



1     **Towards ensemble assimilation of hyperspectral satellite observations with data**  
2     **compression and dimension reduction using principal component analysis**

3  
4                     Yinghui Lu and Fuqing Zhang

5  
6     Department of Meteorology and Atmospheric Science, and Center for Advanced Data  
7     Assimilation and Predictability Techniques, The Pennsylvania State University, University  
8     Park, Pennsylvania

9  
10  
11  
12     Submitted to *Monthly Weather Review* for publication as an article

13     December 2018

14  
15  
16  
17  
18  
19  
20  

---

21     *Corresponding author:* Professor Fuqing Zhang, Department of Meteorology and  
22     Atmospheric Science, The Pennsylvania State University, University Park, PA 16802.

23     E-mail: [fzhang@psu.edu](mailto:fzhang@psu.edu)

**Early Online Release:** This preliminary version has been accepted for publication in *Monthly Weather Review*, may be fully cited, and has been assigned DOI 10.1175/MWR-D-18-0454.1. The final typeset copyedited article will replace the EOR at the above DOI when it is published.

24 **Abstract**

25 Satellite-based hyperspectral radiometers usually have thousands of infrared channels  
26 which contain atmospheric state information with higher vertical resolution compared to  
27 observations from traditional sensors. However, the large numbers of channels can lead to  
28 computational burden in satellite data retrieval and assimilation. Furthermore, most of the  
29 channels are highly correlated and the pieces of independent information contained in the  
30 hyperspectral observations are usually much smaller than the number of channels. Principal  
31 Component Analysis (PCA) was used in this research to compress the observational  
32 information content contained in the Atmosphere Infrared Sounder (AIRS) channels to a  
33 few leading principal components (PCs). The corresponding PC scores were then  
34 assimilated into a PCA-based ensemble Kalman filter (EnKF) system. In this proof-of-  
35 concept study based on simulated observations, hyperspectral brightness temperatures  
36 were simulated using the atmospheric state vectors from convection-permitting ensemble  
37 simulations of Hurricane Harvey (2017) as input to the Community Radiative Transfer  
38 Model (CRTM). The PCs were derived from a pre-existing training dataset of brightness  
39 temperatures calculated from convection-permitting simulation over a large domain in the  
40 Indian Ocean representing generic atmospheric conditions over tropical oceans. The EnKF  
41 increments from assimilating many individual measurements in the brightness temperature  
42 space were compared to the EnKF increments from assimilating significantly fewer  
43 numbers of leading PCs. Results showed that assimilating about 10 to 20 leading PCs could  
44 yield increments that were nearly indistinguishable to that from assimilating hyperspectral  
45 measurements from orders of magnitude larger number of hyperspectral channels.

46

47

## 48 **1. Introduction**

49 Hyperspectral infrared radiometers such as the Atmosphere Infrared Sounder (AIRS;  
50 Aumann et al. 2003) on the Earth Observing System (EOS) Aqua platform and the Infrared  
51 Atmospheric Sounding Interferometer (IASI; Hilton et al. 2012) on MetOp series provide  
52 retrieval profiles with better accuracy and higher vertical resolution than traditional  
53 broadband radiometers such as the Advanced Baseline Imager (ABI) onboard GOES-16  
54 (Schmit et al. 2005, 2017) because of their higher spectral resolution and broader spectral  
55 coverage compared to traditional sounders (Huang et al. 1992). Assimilating selected  
56 temperature sounding channels and humidity sounding channels from these hyperspectral  
57 observations provides positive impact on both global models and regional models  
58 (McNally et al. 2006; Collard and McNally 2009; Guidard et al. 2011; Xu et al. 2013).  
59 Although it is tempting to assimilate all the channels to exploit the full information content  
60 in these hyperspectral observations, its high computational cost can be prohibitive,  
61 especially in the near-real time operational environment. Moreover, although thousands of  
62 channels are available on these instruments (e.g., 2378 channels on AIRS and 8461  
63 channels on IASI), most channels are highly correlated, resulting in much fewer pieces of  
64 independent information than the number of channels (Huang et al. 1992).

65 Carefully selecting a subset of all the channels that contains as much independent  
66 information as possible is one way to balance computational cost and information content  
67 obtained (e.g., Goldberg et al. 2003; Susskind et al. 2003). Alternatively, the full  
68 information content (or part of it) of hyperspectral observations can be compressed using  
69 dimension reduction techniques such as the Principal Component Analysis (PCA;

70 Goldberg et al. 2003; Huang and Antonelli 2001; Aires et al. 2002), which is also called  
71 Empirical Orthogonal Function (EOF; Hannachi et al. 2007; Monahan et al. 2009; North  
72 1984). Besides data compression, another benefit from these PCA-based approaches is  
73 noise reduction. It can be argued that variations of atmosphere signals are more correlated  
74 across the hyperspectral channels while variations of random noises are less correlated.  
75 Thus, after PCA, the atmosphere signals are more likely to be represented by the leading  
76 principal components (i.e., those with larger eigenvalues) while noises are more likely to  
77 be represented by the lower-rank principal components (i.e., those with smaller  
78 eigenvalues). A truncated number of leading principal components (PCs), with  
79 significantly fewer variables than the number of channels, can be used to reproduce the  
80 original hyperspectral observations with reduced noise (Huang and Antonelli 2001;  
81 Antonelli et al. 2004; Turner et al. 2006).

82 PCA-compressed data have been used in data assimilation approaches in various ways.  
83 Collard et al. (2010) assimilated reconstructed IASI radiances from the principal  
84 components. Matricardi and McNally (2014) directly assimilated PC scores instead of  
85 radiances over clear-sky condition in the European Centre for Medium-range Weather  
86 Forecasts (ECMWF) 4D-Var data assimilation system. They showed that directly  
87 assimilating 20 PC scores instead of 165 IASI radiances leads to significant computational  
88 savings with no detectable loss of skills.

89 PCA-based radiative transfer forward models have been developed to efficiently  
90 simulate hyperspectral radiances or brightness temperatures (BTs) from atmospheric states  
91 (temperature, water vapor, and trace gas profiles, etc.). Examples are the PCA-based  
92 version of Radiative Transfer for TOVS (PC\_RTTOV, Matricardi, 2010) and the Principal

93 Component-based Radiative Transfer Forward Model (PCRTM) developed by Liu et al.  
94 (2006). These PCA-based radiative transfer models directly simulate PC scores using pre-  
95 trained PCs. If channel radiances or brightness temperatures are needed, they can be  
96 reconstructed using PC scores and the pre-trained PCs provided by the PCA-based forward  
97 models. The overall execution speed is several to tens of times faster than channel-based  
98 fast forward models. If PC scores are directly used in subsequent retrieval or data  
99 assimilation (e.g., Matricardi and McNally 2014), extra time savings can be obtained since  
100 the reconstruction of brightness temperatures or radiances is not necessary.

101 The typical approach of infrared satellite data assimilation used by operational centers  
102 is to only assimilate radiances not affected by clouds (Geer et al. 2018), which includes not  
103 only cloud-free conditions, but also channels unaffected by clouds (e.g., stratospheric  
104 sounding channels in cloudy conditions) For cases where clouds and precipitation are  
105 ubiquitous, all-sky data assimilation shows positive impact on the prediction skill especially  
106 for short-range forecasting (Zhang et al. 2016; Minamide and Zhang 2019, 2018). Wu et  
107 al. (2017) showed that PCRTM could be used in physical retrieval for IASI cloudy scene  
108 analysis, illustrating possible use of PCA-based radiative transfer models in data  
109 assimilation under all-sky conditions. These findings indicate that assimilating  
110 hyperspectral observations using PCA-based data assimilation method under all-sky  
111 conditions is worth exploring, in particular for the ensemble-based data assimilation  
112 techniques.

113 In this paper, we explore the possibility of directly assimilating PC scores instead of  
114 original hyperspectral observations under the Ensemble Kalman Filter (EnKF) framework  
115 for all-sky conditions. The outline of this paper is as follows. Section 2 discusses the

116 methodology including a brief review of PCA and PC scores in satellite data applications,  
117 and the proposed approach in direct assimilation of PC scores with the EnKF. Section 3  
118 describes data compression and dimension reduction of hyperspectral measurements with  
119 PCA. Section 4 shows data assimilation experiments with EnKF update using PC scores  
120 as innovation vector. Section 5 provides summary and discussion.

121

## 122 **2. Methodology**

### 123 *2.1 Brief review of PCA and PC scores in satellite data applications*

124 Principal component analysis is a widely used method for dimension reduction, data  
125 compression, and noise reduction. Its application in data application has been described in  
126 a number of papers with slightly different approaches and notations (e.g., Collard et al.  
127 2010; Matricardi and McNally 2014). Here we briefly summarize PCA in the context of  
128 hyperspectral infrared observations as follows.

129 Consider a training dataset consists of  $n_{obs}$  hyperspectral observation samples  $\mathbf{y}_i (i =$   
130  $1, 2, \dots, n_{obs})$ , each with  $n_{ch}$  channels, where the boldface non-italic font represents a  
131 multidimensional vector. In principle, the hyperspectral observations can be any  
132 multidimensional variables, such as radiances or normalized radiances respect to  
133 instrument noise. For frequencies where the majority of observed radiances are from  
134 thermal emissions of the atmosphere and the Earth surface, such as longwave infrared and  
135 microwave, radiances are often converted to brightness temperatures. Usually, PCA are  
136 trained using a large dataset that represents the range of variations in atmospheric  
137 conditions, including absorber amounts (e.g., water vapor and trace gas), temperature  
138 profiles, and surface parameters (Matricardi, 2010), and it can be assumed that  $n_{obs} > n_{ch}$ .

139 The mean of all hyperspectral observations of the training dataset is  $\overline{\mathbf{y}^{tr}} = \frac{1}{n_{obs}} \sum_{i=1}^{n_{obs}} \mathbf{y}_i$ ,  
 140 which is a vector with dimension  $n_{ch}$ , and each element of  $\overline{\mathbf{y}^{tr}}$  is the mean value of all  $n_{obs}$   
 141 observations for a given channel. The deviations of all hyperspectral observations from the  
 142 mean are arranged into a  $n_{ch} \times n_{obs}$  matrix  $\mathbf{Y} = (\mathbf{y}'_1, \mathbf{y}'_2, \dots, \mathbf{y}'_{n_{obs}})$ , where each column of  
 143  $\mathbf{Y}$  is the deviation of one hyperspectral observation sample from the training dataset mean  
 144  $\overline{\mathbf{y}^{tr}}$ , i.e.,  $\mathbf{y}'_i = \mathbf{y}_i - \overline{\mathbf{y}^{tr}}$ , and each row of  $\mathbf{Y}$  contains  $n_{obs}$  observation deviations at each  
 145 channel. The covariance matrix  $\mathbf{C}$  of matrix  $\mathbf{Y}$  can be written as

$$146 \quad \mathbf{C} = \frac{1}{n_{obs}-1} \mathbf{Y}\mathbf{Y}^T, \quad (1)$$

147 where the superscript T denotes matrix transpose. Through eigenvalue decomposition,  
 148 covariance matrix  $\mathbf{C}$  can be decomposed as

$$149 \quad \mathbf{C} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T, \quad (2)$$

150 where  $\mathbf{\Lambda}$  is a diagonal matrix with all the eigenvalues  $\lambda_i$  of  $\mathbf{C}$  arranged by their magnitude,  
 151 and  $\mathbf{U}$  is an orthogonal matrix consists of corresponding eigenvectors of  $\mathbf{C}$  as column  
 152 vectors. These column vectors of matrix  $\mathbf{U}$  are the *principal components* (PCs) that lay  
 153 along the directions of maximum variance of the dataset. With the assumption that  $n_{obs} >$   
 154  $n_{ch}$ , the dimension of  $\mathbf{U}$  is  $n_{ch} \times n_{ch}$  and the eigenvectors can be expected to be nonzero.

155 The *PC score*  $\mathbf{z}$  of a multidimensional vector  $\mathbf{y}$  that represents one hyperspectral  
 156 observation, either in the training dataset or not, can be obtained by first subtracting the  
 157 mean hyperspectral observation of *the training dataset*  $\overline{\mathbf{y}^{tr}}$  and then projecting the  
 158 deviation  $\mathbf{y} - \overline{\mathbf{y}^{tr}}$  on the directions of the PCs obtained from the training dataset, i.e.,

$$159 \quad \mathbf{z} = \mathbf{U}^T(\mathbf{y} - \overline{\mathbf{y}^{tr}}). \quad (3)$$

160 From Eq. (3), the PCA process can be interpreted as a linear combination of the original

161 hyperspectral channels into a set of synthesized “super channels”, where each PC (each  
162 column of matrix  $\mathbf{U}$ ) contains the weights of all hyperspectral channels for the  
163 corresponding “super channel”. Then the PC scores are the “measurements” at these super  
164 channels. Note that the PCs and mean hyperspectral observation  $\overline{\mathbf{y}^{tr}}$  trained from the  
165 training dataset should always be used together. Consequently, the hyperspectral  
166 observation can be reconstructed from the PC scores and the PCs by

$$167 \quad \mathbf{y} = \mathbf{U}\mathbf{z} + \overline{\mathbf{y}^{tr}}. \quad (4)$$

168 The variance explained by the  $i^{th}$  PC is given by the eigenvalue associated with it,  $\lambda_i$ . The  
169 proportion of variance explained by the  $i^{th}$  PC, or the *variance ratio*, is given by  $\lambda_i / \sum_j \lambda_j$ .  
170 As will be shown later,  $\lambda_i$  decreases rapidly with increasing  $i$ , indicating the majority of  
171 variance in the hyperspectral observations can be expressed using a smaller number of  
172 leading PCs. Using only the first  $n_{pc}$  PCs and corresponding PC scores leads to a lossy  
173 compression of the original hyperspectral observation:

$$174 \quad \hat{\mathbf{y}} = \mathbf{U}_{npc}\mathbf{z}_{npc} + \overline{\mathbf{y}^{tr}}, \quad (5)$$

175 where  $\mathbf{U}_{npc}$  is a  $n_{ch} \times n_{pc}$  matrix with the first  $n_{pc}$  PCs as its columns and  $\mathbf{z}_{npc}$  is a  
176 column vector with the first  $n_{pc}$  PC scores. When more PCs are used, the reconstruction  
177 error, i.e., the difference between the reconstructed and original hyperspectral observation,  
178 is smaller. However, noises in the lower rank PCs will also be included in the reconstructed  
179 hyperspectral observation so that smaller reconstruction error is not necessarily better. The  
180 optimal choice of  $n_{pc}$  may differ from application to application since it depends on the  
181 noise characteristics of the original hyperspectral observations as well as storage and  
182 transmission limitations. For example, Turner et al. (2006) discussed the number of PCs



183 used for noise reduction of the ground-based Atmospheric Emitted Radiance  
 184 Interferometer (AERI).

185

## 186 2.2 Proposed approach in direct assimilation of PC scores with the EnKF

187 The EnKF update equation at the analysis step is adopted from the formulations  
 188 presented in Houtekamer and Zhang (2016, their Eqs. 1 and 2):

$$189 \mathbf{x}^a = \mathbf{x}^f + \mathbf{K}[\mathbf{y}^o - \mathcal{H}\mathbf{x}^f], \quad (6)$$

$$190 \mathbf{K} = \mathbf{P}^f \mathcal{H}^T (\mathcal{H} \mathbf{P}^f \mathcal{H}^T + \mathbf{R})^{-1}, \quad (7)$$

191 where  $\mathbf{x}^a$  is the updated estimate of atmosphere state from the prior estimate of the  
 192 atmosphere state  $\mathbf{x}^f$  using the extra information from the new observation  $\mathbf{y}^o$ ,  $\mathbf{K}$  is the  
 193 Kalman gain matrix,  $\mathcal{H}$  is the forward operator that performs mapping of the model state  
 194 to observation state,  $\mathbf{R}$  is the observation error covariance, and  $\mathbf{P}^f$  is the background error  
 195 covariance. In EnKF, ensemble-based approximations of  $\mathbf{P}^f \mathcal{H}^T$  and  $\mathcal{H} \mathbf{P}^f \mathcal{H}^T$  are used  
 196 (Houtekamer and Zhang 2016, their Eqs. 6 and 7):

$$197 \mathbf{P}^f \mathcal{H}^T \equiv \frac{1}{N_{ens}-1} \sum_{i=1}^{N_{ens}} (\mathbf{x}_i^f - \overline{\mathbf{x}^f}) (\mathcal{H} \mathbf{x}_i^f - \overline{\mathcal{H} \mathbf{x}^f})^T, \quad (8)$$

$$198 \mathcal{H} \mathbf{P}^f \mathcal{H}^T \equiv \frac{1}{N_{ens}-1} \sum_{i=1}^{N_{ens}} (\mathcal{H} \mathbf{x}_i^f - \overline{\mathcal{H} \mathbf{x}^f}) (\mathcal{H} \mathbf{x}_i^f - \overline{\mathcal{H} \mathbf{x}^f})^T, \quad (9)$$

199 where

$$200 \overline{\mathbf{x}^f} = \frac{1}{N_{ens}} \sum_{i=1}^{N_{ens}} \mathbf{x}_i^f \text{ and } \overline{\mathcal{H} \mathbf{x}^f} = \frac{1}{N_{ens}} \sum_{i=1}^{N_{ens}} \mathcal{H} \mathbf{x}_i^f.$$

201 In the case of hyperspectral observations,  $\mathcal{H} \mathbf{x}_i^f$  is the simulated hyperspectral  
 202 observations from the  $i^{th}$  ensemble member, i.e.,  $\mathbf{y}_i^f$ .  $\overline{\mathcal{H} \mathbf{x}^f}$  is the mean simulated  
 203 hyperspectral observation of *the ensemble*. If the PC scores corresponding to  $\mathbf{y}_i^f$  calculated  
 204 using Eq. (3) is  $\mathbf{z}_i^f$ , with the help of Eq. (4), Eqs. (8) and (9) can be written as

205  $\mathbf{P}^f \mathcal{H}^T = \left[ \frac{1}{N_{ens}-1} \sum_{i=1}^{N_{ens}} (\mathbf{x}_i^f - \overline{\mathbf{x}^f}) (\mathbf{z}_i^f - \overline{\mathbf{z}^f})^T \right] \mathbf{U}^T, \quad (10)$

206  $\mathcal{H} \mathbf{P}^f \mathcal{H}^T = \mathbf{U} \left[ \frac{1}{N_{ens}-1} \sum_{i=1}^{N_{ens}} (\mathbf{z}_i^f - \overline{\mathbf{z}^f}) (\mathbf{z}_i^f - \overline{\mathbf{z}^f})^T \right] \mathbf{U}^T. \quad (11)$

207 Comparing the right-hand-side of Eqs. (8) and (10), as well as the Eqs. (9) and (11), if  
 208 a PCA-based forward model  $\mathcal{H}_{pc}$ , which is based on the same eigenvector basis  $\mathbf{U}$  as in  
 209 Eq. (3), is used to perform mapping from the model state to PC scores such that  $\mathbf{z}_i^f =$

210  $\mathcal{H}_{pc} \mathbf{x}_i^f$ , Eqs. (10) and (11) can be written as

211  $\mathbf{P}^f \mathcal{H}^T = \mathbf{P}^f \mathcal{H}_{pc}^T \mathbf{U}^T, \quad (12)$

212  $\mathcal{H} \mathbf{P}^f \mathcal{H}^T = \mathbf{U} \mathcal{H}_{pc} \mathbf{P}^f \mathcal{H}_{pc}^T \mathbf{U}^T. \quad (13)$

213 Since the matrix  $\mathbf{U}$  is an orthogonal matrix, the equation  $\mathbf{U}^T \mathbf{U} = \mathbf{U} \mathbf{U}^T = \mathbf{I}$  holds, where  
 214  $\mathbf{I}$  is identity matrix. The Kalman gain matrix can be written as

215  $\mathbf{K} = \mathbf{P}^f \mathcal{H}_{pc}^T \mathbf{U}^T (\mathbf{U} \mathcal{H}_{pc} \mathbf{P}^f \mathcal{H}_{pc}^T \mathbf{U}^T + \mathbf{R})^{-1} = \mathbf{P}^f \mathcal{H}_{pc}^T \mathbf{U}^T [\mathbf{U} (\mathcal{H}_{pc} \mathbf{P}^f \mathcal{H}_{pc}^T + \mathbf{U}^T \mathbf{R} \mathbf{U}) \mathbf{U}^T]^{-1}$   
 216  $= \mathbf{P}^f \mathcal{H}_{pc}^T (\mathcal{H}_{pc} \mathbf{P}^f \mathcal{H}_{pc}^T + \mathbf{U}^T \mathbf{R} \mathbf{U})^{-1} \mathbf{U}^T, \quad (14)$

217 and Eq. (6) can be written as

218  $\mathbf{x}^a = \mathbf{x}^f + \mathbf{K} \mathbf{U} [\mathbf{z}^o - \mathbf{z}^f]. \quad (15)$

219 Substituting Eq. (14) into Eq. (15) yields

220  $\mathbf{x}^a = \mathbf{x}^f + \mathbf{K}_{pc} [\mathbf{z}^o - \mathbf{z}^f], \quad (16)$

221  $\mathbf{K}_{pc} = \mathbf{P}^f \mathcal{H}_{pc}^T (\mathcal{H}_{pc} \mathbf{P}^f \mathcal{H}_{pc}^T + \mathbf{U}^T \mathbf{R} \mathbf{U})^{-1}. \quad (17)$

222 Comparing the right-hand-side of Eq. (10) and Eq. (12), as well as Eq. (11) and Eq.  
 223 (13), and using  $\mathbf{z}_i^f = \mathcal{H}_{pc} \mathbf{x}_i^f$ ,  $\mathbf{P}^f \mathcal{H}_{pc}^T$  and  $\mathcal{H}_{pc} \mathbf{P}^f \mathcal{H}_{pc}^T$  can also be calculated from the  
 224 ensemble

225  $\mathbf{P}^f \mathcal{H}_{pc}^T = \left[ \frac{1}{N_{ens}-1} \sum_{i=1}^{N_{ens}} (\mathbf{x}_i^f - \overline{\mathbf{x}^f}) (\mathcal{H}_{pc} \mathbf{x}_i^f - \overline{\mathcal{H}_{pc} \mathbf{x}^f})^T \right], \quad (18)$

226 
$$\mathcal{H}_{pc} \mathbf{P}^f \mathcal{H}_{pc}^T = \left[ \frac{1}{N_{ens}-1} \sum_{i=1}^{N_{ens}} (\mathcal{H}_{pc} \mathbf{x}_i^f - \overline{\mathcal{H}_{pc} \mathbf{x}^f}) (\mathcal{H}_{pc} \mathbf{x}_i^f - \overline{\mathcal{H}_{pc} \mathbf{x}^f})^T \right]. \quad (19)$$

227 Equation (16-19) are very similar to the original EnKF update Eqs. (6-9), with the  
 228 observation variable changed from hyperspectral observations to PC scores, the forward  
 229 operator changed to a PCA-based radiative transfer model that directly calculate PC scores,  
 230 and the observation error covariance changed to  $\mathbf{R}_{pc} = \mathbf{U}^T \mathbf{R} \mathbf{U}$ . This indicates assimilating  
 231 all the PC scores using Eqs. (16-19) and assimilating hyperspectral observations using Eqs.  
 232 (6-9) should give the same analysis increment  $\mathbf{x}^a - \mathbf{x}^f$ .

233 Goldberg et al. (2003) chose to save and distribute 200 PC scores to assure that there  
 234 was sufficient information to reconstruct observed AIRS observations with over 2000  
 235 channels. As will be shown later in this paper, the variance of hyperspectral observations  
 236 explained by the leading PCs are orders of magnitude larger than the lower-rank PCs. This  
 237 means the leading PC scores (leading dimensions of  $\mathbf{z}$ ) are expected to be much larger in  
 238 value than the lower-rank PC scores (lower-rank dimensions of  $\mathbf{z}$ ). The matrix elements  
 239 corresponding to the lower-rank PC scores of the variance term in the Kalman gain matrix,  
 240 i.e.,  $\mathcal{H}_{pc} \mathbf{P}^f \mathcal{H}_{pc}^T + \mathbf{U}^T \mathbf{R} \mathbf{U}$ , are expected to be dominated by instrument noise. As a result,  
 241 assimilating the lower-rank PC scores have much less contribution to the EnKF increment  
 242 compared to the leading PC scores. Thus, assimilating only the leading PC scores can  
 243 provide similar result to that when all the channels of the hyperspectral observations are  
 244 assimilated.

245 The most obvious benefit from this PCA-based EnKF method that assimilates truncated  
 246 PC scores is the huge savings in computation time. The speedup includes both time savings  
 247 in the forward modeling step and the EnKF analysis step. AIRS have over 2000 channels.

248 If assimilating less than 20 leading PCs can provide similar result compared to assimilating  
249 over 2000 channels, tens of times speedup can be expected.

250

### 251 3. Data compression and dimension reduction of hyperspectral satellite 252 measurements with PCA

253 The focus of this work is to demonstrate that directly assimilating a number of leading  
254 PC scores yields similar increments to the model states compared to assimilating a much  
255 larger number of hyperspectral channels. In other words, it is to show that the truncated  
256 version of Eqs. 16 and 17 yields similar increments to Eqs. 6 and 7. As such, full data  
257 assimilation cycle was not performed. Instead, EnKF priors from previous data assimilation  
258 experiments for Hurricane Harvey (2017) described in Chapter 6 of Minamide (2018) was  
259 used, which is described in Section 3.1.

260 When directly assimilating PC scores in an EnKF system, there are generally two  
261 approaches to map atmosphere states into PC scores. One way is to simulate the  
262 hyperspectral observations using traditional channel-based radiative transfer models, then  
263 project the simulated hyperspectral observations to pre-trained PC components. By  
264 following this approach, the data assimilation system benefits from the computational  
265 savings in the analysis step due to reduction in dimensionality and from noise reduction by  
266 discarding the noise-affected lower-rank PCs. The drawback of this approach is that the  
267 computational cost can be prohibitive if large number of channels are need, especially for  
268 operational systems. Another way to obtain PC scores is to directly map from atmospheric  
269 states to PC scores by using a PCA-based radiative transfer model as forward operator  $\mathcal{H}_{pc}$   
270 such as PCRTM and PC\_RTTOV. Additional computational savings in radiative transfer

271 calculations may be achieved with the latter approach. Since simulating hyperspectral  
272 observations was necessary in this research (i.e., used in Eqs. 6 and 7), we chose to follow  
273 the first approach. Another benefit of this approach is that the result can be independent of  
274 the choice of PCA-based radiative transfer model and its underlying assumptions and  
275 constraints.

276 AIRS brightness temperatures (BTs) were simulated using the Community Radiative  
277 Transfer Model (CRTM; Han et al. 2006) using modeled atmospheric and surface states as  
278 input. Principal components should be trained using a dataset that is independent of the  
279 EnKF assimilation experiment. In this work, simulated AIRS BTs from the simulation of  
280 Dynamics of the MJO (DYNAMO) field campaign described in Ying and Zhang (2017) is  
281 used to train the PCA, which is also described in details in Section 3.1; the selection of  
282 AIRS channels used in this study is described in Section 3.2; and the training of PCA and  
283 calculation of PC scores is described in Section 3.3.

284

### 285 3.1. *Generating hyperspectral observations*

286 EnKF inputs from previous data assimilation experiments for Hurricane Harvey (2017)  
287 described in Chapter 6 of Minamide (2018) was used to evaluate the PCA-based EnKF  
288 data assimilation method. The ensemble data assimilation experiment had 60 ensemble  
289 members, which were initialized at 0000 UTC 22 Aug 2017 with perturbations added using  
290 WRFDA CV3 option and integrated for 12 hours for spin-up. Minimum sea level pressure  
291 and all-sky GOES-16 channel 8 observations were assimilated every hour after 1200 UTC  
292 22 Aug 2017 for three days till 1200 UTC 25 Aug 2017. The intermediate resolution (9  
293 km) domain was used because it directly predicted hydrometer mixing ratios and had a

294 relatively large geographical coverage. The WRF single-moment 6-class mixed-phase  
295 microphysics scheme was used (WSM6; Hong and Lim 2006). EnKF priors at two different  
296 analysis times were used: one at 1400 UTC 22 August 2017 and the other one at 0200 UTC  
297 25 August 2017 representing the initial stage and the well-developed stage of Harvey (2017)  
298 respectively. The WRF model used had 42 vertical levels with model top at 10 hPa. More  
299 details about model configuration was described in Chapter 6 of Minamide (2018).

300 Principal components were trained using simulated hyperspectral observations from  
301 the DYNAMO period, which represents general atmospheric conditions over tropical  
302 oceans with different types of cloudy and clear conditions, different moisture profiles, and  
303 a wide range of surface temperatures. The WRF simulation for DYNAMO had a large  
304 domain covering 20°N to 20°S and 50°E to 120°E with 9 km horizontal resolution and 44  
305 vertical levels with model top at 20 hPa. Hydrometer properties were directly predicted by  
306 the WRF double-moment (WDM) scheme (Lim and Hong 2010). More detailed  
307 information in model configuration was described by Ying and Zhang (2017).  
308 Hyperspectral BTs simulated using CRTM from the WRF model output at 0000 UTC 21  
309 October 2011 were calculated. Then Gaussian random noise was added to each channel  
310 with standard deviation equal to instrument noise level of that channel. This noise-added  
311 hyperspectral BTs were used to train PCA.

312 When using modeled atmospheric states as inputs to CRTM, typical effective radius  
313 values of hydrometers were used: 16.8  $\mu\text{m}$  for cloud water drops, 25  $\mu\text{m}$  for cloud ice, 1000  
314  $\mu\text{m}$  for rain, 500  $\mu\text{m}$  for snow, and 500  $\mu\text{m}$  for graupel. Only ocean grid points were used  
315 in this study to avoid additional uncertainties in land surface temperature and emissivity  
316 (that will be explored in future studies). The zenith angles for AIRS cross-track footprints

317 range from near  $0^\circ$  (footprint close to nadir) to about  $57^\circ$  (outermost footprint) in each scan  
318 line. This angle is important for AIRS TB calculations since it influences ocean surface  
319 emissivity and effective path length of the atmosphere. To represent this variability, the  
320 zenith angles were randomly drawn from a uniform distribution from  $0^\circ$  to  $57^\circ$ . For the 60  
321 ensemble members for Hurricane Harvey (2017), the same zenith angles were used at the  
322 same location for all ensemble members, while different zenith angles were used at  
323 different locations.

324 Figure 1 shows simulated BTs of AIRS channel with wavelength  $12.183\ \mu\text{m}$  within the  
325 atmosphere window. Panel (a) shows simulated BTs of the DYNAMO case; panel (b)  
326 shows the 1<sup>st</sup> ensemble member of the early stage Harvey (at 1400 UTC 22 August 2017);  
327 and panel (c) shows the 1<sup>st</sup> ensemble member of the developed stage Harvey (at 0200 UTC  
328 25 August 2017). It is clearly shown that the largest variability in window channel BTs is  
329 due to variations in clouds. Warmer BTs usually correspond to clear sky while colder BTs  
330 usually correspond to clouds of different levels. The three circles in panel (b) labeled with  
331 B1 through B3 are three locations with cold, medium, and warm BTs for the 1<sup>st</sup> ensemble  
332 member of the early stage Harvey case respectively, indicating high-cloud, mid-cloud, and  
333 clear-sky conditions. Note that other ensemble members at these locations could have  
334 different sky conditions than the 1<sup>st</sup> ensemble member. Figure 2 shows histograms of  
335 simulated BTs of AIRS channel with wavelength  $12.183\ \mu\text{m}$  for all the 60 ensemble  
336 members at the three locations B1 through B3. At location B1, most ensemble members  
337 had cold BTs smaller than 240 K, indicating most members had high cloud. At location B2,  
338 most ensemble members had warm BTs (larger than 270 K) while some members had

339 medium BTs (between 240K and 270K). Location B3 is similar to B2, but with more warm  
340 BTs and less medium BTs.

341

### 342 3.2. Channel selection

343 Although using all the channels can exploit the full information content in hyper-  
344 spectral observations, not all channels were used in this demonstrative work. Firstly,  
345 channels with noise equivalent temperature difference (NEDT) larger than 1K were  
346 excluded. Secondly, the model tops of the WRF simulations used in this study only  
347 extended to 10 hPa or 20 hPa. Simulated BTs for channels sensitive to atmosphere layers  
348 close to or above the model tops are not accurate and thus should be excluded, which  
349 includes some temperature and water vapor sounding channels. As a result, only 1670  
350 channels with wavelengths between 4.44  $\mu\text{m}$  to 14.503  $\mu\text{m}$ , with similar wavelength range  
351 to Lin et al. (2017), were used in the current study.

352 Another choice to make is that whether all channels within this wavelength range  
353 should be used in data assimilation. It is known that the number of pieces of independent  
354 information in hyperspectral observations is much smaller than the number of channels.  
355 Figure 3 shows the inter-channel correlations between different channels. Channels with  
356 wavelength 8~13  $\mu\text{m}$  (roughly corresponding to channel number 400~1330) directly sense  
357 the earth surface or cloud top and are highly correlated with each other. They have largest  
358 variations in BTs among various sky conditions because of large differences between cloud  
359 top temperatures and earth surface temperatures. The water vapor channels with  
360 wavelengths about 6~8  $\mu\text{m}$  (channel number 1300~1864) and the temperature sounding  
361 channels with wavelengths about 13~14.5  $\mu\text{m}$  (channel number 137~400) close to  $\text{CO}_2$



362 absorption bands have smaller BTs than the window channels but are also highly correlated  
363 with the window channels. One reason is that when the cloud top is higher than the peak  
364 of the weighting functions of these sounding channels, radiances emitted by the cloud also  
365 have large contribution to the measured radiances by these sounding channels. Another  
366 reason is that correlations exist between atmospheric state variables, such as surface  
367 temperature and atmospheric temperature. As such, correlations between BTs of window  
368 channels and sounding channels also exist over clear-sky condition.

369 In this research, classical fast forward radiative transfer model is used. Simulating BTs  
370 for all AIRS channels and then calculating PC scores require significant computational  
371 resources. Since BTs of AIRS channels were highly correlated, the 324-channel subset  
372 suggested by Goldberg et al. (2003) and Susskind et al. (2003) may be used to balance the  
373 information content and computational requirement, which has 199 of the selected  
374 channels within the frequency range specified in this work. We performed PCA on both  
375 the full-spectral-resolution BTs (1670 channels) and the subset BTs (199 channels) over  
376 the DYNAMO domain. Figure 4(a) shows the variance explained by each PC for both the  
377 full-spectral-resolution BTs (red line) and the subset BTs (blue line). For both full-spectral-  
378 resolution and selected-channel BTs, variance explained by each PC decreases rapidly,  
379 especially for the first few PCs. If  $1 K^2$  was used as variance threshold, corresponding to 1  
380  $K$  white noise in the observation, the full-spectral-resolution BTs had 19 principal  
381 components with variance larger than the threshold, while the subset BTs had 14 PCs.

382 Figure 4(b) shows that the explained variance ratio of both the full-spectral-resolution  
383 BTs (red) and the subset BTs (blue) corresponding to each PC. The full-spectral-resolution  
384 BTs had less PCs than the subset BTs with explained variance ratio above a certain level,

385 indicating variance in the full-spectral-resolution BTs were more strongly concentrated to  
386 the first few leading PCs, possibly because that many highly correlated window channels  
387 are excluded in the 324-channel subset. Figure 4(c) shows the first PCs of the full-spectral  
388 resolution BTs (red) and subset BTs (blue), while Figure 4(d) shows the second PCs. Note  
389 that the amplitude difference was because the PCs were scaled to have unit vector length.  
390 Although the number of channels were different, the shapes of the corresponding PCs were  
391 similar for the full-spectral-resolution and subset BTs. As such, the subset BTs is used in  
392 this study instead of full-spectral-resolution BTs to save computational time in the forward  
393 radiative transfer calculations while retain most information in the hyperspectral  
394 observations.

395

### 396 3.3. *The training of PCA and calculation of PC scores*

397 When selecting AIRS channels used in this study, noisy channels with NEDT values  
398 larger than 1 K were already excluded. Most of the 199 selected channels have similar  
399 noise level lower than 0.5 K. As such, unlike other research that performs PCA on noise  
400 normalized radiances to avoid errors from noisy channels, we performed PCA on BTs  
401 directly. The simulated hyperspectral BTs of the DYNAMO case were used as training  
402 dataset to train the PCA. The domain of DYNAMO case had 444 grids along latitudinal  
403 direction and 777 grids along longitudinal direction. Hyperspectral BTs from every fourth  
404 ocean grid point were used to train the PCA, resulting in a training dataset with 74,612  
405 hyperspectral samples, each with BTs of the 199 selected channels. The choice of the  
406 training dataset in the current demonstrative work was not intended to be optimal. With all  
407 hyperspectral BTs simulated from a single time at 0000 UTC 21 October 2011, they could

408 not represent seasonal or diurnal variations of the atmospheric states. However, this  
409 deficiency could be compensated to some extent by the large geographical coverage of the  
410 domain.

411 PC scores of the hyperspectral BTs simulated from the 60-member ensemble WRF run  
412 for Hurricane Harvey (2017) were calculated using the mean hyperspectral BTs and PCs  
413 trained from the training dataset following Eq. 3, where  $\mathbf{y}$  was the hyperspectral BTs  
414 simulated from the two analysis times of the Harvey case mentioned in Section 3.1, and  
415  $\overline{\mathbf{y}^{tr}}$  was the mean hyperspectral BTs of *the training dataset* from DYNAMO. Then these  
416 differences in hyperspectral BTs were projected onto the PCs trained from *the training*  
417 *dataset* to calculate PC scores.

418 Using more PCs leads to higher reconstruction accuracy but higher computational cost.  
419 It is beyond the scope of this work to find the optimal number of PCs to use. Rather, we  
420 explore the influence of PC truncation on reconstruction error and EnKF increments.  
421 Figure 5 shows the reconstruction error when the first 10 (panels a-c), 15 (panels d-f), or  
422 20 (panels g-i) leading PCs were used to reconstruct the hyperspectral BTs simulated from  
423 the 1<sup>st</sup> member of Harvey (2017) at 0200 UTC 25 August 2017, i.e., the well-developed  
424 stage. The three panels (a, d, g) in the first column show the residual error of a water vapor  
425 sounding channel with wavelength 7.493  $\mu\text{m}$ , panels (b, e, h) in the middle column show  
426 the residual error of a window channel with wavelength 10.213  $\mu\text{m}$ , and panels (c, f, i) in  
427 the right column show the residual error of a temperature sounding channel with  
428 wavelength 13.796  $\mu\text{m}$ . Reconstruction error values for 1400 UTC 22 August 2017 was  
429 similar and not shown.

430 Root mean square of reconstruction error over the entire domain of the DYNAMO case  
431 and the 1<sup>st</sup> member of Harvey (2017) case at the two aforementioned analysis times for all  
432 the selected channels are shown in Figure 6, with the three panels corresponding to using  
433 10, 15, and 20 PCs respectively. When more PCs were used, the reconstruction errors were  
434 smaller. The error level difference between 15 PCs and 20 PCs was small, indicating  
435 information gained from using even more PCs can be limited.

436 It is shown in Figure 1(c) that the range of variation of the window channel was about  
437 80 K, with BTs about 220 K over hurricane eyewalls and about 300 K over clear-sky ocean.  
438 This large variation was successfully captured by the 20 leading PCs with residuals similar  
439 to or below instrument noise levels. Residuals in water vapor sounding channels and  
440 temperature sounding channels showed similar characteristics, indicating the majority of  
441 information content in the hyperspectral BTs could be captured using about 20 leading PC  
442 components.

443

#### 444 **4. Data assimilation experiments with the EnKF update using PC scores as the** 445 **innovation vector**

446 The PCA-based EnKF assimilation method was evaluated by comparing the EnKF  
447 increment  $\mathbf{x}^a - \mathbf{x}^f$  obtained from assimilating PC scores and from directly assimilating  
448 simulated AIRS hyperspectral BTs. The 40<sup>th</sup> ensemble member of the two analysis times  
449 of the Harvey case was used as ‘truth’ because its temperature profiles and water vapor  
450 profiles had relatively large departure from the ensemble mean at the three locations B1  
451 through B3 in Figure 1(b). Gaussian noises with standard deviations equal to instrument  
452 noise levels for all channels were added to the simulated hyperspectral BTs from the 40<sup>th</sup>

453 ensemble member to generate synthetic hyperspectral observations. The other 59 members  
454 were used as ensemble members of the EnKF system.

455 Figure 7(a-c) shows the simulated hyperspectral BTs from all the ensemble members  
456 at the locations labeled as B1, B2, and B3 in Figure 1(b) respectively. Each colored thin  
457 line shows simulated hyperspectral BTs from one ensemble member. The color of each  
458 line was determined by the brightness temperature at the same window channel used in  
459 Figure 1 from each ensemble member at the specific location. As a result, lines with the  
460 same color in Figure 7 (a-c) do not necessarily correspond to the same ensemble member.  
461 Although the three locations were selected as example of high-cloud, mid-cloud, and low-  
462 cloud in the 1<sup>st</sup> member respectively, different ensemble members could have different  
463 scene types, as is shown in Figure 2. Nevertheless, most ensemble members at B1 had  
464 relatively higher clouds, while more members at B2 and B3 had clear sky conditions. The  
465 black lines show the mean hyperspectral BTs of the *training dataset*. It should not be  
466 confused with the BTs calculated using ensemble mean atmospheric states ( $\mathcal{H}\mathbf{x}^f$ ; black  
467 dashed lines, used in EnKF), or the mean hyperspectral BTs of the *ensemble* (not shown),  
468 since it is the mean hyperspectral BTs of the *training dataset* that should be subtracted  
469 from the hyperspectral BTs from the ensembles when calculating their PC scores.

470 Figure 7(d-f) show the corresponding 20 leading PC scores of the ensemble members  
471 at the three locations respectively. Lines with the same colors in Figure 7(a & d), Figure  
472 7(b & e), and Figure 7(c & f) correspond to the same ensemble members respectively. The  
473 black dashed line in Figure 7 (d-f) show the 20 leading PC scores of the ensemble mean.  
474 The absolute values of PC scores decreased rapidly with increasing PC index, which agreed  
475 well with the rapid decrease in explained variance shown in Figure 4(a).

476 When AIRS hyperspectral BTs were assimilated, Eqs. (6)-(9) were used to calculate  
477 EnKF increments. Although instrument noise levels for most AIRS channels are smaller  
478 than 0.5K, observation error for all the channels is assumed as 1K in this experiment to  
479 account for other possible error sources including representative error and forward  
480 radiative transfer model error. As such, the observation error covariance  $\mathbf{R}$  was a diagonal  
481 matrix with all its diagonal elements equaling to  $1 \text{ K}^2$ .

482 When PC scores were assimilated, EnKF increments were calculated using a truncated  
483 version of Eqs. (16-19) that only a number of leading PCs and PC scores were used. The  
484 observation error covariance changed to  $\mathbf{R}_{pc} = \mathbf{U}_{npc}^T \mathbf{R} \mathbf{U}_{npc}$ , where  $\mathbf{U}_{npc}$  was a truncated  
485 version of matrix  $\mathbf{U}$  as in Eq. (5). Since  $\mathbf{U}_{npc}$  was a  $n_{ch} \times n_{pc}$  matrix, the size of matrix  
486  $\mathbf{R}_{pc}$  was  $n_{pc} \times n_{pc}$ . When  $\mathbf{R}$  was a diagonal matrix with all its diagonal elements equaling  
487 to  $1 \text{ K}^2$  as assumed above,  $\mathbf{R}_{pc}$  was also a diagonal matrix with all its diagonal elements  
488 equaling to  $1 \text{ K}^2$  but with a much smaller size  $n_{pc} \times n_{pc}$ .

489 Figure 8 shows the EnKF increments of temperature profiles (panels a-c) and of water  
490 vapor profiles (panels d-f) when AIRS hyperspectral BTs were assimilated (black lines)  
491 and when PC scores were assimilated (colored lines) at 1400 UTC 22 August 2017.  
492 Different choices of  $n_{pc}$  values from 2 PCs to 30 PCs were tested. With increasing number  
493 of PC scores assimilated, the increments from assimilating PC scores became closer to that  
494 from assimilating AIRS hyperspectral BTs. When more than 15 PC scores were assimilated,  
495 the increments from assimilating PC scores and from assimilating AIRS hyperspectral BTs  
496 were almost indistinguishable for both the temperature profiles and water vapor profiles.

497 Figure 9 shows the root-mean-square (RMS) of the difference between increments from  
498 assimilating AIRS hyperspectral BTs and PC scores calculated over all layers at the three

499 locations when different number of PCs were assimilated at 1400 UTC 22 August 2017 at  
500 the three locations B1 through B3. Panel (a) shows the RMS of temperature increments  
501 and panel (b) shows the RMS of water vapor increments. For all of the three locations,  
502 RMS decreased until about 16 leading PCs were used. This number was close to that used  
503 by Matricardi and McNally (2014) where 20 leading PCs were used for sounding channels  
504 under clear-sky condition. This result indicated that most information contents in the  
505 hyperspectral BTs could be captured by a smaller number of leading PCs. Also,  
506 assimilating the leading PC scores instead of assimilating AIRS hyperspectral BTs could  
507 provide significant computational savings with satisfactory accuracy.

508

## 509 **5. Summary and discussion**

510 Satellite-based hyperspectral observations such as those from AIRS and IASI have  
511 thousands of infrared channels that contain information on atmospheric state with much  
512 higher vertical resolution compared to observations from traditional sensors. However, the  
513 large numbers of channels also lead to computational burden in retrieval and data  
514 assimilation. Furthermore, most of the channels are highly correlated and the number of  
515 pieces of independent information contained in the hyperspectral observations are usually  
516 much smaller than the number of channels. Principal component analysis (PCA) was used  
517 in this research to compress the observational information content contained in these  
518 hyperspectral channels to a few leading principal components (PC). The corresponding PC  
519 scores can then be assimilated into a PCA-based ensemble Kalman filter (EnKF) system.

520 In this proof-of-concept study using simulated observations, PCA was trained from  
521 AIRS hyperspectral BTs simulated using the Community Radiative Transfer Model

522 (CRTM) and a large-domain convection-permitting simulation over the Indian Ocean that  
523 represents generic tropical ocean conditions. Brightness temperatures of 1670 AIRS  
524 channels were simulated over the domain, which showed large inter-channel correlations.  
525 Since PCA-based fast radiative transfer model is not used in this study, a subset of 199  
526 channels were selected for subsequent analysis, which contains most of the variabilities in  
527 all AIRS channels, to balance information content and computational requirement. Then  
528 principal components were trained using 74,612 hyperspectral BTs samples.

529 AIRS hyperspectral BTs were simulated from the convection permitting ensemble  
530 simulations of Hurricane Harvey (2017) with CRTM. These hyperspectral BTs were  
531 converted to PC scores using the mean hyperspectral BTs of the training dataset and the  
532 PCs. The EnKF increments from assimilating AIRS hyperspectral BTs and from  
533 assimilating different numbers of leading PCs were compared. Result showed that  
534 assimilating about 10 to 20 leading PCs could yield increments that were nearly  
535 indistinguishable to that from assimilating hyperspectral measurements with 199 channels.

536 In this proof-of-concept study, we chose not to use PCA-based radiative transfer  
537 models so that the results shown in this work can be independent of their underlying  
538 assumptions and constraints. The drawback of this approach is that we had to use a subset  
539 of all AIRS channels which contains most of the information content. In a real PCA-based  
540 EnKF system, PCA-based radiative transfer model can be used to directly simulate PC  
541 scores at the full-spectral-resolution.

542 The current proof-of-concept study is based on simulated observations for both the  
543 training dataset and the reference truth for a tropical cyclone event; research is ongoing  
544 and/or planned to further explore the use of this approach for cycled real-data assimilation



545 under different atmospheric conditions including those over various land surfaces and/or  
546 over higher latitudes.

547

548 **Acknowledgements:** This research is partially supported by NASA Grants  
549 NNX16AD84G and NNX12AJ79G, ONR Grant N000140910526 and NOAA funding  
550 under HFIP and NGGPS. The authors would like to thank Masashi Minamide, Scott Siron  
551 and Yue Ying for providing the WRF model output. Discussions with Eugene Clothiaux,  
552 Scott Siron and Xianglei Huang are beneficial for this study. Computing are performed at  
553 the Texas Advanced Computing Center.

554

555 **6. References**

556 Aires, F., W. B. Rossow, N. A. Scott, and A. Chédin, 2002: Remote sensing from the  
557 infrared atmospheric sounding interferometer instrument 1. Compression, denoising,  
558 and first-guess retrieval algorithms. *J. Geophys. Res.*, **107**, 4619,  
559 doi:10.1029/2001JD000955.

560 Antonelli, P., and Coauthors, 2004: A principal component noise filter for high spectral  
561 resolution infrared measurements. *J. Geophys. Res. Atmos.*, **109**, D23102,  
562 doi:10.1029/2004JD004862.

563 Aumann, H. H., and Coauthors, 2003: AIRS/AMSU/HSB on the aqua mission: design,  
564 science objectives, data products, and processing systems. *IEEE Trans. Geosci.*  
565 *Remote Sens.*, **41**, 253–264, doi:10.1109/TGRS.2002.808356.

566 Collard, A. D., and A. P. McNally, 2009: The assimilation of Infrared Atmospheric  
567 Sounding Interferometer radiances at ECMWF. *Q. J. R. Meteorol. Soc.*, **135**, 1044–  
568 1058, doi:10.1002/qj.410.

569 ———, ———, F. I. Hilton, S. B. Healy, and N. C. Atkinson, 2010: The use of principal  
570 component analysis for the assimilation of high-resolution infrared sounder  
571 observations for numerical weather prediction. *Q. J. R. Meteorol. Soc.*, **136**, 2038–  
572 2050, doi:10.1002/qj.701.

573 Geer, A. J., and Coauthors, 2018: All-sky satellite data assimilation at operational weather  
574 forecasting centres. *Q. J. R. Meteorol. Soc.*, **144**, 1191–1217, doi:10.1002/qj.3202.

575 Goldberg, M. D., Yanni Qu, L. M. McMillin, W. Wolf, Lihang Zhou, and M. Divakarla,  
576 2003: AIRS near-real-time products and algorithms in support of operational  
577 numerical weather prediction. *IEEE Trans. Geosci. Remote Sens.*, **41**, 379–389,

578 doi:10.1109/TGRS.2002.808307.

579 Guidard, V., N. Fourrié, P. Brousseau, and F. Rabier, 2011: Impact of IASI assimilation at  
580 global and convective scales and challenges for the assimilation of cloudy scenes. *Q. J. R. Meteorol. Soc.*, **137**, 1975–1987, doi:10.1002/qj.928.

581

582 Han, Y., P. van Delst, Q. Liu, F. Weng, B. Yan, R. Treadon, and J. Derber, 2006: JCSDA  
583 Community Radiative Transfer Model ( CRTM ) - Version 1. *NOAA Tech. Rep.*,  
584 **NESDIS 122**, 40.

585 Hannachi, A., I. T. Jolliffe, and D. B. Stephenson, 2007: Empirical orthogonal functions  
586 and related techniques in atmospheric science: A review. *Int. J. Climatol.*, **27**, 1119–  
587 1152, doi:10.1002/joc.1499.

588 Hilton, F., and Coauthors, 2012: Hyperspectral Earth Observation from IASI: Five Years  
589 of Accomplishments. *Bull. Am. Meteorol. Soc.*, **93**, 347–370, doi:10.1175/BAMS-D-  
590 11-00027.1.

591 Hong, S.-Y., and J.-O. J. Lim, 2006: The WRF Single-Moment 6-Class Microphysics  
592 Scheme (WSM6). *J. Korean Meteorol. Soc.*, **42**, 129–151.

593 Houtekamer, P. L., and F. Zhang, 2016: Review of the Ensemble Kalman Filter for  
594 Atmospheric Data Assimilation. *Mon. Weather Rev.*, **144**, 4489–4532,  
595 doi:10.1175/MWR-D-15-0440.1.

596 Huang, H.-L., and P. Antonelli, 2001: Application of Principal Component Analysis to  
597 High-Resolution Infrared Measurement Compression and Retrieval. *J. Appl. Meteorol.*,  
598 **40**, 365–388, doi:10.1175/1520-  
599 0450(2001)040<0365:AOPCAT>2.0.CO;2.

600 —, W. L. Smith, and H. M. Woolf, 1992: Vertical Resolution and Accuracy of

601 Atmospheric Infrared Sounding Spectrometers. *J. Appl. Meteorol.*, **31**, 265–274,  
602 doi:10.1175/1520-0450(1992)031<0265:VRAAOA>2.0.CO;2.

603 Lim, K.-S. S., and S.-Y. Hong, 2010: Development of an Effective Double-Moment Cloud  
604 Microphysics Scheme with Prognostic Cloud Condensation Nuclei (CCN) for  
605 Weather and Climate Models. *Mon. Weather Rev.*, **138**, 1587–1612,  
606 doi:10.1175/2009MWR2968.1.

607 Lin, H., S. S. Weygandt, A. H. N. Lim, M. Hu, J. M. Brown, and S. G. Benjamin, 2017:  
608 Radiance Preprocessing for Assimilation in the Hourly Updating Rapid Refresh  
609 Mesoscale Model: A Study Using AIRS Data. *Weather Forecast.*, **32**, 1781–1800,  
610 doi:10.1175/WAF-D-17-0028.1.

611 Liu, X., W. L. Smith, D. K. Zhou, and A. Larar, 2006: Principal component-based radiative  
612 transfer model for hyperspectral sensors: theoretical concept. *Appl. Opt.*, **45**, 201,  
613 doi:10.1364/AO.45.000201.

614 Matricardi, M., 2010: A principal component based version of the RTTOV fast radiative  
615 transfer model. *Q. J. R. Meteorol. Soc.*, **136**, 1823–1835, doi:10.1002/qj.680.

616 Matricardi, M., and A. P. McNally, 2014: The direct assimilation of principal components  
617 of IASI spectra in the ECMWF 4D-Var. *Q. J. R. Meteorol. Soc.*, **140**, 573–582,  
618 doi:10.1002/qj.2156.

619 McNally, A. P., P. D. Watts, J. A. Smith, R. Engelen, G. A. Kelly, J. N. Thépaut, and M.  
620 Matricardi, 2006: The assimilation of AIRS radiance data at ECMWF. *Q. J. R.  
621 Meteorol. Soc.*, **132**, 935–957, doi:10.1256/qj.04.171.

622 Minamide, M., 2018: On the predictability of tropical cyclones through all-sky infrared  
623 satellite radiance assimilation. The Pennsylvania State University, 201 pp.

624 ———, and F. Zhang, 2018: Assimilation of All-Sky Infrared Radiances from Himawari-8  
625 and Impacts of Moisture and Hydrometer Initialization on Convection-Permitting  
626 Tropical Cyclone Prediction. *Mon. Weather Rev.*, **146**, 3241–3258,  
627 doi:10.1175/MWR-D-17-0367.1.

628 ———, and ———, 2019: An adaptive background error inflation method for assimilating all-  
629 sky radiances. *Q. J. R. Meteorol. Soc.*, **145**, 805–823, doi:10.1002/qj.3466.

630 Monahan, A. H., J. C. Fyfe, M. H. P. Ambaum, D. B. Stephenson, and G. R. North, 2009:  
631 Empirical orthogonal functions: The medium is the message. *J. Clim.*, **22**, 6501–6514,  
632 doi:10.1175/2009JCLI3062.1.

633 North, G. R., 1984: Empirical Orthogonal Functions and Normal Modes. *J. Atmos. Sci.*, **41**,  
634 879–887, doi:10.1175/1520-0469(1984)041<0879:EOFANM>2.0.CO;2.

635 Schmit, T. J., M. M. Gunshor, W. P. Menzel, J. J. Gurka, J. Li, and A. S. Bachmeier, 2005:  
636 Introducing the next-generation advanced baseline imager on GOES-R. *Bull. Am.*  
637 *Meteorol. Soc.*, **86**, 1079–1096, doi:10.1175/BAMS-86-8-1079.

638 ———, P. Griffith, M. M. Gunshor, J. M. Daniels, S. J. Goodman, and W. J. Lehair, 2017:  
639 A closer look at the ABI on the goes-r series. *Bull. Am. Meteorol. Soc.*, **98**, 681–698,  
640 doi:10.1175/BAMS-D-15-00230.1.

641 Susskind, J., C. D. Barnet, and J. M. Blaisdell, 2003: Retrieval of Atmospheric and Surface  
642 Parameters From AIRS / AMSU / HSB Data in the Presence of Clouds. *IEEE Trans.*  
643 *Geosci. Remote Sens.*, **41**, 390–409, doi:10.1109/TGRS.2002.808236.

644 Turner, D. D., R. O. Knuteson, H. E. Revercomb, C. Lo, and R. G. Dedecker, 2006: Noise  
645 reduction of atmospheric emitted radiance interferometer (AERI) observations using  
646 principal component analysis. *J. Atmos. Ocean. Technol.*, **23**, 1223–1238,

647 doi:10.1175/JTECH1906.1.

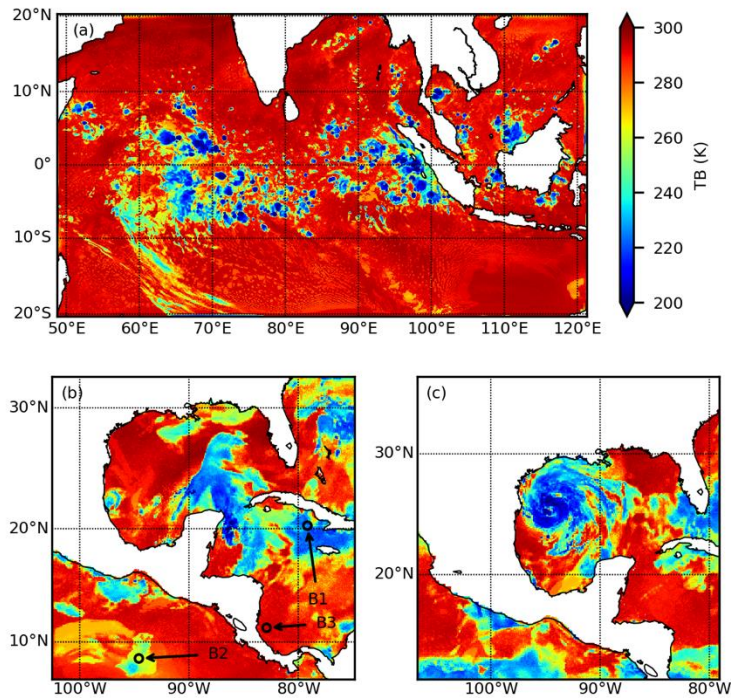
648 Wu, W., X. Liu, D. K. Zhou, A. M. Larar, Q. Yang, S. H. Kizer, and Q. Liu, 2017: The  
649 application of PCRTM physical retrieval methodology for iasi cloudy scene analysis.  
650 *IEEE Trans. Geosci. Remote Sens.*, **55**, 5042–5056,  
651 doi:10.1109/TGRS.2017.2702006.

652 Xu, D., Z. Liu, X. Y. Huang, J. Min, and H. Wang, 2013: Impact of assimilating IASI  
653 radiance observations on forecasts of two tropical cyclones. *Meteorol. Atmos. Phys.*,  
654 **122**, 1–18, doi:10.1007/s00703-013-0276-2.

655 Ying, Y., and F. Zhang, 2017: Practical and Intrinsic Predictability of Multiscale Weather  
656 and Convectively Coupled Equatorial Waves during the Active Phase of an MJO. *J.*  
657 *Atmos. Sci.*, **74**, 3771–3785, doi:10.1175/JAS-D-17-0157.1.

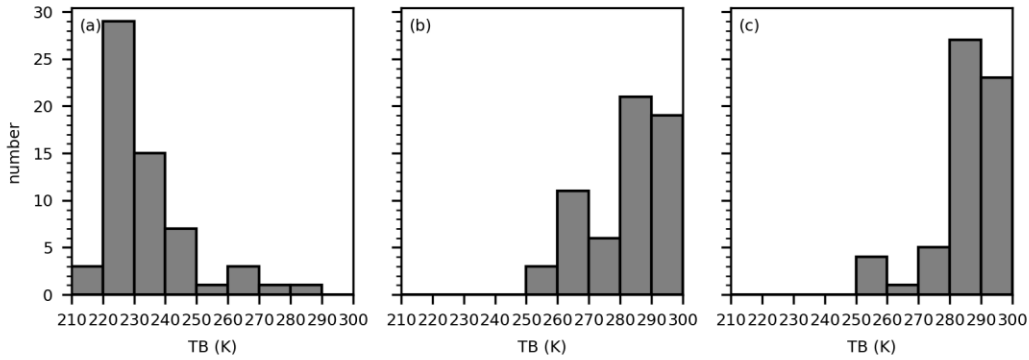
658 Zhang, F., M. Minamide, and E. E. Clothiaux, 2016: Potential impacts of assimilating all-  
659 sky infrared satellite radiances from GOES-R on convection-permitting analysis and  
660 prediction of tropical cyclones. *Geophys. Res. Lett.*, **43**, 2954–2963,  
661 doi:10.1002/2016GL068468.

662



663

664 *Figure 1: Simulated brightness temperatures for AIRS channel with wavelength 12.183  $\mu\text{m}$*   
 665 *for (a) DYNAMO case at 0000 UTC 21 Oct. 2011, (b) early stage Harvey case at 1400*  
 666 *UTC 22 Aug. 2017, and (c) developed stage Harvey case at 0200 UTC 25 Aug. 2017. Small*  
 667 *circles labeled with B1, B2, and B3 in panel (b) indicate high-cloud, mid-cloud, and low/no*  
 668 *cloud location respectively.*

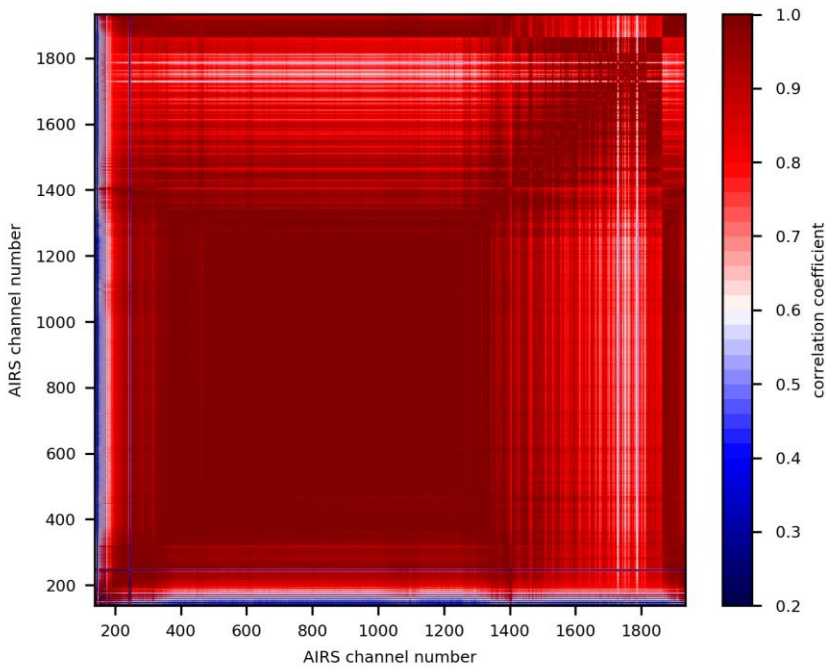


669

670 *Figure 2: Histograms of brightness temperatures of AIRS channel with wavelength 12.183*

671 *μm for all 60 ensemble members at locations labeled with (a) B1, (b) B2, and (c) B3 in*

672 *Figure 1(b).*



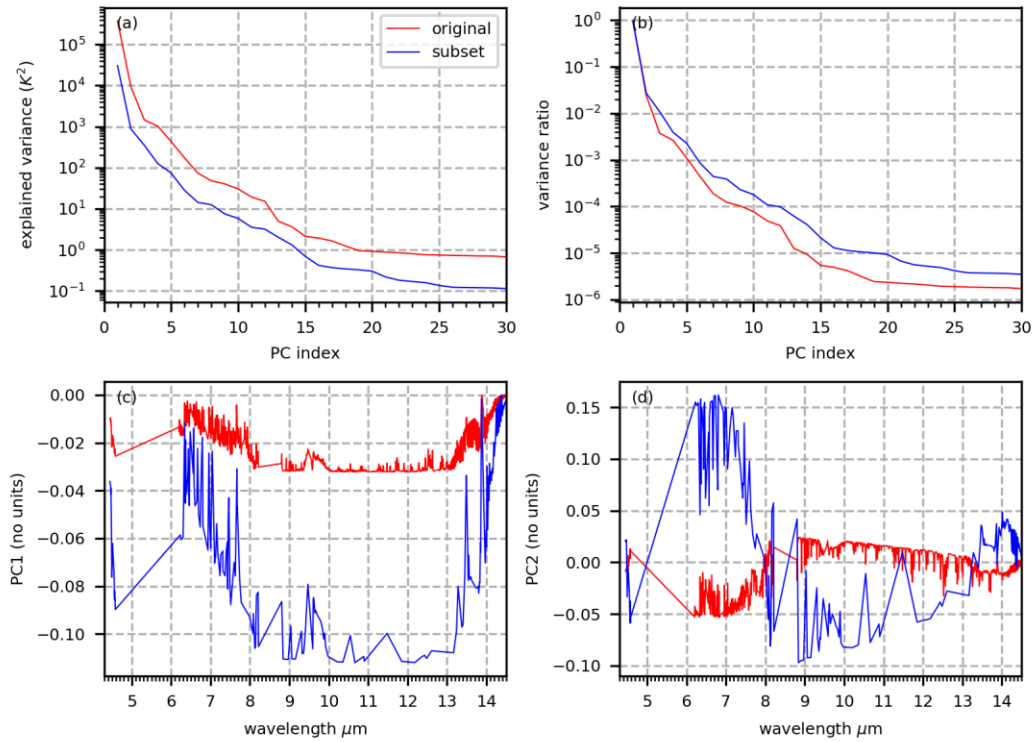
673

674 *Figure 3: Inter-channel correlations for simulated AIRS channels of DYNAMO case. The*

675 *minimum value of the color bar is set to 0.2. The actual minimum correlation value between*

676 *channels is -0.36.*





677

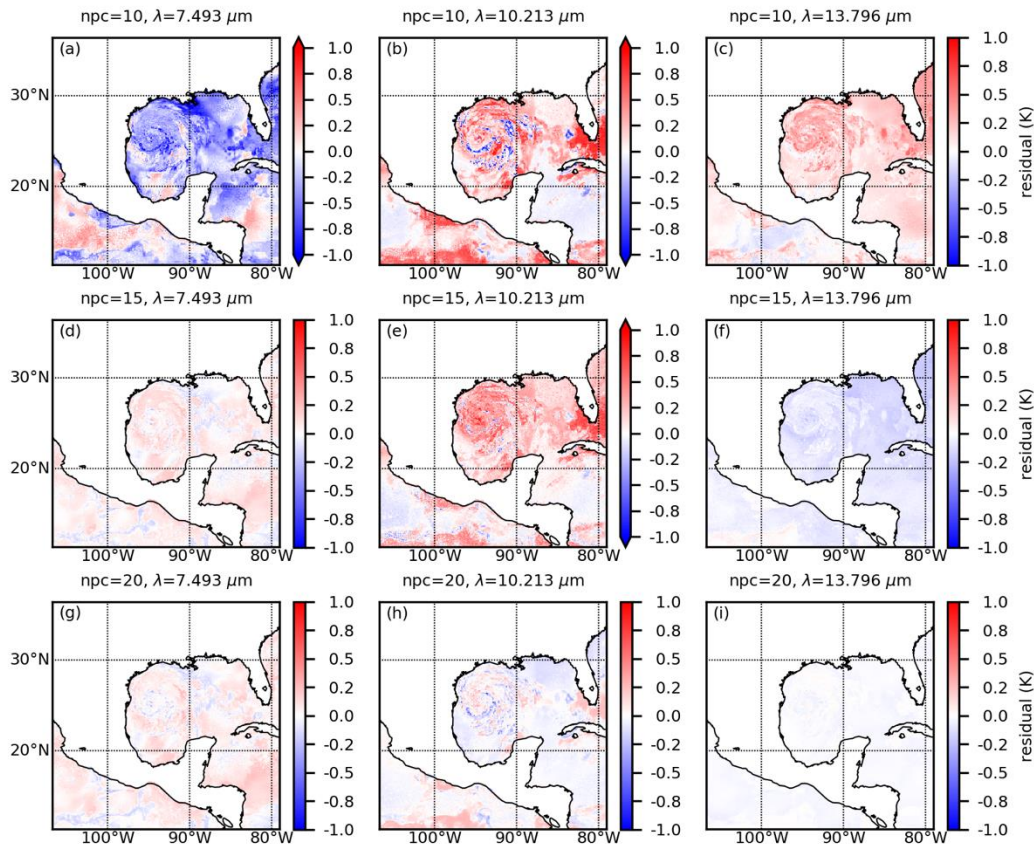
678

679 *Figure 4: (a) variance explained by each PC, (b) variance explained ratio for each PC, (c)*

680 *the 1<sup>st</sup> PC, and (d) the 2<sup>nd</sup> PC of the full-spectral-resolution BTs (red) and the subset BTs*

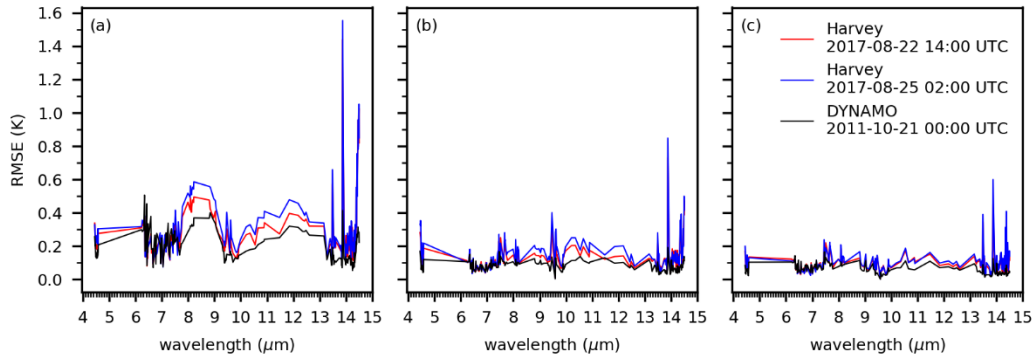
681 *(blue).*

682



683

684 *Figure 5: The reconstruction error when the first 10 (panels a-c), 15 (panels d-f), or 20*  
 685 *(panels g-i) leading PCs were used to reconstruct the hyperspectral BTs simulated from*  
 686 *the 1<sup>st</sup> member of case 3, i.e., the well-developed Harvey (2017) case, for a water vapor*  
 687 *channel with wavelength 7.493 μm (a, d, and g), a window channel with wavelength 10.213*  
 688 *μm (b, e, and h), and a temperature sounding channel with wavelength 13.796 μm (c, f, i).*



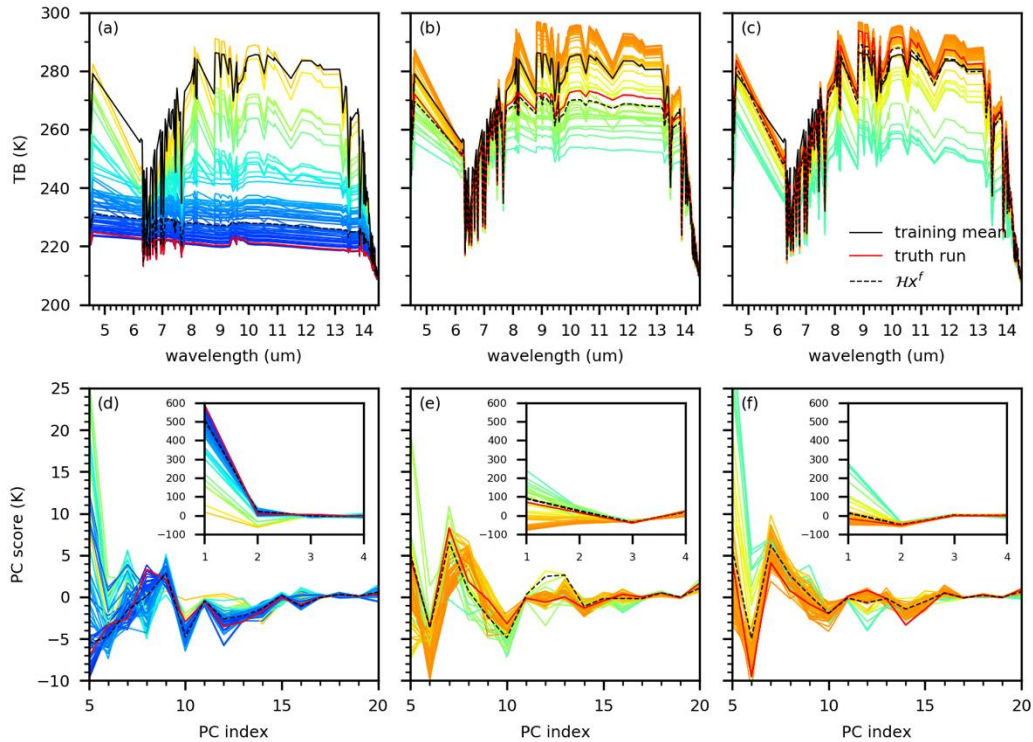
689

690 *Figure 6: Root mean square of reconstruction error over the entire domain of the*

691 *DYNAMO case (black), the 1<sup>st</sup> member of Harvey (2017) case at 1400 UTC 22 August 2017*

692 *(red), and at 0200 25 August 2017 (blue) for all the channels when using (a) 10, (b) 15,*

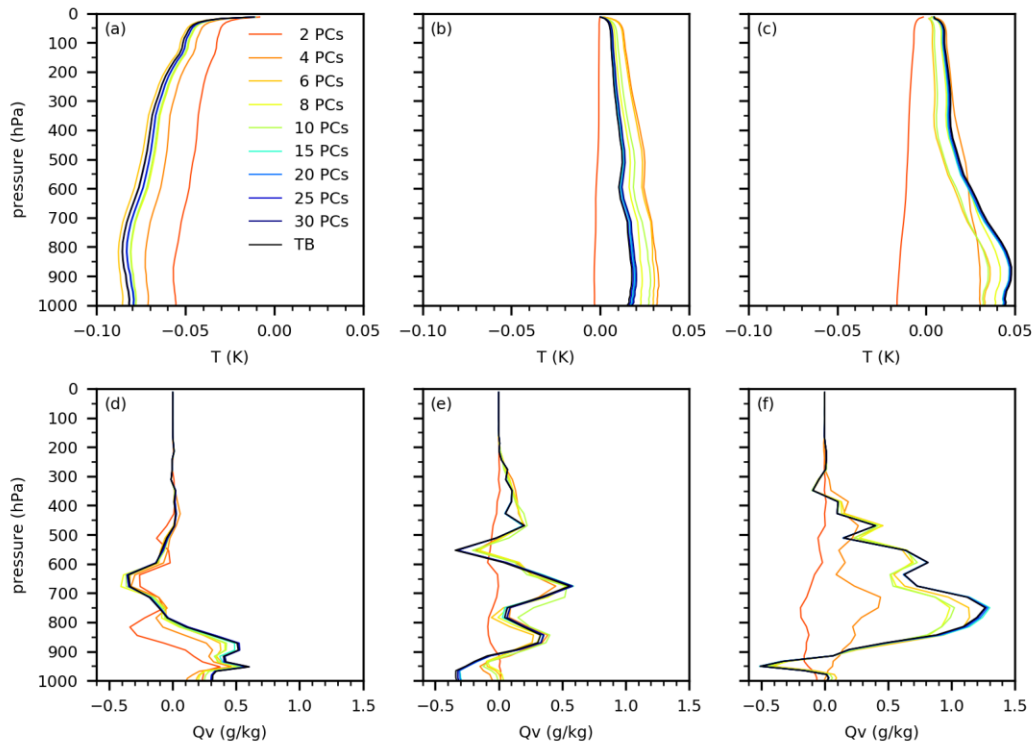
693 *and (c) 20 PCs.*



695

696 *Figure 7: Simulated hyperspectral BTs of the 40<sup>th</sup> ensemble member (red lines), that of the*  
 697 *other 59 ensemble members (other colored lines), hyperspectral BT calculated using*  
 698 *ensemble mean atmospheric states as inputs to CRTM (black dashed lines), and the mean*  
 699 *of hyperspectral BT of the training dataset (black lines) at location (a) B1, (b) B2, and (c)*  
 700 *B3 labeled on Figure 1(b), and the 20 leading PC scores of the 40<sup>th</sup> ensemble member (red*  
 701 *lines), that of the other 59 ensemble members (other colored lines), and that of the*  
 702 *ensemble mean (black dashed lines) at location (d) B1, (e) B2, and (f) B3 labeled on Figure*  
 703 *1(b). The inner plots of panels (d-f) show PC scores corresponding to PC1~PC4 while the*  
 704 *outer plots show that corresponding to PC5~PC20.*

705



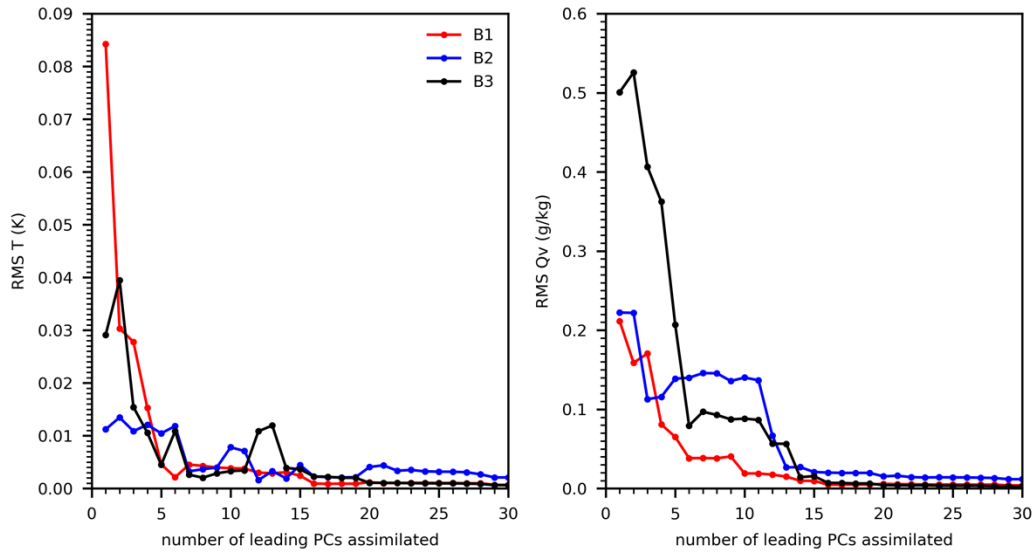
706

707

708 *Figure 8: EnKF increments on (a-c) temperature profiles and (d-f) from assimilating*

709 *hyperspectral BTs (black) and assimilating 2 to 30 leading PC scores (colored lines). The*

710 *legend in panel (a) shows the number of PC scores assimilated according to each color.*



711

712

713 *Figure 9: RMS of the difference between increments from assimilating AIRS brightness*

714 *temperature and PC scores calculated over all layers for (a) temperature profiles and (b)*

715 *water vapor profiles at the locations B1 (red), B2 (blue), and B3 (black) when different*

716 *number of PCs were assimilated.*

717