

## Supplementary Information

# Classifying optical microscope images of exfoliated graphene flakes by data-driven machine learning

Satoru Masubuchi<sup>1,\*</sup> and Tomoki Machida<sup>1,\*</sup>

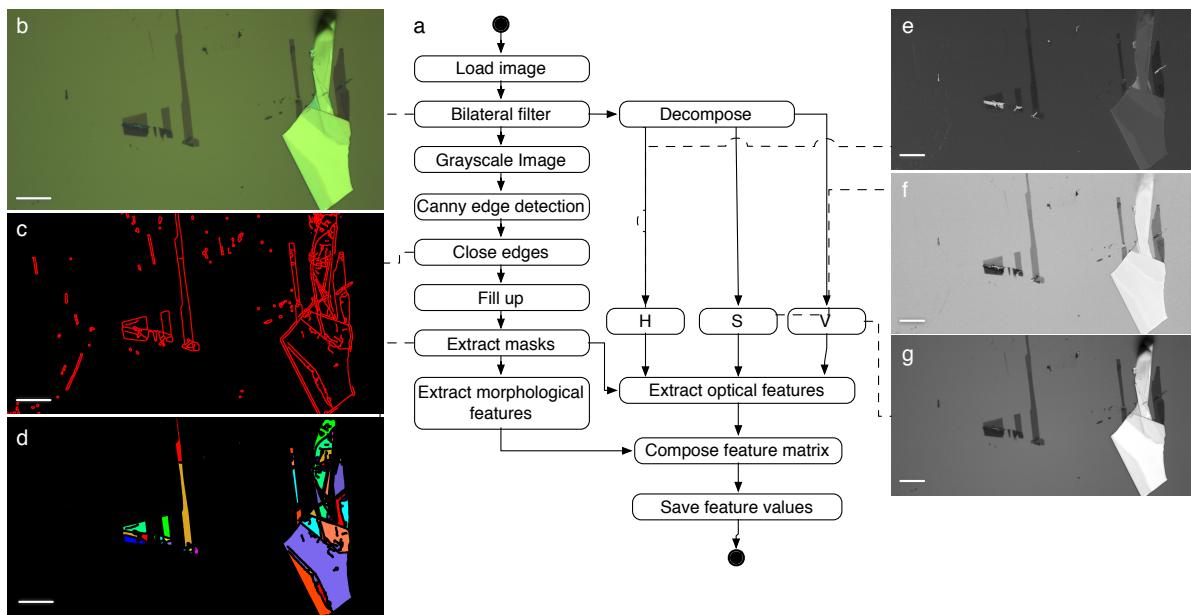
<sup>1</sup>Institute of Industrial Science, University of Tokyo, 4-6-1 Komaba, Meguro-ku, Tokyo 153-8505 Japan

\*Correspondence: [msatoru@iis.u-tokyo.ac.jp](mailto:msatoru@iis.u-tokyo.ac.jp), [tmachida@iis.u-tokyo.ac.jp](mailto:tmachida@iis.u-tokyo.ac.jp)

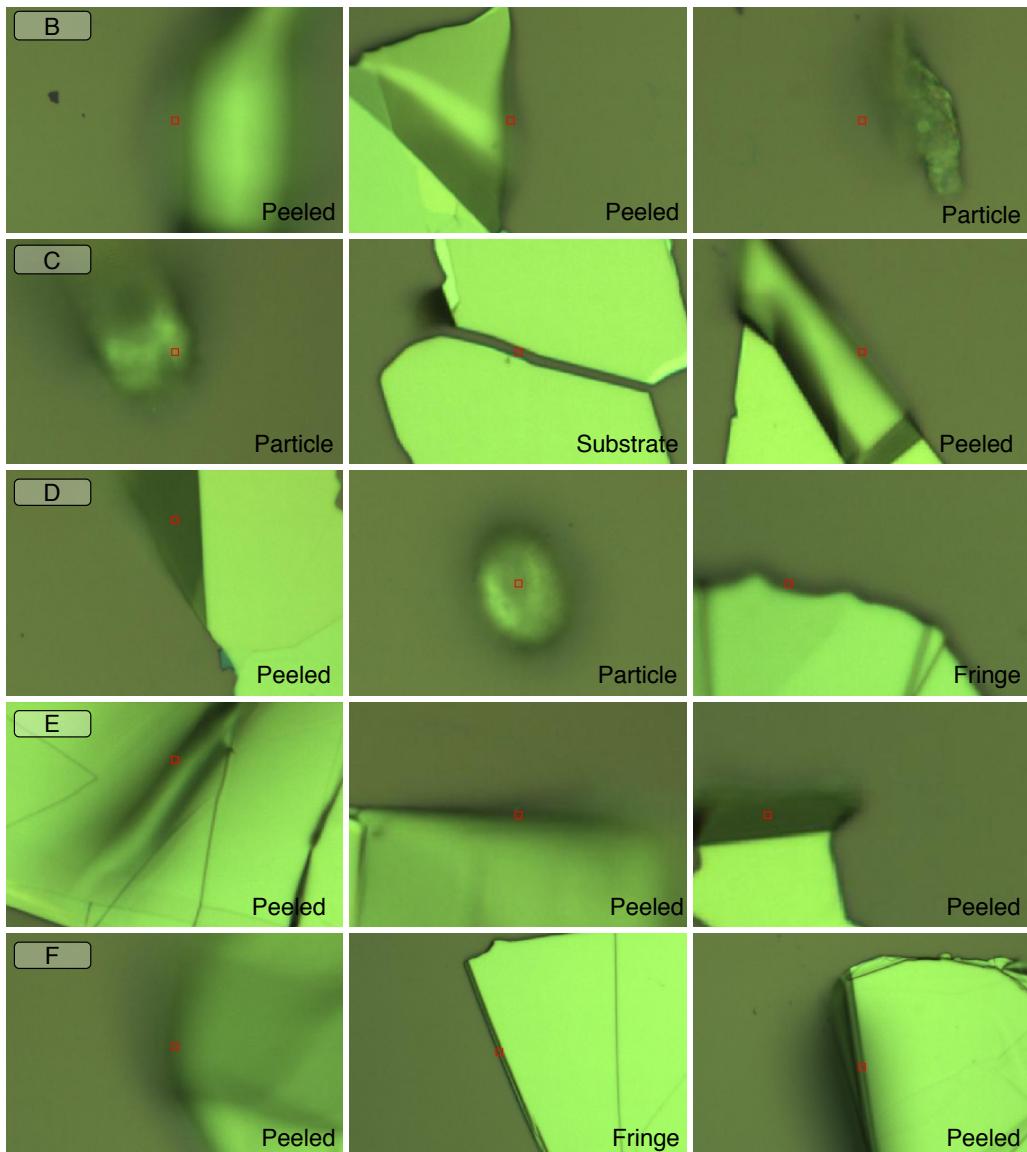
8

9

10 Supplementary Figures



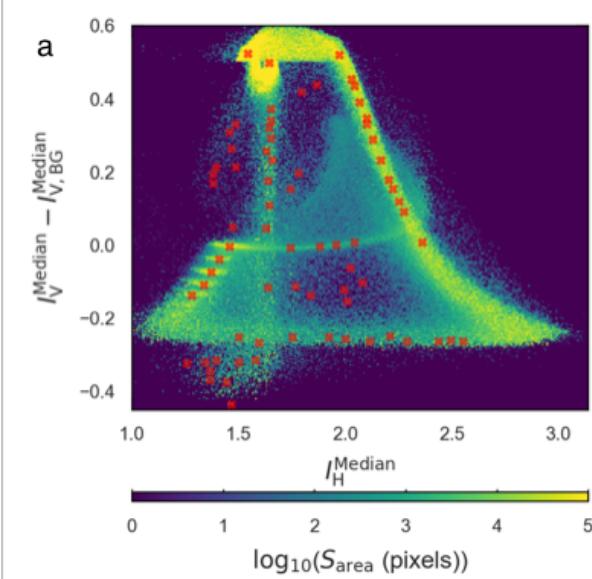
11  
12 **Supplementary Figure 1. Image processing pipeline to extract morphological and optical**  
13 **features. (a)** Schematic of image processing algorithm for extracting features of  
14 **atomically thin 2D crystals from the optical microscope images. (b)** Optical microscope  
15 **image of exfoliated graphene flakes on SiO<sub>2</sub>/Si substrate, where the scale bar represents**  
16 **a length of 20 μm. The image includes bilayer graphene, trilayer graphene, and thick**  
17 **graphite flakes. (c)-(g)** Representative images generated at each step. The correspondence  
18 **between images and blocks are indicated by the dashed lines. The scale bar represents a**  
19 **length of 20 μm.**



20

21 **Supplementary Figure 2.** Typical examples of optical microscope images that lead to  
22 misclassification. For all the clusters B-F, the graphite flakes peeled from the substrate  
23 (Peeled), fringes of graphite (Fringe), the silicon substrate surrounded by the graphite  
24 flakes (Substrate), and particles (Particle) are the main sources of the misclassification.

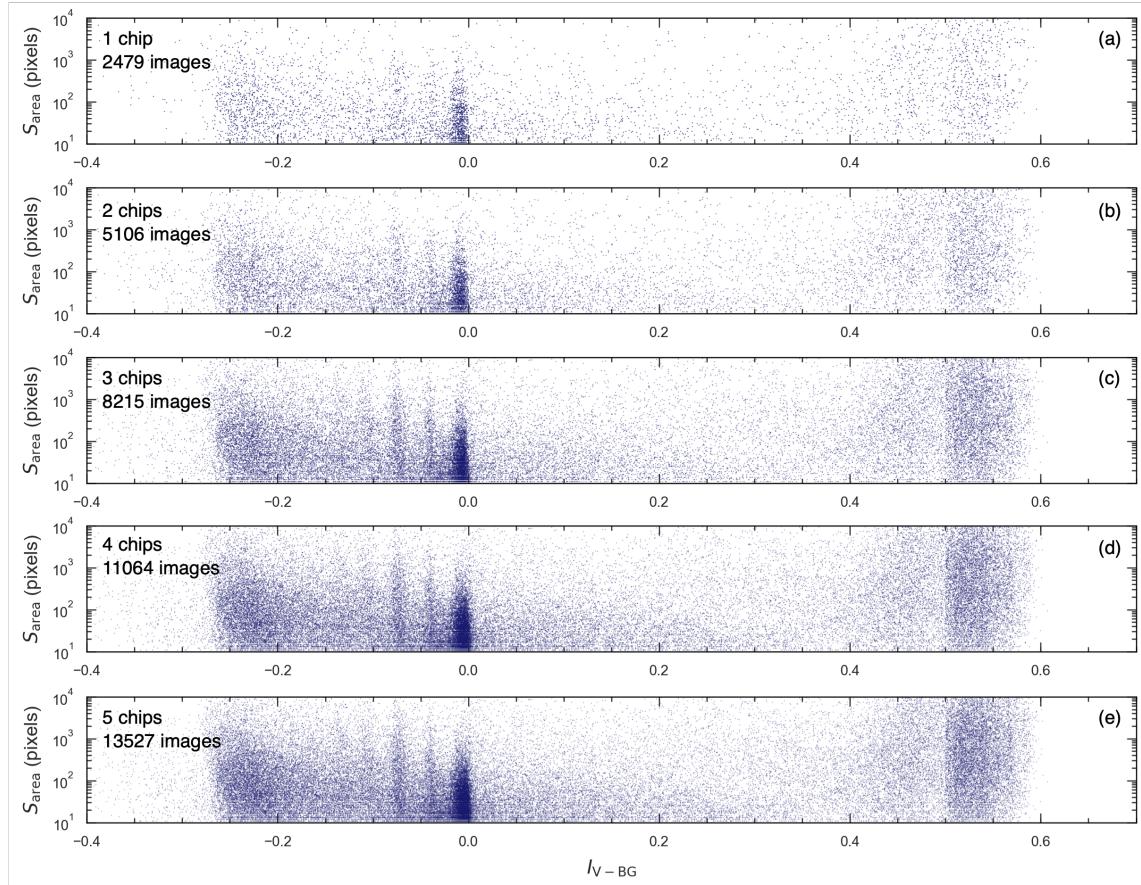
25



26 **Supplementary Figure 3.** The results of peak detection applied to the area histogram  $S_{\text{area}}$ .

27 **(a)** Cumulative area histogram as a function of  $(I_H, I_{V-\text{BG}})$ . The positions of the detected  
28 peaks are indicated by the red crosses.

29

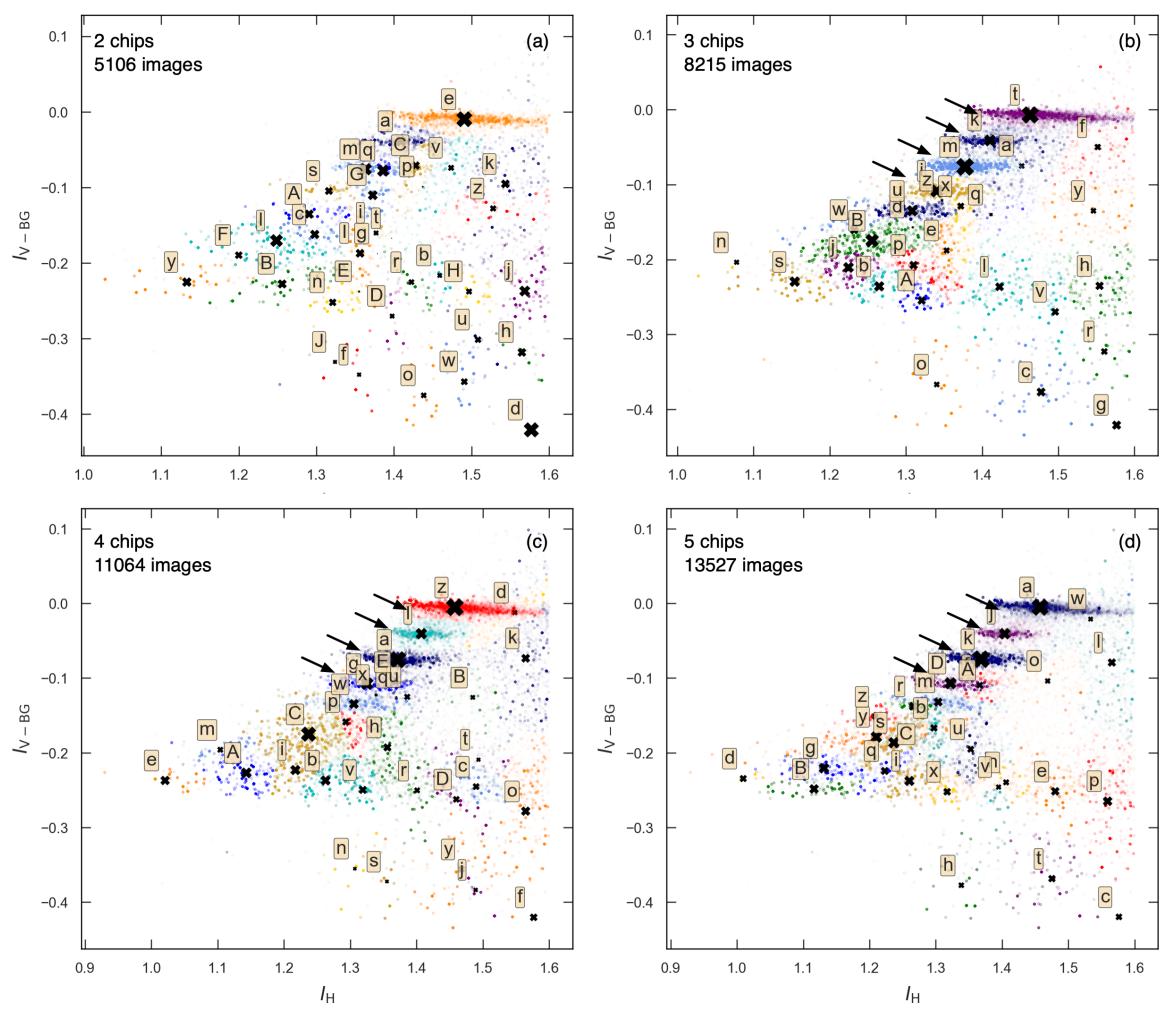


30

31 **Supplementary Figure 4. Scatter plots of the feature values for varying number of**  
 32 **analyzed images.** The scatter plots of  $S_{\text{area}}$  vs.  $I_{V-\text{BG}}$  obtained by analyzing (a) 1 silicon chip  
 33 and 2479 images, (b) 2 silicon chips and 5106 images, (c) 3 silicon chips and 8215 images, (d)  
 34 4 silicon chips and 11064 images, and (e) 5 silicon chips and 13527 images.

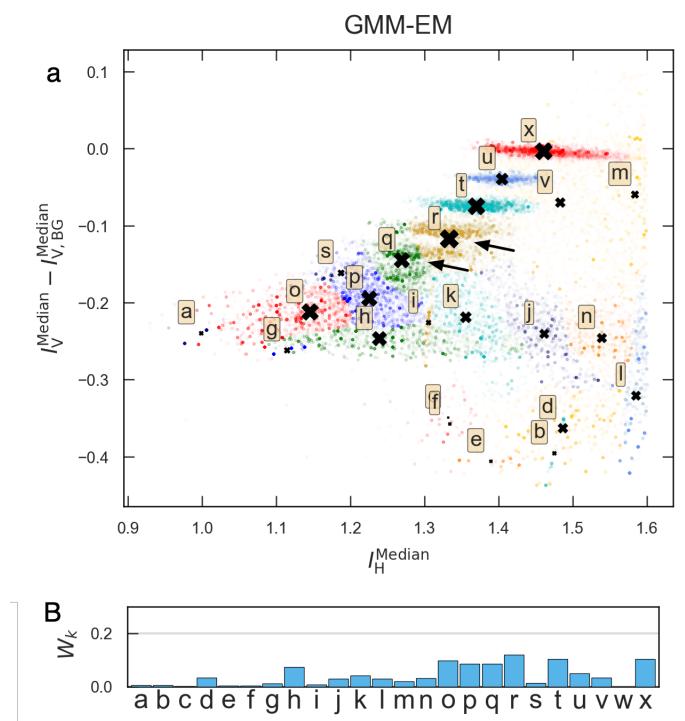
35

36



37  
 38 **Supplementary Figure 5. Clustering results of the Bayesian Gaussian mixture model with**  
 39 **the Dirichlet process for varying number of analyzed images.** (a) 2 silicon chips and 5106  
 40 images, (b) 3 silicon chips and 8215 images, (c) 4 silicon chips and 11064 images, and (d) 5  
 41 silicon chips and 13527 images.  
 42

43



44

45 **Supplementary Figure 6. Clustering results of Gaussian mixture model with expectation  
46 maximization (GMM-EM).** (a) The cluster centers are indicated by filled crosses. The  
47 weights of each cluster are indicated by the sizes of cross marks. (b) The extracted weights for  
48 the clusters of GMM-EM.

49

50 Supplementary tables

Morphology features	Optical features
• Circularity	• BG_mean_(h, s, v, g)
• Compactness	• BG_median_(h, s, v, g)
• Contlength	• Mean_(h, s, v, g)
• Convexity	• Median_(h, s, v, g)
• Rectangularity	• Entropy
• Anisometry	• Anisotropy
• Bulkiness	
• Struct Factor	
• Outer Radius	
• Inner Radius	
• Inner Height	
• Inner Width	
• Dist Mean	
• Dist Deviation	
• Roundness	
• Num Sides	
• Connect Num	
• Holes Num	
• Area Holes	
• Max Diameter	
• Orientation	
• Euler Number	
• Area	
• Column	
• Height	
• Row	
• Width	

51

**Basic information**

- Filenames

52 **Supplementary Table 1. Feature values extracted by the presented algorithm.**

53

54      Supplementary notes

55      Supplementary note 1

56            Here we discuss the minimum number of optical microscope images required to make  
57          good clustering results. In Supplementary Figure 4, we show the scatter plots of the feature  
58          values  $S_{\text{area}}$  vs.  $I_{V-\text{BG}}$  obtained by analyzing the differing number of the silicon chips thereby  
59          that of the optical microscope images. When the analysis was conducted for one silicon chip  
60          and 2479 images (Supplementary Figure 4 (a)), no vertical stripe patterns were emerged. When  
61          the number of silicon chips ( $n$ ) and optical microscope images ( $m$ ) were increased as  
62           $(n, m) = (\text{b})(2, 5106), (\text{c})(3, 8215), (\text{d})(4, 11064)$ , and  $(\text{e})(5, 13527)$ , the vertical stripe  
63          patterns became gradually discernible. From these observations, we could state that around  
64          5000 – 10000 optical microscope images are required for obtaining the vertical stripe patterns  
65          deriving from monolayer-, bilayer-, and trilayer graphene flake as discussed in the main text.

66            When the clustering analysis were conducted using the BGMM-DP, the successful  
67          clustering started from  $(n, m) = (3, 8215)$  [Black arrows in Supplementary Figure 5 (b)],  
68          and the clustering in small number of optical microscope images  $(n, m) = (2, 5106)$  tended  
69          to fall into misclassification. From these observations, we can estimate that at least 10000  
70          optical microscope images are needed to obtain the successful unsupervised clustering results.

71

72      Supplementary note 2

73      **Gaussian mixture model with expectation maximization**

74            We describe the results of the clustering analysis using the Gaussian mixture model  
75          with expectation maximization optimization. In this model, one needs to specify the number  
76          of clusters  $K$  and the initial guess of cluster centers for performing the optimization. Here, to  
77          estimate the cluster numbers and the centers, we applied a local peak detection algorithm to  
78          the cumulative histogram of the area  $S_{\text{area}}$  of the regions [Red crosses in Supplementary Figure  
79          3]. The detected peak positions were supplied to the EM algorithm and optimization iteration

80 was performed. The results are presented in Supplementary Figure 6(a). Compared to the  
81 GMM-DP, the clusters consisting of trilayer-graphene and tetralayer-graphene flakes (clusters  
82 r and q, respectively, indicated by the black arrows in Supplementary Figure 6(b)) were  
83 confused with the different clusters. Even when we changed the initialization conditions, the  
84 GMM-EM results did not exhibit the correct clustering results. These results indicate that when  
85 applied to our feature matrix, the GMM-EM algorithm tend to result in overfitting; therefore,  
86 GMM-DP gives a better performance compared to GMM-EM.