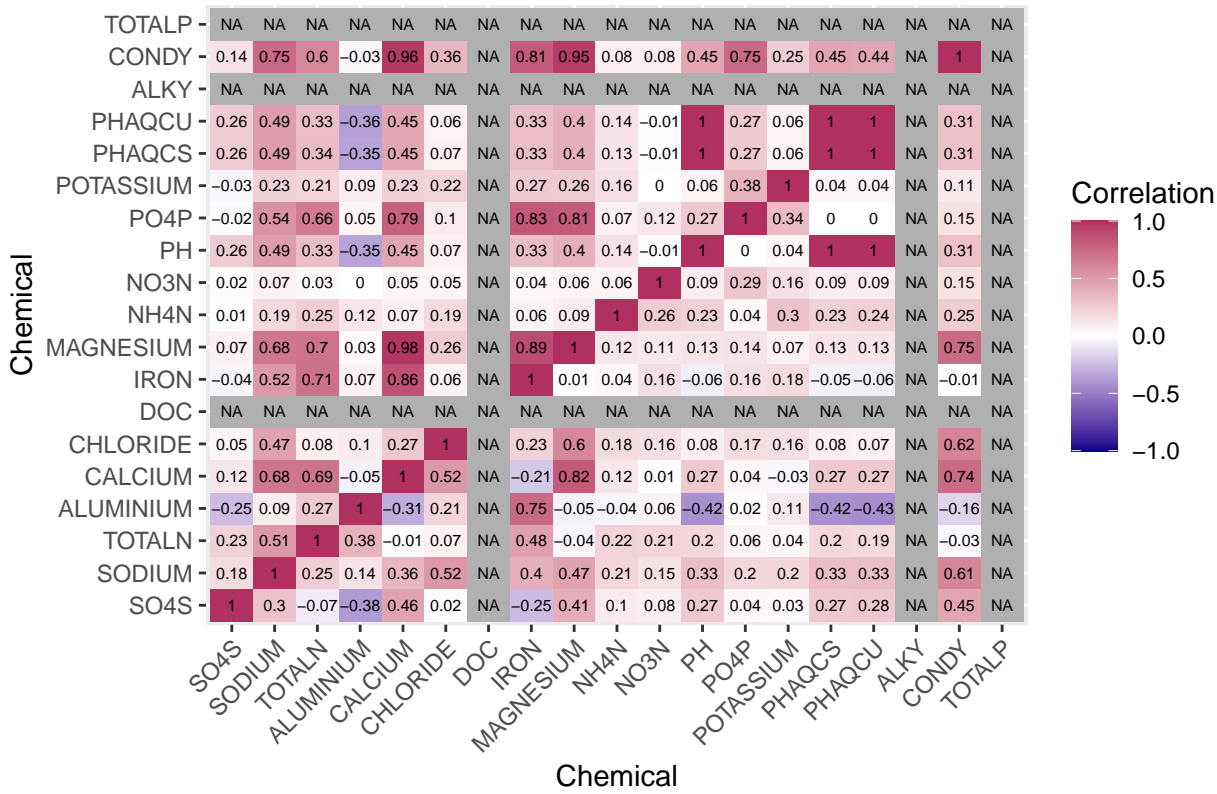## Correlation Heatmap for T02 in Soil Solution dataset (Pearson above



```
print(result_T04_soil$plot)
```

## Correlation Heatmap for T04 in Soil Solution dataset (Pearson above



```
print(result_T06_soil$plot)
```

Correlation Heatmap for T06 in Soil Solution dataset (Pearson above

## Question 3

```r
# View the table of large differences bigger than 0.5 for all sites for each data:
kable(result_T02_stream$differences_table, caption = "Large Pearson-Spearman Differences > 0.5 for T02 for Strea
```

Table 1: Large Pearson–Spearman Differences > 0.5 for T02 for
Stream Water data

| Variables1 | Variables2 | AbsDiff | Site |
|------------|------------|---------|------|

```r
kable(result_T04_stream$differences_table, caption = "Large Pearson-Spearman Differences > 0.5 for T04 for Strea
```

Table 2: Large Pearson–Spearman Differences > 0.5 for T04 for
Stream Water data

| Variables1 | Variables2 | AbsDiff | Site |
|------------|------------|---------|------|

```r
kable(result_T06_stream$differences_table, caption = "Large Pearson-Spearman Differences > 0.5 for T06 for Strea
```

Table 3: Large Pearson–Spearman Differences > 0.5 for T06 for
Stream Water data

| Variables1 | Variables2 | AbsDiff | Site |
|------------|------------|---------|------|

```r
kable(result_T02_precipitation$differences_table, caption = "Large Pearson-Spearman Differences > 0.5 for T02 fo
```

Table 4: Large Pearson–Spearman Differences > 0.5 for T02 for Precipitation data

| Variables1 | Variables2 | AbsDiff | Site |
|---|---|---|---|
| TOTALN | NO3N | 0.71 | T02 |
| NO3N | TOTALN | 0.71 | T02 |
| ALKY | TOTALN | 0.65 | T02 |
| TOTALN | ALKY | 0.65 | T02 |
| TOTALN | PO4P | 0.62 | T02 |
| PO4P | TOTALN | 0.62 | T02 |
| ALKY | CONDY | 0.61 | T02 |
| CONDY | ALKY | 0.61 | T02 |
| TOTALN | NH4N | 0.56 | T02 |
| NH4N | TOTALN | 0.56 | T02 |
| NH4N | NO3N | 0.55 | T02 |
| NO3N | NH4N | 0.55 | T02 |
| CONDY | PO4P | 0.54 | T02 |
| PO4P | CONDY | 0.54 | T02 |

```
kable(result_T04_precipitation$differences_table, caption = "Large Pearson-Spearman Differences > 0.5 for T04 fo
```

Table 5: Large Pearson–Spearman Differences > 0.5 for T04 for Precipitation data

| Variables1 | Variables2 | AbsDiff | Site |
|---|---|---|---|
| PO4P | POTASSIUM | 0.93 | T04 |
| POTASSIUM | PO4P | 0.93 | T04 |
| TOTALP | POTASSIUM | 0.90 | T04 |
| POTASSIUM | TOTALP | 0.90 | T04 |
| NH4N | PO4P | 0.79 | T04 |
| PO4P | NH4N | 0.79 | T04 |
| TOTALP | PO4P | 0.78 | T04 |
| PO4P | TOTALP | 0.78 | T04 |
| TOTALN | PO4P | 0.77 | T04 |
| PO4P | TOTALN | 0.77 | T04 |
| TOTALN | POTASSIUM | 0.68 | T04 |
| POTASSIUM | TOTALN | 0.68 | T04 |
| NH4N | POTASSIUM | 0.66 | T04 |
| POTASSIUM | NH4N | 0.66 | T04 |
| TOTALP | TOTALN | 0.66 | T04 |
| TOTALN | TOTALP | 0.66 | T04 |
| TOTALP | NH4N | 0.65 | T04 |
| NH4N | TOTALP | 0.65 | T04 |
| IRON | ALUMINIUM | 0.62 | T04 |
| ALUMINIUM | IRON | 0.62 | T04 |

```
kable(result_T06_precipitation$differences_table, caption = "Large Pearson-Spearman Differences > 0.5 for T06 fo
```

Table 6: Large Pearson–Spearman Differences > 0.5 for T06 for Precipitation data

| Variables1 | Variables2 | AbsDiff | Site |
|---|---|---|---|
| TOTALN | NO3N | 0.80 | T06 |
| NO3N | TOTALN | 0.80 | T06 |
| NO3N | SO4S | 0.66 | T06 |
| SO4S | NO3N | 0.66 | T06 |
| MAGNESIUM | PO4P | 0.60 | T06 |
| TOTALN | PO4P | 0.60 | T06 |

| Variables1 | Variables2 | AbsDiff | Site |
|---|---|---|---|
| PO4P | MAGNESIUM | 0.60 | T06 |
| PO4P | TOTALN | 0.60 | T06 |
| NH4N | NO3N | 0.56 | T06 |
| NO3N | NH4N | 0.56 | T06 |
| NH4N | PO4P | 0.56 | T06 |
| PO4P | NH4N | 0.56 | T06 |
| CONDY | PO4P | 0.54 | T06 |
| PO4P | CONDY | 0.54 | T06 |
| PO4P | SO4S | 0.52 | T06 |
| SO4S | PO4P | 0.52 | T06 |
| PO4P | POTASSIUM | 0.51 | T06 |
| POTASSIUM | PO4P | 0.51 | T06 |

```r
kable(result_T02_soil$differences_table, caption = "Large Pearson-Spearman Differences > 0.5 for T02 for Soil So
```

Table 7: Large Pearson–Spearman Differences > 0.5 for T02 for Soil Solution data

| Variables1 | Variables2 | AbsDiff | Site |
|---|---|---|---|

```r
kable(result_T04_soil$differences_table, caption = "Large Pearson-Spearman Differences > 0.5 for T04 for Soil So
```

Table 8: Large Pearson–Spearman Differences > 0.5 for T04 for Soil Solution data

| Variables1 | Variables2 | AbsDiff | Site |
|---|---|---|---|
| MAGNESIUM | CALCIUM | 0.69 | T04 |
| CALCIUM | MAGNESIUM | 0.69 | T04 |
| TOTALP | ALKY | 0.63 | T04 |
| ALKY | TOTALP | 0.63 | T04 |
| IRON | CALCIUM | 0.62 | T04 |
| CALCIUM | IRON | 0.62 | T04 |
| CONDY | SO4S | 0.51 | T04 |
| SO4S | CONDY | 0.51 | T04 |

```r
kable(result_T06_soil$differences_table, caption = "Large Pearson-Spearman Differences > 0.5 for T06 for Soil So
```

Table 9: Large Pearson–Spearman Differences > 0.5 for T06 for Soil Solution data

| Variables1 | Variables2 | AbsDiff | Site |
|---|---|---|---|
| MAGNESIUM | IRON | 0.88 | T06 |
| IRON | MAGNESIUM | 0.88 | T06 |
| CONDY | IRON | 0.80 | T06 |
| IRON | CONDY | 0.80 | T06 |
| PO4P | CALCIUM | 0.75 | T06 |
| CALCIUM | PO4P | 0.75 | T06 |
| CALCIUM | TOTALN | 0.68 | T06 |
| IRON | ALUMINIUM | 0.68 | T06 |
| TOTALN | CALCIUM | 0.68 | T06 |
| ALUMINIUM | IRON | 0.68 | T06 |
| PO4P | MAGNESIUM | 0.67 | T06 |
| MAGNESIUM | PO4P | 0.67 | T06 |
| PO4P | IRON | 0.67 | T06 |
| IRON | PO4P | 0.67 | T06 |

| Variables1 | Variables2 | AbsDiff | Site |
|---|---|---|---|
| MAGNESIUM | TOTALN | 0.66 | T06 |
| TOTALN | MAGNESIUM | 0.66 | T06 |
| IRON | CALCIUM | 0.65 | T06 |
| CALCIUM | IRON | 0.65 | T06 |
| PO4P | TOTALN | 0.60 | T06 |
| TOTALN | PO4P | 0.60 | T06 |
| CONDY | PO4P | 0.60 | T06 |
| PO4P | CONDY | 0.60 | T06 |
| CONDY | TOTALN | 0.57 | T06 |
| TOTALN | CONDY | 0.57 | T06 |

# Question 4

| Variables1 | Variables2 | AbsDiff | Site |
|---|---|---|---|
| MAGNESIUM | TOTALN | 0.66 | T06 |
| TOTALN | MAGNESIUM | 0.66 | T06 |
| IRON | CALCIUM | 0.65 | T06 |
| CALCIUM | IRON | 0.65 | T06 |
| PO4P | TOTALN | 0.60 | T06 |
| TOTALN | PO4P | 0.60 | T06 |
| CONDY | PO4P | 0.60 | T06 |
| PO4P | CONDY | 0.60 | T06 |
| CONDY | TOTALN | 0.57 | T06 |
| TOTALN | CONDY | 0.57 | T06 |