# Sentiment Analysis Using Parallel Computing Through GPU

**Harshita Mandloi, Shraddha Masih**

School of Computer Science and IT, DAVV, Indore, Madhya Pradesh, India

## ABSTRACT

Parallel Computing is becoming important in the field of computer science and is proven as a high-performance solution. Over the couple of years, GPU has gained an important place in the field of high-performance computing. Social media is expanding at present and becoming important in society. Social network sites allow users to communicate with people in the network by sharing posts, images, videos, status. The proposed system gathers the information from the social media websites and performs the sentiment analysis on the social media data using GPU. The work concentrates on recognizing the sentiment information from the text reviews and using that to identify the items. The aim of this paper is to do analytics on social media data. Analysis is done on the data using K Nearest Neighbor algorithms and Support Vector Machine algorithm on the GPU.

**Keywords:** Parallel Computing, GPU, Social media, K nearest neighbor, SVM

## I. INTRODUCTION

Nowadays social media, blogs, and other media produce a huge amount of data on the world wide web. This huge amount of data contain a critical opinion related to the information that can be used to benefit business and other aspects of commercial and scientific industries. Analysis of user sentiments can help in a business decisions. Sentiments can be categorized into 3 types:- Positive, Negative and Neutral.

Graphics Processing Unit has shown better computational performance as compared to the current main-stream multicore Central Processing Unit. Apart from graphics and multimedia processing, high-performance GPUs are also mapped to take care of General Purpose Computing.

Machine learning is the current trend in computing that is focused on designing and developing algorithms which allows computers to learn. It can capture characteristics of interest in order to make a prediction for a new data query. The gathered data is considered as training data which illustrate the relationships among the observed variables. Many important patterns can be recognized after applying the learning procedure. Supervised learning is the data mining task of classifying data into labeled data. This work is mainly focused on the classification tools- Support Vector Machine and K Nearest Neighbor. To improve the performance of this algorithm GPU is used with CUDA. Parallel SVM algorithm is implemented on GPU using CUDA framework. The implementation of parallel SVM achieves a great performance using GPU. This software utilizes the parallelism in both data level and task level to maximize the performance of the GPU.

## II. BACKGROUND

**A. Parallel Computing-** Parallel computing is becoming important day by day. Dividing large task into sub task and assigning these tasks to multiprocessor and executing them concurrently is called parallel processing. Parallel computing is used for high performance computing in the field of data analytics. For the high performance computing the

Graphic Processing Unit (GPU) has played an important role because of the low cost and parallel processing power.

**B. Social Media-** Social media is the changing the way of the people find information, share knowledge and communicate with each other. The most popular social networking sites are Facebook, Twitter, yahoo etc. The term social media refers to as:

**Social:** It refers to connect with other people by sharing information with them and receiving information from them.

**Media:** It refers to an instrument of communication like internet, TV, newspaper etc.

**C. Sentiment Analysis-** Sentiment Analysis is also called opinion mining. It is an analysis of the feeling, emotions, attitude. Sentiment analysis is the process of determining whether a piece of script is positive, negative or neutral. It is generally used in social media to check as it allows us to gain an overview of the public opinion behind certain topics. Its application is larges and powerful, the ability to extract insight from social data is a practice that is being widely adopted by organizations across the world.

**Different Classes of Sentiment Analysis**

**i. Positive Sentiment:** It refers to the positive attitude of the speaker about the text.

**ii. Negative Sentiment:** It refers to the negative attitude of the speaker about the text.

**iii. Neutral Sentiment:** In this no emotion are reflected about the text.

**D. Graphic Processing Unit-** Graphic Processing Unit is a solitary chip processor mainly used to manage and boost the performance of videos and graphics. A CPU consist of few cores while GPU consists of massively parallel architecture consisting of thousands of smaller and efficient cores designed to handle multiple tasks at the same time. It is not only

used in a computer on a video card or motherboard, apart it is also used in mobile phones, display adapters, workstations and game consoles. The GPU is also known as Visual Processing Unit (VPU).

**GPU programming Models**

i. CUDA (Compute Unified Device Architecture)
ii. OpenCL (Open Computing Language)
iii. C++ AMP (Accelerated Massive Parallelism)
iv. OpenACC (Open Accelerator)

**E. Support Vector Machine**

It is a supervised machine learning algorithm which can be used for both classification and regression challenges. However, it is mostly used in classification problem. An SVM yields a guide of the sorted data with the margins between the two as far apart as possible. SVMs are utilized in text categorization, image classification, handwriting recognition and in the sciences. It is also known as Support Vector Network (SVN). An SVM creates parallel partitions by generating two parallel lines. Divide two categories by a clear gap that should be as wide as possible, do this partition by a known as hyper plane.

**SVM Algorithm:**

candidateSV = {closest pair from opposite classes }
while there are violating points do
        Find a violator
        candidateSV = candidateSV S
violator
if any $\alpha p < 0$ due to addition of c to S then
        candidateSV = candidateSV \ p
repeat till all such points are pruned
end if
end while

**F. K-Nearest Neighbor**

K-Nearest Neighbor is one of the most basic yet essential classification algorithm in machine learning. It belongs to supervised machine learning domain

and finds intense application in pattern recognition, data mining and intrusion detection.

## KNN Algorithm:

In this algorithm we are assuming n as number of training data set and p as un unknown point.

I.   In this algorithm, store the training data set in an array of data points arr []. This each element of this array represents a tuple (x,y).

II.  for i=0 to m:

III. calculate Euclidean distance d (arr[i],    p)

IV.  Make set S of K smallest distances obtained. Each of these distances corresponds to an already classified data point.

V.   Return the majority label among S.

## III. RELATED WORK

Usually CPU based problem are not suitable for parallel computing due their cost and scalability issues. The GPU is one of the most cost effective solution in the era of big data to solve this problem with massively parallel computing technology[1].

The CPU based algorithm are not fast enough to give a solution in a reasonable amount of time. These problems can become even larger to the point that not even a multi-core CPU based algorithm is fast enough. Today, these problems can be solved faster using massive parallel processing. These problem give a birth to an important hardware for parallel computing known as GPU [2]. The study uncovers the basic and the advanced topics that are required to understand parallel computing and especially GPU computing. The main contribution is in solving computational physics problems that requires massive parallel architecture. In general, data parallel issues are well suited for massively parallel computing in the GPU, that is the problem which divided into several independent sub-problem. This divided problem can be arbitrary distributed or a recursive

divide and conquer strategy. Time dependent problems with small work can lead to inefficient performance and not well suited for GPU computing [3]. In this paper, it is revealed that the performance of two popular GPU integration tools developed in python. The GPU is becoming more popular when the running the machine learning applications including deep learning and computer vision. Cython and PyCUDA are the two most popular python-based framework that supports the high performance computing on the GPU. Cython is the static compiler that wraps C/C++/GPU code in the python. It introduced the Unified Memory feature. PyCUDA allows the user to directly access NVIDIA CUDA driver API and compile the kernel in a just in time fashion and also move data freely between python data objects and GPU memory. As the result Cython achieve comparable performance than PyCUDA.

Though the GPU is becoming the major component of high performance computing, computer architecture has already taken another steps towards more effective utilization of GPUs, in the form of integrated chips [4]. In this chip multicore-CPU and GPU both are integrated in silicon. This both device shares the same physical memory which makes it possible to copy data between device at high speed. The GPU does not have its own memory instead of it, it accesses the system memory with the CPU. Unlike the GPU, which is connected to the host through the PCI bus, the data transfer between memory and both processing units are through the same high speed bus which are controlled by a unified memory controller. The system memory is divided into two parts: the device memory and host memory. These two memories are designed for the GPU and CPU respectively, though both of them can be accessed by either device.

Sentiment Analysis is growing in the field of research with significant application in both industry and academic[5]. The most of the proposed solutions are

centered around supervised, machine learning approaches and review oriented datasets. In this, I am going to focus on the informal communication such as online discussion, tweets and social network comments, unsupervised approaches that estimate the level of emotional intensity contained in the text in order to make predictions. In this experiment were made on the real world datasets, extracted from online social websites, the result will be very robust and reliable solutions for sentiment analysis of informal communication. The sentimental classification is the most useful application of sentiment analysis to classify review using sentiment analysis and classification are technically challenging [6]. The challenging task is to extract the information from the social networking sites and perform sentiment analysis. Using popular social media such as Facebook and Twitter, i have present new perspective to bring out more meaningful information about the social networks.

K-Nearest Neighbors (*k*-NNs) search is a popular and powerful distance based tool for solving classification problem [7]. It is the prerequisite for training local model based classifiers. Fast distance calculation can significantly improve the speed performance of these classifiers and GPUs can be very handy for their accelerations. Meanwhile, several GPU based sorting algorithms are also included to sort the distance matrix and seek for the k nearest neighbors. The speed performances of the sorting algorithms vary depending upon the input sequences. Another research deals with developing parallel processing algorithms for Graphic Processing Unit (GPU) in order to solve machine learning problems for large datasets [8]. In particular, it contributes to the development of fast GPU based algorithms for calculating distance matrix. It also presents the algorithm and implementation of a fast parallel Support Vector Machine (SVM) using the GPU. These application tools are developed using Computing Unified Device Architecture (CUDA), which is a popular software framework for General Purpose Computing using GPU (GPGPU). The GPU version of parallel SVM based on parallel Sequential Minimal Optimization (SMO) implemented in this dissertation is proposed to reduce the time cost in both training and predicting phases. This implementation of GPUSVM is original. It utilizes many parallel processing techniques to accelerate and minimize the computations of kernel evaluation, which are considered as the most time consuming operations in SVM. Although the many core architecture of GPU performs the best in data level parallelism, multi-task processing is also integrated into the application to improve the speed performance of tasks such as multiclass classification and cross-validation.

## IV. METHODOLOGY

The methodology for software development will be as follows:

**A. Acquire Social Media Data-** To gather the social media data through the different websites of the social network sites using the API of the social media. To acquire the twitter data TwitterAPI is used. Tweepy is the python client for the official TwitterAPI.

**B. Preprocess Data**

In preprocessing, remove the unnecessary data like null values, hashtags, repeated letters, and URLs.

    a. **Lower Case:** convert the tweets to lower case.

    b. **URLs:** replace with generic world URLs

    c. **@users:** I can eliminate "@users" via regex matching or replace it with generic core AT_user.

    d. **Hashtag:** Hashtag can give us some useful information, so it is useful to replace them with the exact same words without the hash.

e. **e. Punctuations and Addition whitespace:** remove punctuation at  the start and end of the tweets. It is  also helpful to replace multiple whitespaces with the single whitespace.

## C. Train and Test model using Algorithms-

a.  Support Vector Machine Algorithm
b.  K Nearest Neighbor Algorithm

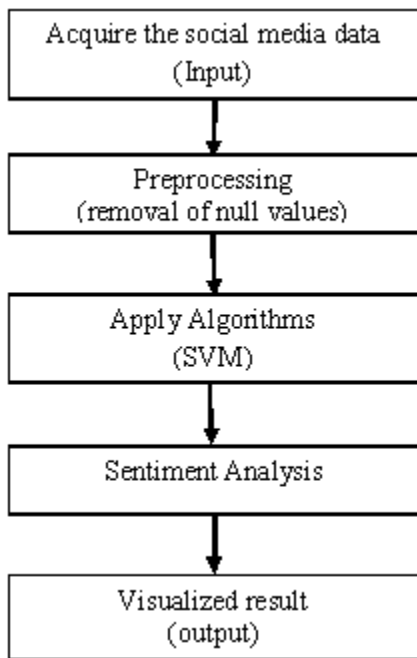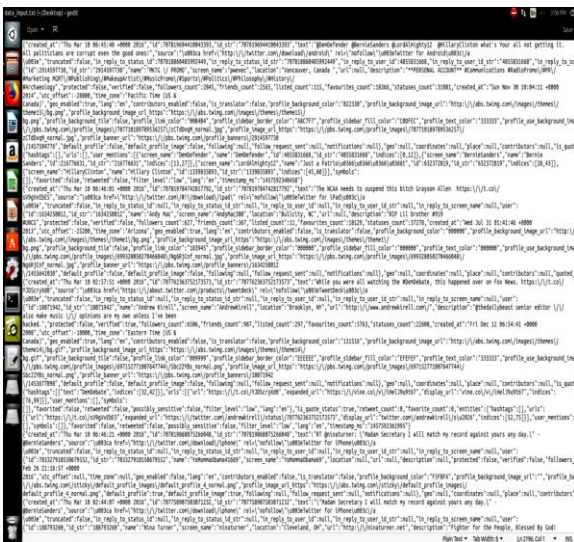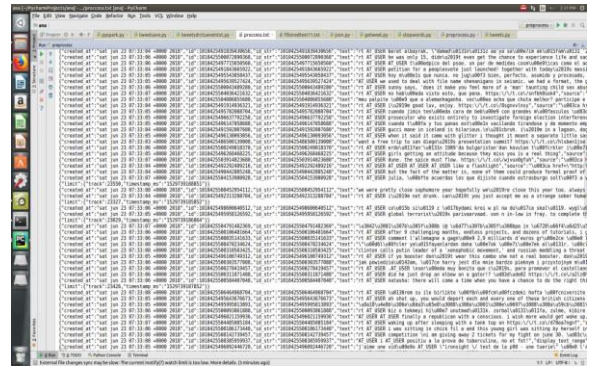**D. Visualized Sentiment Analysis-** Analyze the sentiment of the predicted values.



**Figure 1.** Flow of Methodology

## V.  IMPLEMENTATION

To acquire data, first an application is created to extract data from twitter. To start with the twitter you will need to have a twitter account and obtain credential information that is consumer keys, consumer secret, access token, access secret on the Twitter Developer site to access the Twitter API. If you already have a Twitter account then just go to https//apps.twitter.com/ and login with your twitter user account. This progression gives you a twitter development account under the same name as your user account. Then click to "Create New App" fill the form and click "Create your Twitter Application". Now in the next page, click on "keys and access tokens" and copy your API keys and "API secret", scroll down and click "create my access tokens", and your "access tokens" and "access secret".



**Figure 2.** Create an Application



**Figure 3.** Twitter Application

**A. Data Access:** To gather twitter data using the Twitter API, authentication keys are required to access the Twitter API. The Keys are:

a.  Consumer Key
b.  Consumer Secret Key
c.  Access Token
d.  Access Secret Token

**Figure 4.** Twitter Authentication



**Figure 6.** Preprocessing of Data

### B. Twitter Gathered Data:

Different Packages are used to extract the data:

    a.   Stream Listener: Create a class inheriting from Stream Listener

    b.   Using that class creates a stream object

    c.   Connect the TwitterAPI using the Stream.

### D. Removal of Stops Words:

'a', 'is', 'the', etc. are considered as stop words. These words don't indicate any sentiment and can be removed. The full list of stop words can be found at stop word list.



**Figure 5.** Twitter Data



**Figure 7.** Data without Stopwords

### C. Preprocess Data:

In the preprocessing, remove the un-necessary data like null values, hash tags, repeated letters and URLs.

### E. Applying SVM Algorithm

SVM Algorithm is applied on the dataset using machine learning libraries.

**Figure 8.** Predicted Values

## F. Visualized Output



**Figure 9.** Visualization

**Results:**

42.5% Positive Sentiments

43.3% Negative Sentiments

13.9% Neutral Sentiments

## G. Applying k-NN Algorithm

K-Nearest Search Algorithm is applied on the dataset by importing the machine learning libraries as KNeighborsClassifier. After that, training data is fitted in as a parameter for K-NN. It learns about the data and give the accuracy of the data.



**Figure 10.** Accuracy of KNN

Result:

Accuracy= 85.1666666667

## VI. CONCLUSION

Sentiment Analysis based on social media data is carried out and Positive, negative and neutral sentiment of the Twitter data were found out. Machine learning algorithm is used to classify the data and to analyze the Twitter data. Support Vector Machine is used to analyze the sentiments of the data using GPU. KNN algorithm is used to calculate the accuracy of the data.

## VII. FUTURE WORK

GPU has brought an opportunity of accelerating many applications to various problems. In the future, we can add some extra features like analyzing the data with images or emotions tweets. We can also try to focus on the other sentiments of the data like very positive or very negative.

## VIII. REFERENCES

[1]. Li, Shengren, and Nina Amenta. "Brute-force k-nearest neighbors search on the GPU." International Conference on Similarity Search and Applications. Springer, Cham, 2015.

[2]. Navarro, Cristobal A., Nancy Hitschfeld-Kahler, and Luis Mateu. "A survey on parallel computing and its applications in data-parallel problems using GPU architectures." Communications in Computational Physics 15.2 (2014): 285-329.

[3]. Accelerating Machine Learning Algorithms in Python Patrick Reilly, Leiming Yu and David Kaeli, Department of Electrical and Computer Engineering Northeastern University, Boston, MA

reilly.pa@husky.neu.edu,fylm,kaelig@ece.neu.edu.

[4]. Chen, Linchuan, Xin Huo, and Gagan Agrawal. "Accelerating MapReduce on a coupled CPU GPU architecture." Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis. IEEE Computer Society Press, 2012.

[5]. Li, Peilong, et al. "Heterospark: A heterogeneous cpu/gpu spark platform for machine learning algorithms." Networking, Architecture and Storage (NAS), 2015 IEEE International Conference on. IEEE, 2015.

[6]. Bello-Orgaz, Gema, Jason J. Jung, and David Camacho. "Social big data: Recent achievements and new challenges." Information Fusion 28 (2016): 45-59.

[7]. Huq, Mohammad Rezwanul, Ahmad Ali, and Anika Rahman. "Sentiment analysis on Twitter data using KNN and SVM." Int J Adv Comput Sci Appl 8.6 (2017): 19-25.

[8]. Li, Qi, et al. "GPUSVM: a comprehensive CUDA based support vector machine package." Open Computer Science1.4 (2011): 387-405.

[9]. Thambawita, D. R. V. L. B., Roshan Ragel, and Dhammika Elkaduwe. "To use or not to use: Graphics processing units (GPUs) for pattern matching algorithms." Information and Automation for Sustainability (ICIAfS), 2014 7th International Conference on. IEEE, 2014.