

1     **Validation of Satellite Sea Surface Temperature Analyses in the Beaufort Sea Using UpTempO Buoys**

2  
3  
4                     Sandra L. Castro<sup>1</sup>, Gary A. Wick<sup>2</sup>, and Michael Steele<sup>3</sup>

5  
6             <sup>1</sup>Department of Aerospace Engineering Sciences, University of Colorado at Boulder, Boulder, CO

7                     <sup>2</sup>NOAA Earth System Research Laboratory, Physical Sciences Division, Boulder, CO

8                     <sup>3</sup>Polar Science Center, Applied Physics Laboratory, University of Washington, Seattle, WA

9  
10  
11     **Abstract**

12     Many different blended sea surface temperature (SST) analyses are currently available and exhibit  
13     significant differences in the high latitude regions. It is challenging for users to determine which of  
14     these products is most accurate and best suited for their applications. Nine different SST analyses and  
15     two single sensor satellite products are compared with independent observations from Upper  
16     Temperature of the polar Oceans (UpTempO) buoys deployed in the Beaufort Sea in 2012 and 2013  
17     during the Marginal Ice Zone Processes Experiment (MIZOPEX). The relative skill of the different SST  
18     products is evaluated using a combination of Taylor diagrams and two different verification scores that  
19     weight different statistical measures. Skill thresholds based on satellite accuracy requirements are  
20     chosen to map products with similar performance into three discrete skill categories: excellent, good,  
21     and poor. Results are presented for three subsets of the buoys corresponding to different regimes:  
22     coastal waters, northerly waters, and extreme weather. The presence of strong thermal gradients and  
23     cloudiness posed problems for the SST products, while in more homogeneous regions the performance  
24     was improved and more similar among products. The impact of variations in the ice mask between the  
25     SST products was mostly inconsequential. While the relative performance of the analyses varied with  
26     regime, overall, the best performing analyses for this region and period included the NOAA Optimal  
27     Interpolation SST (OISST), the Canadian Meteorological Centre (CMC) SST, and the Group for High  
28     Resolution SST (GHRSSST) Multi-Product Ensemble (GMPE).

29

30 **1. Introduction**

31           Accurate monitoring of environmental conditions in the Arctic warrants particular attention  
32 both for the unique measurement challenges and the potential for the high latitudes to serve as an early  
33 and strong indicator of potential climate change. Sea surface temperature (SST) is a fundamental  
34 variable critical to weather predictions, climate monitoring, and ship-based operations at high latitudes.  
35 An inherent challenge of high-latitude satellite SST production is persistent cloudiness, which hinders  
36 SST retrievals in the infrared portion of the spectrum, resulting in extended gaps in the satellite imagery.  
37 Microwave sensors, although able to “see” through clouds, have limitations close to land and ice. Multi-  
38 sensor, gridded, satellite-based SST analyses (Level 4 or simply L4 satellite products) offer tremendous  
39 potential for monitoring conditions throughout the Arctic domain over extended periods, since they  
40 combine information available from multiple satellites and types of sensors in an objective analysis to fill  
41 in the coverage gaps, but much remains unknown about the accuracy and representativeness of the SST  
42 analyses at these latitudes.

43           A large number of these gap-free SST products have been developed over recent years (Table 1).  
44 This abundance, however, presents the users with the challenge of choosing the analysis product that  
45 best suits their purpose. L4 SST products are available through the Group for High Resolution Sea  
46 Surface Temperature (GHRSSST; Donlon et al., 2007), and are distributed in a common format for easy  
47 use. A thorough comparison of the GHRSSST L4 SST analyses is described in Martin et al. (2012) and Dash  
48 et al. (2012). A persistent problem, however, is that while the different analyses perform fairly  
49 uniformly globally or in basin-wide regions, there are significant differences at high latitudes. As Dash et  
50 al. (2012) point out, mean analysis differences in excess of 2°C are frequently observed in the Arctic  
51 Ocean.

52           The research presented here was conducted in association with the Marginal Ice Zone Processes  
53 Experiment (MIZOPEX). This was a multi-institutional, multi-instrument Arctic observing campaign  
54 (<http://ccar.colorado.edu/mizopex>) led by the University of Colorado with the support of the National  
55 Aeronautics and Space Administration (NASA )and the National Oceanic and Atmospheric  
56 Administration (NOAA), that involved coordinating different types of unmanned autonomous vehicles  
57 (UAVs), with the simultaneous deployment of in situ instruments and satellite overpasses over the  
58 marginal ice zone (MIZ) to measure SST and sea ice during the 2012 – 2013 melt seasons. Our  
59 intercomparison of L4 SST products was motivated by the need of MIZOPEX management and flight  
60 planners to know which of the satellite SST products provided the most accurate information over the  
61 study area. This area is particularly challenging since persistent cloudiness during the melting season  
62 results in few infrared (IR) retrievals, and proximity to land and ice hinders microwave (MW) SST  
63 retrievals. The highest cloudiness in the Arctic (~80 – 90%) occurs from June to October (Przybylak,  
64 2003), encompassing the duration of the MIZOPEX field campaign. The sparseness of IR SST retrievals  
65 made lower level (Level 2 and Level 3) SST products unfavorable for airborne mission planning, and even  
66 though MW SST retrievals are a valid option under these conditions, the Advanced Microwave Scanning  
67 Radiometer for EOS (AMSR-E) traditionally used for MW products was not operating at the time of this  
68 study. The one remaining option was to choose a satellite product from the plethora of L4 SST analyses.  
69 SST maps of the Beaufort Sea from several widely used analyses exhibited stark differences during the  
70 days leading to the field campaign. These discrepancies made many of the SST analyses unreliable for  
71 determining the locations where experiment activities should be conducted, but motivated devising a  
72 framework for choosing the more skillful L4 products for our particular application.

73           Validation of the SST analyses is especially challenging in the Arctic Ocean and in regions near  
74 sea ice due to the limited number of in situ SST observations. An unprecedented set of high quality  
75 buoys have been deployed in the Arctic Ocean by the Polar Science Center of the Applied Physics

76 Laboratory (APL) at the University of Washington every spring and summer since 2010. Though limited  
77 in number, the quality and high resolution of these Upper layer Temperature of the polar Oceans  
78 (UpTempO) buoys provides a unique opportunity to validate the different SST analyses in the region.

79 In this paper, we employ UpTempO buoys to perform a systematic inter-comparison of  
80 multiple SST analyses in the Beaufort Sea during the Arctic summers of 2012 and 2013. The quality of  
81 the individual L4 analyses, measured relative to the UptempO buoys, is demonstrated using a  
82 combination of performance metrics such as Taylor diagrams and skill scores. The main aims are to  
83 assess which of the products performs best in the Beaufort Sea, and to test a methodology for ranking  
84 the skill of the analyses. The UpTempO buoys uniquely facilitated this study by providing high-quality in  
85 situ observations independent from the analyses (at the time of this study, the UptempO SSTs were not  
86 being reported via the Global Telecommunications System (GTS)). Our focus in this study is on the  
87 seasonally open water of the Beaufort Sea, i.e., we avoid areas of sea ice cover. The UpTempO buoys  
88 are described in detail along with the SST analyses in section 2. The collocation approach and evaluation  
89 methodology are presented in section 3, followed by the results obtained in section 4, and conclusions  
90 in section 5.

## 91 **2. Data Description**

### 92 2.1. Level 4 SST Products:

93 SST L4 analyses are interpolated (gap-free), gridded SST products. These analyses assimilate  
94 both IR and MW satellite SSTs, as they are highly complementary and their error characteristics are  
95 independent of each other. The main passive MW instrument used for SST retrievals, AMSR-E, failed on  
96 October 5, 2011, and data from its successor, AMSR2, was first released in January 2014. During the gap  
97 between AMSR instruments, which coincidentally overlaps with our study period, some of the satellite  
98 SST data producers resorted to an alternative MW data source, WindSat, while others abstained from  
99 using any MW data at all, or temporarily halted production of their MW-based SST products. MW data

100 is especially valuable in regions with persistent cloudiness where the “all-weather” coverage of MW  
101 sensors results in significant improvements in accuracy. As Brasnett (2008) points out, for some of these  
102 analyses, the MW and IR data contribute in equal measure to the analysis quality. It is understood then  
103 that the performance (accuracy) of the L4 products compared here, especially those that rely on MW  
104 data, was greatly compromised during the study period by the special circumstances of not having an  
105 AMSR instrument.

106           Most of the SST analyses used here are available in GHRSSST NetCDF format, and can be  
107 downloaded from the GHRSSST Long Term Stewardship and Reanalysis Facility (LTSRF) at the NOAA  
108 National Centers for Environmental Information (NCEI:  
109 [www.nodc.noaa.gov/sog/GHRSSST/accessdata.html](http://www.nodc.noaa.gov/sog/GHRSSST/accessdata.html)). The following sections include brief introductions  
110 to the L4 products used in this inter-comparison, as they are extensively described elsewhere (see Table  
111 1 for key references). Main features and contributing data sources for all the SST analyses are  
112 summarized in **Tables 1 and 2**, respectively. The different analyses can represent different SST  
113 quantities ranging from a simple daily average temperature to the “foundation” temperature  
114 representing the SST at a depth free from diurnal variability (e.g. Castro et al., 2014).

115

116 Table 1. Characteristics of the SST products considered in this analysis.

Level	Product	Data Producer	Spatial Resolution	SST Type	Ice Mask Source	Reference
<b>L4</b>	<b>CMC</b>	Canadian Meteorological Centre	0.20°	Foundation	CMC	Brasnett (2008)
	<b>FNMO</b>	Naval Research Laboratory	9 km	Skin	FNMO	Cummings and Smedstad (2013)
	<b>GAMSSA</b>	Australian Bureau of Meteorology/ BLUElink	0.25°	Foundation	NCEP	Beggs et al. (2011); Zhong and Beggs (2008)
	<b>GMPE</b>	UK Met Office	0.25°	Median Foundation	OSI-SAF	Martin et al. (2012)
	<b>K10</b>	NAVOCEANO	0.10°	Daily average at depth	N/A	
	<b>MUR</b>	NASA JPL	0.01°	Foundation	OSI-SAF	Chin et al. (1998)
	<b>MWIR</b>	REMSS	9 km	Foundation	OSI-SAF	Gentemann et al. (2006)
	<b>OISST</b>	NOAA/NCDC	0.25°	Daily average at depth	NCEP	Reynolds et al. (2007)
	<b>OSTIA</b>	UK Met Office	0.05° (~6 km)	Foundation	OSI-SAF	Donlon et al. (2012)
<b>L3</b>	<b>LAC</b>	NAVOCEANO	2 km	Skin		May et al. (1998)
	<b>WindSat</b>	REMSS	0.25°	Subskin		

117

118 Table 2. In situ and satellite data sources ingested into the different L4 SST products during 2012–2013.

119 Note that while the Geosynchronous (IR Geo) and TMI sensors are included for completeness, they do

120 not provide coverage in the Beaufort Sea, and thus are not discussed in the text.

Data type	In Situ			IR Polar			IR Geo		MW		
	Argo floats	Buoys GTS	Ships GTS	AVHRR NOAA	AVHRR MetOp	MODIS Aqua,Terra	SEVIRI MSG	GOES	TMI TRMM	WindSat	WindSat Ingest
<b>CMC</b>	✓	✓	✓	✓	✓				✓	✓	01/12
<b>FNMO</b>	✓	✓	✓	✓	✓		✓	✓			
<b>GAMSSA</b>		✓	✓	✓	✓					✓	12/12
<b>K10</b>				✓	✓			✓		✓	01/13
<b>MUR</b>		✓		✓		✓				✓	10/11
<b>MWIR</b>						✓			✓	✓	10/11
<b>OISST</b>		✓	✓	✓	✓						
<b>OSTIA</b>		✓	✓	✓	✓		✓		✓		

121

122 The Canadian Meteorological Centre (CMC) SST analysis is tailored to the needs of the CMC  
123 numerical weather prediction (NWP) system (Brasnett, 2008). It merges the observations listed in Table  
124 2 using optimal interpolation (OI, e.g., Gandin, 1965; Daley, 1991) to provide a daily foundation SST  
125 analysis.

126 Fleet Numerical Meteorology and Oceanography Center (FNMOC) SSTs: The US Office of Naval  
127 Research uses its multivariate OI analysis system, the Navy Coupled Ocean Data Assimilation version 3  
128 (NCODA 3DVAR), run operationally at FNMOC, to produce global SST and sea ice concentration analyses  
129 for GHRSSST. The analyses are executed using a 6-hour update cycle with the U.S. Navy ocean forecast  
130 model, the global Hybrid Coordinate Ocean Model (HYCOM), and are available within 6 hours of real-  
131 time. For the purpose of this intercomparison, only the 12:00 UTC –SST forecast will be used. The  
132 system assimilates satellite SSTs, in situ SSTs, temperature and salinity profiles, altimetric sea surface  
133 heights, and satellite sea ice observations. The analyses have a 12-km resolution at the equator and 9-  
134 km resolution at mid latitudes. The FNOMC GHRSSST analyses are available through the US GODAE  
135 server at <http://www.usgodae.org>.

136 Global Australian Multi-Sensor SST Analysis (GAMSSA): the Australian Bureau of Meteorology  
137 produces this daily foundation SST analysis on a 1/4° grid, and it is used operationally as a boundary  
138 condition in their global NWP system and to initialize their seasonal forecast system. The GAMSSA is an  
139 extension of their 1/12° regional L4 product (RAMSSA: Beggs et al., 2011). The OI system ingests in situ  
140 SST and both IR and MW satellite SST data (Zhong and Beggs, 2008). Data are rejected for low NWP  
141 wind speed thresholds (6 m/s day, 2 m/s night) to reduce effects from diurnal warming on the analysis.

142 The Naval Oceanographic Office (NAVOCEANO) K10 SST analysis uses satellite data only, and  
143 combines the L2 SST products in a weighted average tuned to represent the SST at 1-m depth. This is  
144 one of the few L4s that does not use OI techniques. All the IR inputs are produced by NAVOCEANO  
145 using separate nonlinear regressions trained against quality-controlled GTS drifting buoys from the

146 previous month. To preserve features, the weights decrease exponentially from the center of the  
147 averaging window and the elapsed time from the last observation (B. McKenzie, personal  
148 communication, 2011).

149 The Multi-scale Ultra-high Resolution (MUR) SST analysis is produced daily by the NASA Jet  
150 Propulsion Laboratory (JPL). In contrast to other L4s for which more traditional OI techniques are used,  
151 the MUR system uses a statistical interpolation method based on wavelet decomposition called Multi-  
152 Resolution Variational Analysis (e.g., Mallat 1989). This multiscale signal reconstruction technique is  
153 particularly suitable for dealing with the multiple spatial resolutions of the L2 products entering the  
154 analyses and the irregular swath patterns of the different satellites (Chin et al., 1998). The main  
155 contribution of this product is its fine spatial (horizontal) resolution and capability for resolving high-  
156 resolution SST features such as fronts.

157 The Remote Sensing Systems (REMSS) MW-IR (MWIR) SST product uses an OI analysis and  
158 satellite data only (<http://www.remss.com/measurements/sea-surface-temperature/oisst-description>).  
159 Inputs are diurnally corrected using an empirical diurnal warming model. This foundation SST product  
160 was originally designed for the National Hurricane Center to be used in conjunction with the Statistical  
161 Hurricane Intensity Prediction Scheme (SHIPS) model for hurricane intensity forecasting.

162 The NOAA OISST (Reynolds et al., 2007) is generated by NCEI (formerly the National Climatic  
163 Data Center). While both IR-only and IR-MW products are normally generated, the L4 analysis herein  
164 refers to the IR-only product due to the gap in AMSR data. In addition to the IR SSTs this product uses in  
165 situ data from ships and buoys as well as proxy SSTs, generated from sea ice concentrations, for the MIZ.  
166 The product was designed for applications that target high-resolution features such as fronts and  
167 hurricane forecasting, and to serve as boundary condition for atmospheric models. It represents a daily  
168 average SST (no diurnal warming correction is attempted), that is bias-adjusted using a spatially  
169 smoothed, 7-day in situ SST mean.



170 The Operational Sea Surface Temperature and Sea Ice Analysis (OSTIA) analysis (Donlon et al.,  
171 2012) is produced daily by the UK Met Office and used operationally as a boundary condition in NWP  
172 and Numerical Ocean Forecast systems at the Met Office and the European Centre for Medium-range  
173 Forecasting (ECMWF). Using an OI analysis the OSTIA system normally assimilates MW data from AMSR,  
174 but abstained from ingesting WindSat SSTs during the gap period. Input data are filtered to remove  
175 daytime observations with winds < 6 m/s to eliminate possible instances of diurnal warming. While  
176 provided on a 0.05° (~6 km)-grid the OSTIA SSTs are effectively smoother by design (Donlon et al., 2012).

177 The GHRSSST Multi-Product Ensemble (GMPE) system (Martin et al., 2012), developed and  
178 operated by the UK Met Office, consists of the daily ensemble median and standard deviation, on a  
179 homogenized 0.25°grid, of various GHRSSST L4 operational analyses. Inputs contributing to the ensemble  
180 median include all of the above (with the exception of MUR at the time of this study), plus three other  
181 SST analyses not included in this study (Martin et al., 2012). Although the original purpose of the GMPE  
182 system was to advise GHRSSST data users on the relative performance of the different L4 products by  
183 providing a near-real time global ensemble from a large number of L4 SST analyses, the system has  
184 proven useful for climate-related SST application. The GMPE data is available via the MyOcean project  
185 (<http://myocean.eu.org>).

186 Please note that while most analyses ingest AVHRR (NOAA and/or Metop) L2 SSTs (see Table 2),  
187 there are multiple AVHRR data providers; hence, the source of the AVHRR data can differ among L4  
188 products. Additionally, different AVHRR satellites can be active at any given time. During the study  
189 period, NOAA-16, -18, -19, and Metop A and B were all being used, with NOAA-19 being the designated  
190 operational afternoon orbit and Metop (A in 2012, B in 2013) the morning orbit. For additional  
191 information about the specific AVHRR data sources and sensors being used, please consult the reference  
192 for the appropriate analysis or the metadata source field available in all GHRSSST-compliant SST products.  
193 While Table 2 lists the different L2 products ingested during the study period, another important IR

194 sensor, the Advanced Along-Track Scanning Radiometer (AATSR), was routinely ingested by some of the  
195 analyses (i.e., CMS, FNMOC, GAMSSA, and OSTIA), but ceased operations in April 2012.

## 196 2.2. Ice Masking

197 A main difference among the L4 SST analyses is their treatment of the SSTs near or under ice.  
198 Most of the L4s use independent sea ice concentration (SIC) analyses, generated from space-borne  
199 microwave sensors, to derive an ice mask based on some ice concentration threshold. That is, if the SIC  
200 in an analysis grid exceeds a minimum ice fraction,  $I_o$ , then the corresponding grid cell in the SST  
201 product is flagged as ice. While some SST producers opt for not reporting SSTs once  $SIC \geq I_o$ , others  
202 use ice information to compute proxy SSTs in the range  $[I_o, 1]$ , i.e., the SSTs are relaxed towards the  
203 freezing point temperature of seawater using empirical relationships between SIC and SST. At the time  
204 of this study, the only L4s that simulated proxy SSTs in  $[I_o, 1]$  were the OISST, OSTIA and FNMOC.  
205 Detailed descriptions of the different methodologies used to simulate SSTs under ice can be found in the  
206 product references in Table 1. All other L4s set  $SST = -1.8^\circ\text{C}$  in locations where  $SIC \geq 0.5$ .

207 The two most widely used SIC analyses are the ones produced by the NOAA National Centers for  
208 Environmental Prediction (NCEP)-Marine Modeling and Analysis Branch (MMAB) and the European  
209 Organisation for the Exploitation of Meteorological Satellites (EUMETSAT) - Ocean and Sea Ice Satellite  
210 Applications Facility (OSI-SAF). The newest (starting June 2012) NOAA/NCEP operational SIC analysis  
211 (Grumbine, 1996), distributed daily on 12.7-km polar stereographic projection hemispheric grids, uses  
212 the NASA Team 2 (NT2) sea ice retrieval algorithm (Markus and Cavalieri, 2000) and ice retrievals from  
213 the Special Sensor Microwave Imager (SSM/I) and Special Sensor Microwave Imager/Sounder (SSMIS).  
214 The SST analyses that rely on the NCEP SIC product for ice information are GAMSSA and the OISST. The  
215 OSI-SAF SIC product (Andersen et al., 2007) is based on the SSM/I sensors and is distributed on a 10-km  
216 polar stereographic projection grid every 24 h. The MUR, OSTIA, GMPE, and MWIR use the OSI-SAF SIC  
217 analysis for their ice mask (see Table 1). The FNMOC and CMC, being part of fully integrated NWP

218 systems, construct their own SIC analyses. The FNMOC ice analysis system (Cummings and Smedstad,  
219 2013) assimilates SSM/I and SSMIS ice retrievals, and use the NT2 to calculate SIC, using 6-h forecast  
220 windows. At CMC, the Global Ice Ocean Prediction System (GIOPS) assimilates passive MW satellite  
221 observations together with manual analysis from the Canadian Ice Service to provide a daily, global ice  
222 (and ocean) analysis (Buehner et al., 2013; Smith et al., 2015). The K10 is the only L4 that did not use an  
223 ice mask at the time of this study, relying on sea-ice extent climatologies instead. An ice mask for the  
224 K10 was not introduced until 2016.

### 225 2.3. Single Sensor SST Products

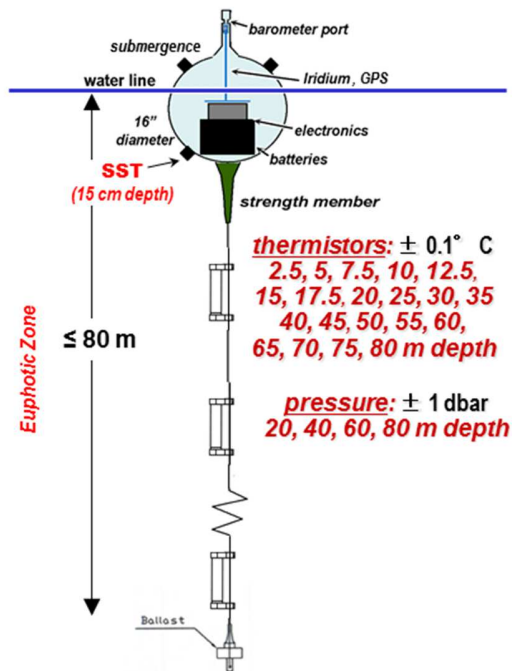
226 Even though the emphasis of this study is on the L4 products, we have also included direct  
227 comparisons against two different single-sensor SST products used as input in some of the analyses  
228 under consideration. These satellite SST products are lower processing-level data, i.e., there has been  
229 no intervention to fill in the gaps, but they have higher spatial resolution. They can be distributed in the  
230 swath of the sensor (Level 2 or L2), or can be gridded for distribution (Level 3 or L3), taking care of  
231 preserving the gaps during the gridding process. These are the NOAA Advanced Very High Resolution  
232 Radiometer (AVHRR) Local Area Coverage (LAC) SSTs produced by NAVOCEANO and the REMSS WindSat  
233 SSTs. The reason for including these products is two-fold: it helps interpret differences among the  
234 multiple L4 products and provides some information on potential limits to the accuracy of the analyses.

235 The LAC SST (obtainable from NAVOCEANO on request) is a 2-km L2 product. The native AVHRR  
236 LAC radiances have 1.1 km resolution at nadir, but retrieved signals from adjacent pixels are averaged  
237 into 2x2 pixel windows before being ingested in the NAVOCEANO SST algorithm. This effectively  
238 reduces the resolution of the L2 SST product to ~2 km. The L2 LAC SSTs were further mapped in house  
239 onto 2 km daily grids to facilitate comparison with the analyses. Hence, technically speaking, this  
240 product was transformed into a “collated level 3,” and hereafter it will be referred to as a L3, despite the  
241 fact that the resolution of the grid did not allow for further aggregation (oversampling) of the data and it

242 is distributed in L2 format. The AVHRR LAC SSTs from NOAA-19 were selected for this study. The  
243 WindSat SSTs (available from [www.remss.com](http://www.remss.com)) is a L3 product, mapped onto a global 0.25°x0.25°  
244 regular grid. WindSat SSTs were largely ignored before the AMSR gap because it was a demonstration  
245 mission sponsored by the U.S. Navy and did not always satisfy timeliness requirements for an  
246 operational mission. Including this product offers a good opportunity to assess the impact of  
247 assimilating WindSat SSTs on the quality of the analyses. This is an important issue for future satellite  
248 SST production since, unlike the diversity of IR spaceborne sensors, there is no redundancy of MW  
249 sensors.

#### 250 2.4. UpTempO Buoys

251 The UpTempO buoys discussed here all consisted of a drifting surface float and a 60-80 m long  
252 string of thermistors. A schematic is provided in **Figure 1**. The buoys reported hourly via Iridium  
253 satellite until individual sensors or the entire buoy failed (often this happened in fall or winter during ice  
254 rafting and ridging). A full description of these buoys will be provided in a forthcoming manuscript  
255 (Steele et al., 2016). Here we provide a brief description of the particular buoys used in this study.  
256 While all of these buoys had a string of thermistors, here we generally only use the uppermost sensor in  
257 order to determine SST. Deeper thermistor data are in fact used initially, but only to show the frequent  
258 presence of an isothermal surface layer. Three different manufacturers supplied the drifting buoys used  
259 in this study:



260

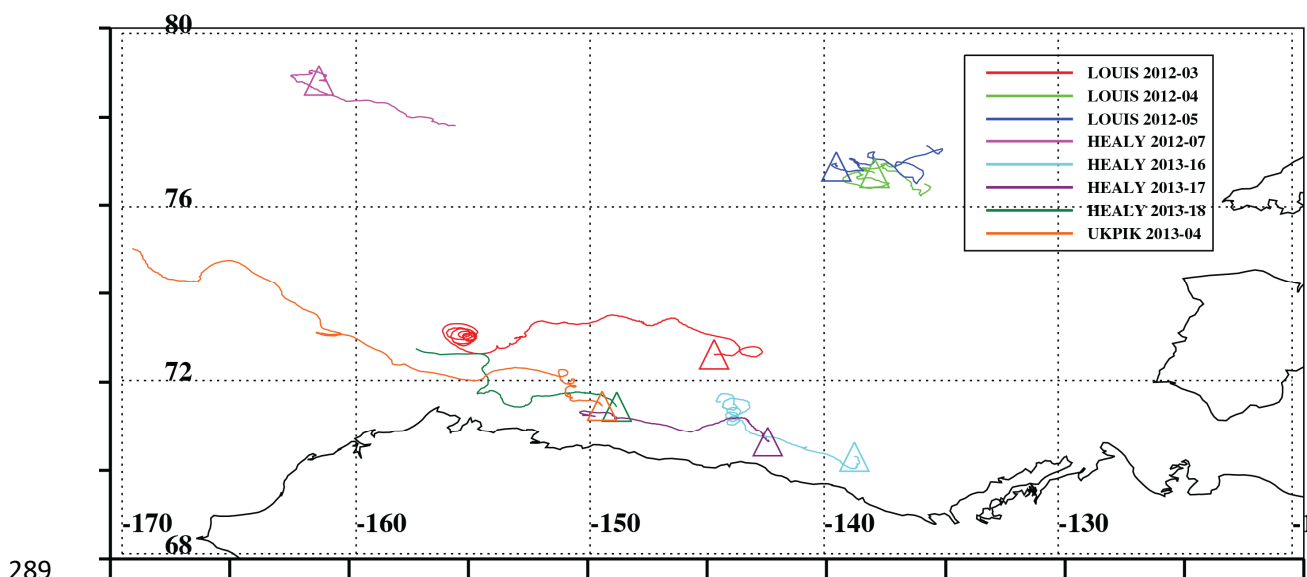
261 Figure 1. UpTempO buoy schematic showing key components and typical depths for temperature  
 262 sensors.

263 (1) Louis 2012-03 and Louis 2012-04 were made by MetOcean Data Systems in Bedford, Nova  
 264 Scotia, Canada. These had a surface hull with electronics, batteries, Iridium satellite antenna, sea level  
 265 pressure sensor, surface atmospheric temperature sensor on a 1-meter mast, but no SST sensor. Below  
 266 the hull was a 60 m long sensor string composed of 12 thermistors and 2 ocean pressure sensors. The  
 267 uppermost thermistor (the only one used in this study) was at 2.5-m nominal depth (i.e., when the  
 268 sensor string was vertical). This sensor was a SBE 39 from Seabird Electronics, with manufacturer-stated  
 269 accuracy of  $\pm 0.002^{\circ}\text{C}$ .

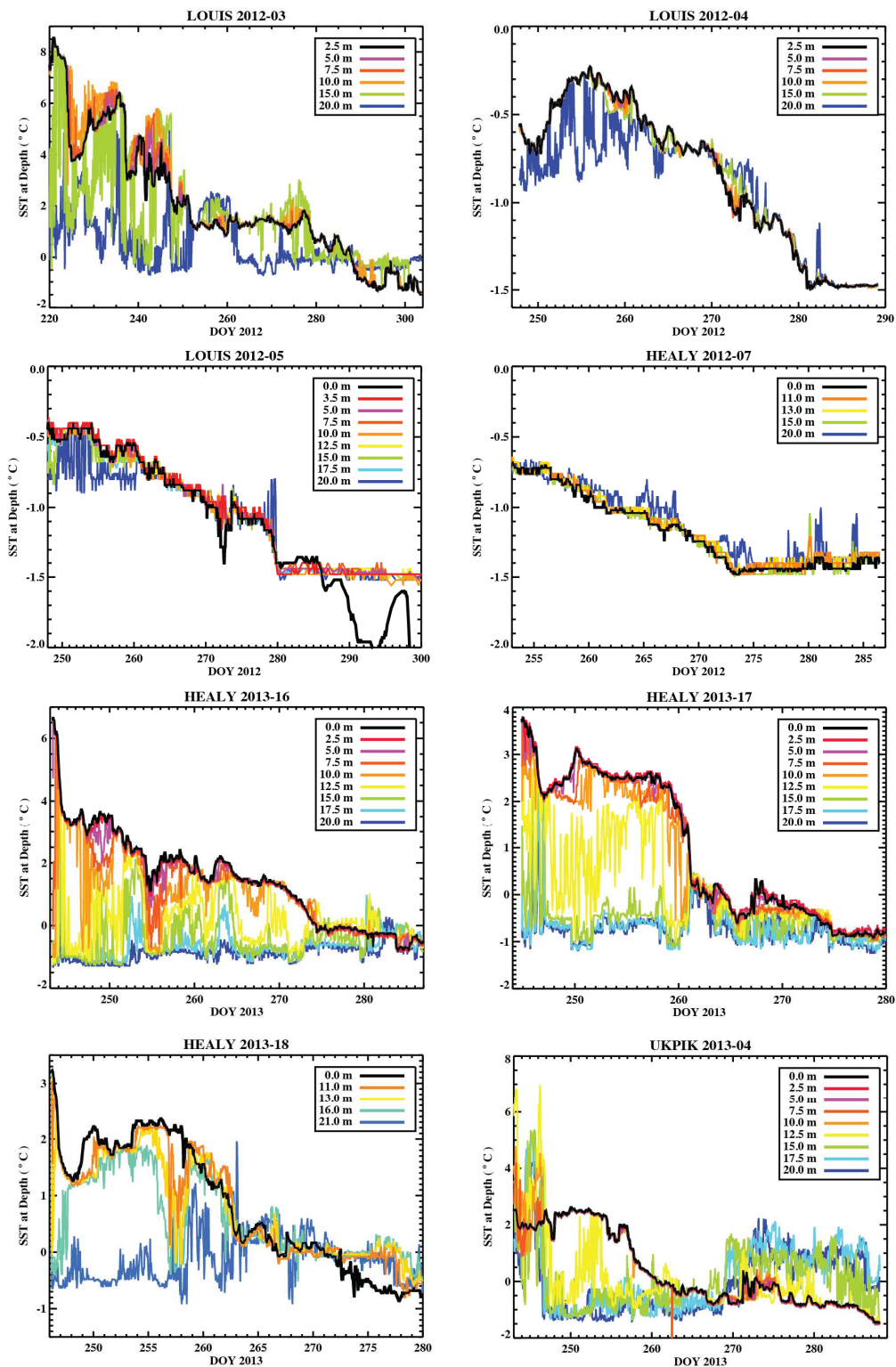
270 (2) Louis 2012-05, Healy 2012-07, and Healy 2013-16, -17, and -18 were made by Marlin-Yug Ltd  
 271 in Sevastopol, Ukraine. The surface float had the same features as the MetOcean buoys, but without  
 272 the atmospheric temperature mast and with a thermistor in the bottom half of the hull to provide SST  
 273 with  $\pm 0.1^{\circ}\text{C}$  accuracy at 15-cm nominal depth, which is what we used in this study. These buoys also  
 274 had a sensor string with 16 thermistors and one ocean pressure sensor.

275 (3) Ukpiik 2013-04 was made by Pacific Gyre Inc. in Oceanside, California, USA. The surface float  
276 had similar features to those made by Marlin-Yug, including an SST sensor with  $\pm 0.1^\circ\text{C}$  accuracy at 15-cm  
277 nominal depth (used in this study), but also including an anemometer. The sensor string had 12  
278 thermistors, 3 ocean pressure sensors, and one Seabird SBE 37-IM for recording ocean temperature and  
279 salinity.

280 A map of the track of the UpTempO buoys deployed in the Beaufort Sea in 2012 and 2013 is  
281 shown in **Figure 2**. The deployment location is indicated by a triangle. Time series of SST-at-depth from  
282 the thermistors in the top 20 m for all the buoys under consideration are presented in **Figure 3**. Note  
283 that the UpTempO temperatures shown in this figure are all fairly uniform throughout the upper 10 m.  
284 The absence of significant thermal gradients near the surface facilitates the comparisons with the  
285 satellite analyses as they characterize different types of SST (Table1), mostly foundation temperatures.  
286 The well-mixed, isothermal surface layer, however, suggests that there were no differences among the  
287 L4s attributable to the SST type. For each buoy, the observations closest to the surface are utilized in  
288 the comparisons.



290 Figure 2. UpTempO buoy tracks deployed in the Beaufort Sea during the Arctic Summers of 2012 and  
291 2013. The triangles indicate the initial deployment location.



292

293

294

Figure 3. Time series of temperatures from the UpTempO thermistors in the top 20 m for all the buoys used in this study. The color coding corresponds to the sensor depth as indicated in the legends.

295

### 296 3. Theoretical Background

#### 297 3.1 Taylor Diagrams

298 Taylor diagrams (Taylor, 2001) are the selected method for comparing the different SST  
299 analyses, since they provide a means of summarizing the relative accuracies of several competing  
300 products graphically, in such a way that they can convey information much more readily and concisely  
301 than the analogous tabular presentation. In its basic form, a Taylor diagram is a two-dimensional scatter  
302 plot in which discrete points give an indication of how closely the various L4 SSTs resemble the  
303 UpTempO observations in terms of their correlation ( $\rho$ ), centered root-mean-square ( $RMS'$ ) error, and  
304 standard deviations ( $\sigma$ ), all at once. These three statistics are related by the equation:  $RMS' =$   
305  $1/N \sum_{n=1}^N [(\sigma_{sat} - \bar{\sigma}_{sat}) - (\sigma_{obs} - \bar{\sigma}_{obs})]^2$ . The diagram is built by plotting a triangle in a rectangular  
306 coordinate system with one point at the origin, and the other two representing the buoy observations  
307 (located along the abscissa), and the corresponding satellite SST matchups, respectively. The radial  
308 distances to the satellite and the buoy SSTs are proportional to their respective standard deviations  
309 ( $\sigma_{sat}$  and  $\sigma_{obs}$ ), whereas the distance from the satellite product to the buoy observations (leg opposite  
310 the origin) is proportional to the  $RMS'$ . The correlation coefficient is given by the cosine of the  
311 azimuthal angle between the radii for the buoy observations and the satellite retrievals. It is important  
312 to know that the means of the SST products are subtracted out before computing  $\sigma_{sat}$  and  $\sigma_{obs}$ , so the  
313 Taylor diagram does not provide information about the overall biases; just about the centered errors.

314 We aim to compare the performance of L3 and L4 products relative to the diagnostic buoy data  
315 set via Taylor diagrams. There is, however, a mismatch in the number of collocated satellite – buoy pairs  
316 between the L4 and L3 comparisons due to the inherent gaps in the L3 products. Because of the  
317 different sampling sizes, the statistics must be standardized before constructing the Taylor diagrams.



318 One way to do this is to normalize the RMS' and all the standard deviations by the standard deviation of  
319 the matched observations, i.e.,  $\widehat{RMS}' = RMS'/\sigma_{obs}$ ,  $\hat{\sigma}_{sat} = \sigma_{sat}/\sigma_{obs}$ , and  $\hat{\sigma}_{obs} = 1$ .

320 The merits of the different SST products can be inferred visually just by looking at their position  
321 in the diagram. The closer the point representing a satellite product is to the buoy observations, the  
322 better the agreement between the two (satellite products lying near the observations have relatively  
323 high correlation and low RMS). In normalized diagrams, these are the points closest to the arc for  $\hat{\sigma}_{obs}$ .  
324 Differences between their respective variances, however, will tend to increase the RMS' pushing the  
325 point farther away from the observations. Point spread from the  $\hat{\sigma}_{obs}$ -arc along the radial distance  
326 from the origin will reflect differences in amplitude between the satellite product and the observations.  
327 The correlation coefficient, on the other hand, remains unaffected by differences in variance but is  
328 sensitive to the relative phasing between the estimates and their associated observations. Hence,  
329 differences in phase will be reflected in azimuthal angle spread (the angle between the point  
330 representing an SST product and the x-axis). As a result, it is possible to have satellite products whose  
331 SST patterns are uniformly too weak or too strong, and still have high correlation with the buoy  
332 measurements or, alternatively, have SST patterns that do not necessarily align with the observations  
333 despite having the right amplitude variability. In other words, the magnitude and the direction of the  
334 scattering of the points representing the SST products in the Taylor diagram helps determine whether  
335 the overall errors in the SST products are attributable to differences in variance or to poor pattern  
336 correlation.

### 337 3.2 Skill Scores

338 Even though Taylor diagrams are very useful in identifying common performance patterns at a  
339 glance, we also desire a verification score that helps us quantify the relative skill of the different L4s at  
340 high latitudes. Under such a scoring system, satellite products with reduced RMS' should be rewarded,  
341 since this indicates close agreement with the observations. The issue remains as to what to prioritize

342 next: pattern similarity or correct amplitude variability. Since there is no universal skill score that  
 343 satisfies all criteria simultaneously, we will consider two different skill scores that emphasize slightly  
 344 different aspects of the product performance.

345 A basic verification score, proposed by Taylor (2001) to evaluate the skill of several precipitation  
 346 models, is given by:

$$347 \quad TS = \frac{4(1+\rho)^4}{(\hat{\sigma}_{sat} + 1/\hat{\sigma}_{sat})^2(1+\rho_{max})^4} \quad (1)$$

348 where  $\rho_{max}$  is the maximum potentially realizable correlation given the uncertainty associated with  
 349 unforced variability. This skill score is defined to vary from one (most skillful) to zero (absence of any  
 350 skill) and fulfills the attributes stated above, i.e.,  $TS \rightarrow 1$  when  $\hat{\sigma}_{sat} \rightarrow \hat{\sigma}_{obs} = 1$  and  $\rho \rightarrow \rho_{max}$ , and  
 351  $TS \rightarrow 0$  as  $\hat{\sigma}_{sat} \rightarrow 0$  or  $\hat{\sigma}_{sat} \rightarrow \infty$  and  $\rho \rightarrow 0$ . Note also that when  $\hat{\sigma}_{sat} \rightarrow 0$ ,  $TS \propto \hat{\sigma}^2$ , and when  $\hat{\sigma}_{sat} \rightarrow$   
 352  $\infty$ ,  $TS \propto 1/\hat{\sigma}^2$ ; i.e, for products with small variance, the skill is proportional to the variance, but when  
 353 the variance is large, the skill is inversely proportional to the variance; thus, the skill always decreases  
 354 when the RMS increases.

355 A potential deficiency of the Taylor verification score and Taylor diagram may arise from the fact  
 356 that centered moments are being used to correct for non-zero mean SST biases. Dash et al. (2012) have  
 357 found differences  $>2^\circ\text{C}$  among the L4 SST products at high latitudes. Consequently, if large  
 358 unconditional biases exist in some of the satellite-derived SST products, the  $TS$  score can substantially  
 359 overestimate their skill. Following Murphy (1988), we propose a complementary skill score that uses  
 360 the un-centered second moments, to account for potential biases in the satellite products being  
 361 evaluated. For a sample of  $N$  pairs of matched SST estimates and observations, the “uncorrected” MSE  
 362 can be expressed as:

$$363 \quad MSE(SST_{sat}, SST_{obs}) \equiv \frac{1}{N} \sum_{n=1}^N (SST_{sat} - SST_{obs})^2. \quad (2)$$

364 A skill score (SS) that takes into account the tradeoffs between possible biases and the variance can  
 365 formulated as the loss function:

366 
$$SS(SST_{sat}, SST_{clim}, SST_{obs}) = 1 - [MSE(SST_{sat}, SST_{obs})/MSE(SST_{clim}, SST_{obs})]. \quad (3)$$

367 Because the climatological reference can be defined, at least hypothetically, based on a sample of  
 368 observations from the experimental period (Murray, 1988), we use the mean of the collocated

369 UpTempO temperatures as our climatological reference, i.e.,  $SST_{clim} = \overline{SST}_{obs}$ . From (2) it follows that

370  $MSE(SST_{clim}, SST_{obs}) = \sigma_{obs}^2$ . Adding and subtracting  $\overline{SST}_{sat}$  and  $\overline{SST}_{obs}$  within the parentheses on

371 the RHS of (3) and expanding the binomial formula, yields

372 
$$MSE(SST_{sat}, SST_{obs}) = RMS'^2 + (\overline{SST}_{sat} - \overline{SST}_{obs})^2. \quad (4)$$

373 Substituting  $MSE(SST_{clim}, SST_{obs})$  and (4) into (3), it follows that

374 
$$SS(SST_{sat}, \overline{SST}_{obs}, SST_{obs}) = 1 - \widehat{RMS}'^2 - (\overline{SST}_{L4} - \overline{SST}_{obs})^2 / \sigma_{obs}^2, \quad (5)$$

375 where the second term on the RHS of (5) is the square of the normalized centered RMS used in

376 constructing the Taylor diagrams; the last term, the square of the mean error normalized by the

377 variance of the observations, is a non-dimensional measure of the overall bias in the SST estimates. This

378 term vanishes only when the satellite estimates are unbiased. Thus,  $SS = 1$  (perfect score) when

379  $\widehat{RMS}'^2 = 0$  and  $\overline{SST}_{sat} = \overline{SST}_{obs}$ . Since the latter two terms in (5) are nonnegative and are preceded

380 by a negative sign, the skill of the SST products tends to decrease as both of these terms increase.

381 Furthermore, because there is no upper bound on the growth of the quadratic terms in (5), negative skill

382 scores can, and will, occur when the analyses have poor or no skill ( $SS$  is defined in the interval  $(-\infty, 1]$ ).

383 This decomposition implies that  $\widehat{RMS}'^2$  is in itself a sort of idealized skill score, attainable when biases

384 are eliminated. In summary, the Taylor diagram and the  $TS$  score can be viewed as measures of

385 potential skill in the absence of any bias (systematic errors of mean bias and standard deviation are

386 intrinsically removed in their computations), while the  $SS$  can be viewed as a measure of actual skill

387 since it incorporates biases (it allows more spread relative to the observations).

388 3.3 Thresholds for Product Performance Classification

389 We will categorize the high-latitude performance of the SST products by ranking the scores from  
390 Equations (1) and (5), into three discrete categories: excellent, good, and poorly skilled. To do so, we  
391 need to establish some basic classification rules for mapping the skill scores into these discrete  
392 categories. This is commonly done via a threshold choice method. Even though there are many  
393 methods for determining decision thresholds (Hernández-Orallo et al., 2012), we opted for the simplest  
394 of all, which is a score-fixed threshold classifier; that is, we use two predefined score thresholds, such  
395 that the satellite products are assigned into one of the three categories mentioned above if their score is  
396 within the limits established by the corresponding thresholds. Because in reality there is no perfect  
397 separation between categories, the threshold, once fixed, can dramatically impact the performance  
398 ranking of the products being compared.

399 A threshold,  $T$ , is usually determined by estimating the cost incurred in misclassifying a product,  
400 and setting the threshold to the value that minimizes the expected loss over different conditions ( $L(t)$ ).  
401 Hernández-Orallo et al. (2012) showed that when dealing with fixed-score thresholds, the accuracy is  
402 the performance metric that minimizes the expected misclassification losses. In this study, we fix the  
403 score thresholds based on the accuracy requirements of operational, real time satellite SST products ,  
404 i.e.,

$$405 \quad T = L(t) = 1 - \sigma^2. \quad (6)$$

406 where the standard deviation,  $\sigma$ , is the adopted GHRSSST convention for SST product accuracy. GHRSSST  
407 has established that for global open ocean products, the user accuracy requirement is  $\sigma < 0.4K$ , with a  
408 tighter user requirement ( $\sigma < 0.3K$ ) for coastal and high-resolution (<2 km) products (Donlon et al.,  
409 2007). Satellite SST L2 products, however, display higher levels of error in the Arctic Ocean with Metop-  
410 A and AVHRR-GAC SSTs showing standard deviations between 0.4 and 0.5°C and MODIS and AMSR-E  
411 between 0.5°C and 0.8°C (Hoyer et al., 2012). In a different study, Martin et al. (2012) found that all the

412 L4 products considered here have  $\sigma < 0.7K$  globally, with the GMPE, CMC, GAMSSA, K10, and OSTIA  
413 having  $\sigma < 0.5K$ .

414           Given the current achievable accuracies of satellite SST products, we assume that  $\sigma$  between  
415  $0.5K - 0.7K$  is a reasonable, if somewhat conservative, absolute accuracy for high-latitude L4 products.  
416 From Equation (6) it follows that the thresholds for  $\sigma = 0.5K$  and  $0.7K$  are  $T = 0.75$  and  $T = 0.51$  (labeled  
417  $T_{75}$  and  $T_{51}$ ), respectively. Note that the equation (6) is independent of score and thus  $T_{75}$  and  $T_{51}$  will  
418 be applied to the two score metrics tested here. In summary, L4 products with skill scores in the  $(0.75,$   
419  $1.0]$  interval will be classified as having excellent high-latitude performance; products with scores in the  
420  $[0.51, 0.75)$  interval will be grouped in the good performance category, and products with scores  $< 0.51$   
421 will be labeled as having poor performance under our set of operating conditions.

422

#### 423 **4. Methodology**

##### 424 4.1 Data Grouping

425           For analysis purposes, the data are divided into three groupings. Among the buoys deployed  
426 during the summer of 2012, Louis 2012-03 is singled out, since this buoy overlapped with a rare weather  
427 event and strong SST gradients, and thus, it is assumed that it was operating under “extreme  
428 conditions.” Louis 2012-04, Louis 2012-05, and Healy 2012-07, on the other hand, were deployed later  
429 in the summer season (September 4, 5, 10 respectively) in cold waters farther offshore (North of  $76N$ , a  
430 more quiescent SST environment), and propagated in a cyclonic direction (opposite to Louis 2012-01)  
431 toward the interior of the Arctic Basin. The latter buoys are grouped together under the “cold northerly  
432 waters” category. The 2013 buoys remained closer to coast, and drifted westward toward the Chukchi  
433 Sea, displaying similar conditions to those sampled by Louis 2012-03 after the extreme weather event  
434 dissipated; hence, the latter portion of Louis 2012-03 and the 2013 buoys are grouped together under  
435 the “coastal waters” category.

## 436 4.2 Collocation Criteria

437           The first step is to construct a matched set of the buoy observations and corresponding satellite  
438 SST products. Collocation of the buoy data with the gridded analyses is straightforward. Buoy  
439 observations are simply compared with the values for the grid cells containing the buoy position on that  
440 day. Multiple buoy observations on a given day were all matched independently with the daily analysis.  
441 Given the gaps in the L3 products, additional steps were taken. Matchups with the AVHRR LAC data  
442 were constructed for satellite observations agreeing within 10 km of the buoy position. Since the  
443 WindSat data are provided in separate ascending and descending grids, available observations for the  
444 grid cell containing the buoy for the orbital segment closest in time to the buoy measurement were  
445 used. This implies a maximum temporal separation of approximately 12 hours.

## 446 4.3 Implementation

447           Collocated satellite-buoy SST pairs were segregated using different UpTempO buoy  
448 combinations as mentioned above, and normalized Taylor diagrams were generated for each of the  
449 classifications. To facilitate labeling in the diagrams, each of the SST products was given an abbreviated  
450 name consisting of the first two letters of the product's name (with the exception of WindSat,  
451 abbreviated WS, and K10 and LAC which stayed the same). The standard deviation of the UpTempO  
452 SSTs is represented in the normalized diagrams by the purple dashed arc depicting  $\hat{\sigma}_{obs} = 1$  and  
453 corresponding purple dot, labeled "observed," at unit distance from the origin along the x-axis.  
454 Additional isolines for  $\sigma$  are drawn at 1°C-intervals as continuous arcs in black. Isolines for the  $\widehat{RMS}'$  are  
455 depicted every 0.25°C, as concentric circles in green, centered on the observations; radial lines in blue,  
456 labeled according to the cosine of the angle made with the abscissa, correspond to 0.1 increments in  
457 correlation. Since  $\sigma_{obs}$ ,  $\sigma_{sat}$  and  $RMS'$  cannot be retrieved from normalized Taylor diagrams, each  
458 diagram has an associated table that includes the sample standard deviations of the collocated SST  
459 products and corresponding observations, as well as the mean bias, the standard deviation (STDEV) and

460 the RMSE of the paired satellite – buoy differences. Numerical values of the skill scores (Equations 5 and  
461 9) for each of the SST products are also included in these tables.

462

## 463 **5. Results**

### 464 5.1 Extreme Weather Event – Strong Gradients: Louis 2012-03

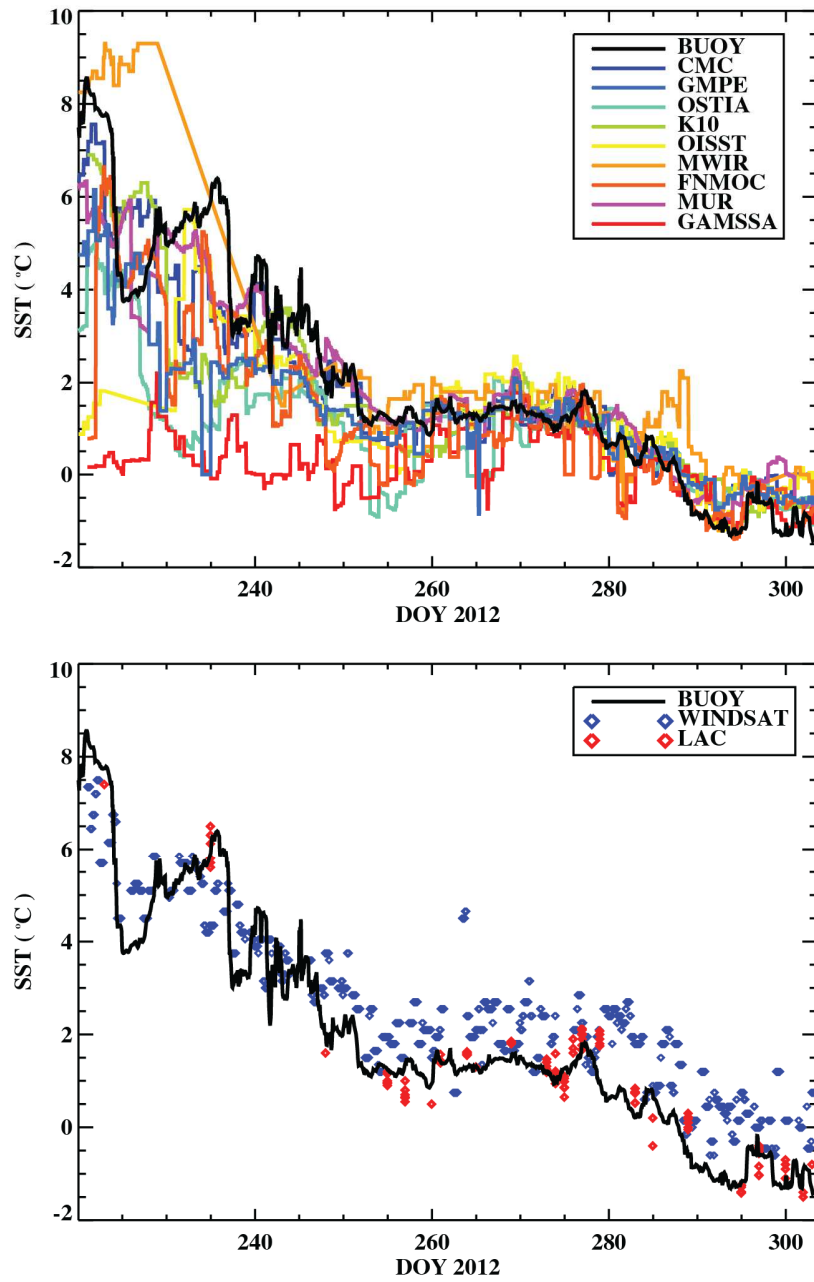
465 The first buoy deployed during the MIZOPEX summers, Louis 2012-03, experienced unique  
466 conditions. A temporal SST animation (not shown) indicates that it drifted along a temperature  
467 gradient, not far from the coast, for several weeks following its deployment. The spiral at the end of the  
468 track in Figure 2, also indicates that the buoy became trapped in a strong anti-cyclonic circulation that  
469 eventually coincided with its demise. The high spatial variability near the thermal gradient is  
470 problematic for satellite SST product comparisons, since buoys make point measurements throughout  
471 the day, whereas satellite SST estimates are representative of spatial averages over much larger areas.  
472 In the case of the L4, the estimates are expected to be smoother and coarser than non-analyzed SSTs.  
473 The grid resolutions specified in Table 1 define the smallest possible SST features that can be resolved by  
474 the different L4 products. Furthermore, oceanic mesoscale fronts promote convection, and hence are  
475 associated with clouds, which in turn, hamper the IR SST retrievals. Martin et al. (2012) have shown that  
476 differences among L4 products tend to be accentuated near coastal and strong gradient regions.

477 To further complicate matters, the Louis 2012-03 deployment coincided with the appearance of  
478 a rare and very powerful summer storm over the Arctic Ocean. The storm, dubbed the Great Arctic  
479 Cyclone of 2012 (GAC-2012, Simmonds and Rudeva, 2012), was first identified over northern Siberia on  
480 2 August 2012. It then crossed into the Arctic basin, intensified off the coast of Alaska on 6 August 2012,  
481 and then tracked into the center of the Arctic basin before slowly dissipating over the next several days  
482 (14 August 2012, day of year 226). Louis 2012-03 was deployed on August 7 off the coast of Alaska  
483 (72.6N, 144.64W) at a time when the storm reached its greatest size and depth. As of summer 2016,

484 this was the most extreme summer storm on record since satellite observations of polar orbiters began  
485 in 1979, and it is believed that the churning action of the cyclone contributed significantly to the rapid  
486 ice melt observed during August 2012 (Simmonds and Rudeva, 2012; Zhang et al., 2013).

487         The thermal structure of the ocean surface sampled by Louis 2012-03 (Figure 3) is characterized  
488 by a rapid cooling of the ocean surface from  $\sim 8.5^{\circ}\text{C}$  to  $1^{\circ}\text{C}$  during the first month of the deployment,  
489 followed by a more gradual cooling for the remainder of its life. Note that the drop in SST during the  
490 storm is about  $4.5^{\circ}\text{C}$ , possibly due to enhanced mixing associated with the storm-generated winds. Time  
491 series of the satellite SST products matched to Louis 2012-01 measurements (**Figure 4**) show extreme  
492 differences among the L4 products for the early part of the deployment (Figure 4a), but after about day  
493 of year (DOY) 260 the various time series quickly converge, indicating renewed agreement among the L4  
494 products. Corresponding comparisons with the L3 products (Figure 4b) show good overall agreement  
495 with Louis 2012-03, although for temperatures below  $2\text{-}3^{\circ}\text{C}$ , the WindSat SSTs appear to have a warm  
496 bias relative to the buoy. Differences in L4 products are further emphasized when comparing their  
497 corresponding images for a randomly chosen day (DOY 226, **Figure 5**) within the period with the  
498 greatest discrepancies. Clearly, significant differences exist for the Beaufort Sea in terms of SST  
499 amplitude, variability, and ice mask coverage, suggesting that some of these products are ill-equipped to  
500 deal with the harsh conditions encountered by Louis 2012-03 during the first month of its deployment.  
501 This clearly motivates determining which, if any, of the analyses accurately represent the actual  
502 conditions.

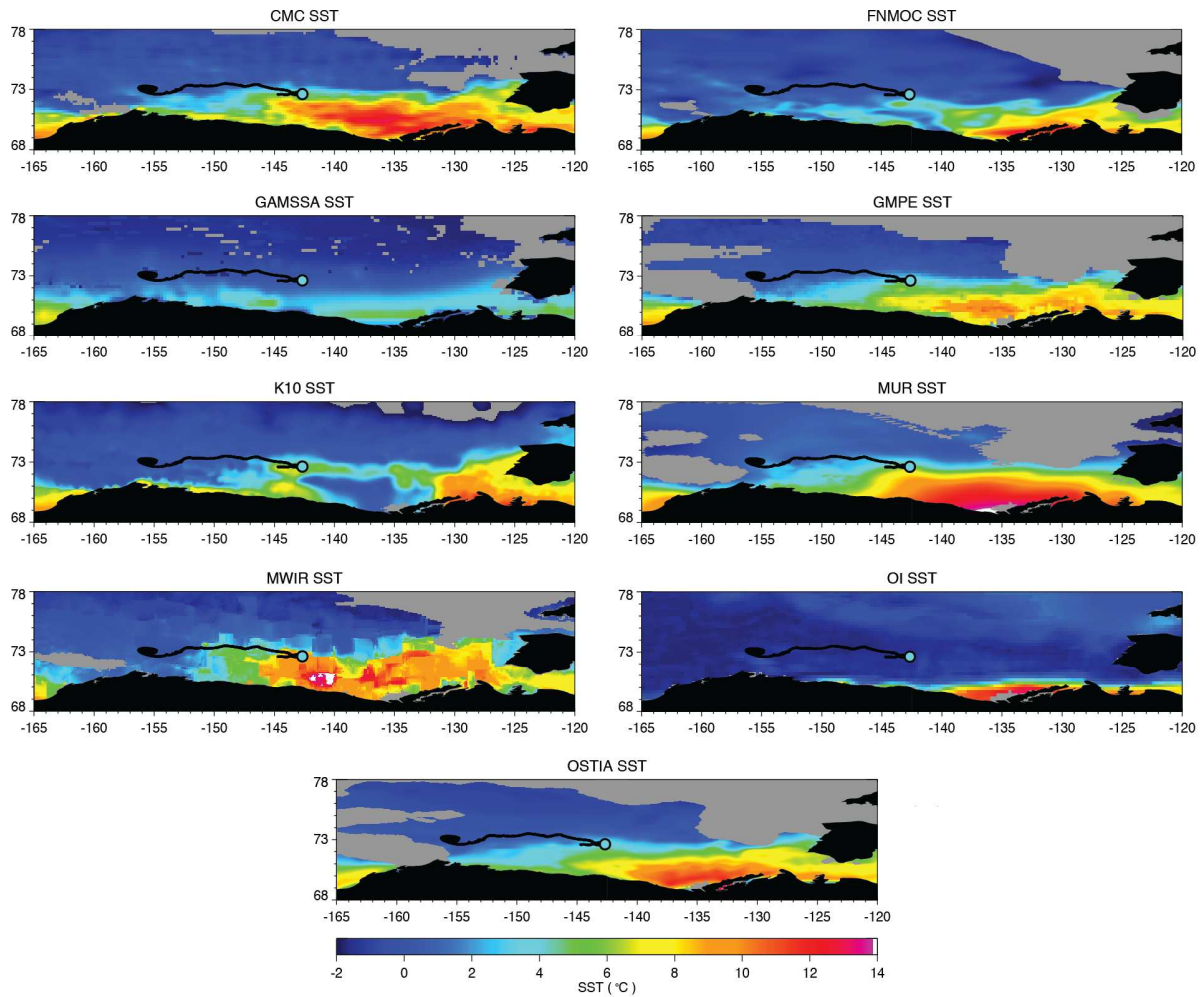




503

504 Figure 4. Time series of the SST products compared to the observations from the Louis 2012-03 buoy.  
 505 The upper panel compares the various L4 products while the lower panel displays the available single  
 506 sensor retrievals.

507



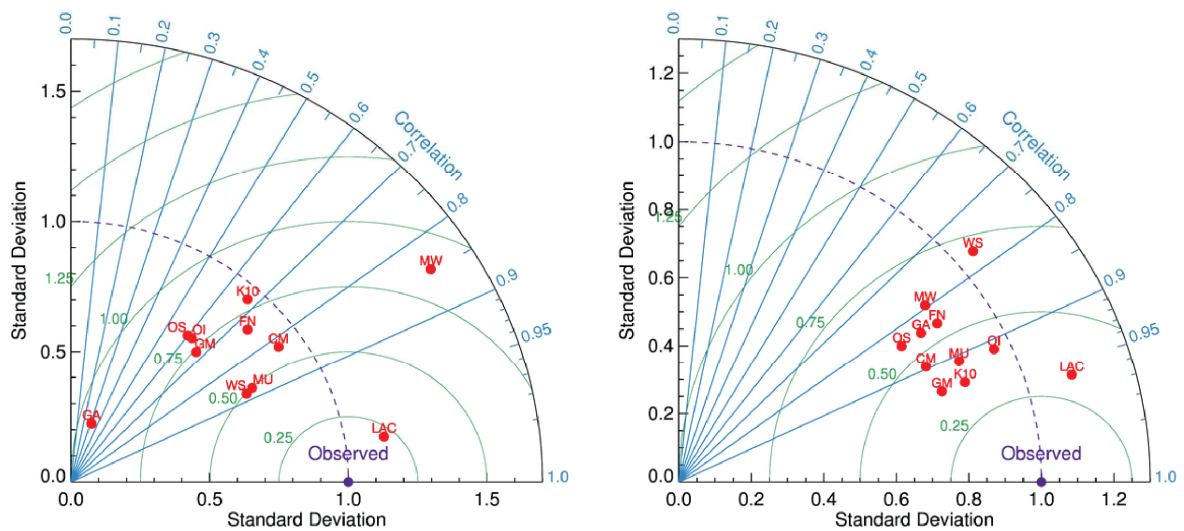
508

509 Figure 5. Comparison of the selected L4 SST analyses in the Beaufort Sea on 13 August, 2012 (DOY 226).  
 510 The color scale has been fixed to facilitate comparisons between products, and the trajectory of Louis  
 511 2012-03 has been plotted over the images with the circle showing the position of the buoy for that  
 512 particular day and the color indicating the corresponding buoy temperature. The gray areas indicate  
 513 that the respective ice mask has been applied, if available. Note that no ice mask is shown in the OISST  
 514 since the ice and water masks were mistakenly inverted during this period and the buoy location was  
 515 inaccurately flagged as ice covered. The anomalously low OISST temperature at the buoy location is a  
 516 result of the improper masking.

517 Because of the two different thermal regimes before and after DOY 260, and because a climate  
 518 event of the proportion of GAC-2012 might have generated oceanic variability that could have  
 519 prevented the analyzed products from agreeing with the buoy observations, we subsampled the Louis  
 520 2012-03 measurements into two datasets for the period prior/post DOY 260. The normalized Taylor  
 521 diagrams showing the performance of the satellite SST products relative to the Louis 2012-03

522 measurements at 2.5-m depth for these two periods are presented in **Figure 6** and associated  
 523 dimensional statistics are summarized in **Tables 3 and 4**, respectively. The  $\sigma_{obs}$  used in the  
 524 normalization of Figure 6 are given in the Tables. These diagrams do suggest two very different regimes,  
 525 as reflected by the re-arrangements of the dots representing the satellite products, with a wide scatter  
 526 in the radial direction for the period overlapping the storm (Figure 6a), and a much closer clustering  
 527 the remaining of the melting season (Figure 6b).

528



529  
 530 Figure 6. Normalized Taylor diagrams showing differences between matched SST from the UpTempO  
 531 buoy Louis 2012-03 2.5 m thermistor and eleven satellite SST products considering matchups before  
 532 DOY 260 (left) and after DOY 260 (right).

533

534 Table 3. Louis 2012-03 statistics for the period before DOY 260. All quantities are as defined in the text.

SST Product	No. Pts	$\sigma_{\text{obs}}$ (K)	$\sigma_{\text{sat}}$ (K)	Bias (K)	STDEV (K)	RMS (K)	$\rho$	SS	TS
CMC	935	1.98	1.80	-0.57	1.14	1.28	0.82	0.58	0.71
FNMOG	934	2.03	1.76	-1.58	1.40	2.11	0.74	-0.08	0.56
GAMSSA	934	2.03	0.48	-3.64	1.93	4.12	0.32	-3.12	0.00
GMPE	935	1.98	1.33	-1.51	1.47	2.11	0.67	-0.13	0.29
K10	934	2.03	1.92	-1.10	1.60	1.94	0.67	0.08	0.51
MUR	935	1.98	1.48	-0.41	0.99	1.07	0.88	0.71	0.59
MWIR	479	2.35	3.60	1.34	2.04	2.44	0.85	-0.08	0.40
OISST	767	2.01	1.42	-1.40	1.59	2.12	0.62	-0.11	0.29
OSTIA	935	1.98	1.39	-2.12	1.60	2.65	0.60	-0.80	0.27
LAC	92	2.02	2.31	-0.30	0.44	0.53	0.99	0.93	0.97
WindSat	715	2.14	1.53	0.06	1.06	1.07	0.88	0.75	0.55

535

536 Table 4. Louis 2012-03 statistics for the period after DOY 260. All quantities are as defined in the text.

SST Product	No. Pts	$\sigma_{\text{obs}}$ (K)	$\sigma_{\text{sat}}$ (K)	Bias (K)	STDEV (K)	RMS (K)	$\rho$	SS	TS
CMC	1045	1.05	0.80	0.22	0.49	0.53	0.90	0.74	0.65
FNMOG	1042	1.02	0.87	-0.18	0.56	0.59	0.83	0.67	0.68
GAMSSA	1045	1.02	0.82	-0.26	0.57	0.62	0.83	0.63	0.62
GMPE	1045	1.05	0.81	0.21	0.40	0.45	0.94	0.82	0.73
K10	1045	1.02	0.86	0.12	0.37	0.39	0.94	0.86	0.84
MUR	1045	1.05	0.89	0.28	0.44	0.52	0.91	0.75	0.79
MWIR	949	1.02	0.87	0.56	0.63	0.84	0.79	0.33	0.62
OISST	1042	1.05	1.00	0.37	0.43	0.57	0.91	0.71	0.88
OSTIA	1045	1.05	0.77	0.07	0.58	0.59	0.84	0.69	0.53
LAC	375	0.99	1.12	0.09	0.32	0.33	0.96	0.89	0.93
WindSat	850	0.99	1.05	0.98	0.70	1.21	0.77	-0.47	0.64

537

538 During the first period (Figure 6a), the satellite SST product in closest agreement with the  
539 observations by far is the LAC (L3) SSTs with  $\widehat{RMS}'_{LAC} = 0.22$ ,  $\hat{\sigma}_{LAC} = 1.14$ , and  $\rho_{LAC} = 0.99$ . The second  
540 and third closest in terms of low RMSE and high correlation are the L3 WindSat and the L4 MUR. Both of  
541 these products lie on the  $0.5 \widehat{RMS}'$ -arc and the  $0.88 \rho$ -radius, but are farther from the observed  
542 variance ( $\hat{\sigma}_{WS} = 0.74$  and  $\hat{\sigma}_{MUR} = 0.78$ ). It is expected that, when available, the lower-level  
543 processing products should outperform the L4 products, particularly in the presence of strong gradients,  
544 as the OI analysis system is in itself a spatial and temporal smoothing filter, damping some of the natural

545 SST variability. The CMC is closer to the arc for  $\hat{\sigma}_{obs} = 1$  ( $\hat{\sigma}_{CMC} = 0.91$ ), but has  $\widehat{RMS}' > 0.5$ . Although  
546 MUR has better overall agreement with the buoy than CMC (smaller RMSE), the fact that  $\hat{\sigma}_{CMC}$  is closer  
547 to the  $\hat{\sigma}_{obs}$ -arc, suggests that, at least in terms of the analyzed SST amplitude variability matching the  
548 observed, the CMC does better than MUR. Which of these two products is appraised over the other  
549 depends on the features valued by the different scoring systems. Of the remaining products, K10 and  
550 FNMOC lie closer to the arc for  $\hat{\sigma}_{obs} = 1$ , whereas GMPE, OSTIA, and the OI align with the arc for  $\hat{\sigma}_{obs} =$   
551  $0.7$ . The positioning of the OI, in particular, follows after careful screening of the data for the period  
552 corresponding to August 10 – August 16, 2012, when the ice and water masks for the Arctic region were  
553 inverted in the OISST product. Failing to remove data from this period results in further degradation of  
554 the OISST statistics. The products whose  $\hat{\sigma}$  are farthest apart from  $\hat{\sigma}_{obs} = 1$  are GAMSSA and the  
555 MWIR. These products display a wide range of SST amplitudes, with GAMSSA being much smoother  
556 ( $\hat{\sigma}_{GAMSSA} = 0.24$ ) and the MWIR being much noisier ( $\hat{\sigma}_{MWIR} = 1.53$ ) than the other L4s. Note also that  
557 their time series exhibit the largest departures from the buoy during the early portion of Figure 4. Since  
558 the spread in the radial direction gives an indication of the degree to which temporal and spatial SST  
559 variability is affecting the SST amplitudes, the uncertainty attributable to variability appears to be  
560 significant for these two products. Our result supports the findings of Reynolds and Chelton (2010)  
561 which found that, since the MWIR attempts to resolve very small SST features based on the 1-km MODIS  
562 (MODerate Resolution Imaging Spectroradiometer) SSTs, when the high-resolution IR data is missing or  
563 the coverage is reduced due to persistent cloud cover over multiple days, as was likely the case during  
564 the storm/proximity to a thermal front, insufficient high resolution IR data results in small-scale noise. It  
565 is noteworthy that the MWIR is highly correlated with the observations ( $\rho_{MWIR} = 0.85$ ) despite having  
566 too much variability.

567           For this period along the strong SST gradient, the skill scores for the L4 comparisons are  
568 generally quite low, confirming once again how ill-equipped the L4s are for extreme conditions. Sorting

569 the  $TS$  scores (Table 3) in descending order of skill (decreasing score magnitude), results in:  $TS_{LAC} \gg$   
570  $TS_{CMC} > TS_{MUR} > TS_{FNMOC} > TS_{WS} > TS_{K10} \gg TS_{MWIR} > TS_{GMPE} > TS_{OI} > TS_{OSTIA} >$   
571  $TS_{GAMSSA}$ . Similarly, the sorted list for the  $SS$  scores (Table 3) indicates:  $SS_{LAC} > SS_{WS} \gg SS_{MUR} >$   
572  $SS_{CMC} \gg SS_{K10} > SS_{FNMOC} > SS_{MWIR} > SS_{OI} > SS_{GMPE} > SS_{OSTIA} > SS_{GAMSSA}$ . The symbols  $\gg$  and  
573  $\gg$  indicate the relative positions of the fixed thresholds for  $T_{75}$  and  $T_{51}$ , delimiting the discrete skill  
574 categories of excellent, good or poor product performance. Only the L3 products can be classified as  
575 excellent during this challenging period. The high-resolution LAC data, when available, clearly can  
576 capture the gradient, but the blending and reduced resolution of the analyses tends to miss or  
577 smear the gradient. MUR, having the highest resolution of the L4 products considered, does among  
578 the best in this case. It can be seen that CMC and MUR, as well as OISST and GMPE, trade positions in  
579 the two scoring systems, since the  $TS$  rewards products with  $\hat{\sigma}$  closer to  $\hat{\sigma}_{obs}$ , whereas  $SS$  rewards  
580 products with small  $\widehat{RMS'}$  and small bias. The dimensional statistics in Table 3 indicate significant  
581 absolute biases ( $> 1.0^\circ\text{C}$ ) with respect to the LOUIS 2012-03 for all satellite products (also evident in  
582 Figure 3), except WindSat ( $0.06^\circ\text{C}$ ), LAC ( $-0.30^\circ\text{C}$ ), MUR ( $-0.41^\circ\text{C}$ ) and CMC ( $-0.57^\circ\text{C}$ ). The K10 and the  
583 FNMOC, which fared well in terms of the  $TS$  classifier, are severely penalized due to their large biases,  
584 and downgraded to poor skill in the  $SS$  scale.

585 During this portion of the Louis 2012-03 deployment it is difficult to assert whether differences  
586 in skill were solely sampling errors associated with the proximity to the strong temperature gradients or  
587 also contained retrieval errors, perhaps associated with the extreme storm conditions. Since Louis  
588 2012-03 followed a front for much of its trajectory, it is plausible to ascribe a fraction of the high RMSE  
589 values in Table 3 to sampling variability. In and around the front, the RMS would have been particularly  
590 sensitive to variations in the analyzed SST amplitudes. Note that in this comparison, we did not attempt  
591 to interpolate the different satellite products to a common grid; instead we opted for working with the  
592 products in their native spatial resolution. It is also possible that the IR SST retrievals were limited by the

593 presence of clouds associated with the storm and strong SST fronts, as suggested by the small number  
 594 of AVHRR LAC matchups obtained before DOY 260 (Figure 4b). The WindSat – buoy matchups, on the  
 595 other hand, were not affected by clouds, giving an apparent advantage to CMC and MUR (two of just  
 596 three L4s that ingested WindSat SSTs at this time; see Table 2). The MWIR, despite ingesting WindSat,  
 597 did not generate many SST retrievals for this period, likely due to the fact that it uses the MODIS SST  
 598 cloud mask in the Arctic, and being an IR sensor, the MODIS instrument has limited coverage under  
 599 cloudy conditions.

600 For the quiescent period following DOY 260, the results change dramatically and the L4 satellite  
 601 products generally show better agreement with the buoy. The L3 SST products in the Taylor diagram  
 602 (Figure 6b) align in an arc to the right of  $\hat{\sigma}_{obs} = 1$  at a radial distance of  $\hat{\sigma}_{L3} = 1.13$ , whereas the L4  
 603 products cluster to the left, at a radial distance of  $\sim 0.80$ , indicating an overall agreement in L4 product  
 604 performance for the second period. Corresponding statistics (Table 4) also indicate that the sampling  
 605 biases, although still present, are significantly smaller ( $< 0.5^\circ\text{C}$ ) for all but the MWIR and WindSat  
 606 products.

607 The convergence in performance among L4 product is substantiated by the *TS* scoring system,  
 608 which finds no poorly skilled products for the latter half of the Louis 2012-03 deployment. Sorted *TS*  
 609 scores in descending order of skill (Table 4), indicate that:  $TS_{LAC} > TS_{OI} > TS_{K10} > TS_{MUR} \gg$   
 610  $TS_{GMPE} > TS_{FNMOC} > TS_{CMC} > TS_{WS} > TS_{MWIR} > TS_{GAMSSA} > TS_{OSTIA}$ , whereas the sorted *SS*  
 611 scores result in:  $SS_{LAC} > SS_{K10} > SS_{GMPE} > SS_{MUR} \gg SS_{CMC} > SS_{OI} > SS_{OSTIA} > SS_{FNMOC} >$   
 612  $SS_{GAMSSA} \gg SS_{MWIR} > SS_{WS}$ . Under quiescent conditions, the satellite products that excel at  
 613 reproducing the observations in the *TS* classifier are LAC, OISST, K10, and MUR. Improvements are  
 614 especially remarkable for the OISST with the closest variance to the observations ( $\hat{\sigma}_{OISST} = 0.95$ ).  
 615 However, despite achieving the correct SST variability and being ranked second best, the OISST is  
 616 demoted one category by the *SS* classifier because of its larger bias ( $0.37^\circ\text{C}$ ) relative to its counterparts

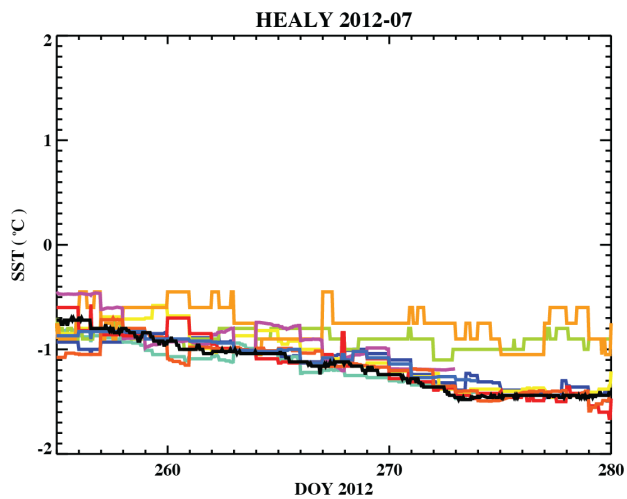
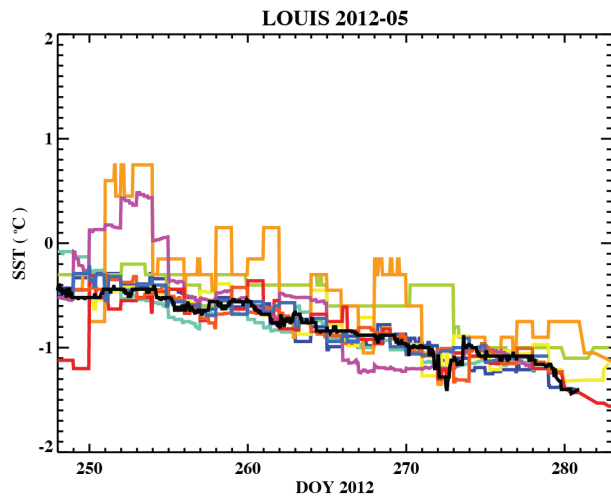
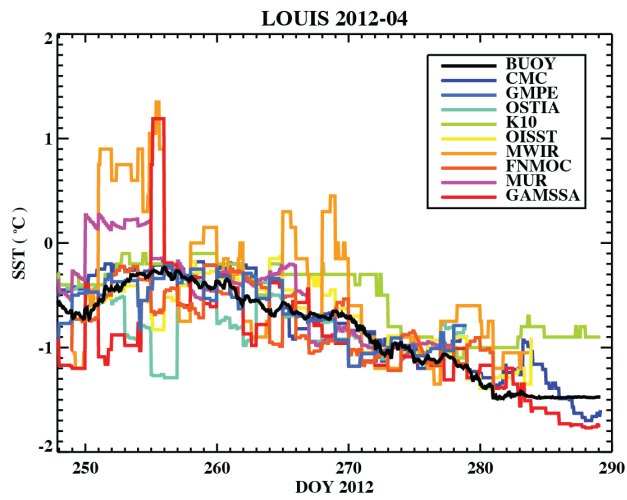
617 at the top of the scale. The GMPE with its small  $\widehat{RMS}'$  and bias, is promoted by the SS classifier in lieu of  
618 the OISST, and placed between the K10 and the MUR. Given the homogeneity in L4 performance under  
619 these quiescent conditions, it is no surprise to find that the median ensemble is among the more skillful  
620 products.

621 It is noteworthy that WindSat drops from second best in the first period according to the SS  
622 scoring to last place in the second period. The good agreement between the MW SSTs from WindSat  
623 and the Louis 2012-03 surface temperatures during the first half (Figure 6a and TS in Table 3) suggests  
624 that, under warmer temperatures, the additional coverage enabled by MW SSTs in the presence of  
625 clouds that obscure the IR is highly beneficial in the objective analyses. After DOY 260, however, Figure  
626 4b shows the appearance of a warm bias, also evident in Table 4, with the L3 WindSat SSTs  
627 overestimating the buoy temperatures by  $\sim 1^\circ\text{C}$ . This contrasts with the fact that Table 3 indicates zero-  
628 bias for WindSat during the first period. This  $1^\circ\text{C}$ -bias persists in all other UpTempO buoy combinations  
629 explored hereafter where the prevailing SSTs were below  $2^\circ\text{C}$ . The WindSat bias is potentially related to  
630 cold SSTs, as a high-latitude bias has also been reported for MW AMSRE SSTs relative to ship-based  
631 observations in the Southern Ocean (Dong et al., 2006), although of a lesser magnitude  
632 (ascending/descending: 0.42/23 K in the summer and -0.21/-0.42 K in the winter).

## 633 5.2. Cold Northerly Waters

634 Louis 2012-04, Louis 2012-05, and Healy 2012-07 were deployed at the beginning of September  
635 2012 in the cold waters further north from the Alaska coast. As illustrated in Figure 2, these buoys  
636 drifted counter-clock wise toward the interior of the Canada Basin. The flow pattern suggests that the  
637 circulation of the Beaufort Gyre had reversed as a result of the low pressure system that persisted for  
638 several days in August 2012. The buoy temperatures are characterized by a narrow dynamic range with  
639 initial temperatures below  $0^\circ\text{C}$ , followed by a gradual cooling (at nearly identical cooling rates) until the  
640 start of fall freeze-up (**Figure 7**).

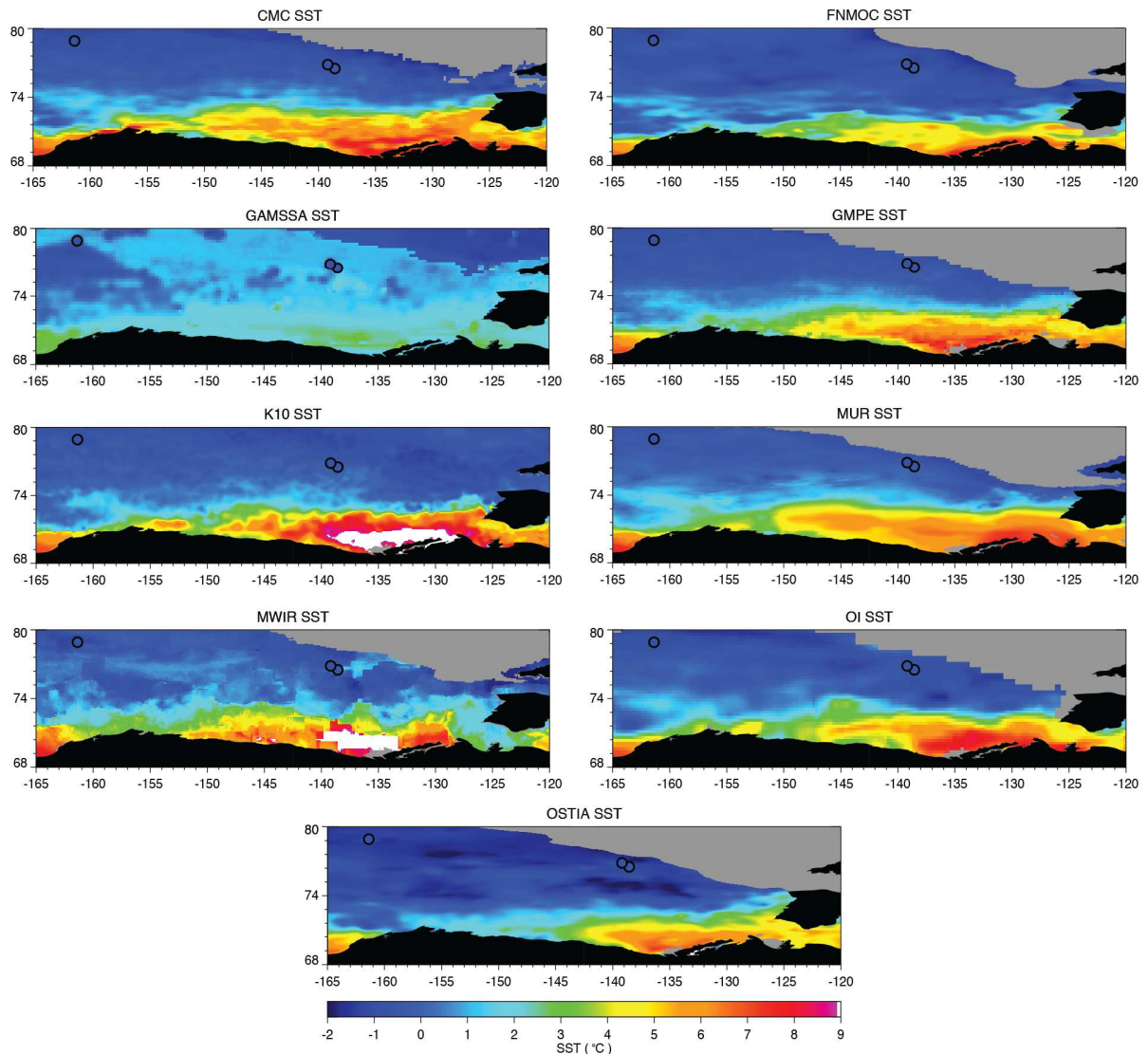




641  
 642 Figure 7. Time series of the selected SST products compared with the remaining UpTempO buoys  
 643 deployed in 2012. These are considered the cold northerly buoys. Separate panels are shown for each  
 644 buoy while the color traces correspond to the different SST products. The buoy observations are always  
 645 shown with the black trace.

646

647           The location of these three buoys (Figure 2) is especially interesting because they allow us to  
648 look at the impact that the different ice masks is having on the L4 SST products. Being farther north,  
649 these buoys were closer to the main ice pack and experienced refreeze earlier than the buoys closer to  
650 the coast. Maps for all the products for a randomly chosen day within this period (DOY 256, 2012) are  
651 shown in **Figure 8**. Buoy positions on this particular day are displayed in Figure 8 as circles, color-coded  
652 by the buoy measured temperature. Although the maps again show pronounced differences among the  
653 satellite products, the buoys are located far off the region with the largest discrepancies, which happens  
654 to be near the coast. The GAMSSA product appears to have experienced some difficulties with the ice  
655 mask around this time, as indicated by the corresponding map in Figure 8, where the region that should  
656 have been flagged as ice is, in fact, shown as water. This issue was corrected after DOY 258.



657

658 Figure 8. Graphical comparison of the selected L4 analyses on DOY 256 of 2012 corresponding to the  
 659 period sampled by the cold northern buoys. The separate buoy positions are indicated with the  
 660 enclosed circles with the color corresponding to the buoy temperature. The buoys, in order from west  
 661 to east, are Healy 2012-07, Louis 2012-05, and Louis 2012-04. The ice masks are again indicated by the  
 662 gray regions. The white areas in the K10 and MWIR analyses correspond to where the temperature  
 663 exceeds the maximum value on the color scale.

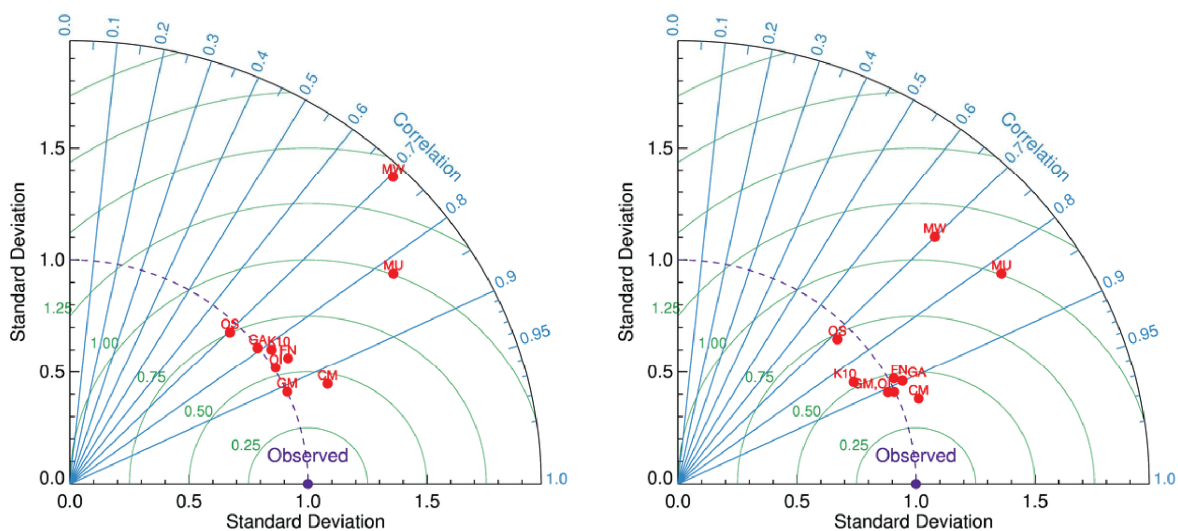
664 In order to explore the impact that a less conservative ice mask (one that delays freezing) had  
 665 on the analyses, we considered two Taylor diagrams: one based on the whole extent of each of the  
 666 products' matchup time series, and a second one in which the most conservative ice mask is used to  
 667 truncate the time series to a common period during which all the products were unambiguously

668 reporting temperatures above the freezing temperature, i.e., a period during which we were highly  
669 confident that the SSTs were not being influenced by ice. A byproduct of the truncation is that, whereas  
670 the former diagram deals with time series of unequal lengths, the latter one is more balanced in terms  
671 of the number of counts. The L3 products are excluded from these comparisons, since the LAC coverage  
672 did not extend that far north, and the WindSat SSTs appeared to be decorrelated from the buoys (not  
673 shown). This could have been the result of cold absolute temperatures and/or the buoys being within  
674 75-km from the ice, where MW SSTs cannot be retrieved.

675 The OSI-SAF ice mask is the most conservative of the ice masks used in the L4s considered here  
676 (see Table 1), with a much earlier freeze cut-off at the end of the melt season. Of the products using the  
677 EUMETSAT OSI-SAF ice mask, the MUR SST matchups defaulted to the freezing temperature on days  
678 278, 279, and 273 for Louis 2012-04, -05, and Healy 2012-07, respectively. OSTIA and GMPE, also using  
679 the OSI-SAF ice mask, followed suit one day later. The other L4s, continued to report SSTs for about 10  
680 more days. Freezing-up last were CMC, K10 and GAMSSA. The ice flagging in the UpTempO buoys  
681 agreed extremely well with the OSI-SAF ice mask indicating possible ice effects on days 279, 280, and  
682 274 for Louis 2012-04, -05, and Healy 2012-07, respectively. It is important to emphasize that while the  
683 UpTempO ice indicator is based on the National Snow and Ice Data Center (NSIDC) SIC, it uses a lower  
684 ice concentration ( $SIC \geq 0.15$ ) and a  $-1.2^{\circ}\text{C}$  threshold on the uppermost thermistor to indicate that the  
685 buoy is in/near ice. Thus, the extent of the time series of the MUR–UpTempO matchups determined the  
686 truncation dates for all the other products, as this period was considered ice-free in both, the satellite  
687 and the UpTempO records.

688 The normalized Taylor diagram for the period free of ice effects is shown in **Figure 9a** and  
689 corresponding statistics are shown in **Table 5**. A visual inspection of this diagram indicates that the  
690 majority of the products tightly align with the arc for  $\hat{\sigma}_{obs} = 1$ , with a spread in azimuthal direction  
691 characterized by an increase in  $\widehat{RMS}'$  (decrease in  $\rho$ ) from 0.42 (0.91) for GMPE to 0.75 (0.70) for OSTIA.

692 The two products that do not conform with the others are the MUR and the MWIR. Their location to the  
 693 right of the  $\hat{\sigma}_{obs}$ -arc, indicates substantially larger  $\hat{\sigma}$  and  $\widehat{RMS}'$  than the rest. For this region with narrow  
 694 dynamic SST range, the scores in Table 5 had a wider range than the scores obtained for the latter part  
 695 of Louis 2012-03 (Table 4), but several products demonstrated good skill. The TS ranking of the L4  
 696 products in decreasing order of skill,  $TS_{GMPE} > TS_{CMC} > TS_{OISST} > TS_{FNMOC} \gg TS_{K10} >$   
 697  $TS_{GAMSSA} > TS_{OSTIA} \gg TS_{MUR} > TS_{MWIR}$ , agrees with the SS ranking except for the placement of the  
 698 K10, which moves from fifth to seventh place, and the threshold delimiters shifting two positions to the  
 699 left, leaving GMPE and CMC at the top of the ranking, followed by OISST, FNMOC, and GAMSSA in the  
 700 intermediate category, and OSTIA, K10, MUR and MWIR at the bottom. The positioning of the GAMSSA  
 701 product follows after careful removal of data for DOY 255, for which this analysis reported unrealistically  
 702 warm SSTs relative to Louis 2012-04 (see Figure 7a). Among the products showing poor skill under the  
 703 sampled conditions, MUR and the MWIR differ from the others mainly in that they are the only ones  
 704 that assimilate MODIS SSTs in the analysis. We speculate that biases due to residual cloud  
 705 contamination at the highest quality MODIS L2 SSTs might have negatively impacted these products.



706  
 707 Figure 9. Normalized Taylor diagrams for the combination of the cold northerly buoys including Louis  
 708 2012-04, Louis 2012-05 and Healy 2012-07. The left panel shows the results for the common truncated

709 time series when the OSI-SAF ice mask indicated the region was ice-free while the right panel shows the  
 710 results for the whole extent of each individual product time series.

711 Table 5. Statistics for satellite SSTs matched to temperatures from the 2012 northerly buoys after  
 712 truncating irregular intervals at the start/end of the summer period to eliminate possible ice effects.

L4 Product	No. Pts	$\sigma_{obs}$ (K)	$\sigma_{sat}$ (K)	Bias (K)	STDEV (K)	RMS (K)	$\rho$	SS	TS
CMC	1761	0.29	0.33	0.04	0.13	0.14	0.92	0.77	0.83
FNMOG	1714	0.27	0.30	-0.04	0.16	0.16	0.85	0.65	0.77
GAMSSA	1736	0.28	0.28	-0.04	0.18	0.19	0.79	0.57	0.68
GMPE	1705	0.28	0.28	0.02	0.12	0.12	0.91	0.82	0.89
K10	1761	0.29	0.30	0.24	0.18	0.30	0.86	-0.08	0.72
MUR	1680	0.27	0.45	0.14	0.28	0.31	0.82	-0.27	0.31
MWIR	1761	0.29	0.55	0.35	0.40	0.53	0.70	-2.49	0.14
OISST	1756	0.29	0.29	0.05	0.15	0.16	0.86	0.68	0.79
OSTIA	1660	0.28	0.27	-0.10	0.21	0.23	0.70	0.31	0.56

713

714 The normalized Taylor diagram using the whole extent of the time series (**Figure 9b**) differs from  
 715 the truncated one (Figure 9.a) in that the products that previously aligned with the arc for  $\hat{\sigma}_{obs} = 1$  re-  
 716 emerge clustered around the arc for  $\widehat{RMS}' = 0.45$  and the radius for  $\rho = 0.90$ . The ranking of the L4s  
 717 based on the TS classifier (**Table 6**) is:  $TS_{CMC} > TS_{OISST} > TS_{GMPE} > TS_{GAMSSA} > TS_{FNMOG} \gg$   
 718  $TS_{OSTIA} > TS_{K10} \gg TS_{MUR} > TS_{MWIR}$ . Note that with the exception of GAMSSA, which benefited  
 719 from more data to compensate for the issues it experienced at the beginning of these buoy surveys, the  
 720 skill categories based on the TS classifier are comprised of the same products in both experiments; i.e.,  
 721 the potential skill of the L4s did not change regardless of the presence of possible ice effects. The SS  
 722 ranking differs from the TS ranking in that K10 and OSTIA get demoted one category because of a large  
 723 bias and large  $\widehat{RMS}'$ , respectively. Consequently, there are no products in the intermediate category  
 724 under the SS classifier when the entire length of the time series is used. The K10 has a bias of 0.29°C  
 725 relative to the northern UpTempO buoys, the second largest after the MWIR. A bias of the same  
 726 magnitude (~0.3°C) has been reported by the authors (Castro et al., 2012) for SSTs < 8°C in validation  
 727 studies of the K10 SSTs using GTS drifting buoys. This warm bias might be a manifestation of the K10

728 reliance on long-term climatologies that show greater ice extent than what was actually observed in a  
 729 year with extreme thinning.

730 Table 6. Statistics for satellite SSTs matched to the temperatures from the 2012 northerly buoys for the  
 731 entire series of matches.

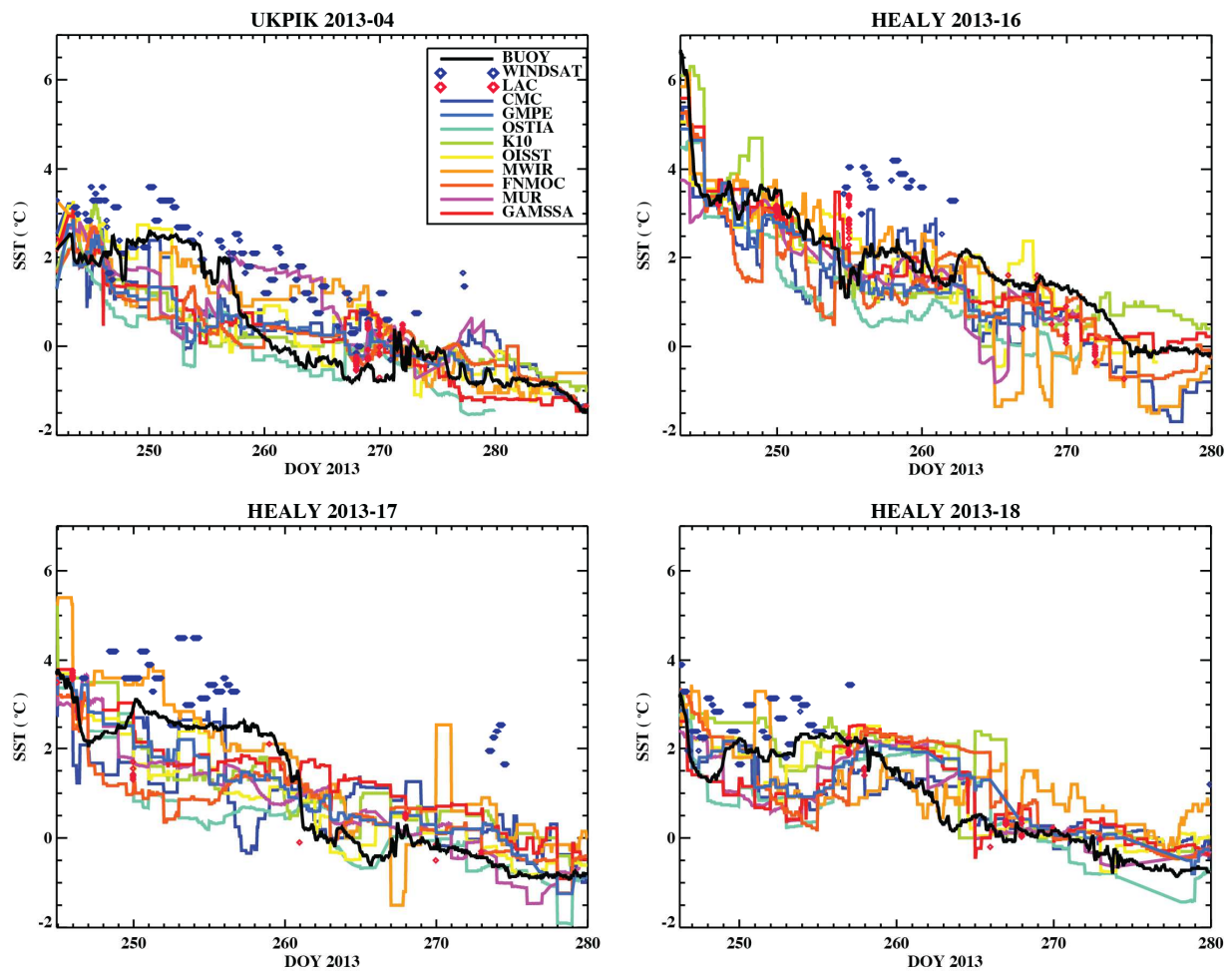
L4 Product	No. Pts	$\sigma_{\text{obs}}$ (K)	$\sigma_{\text{sat}}$ (K)	Bias (K)	STDEV (K)	RMS (K)	$\rho$	SS	TS
CMC	2232	0.36	0.39	0.05	0.14	0.15	0.94	0.84	0.91
FNMOG	2015	0.32	0.33	-0.04	0.15	0.16	0.89	0.76	0.84
GAMSSA	2404	0.37	0.38	-0.05	0.17	0.18	0.90	0.77	0.86
GMPE	1776	0.29	0.29	0.03	0.13	0.13	0.91	0.81	0.88
K10	2394	0.37	0.32	0.29	0.19	0.35	0.85	0.10	0.72
MUR	1680	0.27	0.45	0.14	0.28	0.31	0.82	-0.27	0.31
MWIR	2124	0.34	0.53	0.38	0.38	0.54	0.70	-1.42	0.28
OISST	2148	0.35	0.35	0.06	0.15	0.16	0.91	0.79	0.89
OSTIA	1731	0.29	0.27	-0.09	0.21	0.23	0.72	0.38	0.57

732  
 733 Looking at the small differences between the statistics in Tables 5 and 6, it becomes apparent  
 734 that one statistical measure alone cannot capture the impact that the ice mask is having on SST product  
 735 performance. In order to see a discernible effect, we need to look at the combined effect of all the  
 736 statistics at once. This is captured by scores, which show more drastic changes between the two tables.  
 737 If we ignore the products that rely on the OSI-SAF ice mask (i.e., MUR, OSTIA, GMPE, and MWIR) which  
 738 should have remained mostly unchanged between experiments (the MWIR is the exception since it  
 739 continued to report SSTs for 6 additional days) and look at the differences in actual performance (the SS  
 740 classifier) between Table 5 and Table 6, we find that the scores for all the L4s that use less conservative  
 741 ice masks improved significantly through the availability of more observations. In fact, by using longer  
 742 time series, the OISST, FNMOG, and GAMSSA, which had intermediate skills according to the SS scores in  
 743 Table 5, joined CMC as having excellent skills (Table 6). More data (note that the OISST and FNMOG  
 744 estimate proxy SSTs until SIC = 1) resulted in smaller  $\widehat{RMS'}$  values and better linear fits ( $\rho \rightarrow 1$ ), which  
 745 in turn produced higher scores that rewarded the products in the cluster closer to the observations.  
 746 Thus, product performance was not degraded by using less conservative ice masks in this specific case.

### 747 5.3. Coastal Buoys

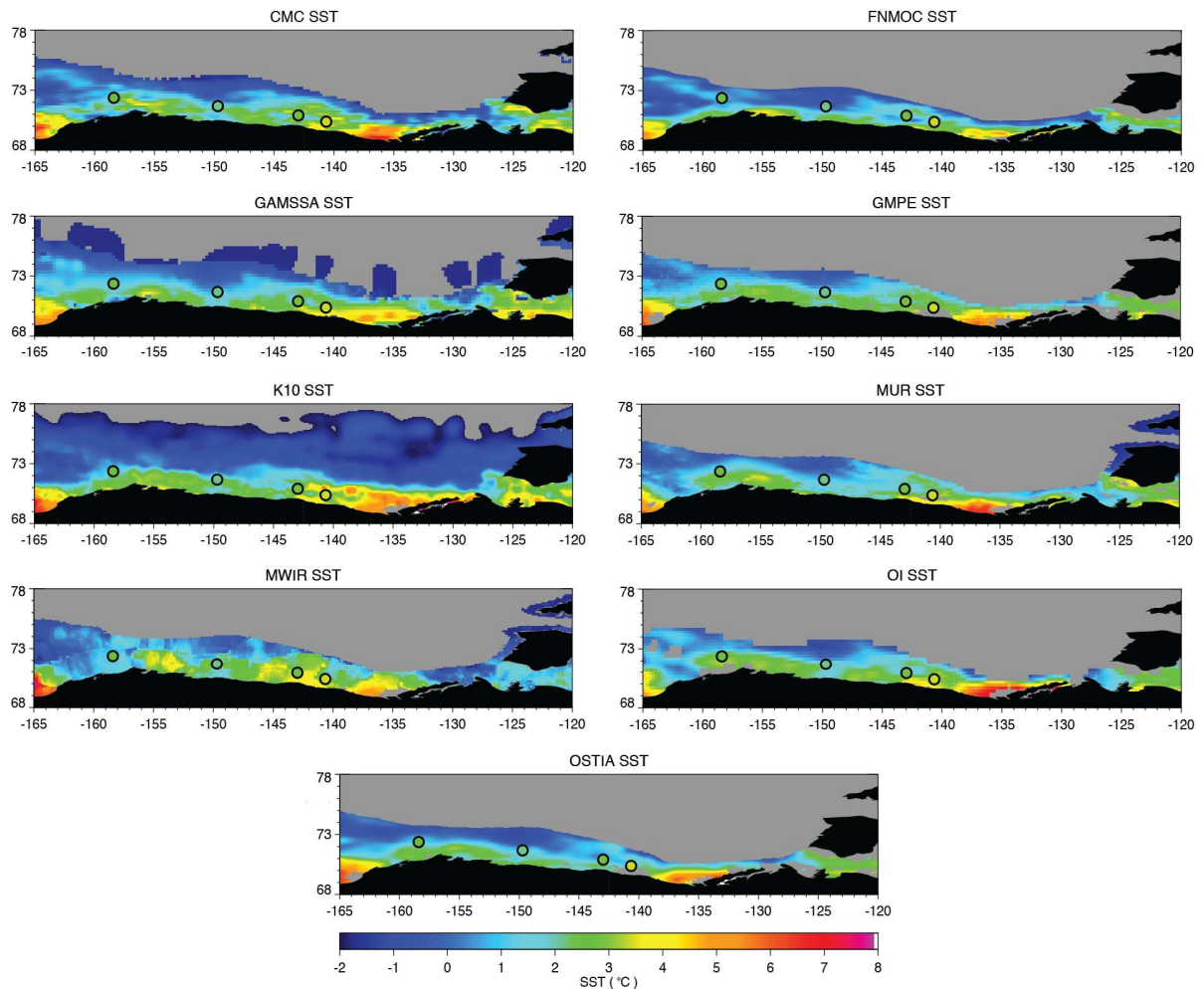
748           The Arctic summer of 2013 (August and September) saw less ice retreat than in the record year  
749 of 2012. The UpTempO buoys deployed in 2013 (Healy 2013-16, Healy 2013-17, Healy 2013-18, and  
750 Ukpik 2013-04) remained closer to the coast (see Figure 2) and moved westward toward the Chukchi  
751 Sea. The time series for the satellite SSTs matched to the observed temperatures from the 2013 buoys  
752 are shown in **Figure 10**. The SST dynamic range is from approximately 4 to -1°C which is rather small  
753 from a global perspective, but is larger than for the northern buoys. In fact, the observed temperatures  
754 here are warmer than for all other groupings other than the storm/gradient period from Louis 2012-03.  
755 A notable difference visible in the time series is that the satellite products tend to show a more  
756 constant, gradual cooling while the buoys suggest a more step-wise drop. Maps for a single, arbitrarily  
757 chosen day (**Figure 11**, DOY 250, 2013), suggest that once again, the buoys were drifting along thermal  
758 fronts, which might explain the choppiness in the buoy measurements, and the smoothing of the SST  
759 analyses along frontal boundaries.





760

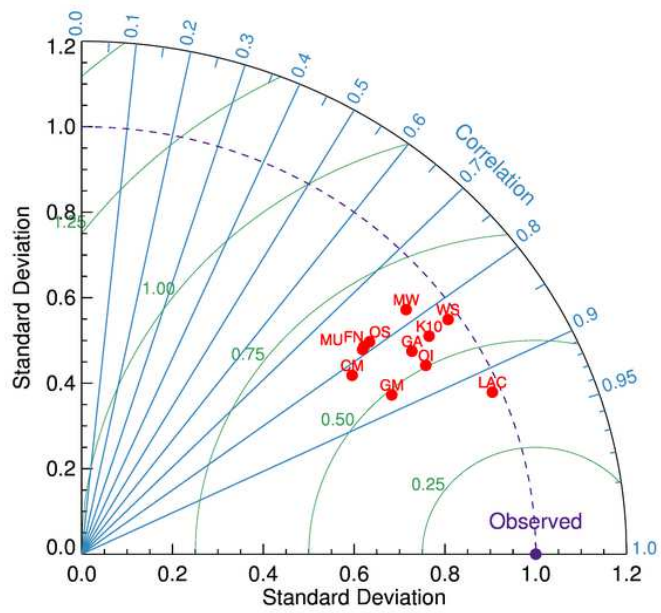
761 Figure 10. Time series of the selected SST products compared with the 2013 UpTempO buoys termed as  
 762 coastal buoys. Separate panels are shown for each buoy while the color traces correspond to the  
 763 different SST products. The buoy observations are always shown with the black trace.



764  
 765 Fig 11. Graphical comparison of the selected L4 analyses on DOY 250 of 2013 corresponding to the  
 766 period sampled by the coastal buoys. The separate buoy positions are indicated with the enclosed  
 767 circles with the color corresponding to the buoy temperature. The buoys, in order from west to east,  
 768 are Ukpik 2013-04, Healy 2013-18, Healy 2013-17, and Healy 2013-16. The ice masks are again indicated  
 769 by the gray regions.

770

771 The normalized Taylor diagram for the coastal buoys, including the latter portion of Louis 2012-  
 772 03, is shown in **Figure 12** and the corresponding statistics are given in **Table 7**. The diagram shows  
 773 similar skill for all SST products, with most products clustering in a narrow region of the parameter space  
 774 delimited by the  $\widehat{RMS}'$ -arcs for 0.5 and 0.65, the  $\hat{\sigma}$ -arcs for 0.8 and 1.0, and  $\rho$ -radii between 0.8 and 0.9.  
 775 The L3 SST products have daily SST amplitudes that are in excellent agreement with the observed (they  
 776 lie on the arc for  $\hat{\sigma}_{obs} = 1$ ), whereas the L4s have slightly smoother amplitudes, as expected.



777

778 Figure 12. Normalized Taylor diagram for the combination of the coastal buoys including all of the 2013  
 779 buoys plus the observations from Louis 2012-03 excluding the period of the storm.

780

781 Table 7. Statistics for satellite SSTs matched to temperatures from the 2012–2013 UpTempO coastal  
 782 buoys.

SST Product	No. Pts	$\sigma_{obs}$ (K)	$\sigma_{sat}$ (K)	Bias (K)	STDEV (K)	RMS (K)	$\rho$	SS	TS
CMC	4215	1.39	1.01	-0.03	0.81	0.81	0.82	0.66	0.50
FNMOC	4211	1.38	1.09	-0.19	0.85	0.87	0.79	0.60	0.55
GAMSSA	4216	1.38	1.20	-0.04	0.76	0.76	0.84	0.70	0.70
GMPE	3990	1.39	1.08	-0.05	0.68	0.68	0.88	0.76	0.65
K10	4215	1.38	1.27	0.09	0.78	0.78	0.83	0.68	0.73
MUR	3979	1.39	1.09	-0.02	0.85	0.85	0.79	0.63	0.54
MWIR	4118	1.38	1.26	0.34	0.88	0.94	0.78	0.53	0.65
OISST	4163	1.39	1.22	0.09	0.70	0.70	0.86	0.74	0.75
OSTIA	4008	1.39	1.12	-0.47	0.86	0.98	0.79	0.51	0.56
LAC	876	1.30	1.27	0.06	0.51	0.51	0.92	0.85	0.91
WindSat	1892	1.22	1.19	1.03	0.71	1.25	0.83	-0.05	0.74

783  
 784 The sorted list of TS scores (Table 7) from high to low is:  $TS_{LAC} > TS_{OISST} \gg TS_{WS} > TS_{K10} >$   
 785  $TS_{GAMSSA} > TS_{GMPE} > TS_{MWIR} > TS_{OSTIA} > TS_{FNMOC} > TS_{MUR} \gg TS_{CMC}$ . Similarly, the sorted list  
 786 for the SS scores, is:  $SS_{LAC} > SS_{GMPE} \gg SS_{OISST} > SS_{GAMSSA} > SS_{K10} > SS_{CMC} > SS_{MUR} >$   
 787  $SS_{FNMOC} > SS_{MWIR} > SS_{OSTIA} \gg SS_{WS}$ . The L3 AVHRR LAC SSTs, once again, agree best overall with  
 788 the smallest centered RMS errors, normalized standard deviation closest to the observed, and maximum  
 789 correlation ( $\widehat{RMS}'_{LAC} = 0.39$ ,  $\hat{\sigma}_{LAC} = 0.98$ , and  $\rho_{LAC} = 0.92$ ). The majority of the L4s have excellent  
 790 to good skills in the coastal region around the Beaufort Gyre. Overall, however, the scores are generally  
 791 lower than for the northern buoys and are also lower than for the latter period of Louis 2012-03. This  
 792 could be a result of the dynamic range and SST variability. The limited LAC observations coincident with  
 793 the buoys could also suggest more issues with cloud coverage during this period. Among the best  
 794 products are the OISST ( $\widehat{RMS}'_{OISST} = 0.50$ ,  $\hat{\sigma}_{OISST} = 0.88$ , and  $\rho_{OISST} = 0.86$ ), and the GMPE  
 795 ( $\widehat{RMS}'_{GMPE} = 0.49$ ,  $\hat{\sigma}_{GMPE} = 0.78$  and  $\rho_{GMPE} = 0.88$ ). GAMSSA and K10, which had difficulties in the  
 796 2012 comparisons, showed good skill for this case study ( $\widehat{RMS}'_{GAMSSA} = 0.55$  and  $\widehat{RMS}'_{K10} = 0.56$ ).  
 797 Both analyses underwent some changes in January 2013, when they started ingesting WindSat SSTs.  
 798 The MUR, FNMOC, OSTIA and MWIR have the largest errors with  $\widehat{RMS}'$  between 0.6 and 0.65. These

799 numbers, however, are consistent with global statistics. A notable exception is the CMC product. This  
800 L4, which was at the top of the rankings for the 2012 buoy combinations, ended up last, according to the  
801 TS classifier ( $TS_{CMC} = 0.5$ , right below  $T_{51}$ ), due to the large difference in variance relative to the  
802 coastal buoys ( $\hat{\sigma}_{CMC} = 0.73$  or 73% of  $\sigma_{obs}$ ). This places the CMC farthest to the left from the dashed  
803 arc in the Taylor diagram, implying that the product is underestimating the SST amplitude variability at  
804 these buoy locations. Note, however, that the SS classifier places the product in fifth position among  
805 the products with good skills ( $\widehat{RMS}'_{CMC} = 0.58$ ). Given this and the fact that we are not aware of  
806 processing changes in the interim period that might have affected the skill of this analysis, we  
807 acknowledge that the CMC was perhaps unjustly penalized by the  $T_{51}$ , as there is some inherent  
808 arbitrariness in how the thresholds used to place products in discrete categories are selected.

#### 809 5.4. Combined Score

810 For a final ranking of the L4 products, we looked at three of the groupings analyzed here, as they  
811 correspond to different operating conditions/regimes: 1) the “weather” system passing, which included  
812 matchups with Louis 2012-03 before DOY 260 (Figure 7b; Table 5); 2) The “northern buoys” (Figure 9b,  
813 Table 6); and 3) The “coastal” buoys included in Section 5.3 (Figure 12; Table 7). Note that these  
814 categories include mutually exclusive data sets that, when combined, add up to all the matchups. We  
815 then considered the average of the TS and SS rankings after removing the L3 products from the sorted  
816 lists; i.e., if, after removing LAC and WindSat from both score rankings in the coastal regime, the OISST  
817 occupies first and second positions (highest rankings) in the TS and SS rankings, respectively, then the  
818 average ranking,  $\bar{S}$ , for this product is  $\bar{S}_{OISST} = 1.5$ . The results are shown in **Table 8**. Since there are  
819 nine L4 products being evaluated,  $\bar{S}$  varies from 1 – 9. In this approach, the lower the value of  $\bar{S}$ , the  
820 better the two skill scores. A natural clustering emerges for each of these regimes in which L4 products  
821 with comparable skills end up having similar average ranking. Table 8 indicates that the best products  
822 have  $\bar{S}$  between 1 and 3.5, the next best group between 4 and 6.5, and the poorest performing

823 products between 7.5 and 9. For instance, the OISST, GMPE, K10, and GAMSSA are all seemingly skillful  
824 in the coastal regime ( $\bar{S}$  between 1.5 and 3.5), whereas CMC, MUR, FNMOC, OSTIA, and MWIR are  
825 equally challenged ( $\bar{S}$  between 7.5 and 9.0). An overall ranking followed from averaging the mean  
826 position in each of the three regimes (see “Overall” in Table 8). In the overall classification the OISST,  
827 CMC, and GMPE, occupy the top three rankings (the most skillful overall), followed by K10, FNMOC,  
828 GAMSSA, MUR, OSTIA, and MWIR. In terms of product resolution, Table 8 indicates that, in general,  
829 high resolution L4 products (<10 km) constituted the best group for frontal regions and extreme  
830 weather conditions; the second best group for the cold Northerly waters, and the poor-performing  
831 group for the coastal regions. Coarser resolution (>10 km) L4 products performed best overall for all but  
832 the frontal regions.

833

834 Table 8. Final ranking of the selected L4 analyses for the individual regimes and overall period. Lower  
835 values indicate the best agreement with the UpTempO observations.

	OISST	CMC	GMPE	K10	FNMOC	GAMSSA	MUR	OSTIA	MWIR
Weather	6.5	1.5	6.5	3.5	3.5	9.0	1.5	8.0	5.0
North	2.5	1.0	2.5	6.5	5.0	4.0	8.0	6.5	9.0
Coastal	1.5	7.5	3.0	3.5	7.5	3.5	7.5	8.0	9.0
Overall	3.5	3.7	4.0	4.5	5.3	5.5	5.7	7.5	7.7

836

## 837 6. Conclusions

838 The Beaufort Sea is an extremely challenging region for SST analyses as evidenced by the very  
839 dramatic differences shown in contemporaneous scenes from the individual SST analyses. Despite the  
840 use of largely similar satellite input products and analysis procedures (they should be highly correlated),  
841 the L4 products exhibit significant differences in the amplitude and the phasing of their spatial patterns.  
842 Products found to be the best performing in other open ocean regions are not necessarily the best here.  
843 The UpTempO buoys provide a unique, and very valuable, verification data set since they are truly  
844 independent from the SST products evaluated in this study, giving great confidence in the results.

845 Taylor diagrams and skill scores provide very useful tools for comparative evaluation of the  
846 different SST analyses. One single statistical measure does not adequately capture all the aspects that  
847 might influence the relative skill of the different SST analyses in this challenging environment. Taylor  
848 diagrams provide a convenient way to graphically summarize the interplay among three different  
849 statistics that gauge different strengths and weaknesses in the products being evaluated. Skill scores  
850 allow emphasis of different measures and permit objective relative ranking of the different products.

851 The products found to be best performing varied with the region and conditions within the  
852 Beaufort Sea. Where available, the IR L3 AVHRR LAC data outperformed all of the analyses due largely  
853 to its high spatial resolution. While AMSR-E data were not available at the time of the comparison,  
854 WindSat SSTs performed well and appeared beneficial to the analyses at warmer temperatures, but  
855 exhibited biases at temperatures below about 2°C. Strong SST gradients in the region, particularly near  
856 the Alaskan coast, posed challenges for the L4 analyses and led to large differences among the products.  
857 SST analyses have lower resolutions that require upscaling the input data streams. This process can  
858 decrease the magnitude of cross-frontal SST gradients and/or slightly change the apparent location of  
859 the front, which increases the likelihood of the satellite product misrepresenting the point buoy  
860 measurement. In dynamic regions, the MUR and CMC analyses exhibited more realistic gradients. For  
861 the Louis 2012-03 buoy, the better gradient representation resulted in very strong product  
862 performance, but for other coastal buoys, slight uncertainties in the positioning of the gradients could  
863 have partially degraded the product scores. In coastal waters, also a region of high spatial variability,  
864 the OISST and the GMPE proved to be very skillful. Where temperatures were more uniform, such as  
865 the cold waters farther to the north, the products performed more similarly. Here, the CMC, OISST, and  
866 GMPE distinguished themselves from the others. The OISST performed best overall, with the best  
867 possible score when considering all buoy observations together, closely followed by the CMC and the

868 GMPE. Those products appear to have the best utility for applications in the challenging Beaufort Sea,  
869 at least during this period.

870 While inclusion of the bias in the SS implies that the score could be influenced by the effective  
871 depth of the analysis, that was not a significant factor here. Of the analyses, only FNMOC is  
872 representative of other than the foundation or daily average temperature. Given the largely isothermal  
873 conditions in the top 10 m, differences between foundation estimates and daily averages are expected  
874 to be small. Any biases associated with the skin layer were insignificant relative to other effects here  
875 such as cloud filtering.

876 Derived uncertainties in the analyses for the Beaufort Sea are generally greater than that  
877 observed globally. Martin et al. (2012) found that, globally, all the analyses they considered (which  
878 included all those considered here except for MUR) had standard deviations less than 0.7 K. In this  
879 study, except for the northern buoys where SSTs were very uniform, the standard deviations commonly  
880 exceeded 0.7 K. For the coastal buoys (Table 7), the analysis uncertainties exceeded the global values  
881 reported by Martin et al. (2012) by about 0.3 K. Interestingly, if we rank the L4s based on the standard  
882 deviations for the coastal buoys alone (column 5, Table 7), we obtain the same performing groupings  
883 reported by Martin et al. (2012) for their global results, despite the fact that those were estimated for  
884 an earlier period when a slightly different mix of sensors was active.

885 A significant difference between the SST analyses is the approach employed for ice masking.  
886 Theoretically, these differences could have an important impact on the merit of the different products  
887 near the ice. The results here suggested that the EUMETSAT OSI-SAF ice mask agreed very well with the  
888 ice flagging included with the UpTempO data, but the choice of ice mask had surprisingly little impact on  
889 the performance of the analyses relative to the buoys for this specific case.

890 The results here were certainly affected by lack of AMSR-E data during the study period. Several  
891 of the analyses that normally incorporate MW data did not do so during this period and others might



892 not have performed as well with the substituted WindSat data. While MW data are of coarser spatial  
893 resolution and cannot retrieve close to land or the ice edge, they provide important independent  
894 observations in cloudy conditions that can commonly obscure IR retrievals in this region. Performing a  
895 similar analysis in a future period when both buoy data and new AMSR2 data are available would be  
896 valuable.

897

898 **Acknowledgments:**

899 This work was supported by the MIZOPEX project with funding from NASA and NOAA. Additional time  
900 for SLC was supported through the Multi-Sensor Improved Sea Surface Temperature (MISST 2) project  
901 funded through the National Ocean Partnership Project (NOPP). The UpTempO buoy SSTs from the Polar  
902 Science Center (PSC) at the Applied Physics Laboratory (APL) in the University of Washington were  
903 obtained from the UpTempO Buoy Project web site at <http://psc.apl.washington.edu/UpTempO/>. The  
904 CMC SSTs were provided by Bruce Brasnett at CMC. JPL MUR SSTs were provided by Mike Chin at NASA  
905 JPL through the Physical Oceanography Distributed Active Archive Center (PODAAC, [podaac.  
906 jpl.nasa.gov/datasetlist](http://podaac.jpl.nasa.gov/datasetlist)). The GAMSSA, K10, OISST, MWIR, and OSTIA analyses were accessed through  
907 the NOAA NCEI Long Term Stewardship and Reanalysis Facility (LTSRF)  
908 ([www.nodc.noaa.gov/sog/GHRSSST/accessdata.html](http://www.nodc.noaa.gov/sog/GHRSSST/accessdata.html)). FNMOC SSTs for GHRSSST were obtained from the  
909 USGODAE server ([www.usgodae.org](http://www.usgodae.org)), hosted by the U.S. Navy at the Naval Research Laboratory Marine  
910 Meteorology Division (<http://www.usgodae.org/pub/outgoing/fnmoc/models/GHRSSST/>). The GMPE SST  
911 data were provided by GHRSSST, the Met Office, and MyOcean ([www.myocean.eu](http://www.myocean.eu)). The MWIR and  
912 WindSat SST retrievals were provided by Remote Sensing Systems, and the AVHRR LAC data were  
913 provided by NAVOCEANO.

914

915 **References**

- 916 Andersen, S., Breivik, L. A., Eastwood, S., Godøy, Ø., Lind, M., Porcires, M. & Schyberg, H. (2007). OSI  
917 SAF Sea Ice Product Manual – v3.5, EUMETSAT OSI SAF – Ocean and Sea Ice Satellite Application  
918 Facility, Tech. Rep. SAF/OSI/met.no/TEC/MA125, (pp. 36), available from  
919 [http://sat.met.no/docs/ss2\\_pmseaice\\_v3p5.pdf](http://sat.met.no/docs/ss2_pmseaice_v3p5.pdf).
- 920 Beggs, H., Zhong, A., Warren, G., Alves, O., Brassington, G., & Pugh, T. (2011). RAMSSA –an operational,  
921 high-resolution, multi-sensor sea surface temperature analysis over the Australian region, *Australian*  
922 *Meteorological and Oceanographic Journal*, *61*, 1 – 22.
- 923 Brasnett, B. (2008). The impact of satellite retrievals in a global sea-surface-temperature analysis,  
924 *Quarterly Journal of the Royal Meteorological Society*, *134*, 1745 – 1760, doi:10.1002/qj.319.
- 925 Buehner, M., Caya A., Pogson L., Carrieres, T., & Pestieau, (2013). A new Environment Canada regional  
926 ice analysis system, *Atmosphere - Ocean*, *51*, 18 – 34, doi: 10.1080/07055900.2012.747171.
- 927 Castro, S. L., Wick, G. A., & Emery, W. J. (2012). Evaluation of the relative performance of sea surface  
928 temperature measurements from different types of drifting and moored buoys using satellite-  
929 derived reference products, *Journal of Geophysical Research*, *117*, doi:10.1029/2011JC007472.
- 930 Castro, S. L., Wick, G. A., & Buck, J. J. H. (2014). Comparison of diurnal warming estimates from  
931 unpumped Argo data and SEVIRI satellite observations, *Remote Sensing of Environment*, *140*, 789 –  
932 799, doi:10.1016/j.rse.2013.08.042.
- 933 Chin, T. M., Milliff, R.F. & Large, W.G. (1998). Basin-scale, high-wavenumber sea surface wind fields  
934 from a multiresolution analysis of scatterometer data, *Journal of Atmospheric and Oceanic*  
935 *Technology*, *15*, 741 – 763.
- 936 Cummings, J. A. & Smedstad, O. M. (2013). Variational Data Assimilation for the Global Ocean, Chapter  
937 13, pp 303 – 343. In *Data Assimilation for Atmospheric, Oceanic and Hydrological Applications* (Vol.  
938 II). S. Park and L. Xu (eds), doi:10.1007/978-3-642-35088-7 13, Springer-Verlag, Berlin, Heidelberg.

939 Daley, R. (1991). Atmospheric Data Analysis, Cambridge University Press, Cambridge, UK.

940 Dash, P., Ignatov, A., Martin, M., Donlon, C., Brasnett, B., Reynolds, R., & Coauthors (2012). Group for  
941 High Resolution Sea Surface Temperature (GHR SST) analysis fields inter-comparisons. Part 2: Near-  
942 real time web-based Level 4 SST Quality Monitor (L4-SQUAM), *Deep-Sea Research II*, 77-80, 31 – 43,  
943 doi:10.1016/j.dsr2.2012.04.002.

944 Dong, S., Guille, S.T., Sprintall, J., & Gentemann, C.L. (2006). Validation of the Advanced Microwave  
945 Scanning Radiometer for the Earth Observing System (AMSR-E) sea surface temperature in the  
946 Southern Ocean, *Journal of Geophysical Research*, 111, C04002, doi:10.1029/2005JC002934.

947 Donlon, C. J., I. Robinson, K. S. Casey, J. Vazquez-Cuervo, E. Armstrong, O. Arino & Coauthors (2007). The  
948 Global Ocean Data Assimilation Project (GODAE) High Resolution Sea Surface Temperature Pilot  
949 Project (GHR SST-PP), *Bulletin of the American Meteorological Society*, 88, 1197-1213.

950 Donlon, C. J., Martin, M., Stark, J. D., Robert-Jones, J., Fiedler, E., & Wimmer, W. (2012). The operational  
951 sea surface temperature and sea ice analysis (OSTIA) system, *Remote Sensing of Environment*, 116,  
952 140 – 158, doi:10.1016/j.rse.2010.10.017.

953 Gandin, L. (1965). Objective Analysis of Meteorological Fields. *Israel Program for Scientific Translations*,  
954 22 pp.

955 Gentemann, C.L., Wentz, F.J., & DeMaria, M. (2006). Near real time global optimum interpolated  
956 microwave SSTs: Applications to hurricane intensity forecasting. Presented at 27<sup>th</sup> Conference on  
957 Hurricanes and Tropical Meteorology, Monterrey, CA.

958 Grumbine, R. W. (1996). Automated passive microwave sea ice concentration analysis at NCEP. NOAA  
959 Tech. Note 120, 13 pp. [available from NCEP/NWS/NOAA, 5200 Auth Road, Camp Springs, MD  
960 20746.]

961 Hernández-Orallo, J., Flach, P., & Ferri, C. (2012). A unified view of performance metrics: translating  
962 threshold choice into expected classification loss, *Journal of Machine Learning Research*, *13*, 2813 –  
963 2869.

964 Hoyer, J. L., Karagali, I., Dybkjaer, G., & Tonboe, R. (2012). Multi sensor validation and error  
965 characteristics of Arctic satellite sea surface temperature observations, *Remote Sensing of*  
966 *Environment*, *121*, 335 – 346, doi:10.1016/j.rse.2012.01.013.

967 Mallat, S.G. (1989). A theory for multiresolution signal decomposition: The wavelet representation. *IEEE*  
968 *Trans. Pattern Analysis Machine Intelligence*, *11*, 674–693.

969 Markus, T., & Cavalieri, D. J. (2000). An enhancement of the NASA Team sea ice algorithm, *IEEE*  
970 *Transactions on Geoscience and Remote Sensing*, *38*, 1387 – 1398.

971 Martin, M., Dash, P., Ignatov, A., Banzon, V., Beggs, H., Brasnett, B., & Coauthors (2012). Group for High  
972 Resolution Sea Surface Temperature (GHRSSST) analysis fields inter-comparisons. Part 1: A GHRSSST  
973 multi-product ensemble (GMPE), *Deep-Sea Research II*, *77-80*, 21 – 30,  
974 doi:10.1016/j.dsr2.2012.04.013.

975 May, D. A., M. M. Parmeter, D. S. Olszewski, and B. D. McKenzie (1998). Operational processing of  
976 satellite sea surface temperature retrievals at the Naval Oceanographic Office, *Bulletin of the*  
977 *American Meteorological Society*, *79*, 397-407.

978 Murphy, A. H. (1988). Skill scores based on the mean square error and their relationships to the  
979 correlation coefficient, *Monthly Weather Review*, *116*, 2417 – 2424.

980 Przybylak, R. (2003). *The Climate of the Arctic*. Series: Atmospheric and Oceanographic Sciences Library,  
981 Vol. 26, Kluwer Academic Publishers, Springer, XI, 272 p. ISBN 1-4020-1134-2

982 Reynolds, R. W., Smith, T. M., Liu, C., Chelton, D. B., Casey, K. C., & Schlax, M. G. (2007). Daily high-  
983 resolution blended analyses for sea surface temperature, *Journal of Climate*, *23*, 5473 – 5496.

984 Reynolds, R. W. & Chelton, D. B., (2010). Comparisons of daily sea surface temperature analyses for  
985 2007 – 08, *Journal of Climate*, 23, 3545 – 3563, doi:10.1175/2010JCLI3294.1.

986 Simmonds, I. & Rudeva, I. (2012). The great Arctic cyclone of August 2012, *Geophysical Research Letters*,  
987 39, L23709, doi:10.1029/2012GL054259.

988 Smith, G.C., Roy, F., Reszka, M., Surcel-Colan, D., He, Z., Deacu, D., & Coauthors (2015). Sea ice forecast  
989 verification in the Canadian Global Ice Ocean Prediction System, *Quarterly Journal of the Royal*  
990 *Meteorological Society*, 142, 659 – 671, doi: 10.1002/qj.2555.

991 Steele, M., Ermold, W., Colburn, K., & Rigor, I. (2016). UpTempO buoy, *Journal of Atmospheric and*  
992 *Oceanic Technology*, in preparation.

993 Taylor, K. E., (2001). Summarizing multiple aspects of model performance in a single diagram, *Journal of*  
994 *Geophysical Research*, 106, 7183 – 7192.

995 Zhang, J., Lindsay, R., Schweiger, A., & Steel, M. (2013). The impact of an intense summer cyclone on  
996 2012 Arctic sea ice retreat, *Geophysical Research Letters*, 40, 720 – 726, doi:10.1002/grl.50190.

997 Zhong, A. & Beggs, H. (2008). Operational Implementation of Global Australian Multi-Sensor Sea Surface  
998 Temperature Analysis, *Analysis and Prediction Operations Bulletin No. 77*, Bureau of Meteorology,  
999 Australia, 2 October 2008, <http://www.bom.gov.au/australia/charts/bulletins/apob77.pdf>.

1000

1001

1002