

**UNIVERSITY OF SCIENCE AND TECHNOLOGY OF HANOI**  
**UNDERGRADUATE UNIVERSITY**



**Group Project**

**Plankton Classification Platform for Supporting  
Research Activities of Environmental Scientists**

**Report 2**

Submitted by

LE Nguyen Khoi	USTHBI7-086
NGUYEN Xuan Bach	USTHBI7-021
LE Minh Anh	USTHBI7-007
DOAN Thanh Huy	USTHBI7-084
PHAM Minh Duc	USTHBI7-041

Under the guidance of:

Dr. NGHIEM Thi Phuong      ICT Lab

Hanoi, January 2019

# Table of Contents

<b>Table of Abbreviations .....</b>	<b>2</b>
<b>Abstract .....</b>	<b>3</b>
<b>I. Introduction.....</b>	<b>4</b>
<b>II. Theoretical Background and Literature Review .....</b>	<b>6</b>
<b>1. Theoretical Background .....</b>	<b>6</b>
<b>2. Literature Review .....</b>	<b>8</b>
<b>III. Scientific Methods .....</b>	<b>8</b>
<b>1. Machine learning and image classification .....</b>	<b>9</b>
<b>2. Deep Learning and Convolutional Neural Network. ....</b>	<b>9</b>
<b>3. Implemented CNN Architecture .....</b>	<b>12</b>
<b>4. Web User Interface.....</b>	<b>13</b>
<b>IV. Experimental Results and Discussion.....</b>	<b>14</b>
<b>1. Experimental setup.....</b>	<b>14</b>
<b>a. Dataset .....</b>	<b>14</b>
<b>b. Data Preprocessing.....</b>	<b>14</b>
<b>c. Training and Validating.....</b>	<b>15</b>
<b>2. Result and Discussion .....</b>	<b>15</b>
<b>V. Conclusion and Future Works .....</b>	<b>17</b>
<b>VI. Management of The Project at The End.....</b>	<b>18</b>
<b>References: .....</b>	<b>20</b>
<b>Acknowledgment.....</b>	<b>22</b>

## Table of Abbreviations

WHOI	Wood Hole Oceanographic Institution
HTML	Hypertext Markup Language
CSS	Cascading Style Sheets
SASS	Syntactically awesome style sheets
ANN	Artificial Neural Network
CNN	Convolutional Neural Network
SGD	Stochastic Gradient Descent
RMSprop	Root Mean Square Propagation
PHP	Hypertext Preprocessor

## **Abstract**

In our daily life, plankton has an important role that we cannot deny it in ecosystem and human life. Because these tiny organisms are the one that is almost controlled the carbon cycle of the Earth. In food cycle, they play the key role of the sustainability as the producer, the primary food for other animals. Even some specific kinds of plankton can be used as medicine for some severe diseases such as cancer. However, even though plankton can cure cancer but for the time being human still not find the solve to it. Moreover, due to the climate change, the population of plankton varies significantly and with controlling. This will cause many problems in the end, such as red tide phenomenon, fish-quantity fall and even economical-relating problems. In this project, the plankton classification platform was built. There are two main part: model development and website development. The model part allows putting unlabeled plankton images to their type of plankton that they belong to. And the website part will provide environmentalists with a friendly web user interface to perform classification task resulted from the model. From the experimental test with a domain in our datasets, we have produces good results after many failed train data, we have the accuracy over 80%.

## **I. Introduction**

Plankton is a diverse group of organisms that live in the oceans and cannot swim against the water current [1]. The existence of plankton is crucial to our ecosystem and also in our human life: they give the most abundant supply of primary productivity and take up to nearly half of the total carbon fixed in global carbon cycle [2]. Moreover, plankton is the major food source of fish and oceanic animals and play significant role in ensuring the health and balance of the oceans. So that we can see, plankton have an important role that we cannot denied.

However, because the earth's temperature is getting warmer from day to day. It causes the amount of plankton in some oceanic areas decreases significantly and while in others oceanic areas, there is an explosion in the number of them. It got up an emerging problem of plankton where the quantity of plankton becomes imbalance. The problem that plankton becomes imbalance can cause us a lot of problems that we cannot anticipate. First of all, loss of plankton populations in some oceanic areas could result in ecological upheaval and negative social impacts. And the uncontrolled rise of some specific plankton in other oceanic areas will lead millions of fish and oceanic animals killed. This problem shows that it is important to study plankton's population to improving the health of the ocean.

In order to study plankton's population that is imbalance nowadays because of the ecosystem changed. Underwater camera systems are developed to capture microscopic, high-resolution images of plankton over large study areas. As a result, it brought us a huge amount of plankton images that we take them underwater. However, it's raises us a key question that how we can analyze the collected plankton images in order to assess the population and distribution of each plankton specie. A natural approach to this question is to manually classify collected plankton images for further human analysis This approach, however, is infeasible since it would be extremely time and cost consuming to manually classify millions of images [3]. Furthermore, this is not all the problem that the question come up. Another common problem is the difficult use of plankton images in research and education. To the detail, we tell that what happened if there is a new picture of plankton came up and a

student or researcher cannot find any information or know what the plankton belong to which domain, or which type.

In order to overcome with this problem, the project has proposed two main goals:

- (1) To develop a classification platform to automatically put new plankton images in their correct categories;
- (2) To build a plankton website so that the environmental scientists can contribute plankton data for easier use of plankton images in research and education.



Figure 1. An Image of Plankton

From working with the project, we had specified our objectives that we must done and learned that listed below:

- This project is about making a website using knowledge that we learned in the university.
- To develop teamwork skill with the use of tools like Github and Google Doc.
- To learn proess of making a project.
- Build a groundwork that can be reused and expanded in the future.
- Understand the basic background about machine learning and image classification.

- Successfully build and develop a part of website that image can be upload to website data, and website can detect that image belong to which type of plankton.
- Apply image classification models for the case of big datasets.
- Apply image classification models to build plankton classification platform.
- Develop web user interface for environmentalist to use plankton classification platform and let people also can contribute plankton data for easier use of plankton images in research and study

The main part of the report contains in 5 sections below. Here is the summary of all the sections.

Section 2: Gives a brief introduction about the background and the state of the art of our work.

Section 3: Describes our method and techniques that we used to build our proposed plankton classification platform for supporting research activities of environmental scientists.

Section 4: Show and explains the result obtained from our experiment. Also, our comments on them.

Section 5: Our conclusion on our work and possible future develop on our project.

Section 6: Shows how we manage our project and what we had done until the end.

## **II. Theoretical Background and Literature Review**

In this section, we will provide the knowledge that we must know in order to build our proposed plankton classification platform. Because it's essential for us to have some knowledge about plankton. So that we can understand about it, and about what kind of plankton belong to which domain, which type.

### **1. Theoretical Background**

Plankton distribution varies horizontally, vertically and seasonally due to two main factors: sunlight and nutrient availability. Sunlight availability is the reason why plankton is abundant in tropical areas and in 200-300m under the surface water. Nutrient such as nitrate, phosphate and silicate govern the growth and development of plankton. There, ocean areas with plenty of nutrients will attract more plankton. Knowing plankton habitat allows capturing plankton images more efficiently.

From what we understand about plankton, it is divided into many subtypes. There are two main ways to classify plankton: by functional groups or by their size. So, in our project we have 6 kind of plankton domain. They are Diatoms, Flagellates, Dinoflagellates, Ciliates, Coccolithophore and Miscellaneous. The detail of each kind plankton we will show below:

About Diatoms, they are a major group of algae, specifically microalgae, found in

the oceans, waterways and soils of the world. Living diatoms number in the trillions: they generate about 20 percent of the oxygen produced on the planet each year, take in over 6.7 billion metric tons of silicon each year from the waters in which they live [4], and contribute nearly half of the organic material found in the oceans. The shells of dead diatoms can reach as much as a half mile deep on the ocean floor, and the entire Amazon basin is fertilized annually by 27 million tons of diatom shell dust transported by east-to-west transatlantic winds from the bed of a dry up lake once covering much of the African Sahara [5].

About Flagellate, we can understand that a flagellate is a cell or organism with one or more whip-like appendages called flagella. The word flagellate also describes a particular construction characteristic of many prokaryotes and eukaryotes and their means of motion. The term presently does not imply any specific relationship or classification of the organisms that possess flagellae. However, the term "flagellate" is included in other terms which are more formally characterized [6].

In Dinoflagellata, the dinoflagellates are a classification subgroup of algae. They are a large group of flagellate eukaryotes that constitute the phylum Dinoflagellata. Most are marine plankton, but they also are common in freshwater habitats. Their populations are distributed depending on sea surface temperature, salinity, or depth. Many dinoflagellates are known to be photosynthetic, but a large fraction of these are in fact mixotrophic, combining photosynthesis with ingestion of prey (phagotrophy) [7]. In terms of number of species, dinoflagellates are one of the largest groups of marine eukaryotes, although this group is substantially smaller than diatoms [8]. Some species are endosymbionts of marine animals and play an important part in the biology of coral reefs. Other dinoflagellates are unpigmented predators on other protozoa, and a few forms are parasitic. Some dinoflagellates produce resting stages, called dinoflagellate cysts or dinocysts, as part of their lifecycles.

Next is the ciliates are a group of protozoans characterized by the presence of hair-like organelles called cilia, which are identical in structure to eukaryotic flagella, but are in general shorter and present in much larger numbers, with a different undulating pattern than flagella. Cilia occur in all members of the group (although the peculiar Suctorians only have them for part of their life-cycle) and are variously used in swimming, crawling, attachment, feeding, and sensation.

Furthermore, we will talk about coccolithophore. A coccolithophore is a unicellular, eukaryotic phytoplankton (alga). They belong either to the kingdom Protista, according to Robert Whittaker's Five kingdom classification, or clade Hacrobia, according to the newer biological classification system. Within the Hacrobia, the coccolithophorids are in the phylum or division Haptophyta, class Prymnesiophyceae (or Coccolithophyceae). Coccolithophorids are distinguished by special calcium carbonate plates (or scales) of uncertain function called coccoliths, which are also important microfossils. However, there are Prymnesiophyceae species lacking



coccoliths (e.g. in genus *Prymnesium*), so not every member of Prymnesiophyceae is coccolithophorid. Coccolithophores are almost exclusively marine and are found in large numbers throughout the sunlight zone of the ocean.

And last but not least, Miscellaneous is a domain that contain a huge amount of kinds of plankton that not belong to any kind of domain that we are not defined or defined but not in any kind of domain that listed above.

Which the detail that we understand from every kind of domains and type plankton that we talk about, we can use it and write the description of each kind of plankton to our plankton classification platform. Also, it will allow us to better characterize plankton images of building classification platform.

## **2. Literature Review**

Nowadays, there are many kinds of methods in the literature to classify plankton around the internet. One of them that we cannot discuss about is an automated platform for classifying and monitoring phytoplankton that using scanning flow cytometer. This cytometer was designed to analyze the natural properties of plankton. For example, the cytometer is capable of analyzing from small to large planktonic particles (from 1 to 700  $\mu\text{m}$  in diameter and a few mm in length) and relatively large water volumes. The platform manipulates and visualizes data using R programming language. The datasets also include original sampled volume, date, time, and depth at which particles were taken.

This platform was proposed by Pomati et al, 2011 [9]. With the idea of a platform that it is able to extract very detail information of planktonic particles for further analysis is very interesting. However, if students or researchers want to have a large amount of plankton datasets for large-scale studies, the platform shows the limit in gaining high efficiency to classify the plankton. This is because the platform has not yet taken into account the advanced data analysis techniques to handle big datasets. So, to solve the problem difficult use of plankton images for research and education, we come up with an idea that we create a website, an online plankton analysis system (PAS) for environmentalist to use plankton classification platform. Basically, the system is a web application allowing marine biologists to classify their plankton images using a web user interface. PAS provides the functions for experts to upload their images, process them, extract features, hand-label image, train classifiers, and use those classifiers to automatically label new images [10]. So that people can contribute plankton data in order to help other people want to research and study. And to make thing much easier than Pomati et al, 2011 [9] platform by we apply the modern machine learning methods to our website, namely convolutional neural network, in order to gain high performance computing when analyzing large datasets.

## **III. Scientific Methods**

From what we mentioned above that our motivation in this work is to build a classification platform to put new plankton images into their correct categories. In the following, we will present in detail the machine learning methods we use to perform this motivation.

## **1. Machine learning and image classification**

According to Michalski et al, 2013 [11] machine learning is a field of study that provides computers with ability to learn from experience. To the detail, that we talked about that a computer program is said to learn from experience if it can teach itself to grow and change having new data. With the detail, we can tell that machine learning focuses on the development of computer algorithms for transforming data into appropriate action. For example, that make it easier to understand how the machine learning environment work. In the environment, there are available data, statistical models and computing power are simultaneously interacted with each other. As a result, from that, growth in data causes the need of additional computing power, which in turn spurs the development of statistical method for analyzing large datasets. This make it become a loop cycle of advancement allowing even larger and more interesting data to be collected.

Image classification refers to the task of classifying images into appropriate classes or categories that they belong to. In other word, we can say that image classification algorithms perform the press of putting unlabeled images into their class, type or categories that they belong to.

In the following, we will present our study on the deep learning and convolutional neural network as the image classification models.

## **2. Deep Learning and Convolutional Neural Network.**

### **Deep Learning**

In practical term, deep learning is the subset of machine learning, from the technical to function, it's similar to machine learning. However, its capabilities are different. Basic machine learning models do become progressively better at whatever their function is, but they still need some guidance. For example, if we train our model base on a machine learning algorithm then put it on prediction phase, and when it returns an inaccurate prediction after the prediction phase finish, we have to step in and make adjustments by our own hand. But deep learning algorithm making more impact than machine learning algorithm, the algorithms can determine on their own if prediction is accurate or not. One of most famous algorithms to continually analyze data with a logic structure how human would draw conclusions is Artificial Neural Network (ANN), commonly known as "neural network", is a machine designed to model the way in brain performs a particular task or function of interest [12]. The design of ANN is inspired by the biological neural network of the human brain. This makes for machine intelligence that's far more capable than that of standard machine

learning models. However, we do not go to the detail in ANN, instead we will go to in the study of Convolutional Neural Networks (CNN), a far more model better than ANN that we used in our model training plankton image dataset.

## Convolutional Neural Networks

Convolutional Neural Networks (CNN) is a variant of multilayer perceptron (MLP) introduced by LeCun et al, 1998 [13] to overcome the weakness of MLP. To the detail, MLP is feedforward of ANN in which each layer is fully connected to the next one. MLP is often used for classification task, where it start with an input layer, transforms through several hidden layer, and finally ends at the output layer. The output layer contains the number of neurons equal to the number of classes required for classification. So, using MLP to classify images often raises the problem of wasteful to connect each neuron in a hidden layer to only a local region of the input layer. This connection is executed by using a convolution filter over the input layer and the hyper-parameter of such region is called the receptive field of the neuron. The convolution filters act as local filters over the input space and are well-suited to exploit the strong spatially local correlation in the natural images.

Compare to other neural network, CNNs have more layers to approach the high dimension data. As you can see in the Figure 2 above is an example that CNNs have more layers than a single line like other neural network. To go deeper in neural network, CNNs have many of neurons which share the same parameterization (weight vector and bias). Replicating neurons in this way allows for features to be detected regardless of their position in the input space. Moreover, weight sharing increases learning efficiency by greatly reducing the number of free parameters being learnt. The constraints on the model enable CNNs to achieve better generalization on the model enable CNNs to achieve better generalization on vision problems. The figure 2 below will show the example about how CNN work.

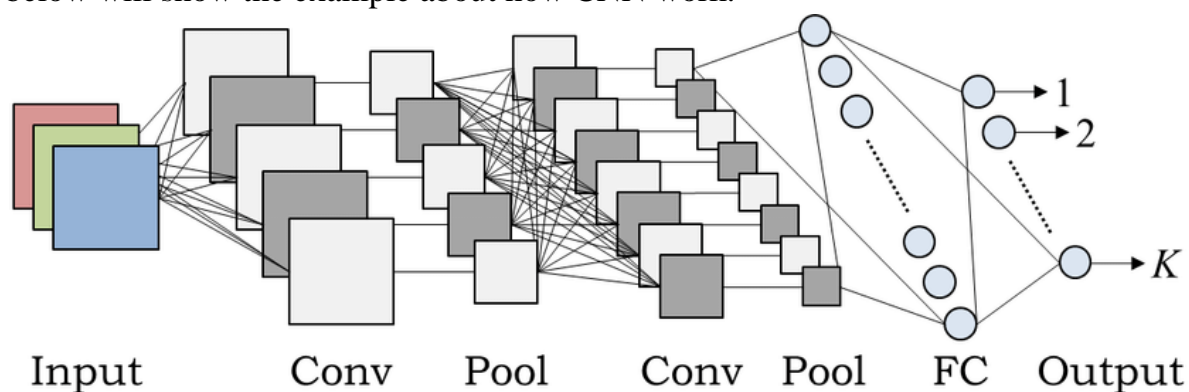


Figure 2. An example of CNN architecture.

Next, we will to the detail about 3 kind of layers that had shown in example of CNN architecture in figure 2 above.

## Convolutional Layer

Convolution is a mathematical operation that's used in single that's used in single processing to filter signals. In a convolutional layer, all the neurons will apply to convolution operations to the inputs, namely convolutional neurons. The most important parameter in a convolutional neuron is the filter size. So that it makes the convolutional layer is the core building block of a CNN. Moreover, from the way that convolutional layers work, each filter is convolved across the width and height of the input volume, computing the dot product between the entries of the filter and the input and producing a 2D activation map. From that, we will have a result that the network learns filters that will activate when it detects some specific type of feature at some spatial position in the input.

## Pooling Layer

Pooling layer is a nonlinear layer, that work with width and height of the image from the dataset that we put on to train. It will perform a down sampling operation on. As we can see, from that we have a result that the image volume that we put in is reduced. This means that if some feature has already been identified in the previous convolution operation, that a detailed image is no longer needed for further processing, and it is compressed to less detailed picture.

So, from what we mentioned above, we can see that the pooling layer is mostly used immediately after the convolutional layer to reduce the spatial size. The most common form that pooling layer always use is Max pooling where we take a filter of size  $F \times F$  and apply the maximum operation over the  $F \times F$  sized part of the image.

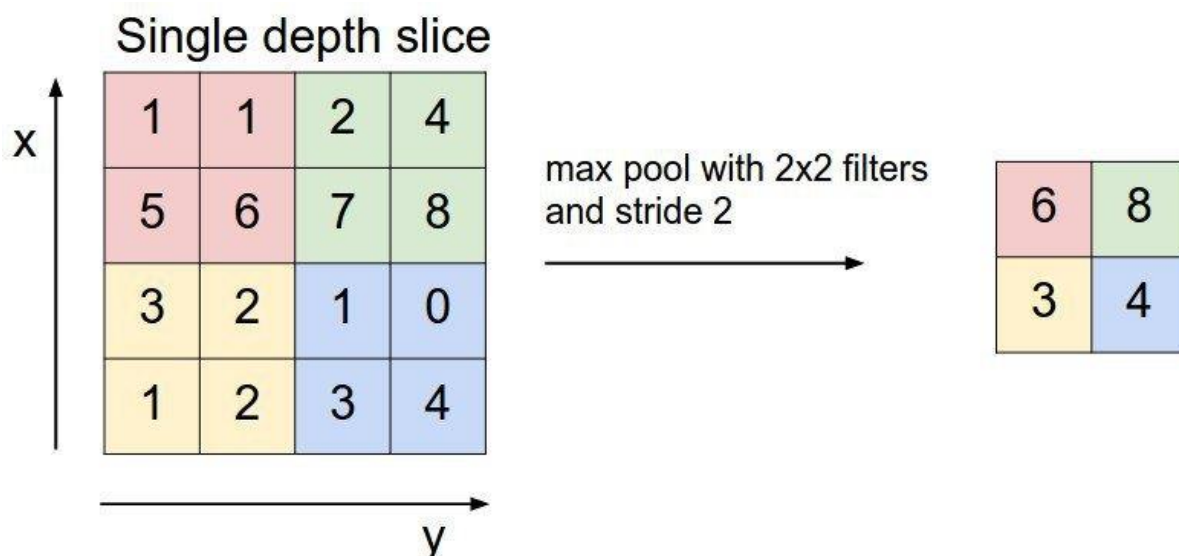


Figure 3. An example about Max pooling

From the figure about we can see that max pool size  $2 \times 2$  with a stride of 2. Because with that we can calculate the formula that we will show below.

For example, that we have our input is of size  $w_1 \times h_1 \times d_1$  and the size of the filter

is  $F \times F$  with stride  $S$ . Then the output sizes  $w_2 \times h_2 \times d_2$  will be:

$$w_2 = (w_1 - f)/S + 1$$

$$h_2 = (h_1 - f)/S + 1$$

$$d_2 = d_1$$

So that with the size  $2 \times 2$  with a stride of 2, it essentially reduces the size of input by half. That is the reason why we not take the average in place of taking maximum, because it will become average pooling and it will not make better result than max pooling.

### **Fully Connected Layer**

From what we saw in figure 2, fully connected layer is the last phase of CNN architecture. It connected every neuron in one layer to every neuron in another layer. It is in principle the same as the traditional multi-layer perceptron neural network (MLP). The flattened matrix goes through a fully connected layer to classify the images, so that we had the output.

### **3. Implemented CNN Architecture**

Our training model was based on tatsuyah's CNN Image Classifier Model that we find out and learn from GitHub (<https://github.com/tatsuyah/CNN-Image-Classifier>). However, in the original model, they used RMSprop optimization algorithms also known as Root Mean Square Propagation to model compile. The RMSprop optimizer restricts the oscillations in the vertical direction. Therefore, with it, we can increase our learning rate and our algorithm could take larger steps in the horizontal direction converging faster. However, RMSprop is not as good as Adam because Adam is a combination of RMSprop and Stochastic Gradient Descent with momentum. It uses the squared gradients to scale the learning rate like RMSprop and it take advantage of momentum optimizer algorithm by using moving average of the gradient instead of gradient itself like SGD with momentum. So, it can process fragmented data more effective than RMSprop. We can see how the different that why we choose optimization more in the section in Result and Discussion section.

Inspired by Maxout Network [14], we configure in our CNN models to five models. In each convolutional layer we set filter are 32, 64, 128, 256 and kernel size is (3,3) which determinate the width and height of the 2D convolution window. So that models can extract more features in the input shape to classify the plankton more correctly, because in plankton, there are many kinds of plankton that look very similar together and we cannot know what plankton belong to what type, so to deal with the problem, we come up with a solution that we must extract more features in the input shape so that computer can detected every kind of plankton and put it to their correctly categories. Like we had mentioned above, with the max pooling have size  $2 \times 2$  reduces the number of parameters within the mode and to generalizes the results from

a convolutional filter - making the detection of features invariant to scale or orientation changes. Also, with these CNN structures for the purpose of reducing data dimension: the higher nodes of CNN will produce the fewer dimension of data before go to Fully connected layer. To the detail of our implemented CNN design, it will be shown on the table below:

Model 1	Model 2	Model 3	Model 4	Model 5
32 nodes convenet	32 nodes convenet	32 nodes convenet	32 nodes convenet	32 nodes convenet
	64 nodes convenet	64 nodes convenet	64 nodes convenet	64 nodes convenet
	128 nodes convenet	128 nodes convenet	128 nodes convenet	128 nodes convenet
		256 nodes convenet	128 nodes convenet	128 nodes convenet
			256 nodes convenet	256 nodes convenet
			256 nodes convenet	256 nodes convenet
1024 fully	1024 fully	1024 fully	1024 fully	1024 fully
				2048 fully

Table 1. Models used in the platform

#### 4. Web User Interface

Since our primary aimed like we mentioned on the title plankton classification platform for supporting research activities of environmental scientists, it is more convenient to provide them with a graphical user interface to use our classification platform. For this, we develop a website acting as a web user interface for our plankton classification platform by divided it to two part: Client-side and Server-side.

A detail of our design and implementation is presented in the next following.

##### Client-Side

In web development, “client-side” refers to everything in a web application that is displayed or takes place on the client. This includes what user sees, such as text, images, and the rest of the UI, along with any actions that an application performs within that user’s browser. So, in order to build our front-end for users to interact directly, we use the skill that we had studied from the course web application. We use HTML, CSS, Javascript, JQuery to design our web framework. However, in CSS part, there are too much style to designed our website and it make our CSS part became complicated. To solve this problem, we decided to use a high level of CSS technique, namely SASS. Because SASS can provide us variable to set same like Javascript and SASS boasts more features, abilities than any other CSS that we used to work with

when we making website interface.

### **Server-Side**

Much similar to “client-side”, “server-side” also has an important in web development. However, not like “client-side”, that everything happened in “server-side” is only happened on the server, instead of on the client. In order to create a back-end to work with our website. We use PHP is our “server-side” scripting language, a general-purpose programming language originally designed for web development. It was originally created by Rasmus Lerdorf in 1994 [15]. It’s suited for web development and can be embedded in HTML part, PHP is integrated with a number of popular databases, and our most use database that we use for our website is MySQL database. It is one of the most popular open source databases in the world nowadays that very fit to use on web development. Also, a one more reason that why we choose MySQL database to work with because MySQL is friendly to PHP which is the server-side language that used in this project.

## **IV. Experimental Results and Discussion**

### **1. Experimental setup**

#### **a. Dataset**

Our data presented here are subset of a larger plankton imagery data set contain millions of images of microscopic marine plankton, organized according to category labels provided by researcher at the Wood Hole Oceanographic Institution (WHOI). The images are currently placed into 6 Domains of plankton: Diatoms, Flagellates, Dinoflagellates, Ciliates, Coccolithophore, Miscellaneous. To the detail, in each domain, there are many types of plankton. There are total 103 types of plankton in 103 folders. These images were collected in situ by automated submersible imaging-in-flow cytometry with an instrument called imaging FlowCytobot (IFCB) [16]. Our data that we use on this project take from 2010 to 2014 (There are total 1,980,000 images in 103 folder)

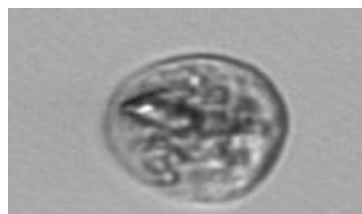


Figure 4. An image of plankton in our data set

In our plankton data set that we extracted from the raw data were sorted into 103 plankton classes, so that we split 6:3:1 into train, validation and test datasets.

#### **b. Data Preprocessing**

In our dataset, there are many images with different resolution, the smallest ones are less than 50 x 50 pixels, whereas the largest ones are up to 1000 x 1000 pixels, even some images will come up with 500 x 1000 pixels or 148 x 200 pixels. Therefore, we decided that we rescaled in fixed dimension: 128 x 128 pixels in order to unified all images on a resolution. To improve the stability and convergences of the model.

However, there are some folder of plankton datasets only contain few images or less than 200 images. With this problem, this will cause us overfitting our training model. So, we used data augmentations to handle randomly during training for preventing overfitting. The parameters of data augmentations are presented below:

- Rotation: random with angle between 0 – 180 degree
- Width shift and height shift: randomly translate pictures 0.2 vertically or 0.2 horizontally
- Rescale: Target values between 0 and 1 by scaling with a 1/255
- Shear range: randomly applying 0.2 shearing transformations
- Zoom range: 0.2 for randomly zooming inside pictures
- Horizontal flip: randomly flipping half of the images horizontally - relevant when there are no assumptions of horizontal asymmetry. We set to True.
- Fill mode: it's a strategy used for filling in newly created pixels, which can appear after a rotation or width/height shift. We set nearest.

### **c. Training and Validating**

Like we just mentioned in section b we divided our data set with the ratio is 6:3:1. The training process with bath size is 8, and samples that we put to train per epoch is 1000 and validation steps is 500. The learning rate is initialized at 0.0001 and a learning rate decay of 0.0001 is applied after each epoch. So that we can gain more accuracy when training plankton datasets.

Accuracy = total true sample / total sample

## **2. Result and Discussion**

### **Effect of Number of Neural Network Layers:**

Overall, the accuracy of all neural network training until now have passed 80% (see the table below) and the time of training also increase with the number of layers as a function like  $n^a$ . As you can see, from the table below, we achieve the best result with model 5 and the best accuracy that we gain from our training data is 83.3% with our training CNN model. From the list table below, through time, we increase the number of neural network layer, it gains more accuracy. As increasing the number of convolutional layers like the solution that we mentioned above, in order to extract more features in the input shape so that computer can detected every kind of plankton



and put it to their correctly categories. However, it will end up overfitting model to the training data and it will not do well with unseen data. So, we did not want try to increase more layer, since we had achieved a good result.

Number of neural network layer	Time (hr)	Accuracy (%)
2	7.1	73.1
3	7.8	76.4
4	8.1	80.1
5	8.5	83.3

Table 2. Result table

Like we mentioned above that we use different algorithm than the original model than we find and base on it on GitHub. We had tested with the RMSprop with learning rate we set is 0.0001, it's only gain us 65% accuracy. However, when we change the optimization algorithm to Adam, it had better result, that is result we shown it on Table 2 – Accuracy reached 83.3% in the last result with 5 neural network layers.

#### **Plankton characteristics classified by the platform:**

As the result is scientific name of plankton, more characteristics of that plankton can be known. According to Jane M. Browles [17] and Jerry G. Johnson [18], the scientific name is in Latin name and consists of two parts: the genus or generic name and the specific epithet. So that, from the name of plankton, the properties of plankton can be almost understood and this is very useful for further research.

#### **Website interaction:**

Our website has two main function, they are a news website about plankton and a huge gallery that people can access and contribute their work to our data through the website. To the detail, with the part news website, people will just use it for update their information about new story, or new discovery about new kind of plankton or something else that relate to plankton. And about the gallery part, people will access to it to find what their plankton images that they want to see, also, people can upload their image about plankton to this website with only few easy steps.

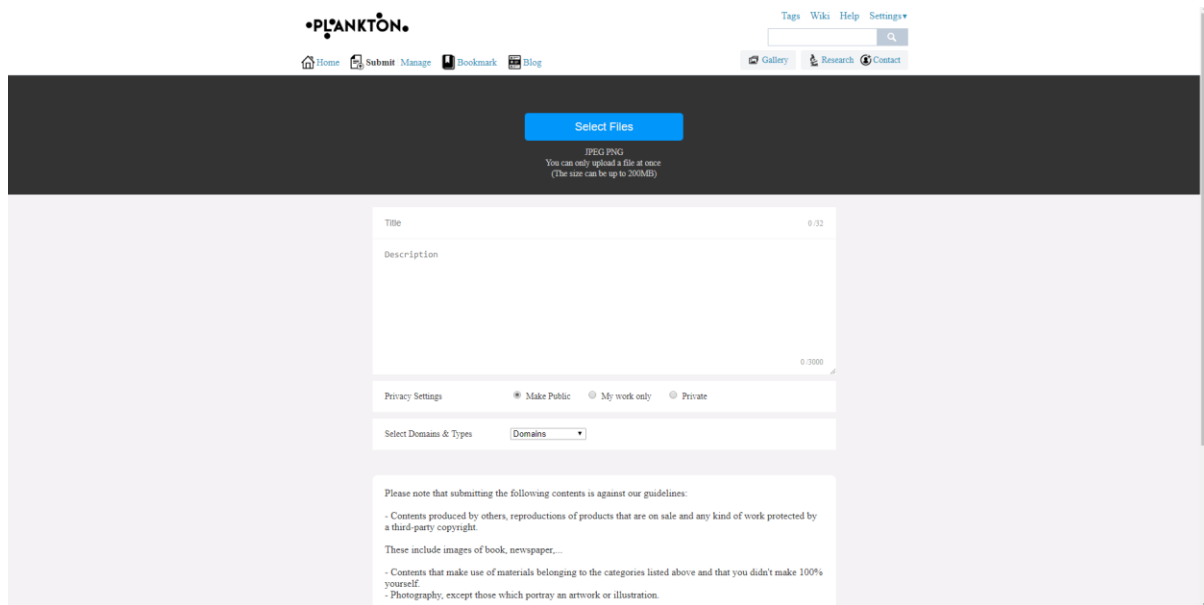


Figure 5. Web interface

If the image that already that exist in the website database, the website will immediately warn the user that this image had already had on our data, so it can be uploaded to our data.

Also base on “intellectual property”, we applied some law to our website to write term of use. In order let people understand that people cannot reup images that already exist in our database.

## V. Conclusion and Future Works

In the project, the accuracy of the experiment that until now, we have is over 80%. However, that only the number that we use for only a domain in our data. We haven't yet used all the domain in WHOI's plankton dataset for training. We are trying to make it work for real at the moment. But to comparing with the classification rate achieved by machine learning method to human accuracy that over 80% is a fine number. And still there still much problem that we do not finish yet. Because our main focus is CNN model, so there is still much more thing to do website.

For the future work, this would be possible to improve the activities of the website. For the time being, we only make it can be uploaded the image to our website and move it to our database so people can used it. In the future, we also want that people from outside also can write article and upload it to blog page in our website. So, people can learn something from the point of view of the people wrote it and uploaded to our website.

## VI. Management of The Project at The End

Activities	Plan	Expected results	Done
Research of Plankton	+) Understand overall characteristic and a bit knowledge of plankton +) Understand 6 domains of plankton +) Prepare the content about plankton	+) Understand overall characteristic and a bit knowledge of plankton +) Understand 6 domains of plankton +) Prepare the content about plankton	+) Understand overall characteristic and a bit knowledge of plankton +) Understand 6 domains of plankton +) Prepare the content about plankton
Preparing detailed information about plankton and code	All required software is installed	Everything is ready for testing	+) Collected and checked all the articles and documents +) Downloaded database (103 types of plankton) from Wood Hole Oceanographic Institution (WHOI)
Model Development	+) Data exportation +) Image preprocessing + Apply machine learning for building a model	The model has accepted classification score	The model reaches up to over 80% (For only one domain that we are now testing, not the whole data) – It's a good result, since there are many mistakes that we had fixed.
Website Development	+ Design website of plankton classification	+) Simple and responsive user interface +) Can use pre-built model to recognize plankton	Website is simple and friendly to user
Writing report and preparing for the final	+) Preparing the slides for presentation. +) Practicing the	+) Completed on time +) Well prepared the slides for	+) Well checked before submitting report. +) Prepared well for

presentation	presentation +) Think and ready for every possible question that may come up +) Writing report	presentation +) All important information is covered +) Come up with many questions that may asked during presentation	presentation +) All important information is covered +) Be sure all members understand well their own task and of others
--------------	--	--	--

## References:

- [1] Castro Pond, Stephen, and George L. Pickard. Introductory dynamical oceanography. Elsevier, 2013
- [2] Paul Falkowski, the power of plankton, Nature, Vol.483, March 2012.
- [3] Culverhouse, Phil F., et al. "Do experts make mistakes? A comparison of human and machine identification of dinoflagellates." Marine Ecology Progress Series 247.17-25 (2003):5.
- [4] Treguer, P.; Nelson, D. M.; Van Bennekom, A. J.; Demaster, D. J.; Leynaert, A.; Queguiner, B. (1995). "The Silica Balance in the World Ocean: A Reestimate". Science. 268 (5209): 375–9.
- [5] Bristow, C.S.; Hudson-Edwards, K.A.; Chappell, A. (2010). "Fertilizing the Amazon and equatorial Atlantic with West African dust". Geophys. Res. Lett. 37 (14): L14807.
- [6] Cavalier-Smith T. (1995). "Zooflagellate phylogeny and classification". Tsitologiya. 37 (11): 1010–29.
- [7] Stoecker DK (1999). "Mixotrophy among Dinoflagellates". The Journal of Eukaryotic Microbiology. 46 (4): 397–401.
- [8] Guiry MD (October 2012). "How many species of algae are there?". Journal of Phycology. 48 (5): 1057–63.
- [9] Pomati, Francesco, et al. "An automated platform for phytoplankton ecology and aquatic ecosystem monitoring. "Environmental science & technology 45.22 (2011): 9658-9665.
- [10] Plankton Analysis System, <http://vis-www.cs.umass.edu/~pas/>, accessed February 2016.
- [11] Michalski, Ryszard S., Jaime G. Carbonell, and Tom M. Mitchell eds. Machine learning; An artificial intelligence approach. Springer Science & Business Media, 2013.
- [12] Patnaik, A., and R. K. Mishra. "ANN techniques in microwave engineering." Microwave Magazine, IEEE 1.1 (2000): 55-60
- [13] LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. Gradient-based learning applied to document recognition. Proceedings of the IEEE, 86(11):2278-2324,

November 1998.

[14] Goodfellow, Ian J., Warde-Farley, David, Mirza, Mehdi, Courville, Aaron, and Bengio, Yoshua. Maxout networks. In ICML, 2013.

[15] History of PHP, php.net

[16] Kate Madin, Building an Automated Underwater Microscope, Oceanus magazine, Woods Hole Oceanographic Institution (WHOI)

[17] Jane M. Browles, The Name of Plants, 1996

[18] Jerry G. Johnson, The Classification of Living things, 2005

## **Acknowledgment**

We would like to express our heartfelt appreciation to Dr. Nghiem Thi Phuong for the patience and enthusiasm that she guided us with during the course of this project. Without your help, would not be able to complete our work successfully.

Our gratitude also to the staff of Information and Communication Technology Department and ICT Lab for their valuable assistance.