1    An Improved Flexible Spatiotemporal DAta Fusion (IFSDAF) method for

2    producing high spatiotemporal resolution NDVI time series

3

4    Meng Liu[a,e], Wei Yang[b], Xiaolin Zhu[c], Jin Chen[a*], Xuehong Chen[a], Linqing Yang[d,e]

5

6    [a] State Key Laboratory of Earth Surface Processes and Resource Ecology, Institute of

7    Remote Sensing Science and Engineering, Faculty of Geographical Science, Beijing

8    Normal University, Beijing 100875, China

9    [b] Center for Environmental Remote Sensing, Chiba University, Chiba 263-8522, Japan

10   [c] Department of Land Surveying and Geo-Informatics, The Hong Kong Polytechnic

11   University

12   [d] State Key Laboratory of Remote Sensing Science, Institute of Remote Sensing

13   Science and Engineering, Faculty of Geographical Science, Beijing Normal

14   University, Beijing 100875, China

15   [e] Department of Ecosystem Science and Management, Texas A&M University,

16   College Station, TX 77843, USA

17

18

19

20    *Corresponding author: Prof. Jin Chen. E-mail address: chenjin@bnu.edu.cn

21  **ABSTRACT**

22      The Normalized Difference Vegetation Index (NDVI) is one of the mostly used

23  vegetation index for ecosystem dynamics monitoring and biosphere process modeling.

24  However, global NDVI products are usually provided with relatively coarse spatial

25  resolutions, being short of important spatial details. Producing NDVI time-series data

26  with a high spatiotemporal resolution is thus indispensable for monitoring land

27  surface and ecosystem changes, especially in heterogeneous areas. An Improved

28  Flexible Spatiotemporal DAta Fusion (IFSDAF) method is developed in this study to

29  overcome the existing issues. In accordance with the distinctive characteristics of

30  NDVI with large data variance and high spatial autocorrelation compared with raw

31  reflectance bands, the IFSDAF method first produces temporal increment with linear

32  unmixing and spatial-dependent increment by thin plate spline (TPS) interpolation,

33  and then obtains final prediction from the optimal integration of two increments by

34  Constrained Least Square (CLS) theory. Moreover, IFSDAF is developed with a

35  capacity of employing all available and partially contaminated fine images. Coarse

36  spatial resolution NDVI (MODIS) and fine spatial resolution NDVI images (Landsat

37  and Sentinel) in areas with great spatial heterogeneity and significant land cover

38  changes were used to test the performance of the new method. The promising results

39  (RMSE 0.0884, rRMSE 22.12% in heterogeneous areas, RMSE 0.0546, rRMSE 25.77%

40  in land cover change areas) demonstrate the strengths and robustness of the proposed

41  method in providing reliable high spatial and temporal resolution NDVI datasets to

42  support research on land surface processes. The proposed IFSDAF method can be

43  further simplified by only using spatial-dependent increment to improve the efficiency

44  to a great extent. It will make IFSDAF a feasible method for applications in large

45  geographical area and has the potential for global studies.

46

47  **Keywords**: Normalized Difference Vegetation Index (NDVI), Spatiotemporal Data

48  Fusion, High Spatial and Temporal Resolution, Constrained Least Square (CLS)

49  method, Weighted Integration

## 1. Introduction

The Normalized Difference Vegetation Index (NDVI) enhances the absorptive and reflective features of vegetation and provides a proxy for measuring canopy greenness and vigor (Rouse et al., 1974; Huete et al., 2002). Accordingly, NDVI time-series data derived from spaceborne sensors are widely employed in ecosystem dynamics monitoring and biosphere process modeling, helping to understand responses of ecosystems to climate change (Pettorelli et al., 2005). As the most significant constraint of the available NDVI time-series products (e.g., GIMMS, MODIS, SPOT VGT), coarse spatial resolutions ranging from 250 m to 8 km prevent these products from capturing spatial details necessary for monitoring land surface and ecosystem changes, especially in geographically heterogeneous areas (Gao et al., 2006; Rao et al., 2015). Producing NDVI time-series data with both high spatial and high temporal resolutions is thus critically required for such applications, raising the need for developing spatiotemporal fusion methods by blending the high frequent but low spatial resolution images (e.g., MODIS images, hereinafter referred to as coarse images) with the high spatial resolution but low frequent images (e.g., Landsat images, hereinafter referred to as fine images) (Zhu et al., 2018). Recently with emerging constellations of CubeSats and new satellite systems (e.g. Sentinel 2 with 5 day NDVI at 10 m resolution observations), new opportunities to alleviate the issue of the classical trade-off between spatial and temporal resolution is becoming hoped, however, spatiotemporal fusion is still necessary for long time series analysis as such

71    data are unavailable before 2015.

72        When using spatiotemporal fusion technology to produce NDVI data with high

73    spatial and temporal resolutions, users need to solve the two puzzles: (I) selecting an

74    appropriate blending strategy: Blend-then-Index (BI) or Index-then-Blend (IB), and

75    (II) selecting a suitable and accurate spatiotemporal fusion method. For the first

76    puzzle, recent studies (Chen et al., 2018; Jarihani et al., 2014; Tian et al., 2013) have

77    demonstrated that the IB strategy consistently yields better or comparable results than

78    the BI, mainly because the IB method has these advantages compared with the BI: (i)

79    less error propagation in the blending process; (ii) less computationally expensive;

80    and (iii) easier to clean the noises (e.g., cloud effects) on NDVI than the raw

81    reflectance bands by the advanced filters (e.g., Chen et al., 2004). Consequently, IB is

82    generally recommended and becomes the dominant blending strategy for producing

83    fused NDVI products.

84        Regarding the second puzzle, a number of spatiotemporal fusion methods have

85    been proposed and validated over past years (Zhu et al., 2018). These methods need at

86    least one pair of cloud-free fine and coarse NDVI images at a base date and a series of

87    coarse NDVI images at the prediction dates as the input. However, the consensus

88    regarding the most suitable method for producing high spatiotemporal resolution

89    NDVI data has not been reached. Generally, as a band combination index for feature

90    enhancement, NDVI enlarges the contrast between vegetated and non-vegetated

91    pixels and therefore displays larger spatial and temporal variance (i.e., larger

92   heterogeneity) than the raw reflectance in most satellite images. Accordingly, a

93   suitable spatiotemporal fusion method for fusing NDVI product is supposed to satisfy

94   the following criteria in practice: (i) obtaining good prediction in areas with large

95   spatial and temporal variance; (ii) requiring only one pair of clear fine and coarse

96   NDVI image at a base date, ensuring its applicability in areas with frequent cloud

97   contamination; (iii) having a capacity to handle land cover change, such as

98   urbanization, deforestation/reforestation, wildfires, floods and land cover transitions

99   caused by other forces. Among the existing spatiotemporal fusion methods, the

100  Flexible Spatiotemporal DAta Fusion method (FSDAF) (Zhu et al., 2016) is the one

101  meeting these criteria and can be considered a potential candidate, while other

102  existing methods fail in at least one criterion, especially the third criterion. For

103  example, the spatial and temporal adaptive reflectance fusion model (STARFM, Gao

104  et al., 2006), the enhanced STARFM (ESTARFM, Zhu et al., 2010), the spatial and

105  temporal adaptive vegetation index fusion model (STAVFM, Meng et al., 2013),

106  unmixing-based spatiotemporal reflectance fusion model (U-STFM, Huang and

107  Zhang, 2014), NDVI linear mixing growth model (NDVI-LMGM, Rao et al., 2015),

108  and spatial and temporal reflectance unmixing model (STRUM, Gevaert and

109  Garcia-Haro, 2015) cannot handle land cover changes occurring between base date

110  and prediction date. The learning-based methods, such as Sparse-representation-based

111  spatiotemporal reflectance fusion model (SPSTFM, Huang and Song, 2012; Song and

112  Huang, 2013), an error-bound-regularized semi-coupled dictionary learning model

113    (EBSCDM, Wu et al., 2015) and an extreme learning machine based fusion method

114    (Liu et al., 2016) can better capture land cover change but their learning step is time

115    consuming, and the accuracy decreases when the spatial heterogeneity is high and

116    scale differences between coarse and fine images are large (Zhu et al., 2016).

117        FSDAF is based on the spectral unmixing analysis and further introduces thin

118    plate spline (TPS) interpolation to capture land cover change if the change is

119    detectable in coarse images (Zhu et al., 2016). Compared with two widely used

120    spatiotemporal fusion methods, STAFRM algorithm (Gao et al., 2006) and an

121    unmixing-based data fusion (UBDF) algorithm (Zurita-Milla et al., 2008), FSDAF

122    needs the same input data as these two methods but is superior in producing more

123    accurate predictions especially in the NIR band of heterogeneous landscapes (Table 3

124    and Table 4 in Zhu et al., 2016). Like NDVI, the NIR band has larger spatial and

125    temporal variances than red band, because the reflectance in NIR band generally has

126    larger difference among different land covers than red band, and it has more significant

127    temporal changes than red band during vegetation growth cycles. Moreover, FSDAF

128    can capture both the gradual and abrupt land cover changes, which is an existing issue

129    with current spatiotemporal fusion methods. Considering many advantages of FSDAF,

130    it could be the appropriate method for producing high spatiotemporal resolution

131    NDVI data. However, there is still space to further improve the FSDAF method. It

132    should be noted that the FSDAF method only relies on the result of TPS interpolation

133    to distribute residuals ($\varepsilon$) between prediction and true values under an assumption that

134 errors mainly depend on the landscape homogeneity. Such an assumption is very

135 empirical and has not been demonstrated by theoretical analysis. It may be not an

136 optimal way to distribute residuals for different scenarios. Furthermore, in practice,

137 many available fine images are partially contaminated by clouds. Clear pixels on

138 these partially contaminated fine images can provide significant information of

139 temporal changes, demonstrated by STAIR method proposed by Luo et al. (2018)

140 with better result of producing daily surface reflectance than STARFM. Consequently,

141 using cloud-free fine images together with partially contaminated fine images will

142 benefit spatiotemporal NDVI fusion and expand its applicability to clouded regions.

143 Unfortunately, the FSDAF method falls short in such a capacity and is not applicable

144 in clouded regions.

145 　　To address the abovementioned limitations, we propose an Improved Flexible

146 Spatiotemporal DAta Fusion (IFSDAF) method for producing high spatiotemporal

147 resolution NDVI time series. The IFSDAF incorporates Constrained Least Square

148 (CLS) theory into FSDAF method, by which temporal prediction derived from

149 unmixing procedure and spatial prediction derived from TPS interpolation are

150 combined, thus ensuring final prediction obtained from the optimal integration of

151 temporal and spatial predictions. Moreover, IFSDAF was developed with the capacity

152 of employing all available and partially contaminated fine images (e.g. maximum

153 cloud coverage is less than 70%). To validate the effectiveness of the proposed

154 method, comparison with three popular NDVI fusion methods (i.e., NDVI-LMGM,

155　STARFM and FSDAF) under the IB strategy were performed in several experiment

156　areas, including a site with a heterogeneous landscape, a site with abrupt land cover

157　changes, and a site where satellite images contain a lot of clouds.
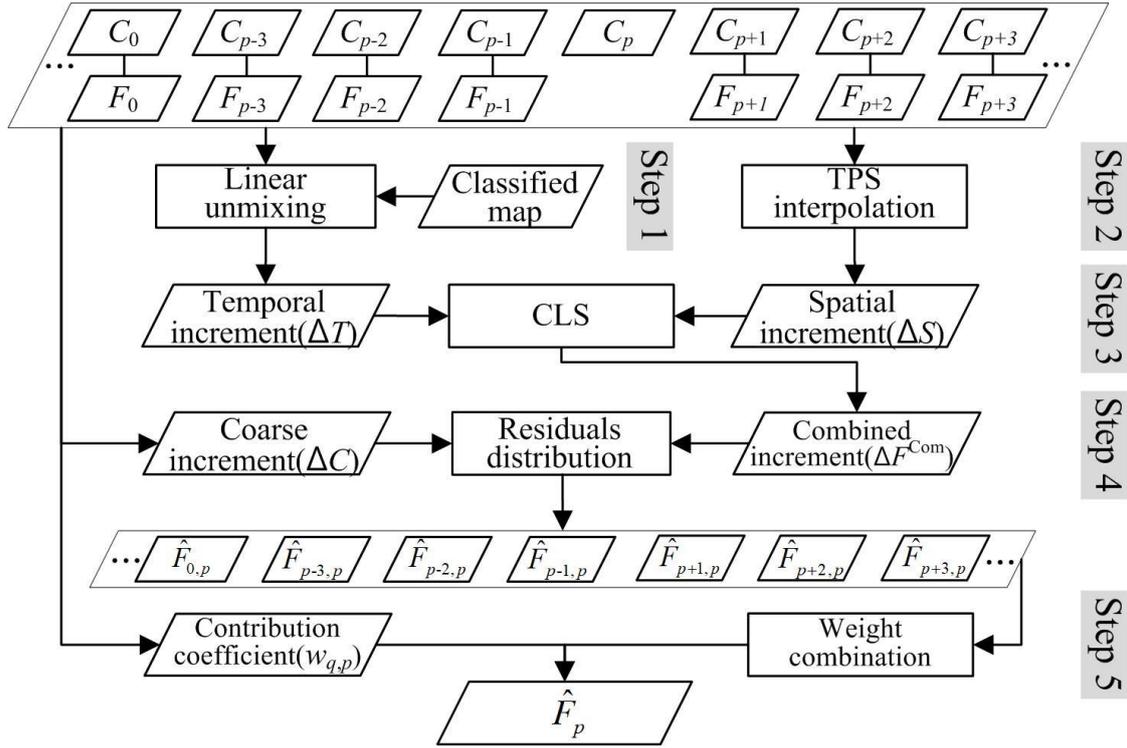
158　**2. Methodology**

159　　Although the principles of existing spatiotemporal fusion methods have a great

160　variety, the main idea can be framed by Eq. (1), in which the fine increment of NDVI

161　($\Delta F$) between the predicting date ($t_p$) and the base date ($t_0$) is firstly estimated, and then

162　fine NDVI values ($F_p$) on the predicting date ($t_p$) are predicted as the sum of the base

163　fine NDVI value ($F_0$) and the increment ($\Delta F$), plus the residuals $\varepsilon$.

$$F_p = F_0 + \Delta F + \varepsilon \qquad (1)$$

165　Given that $F_0$ is known, IFSDAF also follows this unified equation but estimates the

166　increment in two ways: the temporal increments using (i) unmixing analysis and

167　spatial-dependent increments using (ii) the Thin Plate Spline (TPS) interpolation

168　method, and then combine the two increments to obtain final $\Delta F$ through a Constrained

169　Least Square (CLS) method. The CLS method adopted here purifies the original

170　FSDAF because it can adaptively combine the two increments, allowing the final $\Delta F$

171　approaching the one with higher accuracy.

172　　The flowchart of the proposed IFSDAF is shown in Fig. 1. The input data for

173　IFSDAF include coarse NDVI time series images and all available fine NDVI images

174　within the same period. In these images, coarse NDVI and fine NDVI images acquired

175　at the same dates are named as one pair. The pair with minimal cloud contaminations is

176  selected as the base images ($C_0$ and $F_0$) and its acquisition date is the base date $t_0$. The

177  dates of other pairs are denoted as $\cdots$, $p$-3, $p$-2, $p$-1, $p$+1, $p$+2, $p$+3, $\cdots$. The coarse and

178  fine NDVI images of these pairs are denoted as ($\cdots$, $C_{p-3}$, $C_{p-2}$, $C_{p-1}$, $C_{p+1}$, $C_{p+2}$, $C_{p+3}$, $\cdots$)

179  and ($\cdots$, $F_{p-3}$, $F_{p-2}$, $F_{p-1}$, $F_{p+1}$, $F_{p+2}$, $F_{p+3}$, $\cdots$) respectively. The task of IFSDAF is to

180  predict fine NDVI images at any dates whenever a course NDVI image is available, e.g.,

181  the date of $t_p$. In IFSDAF, the input fine NDVI images are not required to be cloud free

182  except $F_0$. Like other spatiotemporal fusion methods, all the coarse and fine NDVI

183  images need to be geo-registered and cropped to become the same image size. Besides,

184  coarse NDVI time-series images need to be smoothed by an algorithm based on

185  Savitzky-Golay filter (Chen et al., 2004), which was designed to reconstruct a

186  high-quality NDVI time-series data by keeping the clear-sky values and interpolating

187  clouded values. And cloud pixels in partially cloud-contaminated fine NDVI images

188  are also masked by the Fmask algorithm (Zhu and Woodcock, 2012). A land cover

189  classification map at a fine resolution, which can be derived from either existing land

190  cover products (e.g. Globeland30, Chen et al., 2015) or the classification result of the

191  input clear fine images, is needed to provide fractional cover for the unmixing process.

192  The output of IFSDAF is synthetic fine NDVI images ($\hat{F}_p$) on the prediction date $t_p$

193  ($p$=1, 2, 3, …). More detailed description for each implementation step of IFSDAF is

194  given below and a list of notations and explanation is given in Appendix.

**Fig.1.** Flowchart of the Improved Flexible Spatiotemporal DAta Fusion method (IFSDAF)

## 2.1 Generation of temporal increment by unmixing method

Following the linear spectral mixing theory, the temporal NDVI change (increment) of a coarse pixel can be considered as the linear combination of NDVI increments of fine pixels within the coarse pixel during a short period (Rao et al., 2015). Accordingly, a linear mixture model is used to unmix the increment of coarse pixels from the base date $t_0$ to the prediction date $t_p$, assuming that fine pixels belonging to the same land cover class have a similar increment within the local region (Busetto et al., 2008; Rao et al., 2015). Neighboring coarse pixels within a moving window centered by coarse pixel $(x, y)$ are used to establish a linear equation system, as shown in Eq. (2).

11

208

$$
\begin{bmatrix} \Delta C(1,1) \\ \mathrm{M} \\ \Delta C(x,y) \\ \mathrm{M} \\ \Delta C(n,n) \end{bmatrix} = \begin{bmatrix} f_1(1,1) & f_2(1,1) & \mathrm{L} & f_l(1,1) \\ \mathrm{M} & \mathrm{M} & & \mathrm{M} \\ f_1(x,y) & f_2(x,y) & \mathrm{L} & f_l(x,y) \\ \mathrm{M} & \mathrm{M} & & \mathrm{M} \\ f_1(n,n) & f_2(n,n) & \mathrm{L} & f_l(n,n) \end{bmatrix} \begin{bmatrix} \Delta F_1 \\ \mathrm{M} \\ \Delta F_c \\ \mathrm{M} \\ \Delta F_l \end{bmatrix},
\tag{2}
$$

209

210 $\quad$ *with s.t.* $\min(\Delta C_{\mathrm{window}}) - \mathrm{std}(\Delta C_{\mathrm{window}}) \leq \Delta F_c \leq \max(\Delta C_{\mathrm{window}}) + \mathrm{std}(\Delta C_{\mathrm{window}})$

211 where $n$ is the number of coarse pixels and $l$ is the number of land cover classes within

212 the moving window. $\Delta C(x, y)$ is the NDVI increment of the coarse pixel $(x, y)$ that can

213 be obtained directly from coarse NDVI time series images. $\Delta F_c$ is the fine NDVI

214 increment of class $c$ within the window. $f_l(x, y)$ is the fraction of class $l$ within the coarse

215 pixel $(x, y)$, which can be obtained from the land cover map at a fine resolution.

216 $\Delta C_{\mathrm{window}}$ is the set of all coarse NDVI increments in the window. $\min(\Delta C_{\mathrm{window}})$,

217 $\max(\Delta C_{\mathrm{window}})$, and $\mathrm{std}(\Delta C_{\mathrm{window}})$ are the minimum value, maximum value and

218 standard deviation of $\Delta C_{\mathrm{window}}$, respectively. A moving window sized at a 7×7 coarse

219 pixel is recommended because the number of coarse pixels in the window, 49, is

220 commonly much larger than the number of land cover classes. This choice of window

221 size ensures the abovementioned overdetermined linear equations less influenced by

222 collinearity and land cover changes. By solving the linear equations, the temporal

223 NDVI increment of each class ($\Delta F_c$) in the moving window can be acquired. Then, the

224 fine temporal increment $\Delta T(x_j, y_j)$, where $(x_j, y_j)$ devotes $j$th fine pixel in the coarse

225 pixel $(x, y)$, is defined by Eq. (3), as following,

226 $\quad\quad\quad\quad \Delta T(x_j, y_j) = \Delta F_c$ if fine pixel $(x_j, y_j)$ belongs to class $c$. $\quad\quad\quad$ (3)

227        The fine-resolution land cover map used to compute the class fractions can be

228    an available land cover product or classification of a cloud-free fine image. In practice,

229    to make the fusion process automatic, existing fusion methods often use unsupervised

230    classifiers (e.g. K-means and ISODATA) to obtain spectral classes rather than real

231    land cover classes (Rao et.al 2015, Zhu et.al, 2017). Users need to set the number of

232    classes in unsupervised classification. According to previous studies, the number of

233    classes ranging from 3 to 6 could get satisfied results for most situations (Rao et.al

234    2015, Zhu et.al, 2017). Accuracy assessment of the classification map is not included

235    in the fusion process because: (1) aggregation of fine-scale class to coarse-scale

236    fraction will average out some errors in classification so it may not cause large

237    problem in solving Eq. (2); (2) temporal change assigned to a pixel with wrong class

238    labels using Eq. (3) will be compensated by the spatial-dependent increment

239    introduced in the next section; and (3) reference samples selection for accuracy

240    assessment introduces more human-computer interaction. Although the proposed

241    method is not sensitive to classification accuracy, including more accurate and robust

242    classification methods in IFSDAF could further improve its performance.

**2.2 Generation of spatial-dependent increment by TPS interpolation**

Coarse NDVI image on $t_p$ contains signals of land cover changes when changes are significant enough to be shown in coarse pixels. Therefore, spatial interpolation of coarse NDVI to fine resolution will retain useful information of land cover changes. Accordingly, coarse spatial resolution NDVI images on $t_p$ and $t_0$ are interpolated to fine spatial resolution respectively, through Thin Plate Spline (TPS) interpolation method (Chen et al., 2014; Zhu et al., 2016). TPS as a spatial interpolation technique for point data based on spatial dependence (Dubrule, 1984), is employed to obtain interpolation result thanks to its high accuracy. Then, another increment from the difference between interpolation results on $t_p$ and $t_0$ can be acquired. As this increment only uses spatial dependence among coarse pixels, it can be referred to as the spatial-dependent increment $\Delta S(x_j, y_j)$, as shown in Eq. (4), where $F_p^{\mathrm{TPS}}(x_j, y_j)$ and $F_0^{\mathrm{TPS}}(x_j, y_j)$ are TPS interpolated values on $t_p$ and $t_0$ respectively, and $(x_j, y_j)$ is the $j$th fine pixel within the coarse pixel $(x, y)$.

$$\Delta S(x_j, y_j) = F_p^{\mathrm{TPS}}(x_j, y_j) - F_0^{\mathrm{TPS}}(x_j, y_j) \qquad (4)$$

Compared with the temporal increment, spatial-dependent increment has two advantages. First, coarse NDVI image on date $t_p$ contains signals of land cover changes if the changes are significant enough to be recorded. By TPS interpolation, such land cover change information can be directly captured at a fine resolution. Second, spatial-dependent increment is independent of classification map and unmixing procedure, thus it has the potential to justify errors in the temporal increment resulted

264   from classification or unmixing. In this study, TPS is used to estimate spatial-dependent

265   increment rather than estimate the NDVI value on $t_p$, a strategy used in FADAF,

266   because the increment reveals the changes of NDVI directly. Zhang et al. (2015) also

267   suggested that using increment yields higher accuracy than predicting the value directly

268   at $t_p$. The use of this spatial-dependent increment will be further discussed in **Section 5**.

269   **2.3 Combination of two increments by CLS**

270   The abovementioned two increments can be considered to be two independent

271   predictions by two different models. Due to the distinct features used by the two

272   predictions, the former uses the information of temporal changes of NDVI, and the later

273   mainly utilizes the spatial dependence. Their prediction accuracies should be different

274   under different scenarios and spatial-dependencies. Therefore, it is natural to expect

275   that a reasonable combination of the two increments can improve the performance and

276   robustness of the fusion method.

277   The simplest and most effective way of combining temporal increment ($\Delta T$) and

278   spatial-dependent increment ($\Delta S$) should be summing them by reasonable weights.

279   Moreover, an ideal combination should be as close to the true fine NDVI increment ($\Delta F$)

280   as possible. Thus, an objective function of weighted combination can be written as,

$$(\hat{w}_S, \hat{w}_T) = \arg \min_{(w_S, w_T) \in (0,1)} \sum_k \left( w_S \Delta S_k + w_T \Delta T_k - \Delta F_k \right)^2, \qquad (5)$$

282   where $\Delta S_k$, $\Delta T_k$, and $\Delta F_k$ are the spatial-dependent increment, the temporal increment,

283   and the true increment of the $k$th fine pixel, respectively. $w_S$ and $w_T$ are weights of the

284   spatial-dependent increment and the temporal increment, respectively. Eq. (5) can be

285  solved by the Constrained Least Square (CLS) method, with constraints of $w_S$ and $w_T$

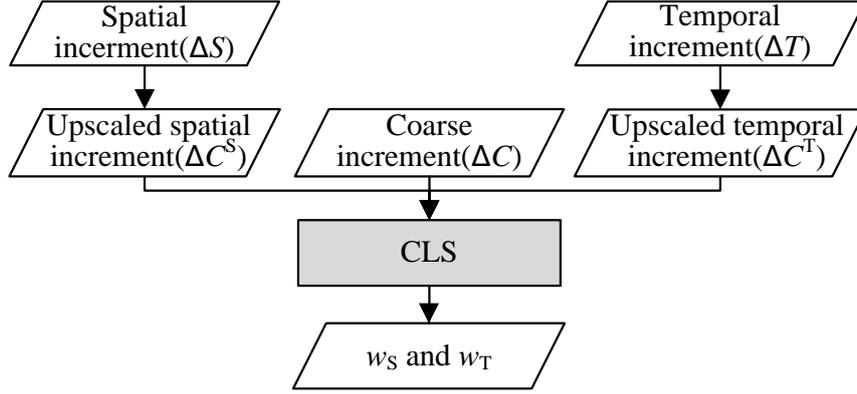286  being nonnegative and summing up to one.

287  However, as fine NDVI values on $t_p$ are unknown, it is impossible to obtain the

288  true fine increment ($\Delta F$). Fortunately, a real NDVI increment of a coarse pixel ($\Delta C$)

289  from $t_0$ to $t_p$ is available because coarse observations are available on two dates.

290  Therefore, both the temporal increment and the spatial-dependent increment are

291  up-scaled to the resolution of coarse pixel ($\Delta C^T$ and $\Delta C^S$), as shown in Fig. 2. Then,

292  $w_S$ and $w_T$ in Eq. (5) can be obtained by solving Eq. (6) alternatively:

293
$$(\hat{w}_S, \hat{w}_T) = \arg \min_{(w_S, w_T) \in (0,1)} \sum_k \left( w_S \Delta C_k^S + w_T \Delta C_k^T - \Delta C_k \right)^2 \tag{6}$$

294  where $\Delta C_k^S$, $\Delta C_k^T$ and $\Delta C_k$ are up-scaled spatial-dependent increment, up-scaled

295  temporal increment and true increment of $k$th coarse pixel, respectively. Here, the

296  average value of all fine NDVI pixels within the coarse pixel are used to produce

297  up-scaled spatial increment ($\Delta C^S$) and up-scaled temporal increment ($\Delta C^T$), and $\Delta C_k$ is

298  calculated as difference between coarse NDVI values on prediction date $t_p$ and $t_0$.

299  Considering that the weights $w_S$ and $w_T$ are spatially-dependent, Eq. (6) is solved in a

300  7×7 moving window at a coarse resolution corresponding to the window size of the

301  unmixing process. Then, with the estimated $w_S$ and $w_T$, the final fine increment can be

302  calculated as following:

303
$$\Delta F^{Com}(x_j, y_j) = w_S \times \Delta S(x_j, y_j) + w_T \times \Delta T(x_j, y_j), \tag{7}$$

304  where $\Delta F^{Com}(x_j, y_j)$ is the combined increment of fine pixel $(x_j, y_j)$. $w_S$ and $w_T$ are

305  supposed to be scale-invariant and its rationality will be discussed **Section 5**.

306

307    **Fig. 2**. Illustration of weighted calibration based on Constrained Least Square (CLS)

308    method.

309    **2.4 Distribution of residuals**

310         After CLS optimization, the combined increment can capture most of the fine

311    NDVI increment. However, residuals are inevitable even though they are minimal. The

312    residuals can be mathematically expressed as Eq. (8),

313    $$R(x, y) = \Delta C(x, y) - \frac{1}{m} \sum_{j=1}^{m} \Delta F^{\text{Com}}(x_j, y_j),$$    (8)

314    where $R(x, y)$ is the residual within a coarse pixel $(x, y)$ and $m$ is the number of fine

315    pixels within the coarse pixel. In order to further improve the accuracy of the combined

316    increment, residual derived above needs to be allocated to each fine pixel $(x_j, y_j)$ within

317    the coarse pixel $(x, y)$. Because the residuals are minimal after the weighted

318    combination of two increments, they can be distributed equally (Chen et al., 2014) as

319    Eq. (9).

320    $$\hat{F}_{0,p}(x_j, y_j) = F_0(x_j, y_j) + \Delta F^{\text{Com}}(x_j, y_j) + R(x, y),$$    (9)

321    where $F_0(x_j, y_j)$ is fine NDVI of pixel $(x_j, y_j)$ on date $t_0$ and $\hat{F}_{0,p}(x_j, y_j)$ is the predicted

322    fine NDVI on date $t_p$. After the residuals distribution, a smoothing process based on

323  similar pixels (Zhu et al., 2016) is applied to remove block effects in the fused image.

324  **2.5 Combination of multi-time predictions**

325  Through **Sections 2.1** to **2.4**, a prediction $\hat{F}_{0,p}$ for date $t_p$ based on the fine

326  NDVI on $t_0$ can be acquired. In the same way, there will be several NDVI predictions,

327  such as …, $\hat{F}_{p-3,p}$, $\hat{F}_{p-2,p}$, $\hat{F}_{p-1,p}$, $\hat{F}_{p+1,p}$, $\hat{F}_{p+2,p}$, $\hat{F}_{p+3,p}$, …, for date $t_p$ based on

328  clear observations at $p+i$ (i=…, -3, -2, -1, 1, 2, 3, …) in other partially clouded fine

329  NDVI images. Recognition of a pixel is either clear or clouded can be performed

330  based on the Fmask algorithm (Zhu and Woodcock, 2012). Generally, the predictions

331  with a base date too far from $t_p$ are excluded considering that the base NDVI images

332  hold weak relationship with the NDVI image on date $t_p$. Operationally, the maximum

333  interval between the base date and the prediction date is set as two months. Then, the

334  NDVI difference of coarse pixels between the base date and the prediction date is

335  used to calculate the contribution of each prediction, as shown in Eq. (10).

336
$$w_{q,p}(x,y) = \frac{1}{\sum_{i=1}^{9}\left|C_q^i(x,y) - C_p^i(x,y)\right|} \tag{10}$$

337  where $C_q^i(x,y)$ and $C_p^i(x,y)$ are coarse NDVI values of the $i$th pixel on base date

338  $q$ and the prediction date $t_p$ in the 3×3 moving window centered by coarse pixel $(x, y)$.

339  $w_{q,p}(x, y)$ is the contribution coefficient of predicted fine NDVI value $\hat{F}_{q,p}(x_j, y_j)$

340  within the center coarse pixel $(x, y)$. Based on the contribution coefficient, the

341  combined prediction of a fine pixel $(x_j, y_j)$ on date $t_p$ is,

342
$$\hat{F}_p(x_j, y_j) = \sum_q [w_{q,p}(x,y) \times \hat{F}_{q,p}(x_j, y_j)] \Big/ \sum_q w_{q,p}(x,y), \tag{11}$$

343  If $C_q^i(x,y)$ equals $C_p^i(x,y)$, $\hat{F}_p(x_j, y_j)$ will be set as $\hat{F}_{q,p}(x_j, y_j)$ since $w_{q,p}(x,y)$ is

344 infinite under this situation. Finally, for each prediction date in the time series, a final

345 prediction in Eq. (11) can be obtained using the same routine described in **Sections**

346 **2.1 through 2.5**.

347 　　To assess the performance of the new method, four accuracy indices, Root Mean

348 Square Error (RMSE), relative RMSE (RMSE divided by averaged observation value),

349 Correlation Coefficient ($r$) and Average Difference (AD) were used. These indices have

350 been widely used to assess the accuracy of fused images in previous studies (e.g. Gao et

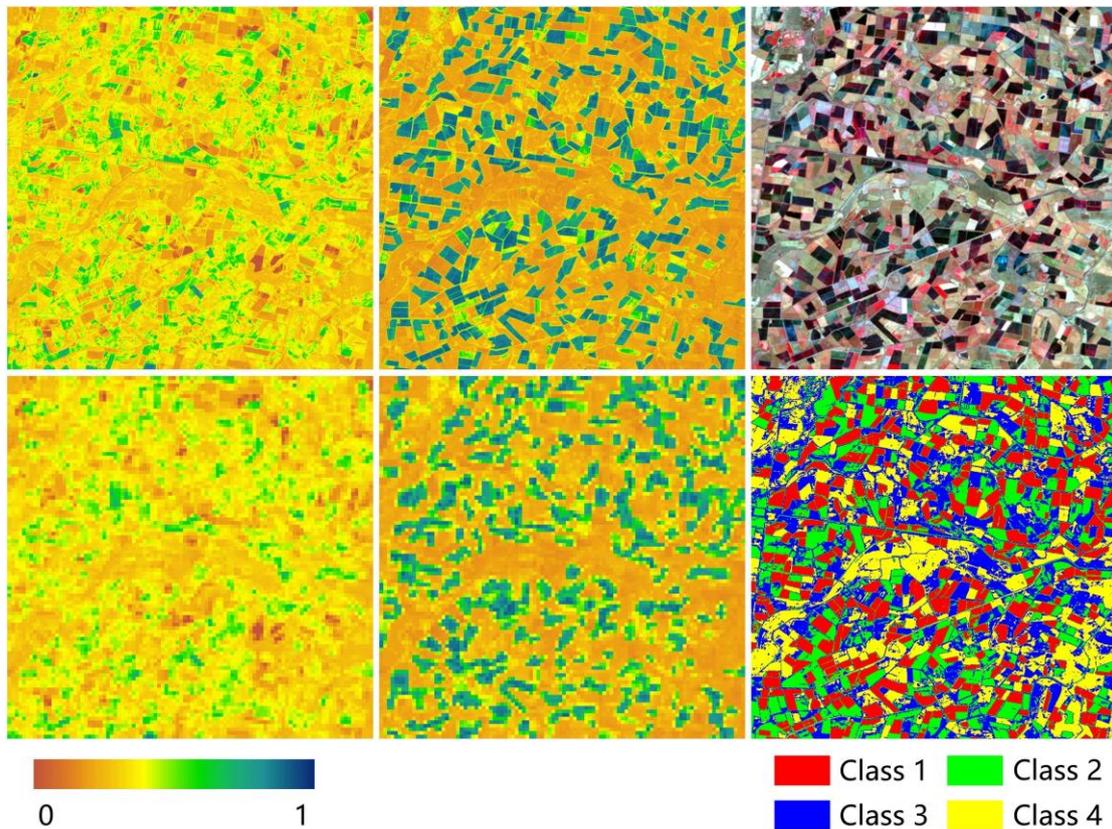351 al., 2006; Rao et al., 2015; Zhu et al., 2016).

352 **3. Data**

353 **3.1 Data for experiments using single cloud-free fine image**

354 　　We used Landsat images without clouded pixels to evaluate the performance of

355 the proposed IFSDAF model at two sites with different land-cover characteristics.

356 Considering that the performance of the existing spatiotemporal fusion methods

357 generally perform well in homogeneous areas (Zhu et al., 2018), the study only tests

358 the performance of the new method in cases with relative complexities (i.e., a

359 heterogeneous site and a site with significant land cover changes.) The Landsat

360 images covering the two sites were shared by Emelyanova et al. (2013) and were also

361 used to test the NDVI-LMGM and FSDAF algorithms (Rao et al., 2015; Zhu et al.,

362 2016).

363 　　This first site is located in the Coleambally irrigated area (34°54′S and 145°57′E),

364 characterized by great heterogeneity in landscape with many small patches of farm

365    land and rapid phenological changes (Fig. 3). Two Landsat ETM+ images (800×800

366    pixels), acquired on November 25th, 2001 ($t_0$) and January 12th, 2002 ($t_p$) during the

367    growing season, were upscaled by the ratio of 8:1 to synthesize MODIS images. In

368    this test, the synthesized MODIS image instead of the real MODIS image was used,

369    because the synthesized MODIS image can exclude the co-registration error (Gevaert

370    and Garcia-Haro, 2015; Wang and Atkinson, 2018; Zhu et al., 2016). This exclusion

371    ensures a fair comparison of different algorithms. The NDVI data were then derived

372    from corresponding reflectance images. Then, the land cover classification map was

373    obtained by the Iterative Self-Organizing Data Analysis Technique (ISODATA)

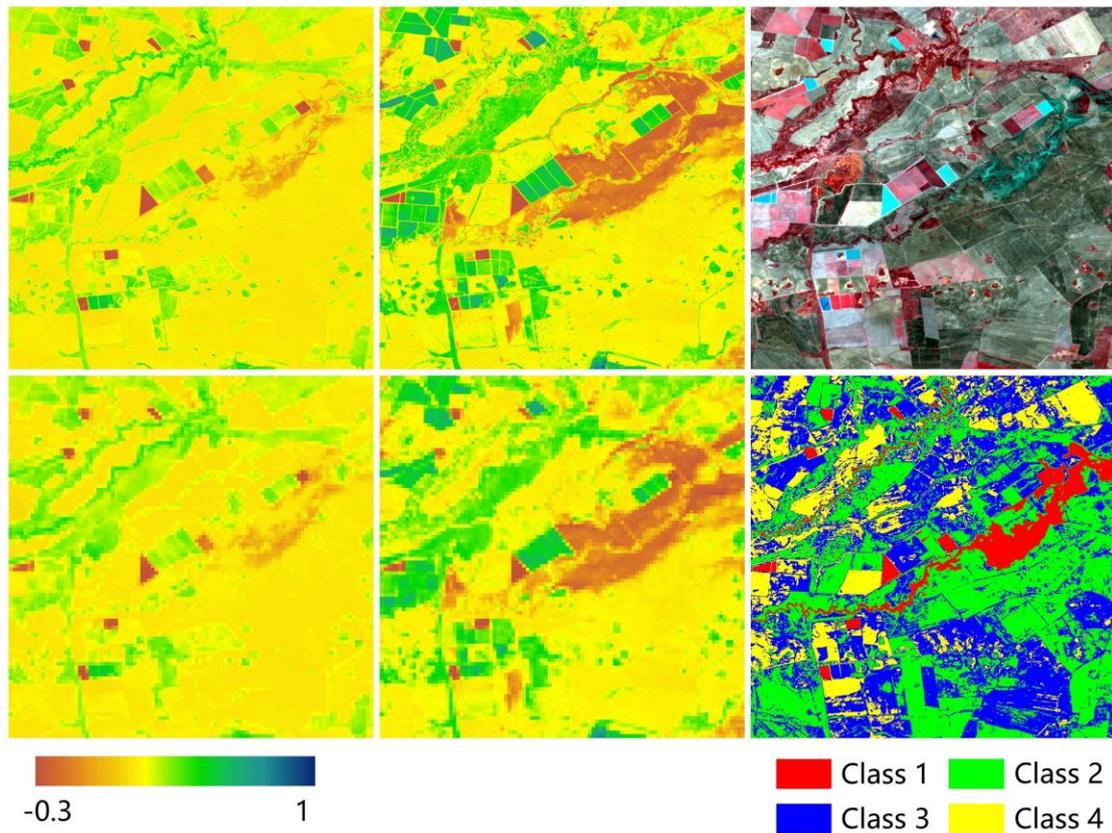374    method based on the Landsat image acquired on November 25th, 2001 ($t_0$).



375

376    **Fig.3.** Test data of the heterogeneous site in Coleambally irrigation area: Landsat

377     NDVI on November 25th, 2001 (a) and January 12th, 2002 (b); false-color-composite

378     Landsat image on November 25th, 2001 (c); MODIS NDVI on November 25th, 2001

379     (d) and January 12th, 2002 (e); and land cover map on November 25th, 2001 by

380     ISODATA (f).

381

382       The second site is located in the Gwydir area (29°07′S and 149°04′E) with a

383     flood event occurred in December 2004. Two Landsat TM images (800×800 pixels)

384     on November 26th, 2004 ($t_0$) and December 12th, 2004 ($t_p$) were used at this site (Fig.

385     4). Abrupt land cover changes can be observed in these two images due to the flood

386     (Emelyanova et al., 2013). Two Landsat images were also upscaled by the ratio of 8:1

387     to synthesize the MODIS images. Then, the NDVI data were derived from all the

388     original images. A land cover classification map was obtained based on Landsat image

389     on November 26th, 2004 ($t_0$) by the ISODATA method.

**Fig.4.** Test data of a site experienced land cover change in Gwydir area: Landsat NDVI on November $26^{th}$, 2004 (a) and December $12^{th}$, 2004 (b); false-color-composite Landsat image on November $26^{th}$, 2004 (c); MODIS NDVI on November $26^{th}$, 2004 (d) and December $12^{th}$, 2004 (e); and classification map on November $26^{th}$, 2004 by ISODATA (f).
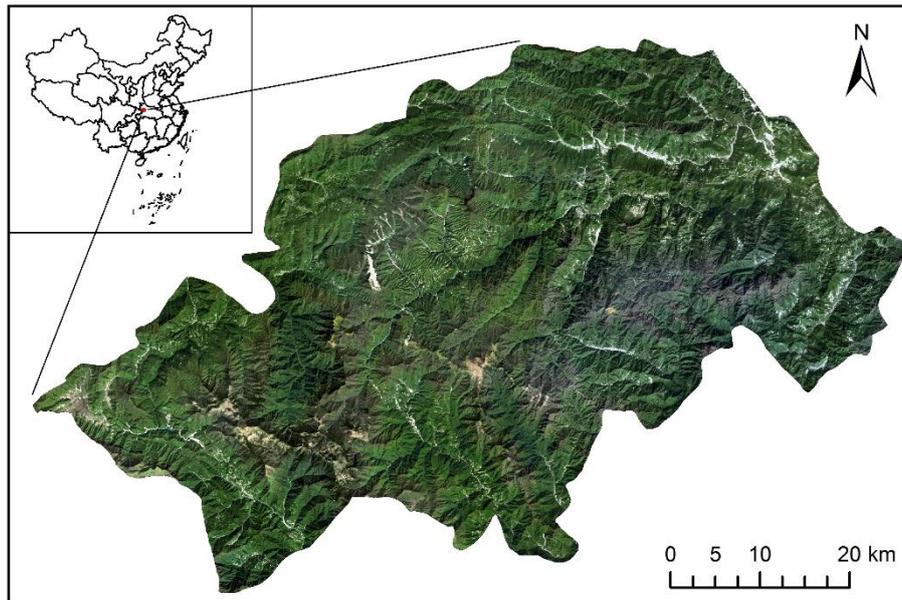
For these two sites, NDVI-LMGM (Rao et al., 2015), STARFM (Gao et al., 2006) and FSDAF (Zhu et al., 2016) were also applied to the same data set for comparison.

**3.2 Data for experiments using multiple clouded fine images**

Experiments using multiple cloudy fine images were implemented to assess the performance of the proposed IFSDAF method for predicting the NDVI time series when the input fine images which were partially contaminated by clouds. To test the
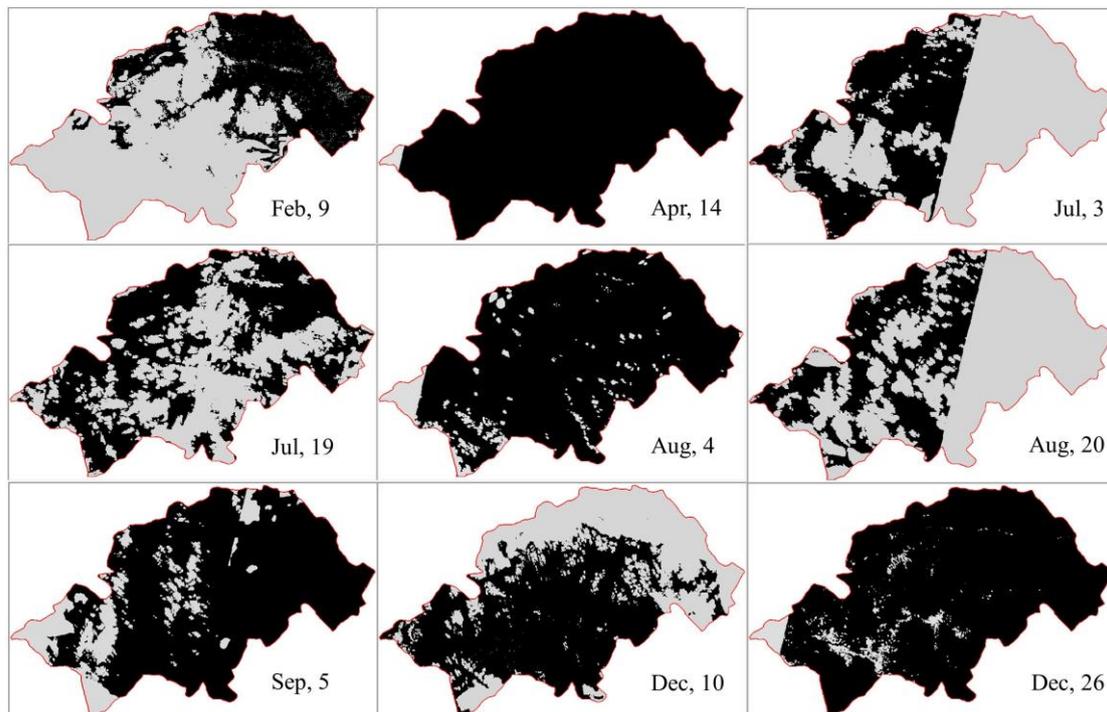
22

402 applicability of IFSDAF for fusing images from diverse sensors, we fused both

403 Landsat and Sentinel-2 images with MODIS in two cloudy sites respectively.

404 The first site is Shennongjia Forestry District (109°59'–110°58' E, 31°15'–31°57'

405 N) located in the western part of Hubei Province, central China (Fig. 5). This area

406 belongs to subtropical monsoon climate; and its elevation ranges from 398 to 3105 m.

407 The vegetation distribution of this area is very heterogeneous with the types of

408 evergreen broadleaf forest, deciduous broadleaf forests, and evergreen coniferous

409 forest. There are also farmlands and artificial surfaces in this area (Wang et al., 2018;

410 Zhao et al., 2005). The 250 m (resampled to 240 m) 16-day composite MODIS NDVI

411 products (MOD13Q1) covering this site in 2015 were acquired from NASA

412 (https://ladsweb.nascom.nasa.gov/search/) and then resampled to the resolution of 240

413 m. Landsat 8 level 2A surface reflectance products and their cloud masks by Fmasks

414 in 2015 were downloaded from the USGS (https://espa.cr.usgs.gov/ordering/new/).

415 All Landsat images were co-registered to MODIS images. M osaic of two adjacent

416 Landsat-8 scenes can cover the whole area of this site. When mosaicking two Landsat

417 8 images with close acquisition dates, pixels in the overlapped part have two NDVI

418 values and the higher one is kept because higher NDVI is less likely affected by poor

419 atmospheric condition. Those Landsat images with clouds, shadows, and snow more

420 than 70% were discarded. Finally, one clear Landsat image ($t_0$, on October 14[th], 2015)

421 and nine partially contaminated Landsat images (Fig. 6) were selected as the input of

422 fine-resolution NDVI images for data fusion.

423

**Fig. 5.** Shennongjia Forestry District in Hubei Province, central China. The image is a

true-color-composite Landsat 8 OLI image acquired on the day of year 287, October

14th, 2015.
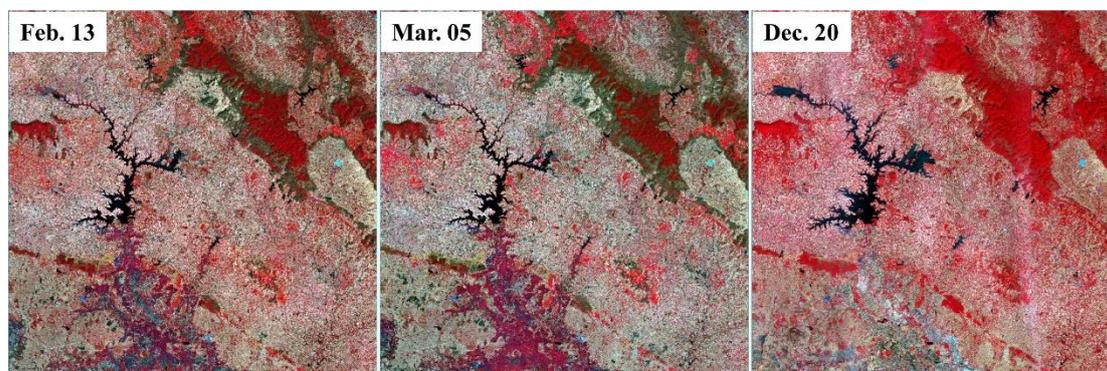


427

**Fig. 6.** Cloud masks of nine partially contaminated Landsat 8 images in Shennongjia

Forestry District by Fmask method, where gray color indicates pixels contaminated by

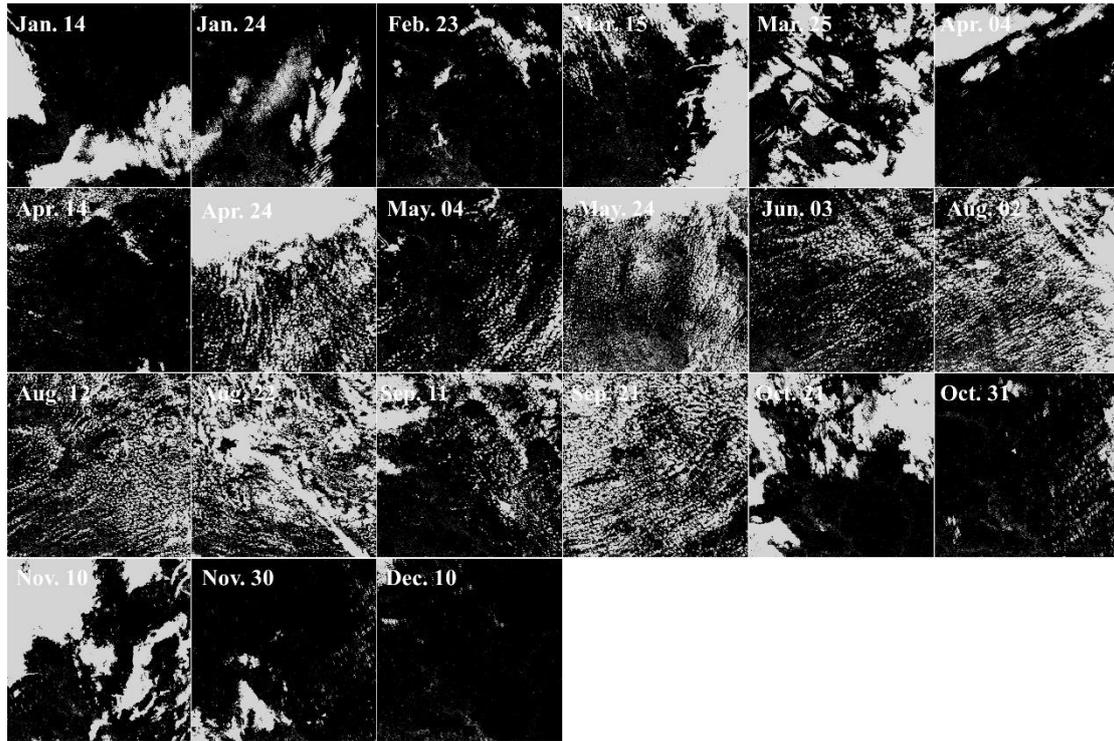clouds and cloud shadows, and black color represents clear pixels.

431

The second site is in Southeast Asia which has complex landscapes with croplands, water, forest and urbans. This site is covered by one Sentinel-2A scene (size 10980×10980 10m pixels) with the tile number of T48QUD. We acquired Sentinel-2A satellite level 1C products from EarthExplorer (https://earthexplorer.usgs.gov/) and 8-day composite MODIS surface reflectance products (MOD09Q1) from NASA. Both images were acquired in 2017. Atmospheric correction of Sentinel-2A images was done with the tool provided by European Space Agency, Sen2Cor (http://step.esa.int/main/third-party-plugins-2/sen2cor/). Cloud masks of sentinel-2A images were produced by the Fmask software (https://github.com/gersl/fmask) and images with cloud cover more than 70% were discarded. Finally, we obtained three clear Sentinel-2A images on Feb. 13th, Mar. 5th and Dec. 20th in 2017 respectively (Fig. 7), and 21 partially cloud contaminated images (Fig. 8). Sentinel-2A NDVI and MODIS NDVI images were then calculated from the surface reflectance data.



**Fig. 7.** False-color-composite of clear Sentinel 2A images on Feb. 13th, Mar. 05th and Nov. 20th in 2017 with the tile number of T48QUD, respectively.

449

**Fig. 8.** Cloud masks of 21 partially cloud contaminated Sentinel-2A images with the

tile number of T48QUD by Fmask method, where black color is clear pixel and gray
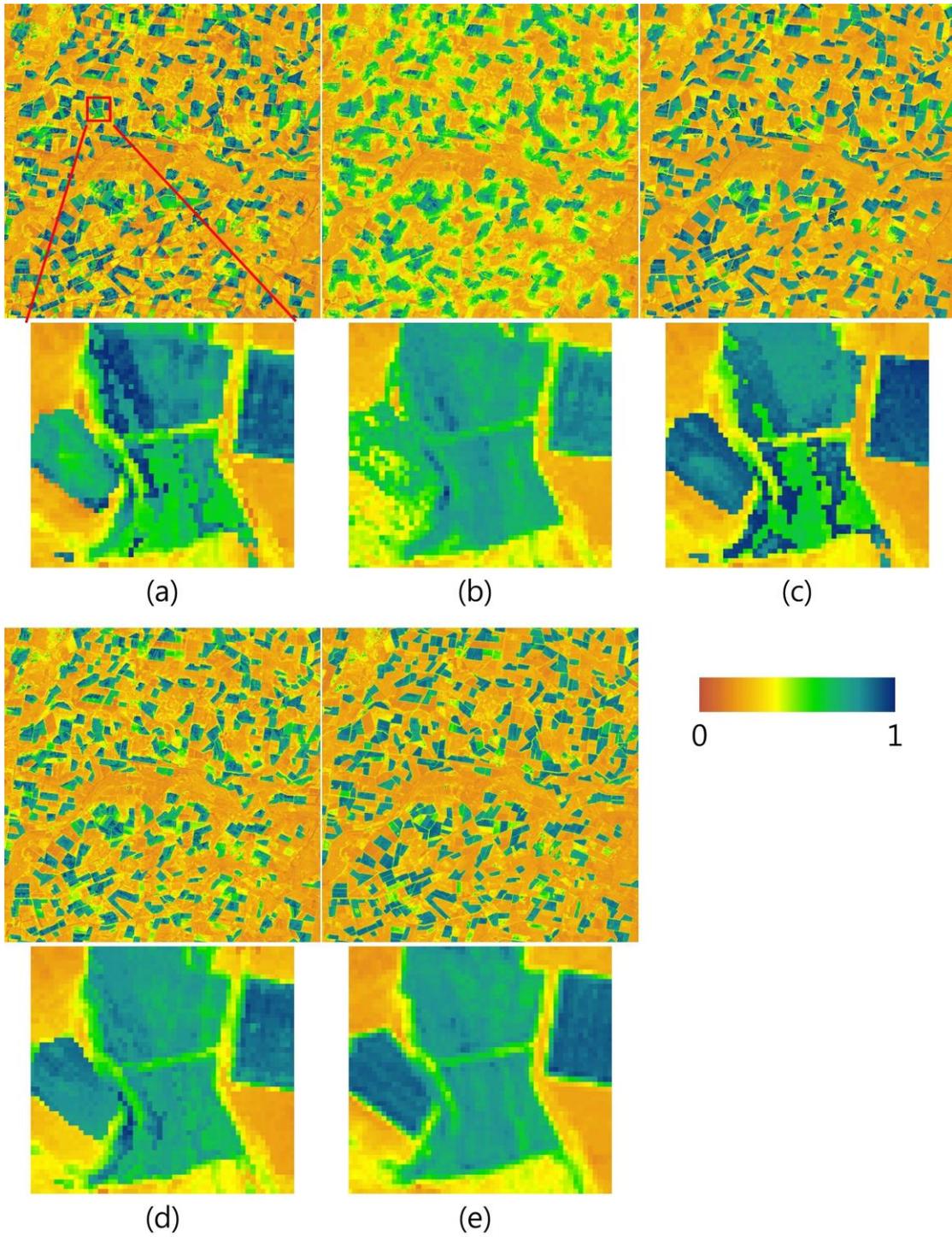
color is cloudy pixel.

## 4. Results

### 4.1 Fusion using single cloud-free fine image

Fig. 9 shows the visual comparison of predicted Landsat NDVI by IFSDAF and

the three existing methods with observed Landsat NDVI on January 12[th], 2002 ($t_p$) for

the farmland site with great heterogeneity and rapid phenological changes. Compared

with the other three methods, the fused image by IFSDAF (Fig. 9d) is more similar to

the actual NDVI image (Fig. 9e) (e.g., the zoomed-in sub-region). On the contrary, the

NDVI-LMGM (Fig.9a) and FSDAF (Fig.9c) methods yield large errors in some

pixels leading to discontinuity in the fused images; and STARFM (Fig.9b) leads to an

462  unsatisfactory blurring effect for small objects. Scatter plots (Fig. 10) and quantitative

463  assessment also confirms that the proposed method obtains the highest accuracy (Table

464  1). The IFSDAF has the lowest RMSE (0.0884), lowest rRMSE (22.12%) and highest *r*

465  (0.9376) for the whole image. Furthermore, the AD (-0.0001) of the newly proposed

466  method is closer to zero, indicating less biased. In addition, the accuracies of the

467  NDVI-LMGM (RMSE= 0.1300 and rRMSE= 32.54%) and STARFM (RMSE= 0.1646

468  and rRMSE= 41.19%) are much lower for the whole image compared with the FSDAF

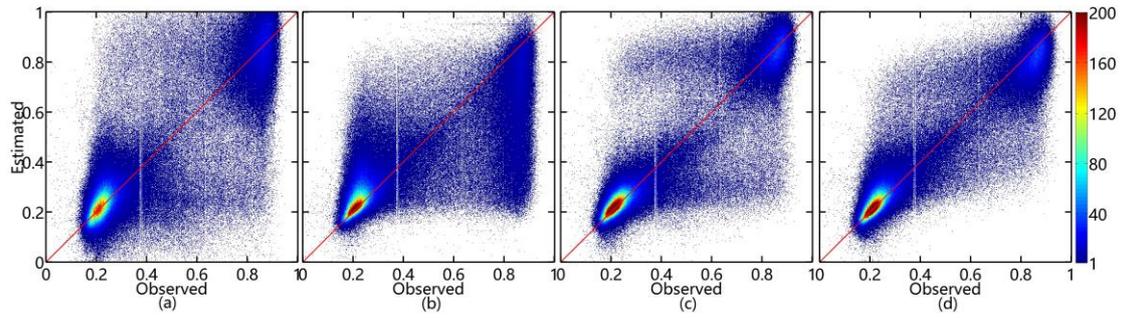469  (RMSE= 0.1002 and rRMSE= 25.06%).

470      As we mentioned, the NDVI normally has a larger variance than raw reflectance

471  bands. Fig. 11a shows histogram distributions of bands Red, NIR, and corresponding

472  NDVI from Landsat image on January 12$^{th}$, 2002 ($t_p$). The NDVI displays two peaks

473  with significantly a greater variance than that of band Red or band NIR due to the

474  amplification of vegetation signals and the suppression of non-vegetation signals. To

475  further investigate the performance of the proposed method in sub-regions with

476  different NDVI magnitudes, based on histogram distribution of NDVI (Fig.11a), the

477  whole image is thus divided into three parts: low NDVI (< 0.4), medium NDVI

478  (0.4-0.7), and high NDVI (>0.7). It can be seen that NDVI-LMGM and STARFM have

479  relatively lower accuracies compared with IFSDAF and FSDAF. And IFSDAF has a

480  better performance than FSDAF in medium NDVI and high NDVI sections.

Fig. 9. Landsat NDVI on January 12th, 2002: predictions by NDVI-LMGM (a), STARFM (b), FSDAF (c), IFSDAF (d) and the actual NDVI (e).

484



485 **Fig. 10.** Scatter plots of estimated results compared with observed value of Landsat

486 NDVI on January 12th, 2002: NDVI-LMGM (a), STARFM (b), FSDAF (c) and
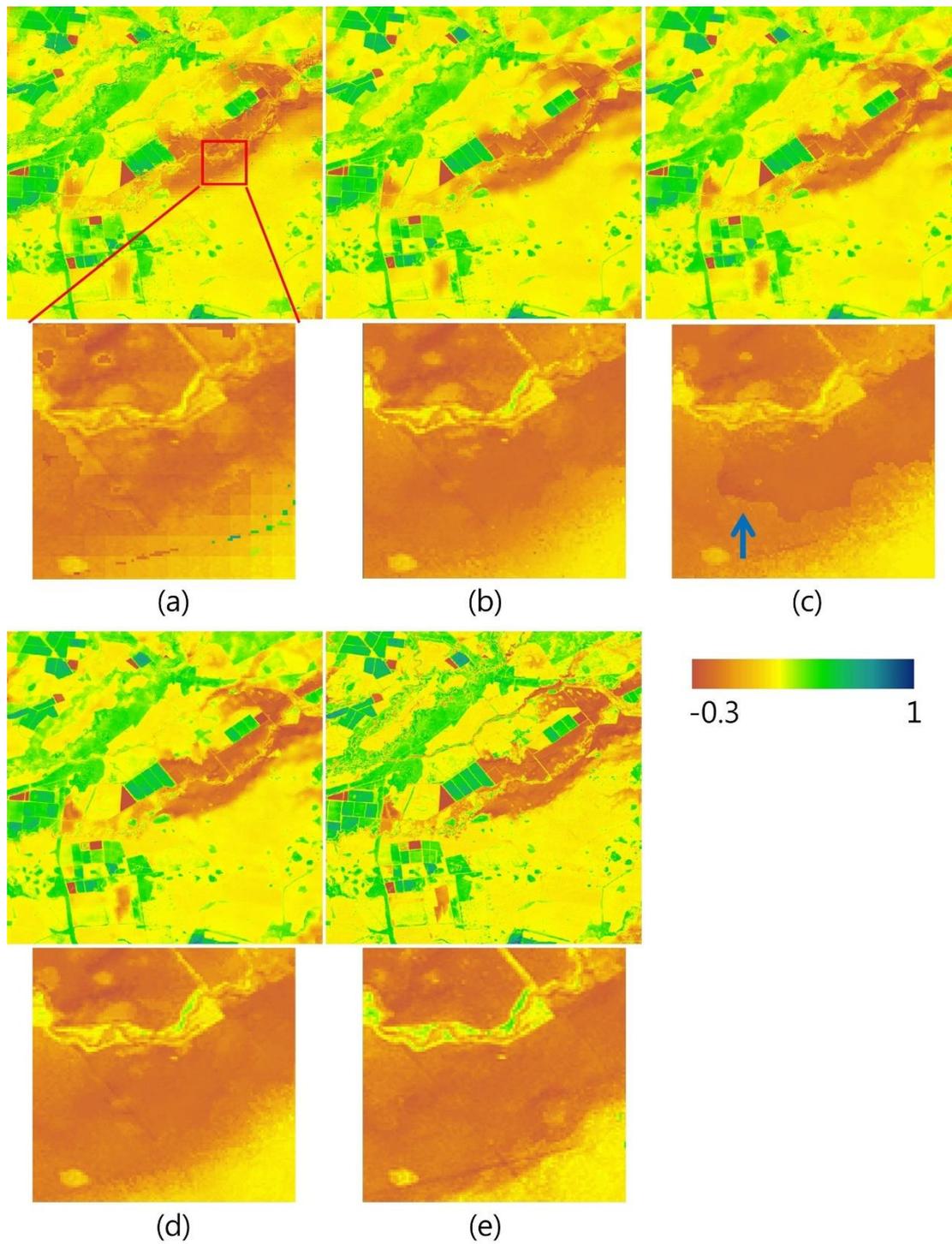
487 IFSDAF (d).



488

489 **Fig. 11.** Histograms of band Red, band NIR and NDVI in Coleambally irrigation area

490 on January 12[th], 2002 (a) and the Gwydir area on December 12[th], 2004 (b).

491 **Table.1.** RMSE, rRMSE, $r$ and AD between predicted NDVI and observed NDVI of NDVI-LMGM, STARFM, FSDAF and the IFSDAF method in the

492 Coleambally irrigation area.

| Method | NDVI-LMGM | | | | STARFM | | | | FSDAF | | | | IFSDAF | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | RMSE | rRMSE | $r$ | AD | RMSE | rRMSE | $r$ | AD | RMSE | rRMSE | $r$ | AD | RMSE | rRMSE | $r$ | AD |
| Low NDVI | 0.0916 | 38.85% | 0.4476 | **0.0148** | 0.1068 | 45.29% | 0.5214 | 0.0564 | 0.0669 | 28.36% | 0.5576 | 0.0160 | **0.0664**[***] | **28.16%** | **0.6334** | 0.0170 |
| Medium NDVI | 0.2415 | 45.11% | 0.2917 | -0.0476 | 0.1482 | 27.69% | 0.2805 | -0.0175 | 0.1962 | 36.64% | 0.3740 | -0.0205 | **0.1473**[***] | **27.51%** | **0.4328** | **-0.0060** |
| High NDVI | 0.1589 | 19.03% | 0.3171 | -0.0477 | 0.2130 | 25.52% | 0.2983 | -0.1656 | 0.1221 | 14.62% | 0.3926 | -0.0426 | **0.1116**[***] | **13.36%** | **0.4493** | **-0.0408** |
| Whole image | 0.1300 | 32.54% | 0.8744 | -0.0053 | 0.1646 | 41.19% | 0.7778 | -0.0295 | 0.1002 | 25.06% | 0.9238 | -0.0012 | **0.0884**[***] | **22.12%** | **0.9376** | **-0.0001** |

493 Note: for t test, * means $p < 0.05$; ** means $p < 0.01$; *** means $p < 0.001$ compared with results of FSDAF.
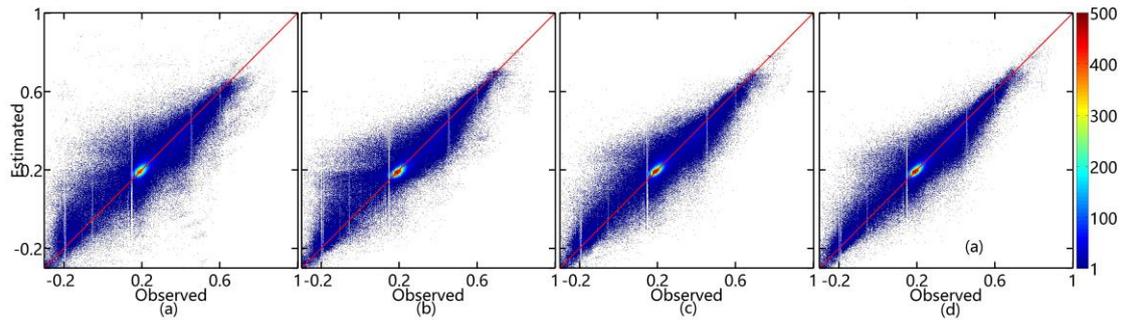
494       For the Gwydir site where a flood event occurred, as shown in Fig. 12, the fusion

495       result of IFSDAF (RMSE = 0.0546) captures the change (Fig. 12d), being more

496       similar to the actual NDVI pattern than the other three methods. NDVI-LMGM gets

497       fused image (RMSE = 0.0794) with significant block effects (Fig. 12a). The result of

498       STARFM (RMSE = 0.0686) is generally similar to the actual NDVI image as shown

499       in Fig. 12b. FSDAF also has high accuracy (RMSE= 0.0617) in fusion but it has

500       abnormal predictions for some pixels shown in the selected area (Fig. 12c). The blue

501       arrow in Fig. 12c indicates the error edges produced by FSDAF. In fact, before the big

502       flood, there is a small river in the zoomed-in area, resulting in the edge (marked by

503       the blue arrow) between water and barren land. However, after the flood, the river

504       overflowed and covered nearby farmland. Thus, the original edge of the river

505       disappeared as shown in the actual Landsat NDVI image of Fig. 12e. IFSDAF is the

506       only among the four methods which can capture this phenomenon. Scatter plots in Fig.

507       13a-d show no obvious bias of these four methods; but points of FSDAF and

508       IFSDAF are closer to the 1:1 line than the other two methods which reveal the

509       comparable capacity of both IFSDAF and FSDAF in capturing land cover changes.

Fig. 12. Landsat NDVI on December 12th, 2004: predictions by NDVI-LMGM (a), STARFM (b), FSDAF (c), IFSDAF (d) and the actual NDVI (e).

**Fig. 13.** Scatter plots of estimated results compared with observed value of Landsat NDVI on December 12$^{th}$, 2004: NDVI-LMGM (a), STARFM (b), FSDAF (c) and IFSDAF (d).

Fig. 11b shows histogram distributions of band Red, band NIR, and NDVI on December 12$^{th}$, 2004 ($t_p$) respectively. The variance of NDVI is also significantly higher than that of band Red or band NIR. Moreover, due to the flood event, there are many negative values in NDVI, causing three peaks around NDVI = -0.2, NDVI = 0.2 and NDVI = 0.4 in the histogram distribution of NDVI. The whole image is divided into three parts (low NDVI < 0, medium NDVI 0-0.3, and high NDVI > 0.3) to quantitatively assess the accuracy (Table 2). It is clear that the new method yields higher accuracy with lower RMSE of 0.0546, higher $r$ of 0.9527 among the whole image than the other three methods. In the three separate NDVI parts, the new method also displays higher accuracy with RMSE = 0.0798, 0.0467, and 0.0584, respectively.

527 **Table.2.** RMSE, rRMSE, $r$ and AD between predicted NDVI and observed NDVI of NDVI-LMGM, STARFM, FSDAF and the IFSDAF method in the Gwydir

528 area.

| Methods | NDVI-LMGM | | | | STARFM | | | | FSDAF | | | | IFSDAF | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | RMSE | rRMSE | $r$ | AD | RMSE | rRMSE | $r$ | AD | RMSE | rRMSE | $r$ | AD | RMSE | rRMSE | $r$ | AD |
| Low NDVI | 0.1267 | -89.59% | 0.6406 | 0.0708 | 0.0964 | -68.16% | 0.7364 | 0.0497 | 0.0906 | -64.03% | 0.7521 | 0.0431 | **0.0798***** | **-56.42%** | **0.7821** | **0.0337** |
| Medium NDVI | 0.0636 | 33.46% | 0.4872 | 0.0061 | 0.0526 | 27.65% | 0.5794 | **0.0042** | 0.0516 | 27.11% | 0.6040 | 0.0066 | **0.0467***** | **24.55%** | **0.6524** | 0.0045 |
| High NDVI | 0.0858 | 19.67% | 0.7425 | -0.0442 | 0.0654 | 15.00% | 0.8410 | -0.0341 | 0.0679 | 15.55% | 0.8302 | -0.0368 | **0.0584***** | **13.38%** | **0.8658** | **-0.0270** |
| Whole image | 0.0794 | 37.45% | 0.8970 | 0.0013 | 0.0686 | 29.53% | 0.9250 | 0.0028 | 0.0617 | 29.09% | 0.9395 | 0.0002 | **0.0546** | **25.77%** | **0.9527** | **0.0001** |

529 Note: for t test, * means $p < 0.05$; ** means $p < 0.01$; *** means $p < 0.001$ compared with results of FSDAF.

**4.2 Fusion using multiple fine images partially covered by clouds**

530

531    In the site of Shennongjia, each of the four Landsat NDVI images captured on

532    Apr 14$^{th}$, Jul 3$^{th}$, Sep 5$^{th}$, and Dec 10$^{th}$ in 2015 serves as reference data for the

533    independent validation. For example, Apr 14$^{th}$ was predicted by IFSDAF using all

534    other eight partially contaminated fine NDVI images as input, and then clear pixels in

535    the true Apr 14$^{th}$ image were used to assess the accuracy of predicted Apr 14$^{th}$ image.

536    For the comparison, FSDAF only used one clear Landsat NDVI image on Oct. 14$^{th}$ of

537    2015 to predict the above four Landsat NDVI images. The accuracies of fusion results

538    of the four images are summarized in Table 3 and the predictions are shown in Fig. 14.

539    For the purpose of simplification, results of NDVI-LMGM and STARFM are not

540    shown in this experiment because they yielded lower accurate results than FSDAF.
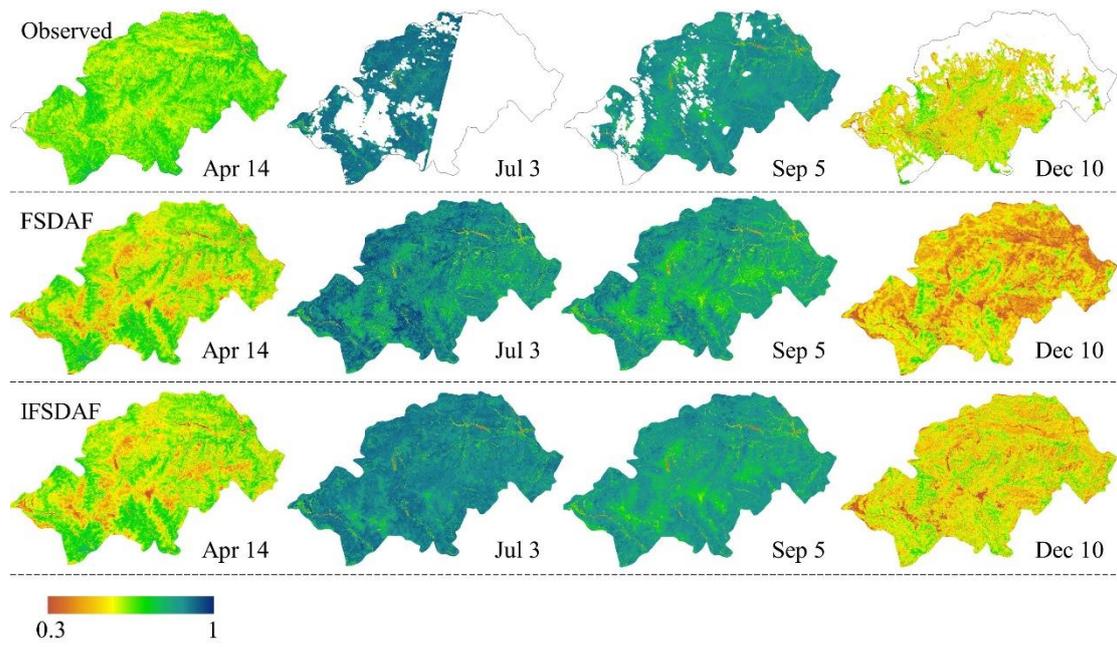
541    It is evident from Fig. 14 that IFSDAF can produce fused images more similar

542    with real Landsat NDVI than FSDAF. In Table 3, RMSE values of IFSDAF on all

543    dates are lower than that of FSDAF. These improvements of accuracy are mainly

544    attributed to the extra information provided by the partially contaminated Landsat

545    images, which can be well used in IFSDAF but not in FSDAF. On the contrary,

546    FSDAF only used one fine image on Oct 14$^{th}$ in 2015 which is far away from some

547    prediction dates, leading to low accuracy on these prediction dates. More important,

548    the improvement of IFSDAF on Jul 3$^{th}$ and Sep 5$^{th}$ during the peak stage of vegetation

549    growth is more significant than other two dates, indicating that IFSDAF may be more

550    effective for fusing images with medium to high NDVI values. This result is similar to

551    the experiment in the Coleambally irrigation area.

552    **Table.3.** RMSE, rRMSE, *r* and AD between the predicted NDVI and observed

553    partially contaminated fine NDVI on Apr 14th, Jul 3th, Sep 5th and Dec 10th in year

554    2015, in the Shennongjia Forestry District.

| Date | Methods | RMSE | rRMSE | *r* | AD |
|---|---|---|---|---|---|
| Apr 14th | FSDAF | 0.0873 | 13.33% | 0.6319 | -0.0481 |
| | IFSDAF | **0.0819***** | **12.51%** | **0.6620** | **-0.0475** |
| Jul 3th | FSDAF | 0.0578 | 6.44% | 0.6504 | -0.0138 |
| | IFSDAF | **0.0368***** | **4.09%** | **0.8508** | **-0.0137** |
| Sep 5th | FSDAF | 0.0671 | 7.86% | 0.7279 | -0.0306 |
| | IFSDAF | **0.0393***** | **4.61%** | **0.8615** | **-0.0173** |
| Dec 10th | FSDAF | 0.1246 | 21.24% | 0.6516 | -0.0729 |
| | IFSDAF | **0.0913***** | **15.57%** | **0.7768** | **-0.0366** |

555    Note: for t test, * means $p < 0.05$; ** means $p < 0.01$; *** means $p < 0.001$ compared

556    with results of FSDAF.

**Fig. 14.** Landsat 8 NDVI in Shennongjia Forestry District on Apr 14th, Jul 3th, Sep 5th, and Dec 10th in year 2015 predicted by FSDAF and IFSDAF, respectively.
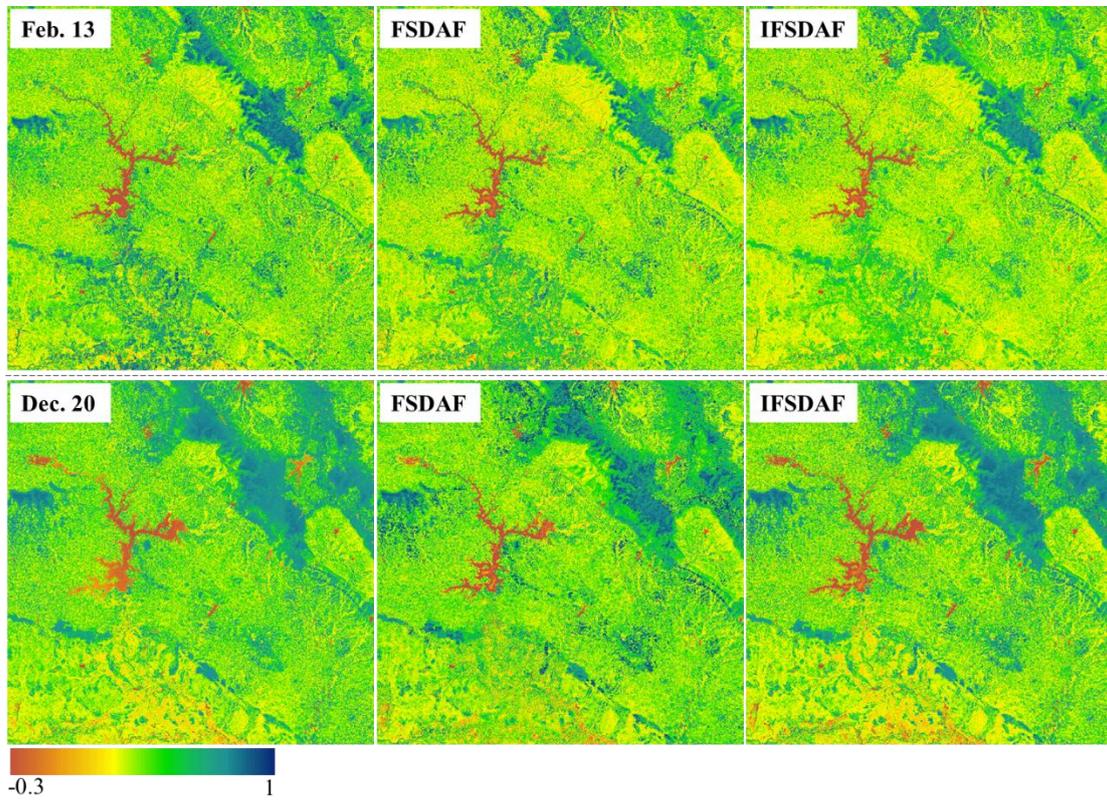
In the South Asia site, clear Sentinel-2A NDVI image on Mar 5th, 2017 was selected as the base image. The other two clear Sentinel NDVI images (Feb. 12th and Dec. 20th) were used as reference data to assess the accuracy of IFSDAF and FSDAF. The base fine spatial resolution NDVI image and the 21 partially cloud contaminated fine NDVI images were used as input for IFSDAF, while only the base fine NDVI image was input to FSDAF. Results in Table. 4 shows that IFSDAF produces more accurate predictions with lower RMSE in both dates (0.0863 and 0.0740) compared with the results by using FSDAF (0.0999 and 0.1469).

**Table. 4.** RMSE, rRMSE, $r$ and AD between the predicted NDVI and observed fine NDVI on Feb. 12th and Dec. 20th, 2017 with Sentinel-2A data.

| Date | Methods | RMSE | rRMSE | $r$ | AD |
|------|---------|------|-------|-----|-----|

| | | | | | |
|---|---|---|---|---|---|
| Feb 13th | FSDAF | 0.0999 | 24.04% | 0.8885 | -0.0463 |
| | IFSDAF | **0.0863***** | **20.76%** | **0.9305** | **-0.0427** |
| Dec 20th | FSDAF | 0.1469 | 33.69% | 0.7401 | **-0.0082** |
| | IFSDAF | **0.0740***** | **17.69%** | **0.9584** | -0.0141 |

Note: for t test, * means $p < 0.05$; ** means $p < 0.01$; *** means $p < 0.001$ compared

with results of FSDAF.



**Fig. 15.** Sentinel-2A NDVI images on Feb. 13th and Dec. 20th (left) and the results

predicted by FSDAF (middle) and IFSDAF (right), respectively.

## 5. Discussion

### 5.1 The way of deriving spatial-dependent increment

In this study, the spatial-dependent increment ($\Delta S$) is acquired based on

difference between interpolation results of coarse NDVI on date $t_p$ and date $t_0$,

579 respectively, as shown in Eq. (4). However, there is also another way of obtaining $\Delta S$

580 in Eq. (12), since the $F_0$ is available on date $t_0$.

$$\Delta S(x_j, y_j) = F_p^{\text{TPS}}(x_j, y_j) - F_0(x_j, y_j) \tag{12}$$

582 where $F_0(x_j, y_j)$ is fine NDVI value of pixel $(x_j, y_j)$ on date $t_0$. However, $\Delta S$ derived from

583 Eq. (4) is a better indicator than that from Eq. (12). A theoretical comparison of these

584 two types of $\Delta S$ is explained as below. Eq. (9) can be simplified as shown in Equation

585 (13), where residuals $R$ are ignored as they are small.

$$\hat{F}_{0,p} = F_0 + \Delta F^{\text{Com}} = F_0 + w_{\text{S}}\Delta S + w_{\text{T}}\Delta T \tag{13}$$

587 For simplification, the notation $(x_i, y_i)$ is removed by replacing Eq. (9) with Eq. (13).

588 And then, replacing $F_0$ by $w_{\text{S}}F_0 + w_{\text{T}}F_0$, as $w_{\text{S}} + w_{\text{T}} = 1$, specifically,

$$\hat{F}_{0,p} = w_{\text{S}}(F_0 + \Delta S) + w_{\text{T}}(F_0 + \Delta T) \tag{14}$$

590 Based on $\Delta S$ in IFSDAF as Eq. (4), Eq. (14) can be written as below,

$$\hat{F}_{0,p} = w_{\text{S}}(F_0 + F_p^{\text{TPS}} - F_0^{\text{TPS}}) + w_{\text{T}}(F_0 + \Delta T) \tag{15}$$

592 Based on $\Delta S$ in Eq. (12), Eq. (14) can also be written as below,

$$\begin{aligned}\hat{F}_{0,p} &= w_{\text{S}}(F_0 + F_p^{\text{TPS}} - F_0) + w_{\text{T}}(F_0 + \Delta T) \\ &= w_{\text{S}}(F_p^{\text{TPS}}) + w_{\text{T}}(F_0 + \Delta T)\end{aligned} \tag{16}$$

594 Difference between Eq. (15) and Eq. (16) is the term $F_0 - F_0^{\text{TPS}}$ in Eq. (15). TPS

595 prediction is a spatially smoothed prediction which loses spatial details to some

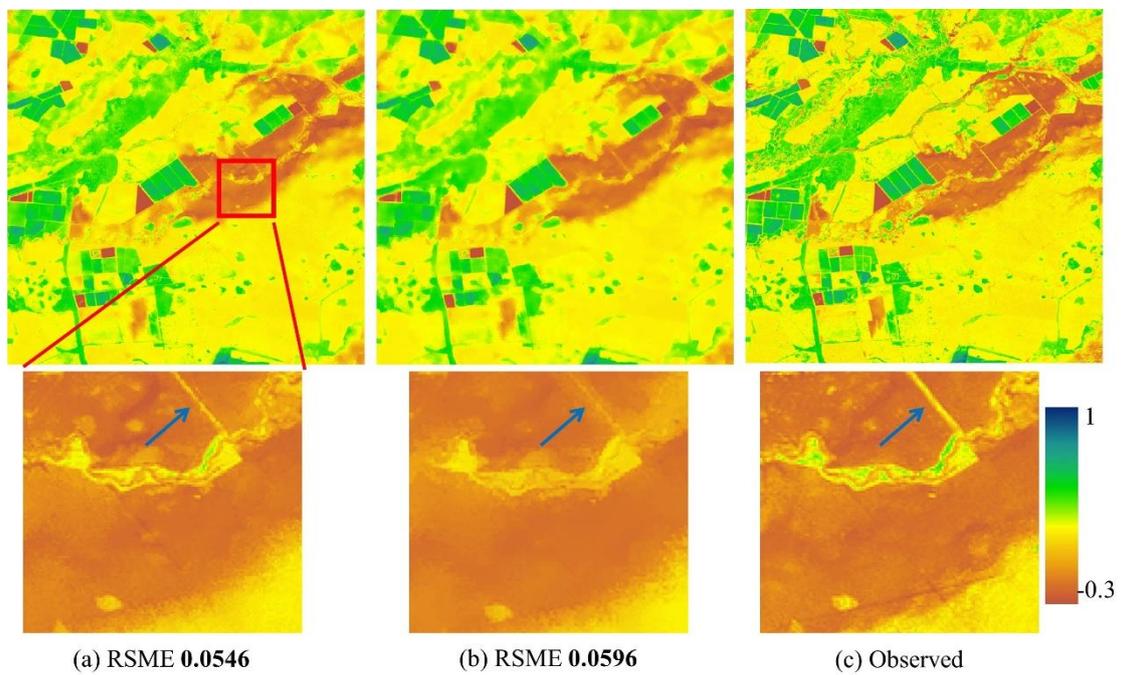596 degree. As a result, $F_0 - F_0^{\text{TPS}}$ functions similarly as a high-pass modulation to model

597 the spatial contrast at $t_0$. This spatial contrast is assumed relatively stable from the $t_0$

598 to $t_p$ in several fusion models (Song and Huang, 2013; Luo et al., 2018), so $F_0 - F_0^{\text{TPS}}$

599 in Eq. (15) can better capture spatial details in the fused image. To demonstrate the

600  abovementioned theoretical analysis, an experiment based on these two types of $\Delta S$

601  was conducted in the Gwydir area. Result shows that the prediction based on $\Delta S$

602  derived from Eq. (4) (Fig. 16a) is more accurate than that from Eq. (12) with RMSEs

603  being 0.0546 and 0.0596, respectively. Moreover, as the zoomed-in pictures in Fig. 16

604  illustrate, the prediction result using Eq. (4) (Fig. 16a) contains more spatial details

605  (e.g., the road marked by blue arrows), providing corroborative support to the

606  analysis.



607  (a) RSME **0.0546**      (b) RSME **0.0596**      (c) Observed

608  **Fig. 16.** Comparison of NDVI value on $t_p$ in the Gwydir area predicted from Eq. (4) (a)

609  Eq. (12) (b), and the real Landsat NDVI image(c).

610  **5.2 The weights scale-invariant assumption**

611  In the proposed IFSDAF method, spatial-dependent increment and temporal

612  increment are combined by optimized weights. As the fine NDVI image on $t_p$ is

613  unknown in the real-world application, coarse NDVI increment ($\Delta C$), upscaled

40

614    spatial-dependent increment ($\Delta C^{S}$), and upscaled temporal increment ($\Delta C^{T}$) are used

615    to derive $w_S$ and $w_T$ in Eq. (5). Such operation assumes that weight is scale-invariant.

616    To verify the assumption, an experiment was conducted in the Coleambally irrigation

617    area, a moving window of 7×7 at a coarse resolution was used to calculate the weights

618    ($w_S$ and $w_T$) for two increments for the center coarse pixel. Because the fine NDVI

619    image $F_p$ actually exists at the site, the two weights at a fine resolution can also be

620    derived based on fine increment ($\Delta F = F_P - F_0$) using the CLS method. Fig. 17a displays

621    the scatter plot of weights derived from the two approaches. All points are close to the

622    1:1 line, where $x$-axis represents the weight of the spatial-dependent increment at the

623    fine resolution and $y$-axis represents the same weight at the coarse resolution,

624    suggesting that the weights derived from both the coarse and fine images are

625    substitutable. Then, combined increments calculated using the two types of weight are

626    very similar (Fig. 17b). RMSE values of combined increments based on weights from

627    the coarse resolution and fine resolution are 0.0941 and 0.0934, respectively. $t$-test

628    shows that there is no significant difference between the two combined increments.

629    Consequently, it can be concluded that the scale effect on the derived weights is

630    minimal, and it will not cause significant errors on the combined increment. Thus, the

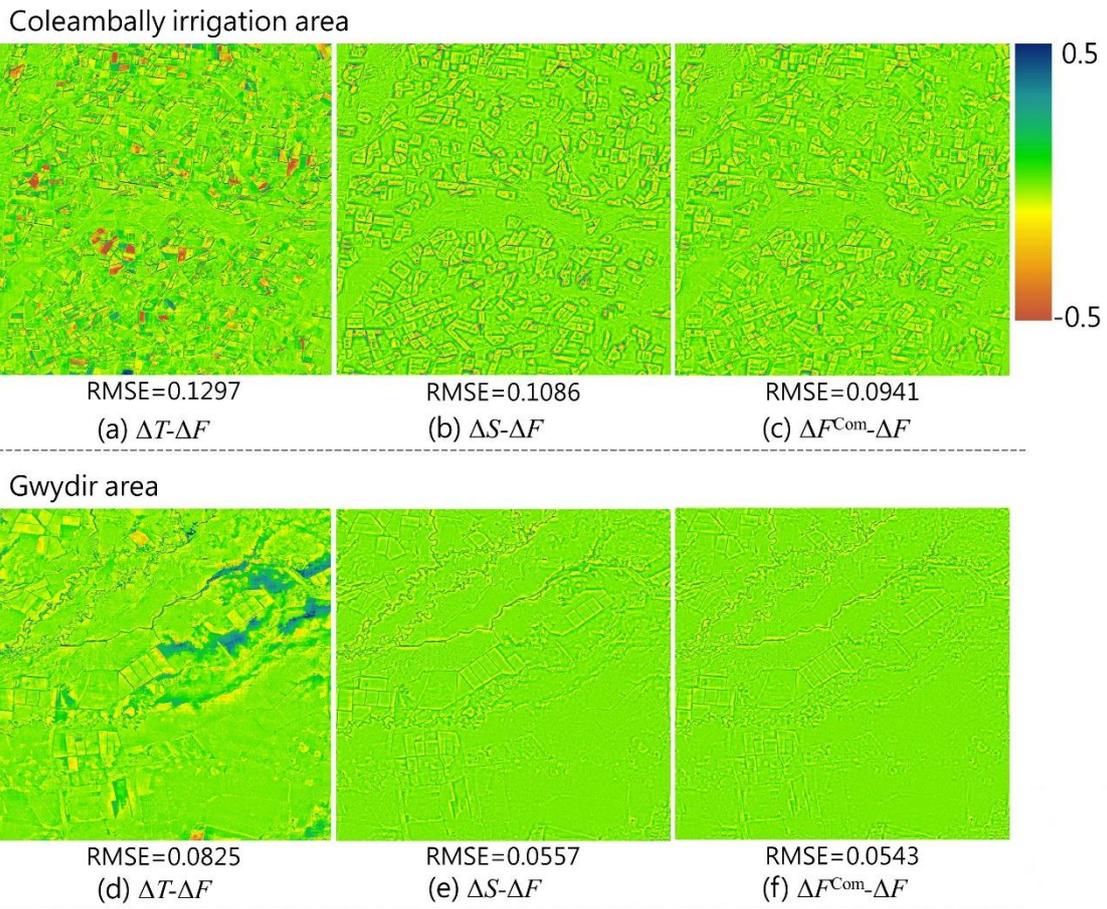631    assumption that weights ($w_S$ and $w_T$) are scale-invariant is reasonable.

**Fig.17.** Scatter plot of weights of spatial-dependent increment based on (a) coarse resolution increment ($\Delta C$) and fine resolution increment ($\Delta F$); (b) comparison of combined increments using weights derived from coarse and fine images.

## 5.3 Combination of temporal and spatial-dependent increments

Temporal increment and spatial-dependent increment are combined in IFSDAF by CLS method in moving windows. Such a combination is based on the assumption that the accuracies of two increment estimations are different under different scenarios; thus the weighted combination is able to improve the accuracy of NDVI prediction through balancing biases in the estimate of two increments. We verified the assumption by comparing the temporal increment, the spatial-dependent increment and the combined increment with the real increment (Fig.18), in which RMSE was used to represent the error of estimation. Theoretically, a good combination is expected to obtain smaller RMSE value than either of the two increment estimations. As shown in Fig.18, the performance of CLS-based combination agrees with our expectation with decreased RMSE values at both study sites, demonstrating the necessity of combining two increments. Moreover, the residual of spatial-dependent

649 increment ($\Delta S$-$\Delta F$) is much more similar to the residual of the combined increment

650 ($\Delta F^{\text{Com}}$-$\Delta F$) than that of the temporal increment ($\Delta T$-$\Delta F$), suggesting that the

651 spatial-dependent increment contributes more to the combined increment than the

652 temporal increment at these two sites.



653

654 **Fig. 18.** Difference between predicted increment and observed increment: difference

655 between (a) temporal increment $\Delta T$, (b) spatial-dependent increment $\Delta S$, (c) combined

656 increment $\Delta F^{\text{Com}}$ and the observed increment $\Delta F$ in the Coleambally irrigation area;

657 difference between (d) temporal increment $\Delta T$, (e) spatial-dependent increment $\Delta S$, (f)

658 combined increment $\Delta F^{\text{Com}}$ and the observed increment $\Delta F$ in the Gwydir area.

659

**5.4 Improvements of IFSDAF compared with FSDAF**

Compared with FSDAF, IFSDAF has improved in the following aspects. First, the increment estimation in FSDAF mainly produced by the unmixing process, and the TPS interpolation result is only used to guide the distribution of residuals rather than producing spatial-dependent increment. However, as shown in Fig.18, the spatial-dependent increment estimated by the TPS interpolation may be more accurate than the temporal increment by the unmixing process. FSDAF underestimates the contribution of the TPS interpolation to some extent. The reason why the spatial-dependent increment is superior to the temporal increment can be found from Table 5, where we calculated the global Moran's I index of the coarse images for band Red, band NIR, and NDVI on the base date $t_0$ and prediction date $t_P$, respectively. The global Moran's I index was used here to measure the spatial autocorrelation of the image, i.e., the relationship of pixel values between neighboring pixels. Larger Moran's I index indicates higher spatial autocorrelation. Table 5 shows that the spatial autocorrelation of NDVI represented by Moran's I index is greater than both the Red and NIR bands, because NDVI, as a feature-enhancing index, can enlarge the data variance compared with the Red and NIR bands (Fig.11a-b). As well known, greater spatial autocorrelation can yield more accurate result in spatial interpolation. Accordingly, the spatial-dependent increment estimated by the TPS interpolation for NDVI should be more accurate than that of Red and NIR bands. Therefore, spatial-dependent increment is more important for fusing NDVI than the raw bands,

681   which greatly benefits the NDVI fusion in IFSDAF. This study implies that without

682   combined with temporal increment, spatial-dependent increment itself may obtain

683   acceptable fusion results. This solution can greatly simplify the fusion process and

684   reduce the computing cost. The simplified fusion model is more effective for the

685   applications in larger areas and when the scale difference between coarse and fine

686   images are not too large, which ensures adequate number of sample points (the center

687   of coarse pixel) to obtain accurate prediction of the fine image by TPS interpolation.

688   **Table. 5**. Moran's I of band Red, band NIR, and NDVI on the base date $t_0$ and

689   prediction date $t_p$ in both the Coleambally irrigation area and Gwydir area at the

690   coarse resolution.

| Band | Base date $t_0$ | | | Prediction date $t_p$ | | |
|---|---|---|---|---|---|---|
| | Red | NIR | NDVI | Red | NIR | NDVI |
| Coleambally irrigation area | 0.5048 | 0.5225 | 0.6439 | 0.5677 | 0.4764 | 0.6840 |
| Gwydir area | 0.5867 | 0.7069 | 0.7881 | 0.6401 | 0.7568 | 0.8584 |

691

692       Second, IFSDAF uses a better way to combine two increments while FSDAF

693   uses only one increment. As we know, the collinearity effect impacts the accuracy of

694   unmixing for temporal increment estimation. Moreover, errors in the classification

695   map and change of land cover also cause uncertainties of temporal increment. In order

696   to correct the potential errors in the temporal increment, FSDAF introduces a

697   homogeneity index $HI(x_j, y_j)$ (over the range of 0-1, derived from the classification

698    map at $t_0$) to help allocate residuals $R(x, y)$ within the coarse pixel. However, when

699    there are land cover changes and misclassification, $HI(x_j, y_j)$ calculated from the

700    classification map at $t_0$ will not be suitable for allocating residuals on date $t_p$. Under

701    this circumstance, the effectiveness of residuals distribution in FSDAF is restricted.

702    Unlike FSDAF, the IFSDAF employs CLS method in moving windows and avoids

703    the use of the homogeneity index; moreover, it allows the final increment estimation

704    with local and adaptive capacity to better combine temporal and spatial-dependent

705    increment.

706        The third improvement of IFSDAF is that it can employ fine NDVI images

707    partially contaminated by clouds, in which clear pixels also provide valuable

708    information. In IFSDAF, the clear pixels in those fine images are also used as base

709    date to estimate the fine NDVI values at the prediction date respectively, and all

710    predictions on date $t_p$ are finally integrated by weights based on the temporal change

711    magnitude in NDVI between base and prediction dates. This weighted prediction can

712    reduce the critical dependency to the clear fine NDVI image and alleviate prediction

713    uncertainties if date of the clear fine NDVI image is far from the prediction date. Of

714    course, the better use of partially contaminated images needs accurate cloud labeling

715    method (e.g., Fmask method). If there are mistakes in cloud labels, estimation results

716    of IFSDAF will be impacted. For instance, land surface with high reflectance (e.g.,

717    sand or snow) is possibly misidentified as clouds (Chen et al., 2016). Moreover,

718    Fmask sometimes omits thin clouds, resulting in cloudy pixels being used in the

719    process of data fusion. Fortunately, effect of the errors in cloud mask can be

720    minimized in IFSDAF because of the weighted combination of predictions from

721    multiple dates. Moreover, with the advance of cloud screening methods (Zhu and

722    Helmer, 2018), the issue can be greatly alleviated. Besides, like other existing

723    spatiotemporal fusion methods, a pure clear image guarantees all pixels to be

724    predicted by IFSDAF. However, in some areas such as tropical areas (e.g., Amazonia),

725    it is difficult to obtain such clear fine image during a long period. Under this condition,

726    using all partially contaminated fine images instead of a purely clear image is a

727    practical choice in IFSDAF although it may leave some pixels not predicted in the

728    fused images if these pixels do not have any one cloud-free observation in the time

729    series.

730    **5.5 Applications to other remote sensing products**

731        Although IFSDAF is designed for the spatiotemporal fusion of NDVI time series,

732    it can also be applied to fuse other vegetation indices like Enhanced Vegetation Index

733    (EVI) and other products such as surface reflectance. To test the applicability of

734    IFSDAF to other products, we assess the performance of IFSDAF in fusion of EVI,

735    Red and NIR bands in Coleambally Irrigation area where have great heterogeneity.

736    RMSEs of the fused images on Jan. 12[th], 2002 (Table. 6) suggest that IFSDAF

737    produces higher accuracy that FSDAF when fusing EVI, while when predicting

738    surface reflectance (Red and NIR bands) the accuracy of IFSDAF and FSDAF does

739    not differ too much. These results confirm that IFSDAF is more suitable than the

740  original FSDAF model for fusing remote sensing products with high spatial

741  autocorrelation.

742  **Table. 6.** RMSE (rRMSE) of predictions by IFSDAF and FSDAF on EVI, Red and

743  NIR bands.

| index/surface reflectance | EVI | Red | NIR |
|---|---|---|---|
| FSDAF | 0.0755 (29.11%) | 0.0271 (19.03%) | 0.0341 (11.02%) |
| IFSDAF | 0.0650 (25.09%) | 0.0245 (17.21%) | 0.0337 (10.91%) |

744  **6.  Conclusions**

745  In this study, we proposed an improved FSDAF method specifically for

746  producing NDVI time series with a high spatiotemporal resolution. Coarse NDVI

747  (MODIS) and fine NDVI images (Landsat and Sentinel) were used to test the

748  performance of the new method for different sensors. Experiments show that the

749  fused NDVI images by IFSDAF is more accurate than FSDAF as well as other two

750  existing methods (NDVI-LMGM and STARFM) in areas with a great degree of

751  spatial heterogeneity and with significant land cover changes. The better performance

752  of IFSDAF can be attributed to producing spatial-dependent increment by the TPS

753  interpolation, employing CLS method in moving windows to adaptively combine the

754  temporal increment and the spatial-dependent increment, as well as the better use of

755  partially contaminated fine images. Such significant improvements are made in

756  accordance with the characteristics of NDVI with larger data variance and spatial

757  autocorrelation compared with raw reflectance bands. Considering the significant

758     contribution of spatial-dependent increment by the TPS interpolation, when the scale

759     difference between coarse and fine images is not very large, the proposed IFSDAF

760     method can be further simplified by only using spatial-dependent increment to

761     improve the efficiency. This result of the study also supports the IFSDAF to be a

762     feasible method for applications in a large area and different sensors. Moreover, it is

763     also applicable to other vegetation index data. We call for more testing of the new

764     method by using other satellite data (e.g. Sentinel and VIIRS data) and in other areas.

765 **Acknowledgement**

771

772    **Reference**

773    Busetto, L., Meroni, M., & Colombo, R. (2008). Combining medium and coarse spatial

774    resolution satellite data to improve the estimation of sub-pixel NDVI time series.

775    *Remote Sensing of Environment, 112*, 118-131

776    Chen, J., Chen, J., Liao, A.P., Cao, X., Chen, L.J., Chen, X.H., He, C.Y., Han, G., Peng,

777    S., Lu, M., Zhang, W.W., Tong, X.H., & Mills, J. (2015). Global land cover mapping at

778    30 m resolution: A POK-based operational approach. *ISPRS Journal of*

779    *Photogrammetry and Remote Sensing, 103*, 7-27

780    Chen, J., Jonsson, P., Tamura, M., Gu, Z.H., Matsushita, B., & Eklundh, L. (2004). A

781    simple method for reconstructing a high-quality NDVI time-series data set based on the

782    Savitzky-Golay filter. *Remote Sensing of Environment, 91*, 332-344

783    Chen, S.L., Chen, X.H., Chen, J., Jia, P.F., Cao, X., & Liu, C.Y. (2016). An Iterative

784    Haze Optimized Transformation for Automatic Cloud/Haze Detection of Landsat

785    Imagery. *IEEE Transactions on Geoscience and Remote Sensing, 54*, 2682-2694

786    Chen, X., Li, W., Chen, J., Rao, Y., & Yamaguchi, Y. (2014). A Combination of

787    TsHARP and Thin Plate Spline Interpolation for Spatial Sharpening of Thermal

788    Imagery. *Remote Sensing, 6*, 2845-2863

789    Chen, X., Wang, D., Chen, J., Wang, C., & Shen, M.G. (2018). The Mixed Pixel Effect

790    in Land Surface Phenology: A Simulation Study. *Remote Sensing of Environment*, *211*,

791    388-344.

792    Chen, X.H., Liu, M., Zhu, X.L., Chen, J., Zhong, Y.F., & Cao, X. (2018).

793 "Blend-then-Index" or "Index-then-Blend": A Theoretical Analysis for Generating

794 High-resolution NDVI Time Series by STARFM. *Photogrammetric Engineering and*

795 *Remote Sensing, 84*, 66-74

796 Dubrule, O. (1984). Comparing Splines and Kriging. *Computers & Geosciences, 10*,

797 327-338

798 Emelyanova, I.V., McVicar, T.R., Van Niel, T.G., Li, L.T., & van Dijk, A.I. (2013).

799 Assessing the accuracy of blending Landsat–MODIS surface reflectances in two

800 landscapes with contrasting spatial and temporal dynamics: A framework for algorithm

801 selection. *Remote Sensing of Environment, 133*, 193-209

802 Gao, F., Masek, J., Schwaller, M., & Hall, F. (2006). On the blending of the Landsat and

803 MODIS surface reflectance: Predicting daily Landsat surface reflectance. *IEEE*

804 *Transactions on Geoscience and Remote Sensing, 44*, 2207-2218

805 Gevaert, C.M., & Garcia-Haro, F.J. (2015). A comparison of STARFM and an

806 unmixing-based algorithm for Landsat and MODIS data fusion. *Remote Sensing of*

807 *Environment, 156*, 34-44

808 Hilker, T., Wulder, M.A., Coops, N.C., Linke, J., McDermid, G., Masek, J.G., Gao, F.,

809 & White, J.C. (2009). A new data fusion model for high spatial- and

810 temporal-resolution mapping of forest disturbance based on Landsat and MODIS.

811 *Remote Sensing of Environment, 113*, 1613-1627

812 Huang, B., & Song, H. (2012). Spatiotemporal Reflectance Fusion via Sparse

813 Representation. *IEEE Transactions on Geoscience and Remote Sensing, 50*, 3707-3716

814 Huang, B., & Zhang, H.K. (2014). Spatio-temporal reflectance fusion via unmixing:

815 accounting for both phenological and land-cover changes. *International Journal of*

816 *Remote Sensing, 35*, 6213-6233

817 Huete, A., Didan, K., Miura, T., Rodriguez, E.P., Gao, X., & Ferreira, L.G. (2002).

818 Overview of the radiometric and biophysical performance of the MODIS vegetation

819 indices. *Remote Sensing of Environment, 83*, 195-213

820 Jarihani, A., McVicar, T., Van Niel, T., Emelyanova, I., Callow, J., & Johansen, K.

821 (2014). Blending Landsat and MODIS Data to Generate Multispectral Indices: A

822 Comparison of "Index-then-Blend" and "Blend-then-Index" Approaches. *Remote*

823 *Sensing, 6*, 9213-9238

824 Lillesaeter, O. (1982). Spectral reflectance of partly transmitting leaves: Laboratory

825 measurements and mathematical modeling. *Remote Sensing of Environment*, 12,

826 247-254

827 Liu, X., Deng, C.W., Wang, S.G., Huang, G.B., Zhao, B.J., & Lauren, P. (2016). Fast

828 and Accurate Spatiotemporal Fusion Based Upon Extreme Learning Machine. *IEEE*

829 *Geoscience and Remote Sensing Letters, 13*, 2039-2043

830 Luo, Y., Guan, K., & Peng, J. (2018). STAIR: A generic and fully-automated method

831 to fuse multiple sources of optical satellite data to generate a high-resolution, daily

832 and cloud-/gap-free surface reflectance product. *Remote Sensing of Environment, 214*,

833 87-99

834 Meng, J.H., Du, X., & Wu, B.F. (2013). Generation of high spatial and temporal

835  resolution NDVI and its application in crop biomass estimation. *International Journal*

836  *of Digital Earth, 6*, 203-218

837  Pettorelli, N., Vik, J.O., Mysterud, A., Gaillard, J.M., Tucker, C.J., & Stenseth, N.C.

838  (2005). Using the satellite-derived NDVI to assess ecological responses to

839  environmental change. *Trends in Ecology & Evolution, 20*, 503-510

840  Rao, Y., Zhu, X., Chen, J., & Wang, J. (2015). An Improved Method for Producing

841  High Spatial-Resolution NDVI Time Series Datasets with Multi-Temporal MODIS

842  NDVI Data and Landsat TM/ETM+ Images. *Remote Sensing, 7*, 7865-7891

843  Rouse, J.W., Jr., Haas, R.H., Schell, J.A., & Deering, D.W. (1974). Monitoring

844  vegetation systems in the Great Plains

845  with ERTS. *In Proceedings of Third ERTS Symposium, Washington, DC, USA, 10–14*

846  *December 1973*, 309–317

847  Song, H., & Huang, B. (2013). Spatiotemporal Satellite Image Fusion Through

848  One-Pair Image Learning. *IEEE Transactions on Geoscience and Remote Sensing, 51*,

849  1883-1896

850  Tian, F., Wang, Y.J., Fensholt, R., Wang, K., Zhang, L., & Huang, Y. (2013). Mapping

851  and Evaluation of NDVI Trends from Synthetic Time Series Obtained by Blending

852  Landsat and MODIS Data around a Coalfield on the Loess Plateau. *Remote Sensing, 5*,

853  4255-4279

854  Wang, P., Teng, M., He, W., Tang, C., Yang, J., & Yan, Z. (2018). Using habitat

855  selection index for reserve planning and management for snub-nosed golden monkeys

856    at landscape scale. *Ecological Indicators, 93*, 838-846

857    Wang, Q.M., & Atkinson, P.M. (2018). Spatio-temporal fusion for daily Sentinel-2

858    images. *Remote Sensing of Environment, 204*, 31-42

859    Wu, B., Huang, B., & Zhang, L. (2015). An Error-Bound-Regularized Sparse Coding

860    for Spatiotemporal Reflectance Fusion. *IEEE Transactions on Geoscience and Remote*

861    *Sensing, 53*, 6791-6803

862    Zhang, H.K.K., Huang, B., Zhang, M., Cao, K., & Yu, L. (2015). A generalization of

863    spatial and temporal fusion methods for remotely sensed surface parameters.

864    *International Journal of Remote Sensing, 36*, 4411-4445

865    Zhao, C.M., Chen, W.L., Tian, Z.Q., & Xie, Z.Q. (2005). Altitudinal pattern of plant

866    species diversity in Shennongjia Mountains, central China. *Journal of Integrative*

867    *Plant Biology, 47*, 1431-1449

868    Zhu, X., Cai, F., Tian, J., & Kay-AnnWilliams, T. (2018). Spatiotemporal Fusion of

869    Multisource Remote Sensing Data: Literature Survey, Taxonomy, Principles,

870    Applications, and Future Directions. *Remote Sensing, 10*, 527

871    Zhu, X., Chen, J., Gao, F., Chen, X., & Masek, J.G. (2010). An enhanced spatial and

872    temporal adaptive reflectance fusion model for complex heterogeneous regions.

873    *Remote Sensing of Environment, 114*, 2610-2623

874    Zhu, X., & Helmer, E.H. (2018). An automatic method for screening clouds and cloud

875    shadows in optical satellite image time series in cloudy regions. *Remote Sensing of*

876    *Environment, 214*, 135-153

877 Zhu, X., Helmer, E.H., Gao, F., Liu, D., Chen, J., & Lefsky, M.A. (2016). A flexible

878 spatiotemporal method for fusing satellite images with different resolutions. *Remote*

879 *Sensing of Environment, 172*, 165-177

880 Zhu, Z., & Woodcock, C.E. (2012). Object-based cloud and cloud shadow detection in

881 Landsat imagery. *Remote Sensing of Environment, 118*, 83-94

882 Zurita-Milla, R., Clevers, J., & Schaepman, M.E. (2008). Unmixing-Based Landsat

883 TM and MERIS FR Data Fusion. *IEEE Geoscience and Remote Sensing Letters, 5*,

884 453-457

885

886 **Appendix**

887 Useful notations.

| | | | |
|---|---|---|---|
| $t_0$ | Base date | $f_l(x, y)$ | Fraction of class $l$ within coarse pixel $(x, y)$ |
| $t_p$ | Prediction date | $\Delta F_c$ | Fine spatial resolution increment of class $c$ within the moving window |
| $(x, y)$ | Location of coarse spatial resolution pixel $(x, y)$ | $F_0^{\text{TPS}}$ | Result of TPS interpolation based on coarse NDVI on $t_0$ |
| $(x_j, y_j)$ | Location of jth fine spatial resolution pixel within coarse pixel $(x, y)$ | $F_p^{\text{TPS}}$ | Result of TPS interpolation based on coarse NDVI on $t_p$ |
| $F_0$ | Fine spatial resolution NDVI on $t_0$ | $w_s$ | Weight of spatial-dependent increment |
| $F_p$ | Fine spatial resolution NDVI on $t_p$ | $w_T$ | Weight of temporal increment |
| $C_0$ | Coarse spatial resolution NDVI on $t_0$ | $\hat{F}_p$ | Fine spatial resolution prediction on date $t_p$ |
| $C_p$ | Coarse spatial resolution NDVI on $t_p$ | $\hat{F}_{0, p}$ | Fine spatial resolution prediction on date $t_p$ based on fine NDVI on date $t_0$ |
| $\Delta F$ | Fine spatial resolution NDVI increment | $\hat{F}_{p+1, p}$ | Fine spatial resolution prediction on date $t_p$ based on fine NDVI on date $p+1$ |

| | | | |
|---|---|---|---|
| $\Delta C$ | Coarse spatial resolution NDVI increment | $\Delta F^{\mathrm{Com}}$ | Combined fine spatial resolution increment based on $\Delta T$ and $\Delta S$ |
| $\Delta T$ | Fine spatial resolution temporal increment | $R(x, y)$ | Residual within the coarse pixel $(x, y)$ |
| $\Delta S$ | Fine spatial resolution spatial-dependent increment | $C_q^i(x, y)$ | $i$th coarse pixel in the moving window centered by coarse pixel $(x, y)$ on date $q$ |
| $\Delta C^{\mathrm{T}}$ | Upscaled fine spatial resolution temporal increment | $C_p^i(x, y)$ | $i$th coarse pixel in the moving window centered by coarse pixel $(x, y)$ on date $t_p$ |
| $\Delta C^{\mathrm{S}}$ | Upscaled fine spatial resolution spatial-dependent increment | $w_{q\text{-}p}(x, y)$ | Contribution coefficient of fine spatial resolution pixels on date $q$ to the final prediction on $t_p$ within coarse pixel $(x, y)$ |

888