

Copyright  
by  
Duo Xu  
2018

The Thesis Committee for Duo Xu  
Certifies that this is the approved version of the following Thesis:

**Assessing the Performance of a Machine Learning  
Algorithm in Identifying Bubbles in Dust Emission**

APPROVED BY

SUPERVISING COMMITTEE:

Stella Offner, Supervisor

Volker Bromm

**Assessing the Performance of a Machine Learning  
Algorithm in Identifying Bubbles in Dust Emission**

by

**Duo Xu**

**THESIS**

Presented to the Faculty of the Graduate School of

The University of Texas at Austin

in Partial Fulfillment

of the Requirements

for the Degree of

**MASTER OF ARTS**

THE UNIVERSITY OF TEXAS AT AUSTIN

December 2018

## Acknowledgments

Reprinted with permission from “Assessing the Performance of a Machine Learning Algorithm in Identifying Bubbles in Dust Emission ” by Xu, D., & Offner, S. S. R., 2017, *Astrophysical Journal*, 851, 149 [49]. Copyright [2017] by American Astronomical Society. D.X. conducted the synthetic observations, trained the machine learning algorithm, carried out the analyses, produced the figures and wrote the paper. S.O. proposed the ideas to harness machine learning to study the stellar feedback bubbles, provided the hydrodynamic simulations, provided input to the analysis and contributed to the paper text.

# Assessing the Performance of a Machine Learning Algorithm in Identifying Bubbles in Dust Emission

Duo Xu, M.A.

The University of Texas at Austin, 2018

Supervisor: Stella Offner

Stellar feedback created by radiation and winds from massive stars plays a significant role in both physical and chemical evolution of molecular clouds. This energy and momentum leaves an identifiable signature (“bubbles”) that affect the dynamics and structure of the cloud. Most bubble searches are performed “by-eye”, which are usually time-consuming, subjective and difficult to calibrate. Automatic classifications based on machine learning make it possible to perform systematic, quantifiable and repeatable searches for bubbles. We employ a previously developed machine learning algorithm, *Brut*, and quantitatively evaluate its performance in identifying bubbles using synthetic dust observations. We adopt magneto-hydrodynamics simulations, which model stellar winds launching within turbulent molecular clouds, as an input to generate synthetic images. We use a publicly available three-dimensional dust continuum Monte-Carlo radiative transfer code, HYPERION, to generate synthetic images of bubbles in three Spitzer bands

(4.5  $\mu\text{m}$ , 8  $\mu\text{m}$  and 24  $\mu\text{m}$ ). We designate half of our synthetic bubbles as a training set, which we use to train *Brut* along with citizen-science data from the Milky Way Project. We then assess *Brut*'s accuracy using the remaining synthetic observations. We find that after retraining *Brut*'s performance increases significantly, and it is able to identify yellow bubbles, which are likely associated with B-type stars. *Brut* continues to perform well on previously identified high-score bubbles, and over 10% of the Milky Way Project bubbles are reclassified as high-confidence bubbles, which were previously marginal or ambiguous detections in the Milky Way Project data. We also investigate the size of the training set, dust model, evolution stage and background noise on bubble identification.

# Table of Contents

<b>Acknowledgments</b>	<b>iv</b>
<b>Abstract</b>	<b>v</b>
<b>List of Tables</b>	<b>viii</b>
<b>List of Figures</b>	<b>ix</b>
<b>Chapter 1. Introduction</b>	<b>1</b>
<b>Chapter 2. Methods</b>	<b>5</b>
2.1 Hydrodynamic Simulations . . . . .	5
2.2 Hyperion . . . . .	6
2.3 <i>Brut</i> . . . . .	11
<b>Chapter 3. Synthetic Observations</b>	<b>13</b>
3.1 Cropped Data . . . . .	13
3.2 Evolutionary Stage . . . . .	15
3.3 Turbulent Realization . . . . .	15
3.4 Noise . . . . .	16
<b>Chapter 4. Results</b>	<b>18</b>
4.1 Retraining <i>Brut</i> with Synthetic Observation . . . . .	18
4.2 Re-Testing <i>Brut</i> on the Milky Way Project Data . . . . .	22
4.3 Application: Bubbles in the Perseus Molecular Cloud . . . . .	36
<b>Chapter 5. Summary</b>	<b>42</b>
5.1 Conclusions . . . . .	42
5.2 Future Work . . . . .	43
<b>Bibliography</b>	<b>46</b>

## List of Tables

2.1	Physical Parameters of the Stellar Sources . . . . .	6
2.2	Model Properties <sup>a</sup> . . . . .	7
2.3	Random Forest Zone . . . . .	12
4.1	Parameters of the Synthetic Images . . . . .	20
4.2	Parameters of the Random Forests . . . . .	20
4.3	Physical Properties of Four Perseus Bubbles . . . . .	40

## List of Figures

2.1	Three-color synthetic images of five sources with 20 different viewing angles adopting the dust model in Koepferl et al. [25]. Red, green and blue represents $24 \mu\text{m}$ , $8 \mu\text{m}$ and $4.5 \mu\text{m}$ emission, respectively. . . . .	9
2.2	Same as Figure 2.1 but adopting the kmh dust model in Kim et al. [24]. . . . .	9
2.3	The SEDs of different dust models compared with the spectra of the Ophiuchus cloud LDN 1688 observed by Rawlings et al. [39]. . . . .	10
3.1	Three-color synthetic images adopting the K16 dust model where the HYPERION input is cropped to 2.2 pc. . . . .	14
3.2	Three-color synthetic images adopting the K16 dust model where the HYPERION input is cropped to 3 pc. . . . .	14
3.3	Three-color synthetic images at an earlier evolution stage with 0.05 Myr (“T2_t0” listed in Table 2.2) where the HYPERION input is cropped to 3 pc. . . . .	15
3.4	Three-color synthetic images with different turbulence where the HYPERION input is cropped to 3 pc. . . . .	16
3.5	Three-color synthetic images with noise. From top to bottom: Figure 2.1, 3.1 and 3.2 with noise added. . . . .	17
4.1	The cumulative distribution function (CDF) of all the scores given by Brut with the original training. The labels are described in Table 4.2. . . . .	23
4.2	The cumulative distribution function (CDF) of all the scores given by Brut retrained on noiseless synthetic images. The labels are described in Table 4.2. . . . .	24
4.3	The CDF of all the scores given by Brut with the original training. The labels are described in Table 4.2. . . . .	25
4.4	The CDF of all the scores given by Brut retrained on synthetic images with and without noise. The labels are described in Table 4.2. . . . .	26

4.5	The CDF of all the scores given by the retrained algorithm without original MWP training set but with only synthetic images. The labels are described in Table 4.2. . . . .	27
4.6	The CDF of all the scores given by the retrained algorithm on half original MWP training set and the synthetic images. The labels are described in Table 4.2. . . . .	28
4.7	The CDF of the scores of 3716 Milky Way Project large bubbles given by algorithm with the original training and the algorithm retrained on several different training sets including synthetic bubbles. . . . .	31
4.8	Distribution of image properties for four different training sets, where the colors indicate the number of bubbles in each bin. The white dashed lines divides the bubbles into four regions. The upper left quadrant indicates images that contain low S/N yellow bubbles. The upper right quadrant indicates images with high S/N yellow bubbles. The lower left quadrant indicates low S/N red bubbles. The lower right quadrant indicates high S/N red bubbles. . . . .	32
4.9	The average hit rate versus the average binned <i>Brut</i> score for the 3716 MWP large bubbles. The red and green lines indicate the average <i>Brut</i> score returned with the original training and the algorithm retrained on synthetic images both with and without noise, respectively. The error bars indicate the standard deviation of the scores and hit rate in each bin. The label indicates the number of bubbles in each bin. . . . .	34
4.10	The distribution of bubble scores returned with the original training and after <i>Brut</i> is retrained. . . . .	35
4.11	One hundred bubbles from MWP. The first number in the title of each panel presents the raw score, which is returned by the original MWP training algorithm. The middle number is the change in score after <i>Brut</i> is retrained with synthetic observations. The last number is the hit rate. . . . .	37
4.12	Nine bubbles, which were previously likely MWP bubbles but are no longer classified as high-confidence bubbles after retraining. The meaning of each number in the title is described in Figure 4.11. . . . .	38
4.13	Four examples of bubbles in the Perseus molecular cloud. The upper right label in each panel corresponds to the bubble name in Arce et al. [2]. The left number in the title of each panel indicates the raw score, which is returned by the original training algorithm. The right number in the title of each panel indicates the new score returned by the retrained algorithm. . . . .	41

# Chapter 1

## Introduction

During the process of star formation, stellar feedback plays a significant role in both physical and chemical evolution of molecular clouds [19, 17]. One of the most important feedback mechanisms is mass-loss [26]. There are two typical manifestations of stellar winds: protostellar outflows, which are often highly collimated, and radiatively driven winds from main sequence stars, which are more isotropic [8, 1, 2, 28]. Both type of stellar winds inject momentum and energy into the environment, and thereby affect the dynamics and structure of the parent molecular cloud.

Recent observational studies have shown that the momentum and energy injected by stellar winds are one or more orders of magnitude larger than those of outflows owing to their larger volume and longer lifetime [2, 28]. Arce et al. [2] found that the energy injection rate from these stellar winds is comparable to the turbulent dissipation rate in the Perseus molecular cloud, which means that in the current epoch, stellar feedback is sufficient to maintain the observed turbulence in Perseus. A similar conclusion was also reached by Li et al. [28] in the Taurus molecular cloud. It is notable that both regions are low-mass star forming regions, and high-mass stars, which generally dominate

feedback energetics are absent.

Simulations confirm the significant kinematic impact due to stellar feedback on the global star formation process. Winds can replenish energy dissipated by turbulence and also trigger star formation by compressing the cloud [9, 30, 10, 11, 12, 33, 48]. Winds can also gradually ablate the molecular material from forming stellar clusters [43]. Offner and Arce [35] quantified the stellar wind mass-loss rates for individual stars, which they found must be greater than  $10^{-7} M_{\odot} \text{ yr}^{-1}$  to be consistent with observations. Additionally, ionizing radiation feedback from O-stars also influences the morphology of clouds and the formation of stars [10, 11, 12, 18, 23].

Despite many observational and theoretical studies, the importance and impact of feedback on molecular clouds remain debated. This is because wind signatures are difficult to identify and quantify. Most bubble searches are done “by eye” [8, 2, 28]. For example, over 35,000 citizen scientists participated in the Milky Way Project [MWP, 46] in order to identify bubbles in Spitzer images. This approach is time-consuming, subjective and difficult to calibrate [4]. Analyzing the completeness of visually identified bubbles, which has a significant effect on the estimation of the injected momentum and energy, remains a great challenge. However, automatic classifications driven by machine learning approaches enable systematic, quantifiable and repeatable searches to identify bubbles [3, 4].

One of the most popular types of machine learning algorithms in astronomical classification is “Random Forests” [e.g. 7, 4, 29], which are based

on decision trees. A decision tree is a data structure which classifies feature vectors by computing a series of constraints, and propagating vectors down the tree based on whether these constraints are satisfied. Compared to other machine learning approaches, the Random Forests approach does well in classifying problems that have a large number of feature dimensions [6]. Beaumont et al. [4] developed an algorithm *Brut* based on Random Forests and applied it to classifying bubbles in the Milky Way. For each bubble, they defined a “score”, which is related to the probability that a given structure is a bubble. After conducting a blind search in the Milky Way, they found a substantial population of low-score bubble candidates not in MWP catalog produced by citizen scientists. In other words, citizen scientists are likely to miss a significant number of bubbles, but machine learning can compensate for some of this incompleteness.

Increasingly rich and detailed data of the local ISM and star-forming regions are available, such as GLIMPSE [Galactic Legacy Infrared Mid-Plane Survey Extraordinaire, 5], Hi-GAL [Herschel infrared Galactic Plane, 31] Survey and GALFA-HI [The Galactic Arecibo L-band Feed Array HI, 36] Survey. Parsing these extensive data visually is prohibitively time-consuming but is possible with the aid of machine learning algorithms.

There are two main types of machine learning algorithms: unsupervised learning and supervised learning. Unsupervised learning algorithms make their own criteria to discover structure in the data. An algorithm that learns from a training dataset and makes decisions based on the input “knowledge” is called

supervised learning. Supervised learning iteratively makes predictions on the training data and is corrected by the input training dataset. Consequently, the training dataset plays a significant role in the performance accuracy.

One fundamental problem with visual identification is that bubbles identified “by eye” are not objective and can be incorrect, such that machine learning approaches trained using flawed visual data will in turn produce defective identifications. In addition, there is no independent, quantitative assessment for completeness or any clear metric to determine how well bubbles are actually identified. One solution is to use realistic simulations, where feedback properties are known and well-defined. Such simulations can evaluate the accuracy of the training data and, in turn, supplement the original training dataset.

In this paper, we assess the performance of *Brut* in identifying bubbles using synthetic observations. We produce synthetic dust observations of bubbles in simulations. We use these as a supplemental training set to retrain *Brut* and test the performance of retrained *Brut* in classifying both synthetic bubbles and observed bubbles. We describe the method we use to construct synthetic observations and the details of the machine learning algorithm in Chapter 2. We compare and discuss several synthetic observation models in Chapter 3. In Chapter 4, we present the performance of retrained *Brut* in classifying both synthetic bubbles and observed bubbles. We summarize our results and conclusions in Chapter 5.

# Chapter 2

## Methods

### 2.1 Hydrodynamic Simulations

We adopt the magneto-hydrodynamics simulations from Offner and Arce [35], which aim to model winds from intermediate-mass stars and explore their impact on cloud morphology and turbulence. The simulations model a piece of a molecular cloud with length of  $L = 5$  pc, mass of  $M = 3762 M_{\odot}$  and periodic boundary conditions. The initial cloud temperature is  $T = 10$  K. The initial density and velocity conditions are set through driving the gas without gravity by adding random large-scale perturbations to the velocity field. These simulations share the same Alfvén Mach number 2.3 but their magnetic field distributions are spatially different at the initial time. Their velocity and density Fourier spectral slopes are comparable to  $S(k) \propto k^{-1.7}$  and  $S(k) \propto k^{-1.3}$ , respectively. The turbulence is initially external driving but ceases when the stellar sources are inserted and the begin feedbacks. Table 2.2 lists the parameters of these models. More details about the simulations can be found in Offner and Arce [35].

We adopt outputs from the strong wind run in which five stellar sources with different mass-loss rates are randomly placed. The number density of

Table 2.1: Physical Parameters of the Stellar Sources

ID	$M (M_{\odot})$	$L (10^3 L_{\odot})$	$T (10^4 \text{ K})$	$\dot{M} (10^{-7} M_{\odot} \text{ yr}^{-1})$
1	3.8	0.19	2.3	0.35
2	10.4	6.3	3.8	9.1
3	12.2	10.3	3.6	17.7
4	13.1	12.8	3.1	12.4
5	12.4	10.8	2.6	2.5

sources is similar to that in Perseus. These sources are all B-type stars with the mass-loss rates ranging from  $2.6 \times 10^{-8} - 1.8 \times 10^{-5} M_{\odot} \text{ yr}^{-1}$ . Table 2.1 lists the physical parameters of each of the five stellar sources. In this work, we explored outputs with different evolution stages and different turbulence realizations.

## 2.2 Hyperion

We use the publicly available three-dimensional dust continuum Monte-Carlo radiative transfer code HYPERION [42] to generate synthetic observations of the simulations described in Section 2.1. We adopt the gas density and temperature distributions from the outputs listed in Table 2.2 and the stellar properties from Table 2.1 as inputs. HYPERION assumes stars radiate as a blackbody.

Assumptions about the dust properties strongly influence the resulting emission. A variety of models for ISM dust have been proposed in the literature [e.g. 24, 13, 25], and we explore four different models in this work. Following Koepferl et al. [25], we combine three different dust grain models with 80.63%

Table 2.2: Model Properties<sup>a</sup>

Model	$t_i$ ( $t_{\text{cross}}$ )	$t_{\text{run}}$ (Myr)
T1_t1 <sup>b</sup>	1.6	0.1
T2_t1 <sup>c</sup>	2.0	0.1
T2_t0	2.0	0.05

Notes:

<sup>a</sup> Model name, the initial start time in crossing times and the evolutionary time. All models have  $L = 5$  pc,  $M = 3762 M_{\odot}$ ,  $T_i = 10$  K and initial  $B = 13.5 \mu\text{G}$ .

<sup>b</sup> Output corresponding to the model “W1\_T1” with an evolutionary time of 0.1 Myr in Offner and Arce [35].

<sup>c</sup> Output corresponding to the model “W1\_T2” in Offner and Arce [35].

big grains ( $>200 \text{ \AA}$ ), 13.51% smaller dust species, called very small grains (20–200  $\text{\AA}$ , vsg), and 5.86% PAH molecules, called ultra-small grains ( $<20 \text{ \AA}$ , usg). We label this dust model “K16” in the following discussion. We assume a moderate gas-to-dust ratio of 100 [44] and adopt a regular Cartesian grid with young stars embedded within. We calculate the emission for 20 different angular views and convolve the spectra with the Spitzer transmission curve [38, 21] to generate synthetic images in three Spitzer bands (4.5, 8, 24  $\mu\text{m}$ ). Figure 2.1 shows synthetic bubble images of the five sources with 20 different viewing angles.

In addition to the K16 dust model above, we adopt three other commonly used dust models to produce synthetic observations:

- (1). “kmh” dust model [24], which consists of astronomical silicates, graphite,

and carbon with full scattering properties,

- (2). “Draine” dust model [13], which is mainly Milky-Way carbonaceous-silicate grains,
- (3). “IPS” dust model [45], which represents “iron-poor” silicate dust.

Figure 2.2 shows synthetic images adopting the kmh dust model. The synthetic observations adopting the Draine and IPS dust models are similar to those adopting the kmh dust model, so we only include images with the kmh model.

The SEDs of different dust models show distinct differences, especially at  $8\ \mu\text{m}$  where PAH emission dominates. We extract the observed spectra of the main molecular cloud of Ophiuchus, LDN 1688 [39] and compare the SEDs of the different dust models as shown in Figure 2.3. The K16 dust model appears to be more realistic since it includes PAH emission while the other models lack PAH emission around  $8\ \mu\text{m}$ . Since the SEDs of the kmh model, Draine model, and the IPS model have a similar intensity at  $4.5\ \mu\text{m}$ ,  $8\ \mu\text{m}$  and  $24\ \mu\text{m}$ , the Draine and IPS three-color synthetic images look similar to the kmh model shown in Figure 2.2. The interiors of the bubbles in Figure 2.2 appear to be redder. This is because the  $24\ \mu\text{m}$  emission is stronger, but they lack  $8\ \mu\text{m}$  emission, compared to Figure 2.1. Consequently, we adopt the K16 dust model for the remainder of the analysis.

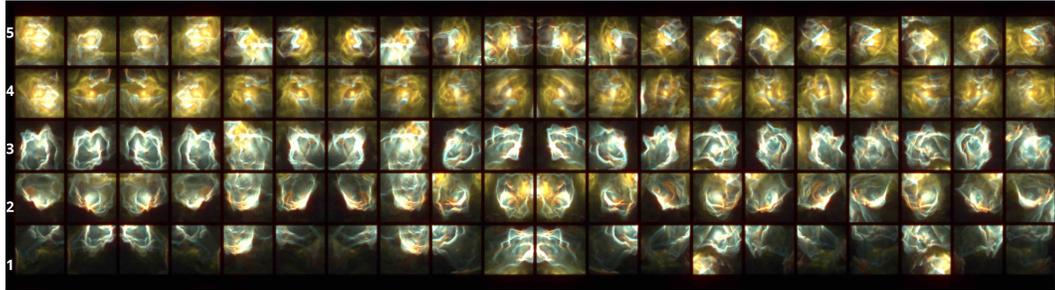


Figure 2.1: Three-color synthetic images of five sources with 20 different viewing angles adopting the dust model in Koepferl et al. [25]. Red, green and blue represents  $24 \mu\text{m}$ ,  $8 \mu\text{m}$  and  $4.5 \mu\text{m}$  emission, respectively.

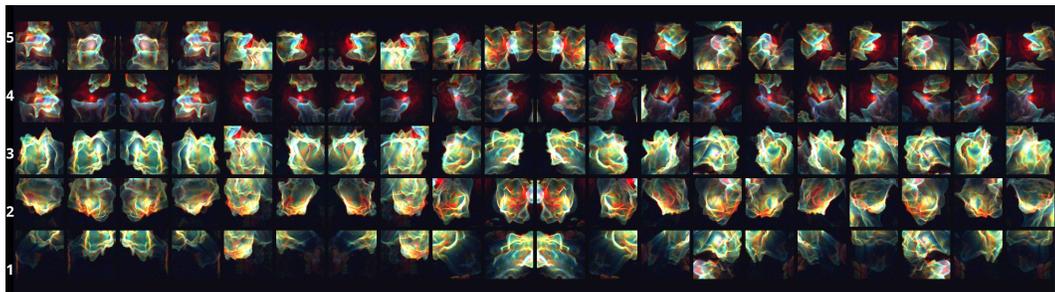


Figure 2.2: Same as Figure 2.1 but adopting the kmh dust model in Kim et al. [24].

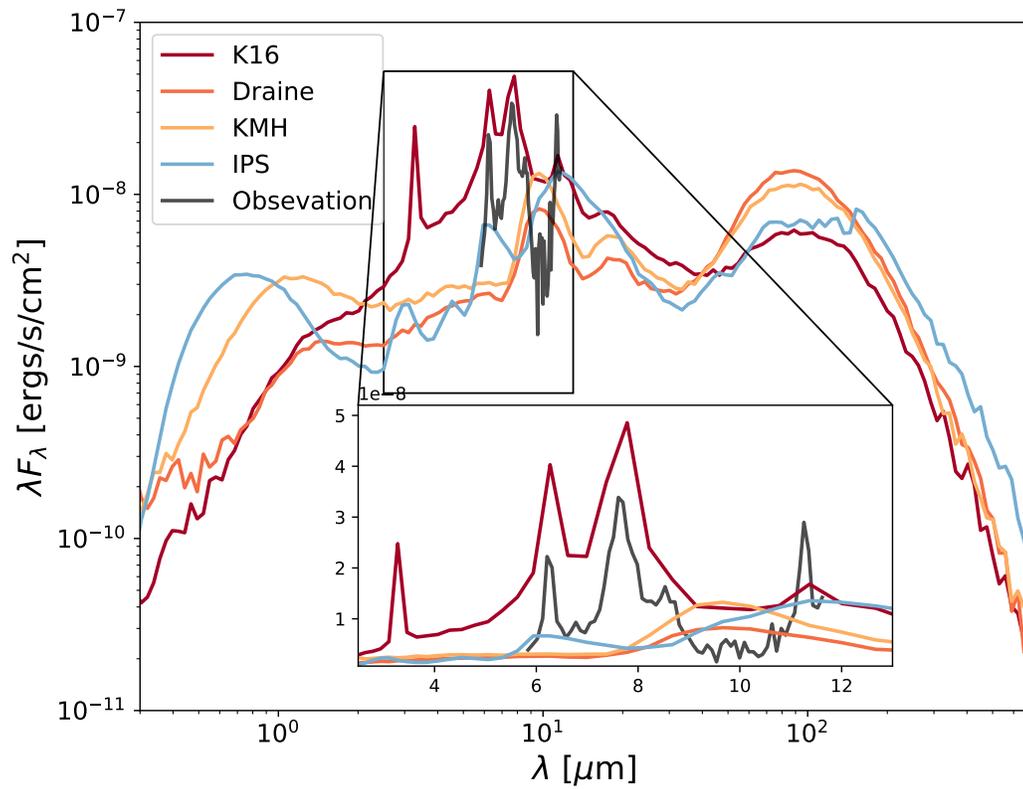


Figure 2.3: The SEDs of different dust models compared with the spectra of the Ophiuchus cloud LDN 1688 observed by Rawlings et al. [39].

## 2.3 *Brut*

*Brut* is a machine learning algorithm developed to identify bubbles in infrared images of the Galactic midplane [4]. *Brut* uses a Random Forest approach that is based on decision trees. A decision tree is a data structure that classifies a set of features, i.e., a numerical vector that describes the properties of each region. *Brut* computes a series of constraints and propagates the features down the tree based on whether these constraints are satisfied. *Brut* defines four features, which extract the most useful information about the difference between bubble and non-bubble images. It concatenates them into a single feature vector to carry out the classification.

Beaumont et al. [4] adopted bubbles identified by citizen scientists from the Milky Way Project as a training set. We include this same data for our analysis. The training set consists of 468 visually identified bubbles and 2289 random fields that are not centered on a bubble. *Brut* has three forests on different subsets of the sky, which we denote r1, r2 and r3. Each forest is trained using examples from two-thirds of the survey area and then tested using the remaining one-third area, as shown in Table 2.3. The illustration of the zones can be found in Figure 7 in Beaumont et al. [4].

After training, *Brut* returns a score related to the probability that a given structure is a bubble. If  $P$  is the probability that a given structure belongs to the bubble set, the *Brut* score is defined as  $2P - 1$ , where -1 is unlikely to be a bubble and +1 is very likely. To find the threshold score for true bubbles, Beaumont et al. [4] conduct a survey using experienced astronomers.

Table 2.3: Random Forest Zone

Random Forest Name	Training Zone ( $l$ ) <sup>a</sup>	Test Zone ( $l$ )
r1	$3n + 0.5^\circ \leq l < 3n + 1.5^\circ$ <sup>b</sup>	$3n + 1.5^\circ \leq l < 3n + 3.5^\circ$ <sup>d</sup>
r2	$3n + 1.5^\circ \leq l < 3n + 2.5^\circ$	$3n - 0.5^\circ \leq l < 3n + 1.5^\circ$
r3	$3n - 0.5^\circ \leq l < 3n + 0.5^\circ$ <sup>c</sup>	$3n + 0.5^\circ \leq l < 3n + 2.5^\circ$

Notes:

<sup>a</sup> The training zones are interleaved across all longitudes.

<sup>b</sup>  $n$  is an integer ranging from 0 to 119.

<sup>c</sup> When  $n$  is 0, the training zone is  $359.5^\circ(-0.5^\circ) \leq l < 0.5^\circ$

<sup>d</sup> When  $n$  is 119, the test zone is  $358.5^\circ \leq l < 0.5^\circ(360.5^\circ)$

They find about 50% of astronomers are likely to judge a region with a *Brut* score of 0.2 as a bubble. Consequently, they set 0.2 as the minimum acceptable score.

# Chapter 3

## Synthetic Observations

We adopt models with different evolutionary stages and turbulence properties as listed in Table 2.2 and consider different dust models in the synthetic observations.

### 3.1 Cropped Data

When carrying out the synthetic observations, we exploit the periodic nature of the simulation domain and wrap the data so all views have complete  $N^3$  voxels, where  $N$  is the number of pixels in one dimension. However, for large image sizes ( $L \geq 3$  pc), the Monte Carlo calculation becomes prohibitively expensive at the resolution we require. Instead, we crop the data into cubes of length 2.2 pc and 3 pc with each individual stellar object at the center.

Figure 3.1 and 3.2 show the cropped synthetic bubble images. Compared with Figure 2.1, in which the bubbles are embedded in the molecular cloud, the synthetic bubble images of the cropped data (Figure 3.1 and 3.2) are less extinguished but the morphology does not change significantly. They appear to be a little bit brighter and bluer, which means the shorter wavelength

emission is less attenuated. Another advantage of this strategy is that the synthetic bubble images are not contaminated by as much foreground or background emission. For example, the bottom row of Figure 2.1 is contaminated by the bubble from the source in the third row. Although observational data likely have overlapping bubbles, most bubbles identified by citizen scientists in MWP tend to be isolated. Consequently, we adopt the cropping strategy to generate the synthetic bubble images in the following discussion.

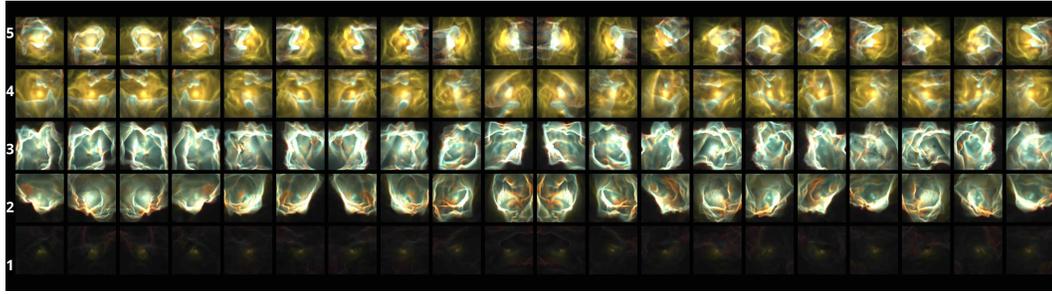


Figure 3.1: Three-color synthetic images adopting the K16 dust model where the HYPERION input is cropped to 2.2 pc.

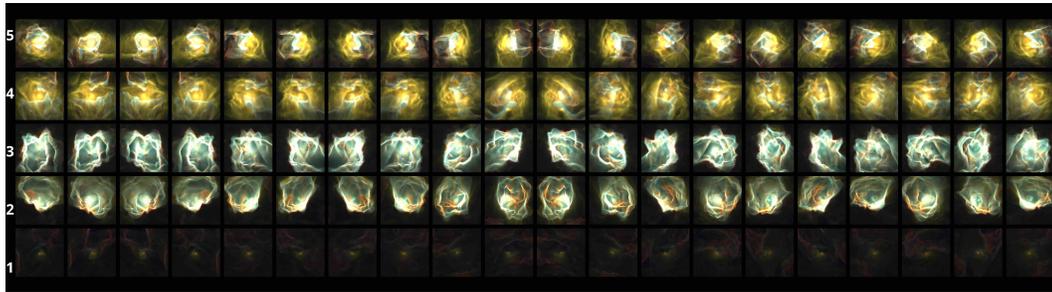


Figure 3.2: Three-color synthetic images adopting the K16 dust model where the HYPERION input is cropped to 3 pc.

## 3.2 Evolutionary Stage

The morphology of the bubbles changes with time as the winds expand into the cloud and interact with the surrounding gas. At earlier evolutionary stages, the bubbles are more compact compared with those at later stages, which have undergone additional expansion driven by the stellar wind. Figure 3.3 shows younger bubbles (“T2\_t0” listed in Table 2.2). The bubbles at the earlier time appear brighter in the center, owing to their compact and concentrated structure.

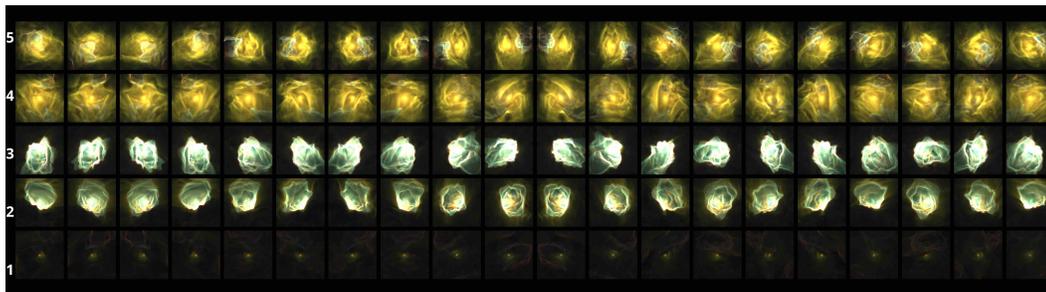


Figure 3.3: Three-color synthetic images at an earlier evolution stage with 0.05 Myr (“T2\_t0” listed in Table 2.2) where the HYPERION input is cropped to 3 pc.

## 3.3 Turbulent Realization

We also analyze a simulation with different initial turbulence. The synthetic observation process remains the same as described above, where we crop the HYPERION input data cube and use the K16 dust model. Figure 3.4 shows the synthetic images with different initial turbulence (“T1\_t1” listed in Table 2.2).

“T1\_t1” and “T2\_t1” have the same initial mean magnetic field, ratio of thermal to magnetic pressure, mean density, and stellar properties, but the shape of the bubbles are distinctly different owing to the different density distribution of the cloud material. Since the turbulent structure of real molecular clouds is varied, we adopt different initial turbulence to explore the diversity of bubble morphology and enrich our training dataset.

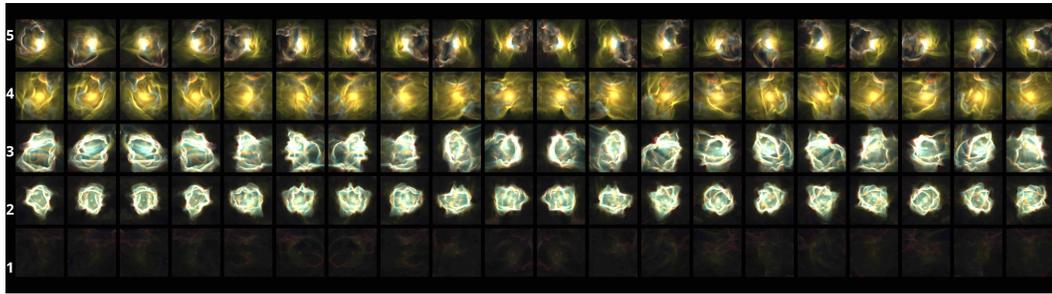


Figure 3.4: Three-color synthetic images with different turbulence where the HYPERION input is cropped to 3 pc.

### 3.4 Noise

The synthetic images are smooth, which is distinct from real observational images, which have fluctuations produced by noise. It is important that the training data be as close as possible to the observational data to reduce bias in detection caused by differences. To make the synthetic images more realistic, we identify patches of GLIMPSE data that are removed from the Galactic plane and have low signal to noise (S/N). We add these “stamps” to the synthetic images using the same S/N as the GLIMPSE data, where  $S/N \sim 8$ . Figure 3.5 shows the synthetic bubble images with noise.

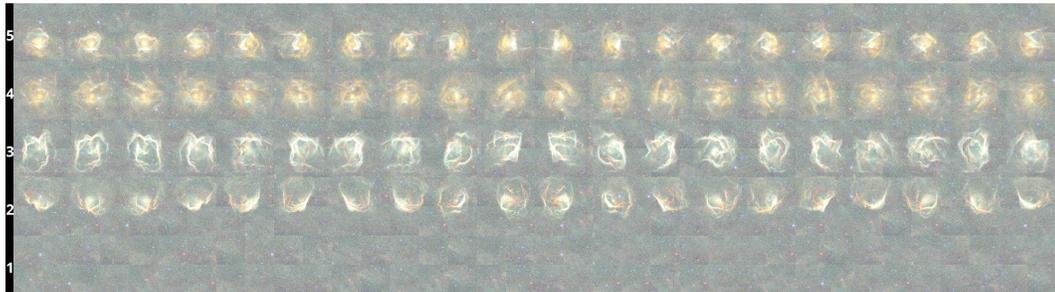


Figure 3.5: Three-color synthetic images with noise. From top to bottom: Figure 2.1, 3.1 and 3.2 with noise added.

## Chapter 4

### Results

#### 4.1 Retraining *Brut* with Synthetic Observation

We divide all the synthetic images into two equal parts. One half acts as a training data set, which we use to supplement the original MWP bubble set. The remainder serve as a test set, which allows us to assess the performance of the retrained algorithm. We summarize all the synthetic images we use in the training and testing sets in Table 4.1.

We analyze the performance of the three Random Forests before and after supplementing with the new training data. First, we retrain *Brut* using the synthetic images without noise (IDs 1-7 in Table 4.1). Figure 4.1 shows the performance with the original training and the algorithm retrained on noiseless synthetic images on the test bubbles. Table 4.2 briefly describes the meaning of labels in Figure 4.1. The scores returned after retraining on the noiseless data are significantly higher than those given by the original training. After retraining, the feature vector more accurately represents the synthetic bubbles and *Brut* does a better job identifying them.

We then augment the training set by adding the bubbles with IDs 7-14 in Table 4.1, so that the new training set consists of half of the bubbles with

and without noise. Figure 4.3 shows the performance with the original training and the algorithm retrained on synthetic images with and without noise on the second half of the noisy data. The scores returned by the retrained algorithm are significantly higher than those given with the original training. Compared with the scores retrained using noiseless data in Figure 4.1, the scores given by the retrained algorithm including some noisy images are more concentrated. This is likely because delicate bubble structure is reduced, i.e., there is less variation in bubble appearance since the noise hides small-scale sub-structure.

We next explore the impact of the training set size and composition on the performance of retrained algorithm. We retrain the algorithm with only synthetic images and retrain the algorithm with a set containing half the number of MWP bubbles and all the synthetic images. Figure 4.5 shows the performance of the algorithm trained with only synthetic images and the algorithm trained with fewer MWP images+synthetic images on the noisy data. Compared with the scores returned when training with all the MWP data and synthetic images in Figure 4.3, the scores returned by different random forests are similar but more concentrated. This is likely caused by the larger fraction of synthetic bubbles, which are similar to the test set, in the training set. The synthetic images are responsible for the better performance of the retrained algorithm on the synthetic images test set.

The increased scores after retraining suggest the original training dataset is incomplete, especially lacking bubbles driven by intermediate or low-mass stars. We further examine the performance of retrained *Brut* on observational

Table 4.1: Parameters of the Synthetic Images

ID	Label <sup>a</sup>	Turbulence <sup>b</sup>	Evolutionary Stage (Myr)	Image Size (pc)	Crop Crop	Noise Noise
1	T1_t1_c2	T1	0.1	2.2	✓	X
2	T1_t1_c3	T1	0.1	3	✓	X
3	T2_t1_2	T2	0.1	2.2	X	X
4	T2_t0_c2	T2	0.05	2.2	✓	X
5	T2_t0_c3	T2	0.05	3	✓	X
6	T2_t1_c2	T2	0.1	2.2	✓	X
7	T2_t1_c3	T2	0.1	3	✓	X
8	T1_t1_c2n	T1	0.1	2.2	✓	✓
9	T1_t1_c3n	T1	0.1	3	✓	✓
10	T2_t1_2n	T2	0.1	2.2	X	✓
11	T2_t0_c2n	T2	0.05	2.2	✓	✓
12	T2_t0_c3n	T2	0.05	3	✓	✓
13	T2_t1_c2n	T2	0.1	2.2	✓	✓
14	T2_t1_c3n	T2	0.1	3	✓	✓

Notes:

<sup>a</sup> The label with “n” indicates the synthetic image with noise.

<sup>b</sup> Turbulent distributions listed in Table 2.2.

data in Section 4.2 and 4.3.

Table 4.2: Parameters of the Random Forests

Random Forests	Training Set		Test Set
	Label <sup>a</sup>	MWP Zone <sup>b</sup> / Synthetic Image (ID) <sup>c</sup>	
	Label <sup>a</sup>	Number	Number
T1_t1_c2_r1		r1 / 314	no / 0
T1_t1_c2_r1s		r1 / 314	1-7 / 280
T1_t1_c3_r1		r1 / 314	no / 0
T1_t1_c3_r1s		r1 / 314	1-7 / 280
T2_t1_2_r1		r1 / 314	no / 0
T2_t1_2_r1s		r1 / 314	1-7 / 280
T2_t0_c2_r1		r1 / 314	no / 0

Table 4.2 Continued:

T2_t0_c2_r1s	r1 / 314	1-7 / 280	T2_t0_c2
T2_t0_c3_r1	r1 / 314	no / 0	T2_t0_c3
T2_t0_c3_r1s	r1 / 314	1-7 / 280	T2_t0_c3
T2_t1_c2_r1	r1 / 314	no / 0	T2_t1_c2
T2_t1_c2_r1s	r1 / 314	1-7 / 280	T2_t1_c2
T2_t1_c3_r1	r1 / 314	no / 0	T2_t1_c3
T2_t1_c3_r1s	r1 / 314	1-7 / 280	T2_t1_c3
T1_t1_c2n_r1	r1 / 314	no / 0	T1_t1_c2n
T1_t1_c2n_r1s	r1 / 314	1-14/ 560	T1_t1_c2n
T1_t1_c3n_r1	r1 / 314	no / 0	T1_t1_c3n
T1_t1_c3n_r1s	r1 / 314	1-14/ 560	T1_t1_c3n
T2_t1_2n_r1	r1 / 314	no / 0	T2_t1_2n
T2_t1_2n_r1s	r1 / 314	1-14/ 560	T2_t1_2n
T2_t0_c2n_r1	r1 / 314	no / 0	T2_t0_c2n
T2_t0_c2n_r1s	r1 / 314	1-14/ 560	T2_t0_c2n
T2_t0_c3n_r1	r1 / 314	no / 0	T2_t0_c3n
T2_t0_c3n_r1s	r1 / 314	1-14/ 560	T2_t0_c3n
T2_t1_c2n_r1	r1 / 314	no / 0	T2_t1_c2n
T2_t1_c2n_r1s	r1 / 314	1-14/ 560	T2_t1_c2n
T2_t1_c3n_r1	r1 / 314	no / 0	T2_t1_c3n
T2_t1_c3n_r1s	r1 / 314	1-14/ 560	T2_t1_c3n
T1_t1_c2n_r1Ns	no positive training set <sup>d</sup> / 0	1-14/ 560	T1_t1_c2n
T1_t1_c2n_r1Hs	half r1 / 159	1-14/ 560	T1_t1_c2n
T1_t1_c2_r2	r2 / 311	no / 0	T1_t1_c2
T1_t1_c2_r3	r3 / 311	no / 0	T1_t1_c2
T1_t1_c2n_r2Hs	half r2 / 161	1-14/ 560	T1_t1_c2n
T1_t1_c2n_r3Hs	half r3 / 158	1-14/ 560	T1_t1_c2n
:	:	:	:

Notes:

Table 4.2 Continued:

---

<sup>a</sup> We list the random forest “r1” and “r1s” for example, where suffix “s” means adding synthetic images into the training set. The random forests label with only “s” adopts the synthetic images without noise as part of the training set. The label with both “n” and “s” means adding the synthetic images with and without noise into the training set. There is a similar set of cases for random forest “r2”, “r2s”, “r3” and “r3s”. “r1Ns” indicates the training set only includes the synthetic images without any MWP bubbles in the positive training set. “r1Hs” means the training set consists of half the MWP bubble in r1 and all the synthetic images.

<sup>b</sup> The random forest zone listed in Table 2.3.

<sup>c</sup> The synthetic images listed in Table 4.1. The synthetic images are divided into two equal parts. One half acts as a training data set, and the second half serves as a test set.

<sup>d</sup> The training set does not have any MWP bubbles in the positive training set but contains the MWP images without bubbles in the negative training set.

---

## 4.2 Re-Testing *Brut* on the Milky Way Project Data

We adopt all 3716 large bubbles found by the citizen scientists in Simpson et al. [46] as a test set to assess the performance of *Brut* after retraining. We ignore the objects contained in the “small bubble” catalogue, which are mainly green knots, dark nebulae, star clusters, galaxies or fuzzy red objects. We compare the performance for the original training, the retrained algorithm (using both noisy and noiseless synthetic bubbles), the algorithm trained with only synthetic images and the algorithm trained with fewer MWP images+synthetic images in classifying MWP bubbles, as shown in Figure 4.7. The scores returned by the retrained algorithm are significantly higher com-

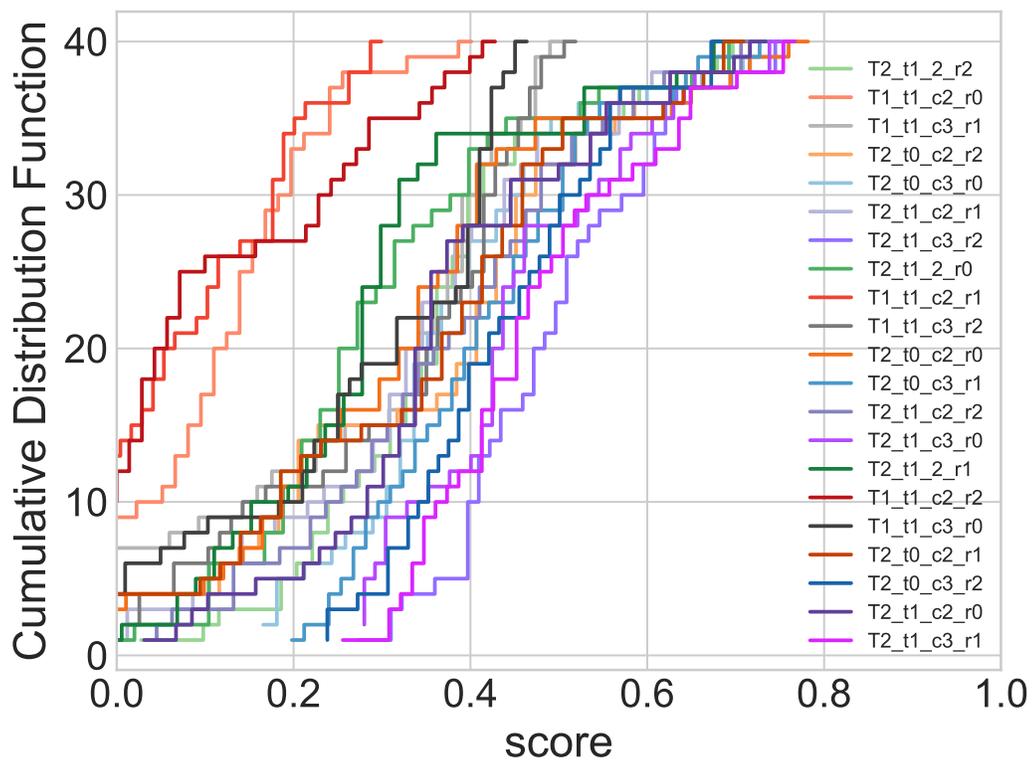


Figure 4.1: The cumulative distribution function (CDF) of all the scores given by Brut with the original training. The labels are described in Table 4.2.

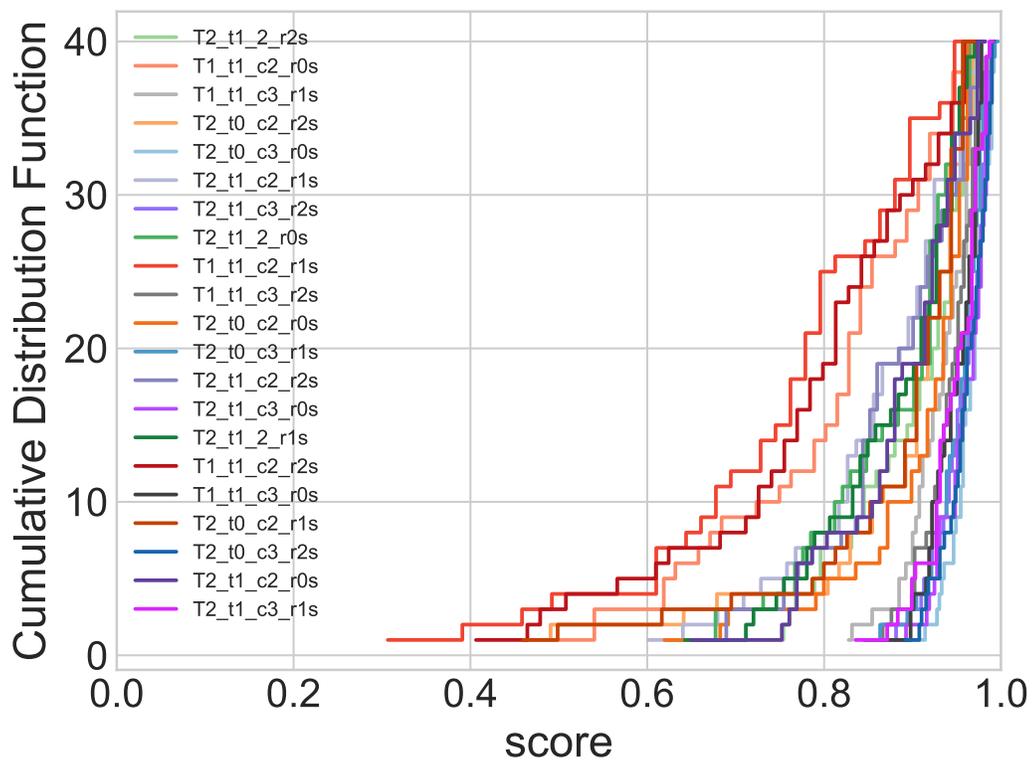


Figure 4.2: The cumulative distribution function (CDF) of all the scores given by Brut retained on noiseless synthetic images. The labels are described in Table 4.2.

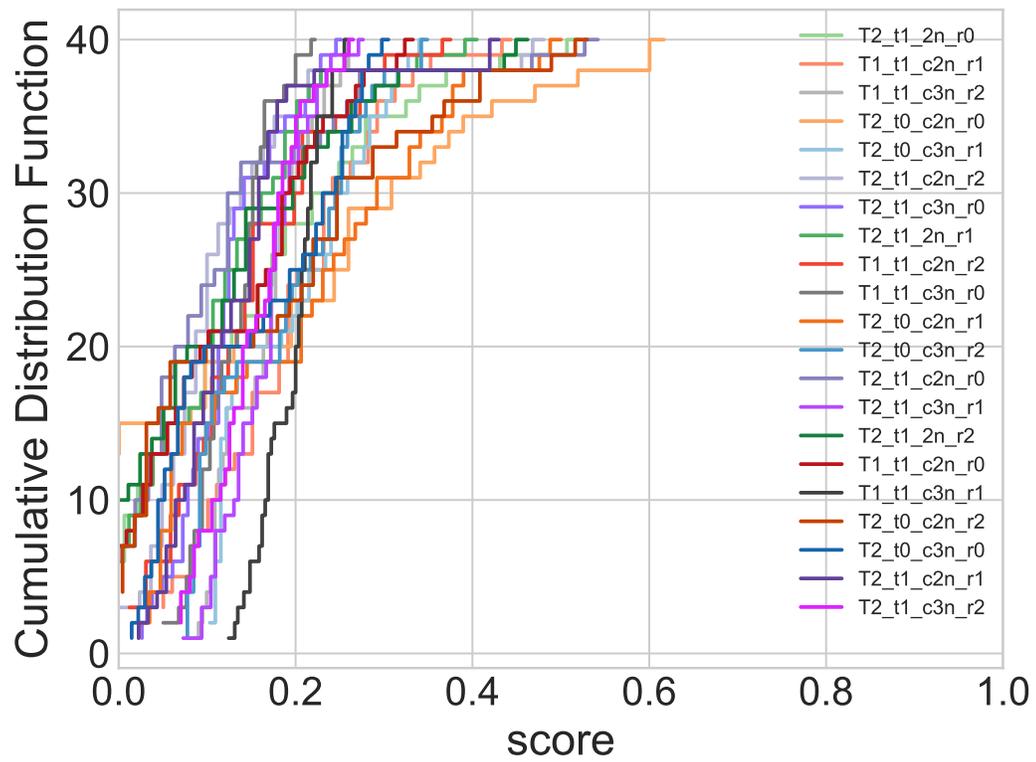


Figure 4.3: The CDF of all the scores given by Brut with the original training. The labels are described in Table 4.2.

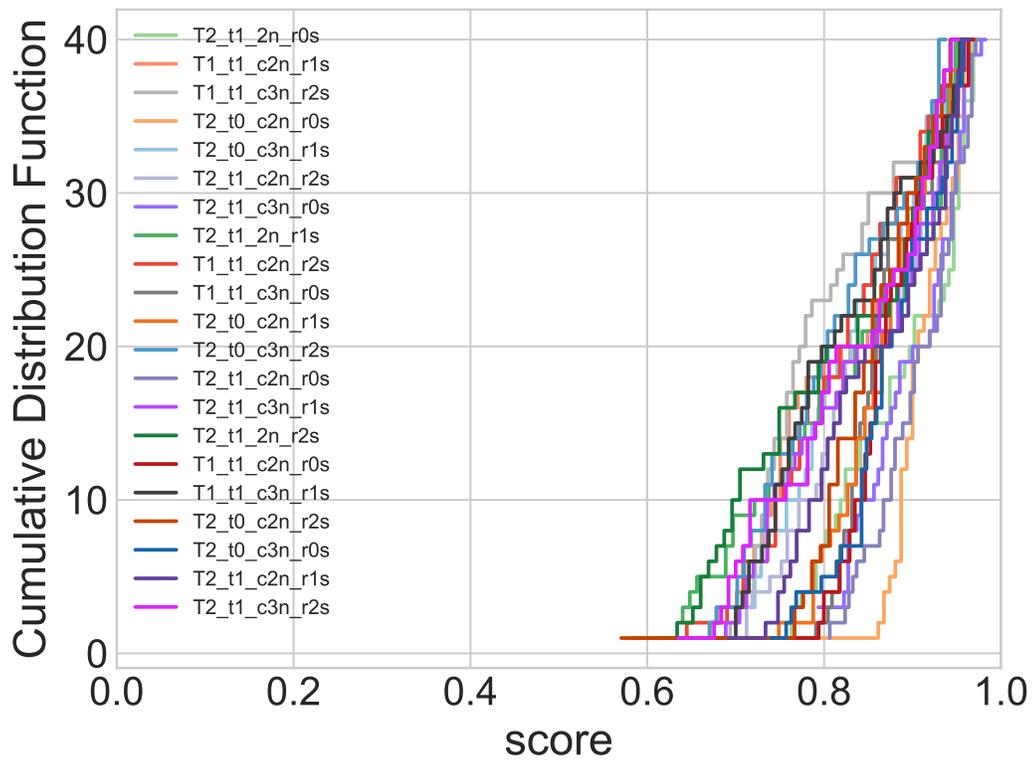


Figure 4.4: The CDF of all the scores given by Brut retrained on synthetic images with and without noise. The labels are described in Table 4.2.

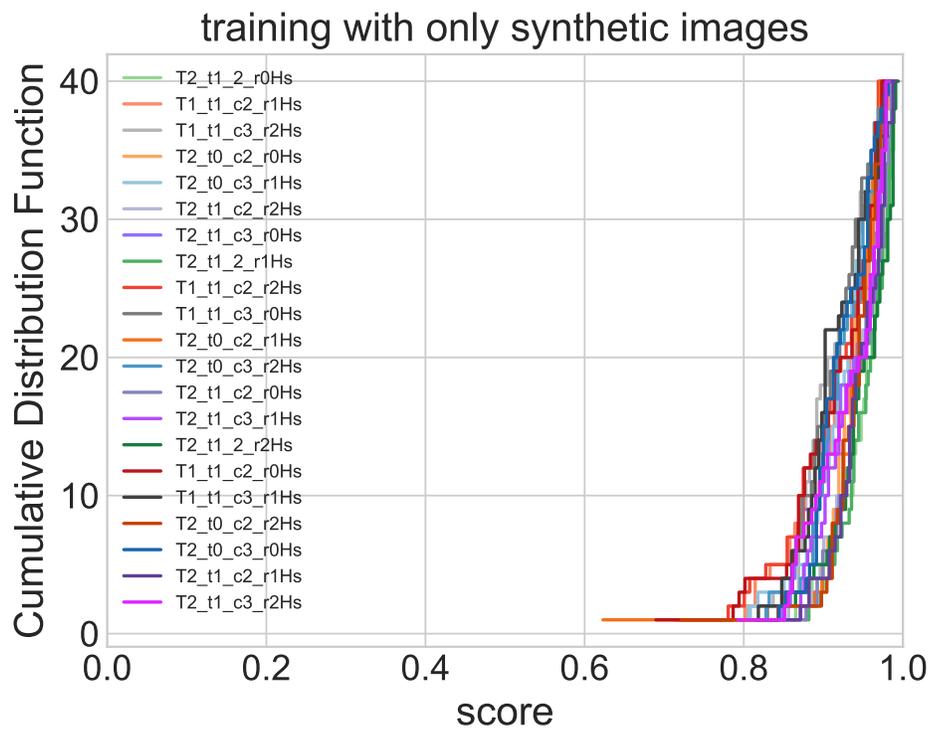


Figure 4.5: The CDF of all the scores given by the retrained algorithm without original MWP training set but with only synthetic images. The labels are described in Table 4.2.

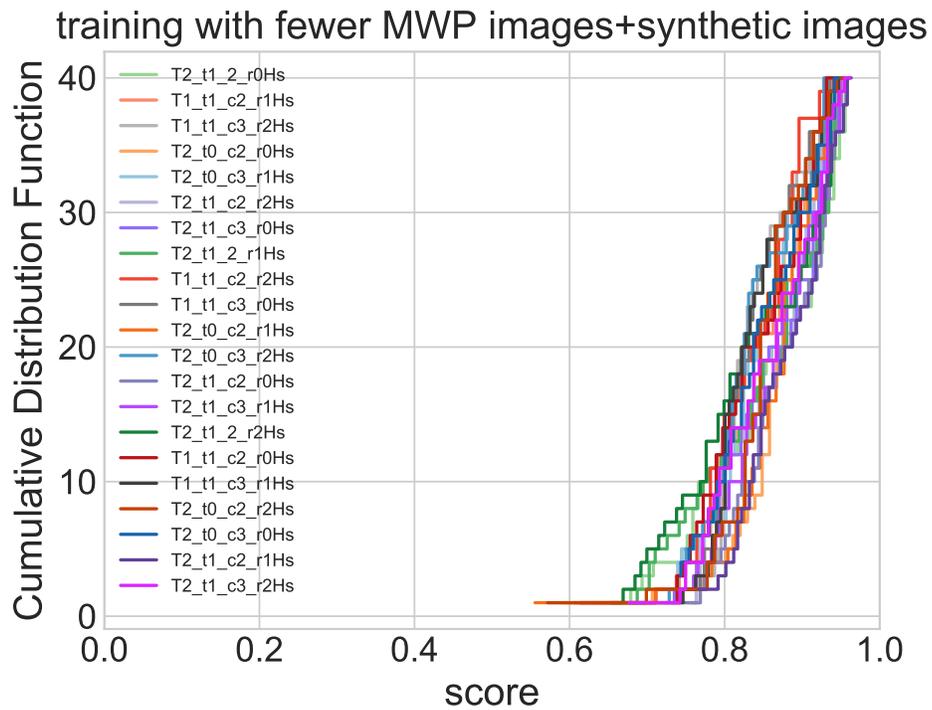


Figure 4.6: The CDF of all the scores given by the retrained algorithm on half original MWP training set and the synthetic images. The labels are described in Table 4.2.

pared with those returned by *Brut* without additional training. When we retrain the algorithm with only synthetic images, the scores under 0.55 show a dramatic improvement. After investigating the high and low score bubble images, we find the algorithm trained with only synthetic images improves the scores of ambiguous bubbles with low S/N and reduces the scores of red bubbles with high S/N.

To explain the performance of the algorithm retrained on several different training sets, we characterize the bubble properties that compose each training set as shown in Figure 4.8. The “Normalized S/N” quantifies the contrast and S/N of the image. We define it as

$$\text{Normalized } S/N = C(\bar{I}_{95} - \bar{I}_{30})(\bar{I}_{95} - \bar{I}_{50})f_{\geq 8\sigma}, \quad (4.1)$$

where  $(\bar{I}_{95} - \bar{I}_{30})$  is the difference between the top 5% and the bottom 30% of values,  $(\bar{I}_{95} - \bar{I}_{50})$  is the difference between the top 5% values and the median value,  $f_{\geq 8\sigma}$  is the fraction of bright pixels ( $\geq 8\sigma$ ), and  $C$  is a constant to normalize the values to unity. In most of the high S/N bubble images, the bubble rim structures occupy the top 5% of the image values, and the noise occupies the bottom 30%. The average of the diffuse emission is well represented by the median image value. We use this product to indicate the contrast of the image. In a random noisy image, the normalized S/N is close to 0. The  $x$ -axis in Figure 4.8 indicates the “Yellow Index,” which describes the color of the bubble. We define it as the ratio between the number of yellow pixels and the number of red pixels. Although the original training set spans a wide range of

color and S/N, it is concentrated in the red domain. The large representation of red bubbles in the training set means that *Brut* will more easily identify red bubbles than yellow bubbles. In contrast, the synthetic bubbles are located in the yellow part of the parameter space. The MWP bubbles are mostly low S/N red bubbles, with some low S/N yellow bubbles and high S/N red bubbles, but there are very few high S/N yellow bubbles. Consequently, the algorithm trained with only synthetic images mainly captures bubbles with low S/N. This explains why a training set with only synthetic images improves the scores of ambiguous bubbles with low S/N and reduces the scores of bubbles that are red and have high S/N.

Figure 4.7 also shows the result when we randomly remove half of the bubbles in the original MWP training set. The score distribution returned by the algorithm trained with fewer MWP images+synthetic images compared to when the results of the algorithm trained with only synthetic images is surprising. When including half of the original MWP bubbles in the training set, the performance of the algorithm dramatically decreases. The original MWP training bubbles are mostly red, while the synthetic images nearly all contain yellow bubbles. Consequently, these sets inhabit two different color domains. The reduction of red bubbles in the training set lowers the scores of these types of bubbles in the test set. When including all the original and synthetic images in the training set, the performance of the retrained algorithm significantly and steadily improves. Consequently, this demonstrates the composition and size of the training set significantly impacts the performance of the algorithm.

Following our comparison of the algorithm performance after retraining with several different training sets, we adopt the training that includes all the original MWP bubbles and synthetic images. This training set significantly improves the scores of most bubbles with little change in the number of high-score objects. Although the algorithm trained with only synthetic images improves the scores of a large number of bubbles, it no longer returns any high score bubbles, which were previously assigned to images with red bubbles.

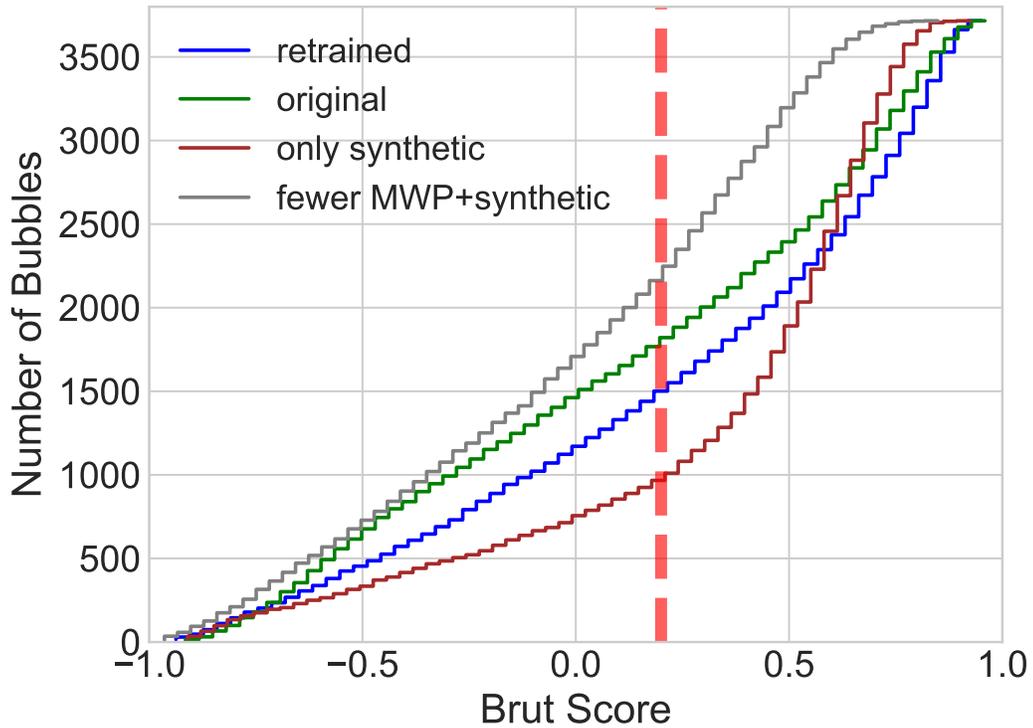


Figure 4.7: The CDF of the scores of 3716 Milky Way Project large bubbles given by algorithm with the original training and the algorithm retrained on several different training sets including synthetic bubbles.

The MWP characterizes the consensus among users that an image con-

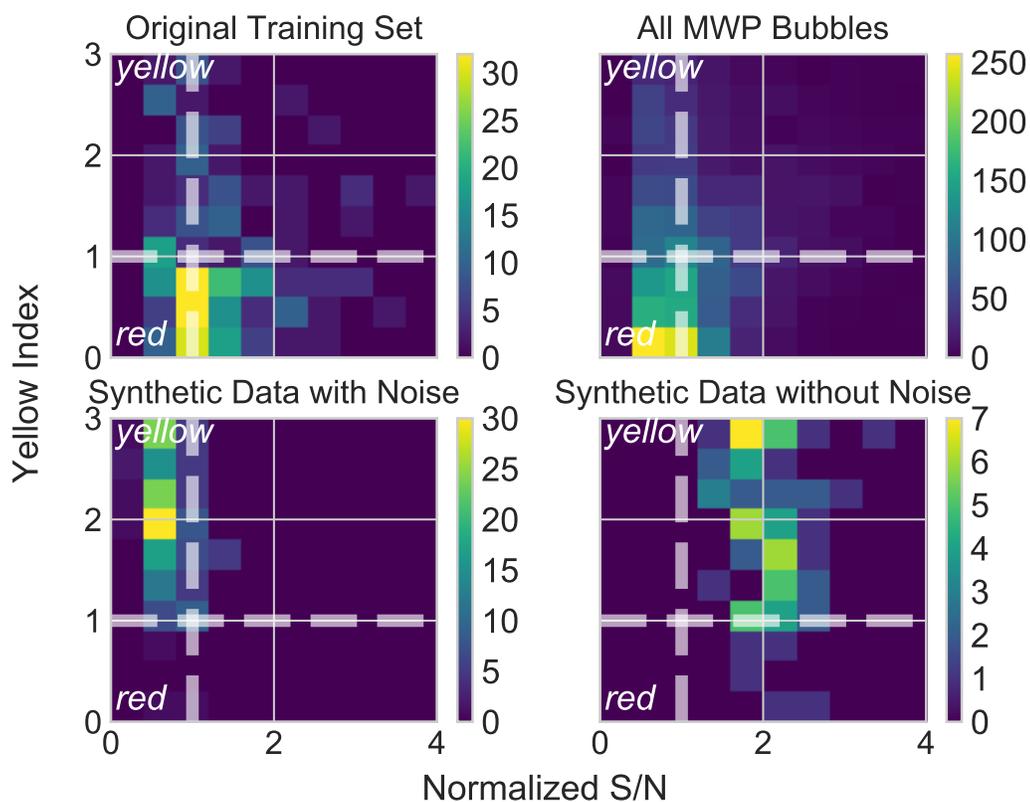


Figure 4.8: Distribution of image properties for four different training sets, where the colors indicate the number of bubbles in each bin. The white dashed lines divides the bubbles into four regions. The upper left quadrant indicates images that contain low S/N yellow bubbles. The upper right quadrant indicates images with high S/N yellow bubbles. The lower left quadrant indicates low S/N red bubbles. The lower right quadrant indicates high S/N red bubbles.

tains a bubble in terms of the “hit rate”, which is the fraction of citizen scientists who identified a bubble in the image. They define hit rates above 0.1 as being high-confidence bubble candidates.

We further compare the scores given by *Brut* with the original training, the scores after it is retrained, and the MWP hit rate as shown in Figure 4.9 and 4.10. The average *Brut* score in each bin with the original training and after retraining both show a clear trend with the hit rate. The error bars indicate the standard deviation of the scores and hit rate in each bin. The higher the hit rate, the higher the score *Brut* returns, which is consistent with our expectations. In other words, the retrained algorithm preserves the hit-rate distribution, where bubbles with low hit rates continue to have low scores.

Moreover, over 10% of the MWP bubbles, which were previously marginal or ambiguous detections, are reclassified as high-confidence bubbles after retraining. Their average *Brut* score increases from -0.07 to 0.39. About 2% of the previously identified MWP bubbles are no longer classified as high-confidence bubbles, and their average *Brut* score drops from 0.31 to 0.06.

Figure 4.11 shows one hundred bubbles, whose score significantly increases after retraining. Most of these bubbles are yellow, indicating the  $8\ \mu\text{m}$  and  $24\ \mu\text{m}$  emission are similar. These yellow bubbles are likely ultra-compact and compact H II regions or analogous regions for less massive B-type stars [22]. The performance of the retrained algorithm is consistent with our training set, in which bubbles are created by the stellar winds of B-type stars. For these

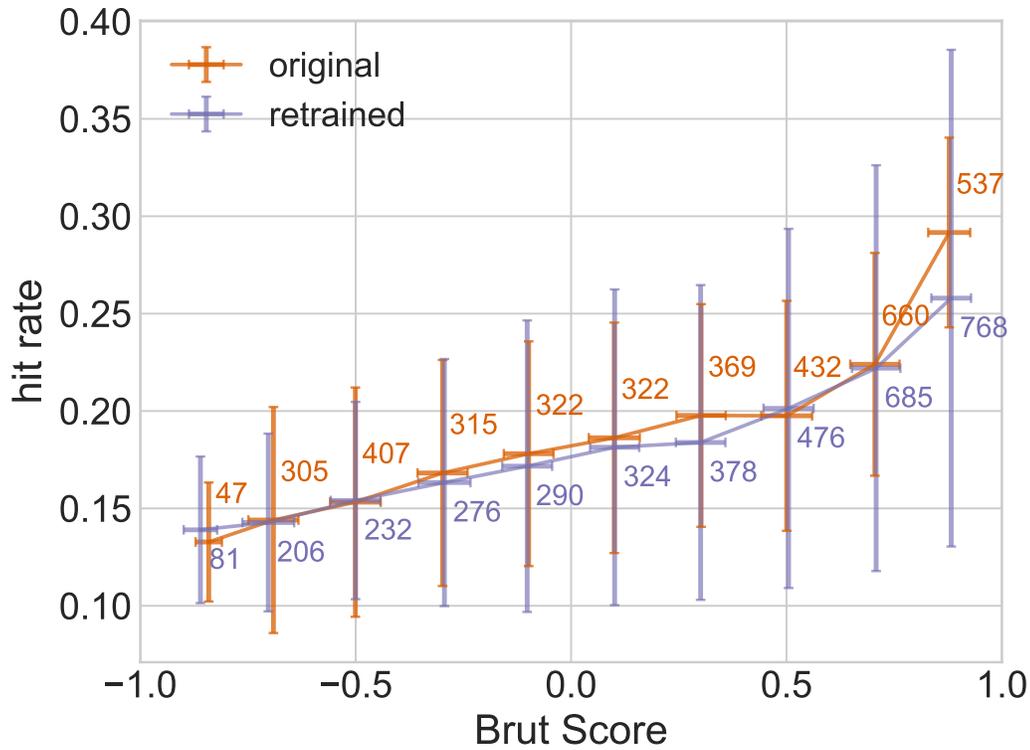


Figure 4.9: The average hit rate versus the average binned *Brut* score for the 3716 MWP large bubbles. The red and green lines indicate the average *Brut* score returned with the original training and the algorithm retrained on synthetic images both with and without noise, respectively. The error bars indicate the standard deviation of the scores and hit rate in each bin. The label indicates the number of bubbles in each bin.

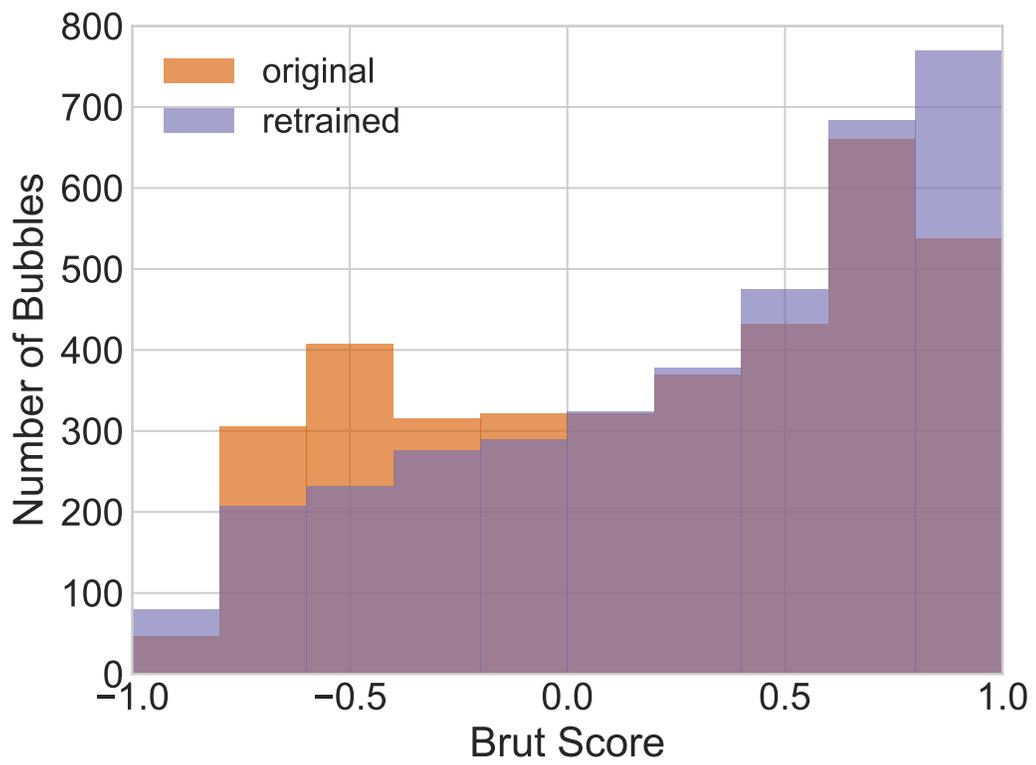


Figure 4.10: The distribution of bubble scores returned with the original training and after *Brut* is retrained.

type of stars, the amount of ionizing radiation is small, so the bubbles are predominantly cleared by the wind (or earlier protostellar outflows) and then illuminated by the stellar radiation field.

Figure 4.12 shows nine bubbles, which were previously identified MWP bubbles but are no longer classified as high-confidence bubbles after retraining. These bubbles are very red and, thus, quite distinct from our yellow bubbles, and their morphology does not show a distinct shell rim. Consequently, since we supplemented the training set with synthetic yellow bubbles, the decline of these bubbles *Brut* scores is unsurprising.

In summary, the performance of the retrained algorithm in classifying yellow bubbles significantly increases when synthetic observations are added to the training set.

### 4.3 Application: Bubbles in the Perseus Molecular Cloud

Perseus is located in the larger Taurus-Auriga-Perseus dark cloud complex with a distance of  $250 \pm 50$  pc, spanning a total area of about  $70 \text{ pc}^2$  [15, 16]. With a mass of  $10^4 M_{\odot}$ , the Perseus cloud is often considered to be an intermediate case between low-mass star forming regions such as Taurus and turbulent, high-mass regions such as Orion [27], making it an ideal location to study low and intermediate-mass star formation. The feedback of young stars makes Perseus a “bubbly” cloud [2].

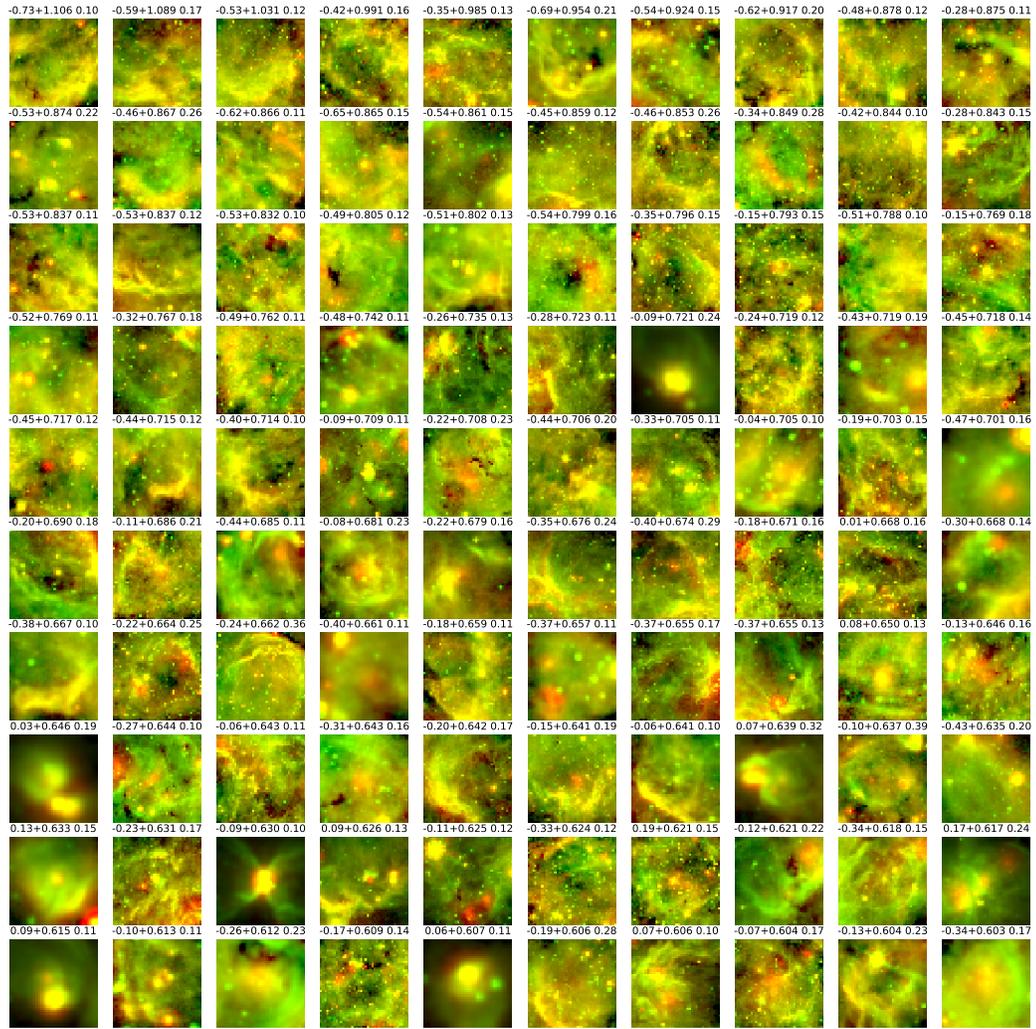


Figure 4.11: One hundred bubbles from MWP. The first number in the title of each panel presents the raw score, which is returned by the original MWP training algorithm. The middle number is the change in score after *Brut* is retrained with synthetic observations. The last number is the hit rate.

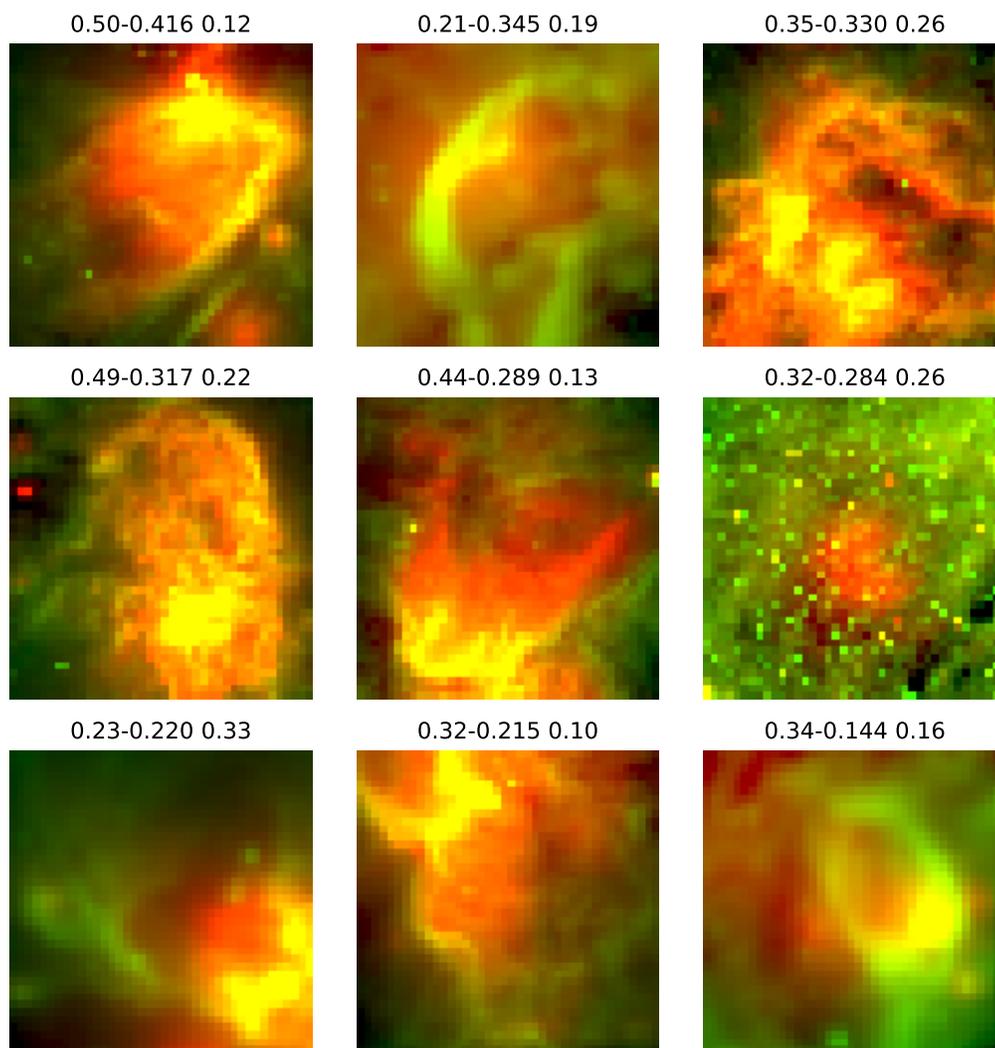


Figure 4.12: Nine bubbles, which were previously likely MWP bubbles but are no longer classified as high-confidence bubbles after retraining. The meaning of each number in the title is described in Figure 4.11.

Arce et al. [2] identified 12 bubbles using CO spectral data. We extract the Spitzer image of Perseus in 4.5  $\mu\text{m}$ , 8  $\mu\text{m}$  and 24  $\mu\text{m}$  bands (Gutermuth priv. comm.) and apply *Brut* to this data. Figure 4.13 shows four examples of bubbles in the Perseus molecular cloud. These bubbles are associated with shells CPS6, CPS8, CPS10 and CPS11 in the CO data, which were visually identified by Arce et al. [2]. Table 4.3 lists the physical properties of these bubbles. All these bubbles are probably driven by relatively low or intermediate-mass young stars such as B type or F type stars. Figure 4.13 shows these bubbles and their associated *Brut* scores before and after retraining. These four cases show a significant improvement in score, most from a negative score (non-bubble) to a positive score (likely bubble).

CPS6 and 8 are similar to the synthetic bubbles and the MWP yellow bubbles. They are the best examples of the good performance produced by retraining *Brut*. CPS11 is a partial bubble, which is probably why its score is still  $< 0.2$ . Table 4.3 shows that CPS10 is driven by a B5V star, but there is no distinct evidence of the existence of the star in the infrared images. However, in the optical data, the star is bright and is clearly visible.

The dust emission exceeds that of the star, so the B5V star becomes invisible when embedded in the cloud. Although CPS10 is not a yellow bubble, it is nonetheless consistent with the bubble model in Figure 1 in Beaumont et al. [4], where green shell structure is produced by PAH emission and the red interior is dominated by hot dust. The bubble score is low, which is likely due to the contamination by other emission at the upper right corner.

Table 4.3: Physical Properties of Four Perseus Bubbles

Bubble Name	Cloud Region	Center ( $\alpha_{2000}, \delta_{2000}$ )	Candidate Source	Source Type
CPS-6	L1468	03 41 24, 31 54 10	IRAS 03382+3145	unknown <sup>a</sup>
CPS-8	IC 348	03 44 10, 32 17 20	omi Per	B1III <sup>b</sup>
CPS-10	IC 348	03 44 35, 32 10 10	HD 281159	B5V <sup>c</sup>
CPS-11	IC 348	03 44 50, 32 18 10	V* 695 Per & IC 348 LRL 30	M3.75 & F0 <sup>d</sup>

Notes:

<sup>a</sup> This is likely not a typical main-sequence star; it is probably a young pre-main-sequence star [2]. There is no reliable mid-IR detection in the MIPS bands because the source is confused with the bright nebulosity in this region [40].

<sup>b</sup> This bright source is observed both in the optical and the infrared, and it is classified as a YSO candidate in the c2d point-source catalog [16].

<sup>c</sup>. HD 281159 is a binary system with two massive B5 main sequence stars, which has a disk around the binary pair. And they have an age  $\leq 10$  Myr [32].

<sup>d</sup> Two possible candidates are V\* 695 Per and IC 348 LRL 30, which are an M3.75 star and F0 star, respectively. Both stars are classified as a YSO candidates in the c2d catalog and both have been detected in X-ray [37].

These results indicate the retrained algorithm can perform well for molecular cloud data not included in the MWP. The synthetic observations are able to improve *Brut* performance in classifying bubbles produced by relatively low or intermediate-mass young stars such as B-type stars.

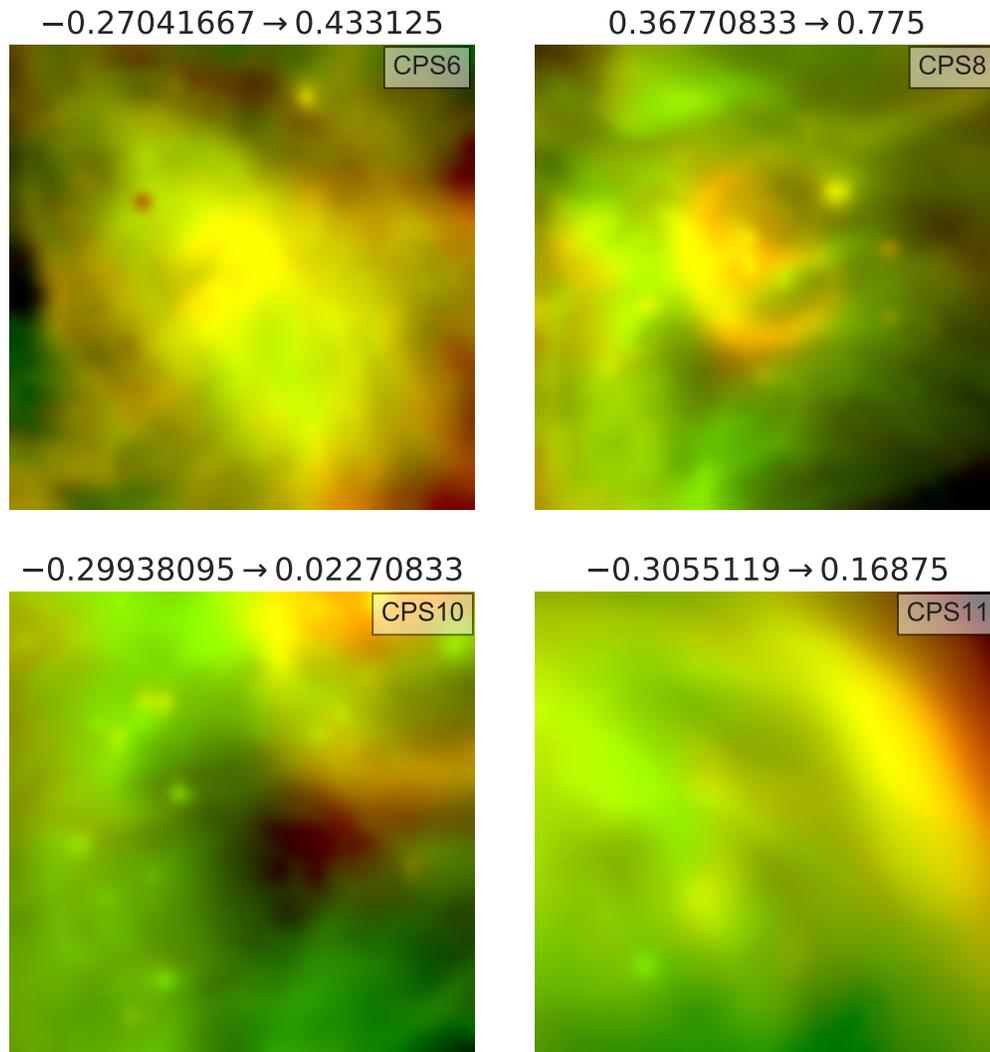


Figure 4.13: Four examples of bubbles in the Perseus molecular cloud. The upper right label in each panel corresponds to the bubble name in Arce et al. [2]. The left number in the title of each panel indicates the raw score, which is returned by the original training algorithm. The right number in the title of each panel indicates the new score returned by the retrained algorithm.

# Chapter 5

## Summary

### 5.1 Conclusions

We adopt magneto-hydrodynamics simulations of stellar winds interacting with a molecular cloud and post-process them using a three-dimensional dust continuum Monte-Carlo radiative transfer code. We generate synthetic observations of bubbles in the Spitzer bands ( $4.5\ \mu\text{m}$ ,  $8\ \mu\text{m}$  and  $24\ \mu\text{m}$ ). We employ a previously developed machine learning algorithm, *Brut* and quantitatively evaluate its performance in identifying bubbles using synthetic dust observations. Our main findings are the following:

1. Synthetic observations in combination with visually identified sources can be used to significantly improve machine learning classification.
2. After retraining with synthetic images, *Brut* better identifies yellow bubbles, which are likely associated with H II regions for less massive B-type stars or cavities evacuated by stellar winds.
3. The completeness of the training set significantly impacts the performance of the algorithm. We suggest that the number of yellow bubbles in the current MWP bubble catalog is incomplete, and we expect a random search

of the full GLIMPSE dataset with *Brut* would return many more yellow bubble candidates.

4. Some of the bubbles with improved scores are associated with lower confidence sources in the MWP. These would likely be identified as bubbles by an expert, and thus the simulations provide an efficient means to enhance machine learning training sets.
5. Turbulent structures greatly affect the morphology of bubbles, yielding a variety of bubble shapes. Different evolutionary stages and different cropped image sizes further enhance the bubbles contained in the training set. Adding noise similar to that in the GLIMPSE data makes the synthetic observations more realistic. In combination, these modifications create a more complete training set to improve the machine learning classifications.
6. The retrained algorithm performs well classifying bubbles associated with more embedded sources located in Perseus. Thus, retraining with synthetic observations expands the parameter space of the training set beyond the less embedded and more distant regions with massive stars covered by the MWP.

## 5.2 Future Work

*Brut* is sensitive to the position of bubbles in the image. This makes it computationally expensive to identify bubbles in a large sky survey map because it needs to crop the map into small chunks at different positions with

different sizes. Recently developed deep learning methods are more powerful in image recognition than *Brut*. Ntampaka et al. [34] develop a deep machine learning tool based on Convolutional Neural Networks (CNNs) to estimate the mass of galaxy clusters in X-ray emission. The CNN is not sensitive to the position of galaxy clusters, making it straightforward to apply on large sky survey maps. Van Oort et al. [47] develop an “Encoder-Decoder” CNN to identify stellar wind bubbles in density slices and 2D CO emission. This CNN approach achieves a 98% accuracy. However, one caveat of these models is that they are limited to 2D integrated intensity maps. These algorithms do not take the information along the velocity axis into consideration, which may lead to a high false detection rate. In other words, this technique may identify a clear ring structure as a bubble even though this structure is caused by a turbulent pattern without any evidence of expansion in the spectra.

In the future, we are planning to extend the algorithm from Van Oort et al. [47] to 3D CNNs in order to exploit the full 3D CO data information (position-position-velocity). We will apply the publicly available radiation transfer code RADMC-3D [14] to model the  $^{12}\text{CO}$  and  $^{13}\text{CO}$  ( $J=1-0$ ) line emission of the MHD simulations. We will train the 3D CNN model with synthetic observations and apply it to the observational CO emission data to identify stellar feedback bubbles in molecular clouds.

Furthermore, substantial surveys and archival observational data are available to study stellar feedback, such as the GPS [Spitzer Galactic Plane Survey, 5], Hi-GAL [Herschel infrared Galactic Plane, 31] Survey, GALFA-HI

[Galactic Arecibo L-band Feed Array HI, 36] Survey and FCRAO Gould Belt Survey[41]. An upcoming survey by the Large Millimeter Telescope (LMT [20]) will provide considerable dust emission and molecular emission data. It is almost impossible for humans to identify feedback features visually given the exponentially increasing amount of data. However, systematic and repeatable identification is possible with the aid of machine learning approaches [3, 4, 47]. An unbiased sample of stellar feedback features will enable a better understanding of the origin of turbulence and the energy budget in molecular clouds.

## Bibliography

- [1] Héctor G. Arce, Michelle A. Borkin, Alyssa A. Goodman, Jaime E. Pineda, and Michael W. Halle. The COMPLETE Survey of Outflows in Perseus. *Astrophysical Journal*, 715:1170–1190, June 2010. doi: 10.1088/0004-637X/715/2/1170.
- [2] Héctor G. Arce, Michelle A. Borkin, Alyssa A. Goodman, Jaime E. Pineda, and Christopher N. Beaumont. A Bubbling Nearby Molecular Cloud: COMPLETE Shells in Perseus. *Astrophysical Journal*, 742:105, December 2011. doi: 10.1088/0004-637X/742/2/105.
- [3] Christopher N. Beaumont, Jonathan P. Williams, and Alyssa A. Goodman. Classifying Structures in the Interstellar Medium with Support Vector Machines: The G16.05-0.57 Supernova Remnant. *Astrophysical Journal*, 741:14, November 2011. doi: 10.1088/0004-637X/741/1/14.
- [4] Christopher N. Beaumont, Alyssa A. Goodman, Sarah Kendrew, Jonathan P. Williams, and Robert Simpson. The Milky Way Project: Leveraging Citizen Science and Machine Learning to Detect Interstellar Bubbles. *The Astrophysical Journal Supplement Series*, 214:3, September 2014. doi: 10.1088/0067-0049/214/1/3.
- [5] Robert A. Benjamin, E. Churchwell, Brian L. Babler, T. M. Bania, Dan P.

- Clemens, Martin Cohen, John M. Dickey, Rémy Indebetouw, James M. Jackson, Henry A. Kobulnicky, Alex Lazarian, A. P. Marston, John S. Mathis, Marilyn R. Meade, Sara Seager, S. R. Stolovy, C. Watson, Barbara A. Whitney, Michael J. Wolff, and Mark G. Wolfire. GLIMPSE. I. An SIRTFF Legacy Project to Map the Inner Galaxy. *Publications of the Astronomical Society of the Pacific*, 115:953–964, August 2003. doi: 10.1086/376696.
- [6] Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- [7] Samuel Carliles, Tamás Budavári, Sébastien Heinis, Carey Priebe, and Alexander S. Szalay. Random Forests for Photometric Redshifts. *Astrophysical Journal*, 712:511–515, March 2010. doi: 10.1088/0004-637X/712/1/511.
- [8] E. Churchwell, M. S. Povich, D. Allen, M. G. Taylor, M. R. Meade, B. L. Babler, R. Indebetouw, C. Watson, B. A. Whitney, M. G. Wolfire, T. M. Bania, R. A. Benjamin, D. P. Clemens, M. Cohen, C. J. Cyganowski, J. M. Jackson, H. A. Kobulnicky, J. S. Mathis, E. P. Mercer, S. R. Stolovy, B. Uzpen, D. F. Watson, and M. J. Wolff. The Bubbling Galactic Disk. *Astrophysical Journal*, 649:759–778, October 2006. doi: 10.1086/507015.
- [9] J. E. Dale and I. A. Bonnell. The effect of stellar winds on the formation of a protocluster. *Monthly Notices of the Royal Astronomical Society*, 391: 2–13, November 2008. doi: 10.1111/j.1365-2966.2008.13802.x.

- [10] J. E. Dale, I. A. Bonnell, C. J. Clarke, and M. R. Bate. Photoionizing feedback in star cluster formation. *Monthly Notices of the Royal Astronomical Society*, 358:291–304, March 2005. doi: 10.1111/j.1365-2966.2005.08806.x.
- [11] J. E. Dale, J. Ngoumou, B. Ercolano, and I. A. Bonnell. Massive stars in massive clusters - IV. Disruption of clouds by momentum-driven winds. *Monthly Notices of the Royal Astronomical Society*, 436:3430–3445, December 2013. doi: 10.1093/mnras/stt1822.
- [12] J. E. Dale, J. Ngoumou, B. Ercolano, and I. A. Bonnell. Before the first supernova: combined effects of H II regions and winds on molecular clouds. *Monthly Notices of the Royal Astronomical Society*, 442:694–712, July 2014. doi: 10.1093/mnras/stu816.
- [13] B. T. Draine. Interstellar Dust Grains. *Annual Review of Astronomy and Astrophysics*, 41:241–289, January 2003. doi: 10.1146/annurev.astro.41.011802.094840.
- [14] C. P. Dullemond, A. Juhasz, A. Pohl, F. Sereshti, R. Shetty, T. Peters, B. Commercon, and M. Flock. RADMC-3D: A multi-purpose radiative transfer tool. Astrophysics Source Code Library, February 2012.
- [15] Melissa L. Enoch, Kaisa E. Young, Jason Glenn, II Evans, Neal J., Sunil Golwala, Anneila I. Sargent, Paul Harvey, James Aguirre, Alexey Goldin, Douglas Haig, Tracy L. Huard, Andrew Lange, Glenn Laurent, Philip

- Maloney, Philip Mautskopf, Philippe Rossinot, and Jack Sayers. Bolocam Survey for 1.1 mm Dust Continuum Emission in the c2d Legacy Clouds. I. Perseus. *Astrophysical Journal*, 638:293–313, February 2006. doi: 10.1086/498678.
- [16] II Evans, Neal J., Michael M. Dunham, Jes K. Jørgensen, Melissa L. Enoch, Bruno Merín, Ewine F. van Dishoeck, Juan M. Alcalá, Philip C. Myers, Karl R. Stapelfeldt, Tracy L. Huard, Lori E. Allen, Paul M. Harvey, Tim van Kempen, Geoffrey A. Blake, David W. Koerner, Lee G. Mundy, Deborah L. Padgett, and Anneila I. Sargent. The Spitzer c2d Legacy Results: Star-Formation Rates and Efficiencies; Evolution and Lifetimes. *The Astrophysical Journal Supplement Series*, 181:321–350, April 2009. doi: 10.1088/0067-0049/181/2/321.
- [17] A. Frank, T. P. Ray, S. Cabrit, P. Hartigan, H. G. Arce, F. Bacciotti, J. Bally, M. Benisty, J. Eislöffel, M. Güdel, S. Lebedev, B. Nisini, and A. Raga. Jets and Outflows from Star to Cloud: Observations Confront Theory. In Henrik Beuther, Ralf S. Klessen, Cornelis P. Dullemond, and Thomas Henning, editors, *Protostars and Planets VI*, page 451, January 2014. doi: 10.2458/azu\_uapress\_9780816531240-ch020.
- [18] Sam Geen, Joakim Rosdahl, Jeremy Blaizot, Julien Devriendt, and Adrienne Slyz. A detailed study of feedback from a massive star. *Monthly Notices of the Royal Astronomical Society*, 448:3248–3264, April 2015. doi: 10.1093/mnras/stv251.

- [19] D. J. Hollenbach and A. G. G. M. Tielens. Photodissociation regions in the interstellar medium of galaxies. *Reviews of Modern Physics*, 71: 173–230, January 1999. doi: 10.1103/RevModPhys.71.173.
- [20] D. H. Hughes, F. P. Schloerb, and LMT Project Team. The Large Millimeter Telescope. In *Revista Mexicana de Astronomia y Astrofisica Conference Series*, volume 35, pages 251–256, May 2009.
- [21] MIPS Instrument. Mips instrument support teams, 2011, mips instrument handbook, version 3. *SSC, Pasadena*.
- [22] C. R. Kerton, G. Wolf-Chase, K. Arvidsson, C. J. Lintott, and R. J. Simpson. The Milky Way Project: What are Yellowballs? *Astrophysical Journal*, 799:153, February 2015. doi: 10.1088/0004-637X/799/2/153.
- [23] Jeong-Gyu Kim, Woong-Tae Kim, and Eve C. Ostriker. Disruption of Molecular Clouds by Expansion of Dusty H II Regions. *Astrophysical Journal*, 819:137, March 2016. doi: 10.3847/0004-637X/819/2/137.
- [24] Sang-Hee Kim, P. G. Martin, and Paul D. Hendry. The Size Distribution of Interstellar Dust Particles as Determined from Extinction. *Astrophysical Journal*, 422:164, February 1994. doi: 10.1086/173714.
- [25] Christine M. Koepferl, Thomas P. Robitaille, James E. Dale, and Francesco Biscani. Insights from Synthetic Star-forming Regions. I. Reliable Mock Observations from SPH Simulations. *The Astrophysical Jour-*

*nal Supplement Series*, 233:1, November 2017. doi: 10.3847/1538-4365/233/1/1.

- [26] C. J. Lada. Cold outflows, energetic winds, and enigmatic jets around young stellar objects. *Annual Review of Astronomy and Astrophysics*, 23: 267–317, January 1985. doi: 10.1146/annurev.aa.23.090185.001411.
- [27] E. F. Ladd, P. C. Myers, and A. A. Goodman. Dense Cores in Dark Clouds. X. Ammonia Emission in the Perseus Molecular Cloud Complex. *Astrophysical Journal*, 433:117, September 1994. doi: 10.1086/174629.
- [28] Huixian Li, Di Li, Lei Qian, Duo Xu, Paul F. Goldsmith, Alberto Noriega-Crespo, Yuefang Wu, Yuzhe Song, and Rendong Nan. Outflows and Bubbles in Taurus: Star-formation Feedback Sufficient to Maintain Turbulence. *The Astrophysical Journal Supplement Series*, 219:20, August 2015. doi: 10.1088/0067-0049/219/2/20.
- [29] Frank J. Masci, Douglas I. Hoffman, Carl J. Grillmair, and Roc M. Cutri. Automated Classification of Periodic Variable Stars Detected by the Wide-field Infrared Survey Explorer. *Astronomical Journal*, 148:21, July 2014. doi: 10.1088/0004-6256/148/1/21.
- [30] Christopher D. Matzner. On the Role of Massive Stars in the Support and Destruction of Giant Molecular Clouds. *Astrophysical Journal*, 566: 302–314, February 2002. doi: 10.1086/338030.

- [31] S. Molinari, B. Swinyard, J. Bally, M. Barlow, J. P. Bernard, P. Martin, T. Moore, A. Noriega-Crespo, R. Plume, L. Testi, A. Zavagno, A. Abergel, B. Ali, P. André, J. P. Baluteau, M. Benedettini, O. Berné, N. P. Billot, J. Blommaert, S. Bontemps, F. Boulanger, J. Brand, C. Brunt, M. Burton, L. Campeggio, S. Carey, P. Caselli, R. Cesaroni, J. Cernicharo, S. Chakrabarti, A. Chrysostomou, C. Codella, M. Cohen, M. Compiègne, C. J. Davis, P. de Bernardis, G. de Gasperis, J. Di Francesco, A. M. di Giorgio, D. Elia, F. Faustini, J. F. Fischera, Y. Fukui, G. A. Fuller, K. Ganga, P. Garcia-Lario, M. Giard, G. Giardino, J. Glenn, P. Goldsmith, M. Griffin, M. Hoare, M. Huang, B. Jiang, C. Joblin, G. Joncas, M. Juvela, J. Kirk, G. Lagache, J. Z. Li, T. L. Lim, S. D. Lord, P. W. Lucas, B. Maiolo, M. Marengo, D. Marshall, S. Masi, F. Massi, M. Matsuura, C. Meny, V. Minier, M. A. Miville-Deschênes, L. Montier, F. Motte, T. G. Müller, P. Natoli, J. Neves, L. Olmi, R. Paladini, D. Paradis, M. Pestalozzi, S. Pezzuto, F. Piacentini, M. Pomarès, C. C. Popescu, W. T. Reach, J. Richer, I. Ristorcelli, A. Roy, P. Royer, D. Russeil, P. Saraceno, M. Sauvage, P. Schilke, N. Schneider-Bontemps, F. Schuller, B. Schultz, D. S. Shepherd, B. Sibthorpe, H. A. Smith, M. D. Smith, L. Spinoglio, D. Stamatellos, F. Strafella, G. Stringfellow, E. Sturm, R. Taylor, M. A. Thompson, R. J. Tuffs, G. Umana, L. Valenziano, R. Vavrek, S. Viti, C. Waelkens, D. Ward-Thompson, G. White, F. Wyrowski, H. W. Yorke, and Q. Zhang. Hi-GAL: The Herschel Infrared Galactic Plane Survey. *Publications of the Astronomical Society of*

*the Pacific*, 122:314, March 2010. doi: 10.1086/651314.

- [32] A. Mora, B. Merín, E. Solano, B. Montesinos, D. de Winter, C. Eiroa, R. Ferlet, C. A. Grady, J. K. Davies, L. F. Miranda, R. D. Oudmaijer, J. Palacios, A. Quirrenbach, A. W. Harris, H. Rauer, A. Collier Cameron, H. J. Deeg, F. Garzón, A. Penny, J. Schneider, Y. Tsapras, and P. R. Wesselius. EXPORT: Spectral classification and projected rotational velocities of Vega-type and pre-main sequence stars. *Astronomy & Astrophysics*, 378:116–131, October 2001. doi: 10.1051/0004-6361:20011098.
- [33] Fumitaka Nakamura and Zhi-Yun Li. Protostellar Turbulence Driven by Collimated Outflows. *Astrophysical Journal*, 662:395–412, June 2007. doi: 10.1086/517515.
- [34] M. Ntampaka, J. ZuHone, D. Eisenstein, D. Nagai, A. Vikhlinin, L. Hernquist, F. Marinacci, D. Nelson, R. Pakmor, A. Pillepich, P. Torrey, and M. Vogelsberger. A Deep Learning Approach to Galaxy Cluster X-ray Masses. *ArXiv e-prints*, art. arXiv:1810.07703, October 2018.
- [35] Stella S. R. Offner and Héctor G. Arce. Impact of Winds from Intermediate-mass Stars on Molecular Cloud Structure and Turbulence. *Astrophysical Journal*, 811:146, October 2015. doi: 10.1088/0004-637X/811/2/146.
- [36] J. E. G. Peek, Carl Heiles, Kevin A. Douglas, Min-Young Lee, Jana Grcevich, Snežana Stanimirović, M. E. Putman, Eric J. Korpela, Steven J.

- Gibson, Ayesha Begum, Destry Saul, Timothy Robishaw, and Marko Krčo. The GALFA-HI Survey: Data Release 1. *The Astrophysical Journal Supplement Series*, 194:20, June 2011. doi: 10.1088/0067-0049/194/2/20.
- [37] Thomas Preibisch and Hans Zinnecker. Deep Chandra X-Ray Observatory Imaging Study of the Very Young Stellar Cluster IC 348. *Astronomical Journal*, 122:866–875, August 2001. doi: 10.1086/321177.
- [38] Manuel A. Quijada, Catherine T. Marx, Richard G. Arendt, and Samuel H. Moseley. Angle-of-incidence effects in the spectral performance of the infrared array camera of the Spitzer Space Telescope. In John C. Mather, editor, *Optical, Infrared, and Millimeter Space Telescopes*, volume 5487 of *Society of Photo-Optical Instrumentation Engineers (SPIE) Conference Series*, pages 244–252, October 2004. doi: 10.1117/12.552061.
- [39] M. G. Rawlings, M. Juvela, K. Lehtinen, K. Mattila, and D. Lemke. Observations of 6–200  $\mu\text{m}$  emission of the Ophiuchus cloud LDN 1688. *Monthly Notices of the Royal Astronomical Society*, 428:2617–2627, January 2013. doi: 10.1093/mnras/sts233.
- [40] L. M. Rebull, K. R. Stapelfeldt, II Evans, N. J., J. K. Jørgensen, P. M. Harvey, T. Y. Brooke, T. L. Bourke, D. L. Padgett, N. L. Chapman, S. P. Lai, W. J. Spiesman, A. Noriega-Crespo, B. Merín, T. Huard, L. E. Allen, G. A. Blake, T. Jarrett, D. W. Koerner, L. G. Mundy, P. C. Myers, A. I. Sargent, E. F. van Dishoeck, Z. Wahhaj, and K. E. Young. The Spitzer c2d Survey of Large, Nearby, Interstellar Clouds. VI. Perseus Observed

- with MIPS. *The Astrophysical Journal Supplement Series*, 171:447–477, August 2007. doi: 10.1086/517607.
- [41] Naomi A. Ridge, Scott L. Schnee, Alyssa A. Goodman, and Jonathan B. Foster. The COMPLETE Nature of the Warm Dust Shell in Perseus. *Astrophysical Journal*, 643:932–944, June 2006. doi: 10.1086/502957.
- [42] T. P. Robitaille. HYPERION: an open-source parallelized three-dimensional dust continuum radiative transfer code. *Astronomy & Astrophysics*, 536:A79, December 2011. doi: 10.1051/0004-6361/201117150.
- [43] H. Rogers and J. M. Pittard. Feedback from winds and supernovae in massive stellar clusters - I. Hydrodynamics. *Monthly Notices of the Royal Astronomical Society*, 431:1337–1351, May 2013. doi: 10.1093/mnras/stt255.
- [44] B. D. Savage and J. S. Mathis. Observed properties of interstellar dust. *Annual Review of Astronomy and Astrophysics*, 17:73–111, January 1979. doi: 10.1146/annurev.aa.17.090179.000445.
- [45] D. Semenov, Th. Henning, Ch. Helling, M. Ilgner, and E. Sedlmayr. Rosseland and Planck mean opacities for protoplanetary discs. *Astronomy & Astrophysics*, 410:611–621, November 2003. doi: 10.1051/0004-6361:20031279.
- [46] R. J. Simpson, M. S. Povich, S. Kendrew, C. J. Lintott, E. Bressert, K. Arvidsson, C. Cyganowski, S. Maddison, K. Schawinski, R. Sherman,

- A. M. Smith, and G. Wolf-Chase. The Milky Way Project First Data Release: a bubblier Galactic disc. *Monthly Notices of the Royal Astronomical Society*, 424:2442–2460, August 2012. doi: 10.1111/j.1365-2966.2012.20770.x.
- [47] Van Oort et al. Identifying Feedback Signatures Using Convolutional Neural Networks. in prep.
- [48] Peng Wang, Zhi-Yun Li, Tom Abel, and Fumitaka Nakamura. Outflow Feedback Regulated Massive Star Formation in Parsec-Scale Cluster-Forming Clumps. *Astrophysical Journal*, 709:27–41, January 2010. doi: 10.1088/0004-637X/709/1/27.
- [49] Duo Xu and Stella S. R. Offner. Assessing the Performance of a Machine Learning Algorithm in Identifying Bubbles in Dust Emission. *Astrophysical Journal*, 851:149, December 2017. doi: 10.3847/1538-4357/aa9a42.