**Capitol Lake and Puget Sound.
An Analysis of the Use and Misuse of the Budd Inlet Model.**

## 2. THE COMPUTER GETS MANY WRONG ANSWERS.

Appendix G2 of the original TMDL Report presents 38 pages comparing the Budd Inlet Model's output with the observed water quality parameters that were used to calibrate it (TMDL Appendix, 2012). There are three pages for each of the Appendix G stations highlighted in Figure 1-2, portraying observed and calculated conditions at the surface, bottom, and a depth midway between surface and bottom. Figure 2-1 shows a typical example, this one for the dissolved oxygen levels in the bottom water at station BI-6 in West Bay (the station nearest the dam). These pages enable us to estimate how many of the calculations were dead-on accurate.

Two features are evident. First is that the computer's graph (dark line) follows the general trend of the observed data (open circles) quite faithfully between January 25 and September 15, 1997. Second, if every calculation were accurate, the graph would go through every one of the open circles. It does not. It "misses the mark" by a wide margin in some cases, by a narrow margin in others, and in some cases (where it touches the circles) it is accurate.
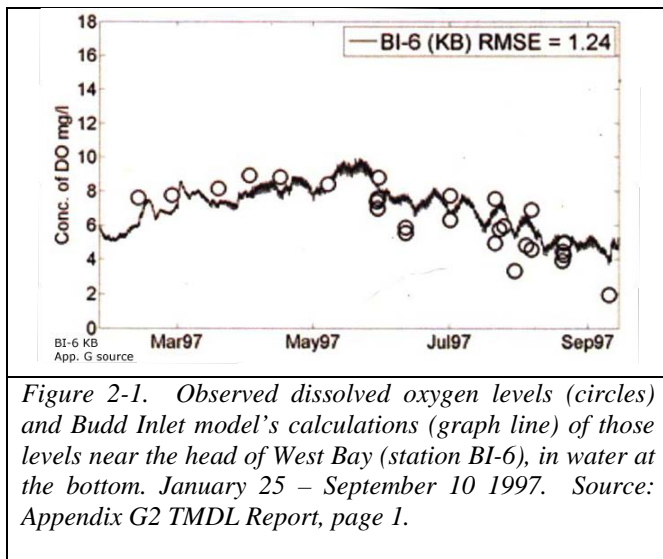
*Figure 2-1. Observed dissolved oxygen levels (circles) and Budd Inlet model's calculations (graph line) of those levels near the head of West Bay (station BI-6), in water at the bottom. January 25 – September 10 1997. Source: Appendix G2 TMDL Report, page 1.*

That is the fact to always bear in mind; *the computer often gets wrong answers.*

I used the following method to estimate the overall accuracy of the Budd Inlet Model's calculations.

## 2a. Methods.  Counting Right Answers.

The data points in the Appendix graphs are at the exact centers of the circles shown there. These circles are about 0.875 mg/L in diameter. If the graph fails to touch ("misses") the circle, the computer's answer in that case is in error by at least 0.44 mg/L (the circle's radius). That is about twice the critical value (0.2 mg/L) used in judging whether a water quality standards violation has been detected, in many cases.

I examined each of the dissolved oxygen graphs in Appendix G2 (36 graphs; 3 depths for each of 12 stations) for visual determination of whether the computer graph missed the observed data point circle, "hit" it, or was undeterminable (not clearly a hit or miss). To qualify as a "hit," the graph had to touch the exact top or bottom of the data circle or pass

through it.  A grazing contact was scored as a "miss;" the graph was close in that case but the top or bottom (over the center) of the data circle was not in contact with the graph on the date of the observation.  An example is shown in Figure 2-2 for station BF-3 surface water (near Boston Harbor).

## 2b. Results. "Hits and Misses."

Figure 2-3 shows the pattern of computer "hits" and "misses" at all stations, three depths per station.  At best over all, the computer's calculations matched observed DO's about 80% of the time in bottom water at sites BI-4 (mouth of West Bay) and BE-2 (center Budd Inlet near the Tamoshan area). At worst, calculations matched the observed values in bottom waters only about 20% of the time at BI-6 and BI-2 (West and East Bays) and BC-2 (Gull Harbor area).  Overall, the calculations were accurate in roughly 40-50% of cases.

## 2c.  Discussion. Hits and Misses.

As a tool for showing broad trends, the Budd Inlet Model is useful.  It is not capable of telling us, however, the exact value of every dissolved oxygen level – every depth, every six minutes[1], every location – for half a year.  Yet the modelers base their most important claims on an apparent assumption



Figure 2-2. Assessment of calculated "hits" and "misses" of observed data circles by the Budd Inlet Model for dissolved oxgyen concentrations in surface water at station BF-3 (near Boston Harbor) by the method described in the text.  Hits ("H" in upper row), misses and undeterminables ("M" and "?" in lower row) show 13 accurate, 13 inaccurate and 1 undeterminable calculation.  Source Appendix G2  p. 36 TMDL Report.



Figure 2-3.  Accuracy of the Budd Inlet model.  Bars show the per cent of calculations that correctly identified observed DO values (counting all "indeterminable" scores as "hits") by stations from south to north in Budd Inlet.  Data from graphs in Appendix G2 TMDL Report.

tion that it is really that accurate.  For example, if the computer finds that a calculated DO level is only 0.2 mg/L below the standard of 6.0 mg/L that prevails over most of Budd Inlet, the real-life water at that site is said by the modelers to be "in violation" of that standard.  This must stem from their assumption that the computer really does "get it
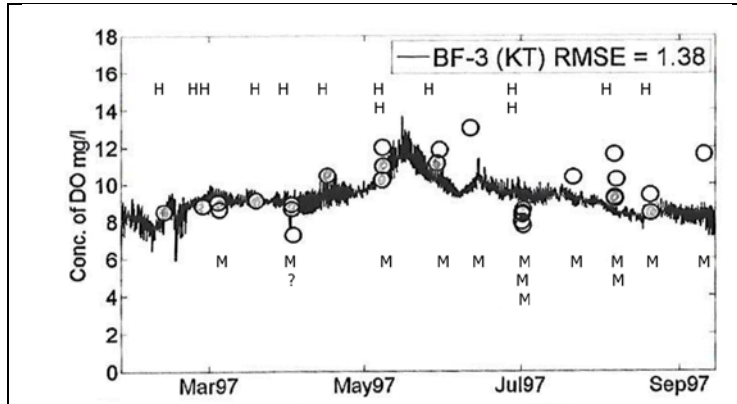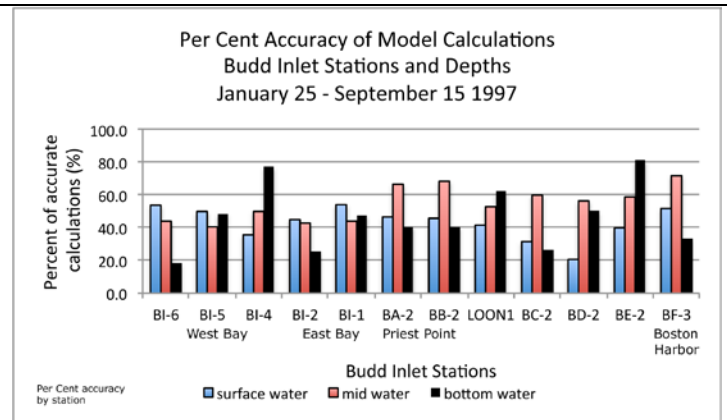
---

[1] Iteration interval given on p. 187 TMDL Report.

right" in every last calculation. (This is also based on their assumption that a theoretical number from their computer is as valid as a measured "violation" in the water.)

In the above analysis, I used only data from the modelers' own graphs in Appendix G2. The inability of the model to "get it right" in every calculation is also evident if data from other sources are used. Figure 2-4 provides an example. That Figure (same as Figure 2-1 above) shows the bottom water at station BI-6 with an overlay of data points from the BISS spreadsheet for that site. The data presented by the modelers (circles) are identical to those from the spreadsheet (triangles) in many instances. The modelers' data include values not found by me in the spreadsheet (for example, two points near July 1 whereas the spreadsheet shows only one) and values found in the spreadsheet that are not shown on the modelers' graph (for ex-ample, the very high data point in mid- September).

Table 2-1 compares all of the bot-tom water DO data from East and West Bay stations for September, 1997 as reported in the Appendix G2 pages and in the BISS spread-sheet data. The correspondence is loose, at best. I haven't attempted to reconcile the two data sources and have taken both at face value throughout this Analysis.
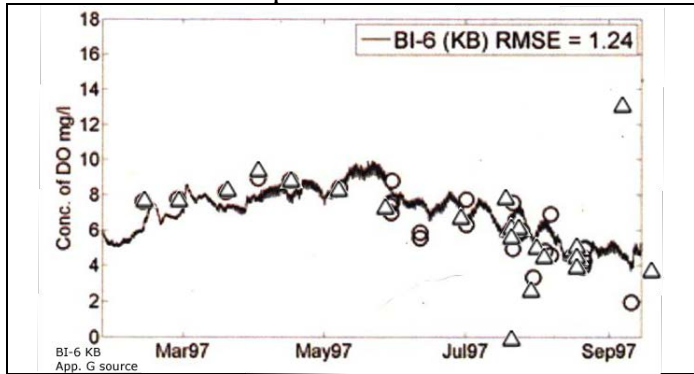


Figure 2-4. Figure from the TMDL Appendix with an overlay of data from the BISS research (triangles). The rightmost triangle is on September 24, a few days after the end of the computer simulation. Source: BISS spreadsheet.

| | Appendix G2 data | | BISS data | | | |
|---|---|---|---|---|---|---|
| Station | Date | Bottom DO (mg/L) | Date | Time | Depth to Bottom (m) | Bottom DO (mg/L) |
| BI-6 | ~Sept 10 | 2.0 | Sept 10 | 12:44 | 9.0 | 12.53 |
| | | | Sept 24 | 12:28 | 10.0 | 3.59 |
| BI-5 | ~Sept 10 | 2.5 | Sept 10 | 12:58 | 13.5 | 2.16 |
| | next day? | 3.5 | Sept 24 | 12:44 | 13.0 | 3.83 |
| | | | Sept 25 | 8:58 | 8.5 | 4.09 |
| BI-4 | ~Sept 10 | 9.5 | Sept 10 | 13:13 | 13.5 | 4.29 |
| | | | Sept 24 | 12:57 | 14.0 | 3.91 |
| BI-2 | -- | ND | Sept 10 | 13:29 | 5.5 | 13.51 |
| | | | Sept 24 | 13:13 | 9.0 | 4.10 |
| BI-1 | ~Sept 10 | 13.5 | Sept 10 | 8:44 | 4.0 | 3.47 |
| | | | Sept 10 | 13:43 | 6.0 | 13.53 |
| | | | Sept 24 | 13:26 | 7.5 | 2.84 |

Table 2-1. Data for all observed bottom water dissolved oxygen levels during Septem-ber, 1997, at all West and East Bay stations. Appendix DO's and date(s) were estimat-ed from the graphs. BISS data were taken from the BISS spreadsheet. Appendix data don't show depth to bottom or sample times. Depths to bottom vary in the BISS data due to tide changes. BISS observations extend past the September 15 end date of the computer simulation interval.

The September 10 BISS value (12.53 mg/L) is startlingly high for bottom water in late summer. Nevertheless it is real. It is not cited in the BISS spreadsheet's "errors" section and similar high bottom- (and midwater- and surface-) values are seen on the same date at East Bay sites BI-2 and BI-1. If the computer were always accurate, it would have "noticed" this high value whether it was portrayed on a graph or not. (The line traced by the computer would have "shot up" to 12.53 mg/L on that date, then back down again by the next day, alerting the modelers to something special happening there.) As we see, the computer "missed by a mile."

The lowest graph value calculated by the computer should have branded the BI-6 site as "in theoretical violation of water quality standards" on September 10 – but it did not.[2] Ironically, an accurate calculation would also have shown (as did the actual measurement) that the bottom water also experienced the highest DO water quality of the entire year for that site -- on that same day.

Figure 2-5 shows an instance in Eld Inlet where the bottom water was dangerously low in oxygen at 2 AM (less than 2.0 mg/L) then astoundingly high in DO (more than 10 mg/L) one hour later. In this case, there is a simple tidal explanation. For a computer making calculations every six minutes and programmed to keep its eye on layers of water (so to speak), this would be easy to detect and report. Station BI-1 in Budd Inlet, comparable in shallow depth and position in the estuary to the Eld Inlet station, shows something similar – a jump from 3.47 to 13.53 mg/L over a period of 5 hours (Table 2-1). The underlying cause at BI-1 in Budd Inlet may be ecological rather than tidal. If so, a more sophisticated simulation than the Budd Inlet Model would be needed to track it.
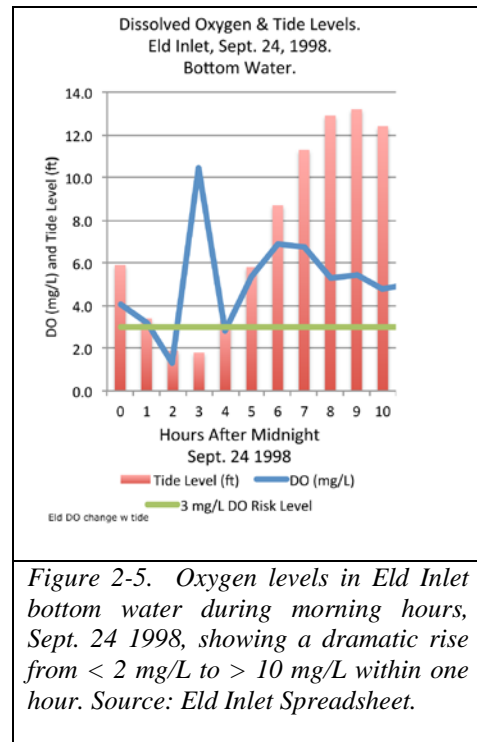


*Figure 2-5. Oxygen levels in Eld Inlet bottom water during morning hours, Sept. 24 1998, showing a dramatic rise from < 2 mg/L to > 10 mg/L within one hour. Source: Eld Inlet Spreadsheet.*

## 2d. Statistics could play a role.

As mentioned, graphic comparisons of the computer's calculations with real observed DO data are presented in Appendix G2 for the surface, bottom and a middle depth at the BISS stations considered by the modelers. Each of these graphs has a box in the upper right corner with the label "RMSE"and a number in it (example; Figure 2-4 above). The number is the "Root Mean Square Error," which is essentially the average distance by which the computer's calculations "miss the mark." In an analogy with bullets fired at a conventional circular target, the RMSE is an approximation of the average distance of all bullet holes from the exact center of the bullseye.

---

[2]   The modelers did not show station BI-6 in violation of water quality standards at any depth on any date (see TMDL Figure 90 and Figure 1-1, Section 1 of this Analysis). A staff member expressed surprise when I pointed that out.

In the DO situation as used by the modelers, the size (radius) of the bullseye is always 0.2 mg/L or more, even though the number at the exact center is not always the same. The average "miss" by the computer (that is, the RMSE) is always larger than 0.52 mg/L at every depth and station, ranging from 0.52 mg/L to 4.72 mg/L (BB-2 and BE-2 surfaces, respectively). This does not mean that *all* "shots" miss the "bullseye" – but where the RMSE is large, half or more of them miss the mark by an amount that obscures the true value of the "target" whose size we would like to estimate. This was the subject of my presentation to the modelers and others on November 14, 2014 (Power Point OK2, 2014).

Statisticians have perfected many reliable tools for overcoming the "misses" in calculations derived from sample measurements and for having confidence that data show (or don't show) what you want to know. One such practice uses "confidence limits" calculated from the data. Here I present an example without burdening unwilling readers with the details.

Figure 2-6 shows a sample of observed DO data obtained in West Bay during a total drainage of Capitol Lake, and two pairs of confidence limits (CL's) that were calculated from that sample. The two sets of CL's are shown at the right edge of the Figure. Each consists of a "data point" linked to 2 horizontal bars (the CL's) above and below it.

Speaking non-statistically, CL's "trap" unknown numbers. Even if you can't know or calculate their exact values, you can still be "pretty sure" that the numbers you're interested in are somewhere between the two CL's. You find CL's by first finding the average of a sample of several measurements (or calculations), then (speaking non-statistically) "take it from there."[3]

In Figure 2-6, what we'd like to know is whether the average real-life DO at that site was low enough to qualify as a Water Quality Standards Violation. The violation threshold is 4.80 mg/L (red line), the average of 17 measurements is 4.06 mg/L – black data point and line – well below the threshold.
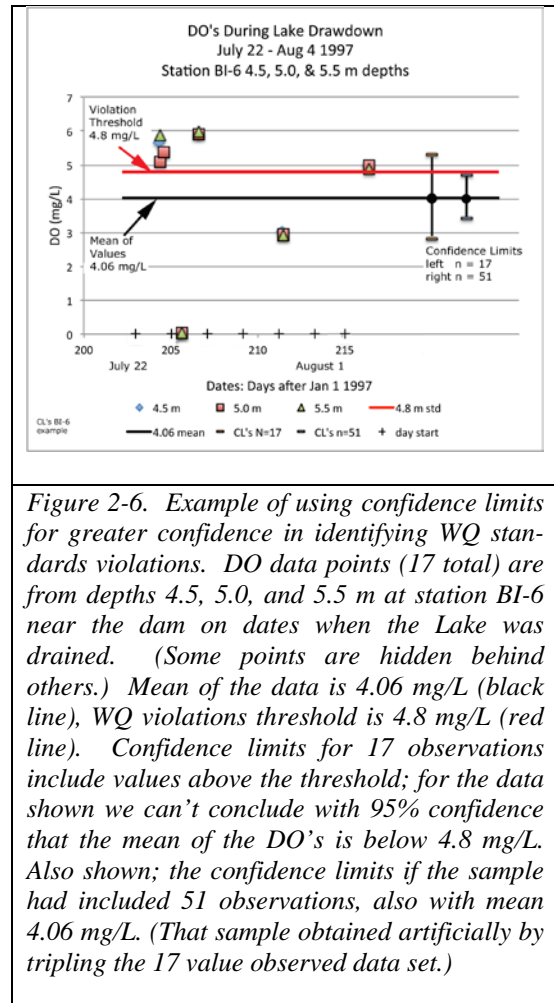


*Figure 2-6. Example of using confidence limits for greater confidence in identifying WQ standards violations. DO data points (17 total) are from depths 4.5, 5.0, and 5.5 m at station BI-6 near the dam on dates when the Lake was drained. (Some points are hidden behind others.) Mean of the data is 4.06 mg/L (black line), WQ violations threshold is 4.8 mg/L (red line). Confidence limits for 17 observations include values above the threshold; for the data shown we can't conclude with 95% confidence that the mean of the DO's is below 4.8 mg/L. Also shown; the confidence limits if the sample had included 51 observations, also with mean 4.06 mg/L. (That sample obtained artificially by tripling the 17 value observed data set.)*

---

[3] The CL's were calculated from 4.06 +/- std dev (of array of 17 values) x $T_{(.95,\ df=15,)}$/sqr root (17). See Keller, 2001.

Because of the vagaries of sampling (or the hit-or-miss nature of the computer's calculations) 4.06 mg/L may or may not be the "real" (= real-life, "population") average DO of the water. But even if it isn't, we can be "pretty sure" ("95% confident," speaking statistically) that the real life average, whatever it is, lies between the CL's.
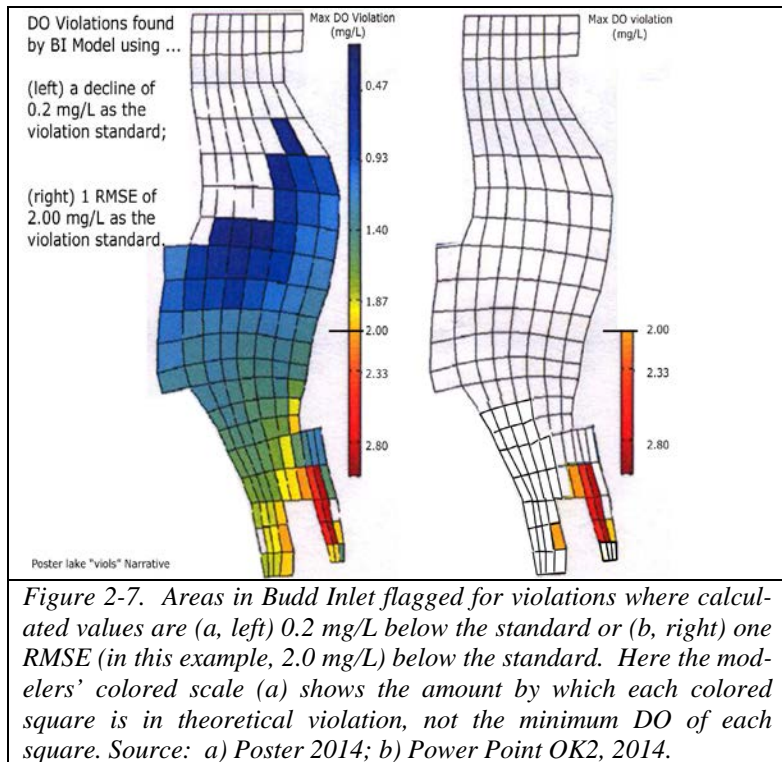
For the sample of 17 measurements, the upper CL is higher than the 4.8 mg/L cutoff threshold. The "real mean," whatever it is, occurs somewhere between the two CL's. If it is in the upper end of this range, it could be above the threshold. We can't be "pretty sure" (that is, "95% confident") with only 17 measurements that a violation really exists.

The way to shrink the CL range is to take more measurements (or include more computer calculations). Figure 2-6 shows the CL range if our sample had 51 measurements (and the same average, 4.06 mg/L). In that case the upper CL would be less than 4.8 mg/L and we could be "95% confident" that the unknown average DO, whether it is really 4.06 or something different, violates the Water Quality Standard.

Caution: the data obtained by the computer do not consist of independent measurements. Each calculated DO value is partially pre-determined by ("dependent upon") its value just a few minutes earlier. That may also be true of measured data, in this instance. That is a complication that introductory statistics are not prepared to deal with. Simple CL's like those for the sample of 17 real-world measurements probably aren't appropriate for data of this kind. *Only a professional statistician can advise on ways of having confidence in calculated answers in such situations.*

If something less complicated than CL's be needed, one possibility might be a simple "rule of thumb" like the one that I suggested to the modelers on November 3, 2014. That is, subtract the RMSE from the WQ standard and compare every calculated DO value with the number thus obtained. If the calculated value is lower, the likelihood is high – maybe 84%?[4] – that a real violation occurred at that time, depth



and place indicated. Figure 2-7 from my presentation illustrates the difference that this rule of thumb would make in understanding Budd

*Figure 2-7. Areas in Budd Inlet flagged for violations where calculated values are (a, left) 0.2 mg/L below the standard or (b, right) one RMSE (in this example, 2.0 mg/L) below the standard. Here the modelers' colored scale (a) shows the amount by which each colored square is in theoretical violation, not the minimum DO of each square. Source: a) Poster 2014; b) Power Point OK2, 2014.*

---

[4] The amount of "confidence" in this "simple" case is beyond the author's statistical comfort level and would need to be calculated by a professional statistician.

Inlet. That is, fewer theoretical violations would be found, but we could have confidence that they really do occur, when and where they are found in this way.

The modelers heard these suggestions in my presentation to them on November 3, 2014. They appear to have taken some notice of it and mention confidence limits in the new SM Report (pp. 27-28). Although their explanation is not easy to follow, they appear to compare the computer's calculations in one scenario ("natural") with its calculations in another ("current conditions") and concentrate on the variability in differences between the corresponding calculations. They find that if there is a difference between a calculated estuary number and the corresponding calculated lake number, that difference is likely to be "real." Nothing appears to be said about comparing the computer's calculations with real data. A better explanation of what they are referring to is needed before knowledgeable readers can evaluate their claims here.

The modelers are not inclined to use averages of DO calculations in their search for theoretical violations. Elsewhere (SPSDOS 2013 Report p. 35) they have said that averages cannot be used to "mask" the fact that a grid cell's DO dropped even briefly below the WQ standard for that area. Their preference is to take each individual calculation at face value and assume that it is accurate enough for real-life policy decisions.

Expressing doubt about the alleged dead-on accuracy of every calculation, personnel of the HDR consulting firm asked the modelers precisely that question in the firm's comments on the draft SPSDOS Report (2013). In their words:

*"Page 19: The DO decreases calculated by the model range from 0.2 to 0.4 mg/L in limited areas due to point sources. These are very modest changes in the DO levels in these locations. Due to these small calculated DO decreases, the following question arises: Is the model sufficiently accurate to predict these DO decreases? And more importantly, is there sufficient confidence in the DO decreases calculated by the model to mandate expensive nitrogen removal upgrades at point source treatment facilities to reduce nitrogen loadings?"*

The Department of Ecology did not respond to the HDR query (Clark, 2016).

## 2e. Hiding the Search for Violations.

One place where the hit-or-miss accuracy of the model makes a huge difference arises from the question: "Did the water quality of Budd Inlet meet modern standards, even in its 'natural' state before human activity began to modify it?" The consequences of that question are explored in this subsection, centering upon the pair of illustrations in the SM Report reproduced here as Figure 2-8.

Figure 2-8a shows the part of Budd Inlet for which the modern water quality standard is 6.0 mg/L (green) and the part where the standard is 5.0 mg/L (orange). If the natural water of times past had DO's that were always higher than these values, then the computer could look for violations of the standards in modern waters by simply comparing the modern waters' DO's with 6.0 or 5.0 mg/L in each grid cell. If the "natural" waters would have violated these standards, however, then the search for violations in modern water becomes very complex. The challenge is to learn or estimate what the DO's in the natural waters of times past really were, to determine the method by which violations are sought in modern waters.

Absent observed data, one way of doing this would be to run the Budd Inlet model with pre-modern conditions of the past – differences in weather, river runoff, and inputs from the Pacific
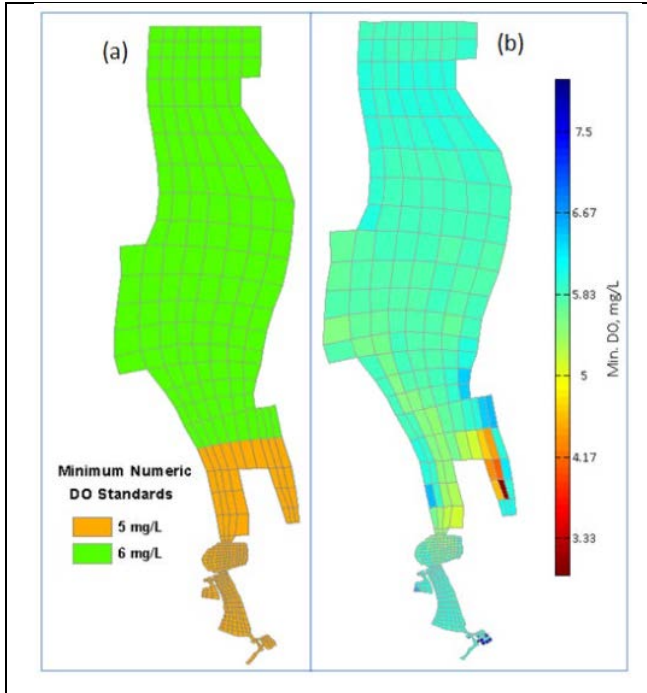


*Figure 2-8(a). Modern water quality standards that apply to Budd Inlet. (b) Minimum dissolved oxygen levels in Budd Inlet as calculated by the modelers for 'natural' waters before they were altered by human activity. ("Capitol Lake" in Fig. 2-8b is an estuarine extension of Budd Inlet, not a dammed impoundment.) Source: Both images make up Figure 7 (p. 32) in the SM Report.*

Ocean to the extent that those can be known or estimated. No dam or impounded Lake would be present. What about tides? Use 1997 tides or those of some year of the past? The modelers are not clear about how they do this.[5]

In any event, the modelers run the model with settings for presumed pre-modern conditions, calculate the DO at every depth underneath every grid location in Budd Inlet every six minutes for 9 months, while comparing each calculation with 6.0 or 5.0 mg/L depending upon the location. Figure 2-8b presents their findings. Rather than show readers the grid cells in which the "natural" waters violate a modern water quality standard, they paint each grid square with a color that represents the lowest DO that they found there during the 9 simulated months.

Are those DO levels above or below the modern standards? It is possible to see that East Bay has several cells clearly in violation, as heads of estuaries naturally do in late

---

[5] The modelers refer (SM Report p. 26) to TMDL Appendix I for 'natural' conditions of the past. Confusingly, Appendix I (p. I-7) says that "current" values of the Deschutes River flow – and temperatures and other properties – were used in their simulations of 'natural' pre-modern waters. This is in stark contrast to their reply to my questions about this (see Section 7, this Analysis). As another example of their typical indifference to consistency and detail, the grid above has three top tiers at Boston Harbor in Figure 2-8a, two in Figure 2-8b.

summers, but aside from there, Figure 2-8b leaves readers no clue. In the following I show how Figure 2-8b can be compared with Figure 2-8a and relate the result to the implications of claiming that every number calculated by the computer is accurate.

## 2f. Methods. Finding the Water Quality Standards Violations in the Pre-Modern ('Natural') Estuary.
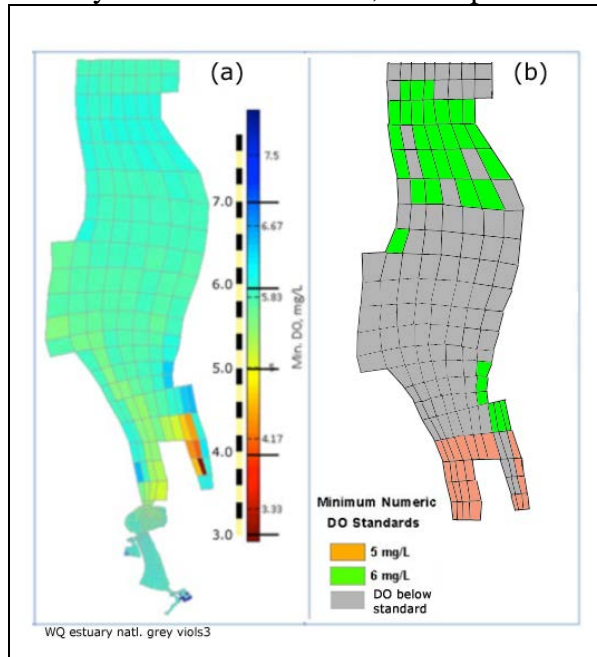
Figure 2-8b can be analyzed using Photoshop Elements 12 software. Using that software, I first added a readable DO scale to the Figure for determining which colors showed DO's less than 5.8 mg/L for the central part of Budd Inlet and which colors showed DO's less than 4.8 mg/L for the southernmost part. (These DO levels are 0.2 mg/L below the respective modern standards. Any DO reading below these levels qualifies as a WQ standard violation.)

I then selected the interiors of various grid squares in Figure 2-8b using Photoshop's "polygonal lasso" tool and clicked "Similar" in the Selection menu. This function identifies every part of Figure 2-8b that has the same color as the selected square. *It also identifies the part of the modelers' scale that has the same color.* By comparing the colors selected on the modelers' color scale with my more readable scale, it was possible

to see whether the 'natural waters' in each selected area of Budd Inlet were calculated to be in violation of modern WQ standards. I assigned every grid square that was in theoretical violation – whether the calculated "violation" was large or small – a grey color and left the squares that were not in violation in the green or orange colors that show the standards. (A detailed description of this technique is in Section 5.)



### 2g. Results. Most of the 'Natural' Estuary Violates Modern Water Quality Standards.

Figure 2-9b shows the result of this exercise. The grey areas in Budd Inlet show locations where the 'natural' waters of some time in the past experienced DO levels lower than the modern standards at least once during the interval January 25 – September 15. Only 57 of the grid locations out of 160 total had 'natural' waters that *always* contained more than 6 (or 5) mg/L of dissolved oxygen. Those grid squares in modern waters can be judged

*Figure 2-9. (a) "Minimum DO" data provided by the modelers for the pre-modern estuary, from which Figure 2-9b was determined (same as Figure 2-8b above with a readable DO scale added by me. (b) Grey areas show where pre-modern ('natural') waters had DO levels 0.2 mg/L or lower below modern water quality standards, as calculated by the Budd Inlet model.*

by the 6.0 or 5.0 mg/L standards. In the other 103 squares, theoretical modern water quality violations must be calculated by the more complex method.

## 2h.  Discussion.  It is Impossible to Check the Calculations when the 'Natural' Estuary is used as the Water Quality Standard.

Once the areas where the "natural" waters violate modern standards are identified (grey grid cells, Figure 2-9b), the modelers use a complex method for those areas to "find" theoretical violations in modern waters.  For each grid square, at times when the 'natural' waters have DO's higher than the standards shown in Figure 2-8a, those numerical standards (6.0 or 5.0 mg/L) are used for the modern waters.  But at times when the DO's of 'natural' waters are lower than the modern standard, then the DO of the 'natural' water itself is used as the standard.

The grey areas are veritable Happy Hunting Grounds for finding theoretical violations in modern waters.  There, compared with calculated DO's of the past whose real values or times of occurrence we can never know, the modelers can assure us that modern-day "violations" of as little as 0.2 mg/L have been identified.  The foundation of this assurance is the assumption that the model gets the exact right answers 100% of the time – first when it calculates the DO's of the "natural waters," then again when it compares the calculated modern DO's with those "natural" DO's.

The modelers have implied elsewhere (SPSDOS 2013, p. 87) that all that is needed to declare a location (= grid square) in violation of modern water quality standards is a single computer calculation of a DO level that is slightly lower than the DO of the mythical 'natural water' at that time and place.  An example described by them (obtained from a model similar to the BI model but expanded to Central and South Puget Sound) is a location with a modern WQ standard of 5.0 mg/L where the calculated DO of the 'natural' water dropped to 4.95 mg/L *for all or part of just one day out of the 302 days simulated by that model.*  The whole grid square was flagged as "in violation."[6]  That is, the 'natural' water's DO fell below 5.0 mg/L just once by an amount so small – 0.05 mg/L -- that it is well-nigh undetectable in real life – an illustration of confidence with which the modelers regard their calculations – namely that they are always dead-on accurate to the second decimal place.

A drawback of the grey zones of Budd Inlet (Figure 2-9b) is that it is impossible for anyone to check up on the numbers used to assign violations to the modern waters in those zones.  The violations originate from the supposed waters of the past, whose exact DO levels we can never know.  If the observed BISS data show "no violations" at the times and places when measurements were made, that reality can always be dismissed by saying "yes, but the computer detected theoretical violations at times other than those hours during which the BISS observers were actually observing."

Despite the impossibility of checking up on the model's calculations over most of Budd Inlet, it is still possible to do so in those few areas in the remaining green and orange zones of Figure 2-9b.  There we know that the DO levels with which the modern water

---

[6] In this case the "violation" – 4.95 mg/L – is only 0.05 mg/L below the standard.  I don't get it.  The violation threshold is supposed to be 0.20 mg/L lower in all other applications. Why this exception?

must be compared to find WQ violations are always 6.0 and 5.0 mg/L, not some unknown/unknowable theoretical DO level of waters past. There are at least two in-stances in which WQ violations occurring in real life were not found and flagged by the computer. In the BISS data, these are West Bay sites BI-5 (observed DO levels 4.74 and 2.16 mg/L on June 12 and September 10 1997) and BI-6 (DO's of 3.83, 4.48 and 4.37 mg/L on August 20-21, 1997). The modelers, on the other hand, found "no violations" at these two sites (see Figures 2-9a and –b).

The modelers' portrayal of where Budd Inlet's modern waters are lower in DO than the 'natural' waters of times past is not credible at face value. It is based on the assumption that every last one of the calculations of DO's in "natural" waters is accurate, and then on the assumption that every one of the corresponding calculations of DO's in modern waters is also accurate, and that therefore the differences between every pair of numbers from the two calculations are also accurate.[7] Even though we can only examine the mar-gins (BI-6 and BI-5) of the 'natural violations' zone, we can still recognize that the com-puter was in error some of the time. That is just one more illustration of the fact that its detailed projections are untrustworthy throughout all of the rest of time and space.

---

[7] Actually, where the DO's of the 'natural' waters are used as the standard, the calculated differences between the 'natural' and 'modern' calculated values are less likely to be accurate than is either individual calculation by itself. If, say the probabilities that the 'natural' and 'modern' values are accurate are (1/2) and (1/3) respectively, the probability that their difference is accurate is only (1/2) x (1/3) =1/6.