

MINUTES
Data Integrity Work Group (DIWG)

The meeting was recorded for note taking purposes.

Wednesday, September 21, 2022
10:00AM-12:50PM

[Meeting Materials Link](#)

ACTIONS

- ✓ At the next DIWG Meeting, Durga Ghosh will demonstrate how to access the data for and view the QA samples.
- ✓ If you have any topics for the next meeting, send to Amy (agoldfischer@chesapeakebay.net) and copy Durga (dghosh@chesapeakebay.net) and Cindy (cindy.johnson@deq.virginia.gov)
- ✓ Please provide any job opportunities within your own organization/agency/jurisdiction to Breck Sullivan (bsullivan@chesapeakebay.net) to distribute within the Chesapeake Bay Program (CBP) and CBP's Scientific and Technical Advisory Committee (STAC) who can share within their networks
- ✓ If interested in attending the Monitoring Kick-Off meeting, which will occur in the second week of January with a background webinar in mid-December, please contact Breck Sullivan (bsullivan@chesapeakebay.net)
- ✓ If you have any further feedback on the cluster analysis work, contact Elgin Perry (eperry@chesapeake.net)

Introductions & Announcements

All

Participants:

Durga Ghosh (USGS), Cindy Johnson (VA DEQ), Breck Sullivan (USGS), Carl Friedrichs (VIMS), Carol Cain (MD DNR), Elgin Perry (consultant), Heather Wright (ODU), Kevin Minga (ODU), Lexis Carter (ODU), Suzanne Doughten (ODU), Jaclyn Mantell (CBL), Jay Armstrong (DCLS), Jerry Frank (CBL), Kim Blodnikar (CBL), Kristen Heyer (MD DNR), Lara Phillips (MDH), Meighan Wisswell (VA DEQ), Mike Mueller (Fairfax County), Mike Mallonee (ICPRB), Pamela Higgins (PA DEP), Renee Karrh (MD DNR), Tammy Zimmerman (USGS – PA), Taylor Hughes (PA DEP), Tyler Shenk (SRBC), Betty Neikirk (VIMS), Mark Brickner (PA DEP), Jake Kilczewski (CBL), Keri Maull (DNREC), Najma Khokhar (MDE), Sergio Huerta (DNREC), Liz Chudoba (Alliance for the Chesapeake Bay)

Monitoring and Laboratory Analysis Updates

All

Cindy Johnson (VA DEQ) reminded the group of the National Water Quality Monitoring Council's 13th [National Monitoring Conference](#) which will take place on April 24-28, 2023 in Virginia Beach, VA.

Cindy gave the monitoring and laboratory updates for Virginia. They had some boat issues but seem to be holding their own for the most part. Jay Armstrong said the Department of Consolidated Laboratory Services (DCLS) are operating in full capacity. Betty Neikirk (VIMS) said things are going well at VIMS and they are on schedule. Their work in the Rappahannock as well as their continuous stations will conclude in October, 2022.

Heather Wright (ODU) said that they missed a few stations for the September cruise due to weather but otherwise things are going smoothly with lab analyses and sampling.

Kristen Heyer (MD DNR) said all of their programs are on schedule. Last week when they were out on the mainstem they took out the new Bay Program director (Dr. Kandis Boyd) so that Dr. Boyd could see monitoring at one of the stations by the Bay Bridge. For nontidal sampling they might be a little short because it hasn't rained everywhere. Bay tributaries, mainstem, data flow, and continuous monitoring are mostly on track for sampling. Jerry Frank said that the Chesapeake Biological Laboratory, one of the Maryland labs, didn't have any serious setbacks.

Lara Phillips gave the update from Maryland Department of Health and said they had some shortage of staff including losing their nutrients lead. They're trying to fill that position. Otherwise, they are meeting turnaround times as usual. Things are going well despite the staff shortage.

Pamela Higgins gave the update for PA. Pam stated they were struggling through staffing vacancies and turnover. They are also working through some supply chain issues. They've had to convert their Gas Chromatography/Mass Spectrometry instruments to running on hydrogen generators because of a helium shortage.

There was no update from West Virginia.

Tammy Zimmerman (PA USGS) reported that in PA the monitoring has been going well. There are a couple sites that they didn't get localized rain in the right places leading to them being a couple samples short. They are also struggling with supply chain issues, particularly with getting certain filters. Cindy Johnson stated their biggest supply chain issue at VA DEQ is getting plastic bottles for sampling.

Tyler Shenk (SRBC) stated Susquehanna River Basin Commission's updates are similar to Tammy's experience. They did get some rain and they're catching up on storms but could use a few more high flow site samples.

Jerry Frank commented in the chat that there were lots of supply chain issues.

Blind Audits

Jerry Frank

Jerry Frank (CBL) stated that the final report for FY22 will be distributed by the end of this month. Samples for the next round will be distributed within the first two weeks of October. Jerry thanked everyone for being timely with their response for FY22.

Updates on field audits

Durga Ghosh

Durga Ghosh (USGS) reported on the status of field audits. Durga said that she was trying to gather feedback for audits. Starting in 2020 there was a disruption to audits. She was able to look at QA samples but that was it. This year, people came up with great ideas for audits, particularly Doug Moyer and Tammy Zimmerman. The tentative plan is to divide audits into external and internal audits. All audits will be done by a given USGS agency, because they use USGS procedures for nontidal monitoring. For internal audits, USGS will audit itself. For external audits, USGS will audit non-USGS partners. For the fiscal year, the first audit will be USGS PA and SRBC. They're doing side by side analysis for 3 of their stations that have been swapped between the two agencies. They were already planning to do a side by side analysis and Tammy asked if Durga could also audit. Cindy is having an internal annual review for their sampling and Durga had wanted to go and watch this training because it would be interesting to incorporate lessons learned from there into other programs. That is tentatively proposed for October 18th. Durga hoped to be able to do an audit on that same day but it may need to be done on a different day due to time constraints. In early November, Durga hopes to do the lab audit for the Upper Potomac Riverkeeper, who are doing IDEX analysis for bacteria. They have two components; the sampling component which is done in DC, and the lab component which is done in MD. Durga is trying to coordinate their sampling audit as well. Durga was hoping to have the Arundel Rivers Federation field audit done this year but there was a major managerial and organizational revamping and they had to resubmit their QAPP. They just got it approved two weeks ago and could not schedule an audit for this year. They will have to push back their audit until the next sampling season but Durga promised that they could do it prior to the season start so that they can be at Tier 3 status starting next year.

Durga brought up bacterial monitoring for the citizen science network and how it's always a bit of an issue because of constrictions that exist for different jurisdictions. MD does things differently from VA and because the monitoring network is spread out through various jurisdictions regions Chesapeake Monitoring Cooperative (CMC) is having concerns. There will have to be further discussions and hopefully starting next year they can come to a consensus. Lastly, Durga said that for the QA samples she typically reviews in summer and fall, she got through the QA samples across the watershed to try to compare and see if there are any consistent outliers to determine any consistent issues. She has not yet looked at this year's data but she gave a shoutout to Mike Mallonee for helping Durga with that and they now have that process streamlined. At the next DIWG meeting, Durga will not only share the data, but also will show how to access the data and view the QA samples.

Coordinated Split Sample Program

Mike Mallonee

[Tributary Split Sample linked here](#) – June was added.

[Mainstem Split Sample linked here](#) – May was added.

For PC, DENRC and VIMS didn't participate in the May sampling for the mainstem.

Cluster analysis work

Elgin Perry

- *Elgin Perry will provide an overview of the Baytrends and GAMs, and describe how the cluster analysis method solves issues with categorizing GAM trends, and variations on the cluster analysis approach.*

Elgin Perry is a statistics consultant to the Bay program and is part of a team called the Bay Oxygen Research Group (BORG) that works on developing data analysis methods for the data that this group generates. It's very important to have the Data Integrity Workgroup take a look at what the BORG is doing and see if they have any concerns. A couple of years ago Elgin presented to the DIWG on baytrends, which use Generalized Additive Models (GAMs).

Since that time, the type of trend analysis Elgin works on has progressed from linear trend analysis to curvilinear lines which can capture things like if the water quality parameters are responding to flow. The three components of the model are a long term trend component, a seasonal component, and the interaction between the trend and seasonal component. The cluster analysis uses predictions from the total model fit, which uses all three of those components. Soon after the start of generating baytrends results, they discovered that it was very tedious to make connections between stations to see if they had similar or dissimilar results. So Elgin automated the process of looking at trend lines with the computer, which brought him to cluster analysis.

This process uses a GAM which fits cubic spline functions which are able to bend with the data. Elgin takes that wavy trend line and gets predictions from the 15th of each month. When the model fits the seasonal term, it does it using day of the year. Elgin has standardized it to the 15th day of the month to get rid of little variations from sampling differences. He does that for each year in the period of record and for each station in a collection of stations in the area of interest, and stores these predictions as a 3D data structure with stations, years and months on different axes. It's necessary to get down to two dimensions before doing cluster analysis. One way to do that is if doing a cluster analysis of stations looking at the long term trend line, you could average over the months and get one number for each station. Elgin showed a table where the stations are in the rows and the columns are different years. He explained a process to compute the Euclidian difference between two rows and that would give a measure of how close any two rows are. That helps find the stations that are closest in trend. Then you start adding stations according to how similar they are to another station, and build up a group called a dendrogram. Elgin uses agglomerative clustering. Another approach is divisive clustering. Agglomerative starts with individual stations and builds them into groups, rather than starting with the whole set and splitting it apart. Elgin uses a method called

Ward's method, which tries to minimize the variance within each cluster group. Once he does the cluster, he displays the results in a dendrogram. He also produces a profile box which allows us to assess what the character of each group is and why the groups differ from one another. The software is set up so the user chooses the years, the months, and the stations that are to be analyzed. It's easy to switch from one tributary to another just by choosing a different set of stations. For parameters like Dissolved Oxygen (DO) you can look at the months where it's primarily a concern (May-September). Elgin calls the things that he clusters items, and he calls the years (the information used to compare the items) the profile. If you wanted with the software, you could cluster the years using the stations as a profile. If you think about the fact that you can switch those things around and that Elgin starts with a 3-Dimensional matrix you can see there's getting to be a large number of different clusters possible. The software is flexible. It also allows the user to decide whether to scale the data or not. The user has to specify a number of groups to interpret because the software starts with stations as individual items and ends with them in one group. New software features include automatic group labeling and color assignments.

Elgin showed a table that shows how you can consider different questions about the data depending on how you specify the item, profile and scaling. For the first row (in the "Variations on Clustering" table which has stations as the item, year as profile, and no scaling), when you specify that type of clustering it's usually the long term mean that differentiates one group from another. Elgin calls that the status cluster. It gives the relative status of the stations that are being grouped by the long term mean. If you want to focus on the shape of the trend over time, it's best to subtract the mean from each station. This is called the mean adjustment. This forces the software to focus more on the shape of the long-term trend curve, and often that means does the water quality respond the same way to flow at two different stations, for example. For the third row on the table (station is the item and month is the profile) if you mean adjust that, it will compare the stations according to the shape of the seasonal signal. Being able to look at seasonality might be of interest to people studying whether climate change has an effect on the data. The software also allows you to use two dimensions of the 3D array to define an item. You can define station by year as an item and let month be the profile. When you do that you can look at whether seasonality changes at a given station over time.

Elgin said you have to choose somewhere along that progression that is the number of groups you want to assess. You can also ask it to cluster the years, according to the stations. When you're looking at a single tributary it's like looking at the longitudinal profile of the estuary from the fall line down to the bay and looking at whether or not you see the same profile year after year. Sometimes you can see things like whether high flow years differentiate from low flow of years in the profile down to the estuary.

Elgin showed different examples in the James River of different types of clustering. He showed a dendrogram output with the stations grouped in a tree formation. In the example, he chose to interpret 4 groups, but also showed what would have happened if he had asked for 3 groups and for 5 groups. The automatic labeler computes the long term mean for each group and it ranks those from 1 to however many groups you have in this

case, 1 to 4. It creates a label that uses the water quality variable, water quality parameter acronym in this case TN, then it appends to that the rank of the long term mean. So we have TN1-TN4 as labels. Elgin showed the plot for this and explained that it also showed an example where the automatic labeler is not working that well because it showed 2 with very similar looking long term means which should be more differentiated. The mean was similar but the shape was different. Elgin showed a map showing that the low TN is mostly at the lower end of the estuary. He pointed out that this trend in the James River is different as compared to the trends that he saw with the other 4 western shore tributaries. The Patuxent, Potomac, Rappahannock and York systems have a monotonically decreasing trend in TN, going from the fall line down to the Bay. In the James, it starts off up at Richmond pretty low and it increases and gets to the highest point right there where the Appomattox comes in and then from there down, it is a decreasing trend. So that's just the way the James is different as compared to others. This difference could be explained by a high TN load coming down the Appomattox and that increases the TN concentration in the lower tidal fresh section of the James. Or it could be that there are some local sources of nitrogen in that area. Downstream from Richmond, that's a unique feature of the James River that can be distinguished with this clustering approach and so that's the complete first set of clustering.

Elgin then went over his second cluster analysis, which still has stations as the items to be clustered and is still clustering according to a year for the long-term trend. However, this time he did the mean adjustment. All of the stations are centered around 0 because the mean has been subtracted out. This rearranges the cluster analysis. The flow information is displayed at the bottom of that figure. The red triangle indicates a low flow year (in the lower 25% of years); the orange indicates the flow was between 25-50% of years, light blue is above average flow (50-75%), and dark blue is high flow (above 75% of years). High TN values tend to coincide with low flow years. The red group is called DEC1 and -F1. The automatic labeler takes the difference between the beginning and end of the period of record, and if it's negative, it calls it declining, and if it's positive, it calls it increasing. Then it ranks those in order. That's the first part of the label for this trends cluster analysis. The second part of the label has either a negative or positive F indicating a negative or positive relationship with flow. This relationship is obtained by doing a Spearman correlation with the mean water quality within each group with the mean log flow of each year. It ranks those as well to show strength of association with flow. This is another case where the labeling process isn't perfect. The red (on the graph) has a more positive relationship with flow – if you look at the flow years 2003-2004, those are high flow years, but the program gave it a negative correlation. Elgin thinks the program gave that a negative correlation because there is an overall decreasing trend in the group, and it got confused sorting the long term trend and the flow relationship. Weak AI doesn't always do the best job of interpreting the data. Elgin stated these imperfections to show the necessity of still keeping a critical eye on the results.

Elgin then showed a seasonal cluster grouping stations across year. The strongest seasonal signal is high in February and low in July. However, there is another group that is low in April and high in September. The labeling process ranks them to the degree of seasonality. It also adds the month of the year where you have the highest value of the

parameter, and the month of the year you have the lowest. It was puzzling to Elgin why there was a group with high TN in the fall and why there would be one at the mouth of the river and the rest in the tidal fresh. If TN has a local source in the upper tidal fresh, that would make sense that concentration would go up in September in low flow conditions.

Elgin showed a graph where he defined station year as item. This graph was very messy. The profile is the month, using mean adjusted data. It also shows one group that has a signal high in February and low in July, and another group that's high in September and low in April. He then showed the data in a different format to demonstrate where seasonality at stations is changing and said the group may be able to interpret the data. He said it has something to do with wastewater treatment plants – with nitrogen, for example, when the advanced nitrogen treatment was first introduced it worked better in the summer than winter, which induced a seasonal cycle in the data. Over time they got it to work better in all seasons, which reduced the seasonal cycle.

Elgin showed another cluster that used year as the item, and station as profile. It is the profile of the estuary. He interpreted four groups. This analysis showed the estuary wide, the greatest improvements in nitrogen in the James River have occurred over the last decade, as well as greatest improvements occurring in the lower tidal fresh stations. The analysis showed high flow years had greater nitrogen flowing to the lower stations.

Elgin wrapped up his presentation by explaining what he liked about the cluster analysis method, his concerns with the method, and plans for the future. A positive is that you can organize complicated GAMs results to tell a story. The method allows the data to self-organize without pre-conceived notions. It allows exploration of new concepts such as changing seasonality and changing spatial trends. His concerns includes the GAMs smoothing, which sometimes is faulty and may overestimate or underestimate to keep a smooth line and de-emphasize differences among neighboring years. Auto-labeling may interfere with critical thinking. Elgin has turned his code over to TetraTech to build it into a package to add onto baytrends. Adding flow adjustments in the GAMs is another next step, as well as adding mixing curves as an interpretive tool. Allowing for 2 term profile is another potential next step, as well as adding more scaling options to see the cluster analysis on a more normal data scale. Another potential next step is adding a statistical testing process for groups, and another is using classification and regression analysis to try to explain differences among the groups. Another future step is adding water quality parameters as a fourth dimension, trying to see if stations are similar or not according to that profile.

Durga asked, do you have any concerns about integrating water quality parameters or would it be very straightforward? Elgin responded that currently he only works with one at a time but you could work with a group of water quality parameters. If you did that the first issue would be scaling because they all come with their own order of magnitude for the values. That would maybe take a z score type scaling where you subtract the mean and divide by the standard deviation in order to do that. Durga responded that it seems like a great thing to do and that it would be visually compelling to see that data if we

could. Elgin said he put that on the table for TetraTech to look at. They are taking that to the next level.

Carl Friedrichs commented that this information benefits from being reviewed more than once. He said he likes the good news plot (slide 18) where it shows the most recent years in purple have the lowest values in total surface nitrogen. Elgin said that's a mixture of high flow and low flow years. Carl responded that he thinks that goes to what he's seen discussed that sometimes that the results from the fall lines in terms of nutrient input isn't as positive as wished for, that the watershed isn't improving as much as wished for but data like this shows that the tidal components of the estuaries, the tributaries, and he thinks it holds for the bay as well, in many ways are getting much better. In terms of explanation, maybe the below fall lines nutrient inputs to the bay are really important and that's where clean-up has been most successful, at major sources beneath fall lines like cities and water treatment plants. There's a lot of good news in the coastal plain that is missed with popular press attention on farther parts of watershed. It seems like the tidal waters are getting better more than you'd expect from the trends in the above fall line watershed where most of the really tight data on nutrient input is coming from.

Elgin responded that he thinks most of that improvement is reflected in the nutrients. Often when you move up to chlorophyll and DO, you don't see those improvements, which is a bit of a conundrum for us. With our bottom-up model of how to improve the Chesapeake Bay, we get improvement at the low level but when we move to the next levels, we don't see that. The Rappahannock and the Patuxent are examples of this.

Carl said he agreed and that complications include climate warming, which is hard on the oxygen, since water can't hold as much oxygen. For chlorophyll, Carl's theory is we're doing so well cleaning the water of solids, TSS is going down, that the light limitation is going down which means we have more productive waters and nutrients can get used up more easily and produce chlorophyll. We're almost hurting ourselves with too much of a good thing.

Elgin said he is open to further feedback on if these models are doing justice to the data.

Breck commented in the chat that the CBP Web Team is currently going through a process to revamp the chesapeakebay.net. They are hoping to post these results and ones for other tributaries on the Integrated Trends and Analysis (ITAT) webpage once that is done, and are also looking to include some of the results in the Tributary Summaries updates. <https://cast.chesapeakebay.net/Home/TMDLTracking#tributaryRptsSection>

Citizen Monitoring Updates

Liz Chudoba

Liz Chudoba gave the citizen science monitoring update because Alex Fries was out. The Arundel Rivers Tier 3 Quality Assurance Project Plan (QAPP) was approved. They postponed scheduling their field audit until the spring of next year, so all of their 2022 data will be Tier 2, and at very beginning of the 2023 sampling season, they will schedule a field audit to get them back up to Tier 3 in 2023.

Topics for Next DI Meeting

All

Breck Sullivan reminded the group of the effort to enhance monitoring networks. This includes lab analysis and all the monitoring samples that this group does. There will be a kick-off meeting in fall (**update: it will be in January**) with monitoring representatives from different jurisdictions and agencies to find resources to address the recommendations provided by this group and others. Breck or another member of the monitoring team can give an update at the next DIWG meeting. Breck also welcomed members of the DIWG group to attend this kick-off meeting. The meeting will discuss different areas people could invest in and what actions monitoring representatives could be interested in and take part in, and what follow up is needed to make sure the recommendations go through. Cindy requested to be added to the meeting.

Breck also said she heard from a couple people that there's shortage in staffing and there are open positions. Is this something that's held up in HR or are there listings available that Breck could send through the CBP network? CBP STAC representatives are interested in finding out how to provide green job opportunities for their students.

12:50 Adjourn