

Vainio et al. have made major changes to their manuscript which have greatly improved it, but there are still a few minor points that should be clarified.

I have only one major point:

I had a question on your previous manuscript version on line 225 (the question about the number and distribution of measurements) which I am still not fully understanding. I am not worried about the temporal or spatial comprehensiveness of your measurements, but it would be important to understand when the soil moisture and CH4 flux measurements that were used to train each model (May-July or Aug-Oct models) were taken and how the measurement date could influence your model performance. For example, ideally, you would do your soil moisture measurements across all the sampling locations on the same day, to make sure your measurements are representing the same conditions, and use this information to train your soil moisture model. Now I think your measurements are distributed quite randomly throughout the model period (either May-July or Aug-Oct), making them less comparable with each other. For example, you might have measured one location on a sunny day in early June and another one during a cold and rainy day in late June. Or, you might have measured one sampling location primarily in May, whereas you measured another location only in July. Even though the topographic position, soil texture, and vegetation conditions would be identical, and consequently soil moisture of these two locations should be the same if they would be measured at the same time, the location sampled only in July might suggest that the soil is drier than the other location because it tends to dry out during the summer. What I mean is that some part of the unexplained variability in your soil moisture/CH4 flux models might be related to the measurement time (date) and how warm and rainy it has been before it. And your current predictors do not consider this variability at all since they describe static topographic properties. I would either try to test how the measurement date or air temperature/precipitation conditions prior/during your measurement explain soil moisture or CH4 flux to check how much this could explain the performance issues in your models, or discuss this as a potential reason for the fact that a relatively large amount of variability in soil moisture/CH4 flux remained unexplained somewhere in discussion. You are kind of discussing this on line 620->, but I think you should also mention the comparability issues of your measurements somewhere, if I have understood the sampling scheme correctly.

Minor points:

29: "which was enough" sounds a bit strange to me

30: CH4 -> CH<sub>4</sub>

31: "upscaling predicted stronger CH4 uptake" -> do you have an idea why? Because you might have sampled more wetter environments than what their aerial extent truly was? Maybe this could be discussed briefly in the discussion.

33-35: I still think these points come a bit out of the blue here in the abstract because they haven't been mentioned before. Maybe you could try to link this sampling strategy comment more specifically (but shortly) to your results listed in the abstract, e.g. to the ones related to the differences in measured vs. predicted fluxes (see my comment above)? The main text mentions the sampling strategy in the introduction and discussion, which is great, but I'm just trying to assure that the reader understands the abstract as well as possible too.

35: "..., and the measured fluxes..." seems a bit redundant to be. Doesn't it always make sense to link fluxes to their environmental drivers?

213: distribution was non-normal

237: "The calculations were made with TopoToolbox (Schwanghart and Kuhn, 2010)." -> can be removed, since this is already mentioned on line 234

242: What was the spatial resolution of the other gridded data sets mentioned on lines 228?

252: You could report here how many measurements and sampling locations you used in these two models. I assume that you had data from 60 sampling locations in all models, but it's unclear to me how many observations you had in total (i.e. what was the sample size in your models).

269: I'm not sure I understand the beginning of this sentence/ whether it's formulated correctly: "In this method, one RF model is developed for each sample point...". To me it says that you developed a model for each individual sample point (i.e. model sample size n=1). Maybe mention that the validation data came from one sample point/you predicted the model developed with the training data to each individual sample point or something similar.

278: Perhaps mention here how you calculated uncertainty. Many studies also use a prediction interval of e.g. 2.5<sup>th</sup>-97.5<sup>th</sup> percentiles. Why did you choose to use standard deviations?

279: I think your uncertainty estimate also describes the effect of the distribution of the sampling points.

281: Aalto et al. (2018)

283: Do you think it's problematic that your soil moisture predictions are not entirely independent from the topographic indices that you use as predictors of fluxes as well? Ideally, you'd use entirely independent data sets to train your CH4 flux models. Now your soil moisture predictions already have some information about the topographic indices too.

438: I'd mention somewhere here that your soil moisture and CH4 flux models seemed to have issues particularly with small and large fluxes, which were over- and underestimated, respectively. This might influence your average predicted fluxes too. E.g., if the area covered by dry conditions which are suitable for net CH4 uptake is high compared to the wet CH4 emitting patches, maybe your predictions are underestimating the net uptake and the difference between measured and predicted fluxes could actually be even greater?

537: I would not repeatedly mention the acronym in the discussion SW-W-3

538: Do you have the data about water table depth somewhere?

547: Or is it possible that SW-W-3 was measured after a rainy day and other on drier conditions?

566-599: Could this discussion be its own section? Now that it's listed under upscaling, I start to think again about the limitations in using 15-20 sampling points to train a random forest model, even though I understand that this discussion is not related to that.

582: You could add a few references here. Also, I guess quite frequently mean flux is calculated for each vegetation type, and then a spatially representative mean flux is calculated based on the spatial extent of the vegetation types? Or is that still rare in boreal forest CH4 flux studies? This method of course has issues related to the vegetation type map used, so I think your method is better!

584: "spatially modelled CH4 flux" – are you referring to the modelled estimate for the whole area?

612: But how about vegetation indices derived from WorldView, Planet, or Sentinel-2 that have a pixel resolution of 2-10 meters?

613-619: A large part of this paragraph could actually be located in Section 4.1. which focuses on the drivers.