

## ***Interactive comment on “A synthesis dataset of permafrost-affected soil thermal conditions for Alaska, USA” by Kang Wang et al.***

### **Anonymous Referee #2**

Received and published: 11 June 2018

The dataset described in this manuscript is certainly useful, the manuscript, however, in its current form is not. There is too much ambiguity and missing information about the dataset, the language and organization are confusing in many places, and the presentation is a bit sloppy (especially the figures). This makes the utility of the dataset difficult to assess. The organization of the results section is strange and the paper leaves it unclear what all variables are actually included in this dataset and which ones are just presented for some type of qualitative validation.

Major points:

1. It is not clear what the dataset actually is. The introduction (p. 2 lines 30-31) says it is \*measured\* air and ground temperatures, snow depth, and soil volumetric water content. But then later other variables like frost numbers, thawing index, and freezing

C1

index are mentioned. It is also later written that the data are provided as interpolated values (p. 5 lines 22-23), which is very different from measured values. At the very least, readers should come away from this paper knowing clearly what the dataset actually is.

2. A shortcoming of the dataset is the monthly timescale. While this may be OK for model validation, this is a very limited audience and most non-modelers would probably prefer the daily data. From a practical standpoint, you are disincentivizing users to turn to your synthesis product, given that the daily data are already readily available from the original UAF, USGS, and NPS sources. For example, based on the data you are providing I could not use it to quantify many processes that occur on shorter timescales, such as the onset of thaw or freeze, snowmelt timing, etc.

3. The dataset is supposed to be a synthesis of near-subsurface ground temperature data (p. 2 line 27). However, the "Overview of this dataset" section focuses on volumetric water content, snow depth, and frost number. It is very strange that you are providing an overview of some of the peripheral and derived variables, but not of the primary variables that make up the dataset. In fact, the entire final paragraph of section 3.1 should be deleted, because it suddenly presents research results, as opposed to describing and showcasing the dataset.

4. A crucial missing piece of this dataset is metadata information about the soil itself (soil type, density). This information is already available because, as stated (p. 5 line 4-5), at least the GI-UAF sensor installation was dependent on the soil profile and texture. Given the heterogeneity of soils and therefore its thermal conditions, adding this to the dataset would make it vastly more useful.

5. I question the linear interpolation method employed. Other studies, for example, Sherstiukov (2009) and Streletskiy et al. (2008, 2015) found that a polynomial fit better captures the exponential attenuation of temperature with depth. Regardless, this interpolation must be described with much greater detail if the entire dataset is based

C2

on this interpolation. How many soil depth observations were required for the interpolation? Did you only interpolate between 'adjacent' soil depths? Or, for example, if a USGS location only had a 5 cm and a 1.2 m observation, did you still interpolate and provide temperature at 0.25, 0.5, 0.75, and 1 m? Was the interpolation done on the raw hourly or daily data, or on the final monthly data? Does the final, published dataset only contain these interpolated values, or also the original ones? Do you distinguish between observed and interpolated data in the official dataset? These are all crucial details that cannot be omitted.

6. Because this dataset is comprised of interpolated data (p. 5 line 22-23) anyway, why not also interpolate to fill in missing observations (or did you)? In cases where there are only one or two days of missing soil temperatures at a certain depth, you could relatively reliably fill in that gap based on the average of the previous and next day's observations. And if a certain soil depth has a missing observation but the immediately adjacent upper and lower depth observation is available, the missing depth could also be interpolated.

7. It is peculiar that a linear regression trend analysis was chosen as a core validation technique, given that you acknowledge how there are probably not enough data available to reliably do so. Using trend analysis and other derived variables like frost numbers and effective snow depth to then qualitatively 'eyeball' whether things generally look right are not robust validation techniques. This may uncover glaring issues, but not the subtle non-climatic artifacts and discontinuities that commonly plague climate data.

Specific points:

Can you quantify or estimate the thermal disturbance caused by drilling the holes at the soil temperature measurement sites?

p. 1 line 9: why report that the dataset consists of 41,667 monthly values (to make it seem really big)? This is not a useful statistic but instead you should provide a

C3

percentage of how complete or how much missing data there are based on the overall date range.

p. 3 line 24-25: is this a personal hunch, or can you provide a reference for this statement?

p. 3 lines 28-29: please clarify how or why the thermistors are designed for a low temperature of  $-30^{\circ}\text{C}$ , yet they record down to  $-35^{\circ}\text{C}$ ?

p. 5 lines 16-20: you explicitly mention the temporal availability of the USGS and GI-UAF data, but why not for NPS?

p. 6 lines 12-13 and 14-15: how did you choose those thresholds (20 days and 90%)? Based on those cited references, or did you perform your own cost-benefit analysis to determine how many missing observations you can allow while still obtaining the most continuous monthly/annual dataset?

p. 7 lines 5-6: what is this statement based on? Is this a visual assessment or did you perform a statistical analysis? Was this the only site that experienced a station move or other non-climatic change? Is there a list of all station moves, instrument changes, and other events that could affect the data quality?

p. 7 line 25: what is effective snow depth and how was it calculated?

p. 8 lines 21-24: you are reporting frost numbers without even explaining what those 0.5 versus 0.6 values mean. How is frost number used to indicate permafrost occurrence? What does a 0.5-magnitude frost number indicate?

p. 9 lines 5-7: how or why were those stations chosen? Are they the only ones with 10 or more years of data?

p. 10: is Smith Lake the only instance of multiple sites for the same (general) location? How or why was this Smith Lake example chosen for discussion?

Figure 5: there is enough room on these figures to write out the actual variable names

C4

in the title.

Figure 6: primary y-axis labels overlap (why not write "Air (°C)," "Surface (°C)," "Ground (°C)" and elaborate in the caption that the top three rows show temperature?); I am unsure about the secondary y-axis labels, are they necessary? What does A, B, C, D refer to? What does the asterisk mean? What does the grey shading mean?

Figure 7: what is being shown here on these modified box and whisker plots? What does the circle versus the range indicate? Are all the trends plotted for the identical time period (the caption implies not)? If not, they are not comparable and this plot is misleading. Even if a location has more than 5 years of data, you cannot compare a 6-year trend to a 10-year trend if they have different beginning and end points. Are all the trends significant? What does A, B, C in the legend mean?

Figure 8: what is the grey shading?

Figure 9: needs actual x and y axis titles (instead of acronyms), and units.

The entire manuscript needs to be carefully edited. I am sure one of the 15 authors is a native English speaker who could do this?

---

Interactive comment on Earth Syst. Sci. Data Discuss., <https://doi.org/10.5194/essd-2018-54>, 2018.