



Advanced Deep Learning-Based Supervised Classification of Multi-Angle Snowflake Camera Images

C. Key, A. Hicks, and B. M. Notaroš[†]

Department of Electrical and Computer Engineering

Colorado State University, Fort Collins, CO, USA

Submitted to **Journal of Atmospheric and Oceanic Technology**

November 18, 2020

Revised Manuscript, 5 March 2021

Second Revision, 1 June 2021

[†]Corresponding Author:
Branislav M. Notaroš
Colorado State University
Department of Electrical and Computer Engineering
1373 Campus Delivery
Fort Collins, CO 80523, USA
Phone: (970) 491-3537, Fax: (970) 491-2249
Web: www.engr.colostate.edu/~notaros
E-mail: notaros@colostate.edu

Early Online Release: This preliminary version has been accepted for publication in *Journal of Atmospheric and Oceanic Technology*, may be fully cited, and has been assigned DOI 10.1175/JTECH-D-20-0189.1. The final typeset copyedited article will replace the EOR at the above DOI when it is published.

13

14

ABSTRACT

15 We present improvements over our previous approach to automatic winter hydrometeor
16 classification by means of convolutional neural networks (CNNs), using more data and improved
17 training techniques to achieve higher accuracy on a more complicated dataset than we had
18 previously demonstrated. As an advancement of our previous proof-of-concept study, this work
19 demonstrates broader usefulness of deep CNNs by using a substantially larger and more diverse
20 dataset, which we make publicly available, from many more snow events. We describe the
21 collection, processing, and sorting of this dataset of over 25,000 high-quality multiple-angle
22 snowflake camera (MASC) image chips split nearly evenly between five geometric classes:
23 aggregate, columnar crystal, planar crystal, graupel, and small particle. Raw images were
24 collected over 32 snowfall events between November 2014 and May 2016 near Greeley,
25 Colorado and were processed with an automated cropping and normalization algorithm to yield
26 224x224 pixel images containing possible hydrometeors. From the bulk set of over 8,400,000
27 extracted images, a smaller dataset of 14,793 images was sorted by image quality and
28 recognizability (Q&R) using manual inspection. A presorting network trained on the Q&R
29 dataset was applied to all 8,400,000+ images to automatically collect a subset of 283,351 good
30 snowflake images. Roughly 5,000 representative examples were then collected from this subset
31 manually for each of the five geometric classes. With a higher emphasis on in-class variety than
32 our previous work, the final dataset yields trained networks that better capture the imperfect
33 cases and diverse forms that occur within the broad categories studied to achieve an accuracy of
34 96.2% on a vastly more challenging dataset.

35

Significance Statement

36 Classification of precipitation, namely, deciding to which of the several typical classes of
37 winter hydrometeors the observed particles belong, can enrich our understanding of polarimetric
38 radar signatures of snow, as well as ice cloud processes and the resulting precipitation
39 production. The high-resolution photographs of snowflakes collected by the multi-angle
40 snowflake camera (MASC) are especially suitable for snowflake classification. However,
41 classifying particle types from MASC photographs by visual inspection is not practical given the
42 typical amounts of MASC data. We present advanced automatic deep machine learning-based
43 classification of MASC images using convolutional neural networks. This study demonstrates
44 broad usefulness of our approach yielding trained networks that achieve extremely high
45 classification accuracy on a large and diverse dataset from many snow events.

46

47 1. Introduction

48 Snowflake classification is important for improved weather radar, assessment of storm
49 structure, and characterization of winter precipitation events from ground sensors (Zhang et al.
50 2011, Straka et al. 2000, Libbrecht 2017). Several types of in-situ image capturing devices used
51 for ground-based collection of data relevant to snowflake classification include the Two-
52 Dimensional Video Disdrometer (Schönhuber et al. 2008), the Precipitation Instrument Package
53 (an improved version of the system in Newman et al. 2009), and the Multi-Angle Snowflake
54 Camera (MASC). We focus on snowflake images collected by MASC systems in the present
55 study. To allow researchers to study the microphysical characteristics of snowfall, relevant to a
56 storm's composition, the MASC captures high resolution images of falling hydrometeors from

57 several angles. These images can be processed to extract images of individual snowflakes from a
58 variety of perspectives, or even used to generate 3D models of hydrometeors automatically
59 (Kleinkort et al. 2017). A MASC system is capable of capturing tens to hundreds of thousands of
60 images during a single winter storm event, leading to datasets too large for manual classification.
61 This has been a major motivation for accurate, automated snowfall classification.

62 Existing approaches to automated snowfall classification from MASC images vary and
63 include the excellent work of Praz et al. (2017), our previous work (Hicks and Notaroš 2019),
64 and an unsupervised technique (requiring no human input) from Leinonen and Berne (2020). The
65 multinomial logistic regression (MLR)-based method described in (Praz et al. 2017) has been
66 demonstrated effective but requires careful definition and algorithmic extraction of several image
67 features from which classifications are made. This approach has achieved an outstanding 95%
68 classification accuracy, but may be somewhat rigid, relying on human-described features such as
69 morphological skeleton statistics, rotational symmetry, and gray-level co-occurrence. Older
70 supervised classification work in Lindqvist et al. (2012), similarly, applies principal component
71 analysis coupled with Bayesian and weighted nearest-neighbor techniques to classify ice-cloud
72 particles, typically achieving accuracies between 80% and 90%. We have previously presented
73 convolutional neural networks (CNNs) as a robust alternative that can easily be applied and
74 generalized in a black-box manner without expert definition of features. Both methods, of
75 course, require manual input to generate training and test data labels. The work of Leinonen and
76 Berne (2020), on the other hand, automatically classifies snowflake images by exploring the
77 latent space of generative, as opposed to predictive, models. Such unsupervised approaches are
78 extremely promising for discriminating and classifying different hydrometeor images in general,
79 but an unsupervised method inherently produces its own categories, rather than directly

80 assigning images to existing, known categories with which researchers are likely already
81 familiar.

82 Accordingly, we offer improvements to our existing CNN-based, supervised approach
83 (Hicks and Notaroš 2019), using more data and improved training techniques to achieve higher
84 accuracy on a more complicated dataset than we had previously demonstrated. As an
85 advancement of our previous proof of concept study, which used a geometric dataset focused on
86 easily identifiable examples of each of the snowflake classes considered, a principal goal of this
87 work is to demonstrate broader usefulness of deep CNNs for automated snowfall classification
88 by using a larger dataset containing wider in-class variety. We present improved training
89 methods and new, automated techniques for detection, cropping, and normalization of snowflake
90 images as well as quality and recognizability preprocessing of image data. From these
91 improvements, we demonstrate higher overall test accuracy on a vastly more challenging dataset
92 than that used in our previous work. Together, these improvements constitute an accurate,
93 efficient, and robust supervised machine learning approach to snowflake classification, using
94 deep neural networks and images collected by the MASC or another image-based particle
95 recording instrument or system.

96

97 **2. Data Collection and Image Processing**

98 This section describes the collection of raw MASC images as well as the automated
99 cropping and normalization performed on raw images to isolate potential snowflakes present in
100 each image.

101

102 **2.1 Raw Image Collection**

103 The 3,458,848 raw images used to generate the training set were collected from several
104 winter weather events between November 2014 and May 2016 using a modified MASC system.
105 The system was located at a surface instrumentation field site established under MASCRAD
106 (MASC + RADar) (Notaroš et al 2016; Bringi et al. 2017; Kennedy et al. 2018). This is the same
107 site and system used for data collection in Hicks and Notaroš (2019). The MASCRAD field site
108 is located at the Easton Valley View Airport in La Salle, near Greeley, Colorado, shown in
109 Figure 1. The MASC system, along with other ground-level instrumentation at the site, is
110 situated within a double fence intercomparison reference (DFIR). Raw images from both winter
111 storm events used in Hicks and Notaroš (2019) constitute a subset of the total raw image set used
112 in the present work. Details of the MASC system used are presented in Hicks and Notaroš
113 (2019). Although the MASC allows for collection of snowflake imagery from multiple angles to
114 help determine three-dimensional shape (Kleinkort et al. 2017), we did not make use of this
115 feature directly for the present work. As described in Leinonen and Berne (2020), it is common
116 that a given snowflake will only be captured at usable quality by a single camera of a multi-
117 camera system, the snowflake often out of focus or occluded in other fields of view, so limiting
118 study to only snowflakes that appear at high quality in all fields of view substantially reduces the
119 number of useable examples. By limiting study to single-view cases, we were able to manually-
120 classify thousands, rather than hundreds of snowflakes at a cost of increased ambiguity due to
121 lack of multi-angle data. Note that we did not explicitly remove cases where a single snowflake
122 was imaged from multiple angles when forming the dataset for the present work.

123

124

125 **2.2 Detection, Cropping, and Normalization**

126 As the MASC produces raw, wide field of view images, typically containing many
127 snowflakes, it is necessary to isolate individual examples for classification. All images were
128 processed in grayscale (single channel). To detect possible flakes in each raw MASC image, we
129 first normalized the entire grayscale image, dividing all pixel values by the maximum brightness
130 value. An example of a normalized raw image is shown in Figure 2. We then converted the
131 grayscale image into a binary image by application of a threshold. Pixels in the grayscale image
132 with brightness greater than or equal to the threshold were assigned value 1, and pixels less than
133 the threshold were assigned value 0. For the present work, this threshold was set to 0.1. We then
134 set any pixels in the binary image with value 0 to 1 if they were within a 2-pixel radius (using
135 Chebyshev distance) of any pixel that had already been assigned value 1 in the previous
136 thresholding step. The example image from Figure 2 is shown after thresholding and application
137 of the 2-pixel radius in Figure 3. This radius was chosen by hand as a reasonable value. Next, we
138 computed sets of connected components in the binary image. A connected component is any
139 group of active (value 1) pixels that form an unbroken group. If a connected component
140 contained fewer than 26 active pixels, it was discarded. For each connected component not
141 discarded, we cropped a rectangular region from the original grayscale image corresponding to
142 its bounding box. Two such examples produced from Figure 3 are shown in Figure 4(a) and (d).
143 Cropped images were then contrast scaled linearly such that the top 1% of brightest pixels were
144 saturated. Figure 4(b) and (e) show cropped image examples from the previous step after contrast
145 scaling. Note that contrast scaling destroys some information theoretically available in the
146 images (by loss of absolute brightness and saturation of some pixels). However, we found that
147 brightness variations between flakes were dominated by differing lighting conditions, rather than

148 useful information like snowflake class, so contrast scaling was performed to give the network
149 input for which pixel brightness variations are dominated by microphysical characteristics rather
150 than lighting conditions. After scaling, any cropped image was rejected if the mean value of its
151 pixels was greater than 0.5. We then centered each remaining cropped, scaled image on a
152 224x224 black background to produce final image chips. Examples are shown Figure 4(c) and
153 (f). Cropped images that exceeded the 224x224 image chip sized were cropped to 224x224
154 pixels after centering. Camera configurations are

155 This approach to cropping and normalization was arrived at for several reasons. In
156 contrast to simply cropping a 224x224 pixel region centered on each connected component in an
157 image (or similar), we found that the above method significantly reduced the number of image
158 chips that contained multiple, physically disconnected snowflakes. In other words, during heavy
159 snowfall events, we found it was common for two or more snowflakes to appear within 224
160 pixels of each other. By cropping a tight bounding box as above, we were able to recover far
161 more closely spaced snowflakes into usable, unambiguous image chips. Rejection of cropped,
162 scaled images with mean pixel value greater than 0.5 rejected most crops of the sky and
163 background that did not actually contain a snow particle. By also rejecting connected
164 components with pixel counts below 26, we avoided cases where a single bright pixel caused a
165 false detection. In general, the described cropping and normalization approach was able to detect
166 far more small particles, and dim, unrimed planar crystals than the approach used for Hicks and
167 Notaroš (2019). Application of this cropping algorithm to all raw MASC images from November
168 2014 to May 2016 produced 8,441,563 image chips.

169
170

171 **3. Hydrometeor Classification Scheme and Training Sets**

172 This section describes how the 8,441,563 224x224 pixel image chips extracted from raw
173 MASC images were automatically sorted to quality classes and how images from the best class
174 were manually sorted into the five geometric categories studied. A total of 25,199 examples were
175 manually sorted for the final geometric dataset covering 32 snowfall events, an event defined
176 here as a period during which no more than 24 hours passed between collection of any two
177 image chips identifies as snowflakes during manual classification). All classification was
178 performed by a single analyst who reviewed each image at least three times. Overall, we are
179 confident the manual classifications used for training accurately represent the opinions of our
180 analyst and have made this dataset available at Key et al. (2021). Note, however, that our use of
181 only one human analyst has potential to introduce more bias relative to other work for which
182 multiple humans performed analysis, such as Praz et al. (2017). We had originally planned to
183 also produce an expanded riming dataset in addition to the presented geometric dataset, but we
184 found that some riming degrees were insufficiently represented for production of a larger,
185 balanced riming dataset from our current pool of raw images. We hope to contribute such a
186 dataset in future work.

187

188 **3.1 Quality and Recognizability Preprocessing**

189 The snowflake detection, cropping, and normalization method described in Section 2.2
190 remains imperfect. Therefore, many of the image chips produced contained bright points from a
191 raw image that are not snowflakes. These included sources like glare, sensor noise, and
192 sky/ground glow. In addition, operators of the MASC system occasionally forgot to turn off data
193 collection while calibrating and testing the system after maintenance and redeployment. This led

194 to captures of test probes, hands, coins, and other objects to occasionally appear in the raw image
195 dataset. Several examples of image chips due to non-flake objects are shown in Figure 5.

196 For image chips that contain snowflakes, there is an inherent range of quality. Some
197 flakes appear out of focus in raw images. Others are poorly cropped, either due to over-cropping
198 by the image processing method in Section 2.2, or because they originally appeared partially out
199 of field of view in a raw MASC image. We considered image chips containing snowflakes to fall
200 into four recognizability categories: Bad-Crop, Bad, Okay, and Good. Image chips in the Bad-
201 Crop category are those where unambiguous recognizability of the imaged snowflake is made
202 difficult due to over-cropping by the processing method described in Section 2.2 or part of the
203 flake appearing out of field-of-view in the raw image, leaving a substantial portion of the flake
204 absent from the image chip. Note that cases where a flake was simply too large to fit in a single
205 image chip were not included in the Bad-Crop category. In our manual exploration of the dataset,
206 such flakes were almost exclusively in the AG class and easily identifiable despite cropping to
207 224x224 pixels. Rather, over-cropping by the processing described in Section 2.2 is typically due
208 to poor or uneven illumination of the flake causing the rectangular bounding box of the resulting
209 connected component to not contain most of the pixels covered by the snowflake. Four examples
210 of Bad-Crop image chips are shown in the first column of Figure 6. Bad image chips are those
211 for which poor focus or poor illumination rendered the target snowflake unrecognizable. Image
212 chips containing more than one disjoint (non-aggregated) snow particle are also included in the
213 Bad category, regardless of lighting and focus. We consider two snow particles disjoint if they
214 were clearly identifiable as discrete, physically unconnected particles by our human analyst. Four
215 such examples are shown in the second column of Figure 6. Okay image chips were those that
216 contained a recognizable snowflake but suffered from mild blur or high background noise that

217 made examination of microphysical characteristics difficult. Four examples of okay image chips
218 are shown in the third column of Figure 6. Good image chips were those that were free of
219 substantial over-cropping and clear enough to identify relevant microphysical features. Column
220 four of Figure 6 shows four examples of Good image chips.

221 To avoid wasting human time visually inspecting images that did not contain flakes or
222 were of quality too poor to use, we trained a preliminary quality and recognizability (Q&R)
223 classifier on a small, manually sorted subset of the 8,441,563 image chips. This classifier was
224 implemented by necessity to reduce the data volume needing manual inspection, and its results
225 were not further analyzed or verified in the present work. To train the Q&R classifier, we
226 collected at least 1,500 examples for each of five categories: Not-Flake, Bad-Crop, Bad, Okay,
227 and Good, with an emphasis on variety within each class. Counts per category for the Q&R
228 dataset are presented in Table 1 along with descriptions. When collecting example images, we
229 included roughly equal numbers of examples from each geometric class in Okay and Good
230 categories to avoid biasing the classifier against a given geometric type. The Q&R classifier was
231 trained using the same methodology used for the geometric classifier in Hicks and Notaroš
232 (2019). For training, 1,500 examples from each Q&R category were drawn randomly. The
233 trained Q&R classifier was then applied to all 8,441,563 image chips to sort each into Not-Flake
234 (3,791,326), Bad-Crop (723,550), Bad (3,062,288), Okay (582,333), and Good (282,001)
235 categories. Only image chips assigned by the Q&R network to the Good category were examined
236 to produce the geometric dataset for the present study.

237

238

239 **3.2 Geometric Classes**

240 A variety of attempts have been made to classify snowflakes (Nakaya and Sekido 1936,
241 Magono and Lee 1966, Korolev and Sussman 2000, Grazioli et al. 2014, Vasquez-Martin et al.
242 2020). As in our previous work (Hicks and Notaroš 2019), we chose to use the scheme adopted
243 by Praz et al. (2017) for training and testing of their multinomial logistic regression snowflake
244 classifier. We summarize this scheme here.

245 The scheme uses the nine categories of snowflakes defined in Magono and Lee (1966),
246 with a few simplifications for data availability. Praz et al. (2017) additionally defined the
247 Aggregate and Small Particle classes. Aggregates are defined as single snowflakes that are the
248 result of in-air collision of two or more particles. Small Particles are snowflakes whose features
249 are too small to categorize. Note that this is based on the subjective opinion of the analyst, rather
250 than a strictly defined size threshold. Simplifications from Magono and Lee (1966) and addition
251 of AG and SP classes resulted in 10 individual categories, of which only six were used in Praz et
252 al. (2017) due to data availability: Aggregates (AG), Small Particles (SP), Columnar Crystals
253 (CC), Planar Crystals (PC), Combination of Columnar and Planar Crystals (CPC), and Graupel
254 (GR). As in Hicks and Notaroš (2019), we chose to exclude the CPC class from the present study
255 due to data availability. We found only a few hundred clear examples of CPC in the Good Q&R
256 class. CPC appeared far less commonly than the next rarest class, GR, which had several
257 thousand Good Q&R examples. Image chips that fell into unconsidered categories, like CPC, we
258 simply omitted from consideration for the present work.

259

260

261 **3.3 Building the Geometric Dataset**

262 Our goal in collecting the geometric dataset for the present work was to establish a large,
263 highly varied collection of image chips in each of the five categories considered. Deep neural
264 networks, like that used in Hicks and Notaroš (2019) and the present work can achieve high
265 accuracies but require substantial training data to avoid over-fitting (Simonyan and Zisserman
266 2015, Szegedy et al. 2015). With tens of millions of parameters, deep CNNs like the ResNet-50
267 architecture (He et al. 2016) can store substantial quantities of information to learn highly
268 complicated associations and trends (Zeiler and Fergus 2014). Care must therefore be taken to
269 train such networks on large enough datasets that they cannot simply memorize associations
270 between specific images and their labels or extract spurious trends.

271 Another important consideration is balance between classes during training. Unless
272 special precautions such as class-specific learning rates are used (not used in the present study),
273 training a neural network on a dataset biased toward a particular class will often bias the network
274 toward that class. As an extreme example, consider a network trained on a dataset of 900 GR
275 images and 100 PC images; the network can attain 90% accuracy on the training set simply by
276 learning to label every image as GR. It is therefore important to present the network with roughly
277 equal numbers of examples in each class during training.

278 To account for these factors, we limited the number of examples in our geometric dataset
279 for each class to the maximum number of Good Q&R examples we could find for the rarest class
280 considered. After CPC (not considered), GR was the rarest class, for which we could only find
281 roughly 5,000 examples. Accordingly, we collected roughly 5,000 examples of each of the other
282 classes considered, for a total of 25,199 examples. Exact image chip counts per class are
283 presented in Table 2. Figures 7 through 11 show representative examples from the final AG, CC,

284 GR, PC, and SP sets, respectively. When collecting examples for each class, we put emphasis not
285 only on archetypical examples, but also examples we considered good counterexamples to
286 possible oversimplifications of each class: e.g. AGs are always large, PCs always have six-fold
287 symmetry, or GR always has a smooth outline. Image chips were not included in the geometric
288 dataset if we could not determine an appropriate label based on information present in the image
289 chip alone, i.e. no multi-angle information was used during manual sorting. We note overall that
290 there is an inherent subjectivity in identification of snowflakes in single-view images, especially
291 for classes like GR (Figure 9), for which distinguishing from other heavily-rimed particles is
292 subjective, and SP (Figure 11), for which deciding unrecognizability of features due to small size
293 is highly subjective. We did not avoid using backlit examples where available, although these
294 were rare, only occurring where a snow particle was imaged while falling in front of a
295 sufficiently bright glare point in the background. Due to their rarity, inclusion of backlit cases
296 likely did not have a substantial impact on accuracy of the trained network. Our analyst
297 recollects seeing at most a dozen backlit cases during manual classification, but such cases were
298 assigned no special designation or identifying information that would make quantification of
299 their impact possible without another manual review of the dataset.

300

301 **4. Convolutional Neural Networks Methodology**

302 A brief discussion of the network architecture is presented in this section. We also
303 present a summary of the training method and hyperparameters used. Note that, although the
304 network architecture remains the same as that in our previous work (Hicks and Notaroš 2019),
305 hyperparameters for training differ.

306

307 **4.1 Neural Network Architecture**

308 We used an identical ResNet-50 architecture to that in Hicks and Notaroš (2019). The
309 ResNet-50 architecture has been demonstrated as an excellent balance between speed and
310 accuracy for image classification tasks and is described in detail in (He et al. 2016). The residual
311 approach, in general, was groundbreaking at the time of its publication, as it presented an elegant
312 solution to the vanishing gradient problem that had previously limited scaling of CNN accuracy
313 with increased depth. The use of residual connections (or similar), as described in (He et al.
314 2016) has since been widely adopted by deep learning researchers and practitioners. As in (Hicks
315 and Notaroš 2019), we used a ResNet-50 model that had been pretrained for general image
316 classification on the ImageNet database (Russakovsky and Deng et al. 2015). We also
317 experimented with randomly initialized (no pretraining) versions of the same architecture but
318 found no substantial benefit. We therefore chose to only focus on the pretrained model for the
319 present work for easy comparison with (Hicks and Notaroš 2019). A necessary change made to
320 the architecture was reduction in the number of outputs of the final, fully connected layer for our
321 substantially lower number of classes (the original ResNet-50 architecture trained on ImageNet
322 had 1,000 classes, not five). Weights in the modified fully connected layer were initialized
323 randomly.

324

325 **4.2 Training Method and Hyperparameters**

326 As in (Hicks and Notaroš 2019), network performance was determined by cross-entropy
327 error, and network weights and biases were optimized by stochastic gradient descent to minimize
328 this loss function. For training, validation, and testing, we again limited the number of examples
329 used in each class to the number of examples available in the smallest class (in the present work,

330 GR with a total of 5,000 hand-classified image chips available). The examples used from classes
331 with raw counts larger than the minimum were drawn randomly. We again used a mini-batch
332 size of 10. Beyond this, we made several changes to the hyperparameters and training method
333 used in (Hicks and Notaroš 2019). Our dataset was also substantially larger; the testing set alone,
334 in this case, was comparable in size to the entire geometric dataset used for (Hicks and Notaroš
335 2019), roughly 1,450 examples. In the present study, we randomly selected 500 examples from
336 each class for a total of 2,500 testing examples. The remaining 22699 examples were randomly
337 partitioned into a training set (~90%) and a validation set (~10%), both evenly distributed among
338 the classes studied. The random partitioning between training and validation was unique to each
339 training run. Only the training and validation sets were used for hyperparameter tuning, which
340 was performed by a mix of expert hand-tuning and small parametric sweeps and included tuning
341 of the mini-batch size, learning rate, and number of training epochs. We also trained for
342 substantially longer than our previous work, training for a total of 20 epochs, as opposed to 10.
343 The training set was shuffled (re-ordered) randomly every epoch. An epoch is defined as one
344 complete pass through the training set, so, the present training dataset containing many more
345 examples than that available in (Hicks and Notaroš 2019), this corresponds to roughly a 30-fold
346 increase in training time. We were able to extend the training time substantially due to
347 prevention of overfitting by the larger training dataset used in the present work. As opposed to
348 the constant learning rate of 0.0003 used in (Hicks and Notaroš 2019), we began with a learning
349 rate of 0.001 which was then scaled by a factor of $1/\sqrt{10}$ every five epochs. We found this led to
350 a small but noticeable improvement in final network accuracy. We expect improvements in
351 network accuracy could be further improved with additional hyperparameter tuning using more
352 compute resources for large parametric sweeps.

353

354 **5. Results and Discussion**

355 This section presents and discusses the performance of the trained classification networks
356 on the test dataset. The final mean test accuracy achieved was 96.23% with a standard deviation
357 of 0.29% across 10 training runs, the individual test accuracies of which are presented in Table 3.
358 Only the order in which images were presented to the network and random partitioning of non-
359 test images between training and validation differed between training runs. We expect we could
360 have achieved even higher accuracy if we had limited our dataset to only archetypal examples,
361 but this would have diminished the usefulness of the dataset and resulting trained model for
362 general snowfall classification tasks.

363 Figure 12 shows accuracy and loss of a typical trained network (test accuracy close to the
364 mean) on the training and validation set with respect to training iteration (and epoch, indicated
365 by alternating vertical bands) for a typical training run. There is no evidence of overfitting, and
366 validation accuracy increased nearly monotonically with iteration count. Overfitting, if present,
367 would be apparent in Figure 12 as divergence of the black validation accuracy and blue training
368 accuracy curves. For the training run shown, the network achieved a validation accuracy of
369 96.1% and a test accuracy of 96.2%. We suspect the much larger size of the geometric dataset is
370 the dominant factor in improving performance over our previous work but did not have sufficient
371 compute time to perform a full parametric sweep to confirm this. We found that network
372 performance on the validation and training sets were comparable, indicating that training, testing,
373 and validation datasets all sampled the underlying distribution of snowflake geometries well. The
374 validation accuracy standard deviation for the 10 example runs shown in Table 3 was 0.42%, and
375 their mean validation accuracy was 96.26%. We attribute the larger validation accuracy standard

376 deviation, as compared to the test accuracy standard deviation, to random selection of the
377 validation set for each training run (the test set did not change between runs). There was little
378 variation between training runs, with the only nominal differences due to this random
379 partitioning of the validation and training sets as well as random re-ordering of the training set
380 during each epoch. Figure 13 shows a confusion matrix for the same network, the training
381 progress of which is shown in Figure 12.

382 In general, trained networks would confuse PC and AG classes most often. We included
383 many difficult examples in the AG class that featured a prominent planar crystal with several
384 less-prominent particles that had adhered due to mid-air collisions, so confusion between the two
385 classes seems understandable to us. Figure 14 presents examples of image chips misclassified by
386 the typical network from Figures 12 and 13. Overall, most misclassifications appear to be blatant
387 errors due to imperfection of the trained model, but several stand out as ambiguous cases or
388 possibly even human error. Figure 14, row 2, column 2, for instance, was assigned by the
389 network to the AG class, having been human labeled as a columnar crystal. Further inspection
390 indicates this snowflake may indeed be a simple aggregate or even a malformed planar crystal,
391 suggesting this misclassification is due human error rather than network error. Figure 14, row 4,
392 column 3 shows a clear planar crystal adhered to a small aggregate of columnar crystals.
393 Although the planar crystal dominates the image chip, the aggregation present indicates the
394 network is correct to assign this image chip to the AG class. Figure 14, row 3, column 4 and row
395 5, column 2, respectively show a GR image chip misclassified as SP and a SP image chip
396 misclassified as GR, respectively. These two cases show the ambiguity of the SP class and the
397 difficulty of drawing a distinction between small GR flakes and relatively large, round SP flakes.
398 Figure 14, row 5, column 3 shows another ambiguous case. Human-classified as SP but network

399 classified as CC, this particle shows possible CC-like features (dominant uniaxial crystal growth)
400 but is barely too small for our analyst to assign confidently to the CC category.

401

402 **6. Conclusions**

403 This paper has presented improvements over our previous approach (Hicks and Notaroš
404 2019) to automated winter hydrometeor classification using deep convolutional neural networks.
405 Using improved training methods and a substantially larger and more complicated dataset from
406 many more snow events than in our previous study, we were able to achieve over 96.2%
407 accuracy on a test set of 2,500 images. We consider this result substantial for several reasons.
408 The MASC is a high-throughput sensor, collecting tens to hundreds of thousands of detectable
409 snowflake images during a winter storm event, so even small accuracy improvements lead to a
410 substantial reduction in the total number of misclassified snowflake images. Namely, this is a
411 ~40% reduction in the fraction of incorrectly classified snowflakes relative to the already very
412 high geometric classification accuracy result reported in our previous work and corresponds to a
413 2.8% increase in overall accuracy. Even more importantly, the dataset of 25,199 image chips
414 sorted by geometric class used in the present study differs substantially from that developed for
415 (Hicks and Notaroš 2019). As a proof of concept study, (Hicks and Notaroš 2019) used a
416 geometric dataset focused on easily identifiable examples of each of the snowflake classes
417 considered. To demonstrate the broader usefulness of deep CNNs for automated snowfall
418 classification, the dataset used in the present study is not only larger but also contains wider in-
419 class variety. In using such a dataset, we have shown that, with a few modifications to the
420 network training process, the geometric classification method described in (Hicks and Notaroš
421 2019) can achieve higher accuracy on a vastly more challenging dataset. Finally, the paper has

422 presented several important components of the CNN-based, supervised approach to snowflake
423 classification, including an improved training method and hyperparameters for training; new
424 automated techniques for snowflake detection, cropping, and normalization of snowflake images;
425 and new quality and recognizability preprocessing of image data. The described methodologies
426 and techniques may be of great use to researchers and practitioners applying the same or similar
427 approaches to hydrometeor classification based on the images collected by the MASC or another
428 image-based particle recording instrument or system.

429

430

Acknowledgment

431 This work was supported in part by the National Science Foundation under Grants AGS-
432 1344862 and AGS-2029806.

433

434

Data Availability Statement

435 The dataset of MASC images generated for and used in this study has been made publicly
436 available at Key et al. (2021).

437

438

439

References

440 Bringi, V. N., P. C. Kennedy, G.-J. Huang, C. Kleinkort, M. Thurai, and B. M. Notaros, 2017:
441 Dual-polarized radar and surface observations of a winter graupel shower with negative Z_{dr}
442 column,” *J. Appl. Meteor. Climatol.*, **56**, 455–470.

443

- 444 Grazioli, J., D. Tuia, S. Monhart, M. Schneebeli, T. Raupach, and A. Berne, 2014: Hydrometeor
445 classification from two-dimensional video disdrometeor data, *Atmos. Meas. Tech.*, **7**, 2869-2882.
446
- 447 He, K., X. Zhang, S. Ren, and J. Sun, 2016: Deep Residual Learning for Image Recognition,
448 *CVPR*.
449
- 450 Hicks, A. and B. M. Notaroš, 2019: Method for Classification of Snowflakes Based on Images
451 by a Multi-Angle Snowflake Camera Using Convolutional Neural Networks. *J. Atmos. Ocean.
452 Tech.*, **36**, 2267– 2282.
453
- 454 Kennedy, P., M. Thurai, C. Praz, V. N. Bringi, A. Berne, and B. M. Notaros, 2018: Variations in
455 Snow Crystal Riming and Z_{DR} : A Case Analysis. *J. Appl. Meteor. Climatol.*, **57**, 695–707.
456
- 457 Key, C., A. Hicks. & B. Notaros. (2021). Colorado State University Geometric Snowflake
458 Classification Dataset (Version 1.0) [Data set]. *Zenodo*. <http://doi.org/10.5281/zenodo.4584200>
459
- 460 Kleinkort, C., G.-J. Huang, V. N. Bringi, and B. M. Notaros, 2017: Visual Hull Method for
461 Realistic 3D Particle Shape Reconstruction Based on High-Resolution Photographs of
462 Snowflakes in Free Fall from Multiple Views. *J. Atmos. Oceanic Technol.*, **34**, 679–702.
463
- 464 Korolev, A. and B. Sussman, 2000: A technique for habit classification of cloud particles, *J.
465 Atmos. Ocean. Tech.*, **17**, 1048-1057.
466

467 Leinonen, J. and A. Berne, 2020: Unsupervised classification of snowflake images using a
468 general adversarial network and K-medoids classification. *Atmospheric Measurement*
469 *Techniques*, **13**, 2949-2964.

470
471 Libbrecht, K. G., 2017: Physical Dynamics of Ice Crystal Growth, *Annu. Rev. Mater. 2017*. **47**,
472 271-295.

473
474 Lindqvist, H., Muinonen, K., Nousiainen, T., Um, J., McFarquhar, G., Haapanala, P., Makkonen,
475 R., and Hakkarainen, H. 2012: Ice-cloud particle habit classification using principal components,
476 *J. Geophys. Res-Atmos.*, 117, 2156–2202.

477
478 Magono, C. and C.W. Lee, 1966: Meteorological classification of natural snow crystals, *J. Fac.*
479 *Sci.*, Hokkaido Univ., Series VII, 2, 321-335.

480
481 Nakaya, U. and Y. Sekido, 1936: General Classification of Snow Crystals and their Frequency of
482 Occurrence. *J. Fac. Sci.*, Hokkaido Univ., Series II, 1, 243-264.

483
484 Newman, A. J., P. A. Kucera, and L. F. Bliven, 2009: Presenting the Snowflake Video Imager
485 (SVI). *J. Atmos. Oceanic Technol.*, **26**, 167–179.

486
487 Notaroš, B. M., V. N. Bringi, C. Kleinkort, P. Kennedy, G.-J. Huang, M. Thurai, A. J. Newman,
488 W. Bang, and G. Lee, 2016: Accurate Characterization of Winter Precipitation Using Multi-
489 Angle Snowflake Camera, Visual Hull, Advanced Scattering Methods and Polarimetric Radar.

490 invited paper, Special Issue on Advances in Clouds and Precipitation, *Atmosphere*, **7**, no. 6, 81–
491 111.

492
493 Praz, C., R. Yves-Alain, and A. Berne, 2017: Solid hydrometeor classification and riming degree
494 estimation from pictures collected with a Multi-Angle Snowflake Camera, *Atmos. Meas. Tech.*,
495 **10**, 1335-1357.

496
497 Russakovsky, O.,*, J. Deng,*, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A.
498 Khosla, M. Bernstein, A.C. Berg, and L. Fei-Fei, (* = equal contribution) ImageNet Large Scale
499 Visual Recognition Challenge. *IJCV*, 2015.

500
501 Schönhuber, M., G. Lammer, and W. Randeu, 2008: The 2D video disdrometer. In *Precipitation:*
502 *Advances in Measurement, Estimation and Prediction*; Michaelides, S., Ed.; *Springer*: Berlin,
503 Germany, 3–31.

504
505 Simonyan, K., and A. Zisserman, (2015): Very deep convolutional networks for large-scale
506 image recognition. *ICLR*, 2015.

507
508 Straka, J., D. S. Zrnić, and A. V. Ryzhkov, 2000: Bulk hydrometeor classification and
509 quantification using polarimetric radar data: Synthesis of Relations. *J. Appl. Meteor.*, **39**, 1341–
510 1372.

511
512 Szegedy, C., W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and
513 A. Rabinovich, 2015: Going deeper with convolutions. *CVPR*.

514

515 Vazquez-Martin, S., T. Kuhn, and S. Eliasson, 2020: Shape Dependence of Falling Snow
516 Crystals' Microphysical Properties Using an Updated Shape Classification. *MDPI Applied*
517 *Sciences*.

518

519 Zeiler, M.D. and R. Fergus, 2014: Visualizing and understanding convolutional neural networks.
520 *ECCV*.

521

522 Zhang, G., S. Luchs, A. Ryzhkov, M. Xue, L. Ryzhkova, and Q. Cao, 2011: Winter Precipitation
523 Microphysics Characterized by Polarimetric Radar and Video Disdrometer Observations in
524 Central Oklahoma. *J. Appl. Meteor. Climatol.*, **50**, 1558–1570.

525

526

527

Tables

528 **Table 1.** Category names, counts, and descriptions for the quality and recognizability dataset, a
 529 balanced subset of which was used to train a presorting network using the methods of Hicks and
 530 Notaroš (2019).

Category Name	Count	Description
Not-Flakes	7,020	Object other than a snowflake present in the image. Examples include sensor noise, glare, sky/ground glow, and calibration probes.
Bad-Crop	1,500	Likely snowflake present, but poor cropping leaves a substantial portion of the snowflake out of the image chip, interfering with geometric classification.
Bad	1,977	Likely snowflake present, but poor lighting or focus prevent identification. Image chips containing more than one disjoint (non-aggregated) snowflake are also assigned to this class, regardless of image quality.
Okay	2,796	Focus and lighting are good enough to identify coarse flake features, and likely geometric class, but are insufficient to capture microphysical characteristics.
Good	1,500	Lighting and focus are good enough to resolve microphysical characteristics and determine snowflake geometric class.

531

532 **Table 2.** Number of examples in each class for the geometric dataset.

Class Name	Count
AG	5,038
CC	5,021
GR	5,000
PC	5,014
SP	5,126

533

534 **Table 3.** Test accuracy results of 10 independent training runs. Note that training runs 5 and 6
535 producing test accuracies identical to two decimal places occurred by chance and was verified
536 not to be a mistake.

Run	Test Accuracy
1	96.56%
2	96.04%
3	96.24%
4	95.88%
5	96.00%
6	96.00%
7	96.20%
8	96.08%
9	96.68%
10	96.64%

537

538

539

540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561

Figure Caption List

Figure 1. MASCRAD Snow Field Site at Easton Valley Airport, near Greeley, Colorado, under the umbrella of CSU-CHILL Radar. MASC (top right), along with other surface instrumentation, is contained in the 2/3-scaled DFIR.

Figure 2. Example normalized raw MASC image. Several snowflakes can be seen in addition to background glare (center left) and subtle ground and sky glow (top and bottom). Note: ground and sky glow may not be visible in all prints or computer monitor settings.

Figure 3. Example binary image produced by application of a brightness threshold and 5-pixel radius to the normalized raw image in Figure 2. Possible snowflake silhouettes are now apparent. Background glare (center left) was rejected due to exceeding the mean brightness threshold. Dimmer glare cases are reliably assigned to the Not-Flakes Q&R class.

Figure 4. Example crops and image chips extracted from the MASC image shown in Figures 2 and 3. (a) Cropped image of a planar crystal. (b) Example crop from (a) after contrast scaling. (c) Final image chip produced from contrast scaled crop in (b). (d) Cropped image of an aggregate. (e) Example crop from (d) after contrast scaling. (f) Final image chip produced from contrast scaled crop in (e).

Figure 5. Examples of image chips in the Not-Flakes quality and recognizability category. A description of this category is given in Table 1. First row (left to right): a coin; background glare;

562 sky glow seen between fence posts; a finger. Second row: a sensor probe; an out of focus sensor
563 probe; part of a pair of calipers; a string. Third row: a metal ball; part of a mitten; background
564 glare; amplified sensor noise. Fourth row: background glare; sky glow seen above fence posts;
565 background glare; background glare.

566
567 **Figure 6.** Examples of Bad-Crop (first column), Bad (second column), Okay (third column), and
568 Good (fourth column) image chips. Category descriptions given in Table 1.

569
570 **Figure 7.** Examples of image chips in the aggregate (AG) class of the final geometric dataset.
571 All image chips in the final geometric dataset had been automatically categorized into the Good
572 Q&R category. We placed emphasis on collecting a wide variety of sizes and forms of aggregate
573 with varying types of constituent particles.

574
575 **Figure 8.** Examples of image chips in the columnar crystal (CC) class of the final geometric
576 dataset. All image chips in the final geometric dataset had been automatically categorized into
577 the Good Q&R category. We included a variety of sizes, forms, and degrees of riming. An
578 example of a backlit snowflake is shown in row 2, column 2. Such cases were rare but were
579 included whenever backlighting did not interfere with recognizability.

580
581 **Figure 9.** Examples of image chips in the graupel (GR) class of the final geometric dataset. All
582 image chips in the final geometric dataset had been automatically categorized into the Good
583 Q&R category. We included a variety of textures and sizes.

584 **Figure 10.** Examples of image chips in the planar crystal (PC) class of the final geometric
585 dataset. All image chips in the final geometric dataset had been automatically categorized into
586 the Good Q&R category. We included difficult examples like row 1 column 2 where possible to
587 help differentiate such PC cases from CC examples.

588
589 **Figure 11.** Examples of image chips in the small particle (SP) class of the final geometric
590 dataset. All image chips in the final geometric dataset had been automatically categorized into
591 the Good Q&R category. As small particles are, by definition, particles with features too small to
592 classify, there is little interesting variety among the collected examples other than various shapes
593 and degrees of riming.

594
595 **Figure 12.** Training progress for an example training run using the methods and hyperparameters
596 described in Section 4.2

597
598 **Figure 13.** Confusion matrix for the network trained in Figure 12 applied to the test set. A final
599 accuracy of 96.2% was achieved. AG and PC were the most confused classes.

600
601 **Figure 14.** Examples of image chips misclassified by a trained network. Misclassified aggregates
602 (first row), misclassified columnar crystals (second row), misclassified graupel (third row),
603 misclassified planar crystals (fourth row), and misclassified small particles (fifth row) are shown
604 with the label assigned by the network overlaid for each image chip.

605

Figures

606



607

608 **Figure 1.** MASCRAD Snow Field Site at Easton Valley Airport, near Greeley, Colorado, under
609 the umbrella of CSU-CHILL Radar. MASC (top right), along with other surface instrumentation,
610 is contained in the 2/3-scaled DFIR.

611



612

613 **Figure 2.** Example normalized raw MASC image. Several snowflakes can be seen in addition to
614 background glare (center left) and subtle ground and sky glow (top and bottom). Note: ground
615 and sky glow may not be visible in all prints or computer monitor settings.

616



617

618 **Figure 3.** Example binary image produced by application of a brightness threshold and 5-pixel

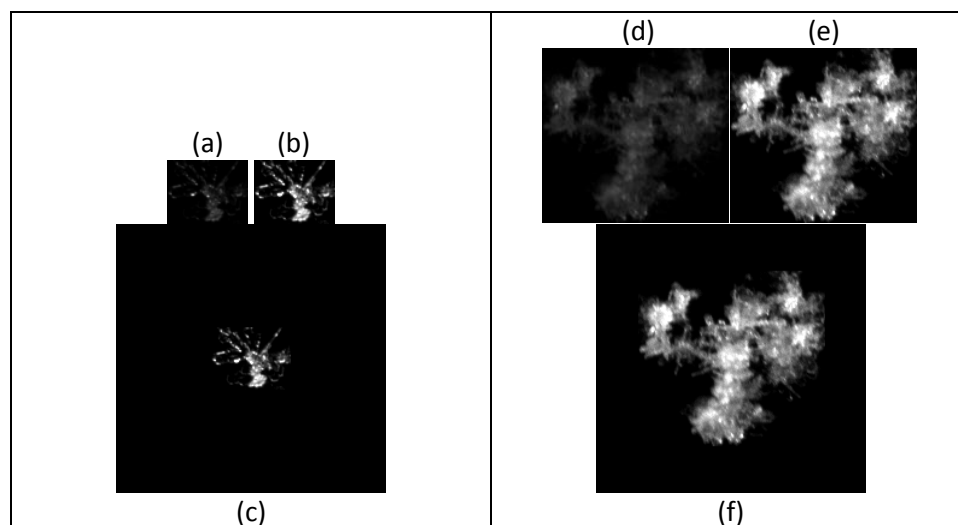
619 radius to the normalized raw image in Figure 2. Possible snowflake silhouettes are now apparent.

620 Background glare (center left) was rejected due to exceeding the mean brightness threshold.

621 Dimmer glare cases are reliably assigned to the Not-Flakes Q&R category.

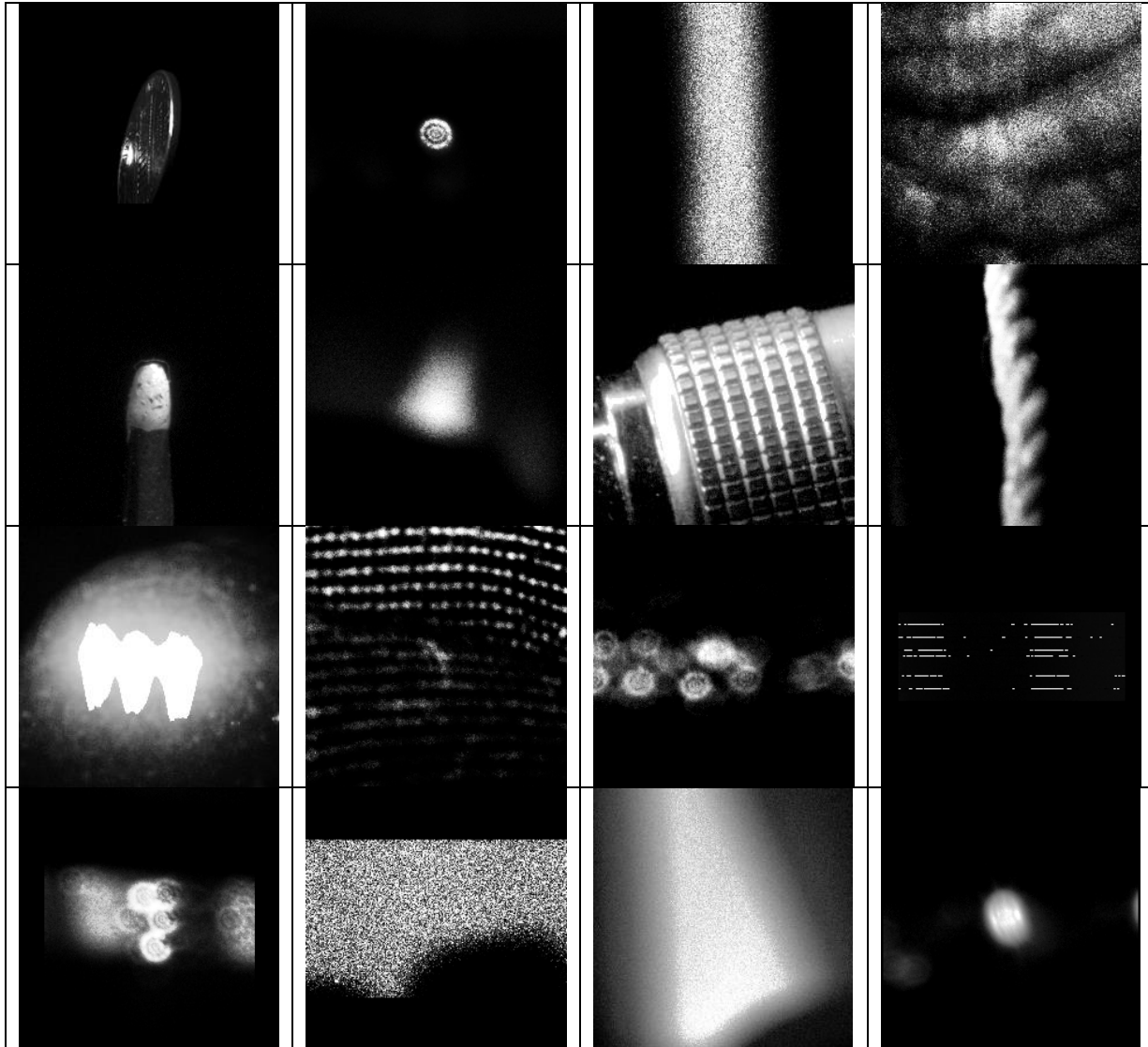
622

623

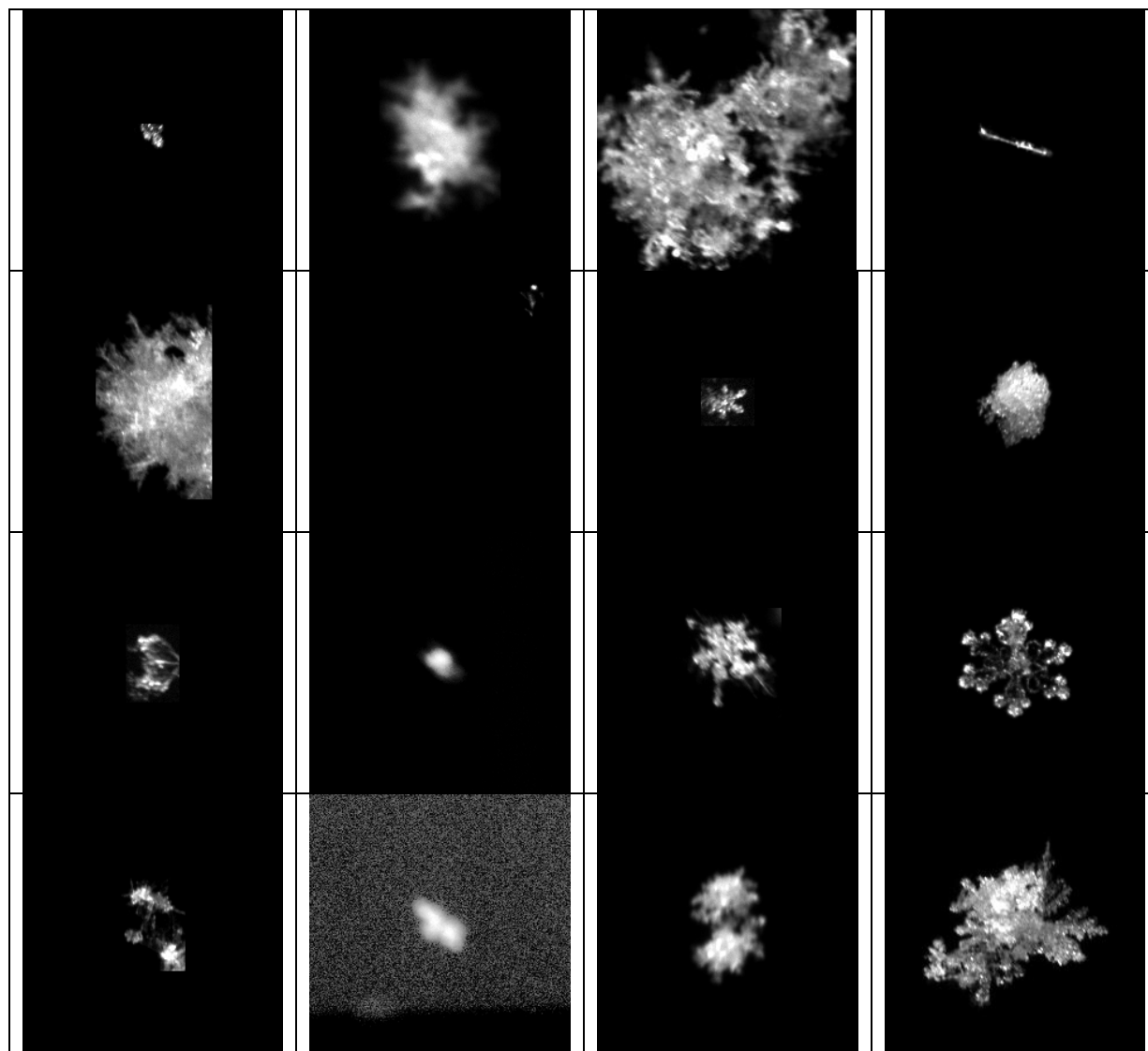


624 **Figure 4.** Example crops and image chips extracted from the MASC image shown in Figures 2
625 and 3. (a) Cropped image of a planar crystal. (b) Example crop from (a) after contrast scaling. (c)
626 Final image chip produced from contrast scaled crop in (b). (d) Cropped image of an aggregate.
627 (e) Example crop from (d) after contrast scaling. (f) Final image chip produced from contrast
628 scaled crop in (e).

629
630
631



632 **Figure 5.** Examples of image chips in the Not-Flakes quality and recognizability category. A
633 description of this category is given in Table 1. First row (left to right): a coin; background glare;
634 sky glow seen between fence posts; a finger. Second row: a sensor probe; an out of focus sensor
635 probe; part of a pair of calipers; a string. Third row: a metal ball; part of a mitten; background
636 glare; amplified sensor noise. Fourth row: background glare; sky glow seen above fence posts;
637 background glare; background glare.

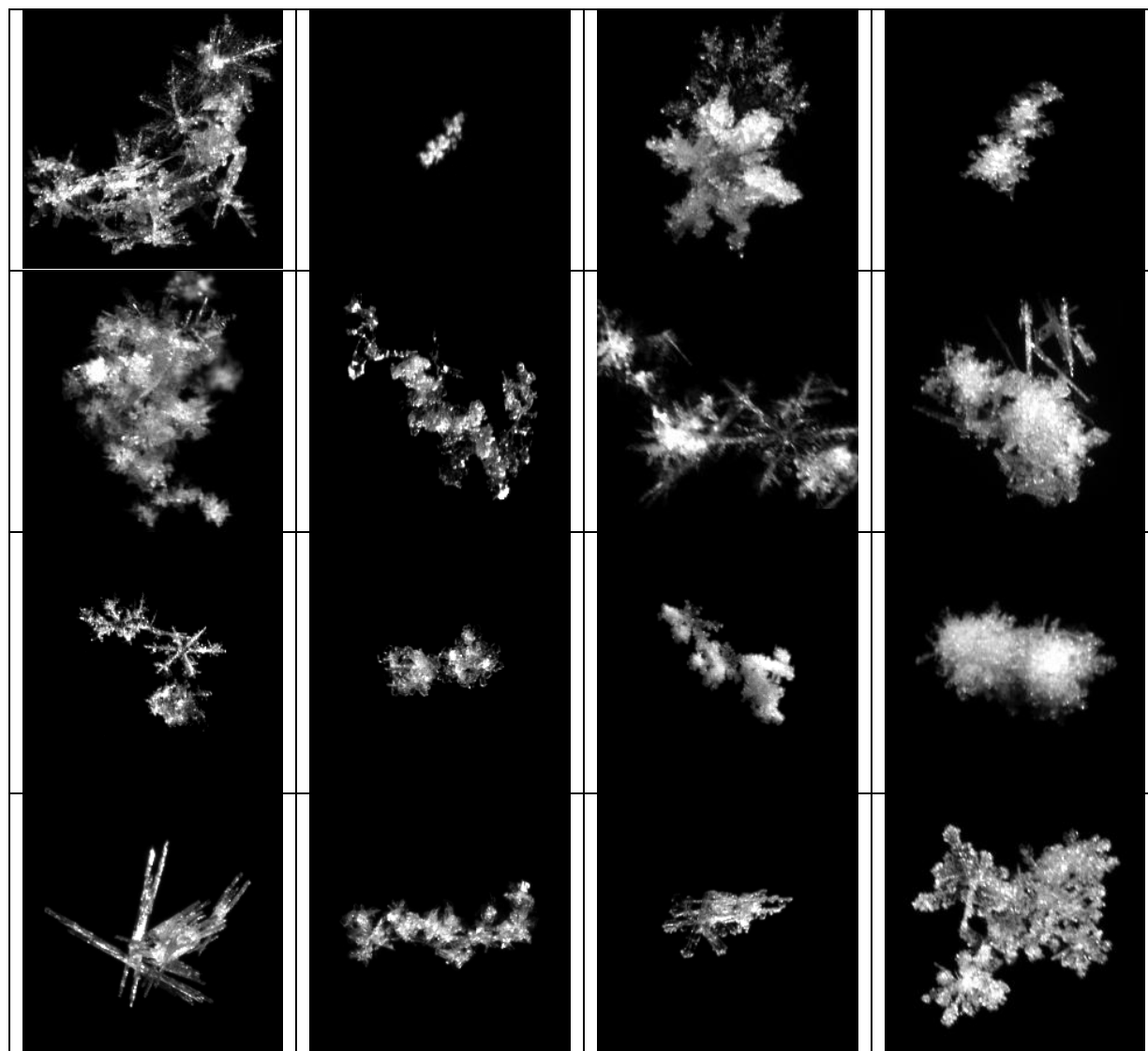


638 **Figure 6.** Examples of Bad-Crop (first column), Bad (second column), Okay (third column), and
639 Good (fourth column) image chips. Category descriptions given in Table 1.

640

641

642



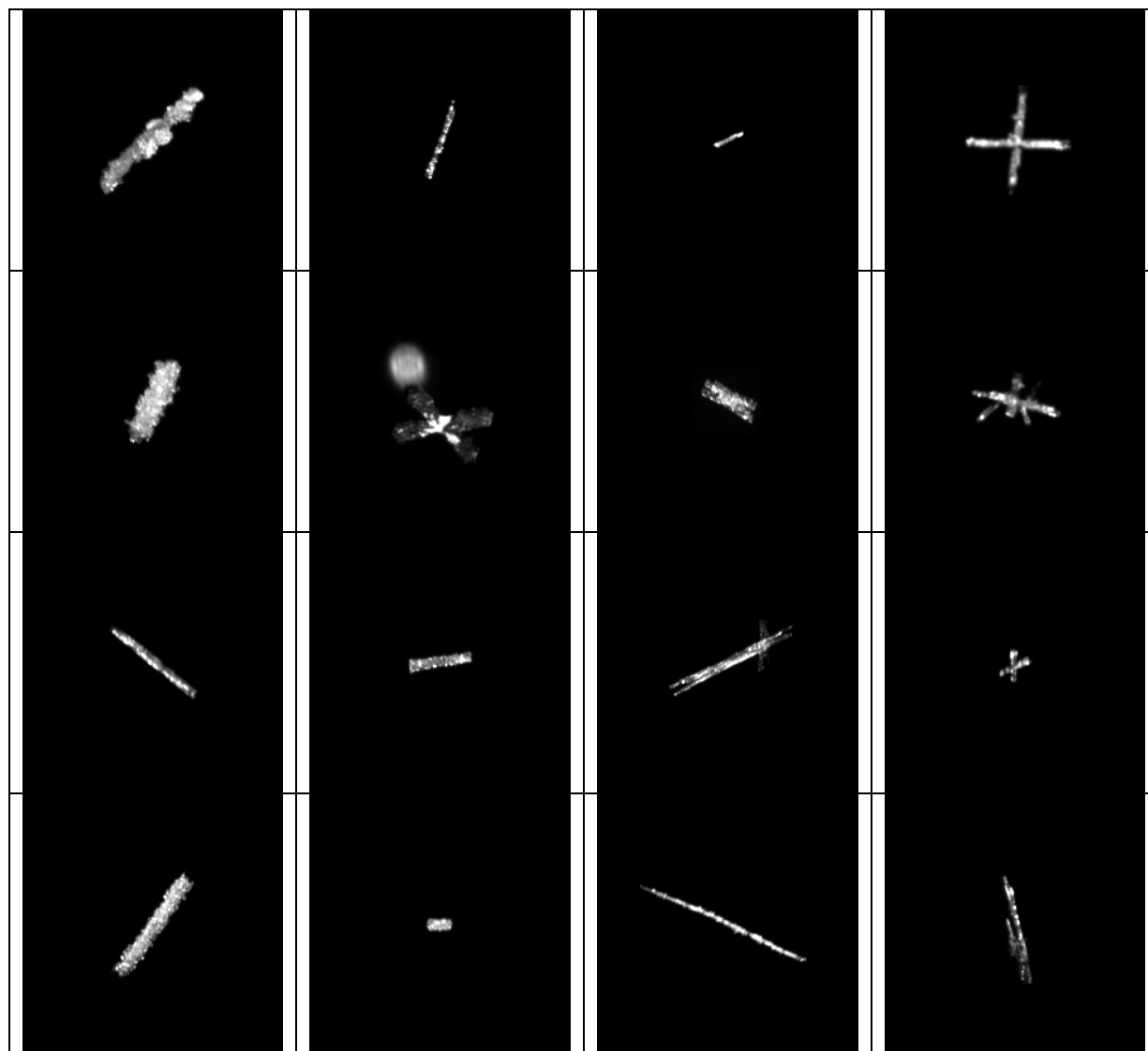
643 **Figure 7.** Examples of image chips in the aggregate (AG) class of the final geometric dataset.

644 All image chips in the final geometric dataset had been automatically categorized into the Good

645 Q&R category. We placed emphasis on collecting a wide variety of sizes and forms of aggregate

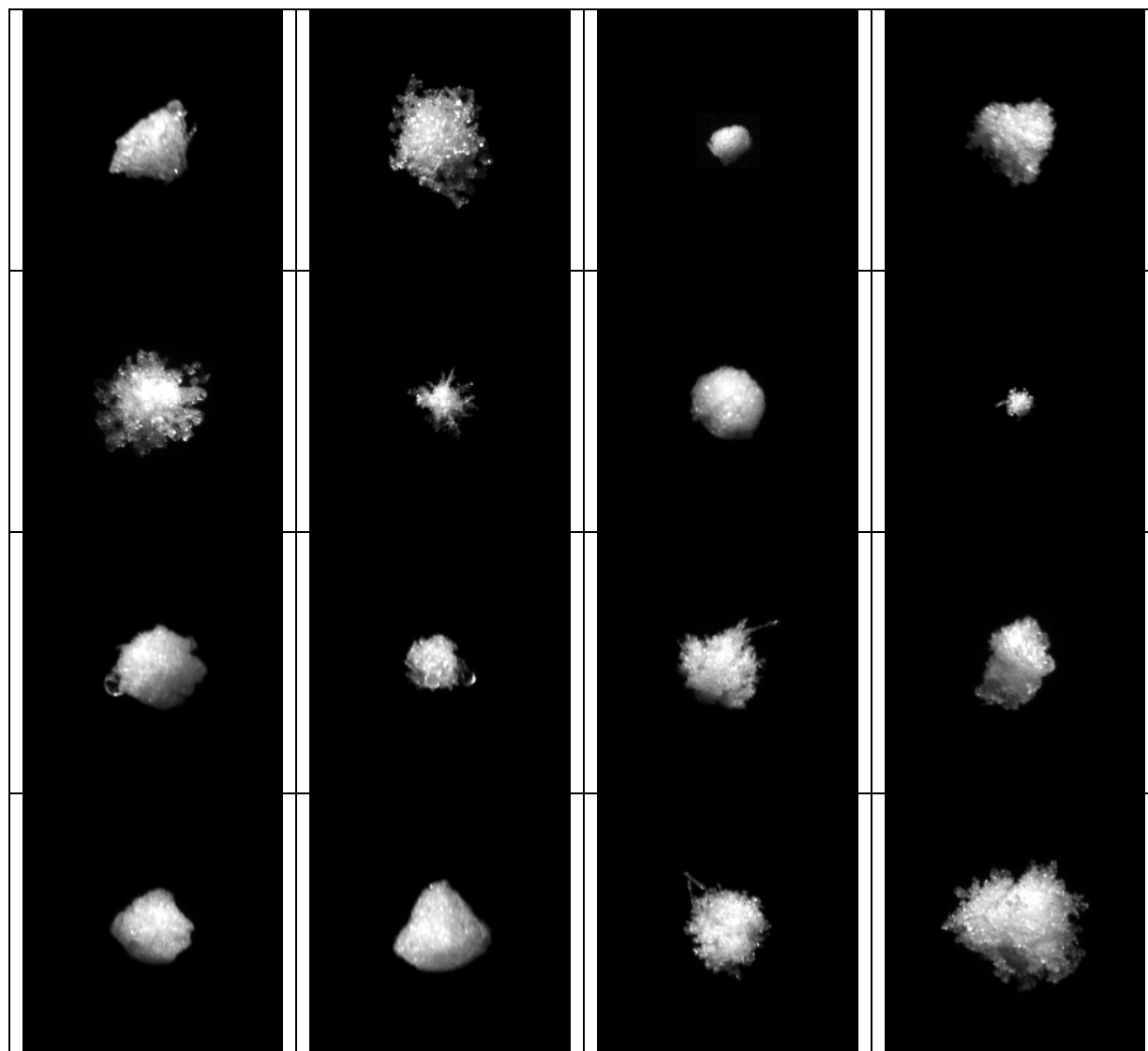
646 with varying types of constituent particles.

647



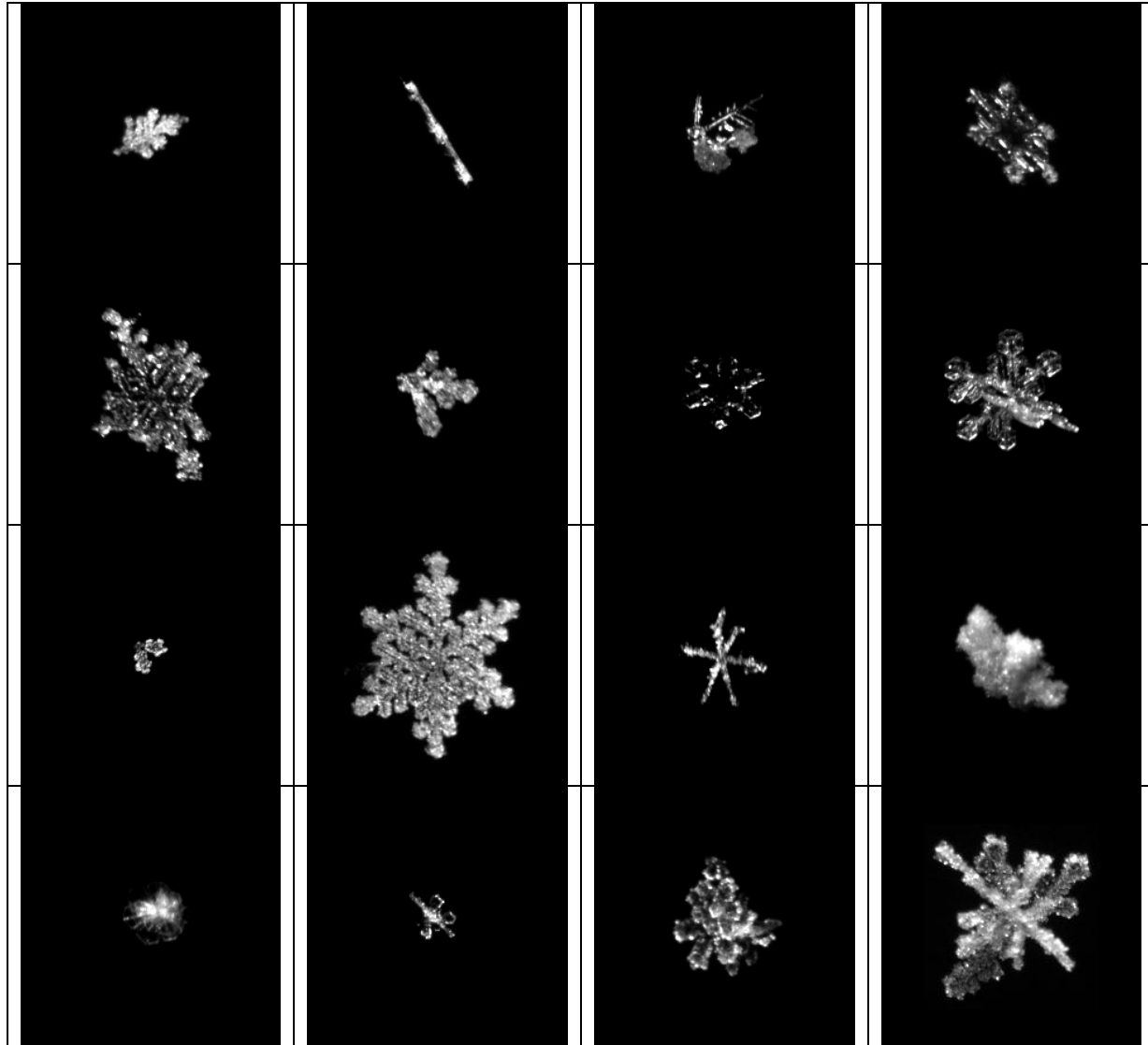
648 **Figure 8.** Examples of image chips in the columnar crystal (CC) class of the final geometric
649 dataset. All image chips in the final geometric dataset had been automatically categorized into
650 the Good Q&R category. We included a variety of sizes, forms, and degrees of riming. An
651 example of a backlit snowflake is shown in row 2, column 2. Such cases were rare but were
652 included whenever backlighting did not interfere with recognizability.

653



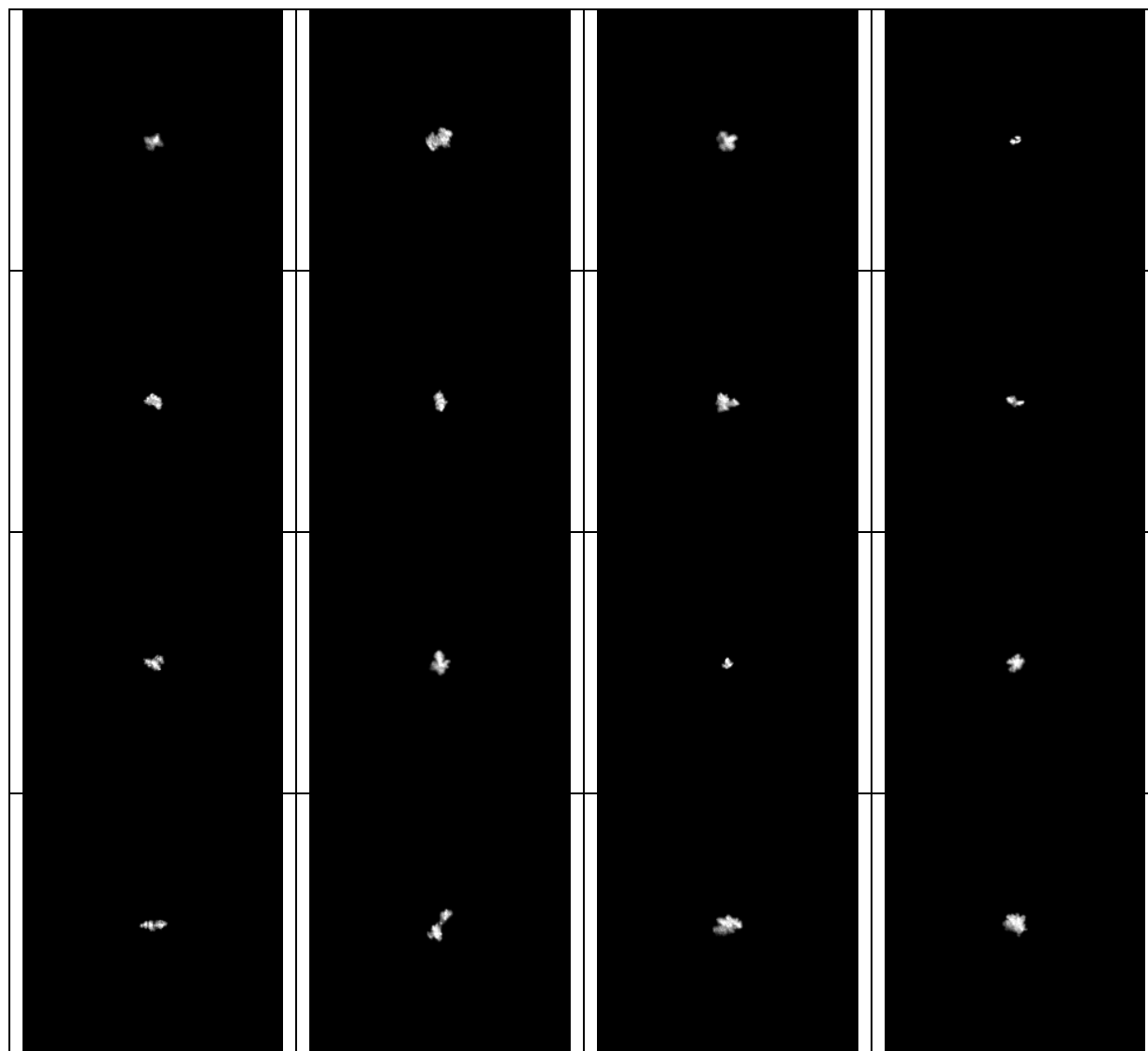
654 **Figure 9.** Examples of image chips in the graupel (GR) class of the final geometric dataset. All
655 image chips in the final geometric dataset had been automatically categorized into the Good
656 Q&R category. We included a variety of textures and sizes and also included melting examples
657 when available.

658



659 **Figure 10.** Examples of image chips in the planar crystal (PC) class of the final geometric
660 dataset. All image chips in the final geometric dataset had been automatically categorized into
661 the Good Q&R category. We included difficult examples like row 1 column 2 where possible to
662 help differentiate such PC cases from CC examples. Emphasis was also placed on including
663 examples that lacked easily identifiable six-fold symmetry.

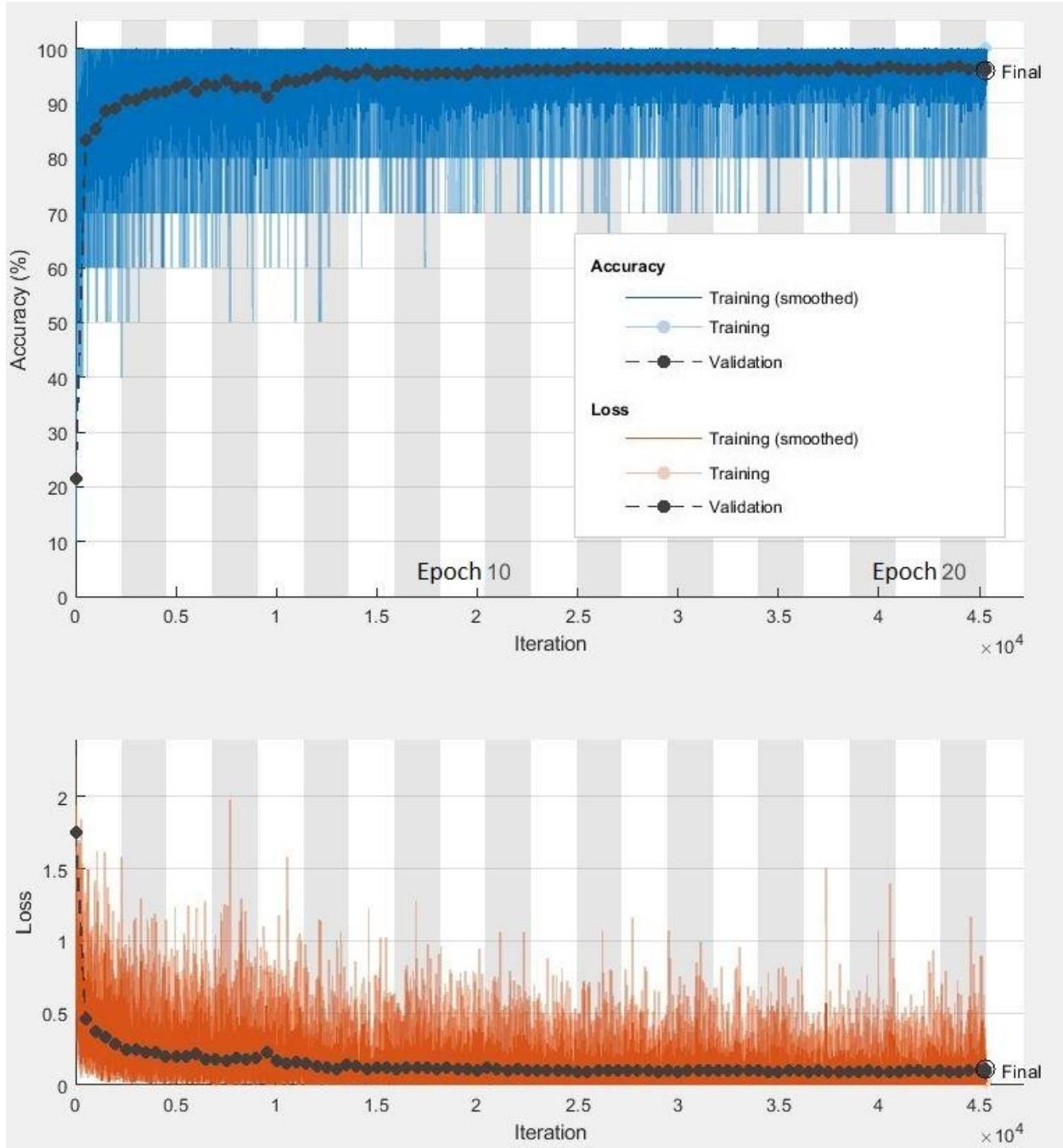
664



665 **Figure 11.** Examples of image chips in the small particle (SP) class of the final geometric
666 dataset. All image chips in the final geometric dataset had been automatically categorized into
667 the Good Q&R category. As small particles are, by definition, particles with features too small to
668 classify, there is little interesting variety among the collected examples other than various shapes
669 and degrees of riming.

670

671



672

673 **Figure 12.** Training progress for an example training run using the methods and hyperparameters

674 described in Section 4.2

675

676

Confusion Matrix

Output Class	AG	471 18.8%	2 0.1%	1 0.0%	14 0.6%	1 0.0%	96.3% 3.7%
	CC	5 0.2%	482 19.3%	0 0.0%	4 0.2%	8 0.3%	96.6% 3.4%
	GR	3 0.1%	0 0.0%	491 19.6%	1 0.0%	1 0.0%	99.0% 1.0%
	PC	20 0.8%	3 0.1%	5 0.2%	478 19.1%	8 0.3%	93.0% 7.0%
	SP	1 0.0%	13 0.5%	3 0.1%	3 0.1%	482 19.3%	96.0% 4.0%
			94.2% 5.8%	96.4% 3.6%	98.2% 1.8%	95.6% 4.4%	96.4% 3.6%
		AG	CC	GR	PC	SP	
		Target Class					

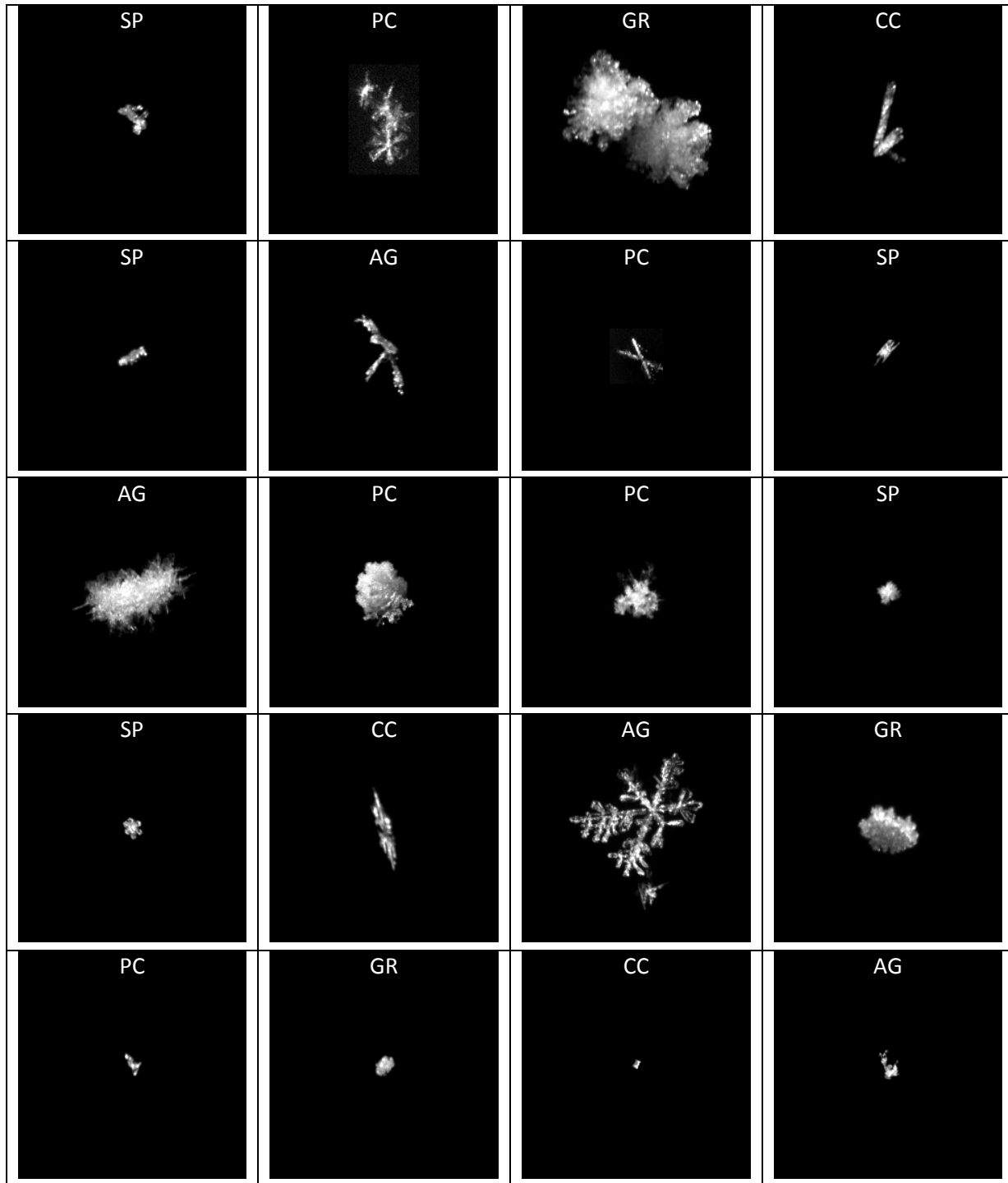
677

678 **Figure 13.** Confusion matrix for the network trained in Figure 12 applied to the test set. Each red
 679 or green cell corresponds to a target class (horizontal) and output class (vertical). Row 2, column
 680 1, for instance, shows that 5 image chips in the test set with target class AG were assigned to the
 681 CC class by the trained network, and this corresponded to 0.2% of the entire dataset. The first
 682 five cells of the bottom row show accuracy (green) and error (red) for each target class. Row 6,
 683 column 1, for instance, shows that, of image chips in the test set with target class AG, 94.2%
 684 were classified correctly by the network while 5.8% were classified incorrectly. The first five

685 cells of the rightmost column similarly show accuracy and error for each output class. Row 1,
686 column 6, for instance, shows that, of image chips assigned by the network to the AG class,
687 96.3% were classified correctly while 3.7% were classified incorrectly. An overall network
688 accuracy (all classes) of 96.2% is shown in the bottom right cell. AG and PC were the most
689 confused classes.

690

691



692 **Figure 14.** Examples of image chips misclassified by a trained network. Misclassified aggregates
693 (first row), misclassified columnar crystals (second row), misclassified graupel (third row),

694 misclassified planar crystals (fourth row), and misclassified small particles (fifth row) are shown
695 with the label assigned by the network overlaid for each image chip.

696