

1 **Disaggregation of Remotely Sensed Land Surface**
2 **Temperature: A Simple yet Flexible Index (SIFI) to Assess**
3 **Method Performances**

4

5 Lun Gao^a, Wenfeng Zhan^{a, b*}, Fan Huang^a, Xiaolin Zhu^c, Ji Zhou^d, Jinling Quan^e,

6 Peijun Du^f, and Manchun Li^f

7

8 a. Jiangsu Provincial Key Laboratory of Geographic Information Science and
9 Technology, International Institute for Earth System Science, Nanjing University,
10 Nanjing, Jiangsu 210023, China

11 b. Jiangsu Center for Collaborative Innovation in Geographical Information Resource
12 Development and Application, Nanjing, 210023, China

13 c. Department of Land Surveying and Geo-Informatics, The Hong Kong Polytechnic,
14 Hong Kong 999077, China

15 d. School of Resources and Environment, University of Electronic Science and
16 Technology of China, Chengdu, Sichuan 610054, China

17 e. Institute of Geographic Sciences and Natural Resources Research, Chinese
18 Academy of Sciences, Beijing 100101, China

19 f. Department of Geographic Information Science, Nanjing University, Nanjing,
20 Jiangsu 210023, China

21

22

23

24

25 **Contact information**

26 ***Corresponding Author:** W. Zhan, Nanjing University at Xianlin Campus, No.163

27 Xianlin Avenue, Qixia District, Nanjing, Jiangsu 210023, P. R. China; Fax:

28 +86-25-89681030.

29 **E-mail Addresses:** gaolun724@foxmail.com (L. Gao), zhanwenfeng@nju.edu.cn (W.

30 Zhan), nju_huangfan@163.com (F. Huang), zhuxiaolin55@gmail.com (X. Zhu),

31 jzhou233@uestc.edu.cn (J. Zhou), quanjinlin_@126.com (J. Quan), dupjrs@126.com

32 (P. Du), and limanchun@nju.edu.cn (M. Li).

33

34

35

36

37 **Abstract**

38 Disaggregation of land surface temperature (DLST), the aim of which is to
39 generate LSTs with fine resolution, has been attracting increasing attention since the
40 1980s. The past three decades have been witness to the emergence of DLST methods
41 in large numbers, the accuracies of which were often assessed by comparing the
42 disaggregated with fine spatial resolution LSTs using error indexes such as the root
43 mean square error (RMSE). However, the majority of previous error indexes are, by
44 their nature, insufficient for assessing the performances of DLST methods. This
45 insufficiency is due in part to their lower competence at distinguishing the DLST error
46 from LST retrieval errors and in part to their inability to remove the process controls
47 resulting from different thermal contrasts, temperature units, and resolution ratios
48 among different scenarios in which DLST is conducted. This is also because they are
49 unable to denote the sharpening statuses of the DLST results (e.g., under- or
50 over-sharpening). This *status quo* has made the evaluation of method performances
51 challenging and sometimes unreliable.

52 To better assess DLST method performances under diversified scenarios, we
53 formulated five protocols, through which a simple yet flexible index (SIFI) was
54 subsequently designed. The establishment of an SIFI includes the following four steps:
55 (1) a detail-based evaluation, which is designed primarily to exclude the impacts of
56 systematic deviations on estimated LSTs; (2) a Gaussian normalization, which is
57 primarily intended to remove the differences in temperature units and thermal
58 contrasts; (3) a triple comparison, with the aim of attenuating the influence of the
59 difference in the resolution ratio in comparisons of method performances; and (4) a
60 piecewise comparison, which is primarily scheduled to distinguish among the three

61 sharpening statuses, under-sharpening, acceptable over-sharpening, and unacceptable
62 over-sharpening. The evaluation ability of SIFI was compared with those of the
63 RMSE, Erreur Relative Globale Adimensionnelle de Synthèse (ERGAS), and image
64 quality index (Q) using simulation tests and actual thermal data. The results illustrate
65 that SIFI generally outperforms the other indexes; it is able to mitigate the impacts
66 from process errors and controls during evaluation and is able to indicate the
67 sharpening statuses accurately. We believe this new index will likely promote the
68 design of future DLST algorithms and procedures.

69 **Keywords**

70 Thermal remote sensing; land surface temperature; disaggregation; model

71 performance; and accuracy assessment.

72

73

74

75

76

77

78 **1. Introduction**

79 The large-scale monitoring of the thermal status of land surfaces was difficult
80 and even impossible until the advent of satellite thermal infrared remote sensing.
81 Thermal sensors enable the generation of the land surface temperature (LST) products,
82 which are instrumental to research in many disciplines (Anderson et al., 2012; Bisht et
83 al., 2005; Jiménez-Muñoz et al., 2016; Sandholt et al., 2002; Sobrino et al., 2007,
84 2012; Teggi, 2012). However, spaceborne sensors are subject to a tradeoff between
85 spatial and temporal resolutions (Zhan et al., 2013), and the spatial resolution of
86 thermal spaceborne-derived LST maps is too coarse for many applications. This
87 challenge has encouraged research on the spatial disaggregation of LST (DLST),
88 which is able to generate LST images with high spatial and temporal resolutions.

89 A quick literature survey shows that DLST has experienced phenomenal growth
90 in the past three decades, and more methods have been proposed, particularly since
91 the 2000s (Zhan et al., 2013). Most of them tried to reconstruct thermal details with
92 the aid of finer resolution data sets (e.g., data in other bands, classification maps, or
93 designed scaling factors), transform these details into thermal ones by statistical
94 inference, and finally add them to the coarse resolution LSTs (Chen et al., 2014). In
95 considering the fast dynamics of LSTs, recent developments in DLST have been
96 focused on the simultaneous disaggregation of the spatial and temporal resolutions
97 (Addesso et al., 2015; Mechri et al., 2014; Moosavi et al., 2015; Weng et al., 2014;
98 Wu et al., 2015; Zhan et al., 2016; Gao et al., 2017). It is anticipated that DLST will
99 continue to be a research focus in the foreseeable future because the resolutions of the
100 current and planned satellite thermal sensors remain far from satisfactory for the
101 relevant applications (Anderson et al., 2012; Lagouarde et al., 2013; Roberts et al.,

102 2012; Teggi & Despini, 2014; Zhou et al., 2013).

103 As more and more DLST methods are being proposed, an index that is more
104 appropriate for assessing their performances under various scenarios is urgently
105 required. Early studies on DLST often employ conventional indexes that measure the
106 similarity between disaggregated and fine resolution LSTs, such as the root mean
107 square error (RMSE) (Agam et al., 2007) and the mean absolute error (MAE) (Nishii
108 et al., 1996; Stathopoulou & Cartalis, 2009). Subsequent studies also use the Erreur
109 Relative Globale Adimensionnelle de Synthèse (ERGAS), which is able to eliminate
110 the resolution difference between pre- and post-disaggregation LSTs (Gevaert &
111 García-Haro, 2015; Pardo-Igúzquiza et al., 2006). In considering the documented
112 similarity between DLST and optical image fusion (OIF) since 2010, a few
113 researchers have resorted to the indexes that were designed to evaluate the OIF
114 algorithms. These indexes include the universal image quality index (Q) (Mukherjee
115 et al., 2014; Zhou et al., 2016) and the structural similarity index (SSIM)
116 (Rodriguez-Galiano et al., 2012), among others.

117 Practitioners may also turn to other advanced indexes in OIF that were recently
118 developed for evaluations, such as the four bands multispectral images fusion index
119 (Q4), the Quality with No Reference (QNR) (Vivone et al., 2014), or the combination
120 of various indexes (Despini et al., 2014). However, although it inherited some traits
121 from the OIF, DLST has differed from its counterpart in the following two regards.
122 First, the DLST process strictly requires that the thermal radiance of a single pixel
123 block at the coarse resolution should be equal to the mean thermal radiance of the
124 corresponding disaggregated fine resolution pixels (Liang, 2005) because its
125 applications are primarily quantitative, whereas spectral distortion is occasionally
126 tolerable in OIF (Pohl & Van Genderen, 1998). Second, fine resolution LSTs, which

127 are either obtained by using thermal data from a different sensor or directly produced
128 by the aggregation-and-then-disaggregation strategy, are indispensable for validation
129 (Zhan et al., 2011). Although references from a different sensor can also help in the
130 evaluation of OIF techniques, they are frequently assessed by comparison with the
131 coarse resolution multispectral images in terms of spectral distortion and with the fine
132 resolution panchromatic images with respect to spatial details (Vivone et al., 2014).

133 Although previous indexes can be used to assess the performances of DLST
134 methods, they are intrinsically flawed in the following three regards. First, their values
135 depend on multiple errors, including those from disaggregation methods but also
136 those due to the preprocessing of LST, which is unrelated to the model performance
137 (e.g., the temperature retrieval error; more clarifications are given in Section 2.1). Any
138 comparison that disregards the errors due to LST preprocessing would no longer be
139 related to the method performances alone. Second, their values depend on multiple
140 controls, including that from the disaggregation method but also those related to the
141 thermal contrast difference and resolution gap. Finally, their values are mostly not
142 indicative of the sharpening statuses including under-sharpening, acceptable
143 over-sharpening, and unacceptable over-sharpening (more clarifications are given in
144 Section 2.3).

145 To address these issues, this work designed a new index able to better assess
146 DLST method performances. Followed by the clarifications of background (Section 2)
147 and the five protocols (Section 3) that an index should comply with, Section 4
148 provides the definition of this index. Sections 5 and 6 exhibit the experiment, the
149 results and discussion, respectively. The conclusions are finally drawn in Section 7.

150

151

152 2. Background

153 An accurate evaluation of the method performances requires researchers to first
154 identify all the possible errors/controls that may affect the associated evaluation.

155 Generally, the overall errors of disaggregated LSTs (given as $err_{overall}$) can be
156 expressed as the function of the temperature retrieval errors (given as e_{LST} , including
157 the errors from both the original low-resolution and the reference fine resolution LST
158 images), the image co-registration error (given as e_{cr}), and the DLST error (given as
159 e_{DLST}). In other words, $err_{overall}$ can be expressed as follows:

$$160 \quad err_{overall} = q_1(\overbrace{e_{LST}, e_{cr}}^{\text{process error}}, \overbrace{e_{DLST}}^{\text{DLST error}}) \quad (1)$$

161 where q_1 is the function between $err_{overall}$ and the three types of errors. Hereafter, we
162 refer to the combination of e_{LST} and e_{cr} as the ‘process errors’ because they primarily
163 stem from the pre-processes that are performed before DLST is conducted (more
164 clarifications are given in Section 2.1).

165 Nevertheless, it remains unsuitable to use e_{DLST} to represent the performances of
166 the DLST methods because e_{DLST} is also dependent on several other controls in
167 addition to the performance control (given as c_{pm}). These controls are involved in
168 scenarios under which the method performance can be distorted; they include
169 scenarios with different thermal contrasts (given as c_{tc}), temperature units (given as
170 c_{tu}), and resolution ratios (given as c_{rr}). Therefore, the DLST error can be given by the
171 following:

$$172 \quad e_{DLST} = q_2(\overbrace{c_{tc}, c_{tu}, c_{rr}}^{\text{process control}}, \overbrace{c_{pm}}^{\text{performance}}) \quad (2)$$

173 where q_2 is the function between the DLST error and the associated controls.
 174 Hereafter we refer to the combination of c_{tc} , c_{tu} , and c_{rr} as the ‘process controls’ (refer
 175 to Section 2.2 for more details).

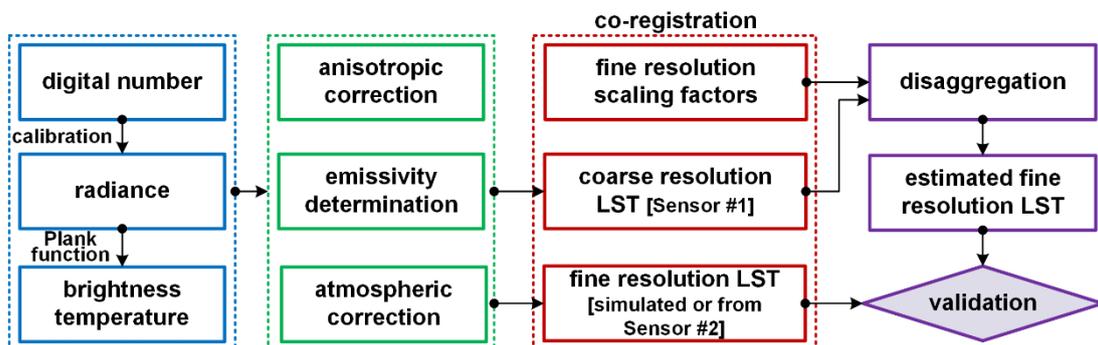
176 The above analysis indicates that the evaluation of method performances should
 177 be conducted by c_{pm} rather than by $err_{overall}$. In other words, the impacts from the
 178 process errors and controls should be excluded before the precise evaluation of
 179 method performances. In addition, the performance evaluations would be further
 180 improved, once the sharpening statuses, including the under-sharpening, acceptable
 181 over-sharpening, and unacceptable over-sharpening, is determined. Elaborate
 182 interpretations of this issue are presented in Section 2.3.

183

184 2.1. Process errors

185 As indicated by Eq. (1) and graphically represented in Fig. 1, the overall errors
 186 for disaggregated LSTs include both the DLST and process errors. The process errors
 187 can be divided into temperature retrieval and image registration errors.

188



189

190 **Fig. 1.** Graphical representation of the combined temperature retrieval and DLST
 191 processes, in which the process errors (including the temperature retrieval and
 192 coregistration errors) and DLST errors are blended.

193

194 The remote retrieval of the surface temperature is a complex process (Fig. 1).
195 Temperature retrieval errors may be directly due to noise-equivalent temperature
196 differences (NE Δ T) (Gillespie et al., 1998), or due to inaccuracies/differences in the
197 conversion from the digital number (DN) to the thermal radiance (i.e., the radiometric
198 calibration process) and then to the brightness temperature (BT) (see Fig. 1). These
199 types of errors depend directly on the calibration coefficients that are estimated from
200 calibration fields (Chander et al., 2009). For LSTs estimated from mature spaceborne
201 thermal sensors that have been carefully calibrated (e.g., MODerate-resolution
202 Imaging Spectroradiometer, MODIS), this type of error is insignificant. However, this
203 error is not trivial for some other spaceborne thermal sensors (e.g., early Landsat
204 series) (Barsi et al., 2007) or when using fine resolution LSTs (e.g., at a meter level)
205 as validation data, which are usually obtained from airborne thermal missions when
206 the calibration may not be adequately accurate (Sobrino et al., 2004).

207 Temperature retrieval errors may also stem from the uncertainties in the surface
208 thermal anisotropy, determination of emissivity, and atmospheric corrections (see Fig.
209 1), such as the approximate parameterizations in the mono-window or single-channel
210 algorithms (Qin et al., 2001; Jiménez-Muñoz & Sobrino, 2003). In this regard, these
211 errors reflect the accuracy of the estimation of the true LST from the thermal radiance.
212 Being subject to the inaccurate parameterization of atmospheric thermal radiation, the
213 problems in estimating the emissivity, and the surface thermal anisotropy, the errors
214 are usually lower over relatively homogeneous surfaces (approximately 0.5 to 2.0 K)
215 but considerably higher over heterogeneous terrains (Li et al., 2013). These errors
216 may be much greater over urban areas (reaching 5.0 K or more) due to significant
217 urban thermal anisotropy (Lagouarde et al., 2010).

218 Errors due to inaccurate temperature retrieval errors are primarily linear and
219 systematic, i.e., the retrieved LSTs compared with the ground truth are systematically
220 higher or lower, with a small portion of errors being nonlinear and random. For
221 example, systematic deviations may occur between the coarse LSTs to be
222 disaggregated and the reference finer resolution LSTs come from other sources
223 (Merlin et al., 2010; Bechtel et al., 2012; Zakšek et al., 2012; Zhou et al., 2015).
224 Errors due to inaccurate calibration are linear because the calibration process itself is
225 often conducted through a linear function (Chander et al., 2009). Errors caused by
226 surface thermal anisotropy are typically systematic for neighboring pixels once the
227 accompanying land cover types are similar, but they may become random for nearby
228 pixels with different land cover types (Lagouarde et al., 2010). Errors caused by
229 inaccurate atmospheric thermal parameterizations are also primarily systematic,
230 because the reflected downward and/or upward atmospheric thermal radiance can be
231 under or over corrected (Li et al., 2013). By comparison, errors that are attributable to
232 inaccurate emissivity depend on the associated estimation process, and these errors
233 can be either systematic or random.

234 We should note in particular that co-registration errors might also be
235 incorporated into the evaluation of the DLST methods. These errors are derived from
236 the mismatch among the coarse resolution LSTs, the fine resolution scaling factors,
237 and the fine resolution LSTs used for validation (see Fig. 1). Errors due to this type of
238 mismatch (i.e., inaccurate coregistration) depend on the practitioner and are highly
239 nonlinear and random. These errors are difficult to quantify and even more difficult to
240 suppress before the evaluation of the method performances. More discussions on this
241 issue are further provided in Section 6.4.

242

243 **2.2. Process controls**

244 As indicated by Eq. (2), process controls can also distort the performance
245 evaluation, and they usually include the thermal contrast control (c_{tc}), temperature unit
246 control (c_{tu}), and resolution ratio control (c_{rr}).

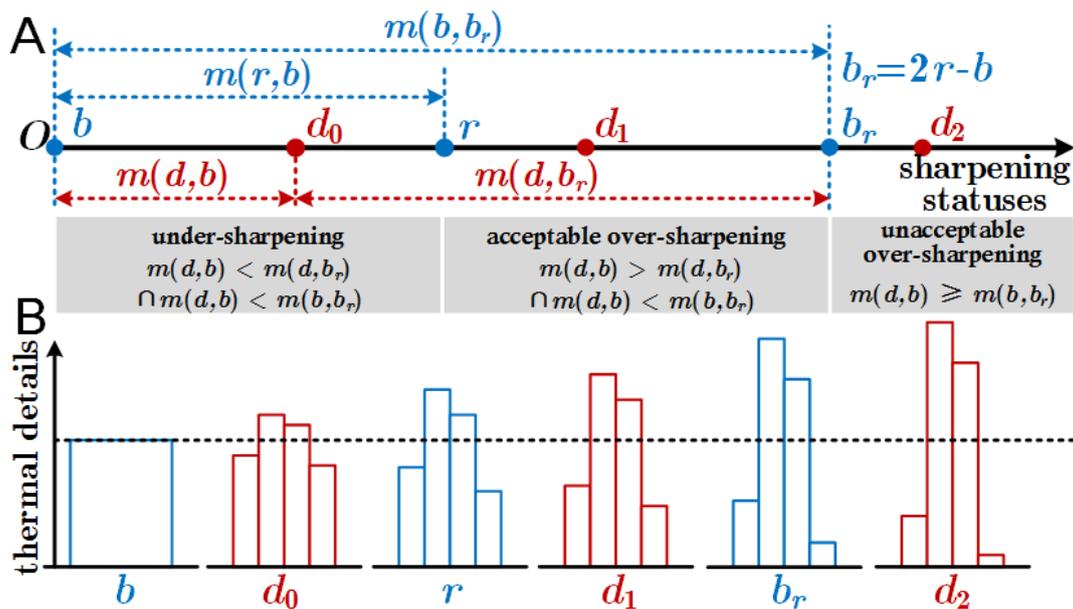
247 First, a DLST may be performed over different areas with a great variety of
248 thermal contrast. For example, an RMSE of 1.5 K for a method that is applied over
249 urban areas with high thermal contrast vs. a value of 1.0 K for another method over
250 vegetated areas does not necessarily indicate that the 1.5 K RMSE method performs
251 worse than the other one. Any index that disregards c_{tc} would no longer be indicative
252 of the method performance. Second, a DLST may be just as well implemented at three
253 levels with different units, including the digital number (DN, no unit) level (Liu and
254 Moore, 1998; Zhukov et al., 1999), the radiance (unit: $W \cdot m^{-2} \cdot \mu m^{-1} \cdot sr^{-1}$) level (Liu and
255 Zhu, 2012), and the temperature level with centigrade ($^{\circ}C$), Fahrenheit ($^{\circ}F$), or Kelvin
256 (K) degrees as its units (Zhan et al., 2013). The index values should be comparable
257 when DLST is performed at all these levels. Third, DLST may be further conducted
258 with different resolution ratios between pre- and post-disaggregation LSTs. The
259 resolution ratio usually ranges from several times (e.g., from ~ 100 to 30 m for the
260 downscaling of Landsat thermal images) (Gao et al., 2017) to as large as several tens
261 of times (e.g., from $\sim 5,000$ to 100 m for the downscaling of geostationary thermal
262 images) (Bechtel et al., 2012). Given the identical RMSEs for these two cases, it is
263 understandable that the method performance for the former case will be worse than
264 the latter. Any index that disregards c_{rr} will be uninformative regarding the method
265 performance.

266

267 **2.3. Sharpening statuses**

268 In general, DLST is used to try to generate fine resolution LST; it can also be
 269 perceived as a process that adds thermal details to the background low-resolution
 270 LSTs (Chen et al., 2014). Through a DLST method, the added thermal details may be
 271 less than or more than needed. In addition, the added thermal details may be ‘much
 272 more’ than needed, making the disaggregated LSTs even further away from the
 273 reference fine resolution LSTs than the background low-resolution ones. This scenario
 274 may sometimes be acceptable for image fusion that aims for target detection, but it is
 275 unacceptable for DLST because the application of LSTs is primarily quantitative (e.g.,
 276 surface flux estimation).

277



278

279 **Fig. 2.** Conceptual description of the three sharpening statuses (A). The coarse
 280 resolution, disaggregated, and fine resolution LST images are denoted by the three
 281 dots at b , d , and r , respectively. In the coordinate axis that starts at O , b and r remain
 282 constant for a single DLST process while d can be located at any point on this axis
 283 depending on the specific DLST method. d_0 , d_1 , and d_2 represent the

284 *under-sharpening, acceptable over-sharpening, and unacceptable over-sharpening*
285 cases, respectively. b_r is the mirror image of b when using r as the center of symmetry,
286 and it can be estimated by finding $2r - b$. $m(\cdot)$ is a distance metric between two LST
287 images, and it corresponds to the RMSE when the Euclidean distance is used. A
288 simple graphical illustration of the sharpening statuses is further provided in (B),
289 where it is assumed that a single LST pixel is divided into four pixels with different
290 LST values.

291

292 We therefore provide the three possible sharpening statuses for the DLST (see
293 Fig. 2), including the *under-sharpening* (corresponding to d_0), *acceptable*
294 *over-sharpening* (corresponding to d_1) and *unacceptable over-sharpening*
295 (corresponding to d_2). Note that the sharpening statuses in Fig. 2 is displayed in a
296 single dimension of thermal details. Please refer to Appendix A for the description of
297 the sharpening statuses at higher dimensions.

298 (1) *Under-sharpening*: This term signifies that generally less thermal details are
299 added to the coarse resolution LSTs than needed. In this case, the distance
300 between d (i.e., the disaggregated LSTs) and b (i.e., the background
301 low-resolution LSTs) is shorter than that between d and b_r (i.e., the mirror image
302 of b) and shorter than that between b and b_r , i.e., $m(d, b) < m(d, b_r) \& m(d, b) <$
303 $m(b, b_r)$ (see Fig. 2A), where $m(\cdot)$ is the distance between the two associated LST
304 images.

305 (2) *Acceptable over-sharpening*: This term implies that more thermal details are
306 added than needed, but these redundant details remain tolerable. In this case, the
307 distance between d and b is greater than that between d and b_r , whereas they are

308 still lower than that between b and b_r : $m(d, b) > m(d, b_r) \& m(d, b) < m(b, b_r)$
309 (see Fig. 2A).

310 (3) *Unacceptable over-sharpening*: This status suggests that DLST fails, not just
311 because there are more added thermal details than necessary, but because they
312 also lead to a situation in which the disaggregated LSTs are further away from
313 the fine resolution ones when compared with the original background
314 low-resolution LSTs. In other words, the post-disaggregation results are even
315 worse than the pre-disaggregation ones, a consequence that is intolerable for the
316 quantitative applications of DLST. In this case, $m(d, b) > m(b, b_r)$ (see Fig. 2A).

317 Note that under-sharpening and acceptable over-sharpening are divided by the
318 reference fine resolution LST image (r), while acceptable and unacceptable
319 over-sharpening are separated by the mirror image of the background LSTs (i.e., b_r).

320 Herein, b_r is defined as the mirror image of b with r as the center of symmetry. In
321 other words, the addition of b and b_r is equal to $2r$ ($b + b_r = 2r$). Physically, the
322 sharpening is no longer tolerable once the disaggregated LSTs possess more thermal
323 details than b_r does because in this case the disaggregated LSTs are less quantitatively
324 accurate even when compared with b (see Fig. 2).

325

326

327

328 3. Protocols and clarifications for designing a DLST index

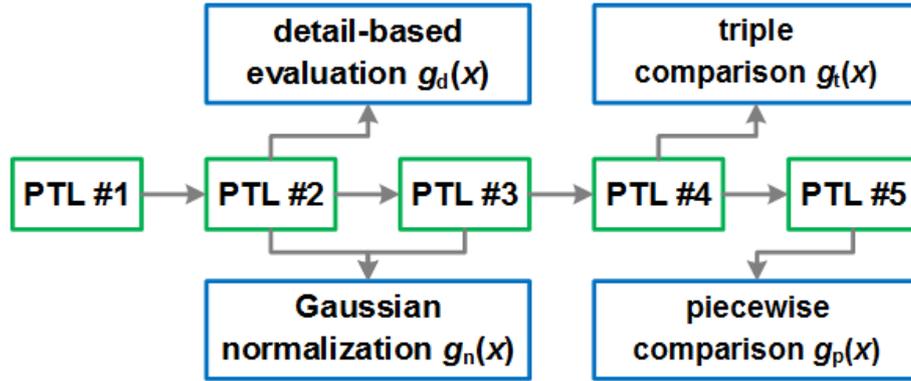
329 A majority of the previous DLST studies used error indexes (e.g., the RMSE),
330 which is commonly directly estimated on the basis of the disaggregated LSTs and the
331 reference fine resolution LSTs to evaluate the performances of the proposed DLST
332 methods (Agam et al., 2007; Zhou et al., 2016). However, the aforementioned analysis
333 shows that, without carefully differentiating the process errors, process controls, and
334 sharpening statuses, the evaluation of the method performances by error indexes such
335 as the RMSE may be inaccurate or even misleading. According to the analysis in
336 Section 2, we propose that a suitable intercomparison index should comply with the
337 following protocols:

338 *PTL #1: The index should simultaneously measure how much the disaggregated LSTs*
339 *are similar to the fine resolution LSTs as well as how much the*
340 *disaggregated LSTs are different from the original (i.e., coarse) LSTs.*

341 PTL #1 is adopted with adaptations from Wald et al. (1997) in which ‘*any*
342 *synthetic image should be as identical as possible to the multispectral set of images*
343 *that the corresponding sensor would observe with the high resolution*’. In PTL #1, the
344 similarity between the disaggregated and reference fine resolution LSTs (i.e., a
345 pairwise comparison) can be evaluated by distance measures (e.g., the RMSE and
346 MAE). Nevertheless, it remains insufficient to only measure this similarity between
347 the disaggregated and reference LSTs using indexes such as the RMSE because the
348 RMSE can be low, indicating that a high accuracy is achieved, even if the actual
349 DLST procedure fails. For instance, the $m(d_0, r)$ may remain small, even if $d_0 = b$ (i.e.,
350 DLST fails or no DLST has been conducted), because the $m(b, r)$ could already be
351 small (see Fig. 2). In other words, the RMSE between the disaggregated and fine
352 resolution LSTs may remain small when no or very few (i.e., under-sharpened)

353 thermal details are added to coarse LSTs because in many cases, the RMSE between
 354 the coarse and fine resolution LSTs is already small (e.g., less than 2.0 K). Therefore,
 355 the dissimilarity between the disaggregated and coarse resolution LSTs also requires
 356 special consideration.

357



358

359 **Fig. 3.** The five protocols (PTLs) and the associated strategies used to design a
 360 suitable index for assessing the DLST results. $g_d(x)$, $g_n(x)$, $g_t(x)$, and $g_p(x)$ represent
 361 four functions (or procedures) that characterize the detail-based evaluation, Gaussian
 362 normalization, triple comparison, and piecewise comparison required by PTLs #2 to
 363 #5.

364

365 **PTL #2:** *The index should be independent of the temperature retrieval errors.*

366 PTL #2 addresses temperature retrieval errors and it demands a detail-based
 367 evaluation as well as Gaussian normalization for the index design, which are given by
 368 the following equations, respectively:

369
$$g_d(x) = x - \bar{x} \quad (3)$$

370
$$g_n(x) = (x - \mu_x) \cdot \sigma_b^{-1} \quad (4)$$

371 where $g_d(x)$ and $g_n(x)$ are the two functions denoting the detail-based evaluation and
 372 Gaussian normalization, respectively; x and \bar{x} are the fine resolution LST images

373 and their aggregated ones at the coarse resolution; and μ_x and σ_b are the mean and
374 standard deviation of LST images. It is notable that σ_b , rather than the standard
375 deviations of d and r , is used uniformly for these three LST images because the latter
376 two standard deviations are more sensitive to outliers.

377 As analyzed in Section 2.1, most of the temperature retrieval errors are linear and
378 systematic. On one hand, the detail-retrieval process expressed by Eq. (3) is able to
379 remove most of the systematic errors due to imprecise atmospheric thermal correction
380 and locally systematic errors due to surface thermal anisotropy as well as a part of the
381 errors due to inaccurate emissivity determinations. It is reasonable that temperature
382 retrieval errors caused by these factors can be suppressed by subtracting the
383 corresponding aggregated LSTs at coarse resolution because LST estimations due to
384 such factors are systematically higher or lower for adjacent pixels. On the other hand,
385 the Gaussian normalization given by Eq. (4) is able to eliminate the linear errors due
386 to inaccurate calibration because this type of normalization is invariant in response to
387 linear transformations. The strict proof is provided in Appendix B.

388 *PTL #3: The index should be comparable when DLST is performed among areas with*
389 *different LST contrasts, and it should be comparable when DLST is*
390 *performed using different types of temperature units.*

391 PTL #3 addresses the thermal contrast control (c_{tc}) and temperature unit control
392 (c_{tu}). In addition to being able to remove a part of the temperature retrieval errors, the
393 Gaussian normalization $g_n(x)$ given by Eq. (4) is expected to be capable of
394 suppressing c_{tu} because the temperature unit conversion (e.g., from K to °C and to °F)
395 is linear. The Gaussian normalization is also able to suppress c_{tc} once the standard
396 deviation of an image is used to represent its thermal contrast (Wang & Bovik, 2002).
397 Note that ERGAS, with its expression given in Section 5.3, also possesses a

398 normalization factor (i.e., the mean of each band). However, it is unable to address the
399 fully linear unit conversion with both the gain and offset (e.g., from K to °C) as well
400 as the dissimilar scenarios with LST contrast differences, e.g., between the highly
401 heterogeneous urban surfaces and relatively homogeneous cultivated lands.

402 *PTL #4: The index should be comparable when DLST is performed with different*
403 *resolution ratios between pre- and post-disaggregation LSTs.*

404 PTL #4 addresses the resolution ratio control (c_{rr}). A triple comparison function
405 among the coarse (b), disaggregated (d), and reference fine resolution (r) LSTs can
406 suppress c_{rr} , which is given as follows:

$$407 \quad g_t(b, d, r) = m(d, r) / m(d, b) \quad (5)$$

408 where $g_t(\cdot)$ denotes the triple comparison function; and $m(\cdot)$ is a distance metric
409 between two LST images, and it corresponds to the RMSE when the Euclidean
410 distance is used. In two disaggregation cases when the $m(d, r)$ remain the same but the
411 in-between resolution gaps are different, it is reasonable that the case with a larger
412 resolution gap indicates a better model performance. Eq. (5) is efficient at suppressing
413 c_{rr} in such cases. This is because the specific case with a larger resolution gap also
414 likely suggests a higher $m(d, b)$, and with the division given by Eq. (5), the resulting
415 $g_t(\cdot)$ will decrease, consequently indicating a better performance. Note that $g_t(\cdot)$
416 physically measures the similarity between d and r as well as the dissimilarity
417 between b and r and it is thus related to PTL #1.

418 *PTL #5: The index should be indicative of the sharpening status.*

419 PTL #5 addresses the sharpening statuses. As illustrated in Section 2.3, the
420 differentiation among the three sharpening statuses requires a piecewise function that
421 considers the position of d on the axis shown in Fig. 2A, with d_0 , d_1 , and d_2 denoting
422 the under-sharpening, acceptable over-sharpening, and unacceptable over-sharpening,

423 respectively. In combining Eq. (5) as required by PTL #4, we provide the following
 424 piecewise function to satisfy PTL #5:

$$425 \quad g_p(b, d, r) = \begin{cases} m(d, r)/m(d, b), & d = d_0, m(d, b) \leq m(d, b_r) \cap m(d, b) < m(b, b_r) \\ -m(d, r)/m(d, b_r), & d = d_1, m(d, b) > m(d, b_r) \cap m(d, b) < m(b, b_r) \\ \text{NaN}, & d = d_2, m(d, b) \geq m(b, b_r) \end{cases} \quad (6)$$

426 where $g_p(\cdot)$ is the piecewise function; NaN indicates that the disaggregated LSTs are
 427 unacceptable for quantitative applications. Note that (1) for the acceptable
 428 over-sharpening, $m(d, b_r)$ rather than $m(d, b)$ is used as the division factor, aiming at
 429 weighting d_0 and d_1 equally once they have the same distance away from r ; and (2)
 430 the minus symbol when $d = d_1$ is used to differentiate the acceptable over-sharpening
 431 from the under-sharpening.

432
 433
 434

435 **4. Definition**

436 **4.1. Standard definition**

437 Using guidance from the proposed protocols, we were able to design a simple yet
438 flexible index (known as the SIFI hereafter) to assess the performances of the DLST
439 methods. Its standard definition is given as follows:

$$440 \quad \text{SIFI} = \begin{cases} m(D, R)/m(D, B), & m(D, B) \leq m(D, B_R) \cap m(D, B) < m(B, B_R) \\ -m(D, R)/m(D, B_R), & m(D, B) > m(D, B_R) \cap m(D, B) < m(B, B_R) \\ \text{NaN}, & m(D, B) \geq m(B, B_R) \end{cases} \quad (6)$$

441 where B , D , and R are the three variables that correspond to the background coarse
442 resolution (b), disaggregated (d), and reference fine resolution (r) LSTs, respectively;
443 and B_R is the mirror image of B that use R as the center of symmetry, i.e., $B_R = 2R - B$.
444 They are obtained by using following equations:

$$445 \quad \begin{aligned} X &= g_n(g_d(x)) \\ &\begin{cases} g_d(x) = x - \bar{x}; \\ g_n(x) = (x - \mu_x) \cdot \sigma_b^{-1} \end{cases} \end{aligned} \quad (7)$$

446 where X denotes B , D , or R , while x denotes b , d , or r ; and $g_d(x)$ and $g_n(x)$ and their
447 associated variables are well explained in the text subsequent to Eqs. (3) and (4). The
448 piecewise functions given by Eq. (6) represent the under-sharpening, acceptable
449 over-sharpening, and unacceptable over-sharpening, respectively.

450 From Eq. (6), one can infer that the SIFI ranges from negative to positive infinity.
451 SIFI approximates to zero once the disaggregated LSTs are close to the fine resolution
452 LSTs (i.e., $m(D, R) \rightarrow 0$). The SIFI becomes very large when few thermal details are
453 added – disaggregated LSTs are relatively close to (but not completely equivalent to)
454 the coarse LST (i.e., $m(D, B) \rightarrow 0$). The SIFI then becomes negatively large when more
455 thermal details than needed have been added – disaggregated LSTs are close to the

456 mirror of the coarse LSTs (i.e., $m(D, B_R) \rightarrow 0$). The SIFI is assigned as ‘NaN’ when
457 redundant details (that would make the DLST fail) have been added (i.e., $m(D, B) \geq$
458 $m(B, B_R)$). In general, the variation in SIFI is continuous for the under-sharpening and
459 acceptable over-sharpening (refer to Fig. 5 for more visual illustrations); it becomes
460 discontinuous for unacceptable over-sharpening by setting its value as NaN. From Eq.
461 (7), one can also infer that for the under-sharpening, the smaller the SIFI values, the
462 better the DLST results, while this phenomenon is reversed for the acceptable
463 over-sharpening. Note that $m(D, B)$ will always be greater than zero because at least
464 ‘some’ details may be added by the DLST process. In addition, this study calculates
465 the distance metric $m(\cdot)$ between two LST images in a global fashion, i.e., only a
466 single distance value is estimated for an entire image. More discussions on the
467 moving window based calculation $m(\cdot)$ for two images are provided in Section 6.4.

468

469 **4.2. Simplifications under specific conditions**

470 The following simplifications can only be justified under particular conditions,
471 and researchers should use Eqs. (6) and (7) to calculate the SIFI when the following
472 conditions/assumptions are not satisfied. First, in considering that the statistical
473 downscaling method is regularly used for DLST (Zhan et al., 2013) and the
474 relationships between the coarse resolution LSTs and scaling factors are usually less
475 represented by the statistical downscaling methods, under-sharpening (i.e., the added
476 thermal details are insufficient) appears more frequently than the other two statuses.
477 Second, the Euclidean distance (i.e., the RMSE) is usually considered the most
478 frequently used similarity metric (More discussions on the use of distance metrics
479 other than the RMSE are provided in Section 6.4). When the over-sharpening does not

480 occur and the Euclidean distance is employed, the SIFI given by Eq. (6) can be
481 simplified into the following equation:

$$482 \quad \text{SIFI} = \text{RMSE}(D, R) \cdot [\text{RMSE}(D, B)]^{-1} \quad (8)$$

483 Once the fine and coarse resolution LSTs are also coming from an identical
484 source, i.e., the disaggregated LSTs are validated by the
485 aggregation-and-then-disaggregation strategy, then the SIFI given by Eq. (8) can be
486 further deduced into the following:

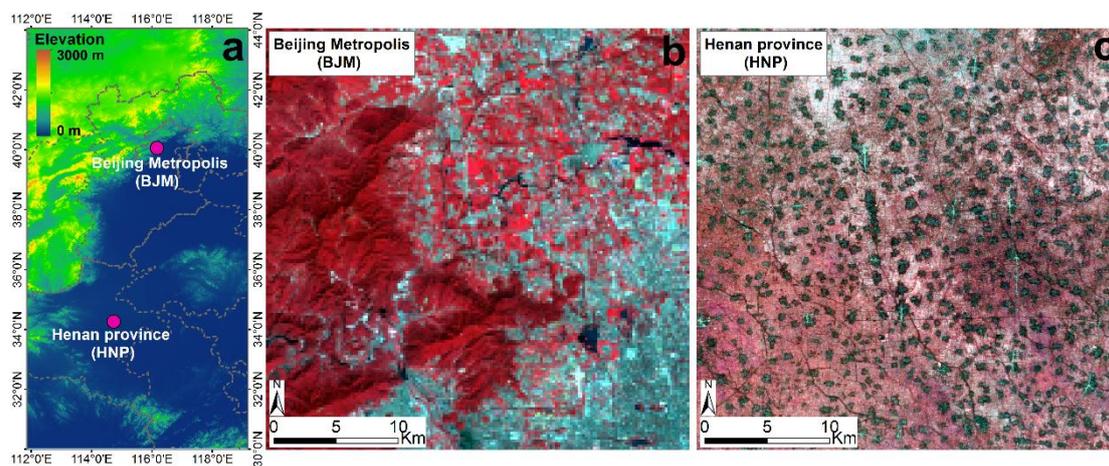
$$487 \quad \text{SIFI} = \frac{\text{RMSE}(d, r)}{\text{RMSE}(b, d)} = \frac{\sqrt{E[(d-r)^2]}}{\sqrt{E[(b-d)^2]}} \quad (9)$$

488 where $E(\cdot)$ denotes the expectation. Please refer to Appendix C for the proof of the
489 simplification from Eq. (8) to (9). Note that the aggregation-and-then-disaggregation
490 strategy may be feasible for the development of new algorithms, but the aim of the
491 DLST is to generate finer-resolution LSTs. For the validation of disaggregated results
492 by fine resolution LSTs from another source, researchers should use the complete
493 form, i.e., Eqs. (6) and (7), rather than Eqs. (8) or (9), to calculate the SIFI.

494 **5. Experiments**

495 **5.1. Datasets and utilities**

496 Two study areas with different surface landcover types were selected (Fig. 4).
497 The first test site, which is labeled BJM, is situated in the northwestern part of the
498 Beijing Metropolis ($39^{\circ}3'6''N - 40^{\circ}4'33''N$; $115^{\circ}5'24''E - 116^{\circ}3'40''E$). The BJM
499 consists of a mixture of urban, rural, and mountainous surfaces. This site was chosen
500 mostly because of its high heterogeneity, which makes it appropriate for testing model
501 performances. The other site, which is labeled HNP, is located in the Henan Province
502 ($34^{\circ}0'44''N - 34^{\circ}7'60''N$, $114^{\circ}0'17''E - 115^{\circ}0'22''E$), and it corresponds to an area
503 covered by fallow field, wheat paddock, and small towns, with a substantially flat
504 terrain. We chose the HNP because the DLST over rural areas is one of its most
505 important applications (Agam et al., 2007; Bindhu et al., 2013).
506



507
508 **Fig. 4.** Geographical location of the two test areas. (a) shows the region over North
509 and East China; (b) demonstrates the northwestern section of the Beijing Metropolis
510 (BJM), as represented by indicating the ASTER bands 3, 2, and 1 as the red, green,
511 and blue channels, respectively; and (c) describes a typical area in central Henan
512 Province (HNP), as provided by indicating TM bands 4, 3, and 2 as the associated

513 channels. The spatial resolutions of (b) and (c) were both aggregated to 200 m from
514 their original resolutions.

515

516 To validate SIFI in its ability to attenuate the impacts, process errors and controls
517 on the evaluation of the method performances, three datasets captured from three
518 satellite sensors were prepared. The first dataset was acquired by Advanced
519 Spaceborne Thermal Emission and Reflection Radiometer (ASTER) at the BJM on
520 August 31, 2004. This dataset includes the spectral reflectance and the associated LST
521 product (AST08, with the spatial resolution of 90 m). The original LSTs (90 m) were
522 aggregated into coarse resolutions, on grids of 100, 200, 400, 800, and 1000 m. The
523 second dataset was obtained by MODIS. It was acquired simultaneously with the
524 ASTER data, and it primarily includes the MODIS LST product (MOD11A1, with a
525 resolution of 1000 m). The third dataset was acquired by the Thematic Mapper
526 (Landsat-5) at the HNP on September 22, 2009. The associated LSTs (with a
527 resolution of 120 m) were retrieved by the mono-window algorithm (Qin et al., 2001)
528 and were further aggregated into coarse resolution datasets, on grids of 200, 400, 800,
529 and 1000 m.

530 To validate the SIFI's ability to remove the temperature retrieval errors (for PTL
531 #2), the upscaled ASTER LSTs at a resolution of 1000 m were systematically shifted
532 by a constant value (including 1.0, 2.0, and 3.0 K) and were then disaggregated into
533 200 m, which was then compared with the upscaled ASTER LSTs at 200 m. In
534 addition, the MODIS LSTs at 1000 m were also disaggregated into 200 m and were
535 referenced to the upscaled ASTER LSTs at 200 m. To illustrate the SIFI's
536 independence from the thermal contrast control (for PTL #3), the upscaled ASTER
537 and TM LSTs at 1000 m over both the BJM and HNP, where the thermal contrasts

538 differ, were disaggregated into 200 m and compared with the corresponding reference
539 LSTs at 200 m. To show the SIFI's competency at excluding the temperature unit
540 control (for PTL #3), the upscaled 1000-m TM (band 6) thermal radiance (unit:
541 $W \cdot m^{-2} \cdot \mu m^{-1} \cdot sr^{-1}$) and LSTs (unit: K) were both disaggregated and compared with the
542 reference 200-m radiance and LSTs, respectively. To show the SIFI's ability to
543 attenuate the resolution ratio control (for PTL #4), the upscaled ASTER LSTs at 1000,
544 800, and 400 m were disaggregated into 200 and 100 m and compared with the
545 reference fine resolution LSTs. To show the SIFI's ability to interpret the sharpening
546 statuses (for PTL #5), the upscaled ASTER LSTs at 1000 m were also disaggregated
547 into 200 m using various DSLT methods.

548

549 **5.2. Generation of a series of DLST methods**

550 The performances of the DLST methods primarily depend on the chosen scaling
551 factors and regression tool as well as the window size used for regression (Zhan et al.,
552 2013). This study employed a series of scaling factors and moving window sizes to
553 generate a large number of DLST methods with different performances, while the
554 regression tool was kept unchanged during the evaluation process, and it was
555 designated the quadratic function (Kustas et al., 2003). We acknowledge that
556 advanced regression tools, such as the support vector machine, are usually able to
557 produce better disaggregation results than simple polynomial functions
558 (Keramitsoglou et al., 2013; Ghosh & Joshi, 2014). Nonetheless, the aim of this
559 article is to evaluate method performances rather than to develop high-accuracy
560 methods.

561 The following scaling factors or their combinations were used: the normalized
562 difference water index (NDWI) (McFeeters, 1996), the panchromatic band (Liu &

563 Moore, 1998), the normalized difference vegetation index (NDVI) (Kustas et al.,
564 2003), the vegetation fraction (f_v) (Agam et al., 2007), the emissivity (ε) (Nichol,
565 2009), the product $f_v \cdot \varepsilon$ (Stathopoulou & Cartalis, 2009), the albedo (Dominguez et al.,
566 2011), the normalized multi-band drought index (NMDI) (Liu & Zhu, 2012), the
567 normalized difference built-up index (NDBI) (Wang et al., 2014), and all
568 multi-spectral bands of ASTER or TM (Ghosh & Joshi, 2014). With the local
569 regression strategy (Gao et al., 2017), the moving-window sizes ranging from 3×3 to
570 21×21 pixels were employed.

571

572 **5.3. Validation strategy**

573 Three indexes commonly used for assessments were employed for comparison.
574 They were the RMSE, ERGAS, and Q, with ERGAS and Q being calculated using the
575 following equations:

$$576 \begin{cases} \text{ERGAS} = 100 \cdot (L_r / L_b) \cdot \sqrt{\frac{\text{RMSE}(d, r)}{\mu_r}} \\ Q = 4\sigma_{dr} \mu_d \mu_r (\sigma_d^2 + \sigma_r^2)^{-1} \cdot (\mu_d^2 + \mu_r^2)^{-1} \end{cases} \quad (10)$$

577 where L_b and L_r are the spatial resolutions of the background coarse resolution LSTs
578 (i.e., b) and the reference fine resolution LSTs (i.e., r); σ_{dr} is the covariance between
579 the disaggregated LSTs (i.e., d) and r ; μ_d and μ_r are the associated means; and σ_d and
580 σ_r are the associated standard deviations. For the RMSE and ERGAS, their values
581 range from zero to positive infinity, in theory. Their values are zero once the best
582 results have been achieved, while their values become greater with poorer results. For
583 Q, its values change from -1.0 to 1.0, with 1.0 indicating the best obtained results
584 (Wang & Bovik, 2002). This study did not consider the mean absolute error (MAE)
585 and the structural similarity index measure (SSIM) that were used for DLST

586 (Rodriguez-Galiano et al., 2012) because these two factors have a parallel
587 performance with the RMSE and Q, respectively.

588 We used three strategies to validate the feasibility of the SIFI. The first was
589 through simple mathematical simulation tests that only include a small number of
590 pixels (refer to Section 6.1); the second was by using actual thermal data (refer to
591 Section 6.2); and the third was through conceptual comparisons of the functionality
592 and design philosophy among different indexes (refer to Section 6.3). Validations
593 based on real thermal data can be further divided into two relatively separated parts.
594 In the first part, different indexes were compared under scenarios that correspond to
595 the proposed protocols (refer to Sections 6.2.1 to 6.2.4). The second part (refer to
596 Section 6.2.5) compared the different indexes through human visual interpretations
597 (HVIs). The HVI has been demonstrated to be plausible and is widely recognized to
598 obtain a relatively accurate image quality from the human visual perspective (Wang &
599 Bovik, 2002). Twenty-two graduate students majoring in remote sensing were
600 recruited and subsequently asked to assign ranks independently for a group of
601 disaggregated LST images (the image with a better quality has a higher score). The
602 quality of a specific LST image was then calculated by averaging all 22 ranks for this
603 specific image. According to the calculated image qualities, the HVI ranks of a series
604 of LST images were finally designated using positive integers, with the higher rank
605 indicating the better disaggregation result.

606

607

608

609

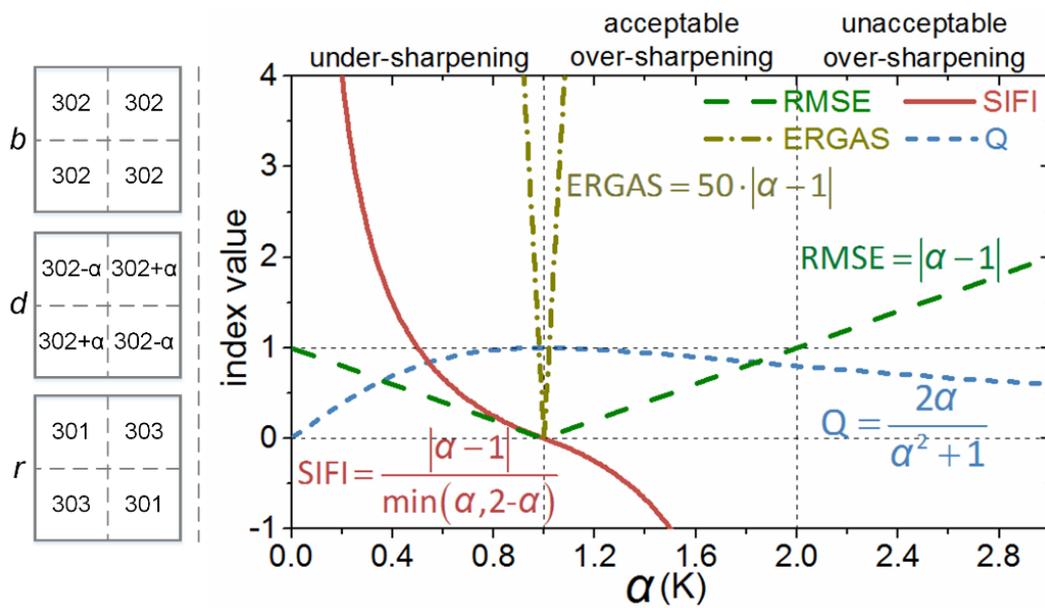
610

611 **6. Results and discussion**

612 **6.1. Comparisons based on simple simulation tests**

613 To explain the differences among the RMSE, ERGAS, Q, and SIFI, we present
 614 two simple simulation tests here. In the first test, let us consider a single pixel to be
 615 disaggregated into half its original resolution, i.e., this pixel is disaggregated into four
 616 subpixels (Fig. 5). Let us further assume that the coarse pixel has an LST value of 302
 617 K, while the values of the four subpixels, from left to right and from above to below,
 618 are $302 - \alpha$, $302 + \alpha$, $302 + \alpha$, and $302 - \alpha$ (unit: K; $\alpha > 0$, and it reflects the added
 619 thermal detail), with 301, 303, 303, and 301 K being the actual values. The variations
 620 in the RMSE, ERGAS, Q, and SIFI as a function of α are provided in Fig. 5.

621



622

623 **Fig. 5.** Variations in the RMSE, ERGAS, Q, and SIFI as a function of the added
 624 thermal detail (quantified by α). b , d , and r represent the original coarse resolution,
 625 disaggregated, and reference fine resolution LSTs (unit: K), respectively; together,

626 they represent a simple DLST process in which a single pixel is disaggregated into
627 four subpixels.

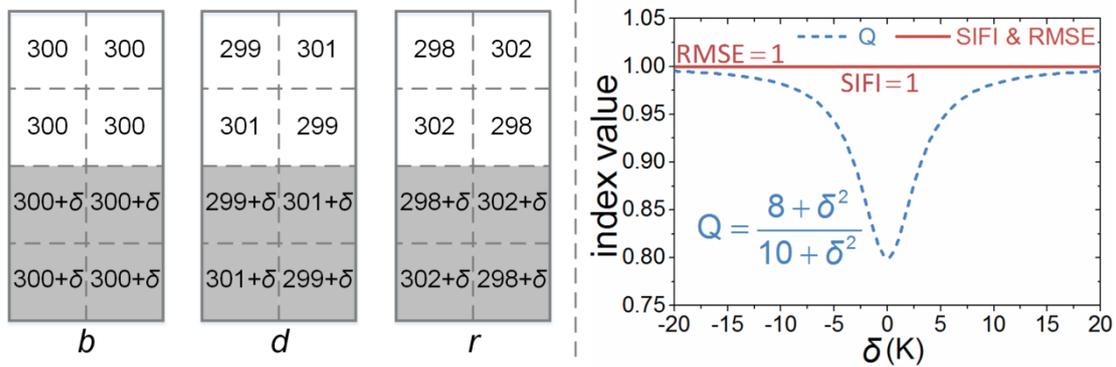
628

629 This simulation test shows that the natures of the RMSE and ERGAS are similar.
630 Their values both decrease when $0.0 < \alpha \leq 1.0$ and increase when $1.0 < \alpha < +\infty$, and
631 they are both axisymmetric with regard to $\alpha = 1.0$. Q increases from zero to 1.0 when
632 $0.0 < \alpha \leq 1.0$, while it decreases, but in a smoother way, from 1.0 to zero when $1.0 <$
633 $\alpha < +\infty$, indicating its asymmetry with regard to $\alpha = 1.0$. By comparison, the SIFI
634 changes from positive infinity to zero ($0.0 < \alpha \leq 1.0$) and then to negative infinity (1.0
635 $< \alpha < 2.0$). SIFI is centrosymmetric when $0.0 < \alpha < 2.0$ but is set as NaN when $\alpha \geq 2.0$.
636 When compared with the RMSE, ERGAS, and Q, the SIFI differ in the following
637 three regards: first, when compared with Q, the symmetry of SIFI, RMSE, and
638 ERGAS shows that SIFI assigns more importance to the quantitative differences
639 between LST images. Second, the values of the three commonly used indexes are
640 unable to indicate the sharpening statuses; however, the calculated SIFI is capable of
641 this type of indication. The background LSTs have been under-sharpened when $0.0 <$
642 $\alpha < 1.0$ (SIFI > 0.0), acceptably over sharpened when $1.0 \leq \alpha < 2.0$ (SIFI < 0.0), and
643 unacceptably over-sharpened when $\alpha \geq 2.0$ (SIFI = NaN). Third, SIFI is especially
644 sensitive to the case in which DLST is poorly performed because its value rapidly
645 increases when α is close to zero or two, indicating that SIFI would be more suitable
646 for differentiating poor disaggregation results from minor in-between differences.

647 In the second test, let us consider two pixels to be disaggregated into half of their
648 original resolutions. Two coarse resolution LSTs are disaggregated into eight fine
649 resolution LSTs. The pixel values of the original coarse resolution, disaggregated, and

650 reference fine resolution LSTs are shown in Fig. 6, where δ is a variable that reflects
 651 the thermal contrast between these adjacent coarse resolution pixels.

652



653

654 **Fig. 6.** Variations in the RMSE, Q and SIFI as a function of the thermal contrast
 655 between two adjacent pixels, which is represented by δ (unit: K). *b*, *d*, and *r* are the
 656 original coarse resolution, disaggregated, and reference fine resolution LSTs,
 657 respectively.

658

659 The simulation results in Fig. 6 illustrate that Q is a function of δ , while RMSE
 660 and SIFI are identically equal to 1.0. Here, the ERGAS is not included due to having
 661 similar properties to the RMSE in this case. These simulations again indicate that SIFI
 662 and RMSE are more highly related to the absolute quantitative differences between
 663 images, while Q varies with the thermal contrast δ . For the assessment of images
 664 quality that is specifically perceived through human visualization, the quantitation
 665 property (the absolute difference between images) is sometimes unimportant because
 666 it will not affect the human interpretation of images (Wang & Bovik, 2002). However,
 667 the DLST as commonly shown is used to assist the detailed analysis of the associated
 668 quantitative applications such as the upscaling of *in situ* data to the pixel level, urban
 669 thermal environment mapping (Zhou et al., 2013), or evapotranspiration estimation
 670 (Anderson et al., 2012). These applications require that an index should remain

671 consistent even if the thermal contrast among adjacent pixels varies. From this
672 viewpoint, an index is applicably better when it is invariant with the thermal contrast,
673 and therefore the RMSE and SIFI in this case are more suitable than the Q for
674 quantitative applications of DLST.

675

676 **6.2. Comparisons based on real thermal data**

677 **6.2.1. Scenario 1 corresponding to PTL #2**

678 Under this scenario, the capabilities of RMSE, ERGAS, Q, and SIFI are
679 compared when there are temperature retrieval errors (corresponding to PTL #2).
680 Table 1 offers the index values for the cases with various systematic LST retrieval
681 errors. The results show that RMSE and ERGAS vary according to the added
682 systematic error (Δ), and specifically, the RMSE increases from 2.17 to 3.73 K for
683 Cases #1 to #4, while the Q and SIFI remain unchanged for these four cases. This
684 finding demonstrates that the RMSE and ERGAS highly depend on systematic
685 temperature retrieval errors, while Q and SIFI are insensitive to such an error. These
686 results reveal that Q and SIFI are better for evaluating model performances than
687 RMSE and ERGAS because model performances should have been unrelated to the Δ
688 (i.e., temperature retrieval error). The less feasibility by RMSE and ERGAS is also
689 evident by comparing Cases #3 and #5, wherein RMSE and ERGAS are consistent. In
690 theory, the performance should have been worse in Case #5 than that of Case #3
691 because the scaling factor used for Case #5 is from ASTER, which has co-registration
692 errors with MODIS LST, making the DLST method for Case #5 not well.

693

694 **Table 1.** Comparisons of index values when biases in temperature retrieval occur.

Cases*	TOD**	RMSE (K)	ERGAS	Q	SIFI
#1	AST + 0	2.17	0.14	0.88	1.06
#2	AST + 1	2.41	0.16	0.88	1.06
#3	AST + 2	2.98	0.19	0.88	1.06
#4	AST + 3	3.73	0.24	0.88	1.06
#5	MOD	2.98	0.19	0.82	2.31

695 * Cases #1 to #5 all used a completely consistent DLST approach with the NMDI as
696 the scaling factor and 7×7 as the moving window size. LSTs were all disaggregated
697 from 1000 to 200 m over the BJM.

698 ** TOD stands for ‘type of data’. For ‘AST + Δ’, the 1000-m LSTs were upscaled
699 from the 200 m ASTER/LSTs, while the validation data were the combination of the
700 200-m ASTER/LSTs and a systematic error of Δ (unit: K). For ‘MOD’, the 1000-m
701 LSTs were the MODIS/LSTs, while the validation data were the 200 m
702 ASTER/LSTs.

703

704 **6.2.2. Scenario 2 corresponding to PTL #3**

705 Under this scenario, the indexes are compared when the thermal contrast and
706 temperature units differ (corresponding to PTL #3). The index values over areas with
707 different thermal contrasts are provided in Table 2. The results include three pairs
708 (Cases #1 and #2, Cases #3 and #4, and Cases #5 and #6), each having an identical
709 RMSE by using different combinations of scaling factors and moving window sizes.
710 For each pair, the RMSE and ERGAS have almost identical values, whereas the Q or
711 SIFI show a different behavior. The Q and SIFI values indicate that the disaggregation
712 over the region with a higher thermal contrast (i.e., the BJM, and its thermal contrast
713 defined by standard deviation is 4.4 K) achieves a better result. This interpretation is
714 reasonable because the specific DLST method should possess a better performance

715 over the regions with higher thermal contrasts once the associated RMSE remains
 716 unchanged, e.g., the RMSE is 1.38 K for Cases #1 and #2. In other words, when using
 717 absolute distances between disaggregated and fine resolution LSTs, RMSE and
 718 ERGAS tend to overestimate the model performance over relatively homogeneous
 719 regions with a lower thermal contrast, while they underestimate the performance for
 720 heterogeneous regions.

721 The temperature levels considered here include the at-sensor radiance (unit:
 722 $W \cdot m^{-2} \cdot \mu m^{-1} \cdot sr^{-1}$) and the LST in Kelvin units. The DN and at-sensor brightness
 723 temperature levels were excluded because they have a very significant linear
 724 relationship with the radiance (Barsi et al., 2007), and the LST values in Celsius and
 725 Fahrenheit were excluded due to their linear relationship with the LST in Kelvin. The
 726 results in Table 3 illustrate that the RMSE and ERGAS vary greatly with the
 727 temperature level (unit), whereas the Q and SIFI values show small differences.
 728 Despite behaving differently at different levels, in practice, the performance of a
 729 certain DLST method should not be significantly altered among these levels. These
 730 results suggest that the RMSE and ERGAS are inappropriate for this type of
 731 assessment, while the latter two indexes are a better option.

732

733 **Table 2.** Comparisons of index values over areas with different thermal contrasts.

Cases*	region	contrast ** (K)	RMSE (K)	ERGAS	Q	SIFI
#1	BJM	4.4	1.38	0.09	0.95	0.69
#2	HNP	1.5	1.38	0.09	0.56	1.30
#3	BJM	4.4	1.44	0.09	0.95	0.71
#4	HNP	1.5	1.44	0.10	0.47	1.40
#5	BJM	4.4	1.58	0.10	0.94	0.79

#6	HNP	1.5	1.58	0.10	0.60	0.93
----	-----	-----	------	------	------	------

734 * Cases #1 to #6 used six different DLST methods considering ASTER band 9, NDVI,
735 NDVI, $f_v \cdot \varepsilon$, ASTER band 5, and TM band 4, respectively, as the scaling factor and
736 considering 11×11, 9×9, 21×21, 7×7, 9×9, and 7×7 as the moving window size,
737 respectively. LSTs were all disaggregated from 1000 to 200 m.

738 ** The thermal contrasts for the BJM and HNP, as represented by the standard
739 deviation (σ), are 4.4 and 1.5 K, respectively.

740

741 **Table 3.** Comparisons of index values with different temperature units.

Cases*	unit	RMSE	ERGAS	Q	SIFI
#1	$W \cdot m^{-2} \cdot \mu m^{-1} \cdot sr^{-1}$	0.13	0.34	0.62	1.98
#2	K	1.15	0.08	0.64	1.85
#3	$W \cdot m^{-2} \cdot \mu m^{-1} \cdot sr^{-1}$	0.14	0.34	0.57	2.60
#4	K	1.19	0.08	0.56	2.40

742 * Cases #1 and #2 used a completely consistent DLST approach with NDVI as the
743 scaling factor and 21×21 as the moving window size, while the DLST approach for
744 Cases #3 and #4 was using vegetation fraction as the scaling factor and 21×21 as the
745 moving window size. The disaggregation was performed for the same LST image
746 over the HNP; and LSTs were all disaggregated from 1000 to 200 m.

747

748 **6.2.3. Scenario 3 corresponding to PTL #4**

749 Under this scenario, the indexes are compared when the resolution gap between
750 the background coarse resolution and reference fine resolution LSTs changes
751 (corresponding to PTL #4). A comparison of the index values when the DLST
752 methods are performed with different initial and target resolutions is presented in
753 Table 4, where the target resolution for Cases #1 to #3 is 100 m, while it is 200 m for

754 Cases #4 to #6. Different DLST methods, each with a particular scaling factor and
 755 moving window size, are considered so that the RMSEs between the disaggregated
 756 and the fine resolution LSTs are approximately equivalent or even identical (see Table
 757 4). This type of setting suggests that the image quality of the disaggregated LSTs is
 758 similar when taking the RMSE as the error index. However, this result is problematic,
 759 when considering that the target resolution is kept unaltered but the initial resolutions
 760 changes. Indeed, we expect that the method with the largest resolution gap should
 761 have the best performance for the DLST methods. Q has a similar behavior with
 762 RMSE in that it hardly changes for all the cases.

763 Instead, the ERGAS is suitable for these model performance evaluations by
 764 considering the spatial resolutions of the pre- and post-disaggregated values (Wald et
 765 al., 1997). This idea is as well evidenced by Table 4; with the same RMSE and target
 766 resolution, the case with a larger resolution gap points to a better model performance,
 767 which is confirmed by the ERGAS values. Note that the SIFI is consistent with the
 768 ERGAS in these tests, showing a similar capability to evaluate the model performance
 769 under this scenario. The SIFI nevertheless differs from the ERGAS in that the
 770 decrease rate of SIFI is not proportional to the resolution ratios (a key variable in
 771 ERGAS) — it decreases in a smoother way (see Table 4).

772

773 **Table 4.** Comparisons of index values with different resolution ratios between pre-
 774 and post-disaggregated LSTs.

Cases*	Resolution (m)**	RMSE (K)	ERGAS	Q	SIFI
#1	200→100	1.77	0.29	0.93	1.16
#2	400→100	1.76	0.14	0.92	0.97
#3	800→100	1.76	0.07	0.93	0.82

#4	400→200	1.76	0.29	0.92	1.04
#5	800→200	1.76	0.14	0.92	1.12
#6	1000→200	1.76	0.11	0.92	0.78

775 * Cases #1 to #6 used the ASTER band 11, vegetation fraction, ASTER band 5,
776 albedo, NMDI, and ASTER band 2, respectively, as the scaling factor and 3×3, 7×7,
777 21×21, 3×3, 9×9, and 7×7 as the moving window size, respectively. DLST was
778 performed over the BJM.

779 ** The numbers on the left and right of the ‘→’ are the resolutions of the original
780 coarse resolution and the disaggregated LSTs, respectively.

781

782 **6.2.4. Scenario 4 corresponding to PTL #5**

783 Under this scenario, the index values are compared when different amounts of
784 thermal details are added to the background LSTs during the DLST process
785 (corresponding to PTL #5). The thermal details were controlled by a multiplier, which,
786 together with the regression coefficients acquired from the relationships between the
787 background LSTs and scaling factors, determines the amount of added details (Zhan et
788 al., 2011). Table 5 illustrates the associated indexed values when different amounts of
789 thermal details were added by varying the multiplier from 0.5 to 1.9. The results show
790 that the RMSE, ERGAS, and Q change with the multiplier, but their values are unable
791 to specify the sharpening statuses. For example, it is very difficult to judge the
792 boundaries among the three sharpening statuses only according the RMSE values,
793 which alter from approximately 1.3 to 3.3 K for Cases #1 to #8. By contrast, the SIFI
794 values are indicative of the sharpening statuses, evidently specifying that the LSTs are
795 under-sharpened for Cases #1 to #4, acceptably over-sharpened for Cases #5 and #6,
796 and unacceptably over-sharpened for Cases #7 and #8. In detail, Case #1 is designated
797 under-sharpened because the $rmse(D, B)$ (the calculated value is 0.28) < $rmse(D, B_R)$

798 (the value is 0.74), while Case #5 shows acceptable over-sharpening because $\text{rmse}(D, B_R)$ (0.62) < $\text{rmse}(D, B)$ (0.74) < $\text{rmse}(B, B_R)$ (0.93), and Case #7 has unacceptable
 799 $\text{rmse}(D, B)$ (0.96) > $\text{rmse}(B, B_R)$ (0.93).
 800
 801

802 **Table 5.** Comparisons of index values with different sharpening statuses.

Cases*	multiplier	RMSE (K)	ERGAS	Q	SIFI	status**
#1	0.5	1.32	0.43	0.96	1.12	USP
#2	0.7	1.31	0.43	0.96	0.80	USP
#3	0.9	1.46	0.48	0.95	0.69	USP
#4	1.1	1.74	0.56	0.93	0.67	USP
#5	1.3	2.08	0.68	0.91	-0.80	AOS
#6	1.5	2.47	0.80	0.88	-0.92	AOS
#7	1.7	2.88	0.94	0.85	NaN	UOS
#8	1.9	3.31	1.08	0.81	NaN	UOS

803 * Cases #1 to #6 all used NDVI as the scaling factor and the statistical regression
 804 between LST and NDVI was conducted in a global window (i.e., the entire image);
 805 the DLST was performed over the BJM from 1000 to 200 m. Note that the method
 806 performances for Cases #1 to #6 are determined by the multiplier coefficient (varying
 807 from 0.5 to 1.9), which determines the amount of thermal details that are added to the
 808 background LSTs.

809 ** USP, AOS, and UOS denote the under-sharpening, acceptable over-sharpening, and
 810 unacceptable over-sharpening, respectively.

811

812 **6.2.5. Comparison reference to human visual interpretations**

813 Under this scenario, the compared index values when referencing human visual
 814 interpretations (HVIs) and with respect to different DLST methods (as specified by

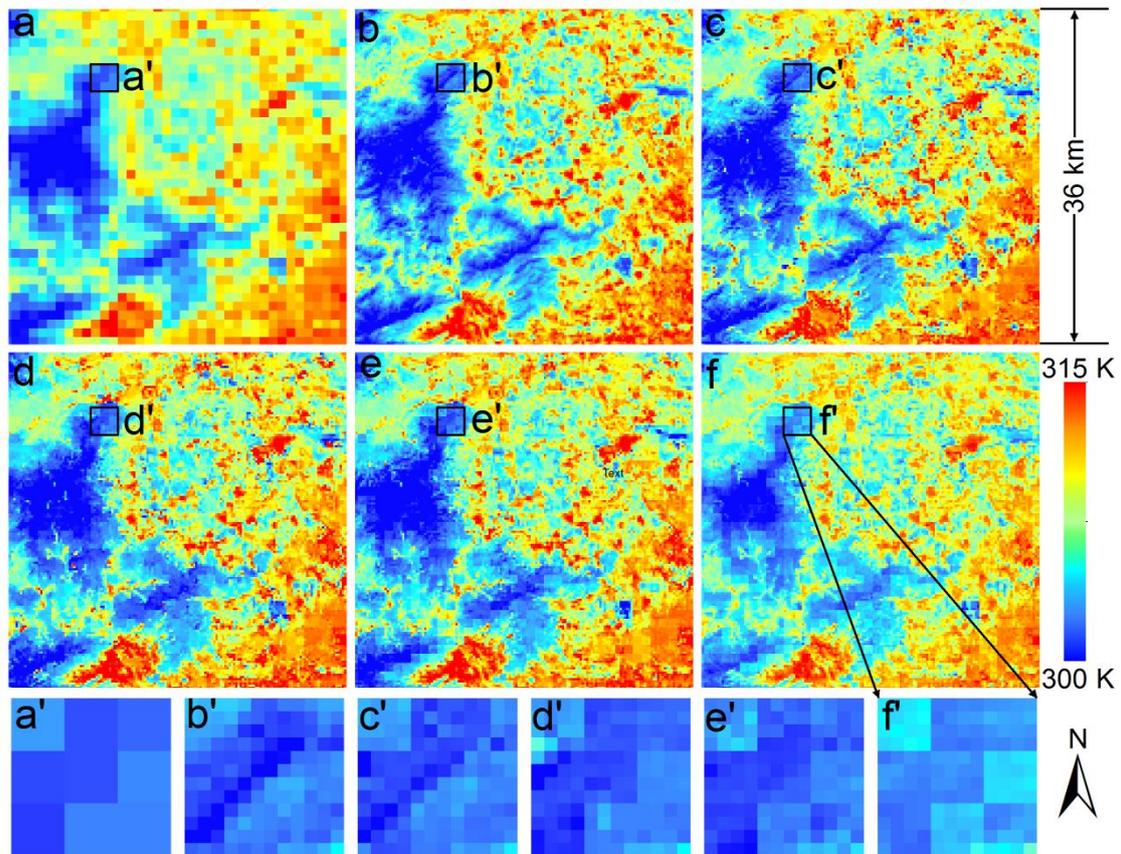
815 dissimilar scaling factors and moving window sizes) are shown in Table 6. The HVI
816 ranks are also reported based on the disaggregated LSTs given in Fig. 7. The HVI
817 ranks are higher once the image quality is better. These results show that RMSE and
818 ERGAS are inaccurate for these assessments. The Q and the HVI ranks are sometimes
819 inconsistent. For example, the disaggregated LSTs for Case #3 possess the highest Q
820 (i.e., the best image quality or method performance). However, the corresponding
821 LST image (see Fig. 7e) is not the best among the four disaggregated LSTs; its block
822 effect (also termed the grid effect) is considerably more distinct than the LST image
823 for Case #1 (Fig. 7c). By contrast, the estimated SIFI values are consistent with the
824 HVI ranks with no exception, with a lower SIFI corresponding to a higher HVI rank.
825

826 **Table 6.** Quantitative comparison between SIFI and other indexes for various DLST
827 methods as represented by different scaling factors and moving window sizes.

Cases*	RMSE (K)	ERGAS	Q	SIFI	HVI Rank
#1	1.530	0.124	0.941	0.774	4
#2	1.611	0.131	0.935	0.791	3
#3	1.496	0.122	0.943	0.840	2
#4	1.491	0.121	0.933	0.868	1

828 * Cases #1 to #4 used four distinct DLST methods when considering the ASTER band
829 1, NDVI, NDWI, and $f_v \cdot \varepsilon$, respectively, as the scaling factor and 7×7 , 5×5 , 9×9 , and
830 21×21 as the moving window sizes. The LSTs were all disaggregated from 800 to 200
831 m. Note that three significant digits are kept for the index values in this table in
832 particular.

833



834

835 **Fig. 7.** Coarse resolution, disaggregated, and fine resolution LSTs used for
 836 comparison. (a) and (b) are the coarse and fine resolution LSTs, respectively; and (c)
 837 to (g) are the disaggregated LSTs corresponding to Cases #1 to #4 in Table 6.

838

839 **6.3. Conceptual comparisons of the index functionality and structure**

840 The aforementioned results were based on mathematical simulations and real
 841 data, and they show that the SIFI assimilates the features of the RMSE, ERGAS, and
 842 Q, and abides by all five protocols (see Table 7). For SIFI, once the Euclidean
 843 distance is used as the metric as employed in this study, it has the same trait as the
 844 RMSE by emphasizing the absolute difference between images (partly corresponding
 845 to PTL #1), which is reflected in their symmetry between the under-sharpening and
 846 acceptable over-sharpening cases. SIFI also incorporates the trait from Q by
 847 combining a normalization process (referred to as $g_d(x)$). This incorporation helps

848 SIFI be independent of a great portion of the temperature retrieval errors
849 (corresponding to PTL #2) and alleviate the thermal contrast and temperature unit
850 controls (corresponding to PTL #3). SIFI further integrates the ERGAS trait by
851 compensating for the resolution difference between the pre- and post-disaggregated
852 LSTs (PTL #4). The former item achieves this objective by using a triple comparison
853 function (refer to $g_t(x)$), while the latter employs the ratio between the coarse and fine
854 resolutions. SIFI is additionally able to identify the three sharpening statuses through
855 a piecewise function with three different constraints (refer to $g_p(x)$ and corresponding
856 to PTL #5). By contrast, indexes such as the RMSE, ERGAS, and Q comply with only
857 a part of the protocols for DLST (see Table 7). For example, RMSE just partly meets
858 the requirement of PTL #1; it portrays the differences between the disaggregated and
859 reference LSTs but disregards those between the disaggregated and background LSTs;
860 ERGAS only complies with PTLs #1 and #4; and Q is less competent when
861 considering the requirements involved in PTLs #4 and #5.

862 In addition, SIFI further incorporates the local mean, i.e., the local details, by
863 combining with the detail-retrieval procedure (refer to $g_d(x)$). In this way, the SIFI
864 provides quantitative assessments on the intensity of the block effect in the
865 disaggregated LSTs (Anderson et al., 2011; Zhan et al., 2013). However, the other
866 three indexes only estimate the global statistical variables of the compared images
867 (e.g., the mean and covariance), which render fewer local details.

868 **Table 7.** Conceptual comparisons on index functionality in reference to the proposed
869 five protocols.

Protocol*	RMSE	ERGAS	Q	SIFI
PTL #1	✓**	✓	✓	✓✓
PTL #2			✓✓	✓✓

PTL #3		✓✓	✓✓
PTL #4	✓✓		✓✓
PTL #5			✓✓

870 * PTL #1 requires that the index should simultaneously measure the similarity
871 between the disaggregated and reference LSTs as well as the dissimilarity between the
872 disaggregated and background LSTs. PTLs #2 to #5 correspond to the indexes that
873 should be independent of the temperature retrieval error (PTL #2), thermal contrast
874 and temperature unit control (PTL #3), and resolution ratio control (PTL #4), and
875 should be indicative of the sharpening statuses (PTL #5).

876 ** ‘✓’ indicates that the index abides by some of the requirements of the specific
877 protocol, while ‘✓✓’ indicates that the index completely complies with the protocol.

878

879 Researchers often use an indirect validation strategy; LSTs are first upscaled to
880 coarser resolutions, which are then disaggregated again to the original fine resolution,
881 at which point an intercomparison becomes possible (Agam et al., 2007;
882 Rodriguez-Galiano et al., 2012). For this type of validation, the temperature retrieval
883 errors vanish because the coarse resolution LSTs and the fine resolution LSTs used for
884 validation are from an identical source. When this validation strategy is used, simple
885 indexes such as the RMSE are mostly feasible for comparing methods with great
886 differences in performance.

887 For communications in the DLST community, one may need to judge the
888 performance of a single method many times simply through a single value, and the
889 SIFI will help for this case. Although the SIFI as illustrated here has shown many
890 advantages, in practice, one may also need to know the absolute differences (e.g., the
891 RMSE) between the disaggregated and reference LSTs for practical applications such

892 as the remote sensing of surface fluxes. For example, the widespread use of the Taylor
893 diagram to evaluate the predicated and the reference geophysical variables underlines
894 the importance of summarizing various aspects of the model performance by plotting
895 the RMSE, standard deviation, and correlation coefficient in a diagram (Taylor, 2001).
896 From this perspective, we therefore recommend these indexes, such as the RMSE,
897 ERGAS, Q, correlation coefficient, and SIFI, even along with the estimated distances
898 for calculating the SIFI (i.e., $m(D, R)$, $m(D, B)$, $m(D, B_R)$, and $m(B, B_R)$), which are
899 used collectively to assess a newly proposed method or compare several methods for
900 DLST. For these reasons, we believe the maximum benefit ultimately lies in this
901 approach.

902

903 **6.4. Problems and prospects**

904 (1) *Problems*

905 First, the procedures used in the design of the SIFI are able to remove (or
906 alleviate) the linear and systematic process errors/controls; these procedures are
907 nevertheless unable to eliminate the random and highly nonlinear process errors, e.g.,
908 the mismatch between the fine resolution scaling factors and the coarse and fine
909 resolution LSTs (i.e., the co-registration error). Practitioners need to be careful to
910 interpret the index values because the co-registration error is fickle. Second, the initial
911 scaling factors always possess a spatial resolution that is much higher than that of fine
912 resolution LSTs. The spatial upscaling of the scaling factors to the resolution of the
913 fine LSTs should consider the point spread function of the sensor (Zhan et al., 2013).
914 The aim of this consideration is to make sure that the true resolution of scaling factor
915 do have the same resolution with the fine LSTs. Otherwise, the RMSE between the

916 disaggregated and reference LSTs will no longer be zero, even if the temperature
917 retrieval as well as the DLST processes are error-free.

918 (2) *Prospects*

919 The SIFI proposed in this study is only one alternative that satisfies the proposed
920 protocols. We need to clarify that other strategies that conform to the protocols can
921 also be applied to help design an index even better than SIFI. First, SIFI employs the
922 strategies given by Eqs. (3) and (4) to remove the linear or the locally/globally
923 systematic process errors and controls. Other normalization schemes that are able to
924 remove these associated errors and controls are also feasible. Second, this study
925 mainly uses the Euclidean distance (i.e., the RMSE) as the metric for calculating SIFI.
926 It is expected that distance metrics such as the general Minkowski distance (Han et al.,
927 2011) may generate a parallel capability for assessments. Nevertheless, metrics such
928 as the Euclidean or Minkowski distances give a high weight to outliers and may make
929 the resultant SIFI less indicative of method performances. Therefore, researchers
930 should try to avoid outliers through setting thresholds during the DLST. One may
931 infer that distance metrics that are insensitive to outliers (e.g., the angle cosine
932 distance) are applicable for performance assessments. However, researchers need be
933 very careful to use such metrics because they more emphasize the structural similarity
934 between two images, that is, they ignore the information retained in the absolute
935 values between images, which is yet important for the quantitative applications of
936 DLST. Third, SIFI is calculated for an entire image, i.e., a single SIFI value is
937 calculated for a single evaluation. SIFI may be modified to be dependent on pixel
938 location – a series of local SIFI values can be obtained by setting moving windows on
939 a LST image. By this modification, method performances can be evaluated for
940 different parts within a single image.

941 SIFI has potential to be further applied to disaggregation/downscaling
942 assessments. The recent rise of the spatio-temporal DLST requires the assessment of
943 sequential fine resolution LSTs rather than a LST image at a single moment. Such
944 assessments, therefore, may be performed by combining the sequential SIFIs within a
945 certain cycle (e.g. the diurnal cycle) (Göttsche & Olesen, 2009). SIFI may be further
946 enhanced to facilitate the assessments of method performances when *in situ*
947 measurements on LST are available. Finally, SIFI and the associated design
948 philosophies may also be used to assess model performances for disaggregation of
949 other satellite products such as precipitation and soil moisture.

950

951

952

953 7. Conclusions

954 At present, the performances of DLST methods are evaluated by simple indexes
955 (e.g., RMSE) or more complicated ones that are adapted from optical image fusion
956 (e.g., ERGAS and Q). These indexes are insufficient because not only do they include
957 all the errors involved in the complete process from thermal radiance to LSTs (termed
958 the process error) but also because they are susceptible to process controls, including
959 differences in the thermal contrast, temperature units, and resolution ratios. In
960 addition, these indexes are unable to differentiate among the three sharpening statuses,
961 the under-sharpening, acceptable over-sharpening, and unacceptable over-sharpening
962 statuses. These deficiencies make evaluating the performance of the DLST methods
963 far from precise under different scenarios. It is therefore of great urgency to design a
964 better index for these evaluations.

965 In considering this issue, five standard protocols were proposed with which a
966 suitable index should be assigned. In being guided under these protocols, a simple yet
967 flexible index (SIFI) was designed. SIFI incorporates the following four procedures:
968 (1) the detail-retrieval procedure $g_d(x)$ that is primarily used to remove the impacts
969 from the temperature retrieval error; (2) the Gaussian normalization $g_n(x)$ primarily
970 aimed at attenuating the controls on the differences in the thermal contrasts and
971 temperature units; (3) the triple comparison $g_t(x)$, which is scheduled to lessen the
972 controls on the difference in resolution ratios; and (4) the piecewise comparison $g_p(x)$
973 and several provisos (given by several inequalities to indicate the three sharpening
974 statuses).

975 Comprehensive evaluations show that indexes that include the RMSE, ERGAS,
976 and Q abide by only part of the requirements denoted in the five protocols. The new

977 index SIFI instead complies with all the proposed protocols. This SIFI is able to
978 capture the model performance more accurately; it can remove the impacts from the
979 process errors and controls on evaluations and can indicate the sharpening statuses
980 such that a disaggregation is under- or over-sharpened. Further analysis illustrates that
981 the SIFI attaches more importance to the scenario in which the DLST is poorly
982 performed and therefore is sensitive to the grid effect in the DLST. Note that it
983 remains difficult for the SIFI to remove the highly nonlinear process errors, such as
984 the mismatch error, and there may be better procedures than those used in this study.
985 Nevertheless, SIFI facilitates the comparison of model performances and therefore
986 helps in the further enhancement of methods for DLST. In addition, we believe that
987 the design philosophies of the SIFI are likely applicable to the model performance
988 comparisons for the disaggregation of other geophysical variables.

989

990

991

992 **Acknowledgements**

993 This work is jointly supported by the Key Research and Development Programs
994 for Global Change and Adaptation under Grant number 2016YFA0600201, the
995 National Natural Science Foundation of China under Grant number 41671420, and the
996 Key Research and Development Programs for Global Change and Adaptation under
997 Grant number 2017YFA0603604. We are also grateful for the financial support
998 provided by the DengFeng Program-B of Nanjing University.

999

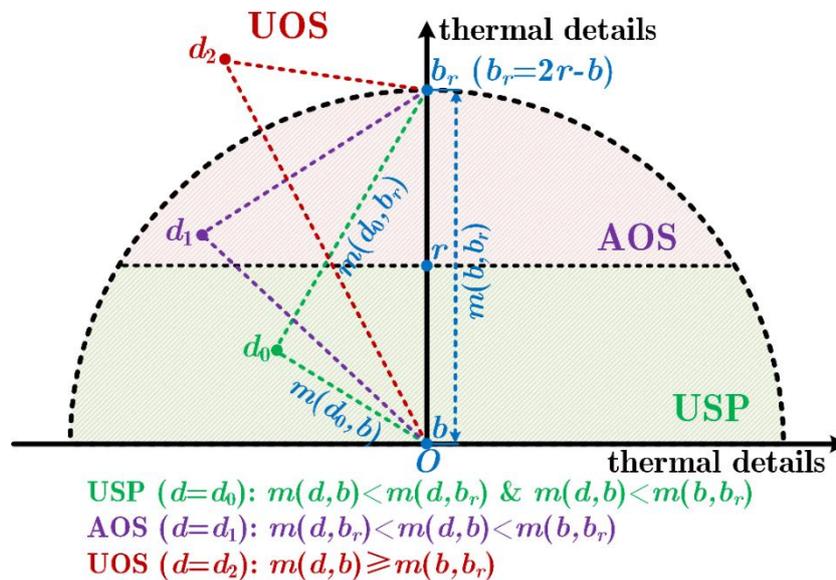
1000

1001

1002 **APPENDIX A: ILLUSTRATION OF THE THREE SHARPENING STATUSES IN**
 1003 **HIGH-DIMENSION**

1004 The conceptual sharpening statuses for a single pixel block (i.e., a single
 1005 dimension) was provided in Section 2.3. However, disaggregation of LSTs is
 1006 conducted for all the pixels of an entire LST image rather than a single pixel block. In
 1007 other words, the precise three sharpening statuses should be displayed in a
 1008 high-dimension. Fig. A1 demonstrates the conceptual illustration of the
 1009 *under-sharpening* (USP), *acceptable over-sharpening* (AOS), and *unacceptable*
 1010 *over-sharpening* (USP) in two dimensions, and illustration of even higher dimensions
 1011 is similar to the two-dimension case.

1012



1013

1014 **Fig. A1.** Two-dimensional conceptual description of the three sharpening statuses. The
 1015 coarse resolution, disaggregated, and fine resolution LST images are denoted by the
 1016 three dots b , d , and r , respectively; b_r is the mirror image of b when using r as the
 1017 center of symmetry. d_0 , d_1 , and d_2 represent the cases of *under-sharpening* (USP),

1018 *acceptable over-sharpening* (AOS), and *unacceptable over-sharpening* (USP),
1019 respectively. $m(d, r)$, $m(d, b)$, $m(d, b_r)$, and $m(b, b_r)$ are the distances between the
1020 associated two LST images.

1021

1022 The constraints for differentiating the statuses in the two-dimension are
1023 consistent with those in the one-dimension. (1) LSTs are under sharpened ($d = d_0$)
1024 when $m(d, b) < m(d, b_r)$ & $m(d, b) < m(b, b_r)$, and therefore the USP can be
1025 geometrically represented by the light green region (refer to Fig. A1). (2) LSTs are
1026 acceptably over sharpened ($d = d_1$) when $m(d, b) > m(d, b_r)$ & $m(d, b) < m(b, b_r)$, and
1027 therefore the AOS corresponds to the light red region. (3) LSTs are unacceptably over
1028 sharpened ($d = d_2$) when $m(d, b) > m(b, b_r)$, which corresponds to the region beyond
1029 the semicircle.

1030

1031

1032

1033 **APPENDIX B: PROOF OF THE SIFI'S INDEPENDENCE OF LINEAR AND SYSTEMATIC**
1034 **ERRORS OR CONTROLS**

1035 The standard definition of SIFI given by Eq. (6) is able to remove the linear
1036 and/or systematic process errors using the Gaussian normalization. Let us consider
1037 two variables T_1 and T_2 , and there is a linear relationship between these two variables,
1038 given as:

$$1039 \quad T_1 = a_1 \cdot T_2 + a_0 \quad (\text{B1})$$

1040 where a_1 and a_0 are the linear coefficient and constant. The Gaussian normalization of
1041 T_1 can be given as follows:

$$1042 \quad g_n(T_1) = \frac{T_1 - u_1}{\sigma_1} = \frac{(a_1 \cdot T_2 + a_0) - (a_1 \cdot u_2 + a_0)}{a_1 \cdot \sigma_1} = \frac{T_2 - u_2}{\sigma_2} = g_n(T_2) \quad (\text{B2})$$

1043 where $g_n(\cdot)$ is a normalization function given by Eq. (4); u_1 , u_2 , σ_1 , and σ_2 are the
1044 means and standard deviations for T_1 and T_2 , respectively. We therefore prove that
1045 SIFI is capable of eliminating the linear and systematic process errors and controls
1046 (e.g., a great portion of the temperature retrieval error).

1047

1048

1049

1050 **APPENDIX C: SIMPLIFICATION OF SIFI UNDER PARTICULAR ASSUMPTIONS**

1051 Once the Euclidean distance (i.e., the RMSE) is employed as the metric and the
 1052 over-sharpening does not occur, the standard definition of SIFI given by Eq. (6) can
 1053 be simplified into the following equation:

1054
$$\text{SIFI} = \text{RMSE}(D, R) \cdot [\text{RMSE}(D, B)]^{-1} \quad (\text{C1})$$

1055 where B , D , and R are given by the following:

1056
$$\begin{cases} B = g_n(g_d(b)) = \sigma_b^{-1}(b - \bar{b} - u_{b-\bar{b}}) \\ D = g_n(g_d(d)) = \sigma_b^{-1}(d - \bar{d} - u_{d-\bar{d}}) \\ R = g_n(g_d(r)) = \sigma_b^{-1}(r - \bar{r} - u_{r-\bar{r}}) \end{cases} \quad (\text{C2})$$

1057 where $g_n(x)$ and $g_d(x)$ are given by Eq. (7); b , d , and r are the background coarse
 1058 resolution, disaggregated, and reference fine resolution LSTs, respectively; σ_x , u_x ,
 1059 and \bar{x} are the standard deviation, mean, and the aggregated coarse resolution LSTs
 1060 of a LST image (i.e., x). Combining Eqs. (C1) and (C2), the following equations can
 1061 be deduced:

1062
$$\begin{aligned} \text{SIFI} &= \text{RMSE}(D, R) \cdot [\text{RMSE}(D, B)]^{-1} \\ &= \frac{\text{RMSE}(\sigma_b^{-1}(d - \bar{d} - u_{d-\bar{d}}), \sigma_b^{-1}(r - \bar{r} - u_{r-\bar{r}}))}{\text{RMSE}(\sigma_b^{-1}(d - \bar{d} - u_{d-\bar{d}}), \sigma_b^{-1}(b - \bar{b} - u_{b-\bar{b}}))} \\ &= \frac{\text{RMSE}(d - \bar{d} - u_{d-\bar{d}}, r - \bar{r} - u_{r-\bar{r}})}{\text{RMSE}(d - \bar{d} - u_{d-\bar{d}}, b - \bar{b} - u_{b-\bar{b}})} \end{aligned} \quad (\text{C3})$$

1063 Once the fine and coarse resolution LSTs are from a same source, i.e., the
 1064 disaggregated LSTs are validated by the aggregation-and-then-disaggregation strategy,
 1065 the aggregated coarse resolution LSTs for b , d , and r will be the identical; and the
 1066 mean of the thermal details for b , d , and r will also be equal to zero. We thus have the
 1067 following equations:

1068
$$\begin{cases} \bar{b} = \bar{d} = \bar{r} \\ u_{b-\bar{b}} = u_{d-\bar{d}} = u_{r-\bar{r}} = 0 \end{cases} \quad (C4)$$

1069 Combining Eqs. (C3) and (C4), we obtain the final equation:

1070
$$\text{SIFI} = \frac{\text{RMSE}(d - \bar{d}, r - \bar{r})}{\text{RMSE}(d - \bar{d}, b - \bar{b})} = \frac{\text{RMSE}(d, r)}{\text{RMSE}(d, b)} \quad (C5)$$

1071 The above proof therefore finally demonstrates that $\text{RMSE}(D, R) \cdot [\text{RMSE}(D, B)]^{-1}$

1072 is equivalent to $\text{RMSE}(d, r) \cdot [\text{RMSE}(d, b)]^{-1}$.

1073

1074

1075

1076 **References**

- 1077 Adesso, P., Longo, M., Maltese, A., Restaino, R., & Vivone, G. (2015). Batch
1078 methods for resolution enhancement of TIR Image sequences. *IEEE Journal of*
1079 *Selected Topics in Applied Earth Observations and Remote Sensing*, 8,
1080 3372-3385, doi:10.1109/JSTARS.2015.2440333.
- 1081 Agam, N., Kustas, W. P., Anderson, M. C., Li, F., & Neale, C. M. (2007). A vegetation
1082 index based technique for spatial sharpening of thermal imagery. *Remote*
1083 *sensing of Environment*, 107, 545-558, doi:10.1016/j.rse.2006.10.006.
- 1084 Anderson, M. C., Allen, R. G., Morse, A., & Kustas, W. P. (2012). Use of Landsat
1085 thermal imagery in monitoring evapotranspiration and managing water
1086 resources. *Remote sensing of Environment*, 122, 50-65,
1087 doi:10.1016/j.rse.2011.08.025.
- 1088 Anderson, M. C., Kustas, W. P., Norman, J. M., Hain, C. R., Mecikalski, J. R., Schultz,
1089 L., González-Dugo, M. P., Cammalleri, C., d'Urso, G., Pimstein, A., & Gao, F.
1090 (2011). Mapping daily evapotranspiration at field to continental scales using
1091 geostationary and polar orbiting satellite imagery. *Hydrology and Earth*
1092 *System Sciences*, 15, 223–239, doi:10.5194/hess-15-223-2011.
- 1093 Barsi, J., Hook, S. J., Schott, J. R., Raqueno, N. G., & Markham, B. L. (2007).
1094 Landsat-5 thematic mapper thermal band calibration update. *IEEE Geoscience*
1095 *and Remote Sensing Letters*, 4(4), 552-555, doi:10.1109/LGRS.2007.896322.
- 1096 Bechtel, B., Zakšek, K., & Hoshyaripour, G. (2012). Downscaling land surface
1097 temperature in an urban area: A case study for Hamburg, Germany. *Remote*
1098 *Sensing*, 4, 3184-3200, doi:10.3390/rs4103184.
- 1099 Bindhu, V. M., Narasimhan, B., & Sudheer, K. P. (2013). Development and

1100 verification of a non-linear disaggregation method (NL-DisTrad) to downscale
1101 MODIS land surface temperature to the spatial scale of Landsat thermal data
1102 to estimate evapotranspiration. *Remote Sensing of Environment*, 135, 118-129,
1103 doi:10.1016/j.rse.2013.03.023.

1104 Bisht, G., Venturini, V., Islam, S., & Jiang, L. (2005). Estimation of the net radiation
1105 using MODIS (Moderate Resolution Imaging Spectroradiometer) data for
1106 clear sky days. *Remote sensing of Environment*, 97, 52–67,
1107 doi:10.1016/j.rse.2005.03.014.

1108 Chander, G., Markham, B. L., & Helder, D. L. (2009). Summary of current
1109 radiometric calibration coefficients for Landsat MSS, TM, ETM+, and EO-1
1110 ALI sensors. *Remote sensing of Environment*, 113, 893-903,
1111 doi:10.1016/j.rse.2005.03.014.

1112 Chen, Y., Zhan, W., Quan, J., Zhou, J., Zhu, X., & Sun, H. (2014). Disaggregation of
1113 remotely sensed land surface temperature: a generalized paradigm. *IEEE*
1114 *Transactions on Geoscience and Remote Sensing*, 52, 5952-5965,
1115 doi:10.1109/TGRS.2013.2294031.

1116 Despini, F., Teggi, S., & Baraldi, A. (2014). Methods and metrics for the assessment
1117 of Pan-sharpening algorithms. In, *SPIE Remote Sensing (pp.*
1118 *924403-924403-924414): International Society for Optics and Photonics,*
1119 doi:10.1117/12.2067316.

1120 Dominguez, A., Kleissl, J., Luvall, J. C., & Rickman, D. L. (2011). High-resolution
1121 urban thermal sharpener (HUTS). *Remote Sensing of Environment*, 115(7),
1122 1772-1780, doi:10.1016/j.rse.2011.03.008.

1123 Gao, L., Zhan, W., Huang, F., Quan, J., Lu, X., Wang, F., Ju, W., & Zhou, J. (2017).
1124 Localization or globalization? Determination of the optimal regression

1125 window for disaggregation of land surface temperature. *IEEE Transactions on*
1126 *Geoscience and Remote Sensing*, 55, 477-490, doi:
1127 10.1109/TGRS.2016.2608987.

1128 Gevaert, C. M., & García-Haro, F. J. (2015). A comparison of STARFM and an
1129 unmixing-based algorithm for Landsat and MODIS data fusion. *Remote*
1130 *Sensing of Environment*, 156, 34-44, doi:10.1016/j.rse.2014.09.012.

1131 Ghosh, A., & Joshi, P. K. (2014). Hyperspectral imagery for disaggregation of land
1132 surface temperature with selected regression algorithms over different land use
1133 land cover scenes. *ISPRS Journal of Photogrammetry and Remote Sensing*, 96,
1134 76-93, doi:10.1016/j.isprsjprs.2014.07.003.

1135 Gillespie, A., Rokugawa, S., Matsunaga, T., Cothern, J. S., Hook, S., & Kahle, A. B.
1136 (1998). A temperature and emissivity separation algorithm for Advanced
1137 Spaceborne Thermal Emission and Reflection Radiometer (ASTER) images.
1138 *IEEE Transactions on Geoscience and Remote Sensing*, 36(4), 1113-1126,
1139 doi:10.1109/36.700995.

1140 Göttsche, F. M., & Olesen, F. S. (2009). Modelling the effect of optical thickness on
1141 diurnal cycles of land surface temperature. *Remote Sensing of Environment*,
1142 113(11), 2306-2316, doi:10.1016/j.rse.2009.06.006.

1143 Han, J., Kamber, M., & Pei, J. (2011). Data mining: concepts and techniques:
1144 concepts and techniques. (3th ed.). Elsevier.

1145 Jiménez-Muñoz, J. C., & Sobrino, J. A. (2003). A generalized single-channel method
1146 for retrieving land surface temperature from remote sensing data. *Journal of*
1147 *Geophysical Research: Atmospheres*, 108, ACL 2-1,
1148 doi:10.1029/2003JD003480.

1149 Jiménez-Muñoz, J. C., Mattar, C., Sobrino, J. A., & Malhi, Y. (2016). Digital thermal

1150 monitoring of the Amazon forest: an intercomparison of satellite and
1151 reanalysis products. *International Journal of Digital Earth*, 9(5), 477-498,
1152 doi:10.1080/17538947.2015.1056559.

1153 Keramitsoglou, I., Kiranoudis, C. T., & Weng, Q. (2013). Downscaling geostationary
1154 land surface temperature imagery for urban analysis. *IEEE Geoscience and*
1155 *Remote Sensing Letters*, 10, 1253-1257, doi:10.1109/LGRS.2013.2257668.

1156 Kustas, W. P., Norman, J. M., Anderson, M. C., & French, A. N. (2003). Estimating
1157 subpixel surface temperatures and energy fluxes from the vegetation
1158 index–radiometric temperature relationship. *Remote sensing of Environment*,
1159 85, 429-440, doi:10.1016/S0034-4257(03)00036-1.

1160 Lagouarde, J.-P., Bach, M., Sobrino, J. A., Boulet, G., Briottet, X., Cherchali, S.,
1161 Coudert, B., Dadou, I., Dedieu, G., & Gamet, P. (2013). The MISTIGRI
1162 Thermal Infrared project: scientific objectives and mission specifications.
1163 *International Journal of Remote Sensing*, 34, 3437-3466,
1164 doi:10.1080/01431161.2012.716921.

1165 Lagouarde, J.-P., Hénon, A., Kurz, B., Moreau, P., Irvine, M., Voogt, J., & Mestayer, P.
1166 (2010). Modelling daytime thermal infrared directional anisotropy over
1167 Toulouse city centre. *Remote sensing of Environment*, 114, 87-105,
1168 doi:10.1016/j.rse.2009.08.012.

1169 Liang, S. (2005). *Quantitative remote sensing of land surfaces*: John Wiley & Sons.

1170 Li, Z.-L., Tang, B.-H., Wu, H., Ren, H., Yan, G., Wan, Z., Trigo, I. F., & Sobrino, J. A.
1171 (2013). Satellite-derived land surface temperature: Current status and
1172 perspectives. *Remote sensing of Environment*, 131, 14-37,
1173 doi:10.1016/j.rse.2012.12.008.

1174 Liu, D., & Zhu, X. (2012). An enhanced physical method for downscaling thermal

1175 infrared radiance. *IEEE Geoscience and Remote Sensing Letters*, 9(4),
1176 690-694, doi:10.1109/LGRS.2011.2178814.

1177 Liu, J., & Moore, J. M. (1998). Pixel block intensity modulation: adding spatial detail
1178 to TM band 6 thermal imagery. *International Journal of Remote Sensing*, 19,
1179 2477-2491, doi:10.1080/014311698214578.

1180 McFeeters, S. (1996). The use of the Normalized Difference Water Index (NDWI) in
1181 the delineation of open water features. *International Journal of Remote*
1182 *Sensing*, 17, 1425-1432, doi:10.1080/01431169608948714.

1183 Mechri, R., Ottlé, C., Pannekoucke, O., & Kallel, A. (2014). Genetic particle filter
1184 application to land surface temperature downscaling. *Journal of Geophysical*
1185 *Research: Atmospheres*, 119, 2131-2146, doi:10.1002/2013JD020354.

1186 Merlin, O., Duchemin, B., Hagolle, O., Jacob, F., Coudert, B., Chehbouni, G., Dedieu,
1187 G., Garatuza, J., & Kerr, Y. (2010). Disaggregation of MODIS surface
1188 temperature over an agricultural area using a time series of Formosat-2 images.
1189 *Remote sensing of Environment*, 114, 2500-2512,
1190 doi:10.1016/j.rse.2010.05.025.

1191 Moosavi, V., Talebi, A., Mokhtari, M. H., Shamsi, S. R. F., & Niazi, Y. (2015). A
1192 wavelet-artificial intelligence fusion approach (WAIFA) for blending Landsat
1193 and MODIS surface temperature. *Remote sensing of Environment*, 169,
1194 243-254, doi:10.1016/j.rse.2015.08.015.

1195 Mukherjee, S., Joshi, P., & Garg, R. (2014). A comparison of different regression
1196 models for downscaling Landsat and MODIS land surface temperature images
1197 over heterogeneous landscape. *Advances in Space Research*, 54, 655-669,
1198 doi:10.1016/j.asr.2014.04.013.

1199 Nichol, J. (2009). An emissivity modulation method for spatial enhancement of

1200 thermal satellite images in urban heat island analysis. *Photogrammetric*
1201 *Engineering & Remote Sensing*, 75, 547-556,
1202 doi:<http://dx.doi.org/10.14358/PERS.75.5.547>.

1203 Nishii, R., Kusanobu, S., & Tanaka, S. (1996). Enhancement of low spatial resolution
1204 image based on high resolution-bands. *IEEE Transactions on Geoscience and*
1205 *Remote Sensing*, 34, 1151-1158, doi:10.1109/36.536531.

1206 Pardo-Igúzquiza, E., Chica-Olmo, M., & Atkinson, P. M. (2006). Downscaling
1207 cokriging for image sharpening. *Remote sensing of Environment*, 102, 86-98,
1208 doi:10.1016/j.rse.2006.02.014.

1209 Pohl, C., & Van Genderen, J. L. (1998). Review article multisensor image fusion in
1210 remote sensing: concepts, methods and applications. *International Journal of*
1211 *Remote Sensing*, 19(5), 823-854, doi:10.1080/014311698215748.

1212 Qin, Z.-H., Karnieli, A., & Berliner, P. (2001). A mono-window algorithm for
1213 retrieving land surface temperature from Landsat TM data and its application
1214 to the Israel-Egypt border region. *International Journal of Remote Sensing*, 22,
1215 3719-3746, doi:10.1080/01431160010006971.

1216 Roberts, D. A., Quattrochi, D. A., Hulley, G. C., Hook, S. J., & Green, R. O. (2012).
1217 Synergies between VSWIR and TIR data for the urban environment: An
1218 evaluation of the potential for the Hyperspectral Infrared Imager (HyspIRI)
1219 Decadal Survey mission. *Remote sensing of Environment*, 117, 83-101,
1220 doi:10.1016/j.rse.2011.07.021.

1221 Rodriguez-Galiano, V., Pardo-Iguzquiza, E., Sanchez-Castillo, M., Chica-Olmo, M.,
1222 & Chica-Rivas, M. (2012). Downscaling Landsat 7 ETM+ thermal imagery
1223 using land surface temperature and NDVI images. *International Journal of*
1224 *Applied Earth Observation and Geoinformation*, 18, 515-527,

1225 doi:10.1016/j.jag.2011.10.002.

1226 Sandholt, I., Rasmussen, K., & Andersen, J. (2002). A simple interpretation of the
1227 surface temperature/vegetation index space for assessment of surface moisture
1228 status. *Remote sensing of Environment*, 79, 213-224,
1229 doi:10.1016/S0034-4257(01)00274-7.

1230 Sobrino, J. A., Jiménez-Muñoz, J. C., El-Kharraz, J., Gómez, M., Romaguera, M., &
1231 Soria, G. (2004). Single-channel and two-channel methods for land surface
1232 temperature retrieval from DAIS data and its application to the Barrax site.
1233 *International Journal of Remote Sensing*, 25(1), 215-230,
1234 doi:10.1080/0143116031000115210.

1235 Sobrino, J. A., Gómez, M., Jiménez-Muñoz, J. C., & Oliso, A. (2007). Application of
1236 a simple algorithm to estimate daily evapotranspiration from NOAA–AVHRR
1237 images for the Iberian Peninsula. *Remote sensing of Environment*, 110(2),
1238 139-148, doi:10.1016/j.rse.2007.02.017.

1239 Sobrino, J. A., Ultra-Carrió, R., Sòria, G., Bianchi, R., & Paganini, M. (2012). Impact
1240 of spatial resolution and satellite overpass time on evaluation of the surface
1241 urban heat island effects. *Remote Sensing of Environment*, 117, 50-56,
1242 doi:10.1016/j.rse.2011.04.042.

1243 Stathopoulou, M., & Cartalis, C. (2009). Downscaling AVHRR land surface
1244 temperatures for improved surface urban heat island intensity estimation.
1245 *Remote sensing of Environment*, 113, 2592-2605,
1246 doi:10.1016/j.rse.2009.07.017.

1247 Taylor, K. E. (2001). Summarizing multiple aspects of model performance in a single
1248 diagram. *Journal of Geophysical Research: Atmospheres*, 106, 7183-7192.

1249 Teggi, S. (2012). A technique for spatial sharpening of thermal imagery of coastal

1250 waters and of watercourses. *International journal of remote sensing*, 33(10),
1251 3063-3089, doi:10.1080/01431161.2011.627888.

1252 Teggi, S., & Despini, F. (2014). Estimation of subpixel MODIS water temperature
1253 near coastlines using the SWTI algorithm. *Remote sensing of Environment*,
1254 142, 122-130, doi:10.1016/j.rse.2013.11.011.

1255 Vivone, G., Alparone, L., Chanussot, J., Dalla Mura, M., Garzelli, A., Licciardi, G. A.,
1256 Restaino, R., & Wald, L. (2014). A critical comparison among pansharpening
1257 algorithms. *IEEE Transactions on Geoscience and Remote Sensing*, 53, 2565 -
1258 2586, doi:10.1109/TGRS.2014.2361734.

1259 Wald, L., Ranchin, T., & Mangolini, M. (1997). Fusion of satellite images of different
1260 spatial resolutions: assessing the quality of resulting images. *Photogrammetric
1261 Engineering and Remote Sensing*, 63, 691-699,
1262 doi:https://hal.archives-ouvertes.fr/hal-00365304.

1263 Wang, Z., & Bovik, A. C. (2002). A universal image quality index. *IEEE Signal
1264 Processing Letters*, 9, 81-84, doi:10.1109/97.995823.

1265 Wang, F., Qin, Z., Li, W., Song, C., Karnieli, A., & Zhao, S. (2014). An efficient
1266 approach for pixel decomposition to increase the spatial resolution of land
1267 surface temperature images from MODIS thermal infrared band data. *Sensors*,
1268 15, 304-330, doi:10.3390/s150100304.

1269 Weng, Q., Fu, P., & Gao, F. (2014). Generating daily land surface temperature at
1270 Landsat resolution by fusing Landsat and MODIS data. *Remote Sensing of
1271 Environment*, 145, 55-67, doi:10.1016/j.rse.2014.02.003.

1272 Wu, P., Shen, H., Zhang, L., & Göttsche, F.-M. (2015). Integrated fusion of
1273 multi-scale polar-orbiting and geostationary satellite observations for the
1274 mapping of high spatial and temporal resolution land surface temperature.

1275 *Remote sensing of Environment*, 156, 169-181, doi:10.1016/j.rse.2014.09.013.

1276 Zakšek, K., & Oštir, K. (2012). Downscaling land surface temperature for urban heat
1277 island diurnal cycle analysis. *Remote sensing of Environment*, 117, 114-124,
1278 doi:10.1016/j.rse.2011.05.027.

1279 Zhan, W., Chen, Y., Zhou, J., Wang, J., Liu, W., Voogt, J., Zhu, X., Quan, J., & Li, J.
1280 (2013). Disaggregation of remotely sensed land surface temperature:
1281 Literature survey, taxonomy, issues, and caveats. *Remote sensing of*
1282 *Environment*, 131, 119-139, doi:10.1016/j.rse.2012.12.014.

1283 Zhan, W., Huang, F., Quan, J., Zhu, X., Gao, L., Zhou, J., & Ju, W. (2016).
1284 Disaggregation of remotely sensed land surface temperature: A new dynamic
1285 methodology. *Journal of Geophysical Research: Atmospheres*, 121(18),
1286 10538–10554, doi: 10.1002/2016JD024891.

1287 Zhan, W., Chen, Y., Zhou, J., Li, J., & Liu, W. (2011). Sharpening thermal imageries:
1288 A generalized theoretical framework from an assimilation perspective. *IEEE*
1289 *Transactions on Geoscience and Remote Sensing*, 49(2), 773-789,
1290 doi:10.1109/TGRS.2010.2060342.

1291 Zhou, J., Chen, Y., Zhang, X., & Zhan, W. (2013). Modelling the diurnal variations of
1292 urban heat islands with multi-source satellite data. *International Journal of*
1293 *Remote Sensing*, 34, 7568-7588, doi:10.1080/01431161.2013.821576.

1294 Zhou, J., Li, M., Liu, S., Jia, Z., & Ma, Y. (2015). Validation and performance
1295 evaluations of methods for estimating land surface temperatures from ASTER
1296 Data in the middle reach of the Heihe River Basin, Northwest China. *Remote*
1297 *Sensing*, 7(6), 7126-7156, doi:10.3390/rs70607126.

1298 Zhou, J., Liu, S., Li, M., Zhan, W., Xu, Z., & Xu, T. (2016). Quantification of the
1299 Scale Effect in Downscaling Remotely Sensed Land Surface Temperature.

1300 *Remote Sensing*, 8, 975, doi:10.3390/rs8120975.

1301

1302

1303

1304

1305