

First, I would like to extend my thanks to both reviewers for their very thorough and thoughtful reviews. Being a newcomer to this particular area of research, the reviews drew my attention to important gaps in my familiarity with relevant prior work as well as giving me occasion to tighten up and clarify my discussions of both my methods and my findings. The revised manuscript that I will submit shortly will be greatly improved as a result.

In the following, the reviewer's comments are indicated by gray shading. My responses appear in the unshaded space below.

- Grant Petty

## 1 General comments

Based on data from a large-eddy simulation for a stratocumulus topped marine boundary layer, the author performed an ensemble of flight measurements and analyzed the convergence of the sampled fluxes towards their truth, as well as investigated the dependence of the random error on the track length. The author compares the observed random error against the theoretically derived expression by Lenschow and Stankov (1986) and found good agreement for track length of 10-30 km. Further, integral length scales of the turbulent quantities were calculated from the LES data for different flight track angles. The author shows that integral length scales depend on the flight track angle and compares these with the proposed, and still commonly-used, approaches by Mann and Lenschow (1994).

The topic and the content of the paper fits well into the journal and is of high interest to the research community. Due to lack of any alternative ways to estimate the random error for flight measurements, almost everybody uses the expressions proposed by Lenschow and Stankov (1986), though it is already known that especially for shorter tracks the estimated random error often does not reflect the true uncertainty. Here, especially for shorter flight tracks, an improved random error estimation is highly desired in order to avoid misinterpretations of observations.

The paper itself is well written and results are sufficiently presented, though at some points in the text more information needs to be given. At several points in the text it is not clear how variables are calculated and from which data, i.e. from the full 3D LES data or from the sampled space-series along the flight tracks. Also, at some points the discussion is rather short and the findings are not well put into the context of previous research. For example, also the study by Schröter et al. (2000) had already analyzed how the sampled flux converges towards its truth with increasing track length. Even though the atmospheric setups are not comparable, the findings presented here should be put into the context of previous work.

Thank you. I will attempt to add that context.

After extensive review I can recommend the manuscript for publication only after major revisions have been done and extended analysis is presented. My major concern is outlined in the following.

## 2 Major comments

(A) I miss one important aspect in the study. The author shows that the predicted random errors according to the Lenschow and Stankov formula matches the observed standard deviation remarkably well. However, in this study the random error using the LS86 formula is calculated based on integral length scales, correlation coefficients, and fluxes, that were inferred from the three-dimensional LES data. Though it is nice to show that the LS86 formula works well in the theoretical case where all these data is available, this approach does not reflect the reality at all. In reality, the integral length scale need to be calculated based on the flight-sampled data itself, hence, it is exposed to the same sampling errors as the flux is. Particularly for short track lengths, the errors in the integral length scales are supposed to be remarkably high. These biased values then propagate into the LS86 formula, with the consequence that the standard deviation of the predicted random error is also very high, with far reaching implications concerning the interpretation of the measurement, as the random error estimation cannot be trusted anymore for short tracks. In some situations the LS86 formula indicates high random error, but in other situations it indicates low random errors (please see the discussion in Sühling and Raasch (2013) about this). Hence, the comparison of the LS86-predicted error against the truly random error in the way you have done it here is actually not fair and does not help much in the interpretation of observed data. The manuscript would strongly profit if you add such an analysis. I would propose to add following analysis: \* How does the integral length scale and the correlation coefficient used in Eq. 5 depend on the track length? Here a direct comparison against the true value is possible since the integral length scales from the 3D LES data are available. \* How does the random error behaves when it is calculated directly from the sampled space-series, and how well does it match with the 'true' random error? Based on this, the performance of the LS86 formula can be directly shown for different track length, which would be extremely helpful for researchers.

This is an interesting point. I had previously been assuming that, for a given regime (e.g., neutrally stratified cloud-topped marine boundary layer), integral lengths could be determined and then assumed to be approximately the same in other similar environments. In that case, the domain-averaged (“global”) integral lengths determine in this paper might be usable by other investigators, and there would be no need to determine integral length scales on a case-by-case basis. However, you raise the possibility of the latter, so I have added a new section specifically looking at this question.

Surprisingly, especially in light of the findings of Durand et al. (2000), I find that the the standard deviation in the error estimation is not as large as expected and, even more surprisingly, it seems to converge quickly to a value that doesn't depend much on track length. (I also plan add some comments cautioning against interpreting the computed magnitudes of the standard deviation too literally, owing to an apples-to-oranges aspect of the comparison, including with regard to your final comments about cross-wind vs. parallel sampling errors.)

### 3 Further comments

Abstract - Line 2: vertical turbulent fluxes Line 26 - 30: The paragraph is quite imprecise, it is not clear to what refinements done by Mahrt (1998) or findings by Hollinger and Richardson (2005) the author refer to. The introduction would profit if the author would elaborate this a bit more.

I have revised the comment about Mahrt (1998) to say that he “further examined the sampling problem in the context of the problems posed by non-stationarity.” I have deleted the reference to Hollinger and Richardson.

As an aside, I will say that as a newcomer to the field of boundary layer meteorology, I'm hesitant to attempt too detailed of a review of prior work on this problem because of the near-certainty that I will miss or misstate something important. My hope is that by providing references to the works that I know about, readers can dive deeper into the subject matter by reference to the original sources rather than depending on my inaccurate interpretation or incomplete knowledge of those earlier works.

line 44: Please remove the word ”exceptional”. Though it is indeed quite a large setup, it is not exceptional any more. In the last one to two years such setups have become already standard in the LES community.

I have modified the statement. I would be happy to cite other studies using a similarly large setup, but I have not seen them yet.

line 42-44: To the reader of the manuscript it might appear unclear what is meant by this sentence ”The single most ...”. I suppose the author mean that the smaller scales in the LES are filtered by the subgrid scheme, as well as by numerical errors, being not sampled adequately in an LES.

I have rewritten and restructured much of the introduction in an effort to be clearer.

line 48-49: I disagree with the last part of the sentence. Emulating turbulence measurements in an LES always suffers from the missing subgrid-scale contributions, numerical errors etc., independent which grid resolution is used. Comparing measured turbulence spectra and spectra derived from LES-sampled data will always show a drop-off spectral energy on the smallest spatial scales (about 10 times the grid spacing as a rule of thumb), as it is also the case here (Fig. 2). The relevant question is if, and when, how much the missing subgrid-scale contributions affect the analysis of sampled turbulence data. I would recommend just to rephrase the part with the subgrid-scale flux here.

Your point is well-taken.

end of introduction: I miss a manuscript outline here, to guide the reader through the manuscript.

I have added an outline.

line 79-88: Beyond the fact that a marine boundary layer was simulated, what is the general atmospheric setup? Is this a boundary layer in the trade-wind zones or in a polar region (which latitude). Of course, according to Fig. 1 it becomes clear, but such information should be also given in the text.

I have added some additional information to clarify the meteorological context, which is basically the prevailing summertime close-cell stratocumulus regime off the coast of California.

line 93: How were these spectra calculated? Were they calculated from the emulated flights or from the 3D LES data (1D or 2D spectra)?

I have clarified in the revision that these were calculated from the entire 2-D horizontal domain with subsequent radial averaging.

Fig. 2: Is it temperature or potential temperature? In Fig. 1 profiles of potential temperature are shown, later on the author only refer to temperature, though I cannot find any statement about a transformation.

The transformation is the standard definition based on Poisson's equation. At the low altitudes considered here, the distinction between  $T$  and  $\theta$ , and especially between the perturbations of these quantities, is numerically insignificant.

line 98: With the phrasing "circular eddies" you indirectly imply that the turbulence is isotropic, which isn't the case as plenty of observations and simulation data show. I would recommend to simply remove the word "circular" here.

Removed.

line 99-103: The author describes the spectra sufficiently here, though some references with respect of the minimum domain size of the LES domain are missing. However, I miss some further discussion about possible implications on further results. It is well known that at smaller wavelength the spectral resolution is bounded due to the subgrid-scale model as well as numerical dispersion and dissipation errors, causing these steep drop-off. In most cases this is no big issue as the smaller scales do not contain much energy. But also at the longer wavelengths, the spectral resolution of LES is bounded by the domain size, especially for humidity (de Roode et al. 2004). From previous studies it is known that structures grow in time, meaning that the spectral peaks move towards larger wavelengths, until the structures cannot grow anymore as they are bounded by the domain size. Somehow the spectra for  $q$  indicates this. In case this happens I suppose the integral length scales of humidity are also affected by this, though it won't change much on the overall results I guess. You already bring this up in line 135-137, but maybe it is worth to already bring this up here too.

I have added some more remarks on this issue as well as a reference to de Roode et al.

line 91/104: Though it is only one author, the we-form is repeatedly used throughout the manuscript.

This touches on a controversial subject. Some journals apparently ban the use of first-person (“I” or “we”) altogether, requiring the use of the passive voice to describe what was done (“the data were analyzed”). Other authorities rail against the passive voice, urging the use of the active first person when possible. Philosophically, I agree with the second position, but it is rare and distracting to see the singular “I” in journal articles, and it feels egotistical. Even for a single-authored publication, it seems appropriate to use “we” if the reader is being included (e.g., “we are concerned here with X”). Some additional discussion can be found here: <https://www.editage.com/insights/is-it-acceptable-to-use-first-person-pronouns-in-scientific-writing>

In the end, I have decided to revert to the passive voice except in the “we” case just mentioned.

Fig. 3: Again, is this absolute or potential temperature?

From the purely numerical perspective, it could be either. From a notational perspective, I will fix the inconsistent usage in the original paper.

line 109-11: I disagree with this. The alignment with the mean wind is also visible in  $q'$  and  $T'$ .  $w'$  and  $q'$  (or  $T'$ ) correlate actually fairly well as it is typical in a convective boundary layer where the updrafts are created by buoyancy (which in turn is related to  $q'$  and  $T'$ ).

My statement was, “This directional anisotropy is far less apparent in  $q'$  and is completely absent from  $T'$  at low levels.” It’s admittedly difficult to make objective statements based on a purely visual interpretation of the images, but my statement seems to be consistent with the crosswind and parallel integral length profiles depicted later in Fig. 11 (original numbering, now Fig. XX).

The near-surface correlation between  $T'$  or  $q'$  and  $w'$  is about 0.5, according to Fig. 12 (original numbering, now Fig. XX). I suspect it would be higher still if the directional anisotropy were the same for both variables, but this is admittedly conjecture.

Eq. 1-2 and Fig. 7-8: Here it would be good to already mention that this quantity will be used for the integral length scale calculation. Without this background, which is not clear at this point in the manuscript, this might puzzle some readers. I am not entirely sure what these quantities are actually represent. To my understanding it is the local vertical flux at one point in time which can only be calculated from simulation data in the special case of a horizontally homogeneous boundary layer. To obtain heat fluxes in the traditional sense, time-averaging would need to be applied on top. Hence,  $w'T'$  is not a real flux but a quantity used to compute integral length scales. I would recommend to make this clear in the manuscript, i.e. that these temporal local fluxes cannot be compared one to one to flux measurements from aircrafts or towers, neither with respect to the spatial pattern nor with respect to the amplitude. But in the context of the integral length scale calculation such equations and plots does make sense.

I think the reviewer and I see Eqs. 1 and 2 differently and apparently disagree on the second-to-last sentence above. Throughout this paper, we are looking at flux determinations from a single time step of the LES. We are calculating fluxes as though an airplane were flying at a high enough speed that we can invoke Taylor’s “frozen turbulence” hypothesis—space can replace time (and vice versa) in the Reynolds averaging of fluxes. So spatial integrals of equations 1 and 2 should be equivalent to temporal averages (e.g., stationary observation point, non-zero wind) and are exactly the fluxes we are purporting to evaluate with our transects through the domain. A spatial integral over the entire horizontal domain gives use the “true” domain-averaged flux, and the linear integrals along the flight tracks give us the imperfect estimates of that “true” flux. In short, we’re not just using (1) and (2) to compute integral lengths; we’re using them for the actual simulated flux estimates along flight tracks. I did revise the notation, because another reviewer pointed out that is is more common to use  $H$  and  $LE$  to denote sensible heat flux and latent heat flux, respectively. I also replaced  $T'$  with  $\theta'$  in (1), because that seems more familiar to some reviewers, even though it makes no significant difference in the numerical results.

That said, I have also been advised by others that I should reduce and consolidate the number of figures, and Figs. 7 and 8 were among those highlighted for deletion.

line 128: one dot too much

Thank you.

line 128-129 and Fig. 9: As the author have already mentioned in the text the momentum-flux profiles look quite noisy which is due to the lack of time averaging. Especially for the momentum flux longer averaging periods are required to obtain smooth profiles. Does the author have access to time-averaged profiles in order to check whether the shown momentum flux profiles changes much? According to the wind profiles shown in Fig. 1, I suppose the momentum flux profiles would decrease linearly with height. Here it would be worth to note this in the text and maybe a statement of how much the comparison against emulated flight measurements will be affected, as the truth in momentum flux is just a first guess.

Again, there may be a fundamental difference of opinion here. I am assuming that time averaging is equivalent to spatial averaging under Taylor's "frozen turbulence" hypothesis in the same way that time averaging *is* effectively spatial averaging when performed by a fast-moving aircraft. I don't think there's anything controversial about that. As for the noisiness of the momentum flux profile, one could just as easily state that that noisiness would be reduced by more spatial averaging (larger domain). My conjecture in this case is that the noisiness of the profile is just the result of the domain not being infinitely large and not in perfect steady state. I don't have access to time-averaged profiles, but I agree that sufficient averaging in time, space, or both would probably result in a linear profile.

In any case, the noisiness (in all three spatial dimensions) of momentum flux is certainly relevant to the flux sampling problem, and it is indeed shown in the original Fig. 16 that sampling errors for this variable are quite large relative to its magnitude.

line 133-134: I am wondering why the lower range for the 90-% contribution changes among the different height levels. To calculate an ogive you actually start at the smallest wavelength (or largest wavenumber) and integrate until the target wavelength. The smallest wavelength, however, is fixed by the cut-off wavelength which depends on the grid spacing, so the intervals should be always something like  $[2.5m-x1]$ ,  $[2.5m-x2]$ , etc.

Unfortunately, I don't understand the comment. Yes, the ogives are calculated starting at the shortest wavelength (left end of the x-axis), which is why they all start out at zero there. And yes, the spectra are obtained at multiples of the Nyquist wavelength, but that interval is not a visible feature of the plots, so I'm not sure what the reviewer is referring to.

Fig. 11 f-h: It is not entirely clear how the integral length scales for fluxes are obtained. I suppose you calculated  $I_f$  from the local fluxes obtained from equation 1-3, but it is nowhere written in the text. Or did you use equation 6-7 to calculate  $I_f$ ?

The integral length scales are calculated using the definition given in Eq. (4). I have slightly expanded the explanation of how this was done. They were calculated for individual variables as well as for the local flux contributions in (1) and (2) (or, more precisely, the products  $w'\theta'$  and  $w'q'$ , since the constant coefficients play no role).

line 196-201: It would be nice to add the resulting length scales from equation 6-7 into Fig. 11 f-h to have this also visually.

I thought about that when preparing the original figure, but since representative values from those equations can be easily computed by hand using point values selected from Fig. 11 (old numbering), I decided to leave that figure as a depiction only of actual integral scales rather than risking confusing readers with the very imperfect estimates of  $I_f$  from (6) or (7). A key finding is that (6) and (7) seem to seriously overestimate the true  $I_f$ .

Section 3.3 and 3.4: The information content is not that big to put it into separate sections, you can easily merge it with section 3.2. But that's of course a matter of taste.

Agree.

line 211: Please add the interval for  $n$  directly after the arctan-formula, this will make it easier to understand what is done here.

I have clarified that  $n$  is a positive integer.

line 208-215: Before you start with your flight-track setup, you should first start to explain this special problem which is just because you use a stationary frozen flow field rather than a moving framework with advection where this problem will not occur unless the domain size is large enough. This will make it easier for the reader to understand what is done here.

I have expanded the introduction to this section.

Fig. 13: Just a minor point, but it would be nice to mention the flight direction - from southwest to the northeast, or was it the other way around?

Because we are "instantaneously" sampling the entire flight track, there is no definable direction of flight, only a track position and orientation.



line 234-240: I am not entirely sure what is the effect of the insufficiently sampled true value due to lack of temporal averaging, especially for the momentum flux which has been clearly not converged. But as you compare the flight-sampled data against these true flux, and both are calculated based on the same data set, I suppose that this won't have any effect on the following results. But interested readers might ask the same question, so it could be worth to spend one or two sentences here to make this more clear.

This comes back to the question of whether pure spatial sampling is an adequate proxy for temporal sampling (e.g., a fixed tower) or temporal-spatial sampling (e.g., a moving aircraft with finite speed). This entire paper depends on the answer being "yes." I'm not sure how to reintroduce that question at this point without risking confusing readers more than I help them. However, I have tried to highlight the philosophical concerns surrounding this issue a little more clearly earlier in the paper.

As an aside, my impression was that the Taylor "frozen turbulence" hypothesis was originally invoked as a way to justify using time averages at a point as an acceptable substitute for spatial averages, rather than the other way around. Regardless, more averaging in either time or space, or ideally both, will always lead to reduced statistical error, as long as the turbulence statistics are stationary.

244-245: It is unclear which bias the author does mean here and how this is connected to the minimum separation distance.

If you don't enforce a minimum separation between parallel segments of a track, then eventually the entire domain is sampled so thoroughly that the track estimate converges prematurely on the "true" value and the apparent error goes to zero. This is likely to start becoming a problem when the track separation is smaller than the autocorrelation distance (i.e. integral length) of the field being measured, in which the adjacent paths are no longer statistically independent. The bottom line is that once that minimum distance is violated, the determination of flux error starts to be biased low relative to what would be observed along a path through a truly infinite domain.

line 262: Though you cannot mention it all here, at least some of these previous studies should be mentioned here directly.

I have added a reference to Grossman (1992). There are probably many other papers with similar findings, but with the deadline to resubmit this revision upon me, I will leave it at that unless more are requested.

Caption Fig. 17: Better to explain here directly what the black dots mean rather than to make a detour via the red curves in Fig. 14.

I have expanded the caption as requested. I have also changed the color to be same red as those other curves so as to eliminate one additional possible source of confusion.

Fig. 17-19: please see my major comment

I have added a new figure and a new section to examine the issue you raised in your major comment.

Fig. 14-16: I assume the plot is based on data from all performed flight tracks? Actually this is not a valid approach to show the convergence of the measured flux towards its true value. For a given track length the data shows a certain spread, though, strictly speaking the data from the flight tracks is not comparable since the alignment of the flight tracks vary. In an atmosphere with fully isotropic turbulence the approach would be fine, but in an anisotropic boundary layer the scatter among the crosswind tracks might be different compared to the scatter among the parallel-wind tracks, i.e. the root mean square might depend on the track angle. This should be stated clearly in the text so that readers are aware of this. For example, you can use different colors for the different track angles and provide the the root-mean square values for each track angle separately (e.g. in an additional table), in addition to the "ensemble" root-mean square value you have already given. This might also indicate the optimal track angle with respect to the mean wind direction, supposed the statistics are sufficient.

This comment is of course correct, but the problem is difficult to satisfactorily address. The creation of periodically continuous tracks allows no choice in their orientation; the angle is  $\phi = \arctan(n)$ , with  $n$  an integer that determines not only the orientation but also the track length. The 90 degree rotations and mirror images ensure a symmetric distribution about the points of the compass, but it remains the case that all possible orientations relative to the wind direction occur, and few are perfectly aligned either parallel or perpendicular to the mean wind. So any segregation of tracks into two groups will inevitably include many that are neither. I have some thoughts on how to take this into account, but in view of the looming deadline to post a response to the reviews, I will simply have to work them into the revised manuscript rather than addressing them here.

Fig. 17-19: Is there a difference between parallel and crosswind tracks and any recommendation?

There is clearly a large difference in the integral length at lower altitudes, implying that a proportionally shorter track length is needed in the crosswind direction to achieve the same sampling error in the flux measurement. I'm new to this field, but I believe that this is already widely understood.

line 274-275: You are right with this when you refer to track length of about 10-30 km. But for even shorter track length the LS86 formula tends to underpredict the error as your data show. Actually, everybody knows that the random error for short track length is remarkably high, but flying only short tracks is sometimes the only way in some situations. I would suggest to be precise when you refer to "short" tracks, the understanding of "short" can be sometimes very different.

Understood.

line 296: Here I miss some final conclusion. From the results shown, which version of integral length scale calculation does the author recommend?

If I'm understanding the question, this is apparently about whether (6) and (7) are adequate substitutes for the directly determined integral length from (4), when the latter is not available. It appears from the results presented that (6) and (7) are prone to overestimate the integral length and thus overestimate the sampling error.