

# GPU Similarity vs Performance

Aren Desai

2024-03-20

## Introduction

Graphics Processing Units (GPUs) are a fundamental part of computer hardware. They are typically tasked with parallel computing for quick matrix calculations, which allows for complex graphics to be rendered quickly. Within the past few decades, it has also been discovered that the parallel computing capabilities of a GPU can also enhance the training and usage of a neural network, which similarly involves a high quantity of matrix calculations. In this report, datasets from alanjo and michaelbryantds are merged by the GPU model and analyzed to determine a few key characteristics revolving around the understanding of what GPUs with a similar architecture have in common in performance. This can be extended with unsupervised learning to cluster GPUs into three categories: low, middling, and high quality.

## Dataset

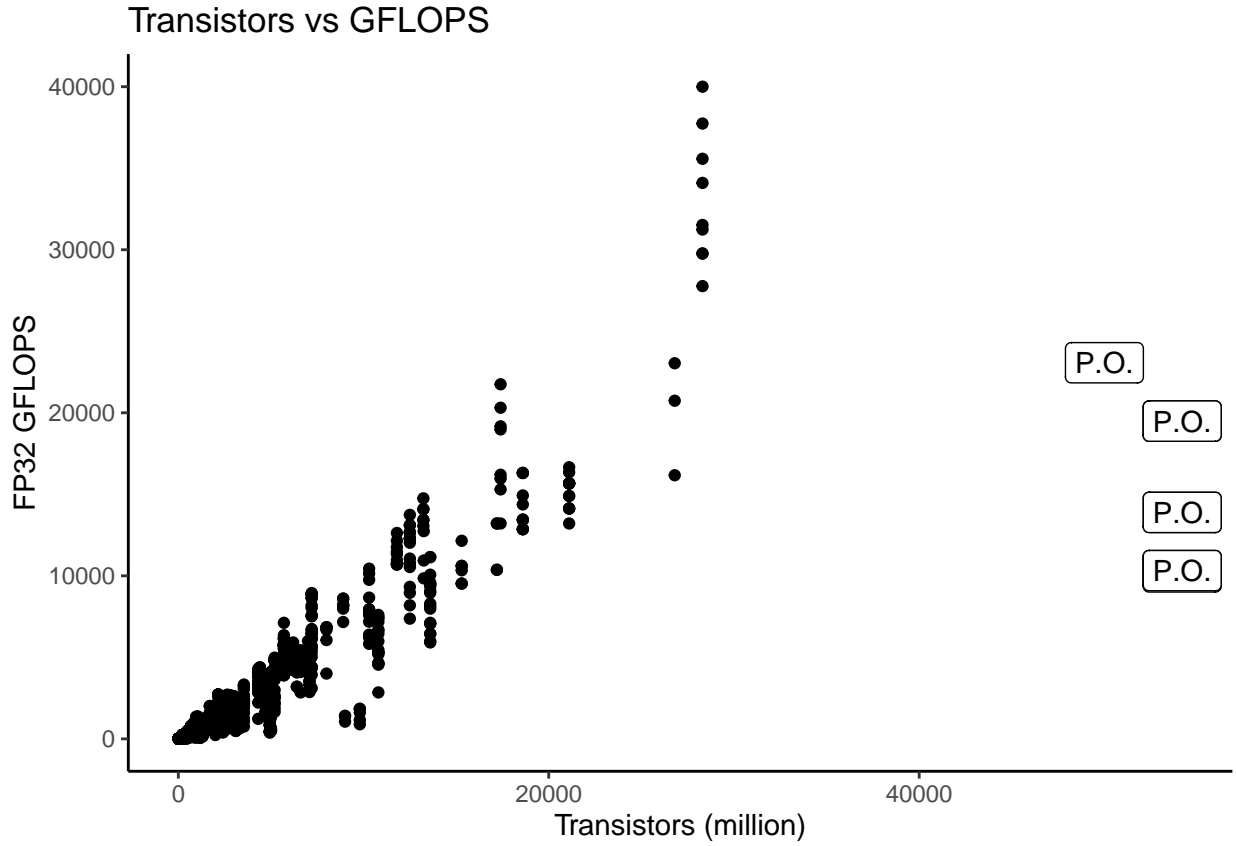
The dataset compiled here is indexed by the names of the models of various GPUs, and includes the memory size, the memory bus width, the gpu clock speed, the memory clock speed, unified shader count, TMU, ROP, the processor size, the die size, the transistor count, the frequency, the g-flops in 16/32/64 bytes, and the manufacturer of the GPU model. All of this is content in the csv file *gpu\_data.csv*, where the manufacturer column is removed and the index is replaced with the model name. Basic statistics of the dataset can be seen in Table 1.

**Table 1: Summary Statistics**

<i>Feature</i>	<i>Mean</i>	<i>Median</i>
<i>memSize</i>	2.618	1.024
<i>memBusWidth</i>	258	128
<i>gpuClock</i>	675	613.5
<i>memClock</i>	860.8	800
<i>unifiedShader</i>	689.3	240
<i>tmu</i>	46.83	24
<i>rop</i>	18.68	8
<i>Process Size(nm)</i>	58.48	40
<i>TDP (W)</i>	82.57	49
<i>Die Size (mm<sup>2</sup>)</i>	205.8	148
<i>Transistors (million)</i>	2541	716
<i>Freq (MHz)</i>	674.3	620
<i>FP16 GFLOPS</i>	2638.24	403.2
<i>FP32 GFLOPS</i>	1691.59	384
<i>FP64 GFLOPS</i>	206.76	33

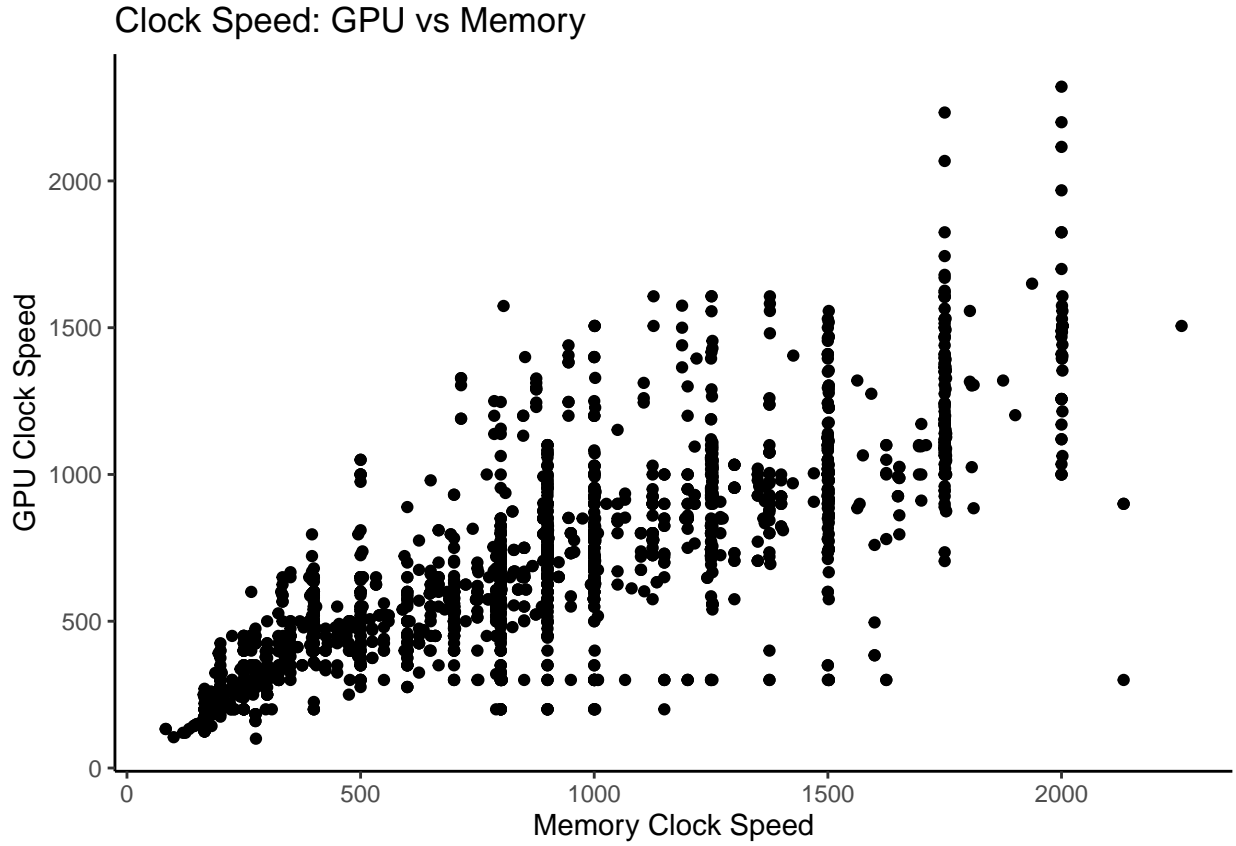
Some meaningful connections can be seen by simple plots. For example, the plot of transistors vs GFLOPS shows a somewhat linear or exponential correlation, with potential outliers labeled as P.O. (potential outlier).

Figure 1: A scatter plot showing the relation between the number of transistors in a GPU and its GFLOPS performance for FP32.



It is clear that a direct upgrade in technology (more transistors, better GPU speed, etc.) indicates a higher quality in other aspects of the GPU. This is to be expected of course, but the linearity of figures 1 and 2 are interesting to note. It appears as a direct correlation with high variance that better GPU architecture leads to better performance.

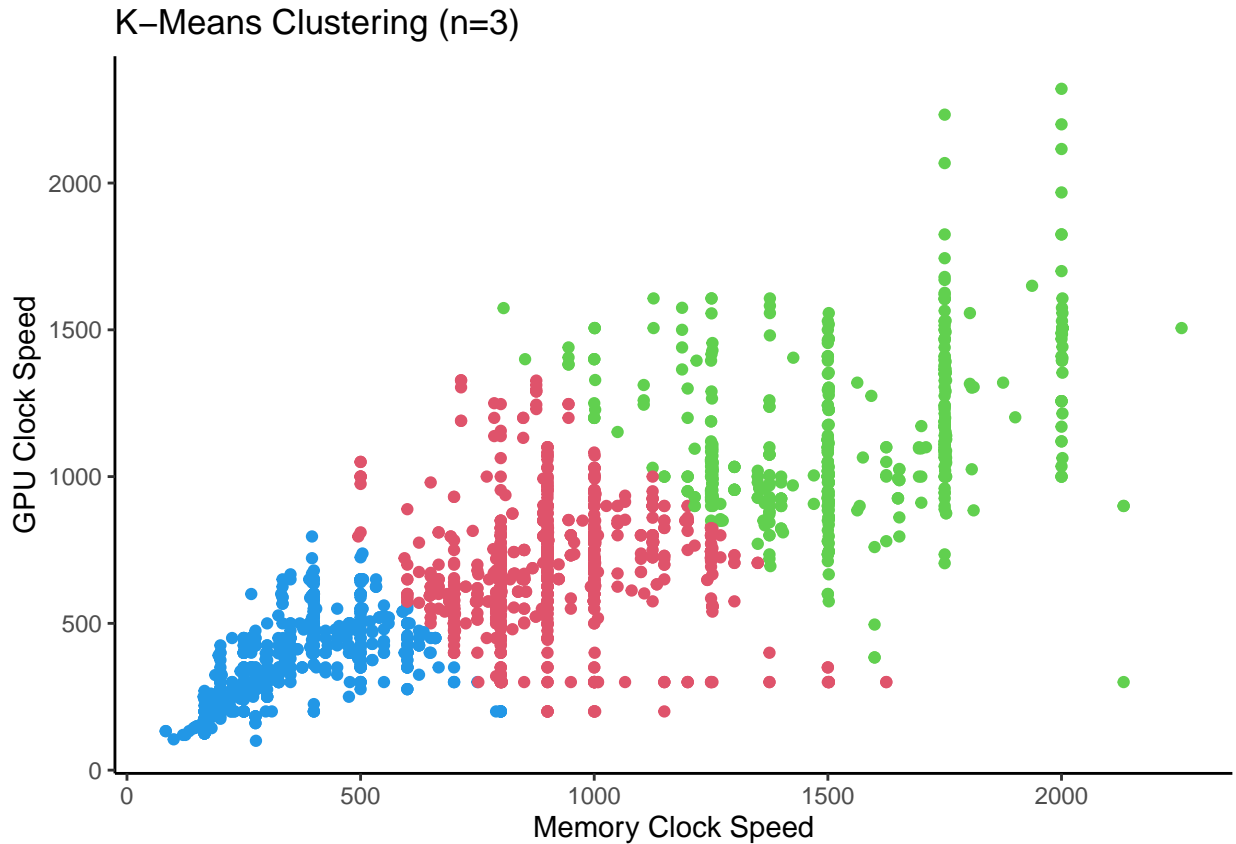
Figure 2: A scatter plot depicting the relation between GPU clock speed and Memory Clock Speed.



### Cluster Analysis

In the pursuit of analyzing similarity in GPUs, a question arose whether GPUs could be grouped into categories based on their attributes. K-Means clustering provided a satisfactory result as shown in Figure 3, a K-Means clustering analysis done on the memClock and gpuClock attributes. Three clear distinctions can be observed in the GPUs, each corresponding to lower, middling, and higher tier GPU build qualities.

Figure 3: Figure 2 separated into 3 clusters depicting low, middling, and high quality.



Increasing K-Means algorithm to include all of the features in the dataset would allow a decent method for classifying whether a new GPU introduced into the dataset would be low, middling, or high quality. In order to do this, spectral clustering was applied for easy visual analysis of the clustering.

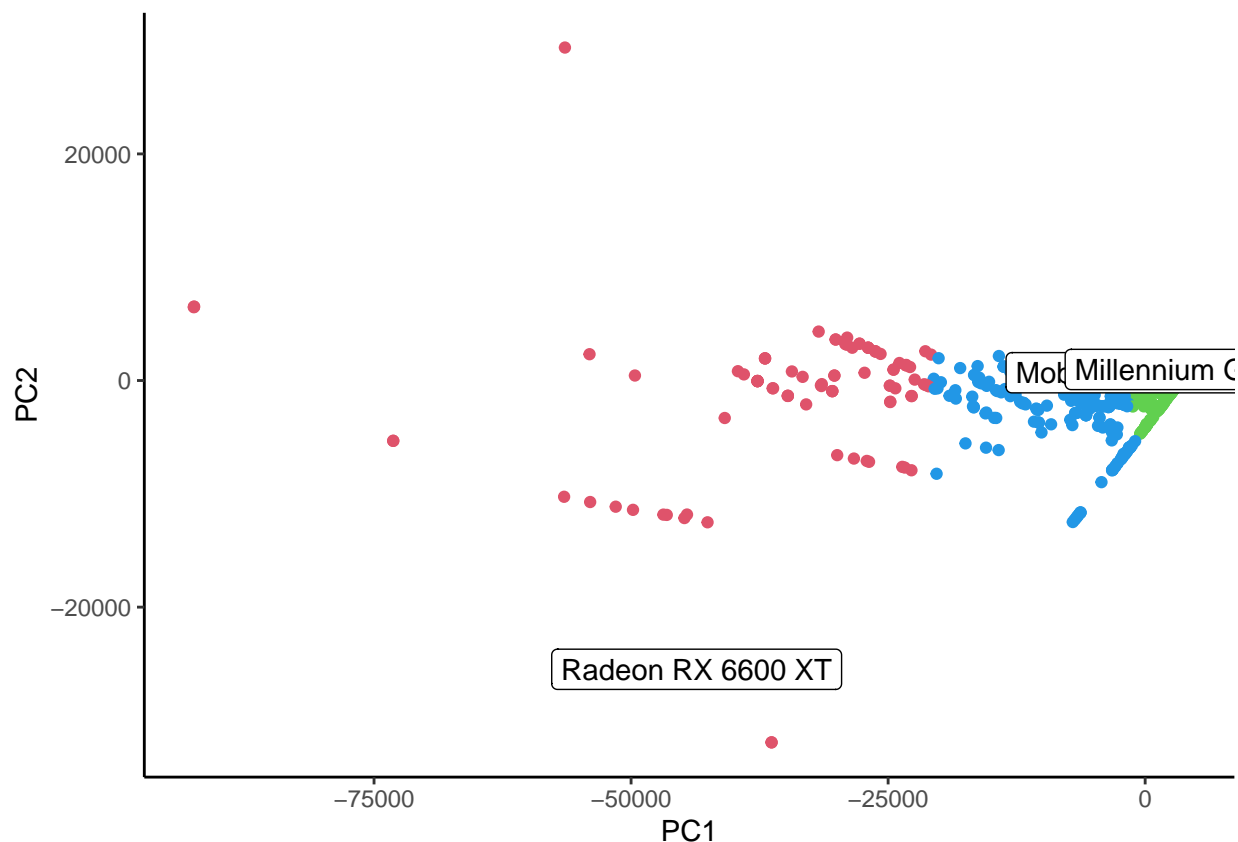
After performing PCA on the GPU dataset, the first two components corresponded to  $\sim 0.78$  cumulative proportion of the variance. While this is slightly on the low side, it is rather good for simplifying from 15 to 2 features. K-Means clustering with 3 centroids provides Figure 4, which depicts the dataset of GPUs clustered by quality according to all available features.

In order to properly identify the groups, these GPUs were isolated by quality and identified on the graph:

- High: Radeon RX 6600 XT
- Middling: Mobility Radeon HD 5145
- Low: Millennium G450

However, the categories as represented on figure 4 were not as clearly defined as expected. Instead, GPUs like the Radeon RX 6600 XT appear largely dissimilar to the majority of other GPUs, while the Mobility Radeon HD 5145 and Millennium G450 appear close together. Most likely, this graph is capturing the similarities in GPU architecture without as much regard to their performance as expected. In this case, GPUs like the Radeon RX 6600 XT are most likely specially developed GPUs that deviate from typical GPU architecture to enhance their niche case.

Figure 4: A scatter plot depicting three GPU clusters with all features, minimized into two dimensions with PCA



## Classification

The GPU clock speed (represented as `gpuClock`) was set as the target for supervised learning, and the other attributes used as training data for our model. The goal was to predict the GPU clock speed based on those other attributes. In the case that a new GPU is being invented, its GPU clock speed could be predicted with this model before its creation.

Formulating parameters for a parametric model seemed both difficult and unreasonable, as the GPUs being created don't necessarily have to follow an underlying distribution (especially with the quick-moving technology sector). Therefore, the non-parametric supervised learning method K-Nearest Neighbors was selected, with  $k=3$ . This resulted in a  $\sim 99\%$  accuracy for all classes, which can be seen in Table 2. This may be considered overfitting in other contexts, but since there are only three classes being detected with 14 features, this can be considered negligible.

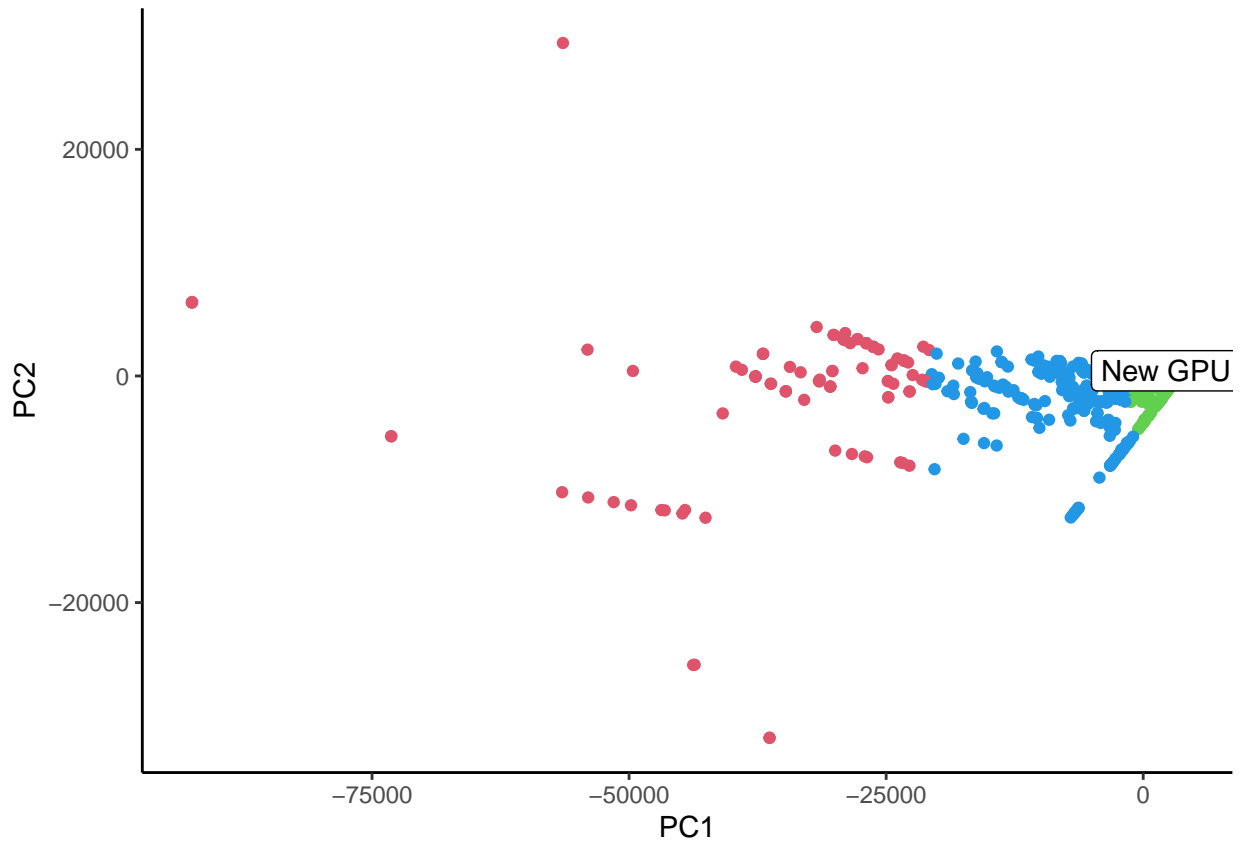
<i>Statistic</i>	<i>Class 1</i>	<i>Class 2</i>	<i>Class 3</i>
<i>Sensitivity</i>	0.9909	1	0.9979
<i>Specificity</i>	0.9980	0.99912	0.9976
<i>Pos Pred Value</i>	0.9880	0.97917	0.9995
<i>Neg Pred Value</i>	0.9985	1	0.9907
<i>Prevalence</i>	0.1404	0.03988	0.8197
<i>Detection Rate</i>	0.1392	0.03988	0.8180
<i>Detection Prevalence</i>	0.1409	0.04073	0.9978
<i>Balanced Accuracy</i>	0.9945	0.99956	0.9978

## New GPU Simulation

To conclude this project, we'll assume a new GPU is being created and simulate its quality and specifications. This new GPU will be based on slightly higher values of the median per feature (with some random variance), and will have the following qualities:

<i>memSize</i>	1.024
<i>memBusWidth</i>	256
<i>gpuClock</i>	800
<i>memClock</i>	700
<i>unifiedShader</i>	345
<i>tmu</i>	38
<i>rop</i>	10
<i>Process Size(nm)</i>	45
<i>TDP (W)</i>	50
<i>Die Size (mm<sup>2</sup>)</i>	200
<i>Transistors (million)</i>	1000
<i>Freq (MHz)</i>	600
<i>FP16 GFLOPS</i>	1200
<i>FP32 GFLOPS</i>	900
<i>FP64 GFLOPS</i>	200

The K-NN categorized this GPU as class 1, and after spectral decomposition, it was placed in the following cluster:



## Conclusion

While the methods presented in this project didn't run into any large errors and, honestly, had better results than expected given the heavy data cleaning done to merge the two initial datasets, it doesn't appear great at classifying GPUs by quality, as the class of the GPU was 1 (which was supposed to indicate low quality) and in the spectral decomposition, the GPU was placed in a high-density area of class 1. However, it appears that the analysis done in this project has revealed a certain "average GPU" architecture, which can identify when a GPU may have been created that deviates from the norm. This is shown in how the GPU prototype achieved a class that places it similar to most other GPUs and in a high-density area, most likely due to its features being around the median.

## Appendix

```
library(dplyr)
library(ggplot2)
library(tibble)
library(MASS)
library(dbSCAN)
library(factoextra)
library(class)
library(caret)

df = read.csv("gpu_data.csv")
dfz = df
df = column_to_rownames(df, "Product")

summary(df)

figure1 = ggplot(df, aes(x = Transistors..million., y = FP32.GFLOPS)) +
  geom_point() +
  theme_classic() +
  xlab('Transistors (million)') +
  ylab('FP32 GFLOPS') +
  ggtitle('Transistors vs GFLOPS') +
  geom_label(data = subset(df, Transistors..million. > 40000), label="P.O.")

figure2 = ggplot(df, aes(y =gpuClock, x = memClock)) +
  geom_point() +
  theme_classic() +
  ylab('GPU Clock Speed') +
  xlab('Memory Clock Speed') +
  ggtitle('Clock Speed: GPU vs Memory')

df1 = data.frame('memClock' = df['memClock'][[1]], 'gpuClock' = df['gpuClock'][[1]])
km.res = kmeans(df1, 3, nstart=25)

figure3 = ggplot(df1, aes(y = gpuClock, x = memClock)) +
  geom_point(col=km.res$cluster+1) +
  theme_classic() +
  ylab("GPU Clock Speed") +
  xlab("Memory Clock Speed") +
```

```

ggtitle("K-Means Clustering (n=3)")

Sigma = cov(df)
gpu_pcacov = princomp(covmat = Sigma, scores=TRUE)
gpu_scores = data.frame(prcomp(df)$x)
# Eliminating outlier
gpu_scores = gpu_scores %>% filter(PC1 > -100000)
km.res2 = kmeans(gpu_scores, 3, nstart=25)

df2 = dfz %>% filter(Product != 'Radeon Instinct MI100') %>% column_to_rownames("Product")
df2$Cluster = km.res2$cluster

figure4 = ggplot(gpu_scores, aes(x = PC1, y = PC2)) +
  geom_point(col=km.res2$cluster+1) +
  theme_classic() +
  geom_label(data = subset(gpu_scores, row.names(gpu_scores) == "Radeon RX 6600 XT"), label="Radeon RX (",
  geom_label(data = subset(gpu_scores, row.names(gpu_scores) == "Mobility Radeon HD 5145"), label="Mobili",
  geom_label(data = subset(gpu_scores, row.names(gpu_scores) == "Millennium G450"), label="Millennium G

predictors = subset(df2, select = -Cluster)
responses = df2$Cluster
knnres = knn(predictors, predictors, responses, k=3)
conf_matrix = confusionMatrix(knnres, factor(df2$Cluster))

df3 = c(
  "memSize" = 1.024,
  "memBusWidth" = 256,
  "gpuClock" = 800,
  "memClock" = 700,
  "unifiedShader" = 345,
  "tmu" = 38,
  "rop" = 10,
  "Process.Size..nm." = 45,
  "TDP..W." = 50,
  "Die.Size..mm.2." = 200,
  "Transistors..million." = 1000,
  "Freq..MH.)" = 600,
  "FP16.GFLOPS" = 1200,
  "FP32.GFLOPS" = 900,
  "FP64.GFLOPS" = 200
)

knn(predictors, df3, responses, k=3)

df4 = rbind(df, df3)
Sigma = cov(df4)
gpu_pcacov = princomp(covmat = Sigma, scores=TRUE)
gpu_scores = data.frame(prcomp(df4)$x)
# Eliminating outlier
gpu_scores = gpu_scores %>% filter(PC1 > -100000)
km.res2 = kmeans(gpu_scores, 3, nstart=25)

figure5 = ggplot(gpu_scores, aes(x = PC1, y = PC2)) +

```



```
geom_point(col=km.res2$cluster+1) +  
theme_classic() +  
geom_label(data = subset(gpu_scores, row.names(gpu_scores) == "2359"), label="New GPU")
```