

1
2
3

4 **Combining Statistical, Physical, and Historical Evidence to Improve Historical
5 Sea-Surface Temperature Records**

6
7

Duo Chan^{†,*}

[†] Department of Earth and Planetary Sciences, Harvard University

ABSTRACT. Reconstructing past sea-surface temperatures (SSTs) from historical measurements containing more than 100 million ship-based observations taken by over 500,000 ships from more than 150 countries using a variety of methodologies creates a wide range of historical, scientific, and statistical challenges. The reconstruction of historical SSTs for studying climate change is particularly challenging because SST measurements are uncertain and contain systematic biases of order 0.1°C to 1°C—these systematic biases are in the range of the historical global warming signal of approximately 1°C. The biases are complicated and have generally been addressed using simplified corrections. In this review, I introduce a history of SST observations, review a statistical method developed for quantifying SST biases, and illustrate scientific insights obtained from adjusted SSTs. This article also documents the scientific journey of my Ph.D. work. As a result, I report personal stories on both successes, difficulties, and setbacks along the way. The statistical method for correcting SSTs (i.e., a linear-mixed-effect intercomparison framework) depends on identifying systematic offsets between intercomparable groups of SST observations. Combining estimated offsets with physical and historical evidence has allowed for correcting discrepancies associated with SSTs, including the North Atlantic warming twice as fast as the North Pacific in the early 20th century and anomalously warm SSTs during World War II. Corrections also permit better hindcasting of Atlantic hurricanes. I conclude with some discussion on how the SST records might be further improved. Given the importance of SSTs for understanding historical changes in climate, I hope that this review can help others appreciate challenges that are present and spark some interest and ideas for further improvement.

8 **Keywords:** climate reconstruction, sea-surface temperature, bias correction, data homoge-
9 nization, linear-mixed-effect model.

10

* duochan@g.harvard.edu

11

MEDIA SUMMARY

12 To better predict what climate change will look like in the future, it is crucial to know how
 13 and why climate has changed in the past. One essential component of climate change is the
 14 ocean, for which we have more than 200 years of ship-based temperature measurements made
 15 at the ocean surface. However, biases in early sea-surface temperatures have limited their usage
 16 in climate studies. These biases are similar in magnitude to historical warming, and they vary
 17 with measurement methods, instruments, protocols, and even postprocessing and data-keeping
 18 practices. The question is, therefore, can we remove the complicated biases and obtain a sea-surface
 19 temperature estimate that is accurate enough to study past climate change?

20 In this article, I review recent progress aimed at correcting sea-surface temperatures for individual
 21 nations and data-collecting groups. I introduce a statistical framework that compares nearby
 22 measurements and estimates systematic offsets in temperatures among groups. Physical and historical
 23 evidence is then combined to understand the origins of significant groupwise differences detected by the statistical method. Correcting data leads to spatially more homogeneous warming
 25 in the early 20th century and removes anomalously warm sea temperatures during World War II,
 26 which reconciles existing model-data discrepancies and brings observations into consistency with
 27 current knowledge of climate forcing and variability. Beyond the ocean itself, adjusted sea-surface
 28 temperatures also allow atmospheric models to simulate more realistic historical variations in North
 29 Atlantic hurricanes, showing potential for improving predictions of these high-impact events. This
 30 review also demonstrates the importance of understanding the social context and history of how
 31 data are collected and postprocessed. When data and models disagree, keeping an awareness of
 32 potential flaws in the quality of data appears to be a necessary practice.

33

1. INTRODUCTION

34 Sea-surface temperature (SST), typically defined at ocean depth of 20–30 cm (Kennedy et al.,
 35 2019), is a crucial quantity for studying the Earth’s climate. Estimates of historical SSTs to an
 36 accuracy of 0.05°C at the global scale and 0.1°C at regional scales are required for a wide range of
 37 climate applications (Kent & Berry, 2008), which include depicting past climate change (Hartmann
 38 et al., 2013), attributing anthropogenic versus internal climate variability (Bindoff et al., 2013), and
 39 understanding changes in climate and weather events that have far-reaching societal impacts, such
 40 as El Niño (Yeh et al., 2009) and hurricanes (Vecchi et al., 2011). Moreover, SSTs are often used
 41 as boundary conditions in numerical models to reproduce or hindcast a variety of meteorological
 42 phenomena (e.g., Gates et al., 1999).

43 Despite their importance for climate sciences, estimates of historical SSTs remain highly uncertain (P. Jones, 2016), with disagreements existing between observational estimates and climate-model simulations. One example of such data-model disagreement would be the recent warming
 44 hiatus, which refers to a slowdown in the increase of the observed global-mean surface temperature
 45 since the late 1990s (Easterling & Wehner, 2009). The hiatus was one of the most popular
 46 climate-related research topics in the first half of the 2010s. At that time, state-of-art climate models
 47 that were used in the latest Intergovernmental Panel on Climate Change report (IPCC AR5,
 48 Taylor et al., 2012) simulated significantly faster warming ($p < 0.05$) than observational estimates
 49 (Fyfe et al., 2013; Medhaug et al., 2017). One possible explanation, as suggested by many studies,
 50 involves natural climate fluctuations that can uptake heat from the surface into the deep ocean
 51 (e.g., X. Chen & Tung, 2014; Kosaka & Xie, 2013). Another plausible explanation, however, is that
 52

54 recent trends in observed SSTs were underestimated by $0.064^{\circ}\text{C}/\text{decade}$ from 2000 to 2014 due to
55 biases in ship-based measurements (Karl et al., 2015). Correcting SSTs reduces the so-called recent
56 warming hiatus, making the estimate of warming rates since the early 2000s consistent with the
57 rapid warming since the late 1970s (Hausfather et al., 2017; Karl et al., 2015). In another example,
58 after removing contributions from major physical modes of climate variations, Thompson et al.
59 (2008) detected a sudden drop of about 0.3°C in global-mean SSTs immediately after World War
60 II, which they attributed to insufficient corrections of instrumental SST biases.

61 When data and models disagree, a common practice is to assume that data reflect reality and
62 to look for new theories to enrich the model and explain the data (e.g., using natural climate
63 fluctuations to explain the recent warming hiatus). However, there is always a second and often
64 overlooked possibility: that data contain undetected problems. As we shall see in detail in later
65 sections, in addition to the recent warming hiatus and the artificial temperature drop at the end of
66 World War II, major data problems also exist in SST estimates in terms of patterns of warming in
67 the early 20th century and temperature evolution during World War II. Observed historical SSTs
68 are particularly likely to contain data problems because of complicated biases associated with using
69 various crude methods to collect early measurements and also because of simplified bias corrections
70 employed when generating SST estimates.

71 **1.1. A Brief History of Measuring SSTs since the 1800s.** Instrumental SSTs have been
72 measured on ships at the ocean surface for more than 200 years, yielding more than 130 million
73 ship-based measurements since the 1850s (Freeman et al., 2017). Such a history is longer than
74 that of studies on anthropogenic climate change; the first estimate of equilibrium warming once
75 doubling atmospheric CO₂ was made by Svante Arrhenius in 1896 (Lapen, 1998). The history
76 of SST measurement is also much longer than that of dedicated scientific efforts to systematically
77 monitor ocean temperatures, which began in the late 1970s.

78 Most of these early SST measurements were made not by dedicated researchers but by voluntary
79 and nonscientist sailors from different countries who sailed, for example, as soldiers, merchants, or
80 fishermen (Kennedy, 2014; Kent et al., 2017). These early SST measurements were made for variety
81 of purposes, including pure scientific interest, facilitating navigation, predicting stormy weather,
82 and mapping a climatological summary of the marine environment (Kennedy, 2014). Although
83 not made for purposes of monitoring climate change, ship-based in-situ observations are the only
84 available source of direct measurements of past states at the ocean surface.

85 Historical SST values were originally recorded in ship logs and were rescued and digitized by a
86 variety of projects and institutions (e.g., Wilkinson et al., 2011). Digitized data from various sources
87 were later put together to construct the International Comprehensive Ocean-Atmosphere Data Set
88 (ICOADS, the most comprehensive modern compilation of available marine meteorological mea-
89 surements since the 1700s). The digitization of ship logs and the construction of ICOADS spanned
90 several decades (Freeman et al., 2017; Woodruff et al., 1998; Woodruff et al., 1987; Woodruff et al.,
91 2011; Worley et al., 2005), and this initiative was accompanied by revolutions in computer and
92 data-storage technologies. ICOADS3 is the latest version, and efforts continue to recover lost his-
93 torical data sets and to include missed metadata during initial digitizations by reprocessing existing
94 data banks (Kent et al., 2017).

95 In addition to the changing purposes of measurements and record-keeping efforts, instruments
96 and associated systematic SST biases during measurement have also experienced major changes.
97 The first instruments to systematically measure SSTs across large spatial scales were buckets and

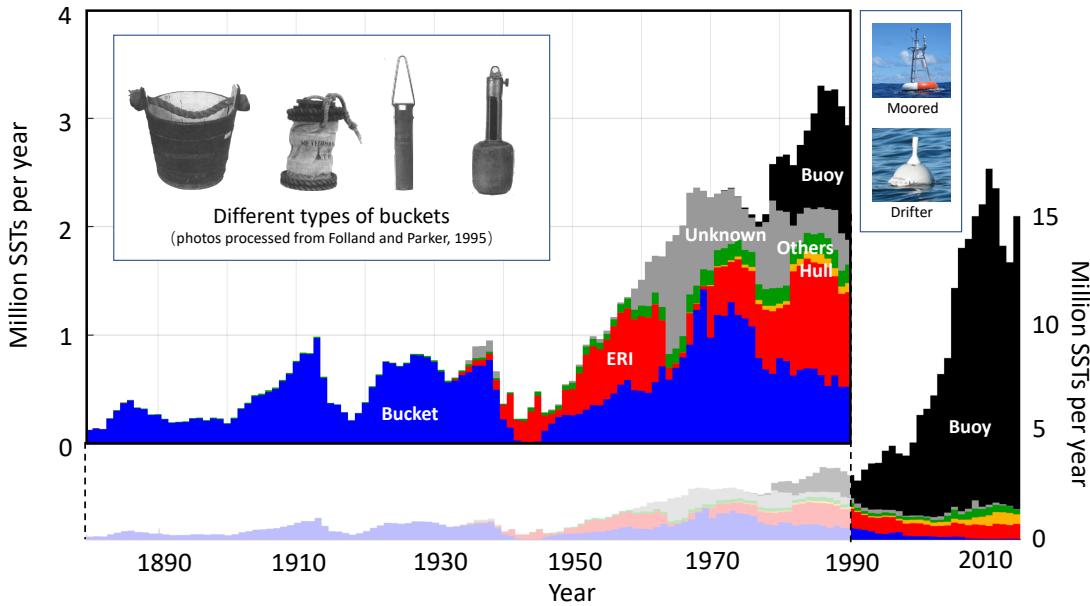


Figure 1. Distinct methods used to take in-situ SSTs compiled under the International Comprehensive Ocean Atmosphere Data Set (ICOADS). The overall number of in-situ SSTs collected by different measurement methods (stacked bars) in individual years from 1880 to 2014. In addition to bucket (blue), engine-room intake (ERI, red), hull sensor (orange), and buoy (black), other methods (green) include radiation thermometer, reversing thermometer, and electronic sensors, which are, however, not thought to be representative of SSTs due to their limited numbers (Kent et al., 2010). Results are based on version 3.0 of ICOADS. Method information is inferred for some unknown measurements, following Kennedy et al. (2011b). For example, SSTs before 1941 come from buckets, if not explicitly indicated otherwise, and U.S. and U.K. Naval SSTs during World War II are assumed to be ERI measurements (Kennedy et al., 2011b). Also shown are photos of some types of buckets used in SST collections, as well as images of moored and drifting buoys that have been widely deployed since the 1980s.

thermometers. The procedure to measure SSTs involved hauling buckets of water from the ocean surface and measuring the temperature of water in buckets on ship decks. SSTs made by this method (hereafter bucket SSTs) are thought to dominate ICOADS before the 1940s, and the number of bucket SSTs gradually decreased after the mid-1970s (Figure 1). During the measurement process, water temperatures in buckets will generally become colder due to wind-induced evaporation, as well as sensible heat loss in the tropics and the subtropics (Folland & Parker, 1995). In midlatitudes, bucket biases are still expected to be cold in winter. During summer, evaporation is suppressed in humid air, and the direction of sensible heat flux can be reversed as air temperatures become warmer than SSTs, leading to less heat loss (Folland & Parker, 1995). Sometimes, bucket water can be heated by the sun, especially on a calm summer afternoon (Kennedy et al., 2019). When averaged annually and over the globe, early bucket SSTs are estimated to be biased cold to an order of 0.4°C (Folland & Parker, 1995; Kennedy et al., 2011b). However, buckets of different materials and designs have been used under different protocols in history, which could lead to distinct biases among groups of bucket SSTs (Folland & Parker, 1995; Kent & Taylor, 2006). For example, a

112 less-insulated canvas bucket can be colder than a more-insulated wooden bucket by around 0.5°C
113 even measured under the same conditions (Folland & Parker, 1995). The time gap between water
114 retrieval and measurement can also affect bucket bias.

115 After the emergence of engine ships in the late-19th century, a second method of measuring
116 SSTs was introduced, which is a byproduct when monitoring the temperature of inlet water before
117 entering and cooling ship engines. SSTs made by the engine-room intake (ERI) method first appear
118 in ICOADS in the 1930s (Figure 1) and are mainly from U.S. ships that dominated the Atlantic
119 and the northeast Pacific. Later, the ERI method was adopted by more nations and gradually
120 became the preferred method because of safety concerns associated with hauling buckets on fast
121 engine ships (Kennedy, 2014). ERI SSTs typically come from a depth of 5–15 m where the ocean
122 is less affected by solar heating and should consequently be colder than SSTs defined at 20–30 cm
123 (Carella et al., 2018; Chan & Huybers, 2020b). However, because of the absorption of heat from
124 ship engines, ERI measurements are estimated to have warm biases of 0.1°C to 0.3°C, depending
125 on ship design and cargo (Kennedy et al., 2011b).

126 In the modern era, a variety of new methods that give more reliable SSTs have been used
127 (Figure 1). Since the 1970s, an increasing number of ships are equipped with specialized digital
128 sensors (Kennedy, 2014). SSTs from hull sensors should be free of engine heating and are therefore
129 expected to be less biased. Scientists have also been deploying drifting and moored buoys since
130 the late 1970s, which has become the dominant data source since the 1990s (Figure 1). Similar
131 to hull sensors, buoys make contact with seawater directly and are expected to give less biased
132 SST measurements, although individual buoys could be problematic due to instrumental drift or
133 biofouling (Kennedy et al., 2012; Kent et al., 2017). Biofouling refers to the accumulation of small
134 ocean organisms on the wet surface of instruments, leading to structural or functional deficiencies.
135 Whereas drifting buoys typically measure at a depth of 20–30 cm, most moored buoys measure at
136 around 1 m deep (Kennedy, 2014). The deployment of drifting buoys substantially increased the
137 spatial coverage of the observing system, especially in the southeastern Pacific and the Southern
138 Ocean, which nonresearch ships rarely traverse. The majority of moored buoys are installed along
139 the coastal United States as marine weather stations and over the tropical Pacific and the Indian
140 Ocean to monitor El Niño evolution (Hervey, 2014). Combining different types of buoys, which
141 sample at different depth, can result in biases due to vertical temperature gradients that often exist
142 near the ocean surface. One cause of these gradients is solar heating in low-wind conditions, which
143 can exceed 3°C in some extreme cases (Kennedy et al., 2007). This depth effect may be damped
144 by ship-induced turbulent mixing for ship-based SSTs and may appear small relative to bucket and
145 ERI biases when averaged over seasons and weather conditions. However, recent ship-based SSTs
146 are reported to be, on average, systematically warmer than collocated buoy SSTs on an order of
147 0.1°C (e.g., Huang et al., 2017; Karl et al., 2015), which needs to be accounted for when combining
148 SSTs from both sources.

149 In addition to in-situ observations collected at the ocean surface, satellite and other remote-
150 sensing techniques became available in the 1980s. Remote-sensing techniques further increase the
151 spatial and temporal coverage of SST measurements. Note that satellites observe the skin tem-
152 perature in the upper several millimeters of the ocean (Kennedy et al., 2007), and this difference
153 in sampling depth, again, needs to be accounted for when homogenizing with in-situ SSTs. Addi-
154 tionally, SSTs retrieved from satellites can be biased due to changes in atmospheric optical depth
155 associated with volcanic and anthropogenic aerosols and, therefore, have to be calibrated and cor-
156 rected against in-situ measurements (T. M. Smith et al., 2008).

157 In addition to instruments dedicated to measuring SSTs, near-surface temperatures are also
 158 available from ocean profiling instruments. Historical profiles have been made on research vessels
 159 for more than a hundred years (Meyssignac et al., 2019). Since the late 1990s, scientists have been
 160 deploying Argo floats that profile temperature and conductivity as functions of pressure (Roemmich
 161 et al., 1999). Since 2006, Argo floats have been able to provide temperatures within the upper five
 162 meters of the ocean with nearly global coverage (Huang et al., 2017). Some of the most recent
 163 Argo floats can provide temperatures at 0.1-meter resolution in the upper 200 meters of the ocean.
 164 Currently, there are approximately 4,000 Argo floats providing nearly global information on near-
 165 surface temperatures at a frequency of once every 10 days.

166 **1.2. Insufficient Metadata and Simplified Corrections.** The shift from measuring SSTs from
 167 buckets to ERI to buoys is, therefore, accompanied by systematic biases varying on the order
 168 of 0.5°C. Such a variation in bias has a similar magnitude to the less-than-1°C global warming
 169 that is thought to have happened in the 20th century. On account of the irreplaceable nature of
 170 these early SST measurements, adjusting biases becomes crucial for quantifying and interpreting
 171 historical climate change. Such a problem, however, is difficult because biases associated with SSTs
 172 coming from the same method can be distinct due to different instrumental designs (e.g., different
 173 bucket materials) and measurement protocols used by various subsets of ships, which will interact
 174 with the uneven sampling to create regionally varying biases.

175 Although biases are complicated, lack of metadata by which to make specific corrections has
 176 necessitated simplifying assumptions regarding the spatial and temporal structure of SST biases,
 177 which inevitably lead to insufficient corrections and SST estimates having higher uncertainty than
 178 land surface temperatures (P. Jones, 2016). SST products from the U.K. Met Office, for example,
 179 assumed that a transition from early wooden buckets to less-insulated canvas buckets happened
 180 with the percentage of canvas buckets increasing linearly, from 35% in 1880 to 100% in 1920 over
 181 the entire ocean (Folland & Parker, 1995; Kennedy et al., 2019; Kennedy et al., 2011b). In other
 182 words, all bucket measurements in the same year are assumed to be biased in the same way, as if
 183 they were measured by the same person using the same bucket. In other SST estimates, corrections
 184 do not distinguish between measurement methods. Rather, biases for SSTs from all methods
 185 are represented using a large-scale fixed pattern, with the amplitude of the pattern estimated by
 186 comparing SSTs with other temperature estimates, for example, nighttime marine air temperatures
 187 (Huang et al., 2015; Huang et al., 2017) or coastal station-based air temperatures (Cowtan et al.,
 188 2018).

189 In addition to biases introduced during measurement, problems may also occur in the record-
 190 keeping and data-processing stage, as information is transferred over time and across technologies.
 191 One example is inaccurate metadata that leads to ERI SSTs being misclassified as coming from
 192 buckets (Carella et al., 2018; Kennedy et al., 2011b). Other postmeasurement problems may also
 193 exist but have not yet been quantified systematically.

194 2. TOWARD REFINED CORRECTIONS FOR INDIVIDUAL NATIONS AND GROUPS OF DATA

195 During my Ph.D. study, I aimed to refine SST corrections by resolving the regional biases that
 196 arise from different measurement and postprocessing characteristics due to distinct physical and
 197 historical reasons. Because ships from the same nation and data-collecting group would have used
 198 similar instruments, followed similar protocols, and experienced similar postprocessing practices, I
 199 corrected biases for individual nations and data-collecting groups, assuming data within the same

group have similar bias characteristics. Specifically, nation information is mainly identified using ICOADS country code (Freeman et al., 2017) and metadata from the World Meteorology Organization No. 47 publication (Kent et al., 2007). Collecting groups are assigned using ICOADS deck number, which is the primary field to track ICOADS data collection and postprocessing (Freeman et al., 2017). Because available metadata is insufficient for physically constraining biases in individual groups, I turned the research question into a big-data problem and developed a statistical method to estimate corrections required for individual groups.

Ship-based SSTs contain information of both physical SST variations and biases. Because nearby measurements tend to have similar physical SST variations, examining differences between nearby SSTs allows for focusing on data heterogeneity associated with distinct biases. SST measurements were, therefore, first paired if they came from different groups and were within 300 km and two days of one another. These scales were chosen to keep expected physical variability with biases on the order of tenths of a degree Celsius. The results are not qualitatively sensitive to scales used in pairing SSTs (Chan & Huybers, 2019). To prevent error covariance between pairs, each measurement was used only once, with an algorithm prioritizing measurements closest in space. Specifically, the method rank-orders all potential pairs within a given month according to distance and selects the closest pair. The next closest pair is selected after removing previously selected measurements. The process repeats until all paired measurements are selected (Chan & Huybers, 2019).

Because measurements in a pair are not perfectly collocated in space and time, we first removed expected physical differences arising from displacements in geographical locations, seasonality, and day-night differences (Chan & Huybers, 2019). For example, SSTs closer to the Equator or during daytime are expected to be warmer. Expected differences were estimated from the 1982–2014 climatology of high-resolution satellite-based retrievals (T. M. Smith et al., 2008). Remaining differences in reported SSTs ($\delta\mathbf{T}$) were represented using a linear-mixed-effect (LME) model,

$$(2.1) \quad \delta\mathbf{T} = \mathbf{X}\boldsymbol{\alpha} + \mathbf{Z}_r\boldsymbol{\beta}_r + \mathbf{Z}_y\boldsymbol{\beta}_y + \boldsymbol{\epsilon},$$

where $\delta\mathbf{T}$ is represented as a fixed-effect term describing offsets between groups ($\boldsymbol{\alpha}$) and random effects describing regional ($\boldsymbol{\beta}_r$) and temporal ($\boldsymbol{\beta}_y$) variations. We constrain $\boldsymbol{\alpha}$ such that the average offset of all compared measurements is zero. \mathbf{X} , \mathbf{Z}_r , and \mathbf{Z}_y are design matrices that specify, respectively, common pairs of groups, years, and regions. See Figure 2 for an element-wise illustration of the LME model. Such a model is similar to an ANOVA approach. It makes use of random effects to give more conservative estimates. As the number of pairs available for constraining a random effect decreases, the estimate is relaxed toward zero such that estimates of regional and yearly variations in groupwise offsets are robust against noise (Chan & Huybers, 2019). The model is flexible in terms of controlling for specific effects and is easily extendable to account for variations in offsets associated with seasonality and day-night differences.

In practice, to reduce computational cost, SST differences are aggregated according to combinations of pairs of groups, regions, and years before estimating offsets (Chan & Huybers, 2019). Pairs in an aggregate are assigned equal weights when taking the average. Errors of aggregated SST differences ($\boldsymbol{\epsilon}$) are budgeted to account for errors from different sources and heteroscedasticity associated with distinct group size (Chan & Huybers, 2019). See Figure 3a for an example of uneven numbers of pairs between different combinations of groups. The error of each aggregated pair (ϵ_k)

$$\begin{aligned}
 \left[\begin{array}{c} \delta T_1 \\ \delta T_2 \\ \delta T_1 \\ \delta T_2 \\ \vdots \\ \delta T_p \end{array} \right]_{p \times 1} &= \left[\begin{array}{ccccc} 1 & -1 & \dots & 0 \\ 0 & 1 & \dots & -1 \\ 1 & 0 & \dots & 0 \\ 0 & 0 & \dots & -1 \\ \vdots & \vdots & \ddots & \vdots \\ -1 & 0 & \dots & 1 \end{array} \right] \left[\begin{array}{c} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_g \end{array} \right] + \\
 &\quad \left[\begin{array}{cccccc} 1 & 0 & \dots & 0 & -1 & 0 & \dots & 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 & 0 & 1 & \dots & 0 & 0 & -1 & \dots & 0 \\ 1 & 0 & \dots & 0 & 0 & 0 & \dots & 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 & 0 & 0 & \dots & 0 & 0 & 0 & \dots & -1 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & -1 & \dots & 0 & 0 & 0 & \dots & 0 & 0 & 1 & \dots & 0 \end{array} \right] \\
 &\quad \downarrow \quad \downarrow \quad \downarrow \\
 &\quad \delta T_{\langle p \times 1 \rangle} \quad X_{\langle p \times g \rangle} \quad \alpha_{\langle g \times 1 \rangle} \\
 \text{Paired SST} & \text{Group-level} & \text{Group-wise} & \text{Regional} & \text{Regional} \\
 \text{difference} & \text{design} & \text{offsets} & \text{design} & \text{offsets} \\
 & \text{matrix} & (\text{fixed effect}) & \text{matrix} & (\text{random effect})
 \end{aligned}$$

Figure 2. An element-wise illustration of the LME model in Eq. 2.1. Also shown is the dimensionality of matrices and vectors (red), where p , g , and r are, respectively, numbers of pairs, groups, and regions. \mathbf{X} , \mathbf{Z}_r , and \mathbf{Z}_y are design matrices whose entries are one, zero, or minus one. Regional effects (β_r) are estimated for individual groups. These regional effects are assigned as random effects and are assumed to follow a Gaussian distribution such that each $\beta_{r_{ij}} \sim N(0, \sigma_r^2)$. Yearly effects, $\mathbf{Z}_y\beta_y$, are also estimated for individual groups and have a similar structure to $\mathbf{Z}_r\beta_r$. Higher-order interactions that involve group, year, and regions are not accounted for in this model to limit the number of free parameters.

is assumed to follow $N(0, \bar{\sigma}_k^2)$, where

$$(2.2) \quad \bar{\sigma}_k^2 = \frac{2\sigma_o^2}{n_k} + \frac{2\sigma_s^2}{n_k^x} + \frac{\sum \sigma_c^2(l)}{n_k^2}.$$

242 $\frac{2\sigma_o^2}{n_k}$ denotes the contribution of random observational errors, where n_k is the number of pairs in the
 243 k^{th} aggregate. Random observational error is denoted by σ_o^2 and is estimated to be $0.86 \pm 0.18^{\circ}\text{C}$
 244 (2 SD, Chan et al., 2019) for individual bucket SSTs. Contributions of partially correlated obser-
 245 vational errors are denoted by $\frac{2\sigma_s^2}{n_k^x}$, with σ_s^2 estimated to be $0.38 \pm 0.14^{\circ}\text{C}$ (Chan et al., 2019). One
 246 possible source of σ_s^2 is systematic errors associated with individual ships. Because ship information
 247 is not always available in ICOADS, n_k^x is used to approximate effective numbers of ships within
 248 the k^{th} aggregate, with x estimated to be 0.57 (Chan et al., 2019). Finally, $\sigma_c(l)$ denotes uncer-
 249 tainties associated with physical SST variations for the l^{th} out of n_k pairs. The estimation of σ_c
 250 accounts for interannual variance and covariance of physical SSTs as a function of location, month,
 251 and displacement, with more details documented in section 5.a.1 of Chan and Huybers (2019). The
 252 robustness of offset estimates to a variety of model formulations and assumptions was explored in
 253 section 5.b of Chan and Huybers (2019).

In the following sections, I will show that the LME method detects significant offsets among groups classified by both nation and deck number. Accounting for these systematic groupwise offsets improves the quality of historical SST estimates at regional and sub-basin scales, resolves

several existing data-model discrepancies, and brings in new opportunities for understanding extreme weather events and climate variations. Beyond the nation-and-deck level, the model in Equation 2.1 can be extended to resolve offsets associated with individual ships for more refined SST corrections. However, metadata of ship information and the algorithm used to fit the LME model need to be improved before estimating ship-level offsets. These steps will be carried out in future works, and associated plans will be discussed in section 7.1.

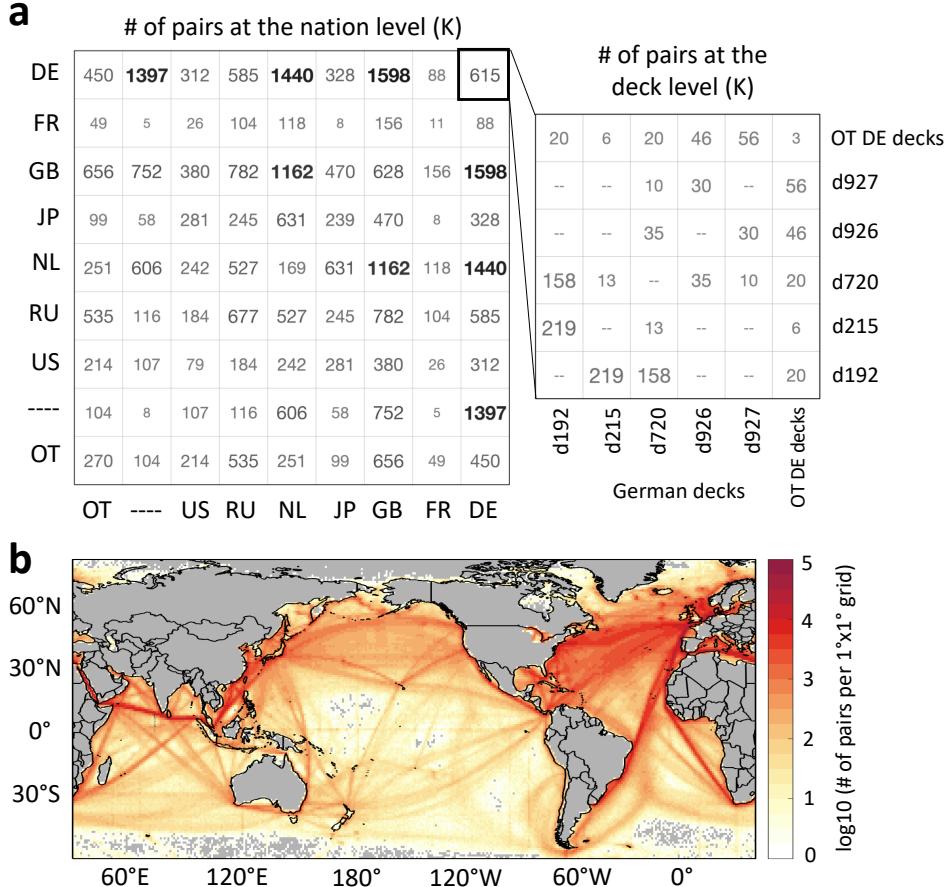


Figure 3. Schematics of an LME intercomparison. (a) Numbers of SST pairs between nation-deck groups (in the unit of one thousand pairs) in the analysis that intercompares SSTs thought to come from buckets. The comparison happens not only between different nations (left) but also between distinct decks from the same nation (e.g., right; zooming in the top-right box in the left and showing the comparison between German decks). The number of pairs can be very different across combinations of groups, with ‘- -’ denoting that no pairs are found between corresponding groups. Nation abbreviations are for Germany (DE), France (FR), Great Britain (GB), Japan (JP), the Netherlands (NL), Russia (RU), the United States (US), and unknown (- -). Nations that contribute to fewer than 500,000 are labeled as “other nations” (OT) for this visualization but are distinguished in the LME analysis. Similarly, Germany decks that contribute to fewer than 50,000 pairs are shown as “OT DE decks.” (b) The spatial distribution of paired measurements follows major ship tracks.

3. OFFSETS AMONG BUCKET GROUPS AND MORE UNIFORM EARLY-TWENTIETH-CENTURY WARMING

When applied to measurements thought to come from buckets, the LME analysis intercompares 17.8 million pairs coming from 162 groups from 1850 to 2014 (Figure 3).¹ The LME methodology detects significant ($p < 0.05$) offsets between major collecting nations and ships sailing for different purposes (Chan & Huybers, 2019). Around 15% of groups remain highly significant after controlling for family-wise error rates using a Bonferroni correction (Chan et al., 2019). A file listing individual offsets and associated uncertainty estimates can be found in the supplement to this article. The identification of significant differences among nation-and-deck groups indicates that we can resolve region- and time-varying SST biases arising from groupwise offsets and varying sampling coverage of individual groups. We, therefore, can perform more detailed corrections that have not yet been accounted for in previous studies (Cowtan et al., 2018; Huang et al., 2017; Kennedy et al., 2011b).

Removing these statistically constrained offsets provides refined SST corrections at regional scales. Central estimates of adjusted SSTs show higher interannual correlations (Pearson's r) with nearby air temperatures from coastal land stations. For example, the correlation in the early 20th century increases from 0.67 to 0.85 after groupwise bucket adjustments over coastal East Asia (Chan et al., 2019). Station-based air temperatures are independently measured using more homogeneous instruments and are expected to contain fewer spatially and temporally varying systematic biases (P. Jones, 2016). Moreover, previous corrections may have missed errors associated with groupwise offsets and, therefore, underestimated SST uncertainties at regional scales. Accounting for uncertainties of groupwise offsets increases the standard error of trend estimates to more than three times at basin scales, which, despite being higher, is a more comprehensive description of our current knowledge of uncertainties.

More importantly, accounting for groupwise bucket SST offsets reconciles a long-standing data-model discrepancy regarding the spatial heterogeneity of the early-20th-century warming. Before groupwise bucket adjustments, whereas the North Pacific warmed by around 0.3°C, the North Atlantic warming exceeded 0.8°C (Figure 4; Hegerl et al., 2018). Such a big difference in warming rate, however, cannot be reproduced by any of the IPCC AR5 models given current knowledge of external forcing and internal climate variability.

The magnitude of adjustments from individual groups depends both on the magnitudes of offsets and on the spatial and temporal coverage of distinct groups. In the early 20th century, one group that determines basin-scale SST estimates is a major subset of Japanese measurements compiled under the Kobe collection, which dominated the North Pacific before World War II (Uwai & Komura, 1992). Interestingly, Japanese Kobe SSTs, compared with nearby measurements from other nations, show a drop of around $0.35 \pm 0.07^{\circ}\text{C}$ (2 SD) in the 1930s (Chan et al., 2019). Such a drop could lead to a significant underestimation of the early-20th-century warming in the Pacific. It is, therefore, crucial to understand why Japanese Kobe SSTs experienced this drop. Is it because Japan used a new type of bucket in the 1930s or did something change in the postprocessing of Japanese data? To disentangle this mystery, I combined historical approaches and physical methods.

302 Retrospectively, figuring out the cause is like detective work: it requires effort, persistence, and
303 some good luck. My first hypothesis was that Japanese sailors measured SSTs on larger ships that

¹Data associated with this analysis are available at <https://doi.org/10.7910/DVN/DXJIGA>.

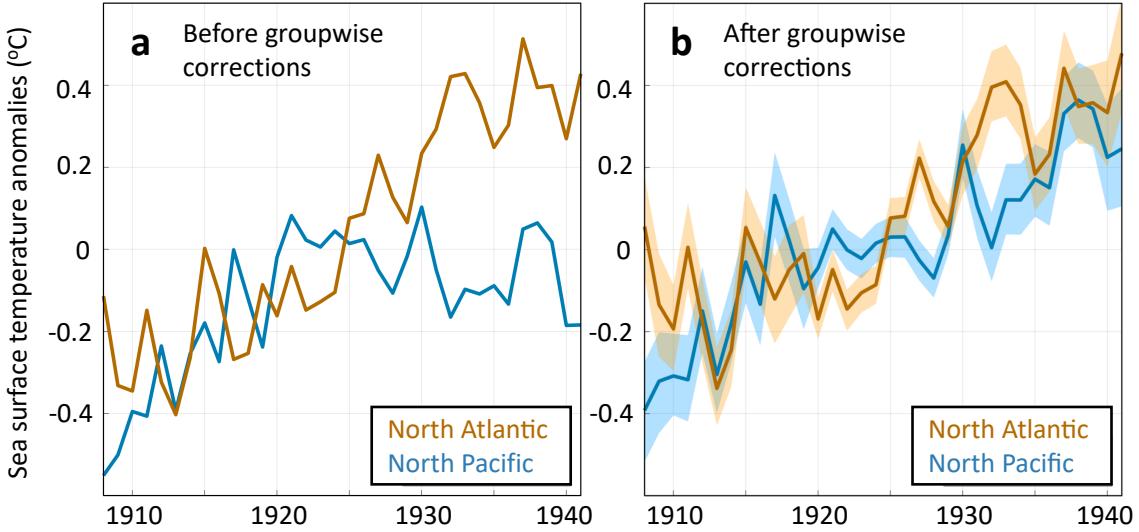


Figure 4. Basin-scale SSTs in the early twentieth century (Chan et al., 2019).

(a) Without groupwise corrections, the annual-mean SST in the North Atlantic (red, 20°N poleward) warms more than twice as fast as that in the North Pacific (blue, 20°N poleward). Shown SSTs are based on ICOADS3.0 bucket measurements with only bulk bucket corrections that do not distinguish groupwise offsets, following Kennedy et al. (2011b). SSTs are shown as anomalies relative to the 1920–1929 average of each basin. (b) As (a) but after adjusting for groupwise offsets in bucket SSTs. Uncertainties (2 SD, shadings) are for annual average SSTs in each basin and are from a 1000-member ensemble of random adjustments that perturb groupwise offsets using their error estimates from the LME analysis in keeping with covariance and spatial structures.

have higher decks. Japanese ships could have increased size as circumstances deviated from the 1922 Washington Naval Treaty that limited ship displacement, and as World War II approached, the Imperial Japanese Navy required large ships for longer voyages across the Pacific. When taking bucket measurements on higher decks, it generally takes longer to haul buckets, and SSTs tend to be collected in stronger winds, leading to colder biases. To test this hypothesis, I first went through individual ships in the Imperial Japanese Navy and mapped out the evolution of the average displacement of Japanese naval ships since the 1920s. The initial result was promising: Japanese ships increased in displacement at a rate that was approximately 40% faster than the U.S. and U.K. ships in the 1930s.²

Despite the confirmation of a rapid increase in the displacement of Japanese naval ships, follow-up analyses served to disprove the initial hypothesis, with two pieces of evidence. First, I used a thermal model of a bucket (Folland & Parker, 1995) to simulate the influence of higher ship decks on bucket bias by increasing the hauling time and the ambient wind. The bucket model indicates that water temperatures will be further biased cold by less than 0.1°C, a magnitude that is insufficient to explain the $0.35 \pm 0.07^\circ\text{C}$ drop seen in Japanese SSTs. Second, most of these large Japanese naval

²This is an unpublished result from my Ph.D. research, and it does not explain the drop in Japanese SSTs in the 1930s. The displacement data are from the World War II database (P. Chen, 2004).

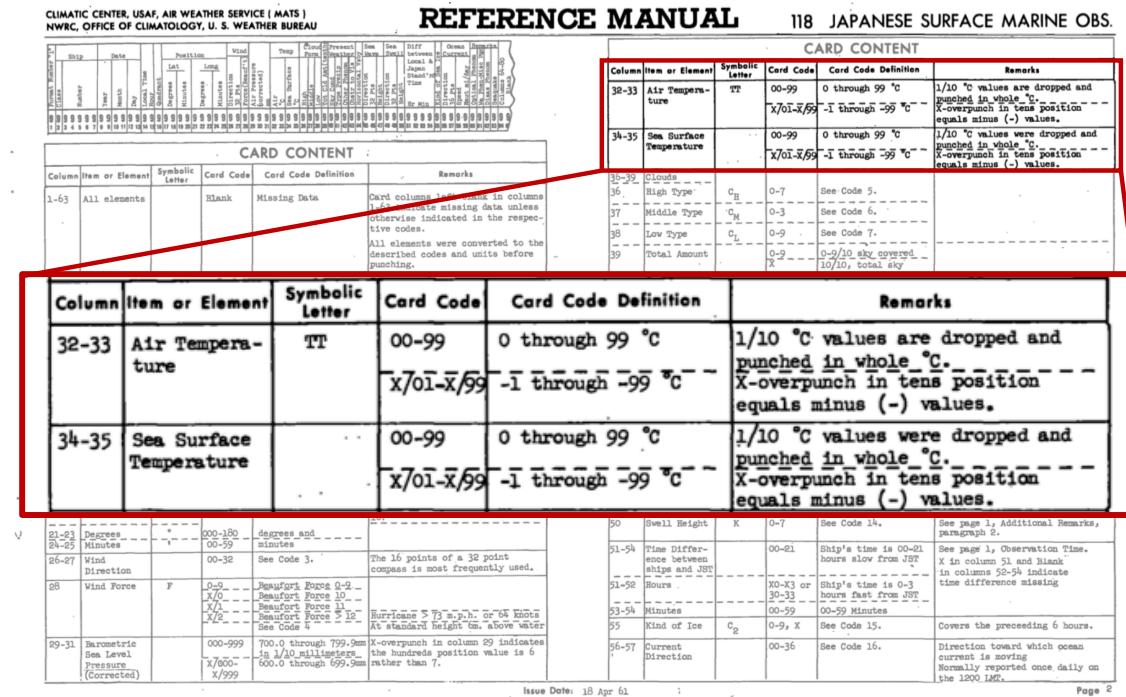


Figure 5. Image of a U.S. Air Force Weather Service document detailing how data from Japanese Kobe Collection deck 118 were digitized (Wilkinson et al., 2011). On the right-upper corner, it indicates that when both SSTs and marine air temperatures from this group were digitized, the data was floored to whole degrees Celsius, with all decimals dropped (<https://icoads.noaa.gov/reclaim/pdf/dck118.pdf>).

ships sank during battles with the U.S. Navy, and the replacement ships were small in size, yet the cold offsets remain in Japanese Kobe SSTs until the 1960s. The failure of my initial hypothesis reflects the difficulty of determining historical fact when faced with many possibilities. Fortunately, I discussed an initial manuscript with Dr. Elizabeth Kent, an expert from the U.K. National Oceanography Center who has been working on ICOADS for more than 30 years. She pointed me to an online library that documents digitization practices of many ICOADS decks. I was not aware of these documents before, and her deep expertise nudged me in the right direction. It turns out that SSTs from the Japanese Kobe collection were digitized during the recovery of logbooks and international marine data (RECLAIM) project (Wilkinson et al., 2011), with the data set divided into three subsets. The U.S. Air Force was in charge of the two parts that span from the 1930s to the early 1960s. During digitization, staff truncated Japanese temperatures and floored values to whole degrees Celsius (Figure 5), leading to the cold offsets that were prevalent in Japanese SSTs since the 1930s. Because SST measurements in ICOADS have a precision of 0.1°C , the expected cold offset due to truncation is -0.45°C when assuming that the 10th of degree digit is uniformly distributed on the values from zero to nine. The detected smaller magnitude of $-0.35 \pm 0.07^{\circ}\text{C}$ could reflect the presence of additional offsets and biases between decks.

335 Accounting for truncation errors in Japanese data, together with removing offsets in other groups,
336 reveals a more homogeneous early-20th-century warming (Figure 4b). The difference in warming
337 rates between the North Atlantic and the North Pacific decreases from 0.54°C to $0.10 \pm 0.07^{\circ}\text{C}$ (2

338 *SD*) and becomes consistent with simulations from IPCC models (Chan et al., 2019). With quantifi-
 339 cation from the statistical method and confirmation from historical documents, we now understand
 340 that the long-standing warming discrepancy is not physical but a result of data problems as simple
 341 as truncation errors. What has happened is more consistent with physics-based expectations of
 342 uniform warming associated with anthropogenic activities (Chan et al., 2019).

343 4. TRACING THE ORIGIN OF BUCKET OFFSETS USING PHYSICAL EVIDENCE

344 The reconciliation of discrepancies in the early-20th-century warming demonstrates the power of
 345 the LME method. It is also a good example of combining statistical methodologies and historical
 346 evidence to make convincing inferences and interpretations. However, despite the cause of the drop
 347 in Japanese Kobe SSTs being explained by limited historical metadata, the origin of offsets remains
 348 unclear for other groups. Due to limited metadata, the problem is approached by using features of
 349 the data. Specifically, we used a physical quantity (i.e., the diurnal cycle of SSTs) to explore the
 350 origin of groupwise offsets (Chan & Huybers, 2020b).

351 The diurnal cycle is variation among individual hours in a day, which can be easily estimated
 352 for each group independent of the LME methodology. Diurnal cycles are used because differences
 353 in diurnal cycles may reflect differences in measurement characteristics and could affect daily mean
 354 SSTs through physical processes. Water in buckets is subject to heat loss from the wind but is
 355 heated additionally by the sun during the daytime. As a result, if a bucket stays longer on the
 356 ship's deck before temperature is measured, it tends to have a higher day-night SST difference
 357 and overall a colder daily mean SST bias. Interestingly, when the amplitude of diurnal cycles and
 358 groupwise SST offsets are plotted against one another, the two quantities scale negatively for data in
 359 the 1980s and 1990s (Figure 6b; Chan & Huybers, 2020b), which strongly indicates the physicality
 360 of groupwise offsets detected by the LME methodology.

361 However, varying time on the ship's deck is not the only reason negative scalings emerge. To
 362 explore other possible origins, I extended the classic thermodynamical bucket model (Folland &
 363 Parker, 1995) to resolve bucket biases at individual local hours and simulate diurnal cycles of bucket
 364 water temperatures. Model simulations show that a negative scaling between diurnal amplitude
 365 and daily mean biases can emerge not only from varying time on deck but also from the type of
 366 bucket insulation or misclassification of ERI measurements (Chan & Huybers, 2020b). The latter
 367 arises because ERI measurements are biased warm by heat from ship engines and the ERI method
 368 samples at a depth of 5–15 m, which is less affected by diurnal variations in radiation from the
 369 sun (Carella et al., 2018). Contribution from each of these origins, however, cannot be determined
 370 from one single slope because expected slopes associated with individual origins can vary with other
 371 unknown factors, including wind and solar exposure.

372 To further trace the origin, we turned to the historical evolution of observed amplitude–offset
 373 relationships, which reveals that negative scalings first emerge in the 1930s (Figure 6c; Chan &
 374 Huybers, 2020b). Data before 1930, however, have a smaller range in both amplitudes and offsets
 375 and show no significant scalings (Figure 6a). Interestingly, the 1930s is also the advent of ERI
 376 measurements in ICOADS. Moreover, groups having the warmest offsets also have amplitudes of
 377 diurnal cycles that are smaller than physical SSTs at a depth of 20–30 cm, which is consistent with
 378 characteristics of ERI measurements. In other words, the misclassification of ERI measurements,
 379 as from buckets, provides the simplest explanation and is also most consistent with the history of
 380 changing data characteristics.

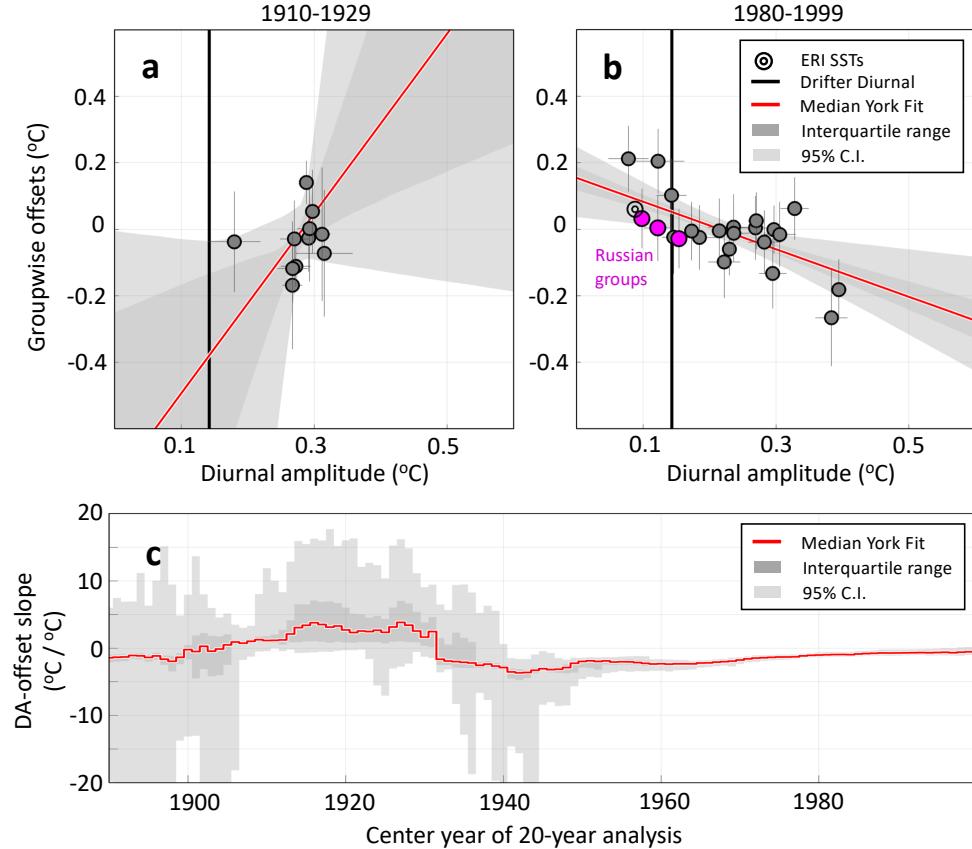


Figure 6. Groupwise SST offsets and diurnal amplitudes for groups thought to contain bucket SSTs (Chan & Huybers, 2020b). Shown results are over the tropics (20°S - 20°N) and are for 20-year periods: (a) 1910–1929 and (b) 1980–1999. Groups (markers) are assigned according to nation (two-letter abbreviations, see the legend of Figure 3) and deck number. Diurnal amplitude is quantified as the amplitude of a once-per-day sinusoid using least-squares fitting. LME analyses shown here include ERI SSTs as a single group such that groupwise bucket offsets are evaluated against ERI measurements (details in Chan & Huybers, 2020b). Slopes between diurnal amplitudes and groupwise offsets (red lines) are based only on bucket groups and are estimated using York regressions (York et al., 2004). In an update to Chan and Huybers (2020b), the uncertainty of regression slopes is estimated using a stratified bootstrapping technique that resamples the entire history of individual groups with replacement (see Appendix for more details). In panel (b), note that the regression intersects the offset and diurnal amplitude of ERI measurements (double circles), indicating that bucket groups on the warm end of the slope (such as Russian groups shown in magenta markers) could contain misclassified ERI measurements. Also note that numerous groups show a diurnal amplitude that is similar to or lower than that of drifter SSTs (vertical black lines), which is consistent with the deep sampling depth of the ERI method. (c) Evolution of the amplitude–offset relationship, which is based on an analysis that uses a 20-year window and slides annually from 1880–1899 to 1990–2009. Results are shown on the center year of each 20-year analysis. Whereas highly uncertain slopes are found before the 1930s (estimates of the 2.5% quantile can be as negative as $-56^{\circ}\text{C}/^{\circ}\text{C}$ before the 1910s), significant negative slopes are found afterward. Shown are median values (red curve), interquartile CI (dark shading), and 95% CI (light shading).

381 Analyzing diurnal cycles reveals that the record-keeping problem of incorrect metadata that
 382 mixes up bucket and ERI measurements is prevalent in data thought to come from buckets after
 383 the 1930s (Chan & Huybers, 2020b). Moreover, examining the decimal distributions of individual
 384 groups indicates that in the whole ICOADS deck, truncation is only seen in the Japanese Kobe
 385 collection. Remaining offsets, including those before the 1930s, could still arise from differences in
 386 physical processes, although, without introducing evidence from further dimensions, the trade-off
 387 among different physical factors makes it hard to attribute offsets to individual processes.

388 5. BEYOND BUCKET-ONLY SSTs — WORLD WAR II WARM ANOMALY

389 The LME methodology detects significant groupwise offsets that arose from both physical pro-
 390 cesses during data collection and record-keeping problems, including truncation and misrecorded
 391 metadata. This method can be easily extended beyond bucket SSTs and provides refined internal
 392 homogenization for measurements coming from various instruments. When applied to all ship-based
 393 SSTs in ICOADS, the LME method resolves another major data-model discrepancy involving excess
 394 warming during World War II (Chan & Huybers, 2020a).

395 Recent SST estimates feature warmer global-mean SSTs during World War II that well exceed
 396 climate-model reproductions (Figure 7a, b). Such warm anomalies are at the end of the early-20th-
 397 century warming and the beginning of the mid-20th-century hiatus and, therefore, have implications
 398 for quantifying decadal climate variations (Hansen et al., 2010; Morice et al., 2012; Vose et al., 2012),
 399 constraining uncertain aerosol forcing (Stevens, 2015), and attributing external anthropogenic forc-
 400 ing and internal climate variability in driving past climate change (Bindoff et al., 2013; Hegerl et al.,
 401 2018; G. S. Jones et al., 2013; Maher et al., 2014). Moreover, the World War II warm anomaly is the
 402 largest remaining data-model discrepancy in the global-mean surface temperature, given current
 403 knowledge of forcing and internal variability (Folland et al., 2018).

404 Such an observed warm anomaly could indicate that current models missed important physical
 405 processes. The anomaly could also arise from a 58% drop in the amount of data collected during
 406 1942–1945 (Figure 1). Another plausible explanation involves the warm anomalies reflecting incom-
 407 plete corrections of SST biases. These biases have been hypothesized to arise from a rapid increase
 408 in the number of warm-biased ERI measurements during the war (Thompson et al., 2008). How-
 409 ever, tracing the origin of biases to specific sets of SST measurements and estimating the amount
 410 of required adjustments has not previously been possible. Existing corrections are limited by not
 411 being able to resolve offsets between groups and also due to the fact that more than 80% of wartime
 412 measurements have missing method information in raw ICOADS.

413 The LME method is, however, suitable in this situation of missing metadata because it allows for
 414 groupwise quantification of data heterogeneity without the need for specifying method information
 415 (although methods can be inferred from offsets and cross-checked using diurnal amplitudes, e.g.,
 416 as in Figure 1b). In a recent work (Chan & Huybers, 2020a), we extended the estimation of
 417 groupwise offsets to all ship-based SST measurements in ICOADS, and the LME model quantifies
 418 that SSTs from some U.S. and U.K. naval ships that dominated data collections in World War II
 419 are, respectively, around 0.45°C and 0.25°C warmer than other groups before and after the war.
 420 These large and fast naval ships were likely to take ERI measurements (Kennedy et al., 2011b).
 421 Moreover, when further extended to resolve diurnal differences, the LME model detects an increase
 422 in nighttime measurements of around 0.3°C for many wartime non-ERI measurements (Chan &
 423 Huybers, 2020a). Such an increase is consistent with warm biases arising from measuring nighttime

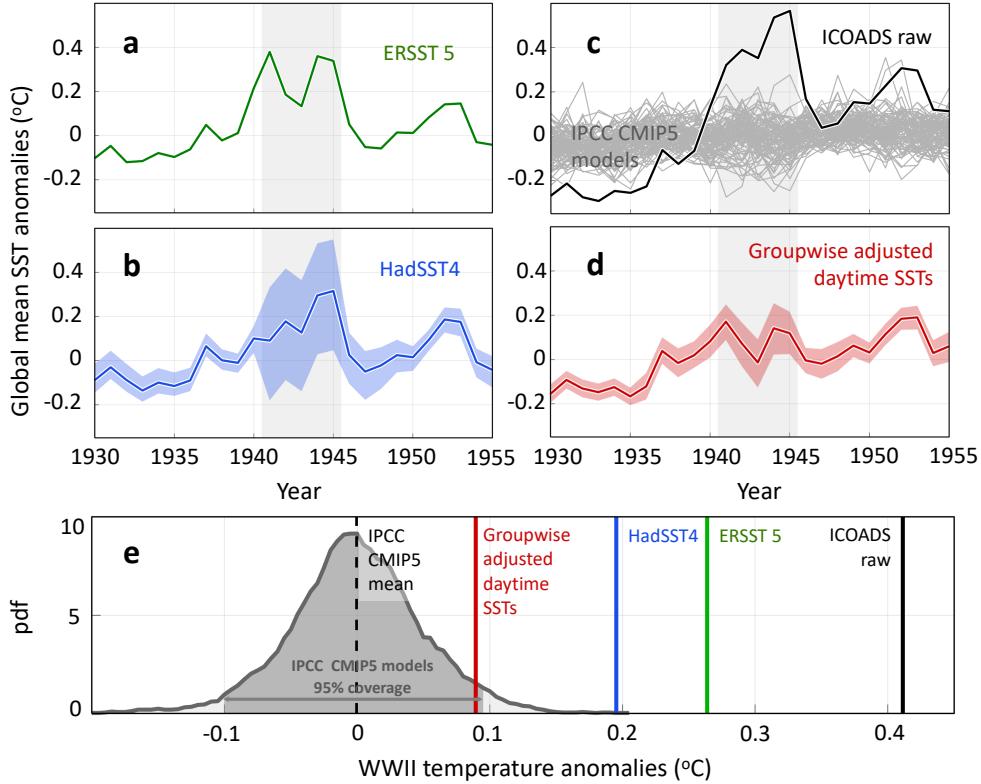


Figure 7. World War II SST anomalies in observational estimates and model simulations (Chan & Huybers, 2020a). Most recent SST estimates from (a) the U.S National Oceanic and Atmospheric Administration (ERSST5, green) and (b) U.K. Met Office (HadSST4, blue) show warm anomalies in global-mean SSTs that exceed 0.2°C during World War II. The axes for (c) and (d) are as in (a), but (c) shows raw ship-based SSTs in ICOADS (black) and (d) shows daytime ship-based SSTs after groupwise adjustments (red). Shown time series are global-averaged SST anomalies relative to the 20-year average over 1931–1940 and 1946–1955, where the global average is taken over grid boxes containing major ship tracks in the early and mid-20th century. Uncertainties are 95% CI (blue and red shading) from ensemble corrections of corresponding estimates. Also shown is an ensemble of 94 historical all-forcing simulations from 39 IPCC models (light gray curves in c; Taylor et al., 2012). In (e), the World War II warm anomaly in groupwise adjusted daytime SSTs (red) becomes consistent with the range of internal variations estimated from IPCC models (gray distribution). The World War II anomaly is quantified as the difference between averages over 1941–1945 and over the surrounding 10 years (1936–1940 and 1946–1950).

424 bucket SSTs inside ships to avoid detection, a wartime practice documented for the U.K. Navy
 425 (Folland et al., 1984).

426 The effect of groupwise offsets and nighttime bucket biases contributes to, respectively, 0.26°C
 427 (95% CI 0.15°C – 0.38°C) and 0.05°C (0.02°C – 0.08°C) warm anomalies in raw ICOADS (Chan &
 428 Huybers, 2020a). Adjustments bring the World War II warm anomaly from 0.41°C in raw ICOADS
 429 to 0.09°C (-0.01°C to 0.18°C), which becomes consistent with the $\pm 0.10^{\circ}\text{C}$ range (95% CI)
 430 of internal variability in IPCC models (Figure 7b, c; Chan & Huybers, 2020a). Groupwise adjustments

431 based on the LME methodology lead to more homogeneous spatial and temporal variations in SSTs
 432 and reconcile the largest remaining data-model discrepancy in global-mean surface temperatures.

433 Fixing problems in the WWII warm anomaly confirms the hypothesis of data biases, as suggested
 434 by Thompson et al. (2008). This piece of work also provides us with a lesson that historical data
 435 may contain not only physical but also social aspects. When interpreting historical data, especially
 436 in the context of comparing with model simulations, it is often implicitly assumed that data reflect
 437 the physical, or broadly, the scientific dimension. This assumption could be valid when a small
 438 amount of data is calibrated carefully for a specific scientific purpose. But for massive data sets
 439 that pool information from heterogeneous sources and take generations to construct, it is crucial to
 440 keep an awareness of the social dimension and the people involved in data collection and processing,
 441 especially over periods having dramatic social changes.

442 6. BEYOND SSTs — HINDCASTING OF NORTH ATLANTIC HURRICANES

443 Adjustments of groupwise offsets have been shown to improve historical SST estimates and
 444 reconcile major data-model discrepancies in surface temperature evolution. On account of the broad
 445 climatic applications of SSTs, the implication of improved SST corrections, however, is not limited
 446 to simple year-to-year variations or linear trends. Beyond surface temperatures, improvements
 447 associated with groupwise SST adjustments could also advance other fields in atmospheric and
 448 ocean sciences.

449 One example is the hindcasting of North Atlantic hurricane activities. Decadal variations in the
 450 frequency of North Atlantic hurricanes are known to depend on patterns of tropical SSTs (Vecchi,
 451 Msadek, et al., 2013; Vecchi et al., 2008). Compared with observational reconstructions of
 452 Atlantic hurricane counts, dynamical climate models prescribed with historical SSTs as boundary
 453 conditions,³ however, reproduce too few hurricanes in the late-19th century and too many in the
 454 mid-20th century ($p < 0.05$, Figure 8a, Chan et al., 2020). The inability for models to skillfully
 455 reproduce a long-term evolution of hurricane counts that are statistically consistent with observa-
 456 tional estimates erodes the credibility of future projections based on these models (Vecchi et al.,
 457 2019). Possible causes for this low reproducibility include inaccurate hurricane reconstructions
 458 (Vecchi & Knutson, 2008) and model deficiency (Zhao et al., 2009). In addition, biases in SSTs
 459 may also undermine simulations of hurricane genesis.

460 The SST value that is thought to affect North Atlantic hurricane frequency is the difference
 461 between the subtropical North Atlantic and the whole tropical ocean. This SST difference is also
 462 known as ‘relative SST’ (RSST) and is thought to influence hurricane genesis through affecting
 463 convective activities and potential energy (Vecchi, Fueglistaler, et al., 2013; Vecchi et al., 2011).
 464 Groupwise bucket SST corrections increase RSST in the late 19th century and decrease RSST in
 465 the mid-20th century. In a recent work in which I collaborated with colleagues from Princeton
 466 University (Chan et al., 2020), we incorporated groupwise bucket corrections to previous SST
 467 estimates and found that simulated hurricane counts also show increases in the late-19th century

³Although SSTs are not entirely independent of hurricanes, these simulations should have partially accounted for the effect that hurricanes lower SSTs by using observed monthly SSTs. Moreover, the response of SSTs to hurricanes is not expected to change much with time (Vecchi et al., 2019). Thus, using monthly SSTs as boundary conditions will not alter the active and inactive phases in the simulated decadal variability of hurricane counts.

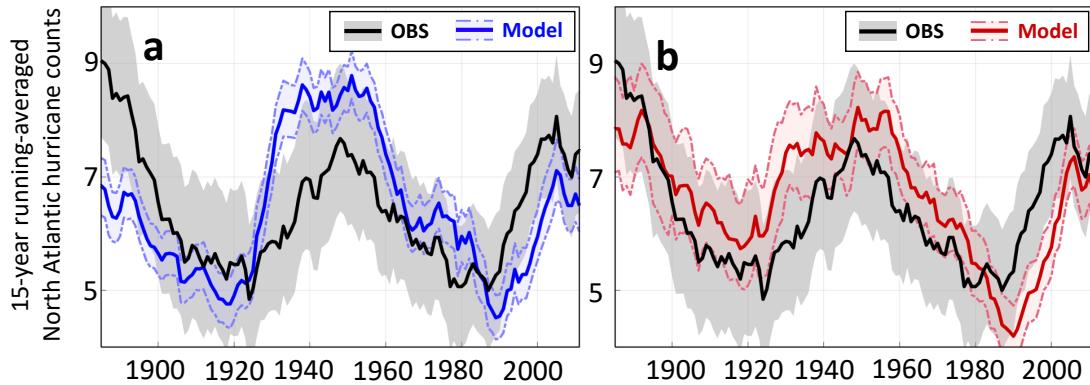


Figure 8. 15-year running-averaged North Atlantic hurricane counts in observational reconstructions and model simulations (Chan et al., 2020). (a) Simulations (blue, average of a five-member ensemble) using SST estimates without groupwise bucket corrections give significantly ($p < 0.05$) lower hurricane counts than observational estimates (black) in the late 19th century and higher counts in the mid-20th century. (b) Simulated (red, average of a five-member ensemble) and observed (black) hurricane counts become consistent using SSTs that include groupwise bucket SST corrections. Shown curves are 15-year running-averaged rather than raw integer counts because we are interested in the decadal variability of North Atlantic hurricane frequency. Uncertainties (95% CI) account for atmospheric internal variability and errors in hurricane adjustments (gray shading), atmospheric internal variability (blue shading in panel a), and atmospheric internal variability and errors in uncertain groupwise SST corrections (red shading). Distinct types of errors are assumed to be independent of one another. For estimates containing errors from two sources (i.e., black and red lines), shown uncertainties are summations of squared errors from both sources. A detailed description of the error model is in the method section of Chan et al. (2020). Note that atmospheric internal variability arises from perturbations to initial conditions, and that variability in observations is not expected to be reproduced by models because of imperfect initial conditions. It is, therefore, necessary to consider atmospheric internal variability as random error in both observation and simulation.

468 and decreases in the mid-20th century, consistent with expectations from adjusting RSST. More
 469 importantly, simulated hurricane counts become statistically consistent with independently recon-
 470 structed observational estimates of hurricane counts after accounting for groupwise SST offsets
 471 (Figure 8b; Chan et al., 2020).

472 Showing that SST biases are the dominant limiting factor for models to recover historical Atlantic
 473 hurricane counts is exciting news for both the SST and hurricane communities. The diminishing
 474 data-model discrepancy in hurricane variability provides dynamical evidence to buttress the im-
 475 proved quality of SSTs after groupwise corrections. On the other hand, the more stable relationship
 476 between observed and simulated hurricane activity increases the credibility of dynamical models in
 477 making accurate predictions of future changes in hurricane activities.

478

7. WHAT IS NEXT?

479 Correcting national and groupwise offsets improves historical SSTs and reconciles a number of
 480 data-model mismatches. Despite these significant improvements, estimates of historical SSTs are
 481 still far from perfect, and there is much scope for further improvements. In addition to recovering
 482 lost data sets and missed metadata (as suggested in, e.g., Kent et al., 2017), opportunities exist to
 483 develop new techniques and further analyze existing data sets.

484 **7.1. Internal Homogeneity to the Level of Individual Ships.** Further improvements could
 485 come from better resolving internal heterogeneity at more refined levels, such as quantifying offsets
 486 associated with distinct measurement characteristics of individual ships. Ship-level biases can lead
 487 to partially correlated errors across space and time as ships passing through different grid boxes
 488 (Kennedy, 2014). Ship-level biases were estimated to have a similar magnitude to random mea-
 489 surement errors by comparing with satellite (Kennedy et al., 2012) or observational-constrained
 490 model simulations (Kent & Berry, 2008) using data in recent decades. Ship-level biases, however,
 491 have not yet been explicitly quantified for data before the 1970s. In version 3 of the HadSST data
 492 set, biases of individual ships were assumed to follow a Gaussian distribution that has a zero mean
 493 and a standard deviation of ship-level biases (Kennedy et al., 2011a), where the ship-level standard
 494 deviation was estimated by comparing ship-based SSTs with collocated satellite retrievals since
 495 the 1990s (Kennedy et al., 2012). Uncertainties associated with ship-level biases were inferred for
 496 gridded data sets after estimating the effective number of ships in individual 5° boxes (Kennedy
 497 et al., 2011a). Such a treatment better accounted for error covariance but could not remove biases
 498 associated with individual ships. In version 4 of HadSST (the latest version), ship-based SSTs were
 499 compared against uppermost temperature measurements from ocean profiles for data after World
 500 War II, with the assumption that profile measurements are free of bias (Kennedy et al., 2019).
 501 Kennedy et al. (2019) assumed that ship-based SSTs and profile temperatures follow a multivariate
 502 normal distribution, which allows for estimating biases of gridded ship-based SST fields at the level
 503 of individual grid boxes. Such a method, to some extent, has implicitly accounted for ship-level
 504 biases.

505 The LME method appears to be a suitable approach for estimating offsets between individual
 506 ships, which further increases internal homogeneity among measurements within ICOADS. Method-
 507 logically, quantifying ship-level offsets can be realized by assigning random effects for individual
 508 ships. A systematic implementation, however, is currently limited by the quality of ship informa-
 509 tion. A total of 44% of paired bucket SSTs from 1850 to 1970 do not have ship identifiers in raw
 510 ICOADS. Moreover, around 85% of ships having IDs have no more than 25 paired measurements,
 511 which is too few for robust offset estimates because random observational errors are estimated to be
 512 0.74°C (Kennedy et al., 2012). To improve ship information, Carella et al. (2017) tried to track mea-
 513 surements with missing ship IDs and combine short tracks into longer ones. The tracking algorithm
 514 of Carella et al. (2017), however, is uncertain at ship crossings on which the LME intercomparison
 515 entirely relies. A careful examination of the suitability of these tracked ships is, therefore, required
 516 before estimating ship-level offsets using the LME method. Other opportunities may come from
 517 redigitizing early ship logs or developing more robust tracking algorithms. In addition to improving
 518 metadata of ship information, the algorithm of fitting the LME model also needs to be modified to
 519 account for hundreds of thousands of additional parameters associated with individual ships.

520 **7.2. External Consistency with Independent Instrumental Measurements and Paleo-**
 521 **proxies.** Potential for further improvements may also come from improving external consistency
 522 with independent temperature measurements or proxies. Measurements from marine air tempera-
 523 tures (Huang et al., 2017; T. M. Smith & Reynolds, 2002) and coastal land stations (Cowtan et al.,
 524 2018) have been used to quantify SST biases at the global scale. Recent studies also use subsurface
 525 temperatures from ocean profiling data to estimate SST biases at both global and regional scales
 526 since the 1940s (e.g., Huang et al., 2018; Kennedy et al., 2019). In addition to these external data
 527 sources, I would like to call attention to the potential value of ocean temperatures at deeper depth
 528 and proxies from coral reefs.

529 The deep ocean communicates with the surface through convective and diffusive processes (Geb-
 530 bie & Huybers, 2011). Deep-ocean temperatures largely reflect SST variations at high latitudes
 531 where ocean convection is most active (Gebbie & Huybers, 2011). Variability in deep-ocean tem-
 532 peratures, however, needs to be interpreted cautiously on account of possible smoothing associated
 533 with eddy diffusion, less constrained variability of ocean circulation, and a time lag between the
 534 surface and interior ocean. Alternatively, SSTs may contain information to constrain and recon-
 535 struct ocean circulation. Furthermore, similar to the SST problem, profiles that contain deep-ocean
 536 temperatures also come from various methods and nations (Meyssignac et al., 2019). Quality con-
 537 trols that involve group- or ship-level examination to profile data using the LME method appear
 538 necessary before calibrating SSTs.

539 There could also be value in paleoclimate proxies, a data source often considered to have higher
 540 noise and be less reliable than instrumental measurements. Mechanistically, heavier isotopes tend
 541 to enrich in the condensed phase due to kinetic fractionation (Urey, 1947). Heavy oxygen isotopes
 542 (e.g., O¹⁸) in coral reefs will, therefore, decrease with water temperature, providing long-term and
 543 homogeneous approximations of SSTs (e.g., Gagan et al., 2000). Pfeiffer et al. (2017) showed that
 544 proxy temperatures from coral reefs in the Indian Ocean do not show abrupt changes during World
 545 War II, which is consistent with our groupwise corrections. Coral reefs have the benefit of a broader
 546 coverage in tropical oceans, including the eastern Pacific, which is not frequently sampled by ships.
 547 Caution is required when interpreting oxygen isotopes. In certain regions that have abundant
 548 rainfall, such as the intertropical convergence zone, the concentration of O¹⁸ in rainwater (thus
 549 seawater and coral reefs) decreases strongly with increasing rainfall, which could mask temperature
 550 signals (Gagan et al., 2000; Lee & Fung, 2008; Pfeiffer et al., 2017).

551 **7.3. New Mapping Techniques.** In addition to correcting SST biases, an equally important
 552 problem in SST reconstructions involves mapping and infilling grids without observations to have
 553 global coverage. Unlike typical kernel functions that have decaying covariance with increasing
 554 displacement, kernels preferred in climate sciences should account for covariance associated with
 555 large-scale variations in ocean and atmosphere (e.g., El Niño and Southern Oscillation). Most
 556 previous SST estimates use principal component analysis (PCA) to learn patterns of covariance
 557 from satellite observations since the 1980s and assume stationarity (e.g., Hirahara et al., 2014;
 558 Huang et al., 2017; Rayner et al., 2003), even though satellite retrievals show that the details of
 559 the SST patterns are distinct across El Niño events in the past 40 years (Timmermann et al.,
 560 2018). Variational Bayesian methods have been proposed to learn patterns from ship-based SSTs
 561 that have a longer history (Ilin & Kaplan, 2009). Reconstructions from this method, however,
 562 contain patterns of ship tracks that we do not expect to exist in physical SSTs (Kennedy et al.,
 563 2013). The most recent advance involves using inpainting techniques in artificial intelligence and

learning patterns from climate model simulations (Kadow et al., 2020), but whether such a method gives reliable error estimates remains questionable. As a result, a statistically rigorous mapping technique that accounts for physical climatic covariance assuming potentially nonstationary spatial covariance is necessary for reconstructing past climate variability and budgeting uncertainties.

7.4. A Unified Statistical Framework. Quantification and correction of SST biases are often treated as separate steps from mapping and infilling for existing SST estimates (e.g., Hirahara et al., 2014; Huang et al., 2017; Rayner et al., 2003). Moreover, the mapping procedure is often further divided into separately performed substeps according to spatial scales of infilling. Such frameworks, however, make it difficult to budget and synthesize uncertainties in SST estimates arising from distinct analyzing steps.

Developing a holistic statistical framework that unifies distinct steps appears to be a solution. Such a framework should incorporate random error, systematic biases, and physical variations of global SSTs simultaneously with fully resolved covariance. Moreover, on account of potentially large uncertainties in estimates of physical SST covariance and observational errors and biases, a Bayesian method may be a more suitable framework for comprehensive quantification of SST uncertainties. Ideally, this framework should also take in biases and uncertainties we have learned from existing works and other pieces of useful information from independent external measurements. The Bayesian framework developed by Tingley and Huybers (2010) appears to be a valid starting point.

8. CONCLUSION

Understanding the history of data is crucial, and one needs to be particularly careful when using data outside their historical context. Most historical SSTs were not collected for studying climate change. These measurements contain various biases due to distinct physical and historical reasons during data collection and postprocessing. Although these SSTs have irreplaceable value for understanding past climate variations, they are undercalibrated to have sufficient accuracy for climatic use. Contrary to complicated biases, previous SST corrections that assumed homogeneous bias structures appear oversimplified, which motivated our scrutinizing historical data. Insufficient corrections have led to substantial remaining errors that result in discrepancies between observations and model simulations of the historical period. When data and models disagree, one can almost always adjust the models so that they better reproduce the data, but being aware of the underlying assumption that data reflect reality and being skeptical about data quality appears to be a good practice.

My Ph.D. work is one step forward toward better resolving complicated SST biases and toward a more accurate depiction of the past climate. Our LME method is ignorant of the existence of data-model discrepancies, but accounting for data heterogeneity among nations and groups of collectors reconciles several data-model discrepancies and provides a more comprehensive estimate of SST uncertainties. These improvements will not be achieved and consolidated without combining evidence from statistical, physical, and historical aspects. Even though not assumed or built-in, the updated SST estimates show simpler spatial and temporal variations and are more in line with expected patterns of warming. Bringing observational estimates into accord with our current knowledge of forcing, climate sensitivity, and internal variability leads to greater confidence in future predictions of global warming made by climate models.

606 **Disclosure Statement.** D.C. is supported by a grant from the Harvard Global Institute and has
607 no conflicts of interest to disclose.

608 **Acknowledgments**

609 I acknowledge the associate editor, the dataviz editor, and two anonymous reviewers for their
610 constructive comments that greatly improved the quality of this review. I thank my advisor Peter
611 Huybers for advice and helpful discussions throughout my Ph.D. study. I also thank Elizabeth
612 C. Kent, Gabriel A. Vecchi, Wenchang Yang, and David I. Berry for their collaboration and discus-
613 sions on specific sections of my PhD project. This review was initiated at the 2020 Harvard Horizons
614 program, and I acknowledge Xiaoli Meng, Edward J. Hall, Pamela Pollock, Hardeep Dhillon, the
615 six other fellow scholars, and the staff at the Derek Bok Center for discussions on an initial version
616 of Sections 1–3 of this review.

617

- 618 Bindoff, N. L., Stott, P. A., AchutaRao, K. M., Allen, M. R., Gillett, N., Gutzler, D., Hansingo,
619 K., Hegerl, G., Hu, Y., Jain, S., et al. (2013). Detection and attribution of climate change:
620 from global to regional. *Climate Change 2013 – the Physical Science Basis: Working Group*
621 *I Contribution to the Fifth Assessment Report of the Intergovernmental Panel on Climate*
622 *Change*. Cambridge University Press. <https://doi.org/10.1017/CBO9781107415324.022>
- 623 Carella, G., Kennedy, J. J., Berry, D., Hirahara, S., Merchant, C. J., Morak-Bozzo, S., & Kent,
624 E. C. (2018). Estimating sea surface temperature measurement methods using characteristic
625 differences in the diurnal cycle. *Geophysical Research Letters*, 45(1), 363–371. <https://doi.org/10.1002/2017GL076475>
- 626 Carella, G., Kent, E. C., & Berry, D. I. (2017). A probabilistic approach to ship voyage reconstruc-
627 tion in ICOADS. *International Journal of Climatology*, 37(5), 2233–2247. <https://doi.org/10.1002/joc.4492>
- 628 Chan, D., & Huybers, P. (2019). Systematic differences in bucket sea surface temperature measure-
629 ments among nations identified using a linear-mixed-effect method. *Journal of Climate*,
630 32(9), 2569–2589. <https://doi.org/10.1175/JCLI-D-18-0562.1>
- 631 Chan, D., & Huybers, P. (2020a). Identifying and correcting the World War 2 warm anomaly in sea
632 surface temperature measurements. *EarthArXiv preprint*. <https://doi.org/10.31223/osf.io/ju26e>
- 633 Chan, D., & Huybers, P. (2020b). Systematic differences in bucket sea surface temperatures caused
634 by misclassification of engine room intake measurements. *Journal of Climate*, 33(18), 7735–
635 7753. <https://doi.org/10.1175/JCLI-D-19-0972.1>
- 636 Chan, D., Kent, E. C., Berry, D. I., & Huybers, P. (2019). Correcting datasets leads to more
637 homogeneous early-twentieth-century sea surface warming. *Nature*, 571(7765), 393. <https://doi.org/10.1038/s41586-019-1349-2>
- 638 Chan, D., Vecchi, G. A., Yang, W., & Huybers, P. (2020). Correcting sea surface temperatures
639 improves simulations of historical hurricane activity. *EarthArXiv preprint*. <https://doi.org/10.31223/osf.io/huz73>
- 640 Chen, P. (2004). World War II Database.
- 641 Chen, X., & Tung, K.-K. (2014). Varying planetary heat sink led to global-warming slowdown and
642 acceleration. *Science*, 345(6199), 897–903. <https://doi.org/10.1126/science.1254937>

- 648 Cowtan, K., Rohde, R., & Hausfather, Z. (2018). Evaluating biases in sea surface temperature
649 records using coastal weather stations. *Quarterly Journal of the Royal Meteorological Society*, 144(712), 670–681. <https://doi.org/10.1002/qj.3235>
- 651 Easterling, D. R., & Wehner, M. F. (2009). Is the climate warming or cooling? *Geophysical Research
652 Letters*, 36(8). <https://doi.org/10.1029/2009GL037810>
- 653 Folland, C. K., Boucher, O., Colman, A., & Parker, D. E. (2018). Causes of irregularities in trends
654 of global mean surface temperature since the late 19th century. *Science Advances*, 4(6),
655 EAAO5297. <https://doi.org/10.1126/sciadv.aa05297>
- 656 Folland, C. K., & Parker, D. (1995). Correction of instrumental biases in historical sea surface
657 temperature data. *Quarterly Journal of the Royal Meteorological Society*, 121(522), 319–
658 367. <https://doi.org/10.1002/qj.49712152206>
- 659 Folland, C. K., Parker, D., & Kates, F. (1984). Worldwide marine temperature fluctuations 1856–
660 1981. *Nature*, 310(5979), 670–673. <https://doi.org/10.1038/310670a0>
- 661 Freeman, E., Woodruff, S. D., Worley, S. J., Lubker, S. J., Kent, E. C., Angel, W. E., Berry, D. I.,
662 Brohan, P., Eastman, R., Gates, L., et al. (2017). ICOADS Release 3.0: a major update
663 to the historical marine climate record. *International Journal of Climatology*, 37(5), 2211–
664 2232. <https://doi.org/10.1002/joc.4775>
- 665 Fyfe, J. C., Gillett, N. P., & Zwiers, F. W. (2013). Overestimated global warming over the past 20
666 years. *Nature Climate Change*, 3(9), 767–769. <https://doi.org/10.1038/nclimate1972>
- 667 Gagan, M., Ayliffe, L., Beck, J. W., Cole, J., Druffel, E., Dunbar, R. B., & Schrag, D. (2000). New
668 views of tropical paleoclimates from corals. *Quaternary Science Reviews*, 19(1-5), 45–64.
669 [https://doi.org/10.1016/S0277-3791\(99\)00054-2](https://doi.org/10.1016/S0277-3791(99)00054-2)
- 670 Gates, W. L., Boyle, J. S., Covey, C., Dease, C. G., Doutriaux, C. M., Drach, R. S., Fiorino, M.,
671 Gleckler, P. J., Hnilo, J. J., Marlais, S. M., et al. (1999). An overview of the results
672 of the Atmospheric Model Intercomparison Project (AMIP I). *Bulletin of the American
673 Meteorological Society*, 80(1), 29–56. [https://doi.org/10.1175/1520-0477\(1999\)080<0029:
AOOTRO>2.0.CO;2](https://doi.org/10.1175/1520-0477(1999)080<0029:
674 AOOTRO>2.0.CO;2)
- 675 Gebbie, G., & Huybers, P. (2011). How is the ocean filled? *Geophysical Research Letters*, 38(6).
676 <https://doi.org/10.1029/2011GL046769>
- 677 Hansen, J., Ruedy, R., Sato, M., & Lo, K. (2010). Global surface temperature change. *Reviews of
678 Geophysics*, 48(4). <https://doi.org/10.1029/2010RG000345>
- 679 Hartmann, D. L., Tank, A. M. K., Rusticucci, M., Alexander, L. V., Brönnimann, S., Charabi,
680 Y. A. R., Dentener, F. J., Dlugokencky, E. J., Easterling, D. R., Kaplan, A., et al. (2013).
681 Observations: atmosphere and surface. *Climate Change 2013 – the Physical Science Ba-
682 sis: Working Group I Contribution to the Fifth Assessment Report of the Intergovernmen-
683 tal Panel on Climate Change*. Cambridge University Press. [https://doi.org/10.1017/
CBO9781107415324.008](https://doi.org/10.1017/
684 CBO9781107415324.008)
- 685 Hausfather, Z., Cowtan, K., Clarke, D. C., Jacobs, P., Richardson, M., & Rohde, R. (2017). Assessing
686 recent warming using instrumentally homogeneous sea surface temperature records. *Science
687 Advances*, 3(1), e1601207. <https://doi.org/10.1126/sciadv.1601207>
- 688 Hegerl, G. C., Brönnimann, S., Schurer, A., & Cowan, T. (2018). The early 20th century warming:
689 anomalies, causes, and consequences. *Wiley Interdisciplinary Reviews: Climate Change*,
690 9(4), e522. <https://doi.org/10.1002/wcc.522>
- 691 Hervey, R. V. (2014). Meteorological and oceanographic data collected from the National Data Buoy
692 Center Coastal-Marine Automated Network (C-MAN) and moored (weather) buoys during

- 693 2014-03 (NODC Accession 0117682). Version 1.1. National Oceanographic Data Center,
694 NOAA. Dataset. <https://data.nodc.noaa.gov/cgi-bin/iso?id=gov.noaa.nodc:0117682>
- 695 Hirahara, S., Ishii, M., & Fukuda, Y. (2014). Centennial-scale sea surface temperature analysis and
696 its uncertainty. *Journal of Climate*, 27(1), 57–75. <https://doi.org/10.1175/JCLI-D-12-00837.1>
- 697 Huang, B., Angel, W., Boyer, T., Cheng, L., Chepurin, G., Freeman, E., Liu, C., & Zhang, H.-M.
698 (2018). Evaluating SST analyses with independent ocean profile observations. *Journal of
700 Climate*, 31(13), 5015–5030. <https://doi.org/10.1175/JCLI-D-17-0824.1>
- 701 Huang, B., Banzon, V. F., Freeman, E., Lawrimore, J., Liu, W., Peterson, T. C., Smith, T. M.,
702 Thorne, P. W., Woodruff, S. D., & Zhang, H.-M. (2015). Extended reconstructed sea surface
703 temperature version 4 (ERSST. v4). Part I: Upgrades and intercomparisons. *Journal of
704 Climate*, 28(3), 911–930. <https://doi.org/10.1175/JCLI-D-14-00006.1>
- 705 Huang, B., Thorne, P. W., Banzon, V. F., Boyer, T., Chepurin, G., Lawrimore, J. H., Menne,
706 M. J., Smith, T. M., Vose, R. S., & Zhang, H.-M. (2017). Extended reconstructed sea
707 surface temperature, version 5 (ERSSSTv5): Upgrades, validations, and intercomparisons.
708 *Journal of Climate*, 30(20), 8179–8205. <https://doi.org/10.1175/JCLI-D-16-0836.1>
- 709 Ilin, A., & Kaplan, A. (2009). Bayesian PCA for reconstruction of historical sea surface tempera-
710 tures. *2009 International Joint Conference on Neural Networks*, 1322–1327. <https://doi.org/10.1109/IJCNN.2009.5178744>
- 711 Jones, G. S., Stott, P. A., & Christidis, N. (2013). Attribution of observed historical near-surface
712 temperature variations to anthropogenic and natural causes using CMIP5 simulations.
713 *Journal of Geophysical Research: Atmospheres*, 118(10), 4001–4024. <https://doi.org/10.1002/jgrd.50239>
- 714 Jones, P. (2016). The reliability of global and hemispheric surface temperature records. *Advances
715 in Atmospheric Sciences*, 33(3), 269–282. <https://doi.org/10.1007/s00376-015-5194-4>
- 716 Kadow, C., Hall, D. M., & Ulbrich, U. (2020). Artificial intelligence reconstructs missing climate
717 information. *Nature Geoscience*, 1–6. <https://doi.org/10.1038/s41561-020-0582-5>
- 718 Karl, T. R., Arguez, A., Huang, B., Lawrimore, J. H., McMahon, J. R., Menne, M. J., Peterson,
719 T. C., Vose, R. S., & Zhang, H.-M. (2015). Possible artifacts of data biases in the recent
720 global surface warming hiatus. *Science*, 348(6242), 1469–1472. <https://doi.org/10.1126/science.aaa5632>
- 721 Kennedy, J. J. (2014). A review of uncertainty in in situ measurements and data sets of sea surface
722 temperature. *Reviews of Geophysics*, 52(1), 1–32. <https://doi.org/10.1002/2013RG000434>
- 723 Kennedy, J. J., Brohan, P., & Tett, S. (2007). A global climatology of the diurnal variations in sea-
724 surface temperature and implications for MSU temperature trends. *Geophysical Research
725 Letters*, 34(5). <https://doi.org/10.1029/2006GL028920>
- 726 Kennedy, J. J., Rayner, N., Atkinson, C., & Killick, R. (2019). An ensemble data set of sea surface
727 temperature change from 1850: The Met Office Hadley Centre HadSST. 4.0.0.0 data set.
728 *Journal of Geophysical Research: Atmospheres*, 124(14), 7719–7763. <https://doi.org/10.1029/2018JD029867>
- 729 Kennedy, J. J., Rayner, N., Smith, R., Parker, D., & Saunby, M. (2011a). Reassessing biases and
730 other uncertainties in sea surface temperature observations measured in situ since 1850: 1.
731 Measurement and sampling uncertainties. *Journal of Geophysical Research: Atmospheres*,
732 116(D14). <https://doi.org/10.1029/2010JD015218>
- 733

- 737 Kennedy, J. J., Rayner, N., Smith, R., Parker, D., & Saunby, M. (2011b). Reassessing biases and
738 other uncertainties in sea surface temperature observations measured in situ since 1850:
739 2. Biases and homogenization. *Journal of Geophysical Research: Atmospheres*, 116(D14).
740 <https://doi.org/10.1029/2010JD015220>
- 741 Kennedy, J. J., Rayner, N., Saunby, M., & Millington, S. (2013). Bringing together measurements of
742 sea surface temperature made in situ with retrievals from satellite instruments to create a
743 globally complete analysis for 1850 onwards, HadISST2. *EGU General Assembly Conference
744 Abstracts*, 15.
- 745 Kennedy, J. J., Smith, R. O., & Rayner, N. A. (2012). Using AATSR data to assess the qual-
746 ity of in situ sea-surface temperature observations for climate studies. *Remote Sensing of
747 Environment*, 116, 79–92. <https://doi.org/10.1016/j.rse.2010.11.021>
- 748 Kent, E. C., & Berry, D. I. (2008). Assessment of the marine observing system (ASMOS): final
749 report.
- 750 Kent, E. C., Kennedy, J. J., Berry, D. I., & Smith, R. O. (2010). Effects of instrumentation changes
751 on sea surface temperature measured in situ. *Wiley Interdisciplinary Reviews: Climate
752 Change*, 1(5), 718–728. <https://doi.org/10.1002/wcc.55>
- 753 Kent, E. C., Kennedy, J. J., Smith, T. M., Hirahara, S., Huang, B., Kaplan, A., Parker, D. E.,
754 Atkinson, C. P., Berry, D. I., Carella, G., et al. (2017). A call for new approaches to
755 quantifying biases in observations of sea surface temperature. *Bulletin of the American
756 Meteorological Society*, 98(8), 1601–1616. <https://doi.org/10.1175/BAMS-D-15-00251.1>
- 757 Kent, E. C., & Taylor, P. K. (2006). Toward estimating climatic trends in SST. Part I: Methods
758 of measurement. *Journal of Atmospheric and Oceanic Technology*, 23(3), 464–475. <https://doi.org/10.1175/JTECH1843.1>
- 760 Kent, E. C., Woodruff, S. D., & Berry, D. I. (2007). Metadata from WMO publication no. 47 and an
761 assessment of voluntary observing ship observation heights in ICOADS. *Journal of Atmo-
762 spheric and Oceanic Technology*, 24(2), 214–234. <https://doi.org/10.1175/JTECH1949.1>
- 763 Kosaka, Y., & Xie, S.-P. (2013). Recent global-warming hiatus tied to equatorial Pacific surface
764 cooling. *Nature*, 501(7467), 403–407. <https://doi.org/10.1038/nature12534>
- 765 Lapen, A. G. (1998). Arrhenius and the Intergovernmental Panel on Climate Change. *Eos, Trans-
766 actions American Geophysical Union*, 79(23), 271–271. <https://doi.org/10.1029/98EO00206>
- 767 Lee, J.-E., & Fung, I. (2008). “Amount effect” of water isotopes and quantitative analysis of post-
768 condensation processes. *Hydrological Processes: An International Journal*, 22(1), 1–8. <https://doi.org/10.1002/hyp.6637>
- 770 Maher, N., Gupta, A. S., & England, M. H. (2014). Drivers of decadal hiatus periods in the 20th
771 and 21st centuries. *Geophysical Research Letters*, 41(16), 5978–5986. <https://doi.org/10.1002/2014GL060527>
- 773 Medhaug, I., Stolpe, M. B., Fischer, E. M., & Knutti, R. (2017). Reconciling controversies about the
774 global warming hiatus. *Nature*, 545(7652), 41–47. <https://doi.org/10.1038/nature22315>
- 775 Meyssignac, B., Boyer, T., Zhao, Z., Hakuba, M. Z., Landerer, F. W., Stammer, D., Köhl, A., Kato,
776 S., L’Ecuyer, T., Ablain, M., et al. (2019). Measuring global ocean heat content to estimate
777 the Earth energy imbalance. *Frontiers in Marine Science*, 6, 432. <https://doi.org/10.3389/fmars.2019.00432>
- 779 Morice, C. P., Kennedy, J. J., Rayner, N. A., & Jones, P. D. (2012). Quantifying uncertainties
780 in global and regional temperature change using an ensemble of observational estimates:

- 781 The HadCRUT4 data set. *Journal of Geophysical Research: Atmospheres*, 117(D8). <https://doi.org/10.1029/2011JD017187>
- 782
- 783 Pfeiffer, M., Zinke, J., Dullo, W.-C., Garbe-Schönberg, D., Latif, M., & Weber, M. (2017). Indian Ocean corals reveal crucial role of World War II bias for twentieth century warming estimates. *Scientific Reports*, 7(1), 1–11. <https://doi.org/10.1038/s41598-017-14352-6>
- 784
- 785
- 786 Rayner, N., Parker, D. E., Horton, E., Folland, C. K., Alexander, L., Rowell, D., Kent, E. C., & Kaplan, A. (2003). Global analyses of sea surface temperature, sea ice, and night marine air temperature since the late nineteenth century. *Journal of Geophysical Research: Atmospheres*, 108(D14). <https://doi.org/10.1029/2002JD002670>
- 787
- 788
- 789
- 790 Roemmich, D., Boebel, O., Desaubies, Y., Freeland, H., King, B., LeTraon, P.-Y., Molinari, R., Owens, B., Riser, S., Send, U., et al. (1999). Argo: The global array of profiling floats. *CLIVAR Exchanges*, 13(4 (3)), 4–5.
- 791
- 792
- 793 Shankar, P. (2020). Tutorial overview of simple, stratified, and parametric bootstrapping. *Engineering Reports*, 2(1), e12096. <https://doi.org/https://doi.org/10.1002/eng2.12096>
- 794
- 795 Smith, T. M., & Reynolds, R. W. (2002). Bias corrections for historical sea surface temperatures based on marine air temperatures. *Journal of Climate*, 15(1), 73–87. [https://doi.org/10.1175/1520-0442\(2002\)015<0073:BCFHSS>2.0.CO;2](https://doi.org/10.1175/1520-0442(2002)015<0073:BCFHSS>2.0.CO;2)
- 796
- 797
- 798 Smith, T. M., Reynolds, R. W., Peterson, T. C., & Lawrimore, J. (2008). Improvements to NOAA's historical merged land-ocean surface temperature analysis (1880–2006). *Journal of Climate*, 21(10), 2283–2296. <https://doi.org/10.1175/2007JCLI2100.1>
- 800
- 801 Stevens, B. (2015). Rethinking the lower bound on aerosol radiative forcing. *Journal of Climate*, 28(12), 4794–4819. <https://doi.org/10.1175/JCLI-D-14-00656.1>
- 802
- 803 Taylor, K. E., Stouffer, R. J., & Meehl, G. A. (2012). An overview of CMIP5 and the experiment design. *Bulletin of the American Meteorological Society*, 93(4), 485–498. <https://doi.org/10.1175/BAMS-D-11-00094.1>
- 804
- 805
- 806 Thompson, D. W., Kennedy, J. J., Wallace, J. M., & Jones, P. D. (2008). A large discontinuity in the mid-twentieth century in observed global-mean surface temperature. *Nature*, 453(7195), 646–649. <https://doi.org/10.1038/nature06982>
- 807
- 808
- 809 Timmermann, A., An, S.-I., Kug, J.-S., Jin, F.-F., Cai, W., Capotondi, A., Cobb, K. M., Lengaigne, M., McPhaden, M. J., Stuecker, M. F., et al. (2018). El Niño–southern oscillation complexity. *Nature*, 559(7715), 535–545. <https://doi.org/10.1038/s41586-018-0252-6>
- 810
- 811
- 812 Tingley, M. P., & Huybers, P. (2010). A Bayesian algorithm for reconstructing climate anomalies in space and time. Part I: Development and applications to paleoclimate reconstruction problems. *Journal of Climate*, 23(10), 2759–2781. <https://doi.org/10.1175/2009JCLI3015.1>
- 813
- 814
- 815 Urey, H. C. (1947). The thermodynamic properties of isotopic substances. *Journal of the Chemical Society (Resumed)*, 562–581. <https://doi.org/10.1039/jr9470000562>
- 816
- 817 Uwai, T., & Komura, K. (1992). The collection of historical ships' data in Kobe marine observatory. *Proceedings of the International COADS Workshop, Boulder, CO, USA*, 13–15.
- 818
- 819 Vecchi, G. A., Delworth, T. L., Murakami, H., Underwood, S. D., Wittenberg, A. T., Zeng, F., Zhang, W., Baldwin, J. W., Bhatia, K. T., Cooke, W., et al. (2019). Tropical cyclone sensitivities to CO₂ doubling: roles of atmospheric resolution, synoptic variability and background climate changes. *Climate Dynamics*, 53(9–10), 5999–6033. <https://doi.org/10.1007/s00382-019-04913-y>
- 820
- 821
- 822
- 823

- 824 Vecchi, G. A., Fueglistaler, S., Held, I. M., Knutson, T. R., & Zhao, M. (2013). Impacts of atmo-
825 spheric temperature trends on tropical cyclone activity. *Journal of Climate*, 26(11), 3877–
826 3891. <https://doi.org/10.1175/JCLI-D-12-00503.1>
- 827 Vecchi, G. A., & Knutson, T. R. (2008). On estimates of historical North Atlantic tropical cyclone
828 activity. *Journal of Climate*, 21(14), 3580–3600. <https://doi.org/10.1175/2008JCLI2178.1>
- 829 Vecchi, G. A., Msadek, R., Anderson, W., Chang, Y.-S., Delworth, T., Dixon, K., Gudgel, R.,
830 Rosati, A., Stern, B., Villarini, G., Wittenberg, A., Yang, X., Zeng, F., Zhang, R., &
831 Zhang, S. (2013). Multiyear predictions of North Atlantic hurricane frequency: Promise
832 and limitations. *Journal of Climate*, 26(15), 5337–5357. <https://doi.org/10.1175/JCLI-D-12-00464.1>
- 833 Vecchi, G. A., Swanson, K. L., & Soden, B. J. (2008). Whither hurricane activity? *Science*, 687–689.
834 <https://doi.org/10.1126/science.1164396>
- 835 Vecchi, G. A., Zhao, M., Wang, H., Villarini, G., Rosati, A., Kumar, A., Held, I. M., & Gudgel,
836 R. (2011). Statistical–dynamical predictions of seasonal North Atlantic hurricane activity.
837 *Monthly Weather Review*, 139(4), 1070–1082. <https://doi.org/10.1175/2010MWR3499.1>
- 838 Vose, R. S., Arndt, D., Banzon, V. F., Easterling, D. R., Gleason, B., Huang, B., Kearns, E., Law-
839 rimore, J. H., Menne, M. J., Peterson, T. C., et al. (2012). NOAA’s merged land–ocean sur-
840 face temperature analysis. *Bulletin of the American Meteorological Society*, 93(11), 1677–
841 1685. <https://doi.org/10.1175/BAMS-D-11-00241.1>
- 842 Wilkinson, C., Woodruff, S. D., Brohan, P., Claesson, S., Freeman, E., Koek, F., Lubker, S. J.,
843 Marzin, C., & Wheeler, D. (2011). Recovery of logbooks and international marine data:
844 the RECLAIM project. *International Journal of Climatology*, 31(7), 968–979. <https://doi.org/10.1002/joc.2102>
- 845 Woodruff, S. D., Diaz, H. F., Elms, J. D., & Worley, S. J. (1998). COADS Release 2 data and meta-
846 data enhancements for improvements of marine surface flux fields. *Physics and Chemistry
847 of the Earth*, 23(5–6), 517–526. [https://doi.org/10.1016/S0079-1946\(98\)00064-0](https://doi.org/10.1016/S0079-1946(98)00064-0)
- 848 Woodruff, S. D., Slutz, R. J., Jenne, R. L., & Steurer, P. M. (1987). A comprehensive ocean-
849 atmosphere data set. *Bulletin of the American Meteorological Society*, 68(10), 1239–1250.
850 [https://doi.org/10.1175/1520-0477\(1987\)068<1239:ACOADS>2.0.CO;2](https://doi.org/10.1175/1520-0477(1987)068<1239:ACOADS>2.0.CO;2)
- 851 Woodruff, S. D., Worley, S. J., Lubker, S. J., Ji, Z., Eric Freeman, J., Berry, D. I., Brohan, P.,
852 Kent, E. C., Reynolds, R. W., Smith, S. R., et al. (2011). ICOADS Release 2.5: exten-
853 sions and enhancements to the surface marine meteorological archive. *International journal of
854 climatology*, 31(7), 951–967. <https://doi.org/10.1002/joc.2103>
- 855 Worley, S. J., Woodruff, S. D., Reynolds, R. W., Lubker, S. J., & Lott, N. (2005). ICOADS re-
856 lease 2.1 data and products. *International Journal of Climatology: A Journal of the Royal
857 Meteorological Society*, 25(7), 823–842. <https://doi.org/10.1002/joc.1166>
- 858 Yeh, S.-W., Kug, J.-S., Dewitte, B., Kwon, M.-H., Kirtman, B. P., & Jin, F.-F. (2009). El Niño in
859 a changing climate. *Nature*, 461(7263), 511. <https://doi.org/10.1038/nature08316>
- 860 York, D., Evensen, N. M., Martinez, M. L., & De Basabe Delgado, J. (2004). Unified equations
861 for the slope, intercept, and standard errors of the best straight line. *American Journal of
862 Physics*, 72(3), 367–375. <https://doi.org/10.1119/1.1632486>
- 863 Zhao, M., Held, I. M., Lin, S.-J., & Vecchi, G. A. (2009). Simulations of global hurricane climatology,
864 interannual variability, and response to global warming using a 50-km resolution GCM.
865 *Journal of Climate*, 22(24), 6653–6678. <https://doi.org/10.1175/2009JCLI3049.1>

868

APPENDIX

869 **Stratified Bootstrapping for Estimating Uncertainties in the Evolution of the Amplitude–
870 Offset Relationship.** Similar to Figure 6c, Chan and Huybers (2020b) also investigated the evo-
871 lution of the amplitude–offset relationship using a sliding 20-year window. The bootstrapping in
872 Chan and Huybers (2020b) resamples available groups in each 20-year analysis independently, which
873 gives a reasonable uncertainty estimate within each 20-year analysis but may not be optimal for
874 intercomparing slopes across 20-year analyses. Here, I supplement earlier estimates by resampling
875 groups with their entire history of diurnal amplitudes and groupwise offsets, which also estimates the
876 path-wise uncertainty. Moreover, to account for the reduced number of groups before the 1950s, a
877 stratified resampling scheme (Shankar, 2020) is used to guarantee that the resampled groups better
878 reflect the prevalence of groups throughout the history of marine observation. Specifically, groups
879 are divided into two strata based on whether they were present in 20-year windows before 1940–
880 1959. The resampling is then performed within each stratum with replacement and repeated 10,000
881 times. On average, updating the bootstrapping technique slightly increases the 95% CI of York fit
882 slopes by 6%, and the interquartile range by 1%. Among the 10,000 time series of bootstrapped
883 York regression slope, 9,306 of them have, on average, positive values over 1910–1929 but negative
884 values afterward, indicating that the relationship between diurnal amplitudes and groupwise off-
885 sets changed sign significantly ($p < 0.1$) in the 1930s. A Matlab script to reproduce this updated
886 bootstrapping analysis is in the supplement to this article.