**Steps and Future Work:**

To identify users that are likely spend money in their game after finishing the tutorial, three main phases are investigated.

First, the data is joined and load, and visualizing is done to understand data. The images are attached in the appendix at the end of this PDF.

Second, Preprocessing and feature selection are done. Duplicates and remaining Nan values are removed. Then, categorical features are changed to numerical. Next, an automatic feature selection is done using random forest algorithm. Then, data is split into train and test. Next, data is normalized using train data. Because data is very unbalanced to determine total spend, downsampling and oversampling using random sampler and SMOTE is done. Applying PCA and drop categorical columns can be done too.

Third, Machine Learning and Deep Learning models are applied. To detect users who likely spend money two approaches are given: classification and regression. Because the game designers' question is including a suggestion for different prices then the total spend can be viewed as a continues for regression problem or it can be categorized into different ranges to deal with classification. Both approaches are done using different classification and regression models. For regression, different methods such as LSTM, XGBoost, Kernel Ridge etc. are applied and, MSE and MAE errors are reported. For classification, different methods such as LSTM, KNN etc. are applied, precision, recall and F-measure are reported.

The results show that imbalance data affect the model very much. For regression, Kernel Ridge results were better than other models. For classification, the minority class identification is a challenge. Precision, recall and F-measure for majority class was very good, however, the minority class criteria were not good generally. This is a very common problem in imbalanced data which needs more investigation with more complex models (for example hybrid models). It is interesting to mention that because the precision, recall and F-measure for the class of users who are not likely to spend money on game are very high, then we can look at the problem reverse (too many False Positive is responsible for low accuracy of negative class). In this scenario, the users who are most likely won't spend will be detected and remain users should be investigated more (a very big chunk of users will be eliminated).

For future work, anomaly detection and hyperparameter tuning using Bayesian optimization or Genetic Algorithm can be done. However, preprocessing is the most important step, and more investigation is needed on preprocessing.
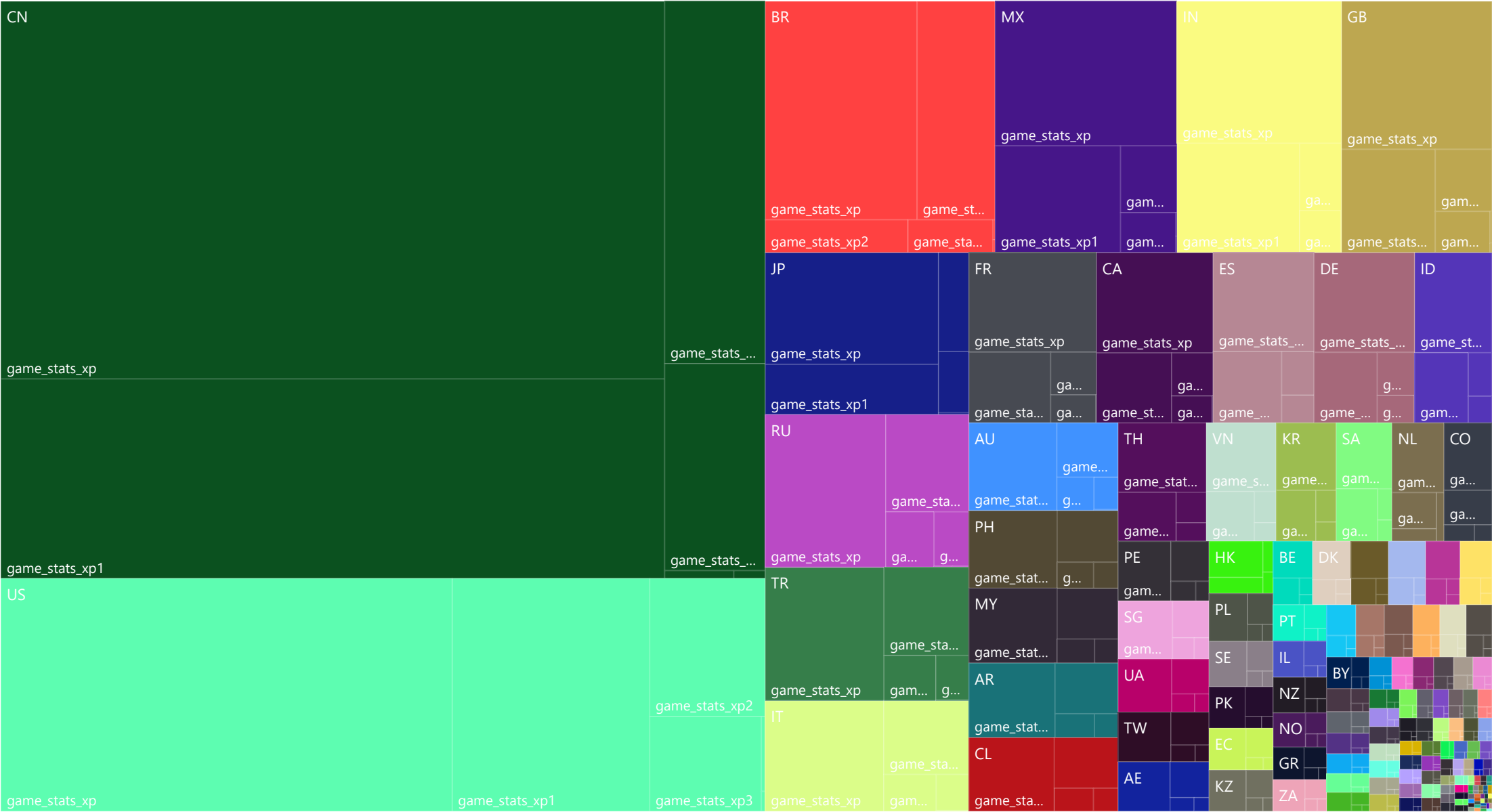
# Appendix:

total_spend, redeemer_actions and scribe_actions by geo_s
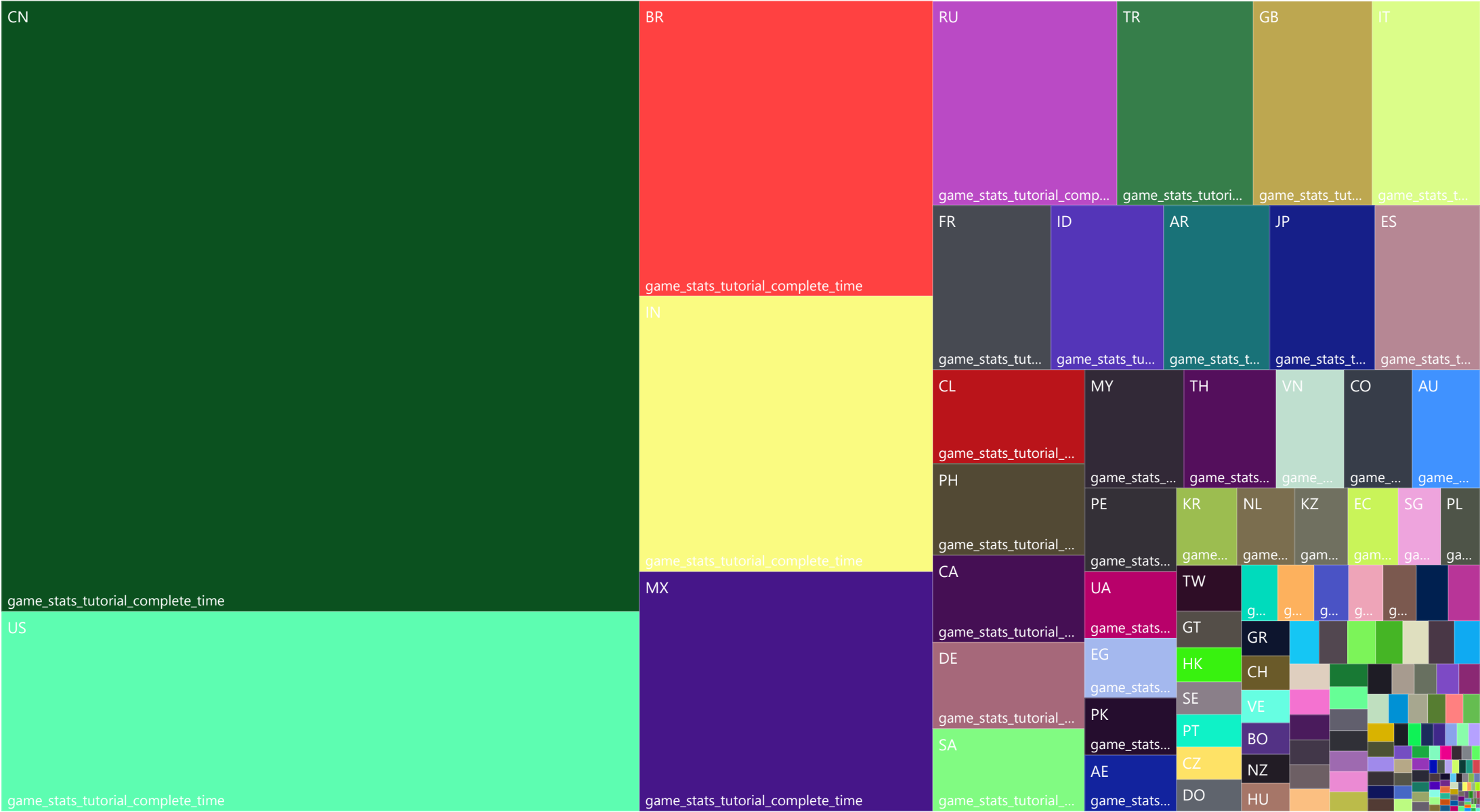


It is not surprising to see most of the people who spend money are from China and US. Brazil, Mexico, India and Great Britain are in the next positions.

# game_stats_tutorial_complete, total_spend, game_stats_xp, game_stats_xp1, game_stats_xp2 and game_stats_xp3 by geo_s
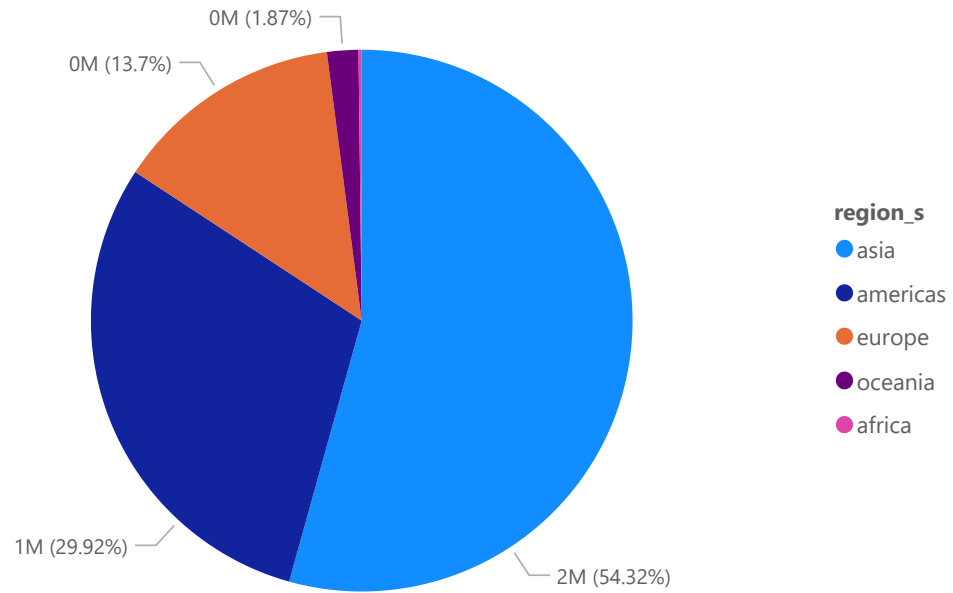


Again the same position based on player who completed tutorial.

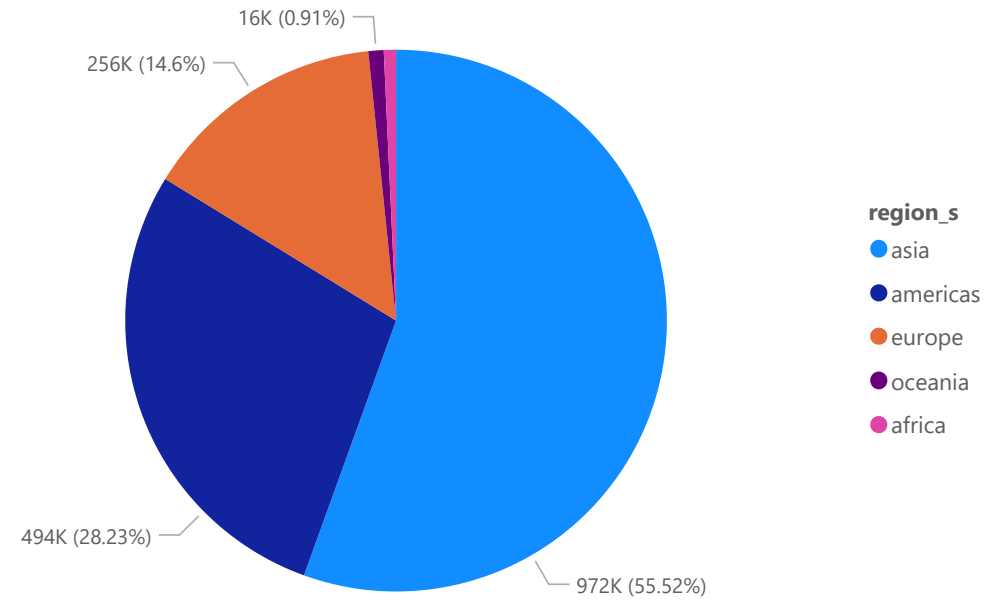**total_spend and game_stats_tutorial_complete_time by geo_s**



For complete time, Brazil, India and Mexico are in the positions 3rd to 5th.
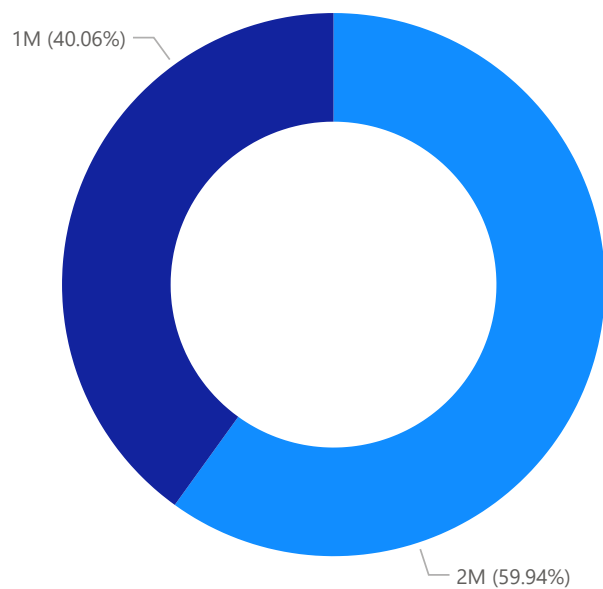
## total_spend by region_s



## game_stats_tutorial_complete by region_s



**total_spend by region_s**

- 0M (1.87%)
- 0M (13.7%)
- 1M (29.92%)
- 2M (54.32%)

region_s
- asia
- americas
- europe
- oceania
- africa

**game_stats_tutorial_complete by region_s**

- 16K (0.91%)
- 256K (14.6%)
- 494K (28.23%)
- 972K (55.52%)

region_s
- asia
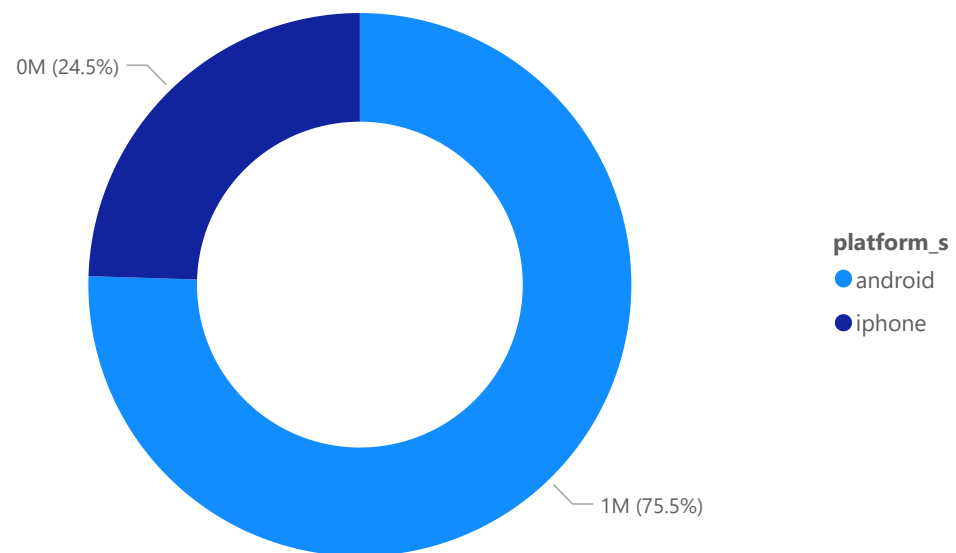- americas
- europe
- oceania
- africa

It is interesting to see total spend and users who completed tutorial are identical by region. 85% from Asia and America.

## total_spend by platform_s
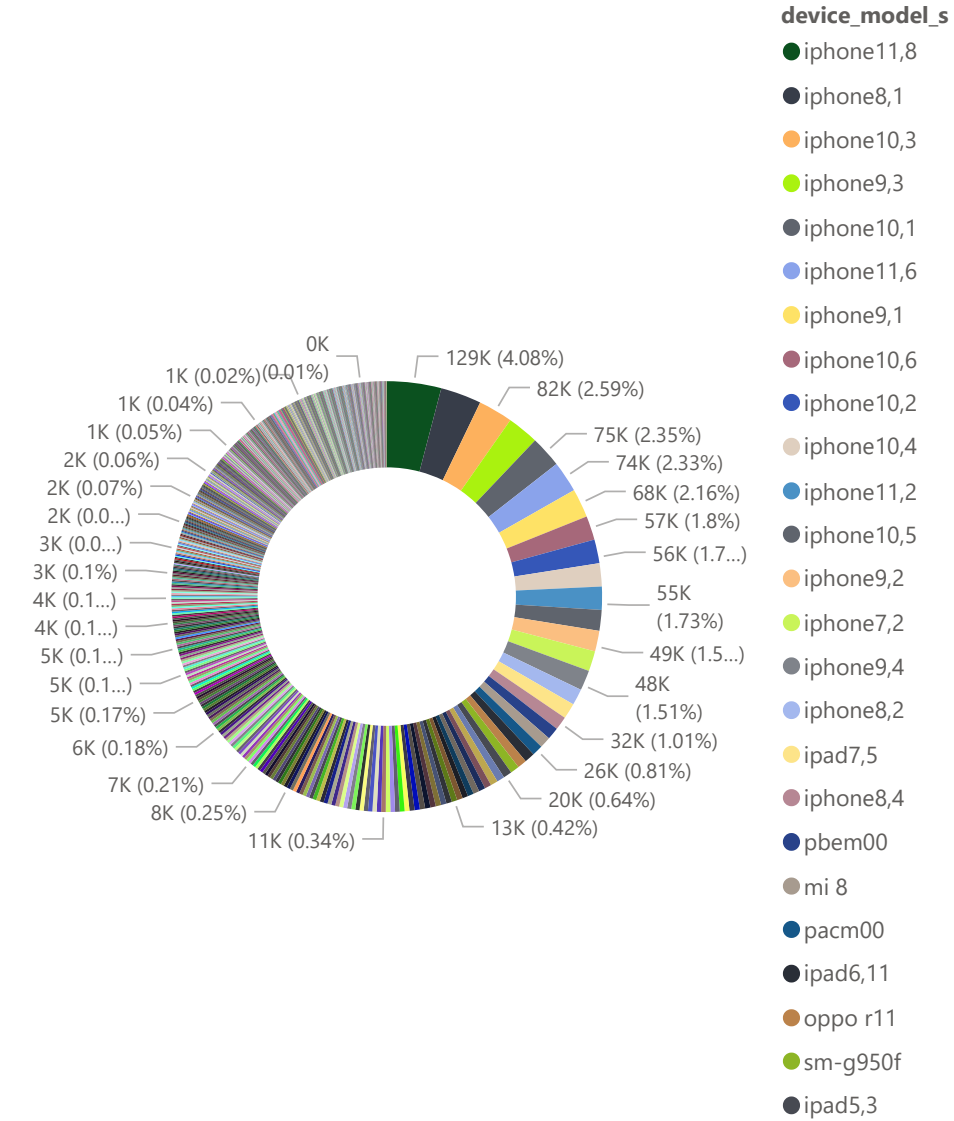


1M (40.06%)

2M (59.94%)

**platform_s**
- android
- iphone

## game_stats_tutorial_complete by platform_s
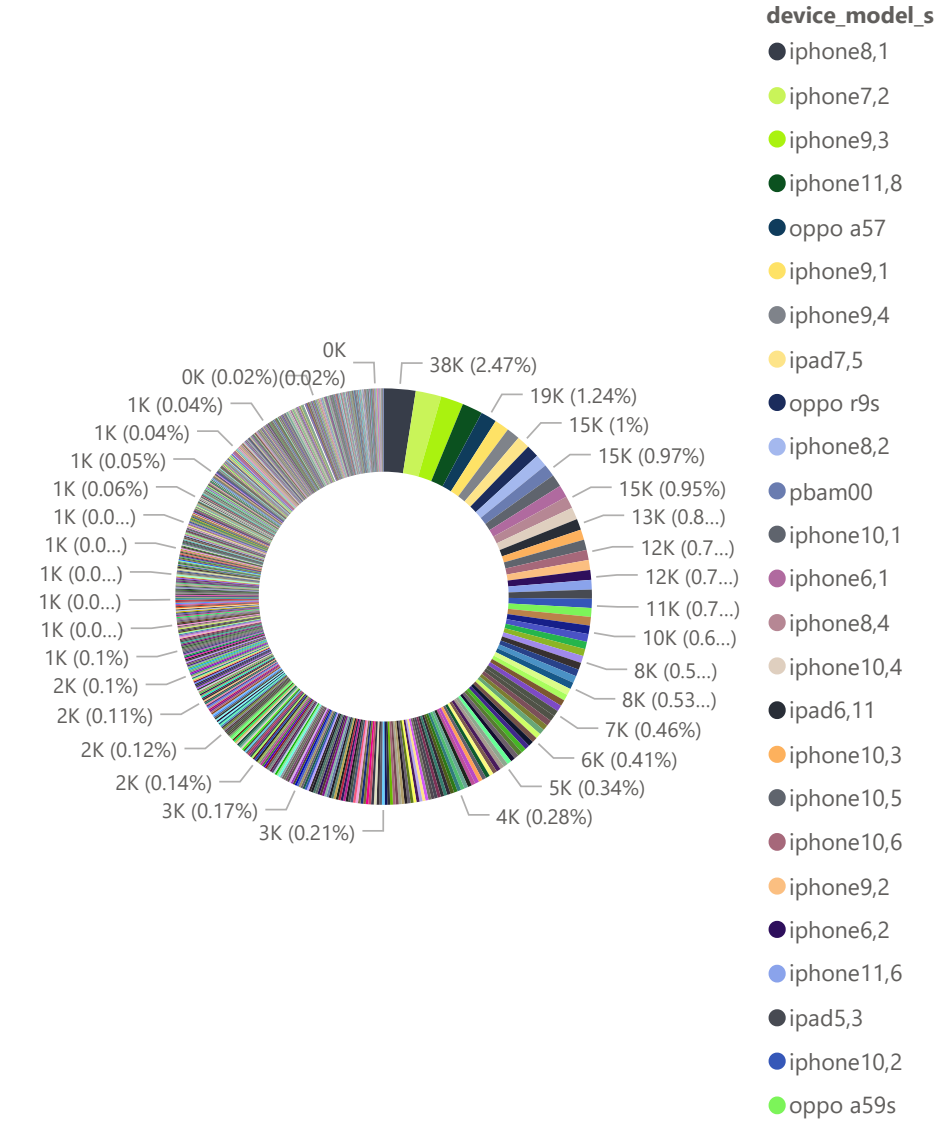


0M (24.5%)

1M (75.5%)

**platform_s**
- android
- iphone

Android is the most popular among people who spend or complete tutorial.

## total_spend by device_model_s



device_model_s
- iphone11,8
- iphone8,1
- iphone10,3
- iphone9,3
- iphone10,1
- iphone11,6
- iphone9,1
- iphone10,6
- iphone10,2
- iphone10,4
- iphone11,2
- iphone10,5
- iphone9,2
- iphone7,2
- iphone9,4
- iphone8,2
- ipad7,5
- iphone8,4
- pbem00
- mi 8
- pacm00
- ipad6,11
- oppo r11
- sm-g950f
- ipad5,3

0K
1K (0.02%)(0.01%)
1K (0.04%)
1K (0.05%)
2K (0.06%)
2K (0.07%)
2K (0.0...)
3K (0.0...)
3K (0.1%)
4K (0.1...)
4K (0.1...)
5K (0.1...)
5K (0.1...)
5K (0.17%)
6K (0.18%)
7K (0.21%)
8K (0.25%)
11K (0.34%)
13K (0.42%)
20K (0.64%)
26K (0.81%)
32K (1.01%)
48K (1.51%)
49K (1.5...)
56K (1.7...)
57K (1.8%)
68K (2.16%)
74K (2.33%)
75K (2.35%)
82K (2.59%)
129K (4.08%)
55K (1.73%)

## game_stats_tutorial_complete by device_model_s



device_model_s
- iphone8,1
- iphone7,2
- iphone9,3
- iphone11,8
- oppo a57
- iphone9,1
- iphone9,4
- ipad7,5
- oppo r9s
- iphone8,2
- pbam00
- iphone10,1
- iphone6,1
- iphone8,4
- iphone10,4
- ipad6,11
- iphone10,3
- iphone10,5
- iphone10,6
- iphone9,2
- iphone6,2
- iphone11,6
- ipad5,3
- iphone10,2
- oppo a59s

0K
0K (0.02%)(0.02%)
1K (0.04%)
1K (0.04%)
1K (0.05%)
1K (0.06%)
1K (0.0...)
1K (0.0...)
1K (0.0...)
1K (0.0...)
1K (0.0...)
1K (0.1%)
2K (0.1%)
2K (0.11%)
2K (0.12%)
2K (0.14%)
3K (0.17%)
3K (0.21%)
4K (0.28%)
5K (0.34%)
6K (0.41%)
7K (0.46%)
8K (0.53...)
8K (0.5...)
10K (0.6...)
11K (0.7...)
12K (0.7...)
12K (0.7...)
13K (0.8...)
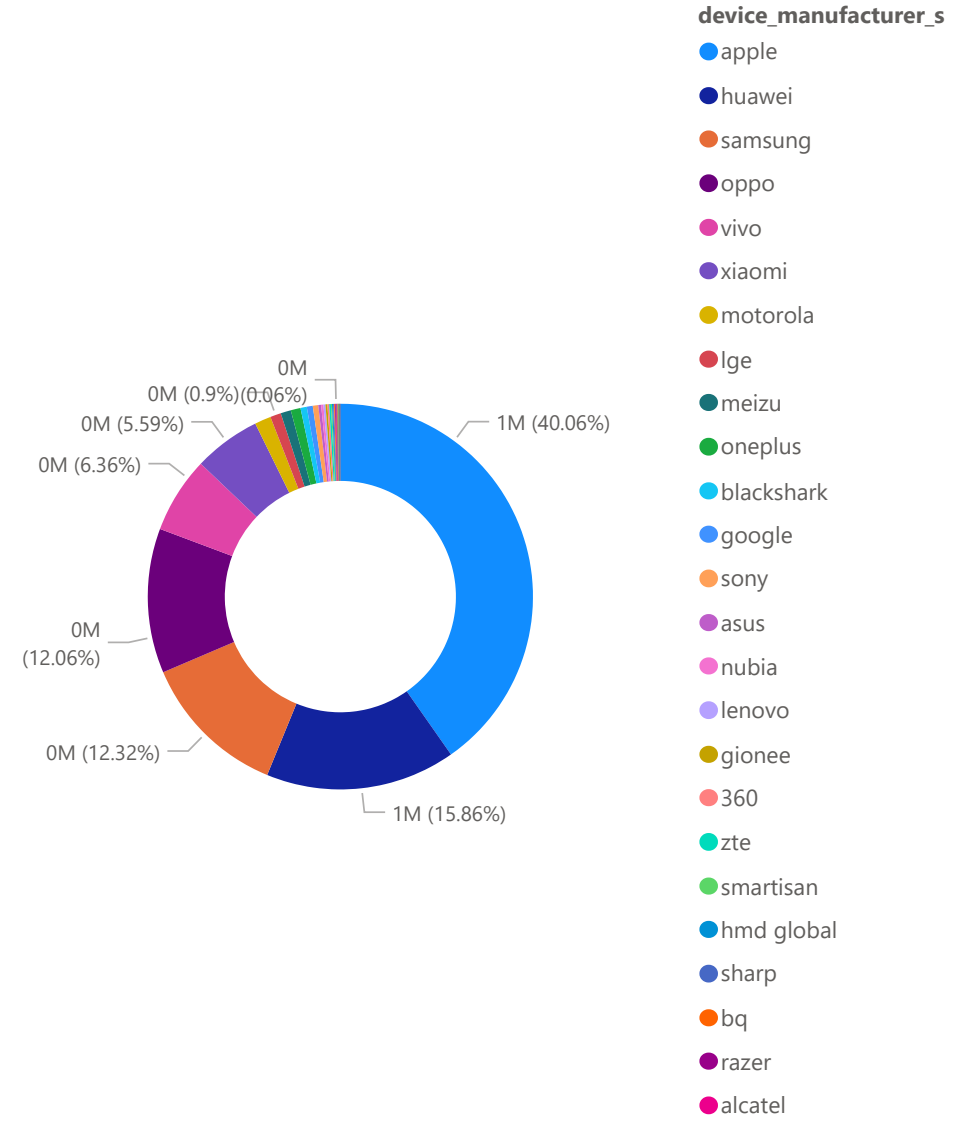15K (0.95%)
15K (0.97%)
15K (1%)
19K (1.24%)
38K (2.47%)

Iphone users are most likely to spend and finish tutorial if we look only at device model?
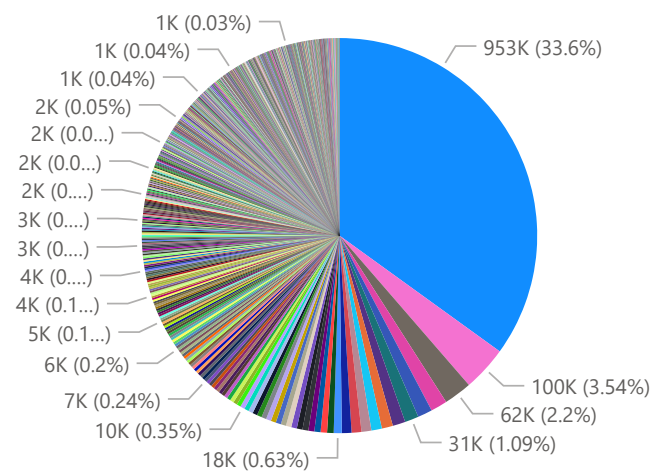
## game_stats_tutorial_complete by device_manufacturer_s



**device_manufacturer_s**
- apple
- samsung
- oppo
- huawei
- xiaomi
- vivo
- motorola
- lge
- meizu
- lenovo
- asus
- oneplus
- sony
- tcl
- zte
- gionee
- hmd global
- google
- htc
- blackshark
- lemobile
- hisense
- wiko
- 360
- general mobile

## total_spend by device_manufacturer_s



**device_manufacturer_s**
- apple
- huawei
- samsung
- oppo
- vivo
- xiaomi
- motorola
- lge
- meizu
- oneplus
- blackshark
- google
- sony
- asus
- nubia
- lenovo
- gionee
- 360
- zte
- smartisan
- hmd global
- sharp
- bq
- razer
- alcatel

Again apple users are mostly to spend. Can we targeted them more? However, most of other phones use only android.

## total_spend by device_os_s



953K (33.6%)
100K (3.54%)
62K (2.2%)
31K (1.09%)
18K (0.63%)
10K (0.35%)
7K (0.24%)
6K (0.2%)
5K (0.1...)
4K (0.1...)
4K (0.0...)
3K (0....)
3K (0....)
2K (0....)
2K (0....)
2K (0.0...)
2K (0.05%)
1K (0.04%)
1K (0.04%)
1K (0.03%)

## game_stats_tutorial_complete by device_os_s



278K (22.57%)
36K (2.95%)
20K (1....)
18K (1.48%)
14K (1.1%)
10K (0.83%)
9K (0.7%)
6K (0.52%)
4K (0.28%)
2K (0.17%)
2K (0.14%)
1K (0.11%)
1K (0.1%)
1K (0.0...)
1K (0....)
1K (0....)
1K (0....)
1K (0....)
1K (0.0...)
1K (0.0...)
1K (0.04%)
0K (0.04%)
0K (0.03%)
0K (0.03%)
0K (0.02%)
0K

ios 12.2 the most popular for users who spend or complete the tutorial.

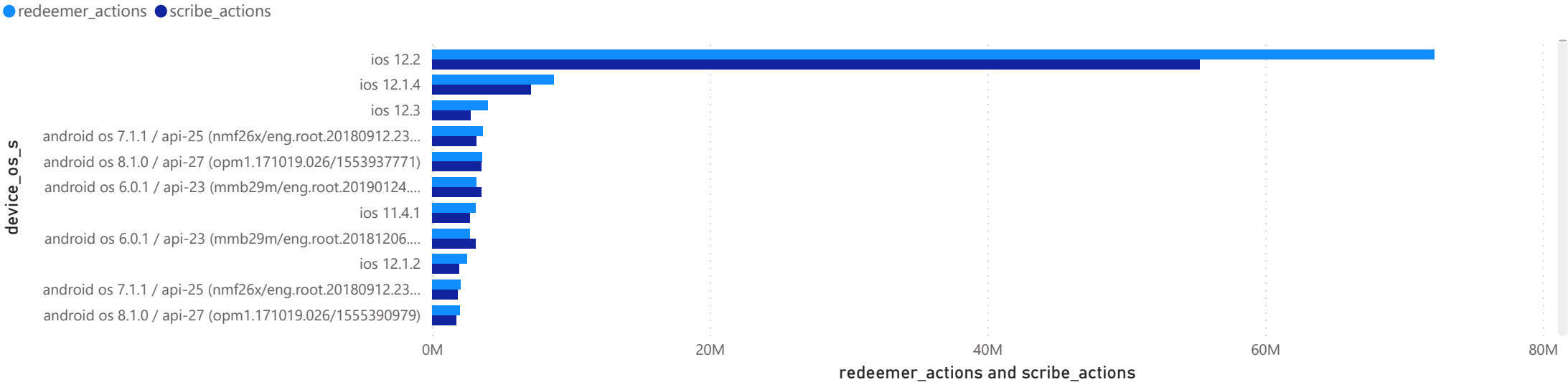## game_stats_xp, game_stats_xp1, game_stats_xp2 and game_stats_xp3 by device_os_s

● game_stats_xp  ● game_stats_xp1  ● game_stats_xp2  ● game_stats_xp3



## game_stats_xp1, game_stats_xp, game_stats_xp2 and game_stats_xp3 by platform_s

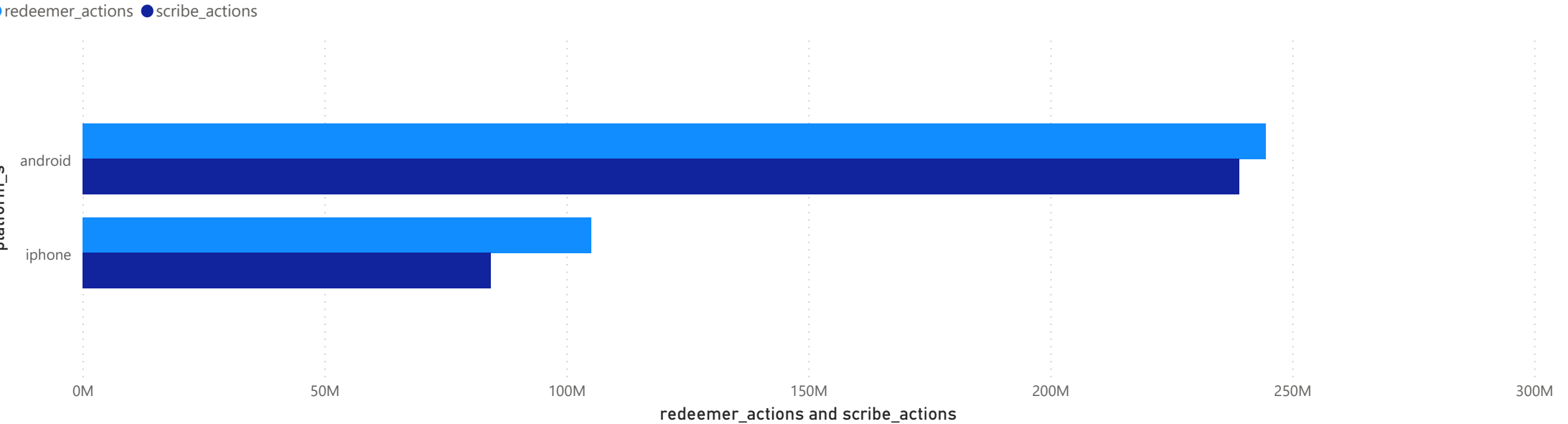● game_stats_xp1  ● game_stats_xp  ● game_stats_xp2  ● game_stats_xp3



For game stats experiences, ios 12.2 again is by far the most used. However, compare to with platform android is far more used. It is interesting to see game_stats_xp has the most points in both graphs.
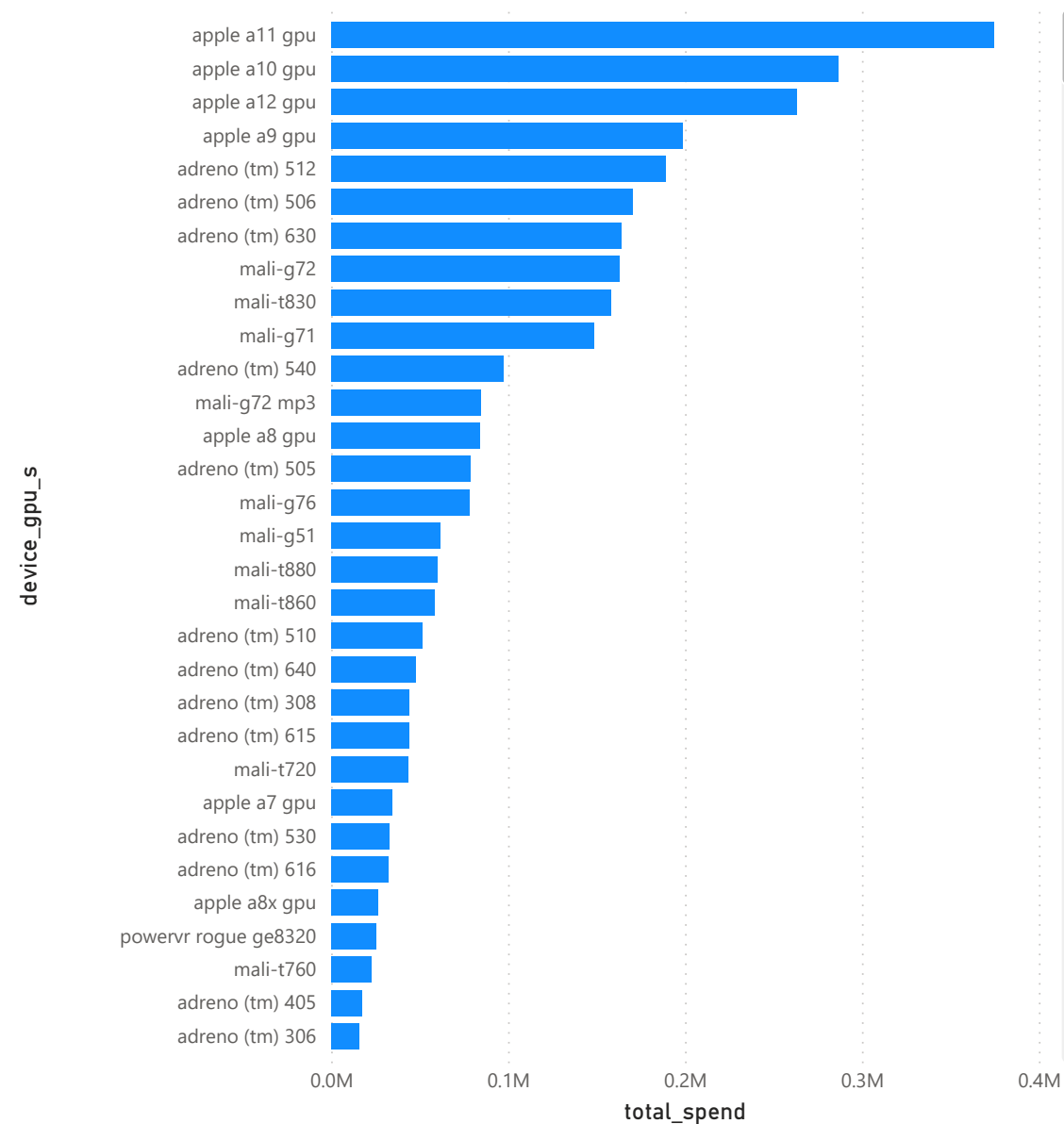
## redeemer_actions and scribe_actions by device_os_s

● redeemer_actions  ● scribe_actions



ios 12.2
ios 12.1.4
ios 12.3
android os 7.1.1 / api-25 (nmf26x/eng.root.20180912.23...
android os 8.1.0 / api-27 (opm1.171019.026/1553937771)
android os 6.0.1 / api-23 (mmb29m/eng.root.20190124....
ios 11.4.1
android os 6.0.1 / api-23 (mmb29m/eng.root.20181206....
ios 12.1.2
android os 7.1.1 / api-25 (nmf26x/eng.root.20180912.23...
android os 8.1.0 / api-27 (opm1.171019.026/1555390979)

device_os_s

redeemer_actions and scribe_actions

## redeemer_actions and scribe_actions by platform_s

● redeemer_actions  ● scribe_actions



android

iphone

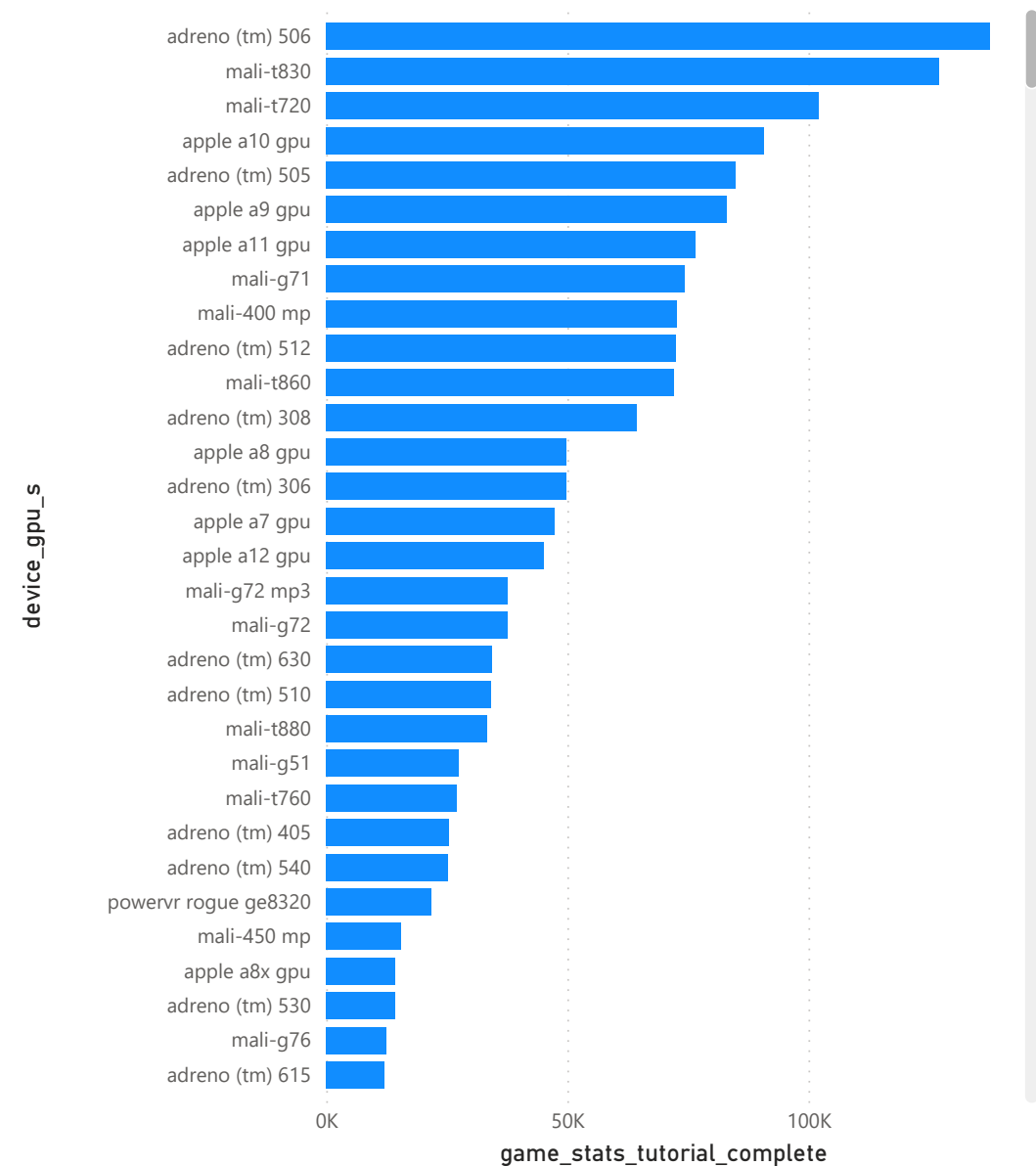platform_s

redeemer_actions and scribe_actions

ios 12.2 far more used for redeemer actions. However, comparing with platform again android is used far more.

## total_spend by device_gpu_s



## game_stats_tutorial_complete by device_gpu_s



The Gpu can be investigated too. Better GPU, most likely to play game?