

Index	Line (initial)	Comment	Reviewer	Response (revision or comment)
1	Section 2.6	For the error propagation described in section 2.6, what is the reasoning for using CO2 flux? Using CO2 flux introduces several other potential biases and errors to the assessment:	1	The initial reasoning for including the CO2 flux was to show implications of error propagation for various imputations during a common use case for DIC time series. However, both reviewers have raised similar concern about the introduction of multiple sources of error when determining CO2 flux. The combined uncertainty for the CO2 flux was initially determined by a Monte Carlo method (n=1000, which was not significantly different than n=10000) and then only the values of imputed DIC and their associated uncertainty were varied as inputs into the calculation. In this way we attributed the percent difference between imputed time series and observed time series to be related only to the gap-filling method because no other input was varied. Similarly the uncertainty of the CO2 flux was determined via this MCM for each method. That said, we understand the concern about multiple sources of error and recognize that this application detracts from the results of just gap-filling the DIC time series. We will remove the CO2 flux aspect of paper and add a focus on long term trend assessment in its place. This will be more consistent with the intentions of the work, enhance the focus of the paper and address multiple comments from both reviewers.
2		· uncertainty in air pCO2,	1	removing CO2 flux aspect
3		· major bias and errors of NCEP winds (see: <a href="https://doi.org/10.5194/bg-15-1701-2018">https://doi.org/10.5194/bg-15-1701-2018</a> , <a href="https://doi.org/10.1029/2018GB006047">https://doi.org/10.1029/2018GB006047</a> , <a href="https://doi.org/10.1002/2017GL073814">https://doi.org/10.1002/2017GL073814</a> )	1	removing CO2 flux aspect
4		· uncertainty in the gas transfer velocity coefficient (resulting in total uncertainty in CO2 flux of ~20%), and	1	removing CO2 flux aspect
5		· uncertainty (~5%) introduced in the calculation of sw pCO2 from DIC and TA.	1	removing CO2 flux aspect
6		How will those biases and errors complicate your assessment of gap filling error propagation? The relative uncertainty for CO2 flux at BATS is reported in line 361 as 3.5%. What does this uncertainty take into account? Not items 1 – 4 above, as this value would be much higher. These issues should be addressed in the error propagation, or another parameter should be used for this assessment.	1	removing CO2 flux aspect
7	NA	Data used in this study need to be cited properly, which is incredibly important to the programs supporting these time series measurements. Those data should be cited in the methods and/or funders noted in the acknowledgements, depending on what each time series program recommends, not recorded as web addresses in the notes of Table 2. For the moorings, if you are accessing original data files via NCEI, those citations can be found at <a href="https://doi.org/10.3334/cdiac/otg.tsm_papa_145w_50n">https://doi.org/10.3334/cdiac/otg.tsm_papa_145w_50n</a> for Papa and <a href="https://doi.org/10.3334/cdiac/otg.tsm_keo_145e_32n">https://doi.org/10.3334/cdiac/otg.tsm_keo_145e_32n</a> for KEO. If you are accessing the mooring data from the synthesis product, the citation can be found at <a href="https://doi.org/10.7289/V5DB8043">https://doi.org/10.7289/V5DB8043</a> . I am not as familiar with the citation requirements of all the ship-based time series, but with a quick search I found this data citation request for HOTS, for example: <a href="https://hahana.soest.hawaii.edu/hot/dataaccess.html">https://hahana.soest.hawaii.edu/hot/dataaccess.html</a>	1	This was a gross oversight on our part. While the sources for data sets were listed in table 1, they were not properly cited as noted. We will cite these as required.

8	NA	Finally, it may be out of the scope to include additional analyses in this paper, but it would be worthwhile discussing future work that can build off these results. For example, what satellite-based products are best suited for the MLR approach? Are there any that can span open ocean and coastal environments, so gap filling methods can be applied consistently across all global ocean and coastal time series? Also, it would be useful to study whether there are discrepancies in calculated trends when using these different gap filling methods (at least the most successful methods) or no gap filling methods at all. Both of these analyses seem like they could have been included in this paper, but I could also understand if those are the next assessments planned using the most promising empirical gap filling methods resulting from this work.	1	These are excellent points some of which we can address in the revision. Firstly, we will include an assessment of impacts on trends in place of the CO2 flux. Secondly, we have already separately performed a cross shelf assessment of the MLR performance using data from the Munida transect and we can included this appliaction. These aspect taken together will also help address Reviewer 2's comment about focusing the scope of the paper more on presenting this MLR method and comparing it to other gap-filling methods, rather than an extensive comparative assessment of techniques since we have only selected a few methods from a very large number of possibilities.
9	31	Use the most recent version of the Global Carbon Budget: <a href="https://doi.org/10.5194/essd-12-3269-2020">https://doi.org/10.5194/essd-12-3269-2020</a>	1	updated reference
10	89	Define DT	1	This was formatting issue only; this is really delta T, but the Greek symbol was lost at somepoint in the upload.
11	89	State the sites that did not measure DIC directly as in line 87 for discrete sampling sites	1	done
12	90	What measured parameters are being used to calculate DIC from the moored data? Measured pCO2 and pH? The measured pCO2 and pH pair has several issues, most importantly in this application is the issue brought up below for line 118, in that data return from pH sensors tend to be poor and data gaps will usually fall at the same time each year. Data return from the pCO2 systems are much better, and you will avoid much of the repeated seasonal gaps if you used established salinity-alkalinity relationships (in the Fassbender references) for those open ocean locations paired with measured pCO2 as discussed in <a href="https://doi.org/10.5194/bg-13-5065-2016">https://doi.org/10.5194/bg-13-5065-2016</a> . This will increase N Years in Table 3 for Papa and KEO	1	We calculated DIC for both KEO and Papa from measured pCO2 and a calculated total alkalinity using the published salinity algorithms from Fassbender et al 2016 and 2017, as you suggest. We will make sure this is more clearly communicated in the methods section for these sites as well as all other sites that do not mearusre DIC directly.
13	96-99	It would be useful to present more information (figure or some statistics like mean diff and standard deviation) about how MODIS and VIIRS compare at this particular site so it is more clear why VIIRS was chosen.	1	We will address this in revision
14	118	"Missing at random" is not a good assumption for many of the moored time series, especially the open ocean sites which tend to be serviced around the same time every year. Sensor failures are more likely late in the deployment, which can be around the same time every year just before servicing. That should be acknowledged here.	1	This has been addressed
15	202	BATS is a different latitude than Mauna Loa, and therefore, has different annual mean and seasonality of air xCO2. xCO2 air from same latitude of BATS should be used from one of these products:	1	removing CO2 flux aspect
16		<a href="https://www.esrl.noaa.gov/gmd/ccgg/obspack/our_products.php">https://www.esrl.noaa.gov/gmd/ccgg/obspack/our_products.php</a>	1	removing CO2 flux aspect
17		<a href="https://www.esrl.noaa.gov/gmd/ccgg/carbontracker/">https://www.esrl.noaa.gov/gmd/ccgg/carbontracker/</a>	1	removing CO2 flux aspect
18	344	What about: <a href="https://doi.org/10.1002/lom3.10232">https://doi.org/10.1002/lom3.10232</a> and <a href="https://doi.org/10.3389/fmars.2017.00128">https://doi.org/10.3389/fmars.2017.00128</a> ?	1	Will review these references during revision
19	358	You should note that the studies cited here do not use ocean DIC time series. Include information on what types of time series these are (soil flux and respiration, etc).	1	This has been addressed
20	406-408	Since trends were not considered in this paper, this statement may be a bit premature?	1	This will be address / revised according to shift from CO2 flux to trend assessment
21	655	What is the note with the "*" referring to?	1	This has been addressed
22	Fig 10	Why aren't the models listed above the top panel? And spline should maybe be presented on the far right or left since it has a diff y axis for the 6 month gap?	1	This was a formatting typo in R, but this figure will be removed per the removal of the CO2 flux assessment

23	Fig 12	Consistent with earlier comments about error propagation for CO2 flux, these results showing higher uncertainty at higher outgassing and uptake values are consistent with increased uncertainty at higher wind speeds. This makes it difficult to understand what is a gap filling uncertainty vs uncertainty in other parameters that impact flux.	1	As stated above, this assessment was meant to illustrate the change related to gap-filling method only since all other input data were held constant during the MCM analyses. However, this will also be removed during revision
24	66	On Line 66 is stated "This study aims to identify the optimal gap-filling methods for carbonate time series by establishing which techniques perform with sufficiently low error and bias to assess seasonal and interannual variability of carbonate biogeochemistry and the biological and physical processes that determine it." The manuscript takes the approach that all gap-filling techniques have been explored and that MLR is recommended as the best performing. While the latter is certainly true of the methods compared, I feel it is not currently possible to say the former while one / a number of machine learning (and other) approaches are absent - these have recently been successfully applied in oceanographic research, and so the manuscript is not fulfilling its own aims by omitting them. Clearly it is not feasible to compare all available methodologies, so I would recommend that you either tone down the aims of the paper (by saying that you present a MLR method for DIC time-series data gap imputation and compare it to other common, computationally inexpensive methods) or a selection of additional methods are included e.g. median as well as mean, machine learning (i.e. neural network, regression trees, random forests that you already mention), curve fitting, exponential moving average, k-nearest neighbours etc.	2	This point is well-taken. Given that we have not here (and could not really) assessed all methods, we will shift the stated focus away from optimization of gap-filling and toward presenting the MLR and comparing it against other common approaches as suggested.
25	NA	When comparing methods a lot of focus is on the magnitude of the RMSE. I feel the reader would benefit from some consideration of the structure of the error e.g. are certain times of the year subject to greater uncertainties, do the models reproduce the timing of the seasonal cycle, and the magnitude of the peaks and troughs or are these far worse than those that vary around annual mean values? Equally, is the error of the preferred MLR technique actually normally distributed, as a lot of its power rests on this assumption. The manuscript would certainly benefit from greater examination of the seasonal cycle, and anomalies from this in the imputation methods.	2	This point is also well taken. With removing the CO2 flux aspect of the paper we can provide more room for showing the distribution of error. As for the structure of the seasonal cycle, we discuss this but had not quantified it. In revision we will provide quantification of the timing and magnitude of the seasonal cycle and some metric(s) for method performance to make this discussion less qualitative.
26	NA	The use of the air-sea CO2 flux for assessing imputation performance is an interesting choice, as it introduces a whole suite of additional uncertainties (wind-speed, piston velocity, K1/K2 equilibrium constants, how missing alkalinity data is filled etc) that are not considered in your error analysis. These uncertainties would also need to be assessed, or another metric/s chosen for comparison. If the air-sea CO2 flux is still the preferred metric, is it not better to calculate pCO2 from DIC/alkalinity first, before imputing missing pCO2 values?	2	See our response to reviewer 1 comments above regarding our initial methods and reasoning; and note that we will be removing this aspect from the paper.
27	NA	I appreciate that this may be being considered in a follow up study, but an assessment of the desired sampling frequency necessary to generate a good representation of the seasonal cycle (1, 1.5, 2, 3 month frequency, only summer and winter etc) would be very interesting/useful.	2	We will add this assessment
28	36	value is singular, so has not have	2	This has been addressed
29	38	40% - This is possibly fossil fuel CO2 emissions? All anthropogenic CO2 (including land-use change and cement) means the ocean component is probably closer to 25% (Global Carbon Project, Friedlingstein et al., 2020)	2	This has been addressed

30	66	"This study aims to identify the optimal gap-filling methods for carbonate time series by establishing which techniques perform with sufficiently low error and bias to assess seasonal and interannual variability of carbonate biogeochemistry and the biological and physical processes that determine it." - see comment above	2	Response as above
31	72	should be principle rather than principal	2	This has been addressed
32	75	(and Table 1) - add citation/references for time-series, possibly through additional column in Table	2	As per response to Adrinne's comment above, this was an oversight and all dataset citations will be properly added.
33	86	Is there an impact on your analyses of averaging data to monthly means?	2	Uncertainty in monthly values was estimated for both single observations and averaged higher frequency measurements from moorings so they could be properly compared. We will make sure this is clearly communicated in the methods during revision
34	89	would be better to use greek delta notation rather than DT	2	fixed per above as well
35	90	What is the uncertainty introduced by the use of estimated DIC values? DIC is only measured at BATS. What do you get if you apply the same techniques to data with DIC, TA and pCO2 e.g. at sea surface?	2	Individual DIC uncertainty budgets were assessed by adding the sources (measurement, natural variability (e.g. monthly averaging), and /or propagation from calculating DIC from other carbonate measurements) in quadrature to determine the combined standard uncertainty for each DIC value in the time series. For DIC calculated from the other variables such as pCO2 and TA, the error function in the R package seacarb was used.
36	122	"The primary goal was imputing timeseries at monthly resolution to investigate variability and trends over seasonal, interannual and decadal timescales" - neither trends nor decadal are covered as far as I can see?	2	See our responses above that indicate we will be removing the CO2 flux aspect and adding an assessment of trends and seasonal structure
37	141	is this not an exponential moving average then, rather than a weighted moving average?	2	I suppose it could be stated both ways. It is a weighted moving average, but the weighting is based on an exponential relation to neighbors
38	148	cite1 and cite2?	2	This was some sort of formatting typo with Endnote - will fix
39	150	does this method also input uncertainty into the fitted values used?	2	I don't believe this inputs uncertainty - rather values are found through convergence of multiple regressions. Uncertainty can be assessed by looking at the spread when the option to have multiple outputs for a given value is selected.
40	190	as above, why this? Is it not better to calculate pCO2 from bottles at the start, then do imputation on pCO2 data set?	2	No imputation of pCO2 data was done. All imputation is on DIC values only. All pCO2 was calculated from the imputed DIC and either measured or estimated alkalinity
41	193	Wanninkhof 2014 recommends to not use Wanninkhof 1992.	2	removing CO2 flux aspect
42	201	why not use Bermuda atmospheric CO2 concentrations?	2	removing CO2 flux aspect
43	215	what were these uncertainties? It would be good to state them here. pCO2 from DIC and TA at their measurement uncertainty is ~6µatm. What is it when DIC is estimated?	2	We will make uncertainties more explicit during revision
44	223	To give a better feeling of interannual variability it would be useful to have the value for n for each month in Figure 2. For example so that a reader doesn't look at FOT and think there is very little variability in months 1-3, when instead n is only 1-2 for these months.	2	will add this info
45	227	& Fig 3. Is this a single MLR encompassing all data from all sites? Or the results of individual MLRs plotted and pooled? I'm don't think this is clear in the text	2	This is pooled results. MLRs must be fit using site specific observations and have unique coefficients. Will update language to clarify
46	229	"worked well"? A RMSE of 12 is beyond the 'weather' goal of measurement quality to assess spatial and short-term variability. I'm not sure stating this metric is useful as it obscures the capability of the method in (primarily) oceanic sites. Instead it might be better to simply focus on individual monitoring station results.	2	This is a good point and we will address during revision

47	234	It would be interesting to hear the thoughts behind why PAPA performs so well	2	The MLR appears to perform best at sites where there is a high correlation to temperature and a large seasonal cycle. The performance of the results follows the trends shown in Figure 5 where there is selective omission of predictor variables. We will investigate and elaborate on this further during revision
48	244	put the numbers in the boxes as well - the colour scale is not the most obvious/immediate to show similarity/disparity	2	We will add this info
49	245	add another line to the bottom of Figure 5 to show mean	2	We will add this info
50	246	Table 5 - change title to Mean model results	2	This has been addressed
51	250	Figure 6 - might be better showing as well / instead the residual (y) versus the measured (x)? - this may better highlight the better performing models, with the distribution of the residual ideally normal about 0.	2	We will explore this suggestion and other possibilities for expressing error distribution across both observed DIC ranges and sites
52	259	I struggle somewhat with this plot (Fig 7) too. The colour scale is not the most obvious/immediate to show similarity/disparity, and seems to be the opposite to Figure 5 where light colours indicate better performance - here they indicate worse performance.	2	We will look to increase the contrast and make these figure gradients consistent for clarity
53	261	I think that showing the performance of the models in recreating the seasonal cycle would be very useful. Whether they get the amplitude and timing correct is important for potential end users of these methods. Showing the anomaly from the observed seasonal cycle may also be useful.	2	This is a great point. As indicated in response above, this was qualified in the discussion but was lacking quantification and we will add that during revision
54	266	Fig 8A I like this plot, but i think it is making false equivalences by using different y scales for the 7 different methods for each monitoring station. It might be worth having this as a standalone figure to give more space to what is an enormous amount of information.	2	With the removal of the CO2 flux aspect and associate figure we will have space to break out this figure. The y scales were held consistent across sites so that methods could be compared. If the scales are held constant for all sites and all methods it will loose significant detail for visual interpretation.
55	275	Assessing error on seasonality and annual sums - not sure these numbers capture this. As mentioned above I'd be interested in seeing the performance of individual methods of capturing the seasonal cycle / amplitude and annual mean, and how they compare to the data, both using the full timeseries, and when there are artificial data gaps. It would certainly be useful to know how critical it is to sample seasonal maxima/minima (or not) in correctly formulating a seasonal cycle, and getting lowering the uncertainty with respect to annual budgets.	2	The aspect of sampling optimization is a good point that is missing here. As noted in responses above, we will quantify the performance of retaining seasonal structure and we will explore assessing sampling optimization.
56	280	and Figure 9A. While these plots are interesting it might be better represented by adding/replacing with anomaly timeseries. Also, I was wondering whether you could comment on how there appears to be a positive bias for the bimonthly and 3 month data gaps towards higher concentrations? Is the reason there are no red dots at the lowest concentrations (particularly in the 3 month timescale) simply the result of random data gaps, or something else? For the 6 month gaps I'd be interested in the performance of the models when only summer data is available, or perhaps completely missing winter data, as this would be a situation facing other time series sites.	2	We can explore representing anomalies here for clarity. As for the gap placement in the 3-month gap series, yes this is just due to randomization. We could explore artificially removing particular seasons and assessing impacts on annual cycles.
57	291	Fig 9b - would it be possible to have the legend across a single row, to aid in identifying models? Or indeed numbering the different box plots.	2	We will address clarifying the identification of methods in this boxplot
58	299	Figure 10 - this plot might be easier to interpret if it was anomalies from observations rather than actual values side-by-side?	2	This figure will be removed along with the CO2 flux aspect of the paper

59		The uncertainty bars also seem particularly low - has the uncertainty from the imputed data been propagated through the calculation? Even a DIC RMSE of 6 $\mu\text{mol/kg}$ would have an impact of 10-25 $\mu\text{atm}$ of $\text{pCO}_2$ depending on temperature. I imagine if there are missing DIC observations, there will also be missing alkalinity observations as well. It will likely be too much to include an estimate from these values as well, but I think you should comment on the fact that the error estimates relating to air-sea $\text{CO}_2$ fluxes presented here will be an underestimate, as there will also be additional uncertainties associated with imputing alkalinity.	2	The uncertainty budget was assessed using a MCM method as noted above in response to other comments, however we will be removing this aspect regardless in place of more focus on assessing trends and seasonal structure
60	328	change 'has a dominant effect the carbonate chemistry' to 'has a dominant effect on carbonate chemistry'	2	This has been addressed
61	333	need to reference these different datasets	2	This will be address as noted above
62	335	missing full stop	2	This has been addressed
63	353	- I don't think you've shown anything about temporal extrapolation.	2	This has been addressed
64	358	either remove the parentheses around the citations, or remove 'in the studies of'	2	This has been addressed - this was an Endnote formatting typo
65	369	This may be so but I don't think the figures you have presented make this obvious. A figure showing the mean seasonal cycle from the full data set compared to those imputed for different percentages of missing data would be necessary to show this.	2	Our quantification of seasonal structure during revision will address this
66	371	it's not clear visually, as you're missing a figure showing it. Figure 9 suggests it's only really obvious for the 6 month gap, while Figure 12 suggests that the mean approach has some of the highest uncertainties for the bi-monthly data gaps.	2	Our quantification of seasonal structure during revision will address this
67	381	I'd again suggest that looking at anomaly plots would be more straightforward to interpret than net flux comparisons	2	Point well taken, and we will explore this for clarity
68	405	change 'In general' to 'Of the methods we tested'	2	This has been addressed
69	408	May and possibly are really not strong enough - the artifice of the mean imputation method introduces bias, and actively removes any trend from the input data.	2	Good point, we will revise language here
70	415	- MLR certainly has the lowest error, but this doesn't necessarily tell the whole story. Showing the residuals of the predicted values will help - would you like to comment on the tendency of MLR methods to revert to the mean, where higher values are typically predicted lower, and lower values are predicted higher. This will have an impact on estimating maxima/minima. And I'd hesitate to recommend best practice until MLR is compared against a fuller suite of gap-filling methods, including machine learning	2	As also noted above in responses to a similar comment, we will revise the focus of the paper to dial back the language for establishing best practices and shift to scoping it as a presentation of this MLR compared to some selected common methods. The expanded seasonal structure assessment will help the discussion about max/min biasing
71	426	(and L433)- can be estimated, but to what uncertainty, and is this the same across all times of the year?	2	We will assess and present seasonal error distribution to further support this; however uncertainty must be assessed on an individual site/data set basis. We will make revisions to the language here to make sure that point is clear
72	432	I sound like a broken record but I think plots of seasonal cycles/anomalies of seasonal cycles/internannual anomalies are really what are needed to help determine this.	2	Noted and will address
73	433	Change "the most robust option for imputing gaps over a variety of data gap scenarios." to "the most robust option from those we compared for imputing gaps over a variety of data gap scenarios."		This has been addressed