# Case Study 2: Paris 2024 Team USA Olympics Gymnastics Team

Aimi Wen, Katie Nash, Nidhi Dhupati, Wendy Shi

## Introduction

**Background Information**

The Olympic Games are an international multi-sport event held every four years, and are largely considered to be the most prestigious sporting competition in the world. These Games are therefore a popular spectator event across the world, with an average viewership of more than 15 million people in past years [1]. The 2024 Olympic Games will be held in Paris, France, and will include both a Men's and Women's Artistic Gymnastics event. Winning a medal in an Olympic gymnastics event is an honor for both the individual athlete and for the country that athlete represents; the medal brings a sense of national pride to citizens across the country and promotes the global recognition of the winning country. The USA is known for its men and women gymnastics teams, which have consistently earned medals in past Olympic games. Because of the high honor and importance of upholding the legacy of the United States in these competitions, it is important to carefully select the men and women that will represent our country as Team USA. The importance of and interest surrounding the Olympics provides a rich foundation for analytics and statistical interpretations of the athletes in the competitions, as well as predictions about who will represent Team USA in the many Olympic sports. There have already been analyses, published in major news outlets, about who Team USA will consist of for non-gymnastics sports in 2024 [2]. The objective of this project is to extend this analysis to the Artistic Gymnastics events, and to select the 5 athletes for each of the men's and women's events that will optimize performance for Team USA.

The format of the Olympic Artistic Gymnastics event includes a qualifying round and final round. Scores from qualifying rounds do not impact medal determinations, but rather decide which athletes advance to the finals. The final round, which is what determines medal placements, consists of team all-around, individual all-around, and individual apparatus events. Men's opportunities for medals are: 1) team all-around, 2) individual all-around, 3) floor exercise, 4) pommel horse, 5) still rings, 6) vault, 7) parallel bars, and 8) high bar. Women's opportunities for medals are: 1) team all-around, 2) individual all-around, 3) vault, 4) uneven bars, 5) balance beam, and 6) floor exercise. The athletes representing Team USA will compete in a pool of 96 men and 96 women, which will include teams of five from other countries as well as the option for individual entries for countries without full team qualifications.

Our hypothesis is:

- Who are the 5 men and 5 women who can make Team USA achieve the highest level of success at the Paris 2024 Artistic Gymnastics Event?

**Data Description**

Gymnastics competitions from 2017-2023 were included in the dataset. These competitions were split into two separate datasets by the seasons leading up to the 2021 Tokyo Games and the seasons leading up to the 2024 Paris Games (the 2022 and 2023 seasons). Competitions were included from around the globe, including world competitions, non-world competitions that included the United States as a competitor, and non-world competitions that did not include the United States as a competitor. The competition data included scoring information, athlete names, athlete demographic information, and athlete scores. Each entry in the dataset

is one event from one athlete, which includes the country affiliation of the athlete, the competition and competition location, the difficulty, execution, penalty, and overall scores, the apparatus that the athlete competed on, and their ranking. Additionally, the countries who qualified as a team, which are the countries that will be allowed to compete in the Team All-Around event for the men's and women's Olympic events, were provided. The dataset we used for analysis was last synced with the central repository on November 2nd.

# Methodology

## Data Cleaning

From the original competition dataset, we standardized variables and added columns. We standardized each athlete's name based on the first instance it appeared and created a column that included the full name of the athlete. We additionally made sure the country codes were consistent with International Olympic Committee Codes (for example, Scotland competes as part of the United Kingdom in the Olympics). We additionally separated the competition data into two datasets by men's and women's competitions, created a variable to track the order that the competitions occurred by their dates, and added a competition scope variable to aid with data visualizations and athlete selection in the analysis approach.

## Analysis Approach

In order to study which athletes could enable Team USA to achieve maximum success at the Paris Olympics, we decided to model possible Olympic scenarios by simulating the medal outcomes of different variations of Team USA competing with variations of teams from other countries. Our metric of success was a weighted medal count: we assigned 3 points to a Gold medal, 2 points to a Silver medal, and 1 point to a Bronze medal. This reflects our belief in the relative prestige of each medal, and will allow the highest score to reflect our interpretation of the highest level of success. We conducted a separate simulation process for the Men's Artistic Gymnastics Event and the Women's Artistic Gymnastics Event. Each simulation can be broken down into the following subsections: athlete selection, creating score distributions, team selection, Team All-Around event simulation, individual event simulation, and finding the top performing USA teams.

### Athlete Selection Criteria

We only included athletes who placed in the top 10 at a World competition or in the top 3 at a non-World competition in our hypothetical team combinations. This reflects an assumption that World competitions are more competitive than non-World competitions, and that each country would only select the highest level of top-performing athletes to represent their country in the Olympics.

### Creating Score Distributions

The next step is to simulate the score that each athlete receives on each apparatus. We created a score distribution for each athlete and apparatus that mimics the characteristics of the observed scores they received in past competitions, while adding noise and allowing for more score selection options than just the exact scores they have received in the past. We used the process of kernel density estimation (KDE), which provides a continuous representation of the score distribution. Mathematical details about the development of our KDEs is included in the appendix. Although we are only simulating scores for the final rounds, since those are the medal-earning events, we decided to include athlete scores from qualification and final rounds in our creation of the KDEs. Figure X confirms that the distribution of scores for qualification and final rounds are similar. Additionally, we only created the KDEs based on total score for the athlete, instead of separating the scores by difficulty score, execution score, and penalty score. Since the total score is difficulty score + execution score - penalty score, it encompasses all of these sub-scores. Figure X additionally demonstrates that most penalty scores are near 0.

**Team Selection**

In order to simulate Olympics outcomes, we had to select a team of five gymnasts for each qualifying country in the Olympics, for each simulation round. Because the Olympics has already released which countries have qualified, we used those countries to filter out athletes who did not belong to a qualifying country. For each qualifying country, we used weighted random sampling to select 5 gymnasts. Weighted random sampling was chosen rather than simple random sampling in order for the highest performing athletes to be sampled more often. The weights also biased selection toward "well-rounded" gymnasts, who have a high score on multiple apparatuses, rather than "specialized" gymnasts. We made the assumption that all-around athletes have the potential to compete in multiple medal events, and that medaling the individual all-around event can bring more media coverage and prestige to the winning country than winning an individual apparatus event. Mathematical details about the weighting process are included in the appendix.

**Team All-Around Simulation**

To simulate the team all-around competition, three gymnasts would need to be selected to compete for each apparatus from the five gymnasts on each team. After selecting the five athletes for each country's team, we simulated how each athlete would score for each individual event by randomly sampling 1 score from the appropriate athlete and apparatus KDE score distribution. Then, for each apparatus, we selected the athletes with the 3 highest scores to compete in the Team All-Around.

Once every qualifying country had three athletes selected for each apparatus, the scores for the actual event with these athletes was simulated. Then, the scores for all apparatuses were summed by country to create the final scores for the Team All-Around event. The top 3 teams were given the appropriate medal counts.

**Individual Events Simulation**

Unlike the Team All-Around, individual events allow for athletes from countries that did not qualify as a team. Therefore, when simulating individual events, our dataset included five athletes that were randomly selected (via the weighted random selection process) for each qualifying country, as well as all of the athletes who passed our athlete selection criteria but did not belong to a qualifying country. For each individual event and athlete, we randomly sampled one score from the appropriate athlete and apparatus KDE score density. For individual events where multiple scores can be recorded (the individual all-around and vault events), we summed each score together to represent the final score for each event. The top 3 athletes for each event were given the appropriate medal counts.

**Finding the Top USA Teams**

We ran 5000 simulations for the Men's Artistic Gymnastics Event, and a separate 5000 simulations for the Women's Artistic Gymnastics Event. Each simulation included a separate selection process for the five athletes who would represent the qualifying countries. After running these simulations, we selected the one simulation where the USA scored the highest weighted medal count for each gender. We selected the athletes selected in this highest-scoring round as Team USA. Then, to address potential variability in our results as well as address the limitation that Team USA may have just scored highly due to poor performance of other countries, we set these athletes as our Team USA. Then, we ran another 100 simulations, where we randomly selected the athletes in other qualifying countries but kept the same athletes in every simulation for the United States. From these simulations, we tracked how many times the United States had the highest weighted medal count of all countries.

# Results

Of the 5000 simulations for both men and women, the simulation with the highest weighted medal score for the United States was taken as Team USA. The results of the simulated Olympics competitions, for both the Men's and Women's Artistic Gymnastics events, are included below.

Table 1: Men's Artistic Gymnastics - Best Simulated Outcome for Team USA

| Event | Gold | Silver | Bronze |
|---|---|---|---|
| Team All-Around | USA | China | Switzerland |
| Individual All-Around | USA (Brody Malone) | China (Boheng Zhang) | China (Wei Sun) |
| Vault | Armenia (Artur Davtyan) | Philippines (Carlos Yulo) | Great Britain (Luke Whitehouse) |
| Still Rings | USA (Javier Alfonso) | USA (Brody Malone) | Italy (Salvatore Maresca) |
| Pommel Horse | Chinese Taipei (Chih Lee) | Kazakhstan (Nariman Kurbanov) | Uzbekistan (Abdulla Azimov) |
| Parallel Bars | China (Jingyuan Zou) | USA (Curran Phillips) | Switzerland (Noe Seifert) |
| High Bar | USA (Frederick Richard) | Kazakhstan (Ilyas Azizov) | USA (Brody Malone) |
| Floor Exercise | USA (Brody Malone) | Israel (Artem Dolgopyat) | India (Satyajit Mondal) |

**Team USA Athletes: Brody Malone, Javier Alfonso, Curran Phillips, Frederick Richard, Joshua Karnes**
**Team USA Medal Count:** 8
**Team USA Weighted Medal Count:** 20 (where Gold = 3, Silver = 2, Bronze = 1)
**Team USA Medal Breakdown:** 5 Gold, 2 Silver, 1 Bronze

Table 2: Women's Artistic Gymnastics - Best Simulated Outcome for Team USA

| Event | Gold | Silver | Bronze |
|---|---|---|---|
| Team All-Around | USA | Brazil | China |
| Individual All-Around | USA (Simone Biles) | Brazil (Rebeca Andrade) | USA (Konnor McClain) |
| Balance Beam | USA (Konnor McClain) | USA (Simone Biles) | China (Qingying Zhang) |
| Vault | USA (Jade Carey) | USA (Simone Biles) | South Korea (Seojeong Yeo) |
| Floor Exercise | USA (Simone Biles) | USA (Jade Carey) | Brazil (Rebeca Andrade) |
| Uneven Bars | Algeria (Kaylia Nemour) | Russia (Viktoria Listunova) | USA (Shilese Jones) |

**Team USA Athletes: Simone Biles, Shilese Jones, Karis German, Konnor McClain, Jade Carey**
**Team USA Medal Count:** 10
**Team USA Weighted Medal Count:** 23 (where Gold = 3, Silver = 2, Bronze = 1)
**Team USA Medal Breakdown:** 5 Gold, 3 Silver, 2 Bronze

These tables show that, for both men's and women's events, the United States earns more medals than any other country.

Simone Biles, Shilese Jones, Karis German, Konnor McClain, and Jade Carey were selected as the women's Team USA. All five of these athletes are proven to be highly accomplished [3]:
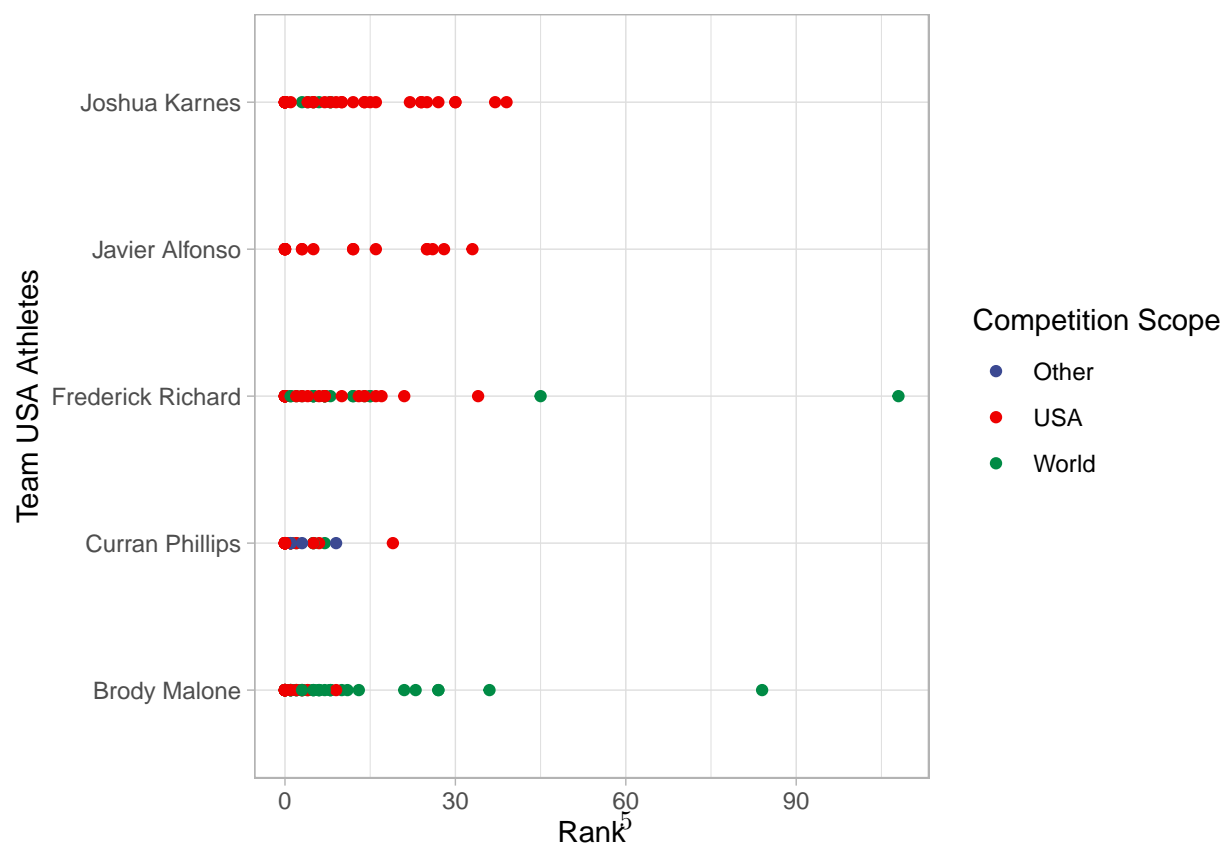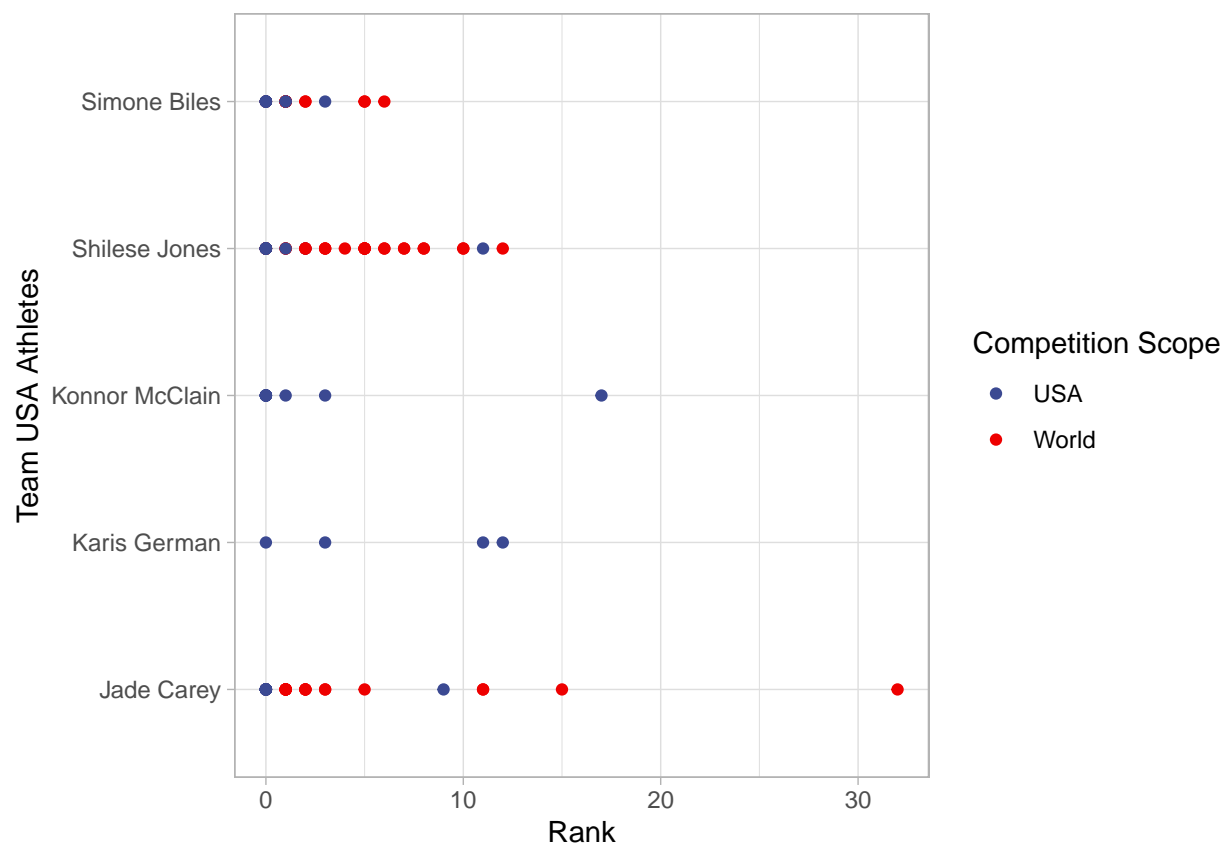
- Simone Biles is the most decorated gymnast in history
- Shilese Jones was part of the gold-winning team at the 2022 World Championships
- Karis German medaled in a floor exercise event at a qualifying event for the 2024 Olympics
- Konnor McClain is the 2022 U.S. National Champion
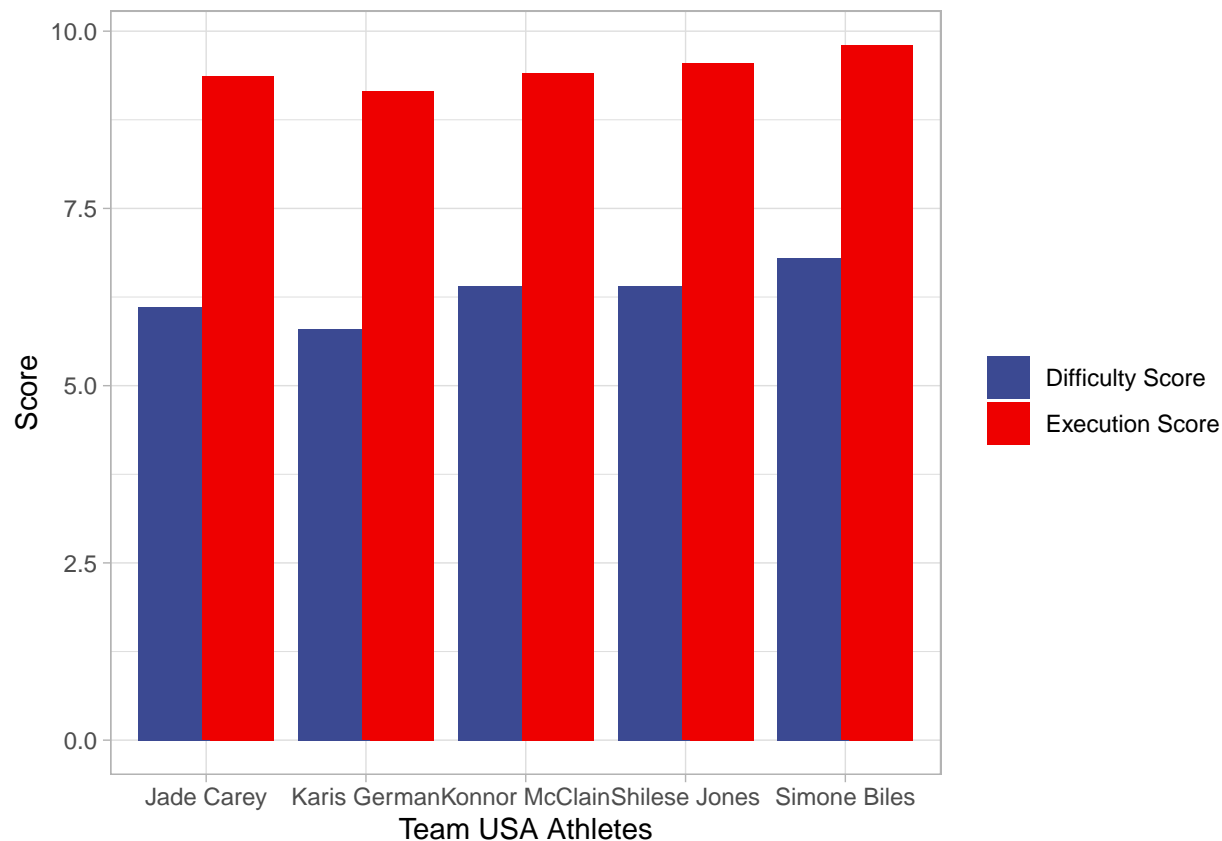- Jade Carey represented Team USA in the Tokyo Olympic Games

Brody Malone, Javier Alfonso, Curran Phillips, Frederick Richard, and Joshua Karnes were selected as the men's Team USA in the simulation where the USA has the highest weighted medal score. All five of these athletes are also proven to be highly accomplished [3]:

- Brody Malone is the 2021 and 2022 U.S. National All-Around Champion
- Javier Alfonso was part of the silver-winning team at the 2023 NCAA competition
- Curran Phillips is the 2023 Pan American Games team champion
- Frederick Richard is the 2023 World all-around bronze medalist
- Joshua Karnes is the 2023 Winter Cup bronze medalist in parallel bars

# Discussion

## Analysis

**Citations**

[1] https://www.cnbc.com/2021/08/04/simone-biles-return-helped-olympics-viewership-average-16point8-million.html

[2] https://www.cbssports.com/nba/news/2024-paris-olympics-projecting-team-usas-12-man-roster-with-joel-embiid-joining-lebron-james-stephen-curry/

[3] https://usagym.org/

**EDA - separated by type of visualization, not gender, figure out where to put (probably appendix, referenced in methodology)**

showing that most penalties are near 0: For both men's and women's competitions, the 95th percentile for the penalties is less than 0.5 (it is 0.5 for men and 0.4 for women). This means that 95% of penalties for both genders are not more than 0.5. The visualizations show that this does not change drastically based on the scope of the competition. Additionally, since the 51st percentile of penalties for both genders is 0.1, this means that the majority of entries had a penalty score of 0.1 or less.

Figure 1: Difference between Difficulty and Exectution Scores for team USA Men

Figure 2: Distribution of the Difficulty Score by Round and Apparatus for Women

Figure 3: Distribution of Difficulty scores by the Scope of Competition and Apparatus for Women

Figure 4: Distribution of Scores by Scope of Competition and Apparatus for Women

## Penalty Distribution for Women Athletes



Why we're using mean for bootstrapping: The distribution of event scores for women and men are both approximately normally distributed, with the mean and median scores less than 0.5 apart within each gender. The mean of women's scores is 12.207, and the median of women's scores is 12.367. The mean of men's scores is 13.048, and the median of men's scores is 13.233.

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

Whether scores vary by type of round in competition: Doesn't seem to vary much; make the assumption that qualification scores are representative of an athlete's scores in the medal rounds

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

Figure 5: Distribution of Penalties by competition scope for Men

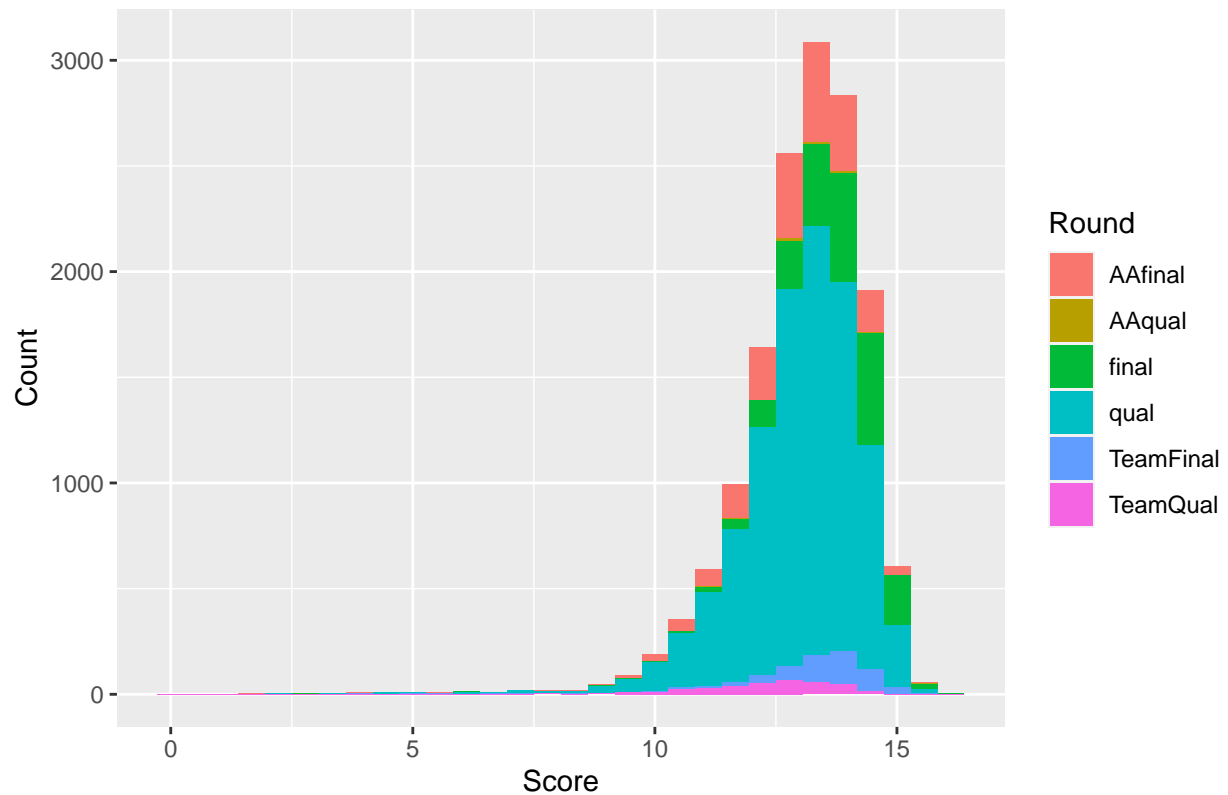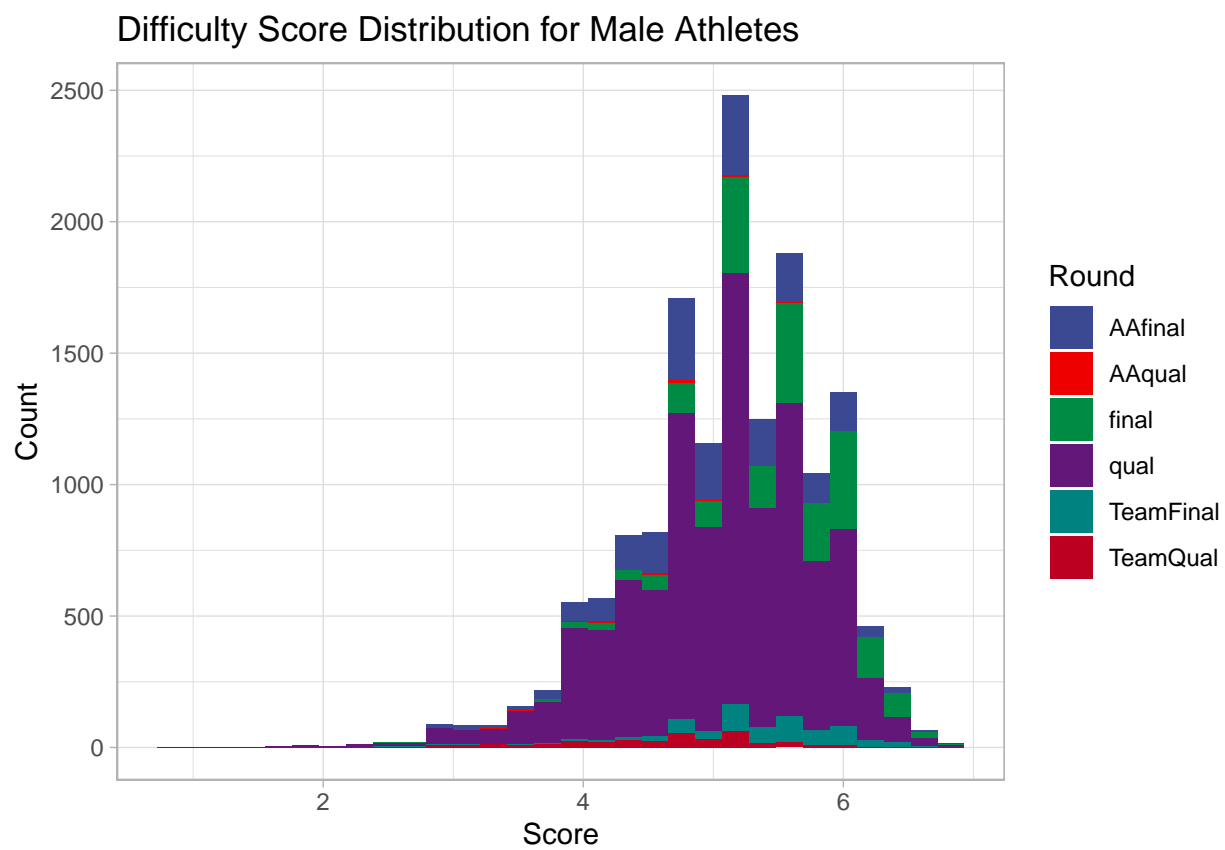Figure 6: Score Distribution for Women Athletes

Figure 7: Score Distribution for Male Athletes

Score Distribution for Women Athletes, by Round

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

## Score Distribution for Men Athletes, by Round



```
## [1] 12.21275
```
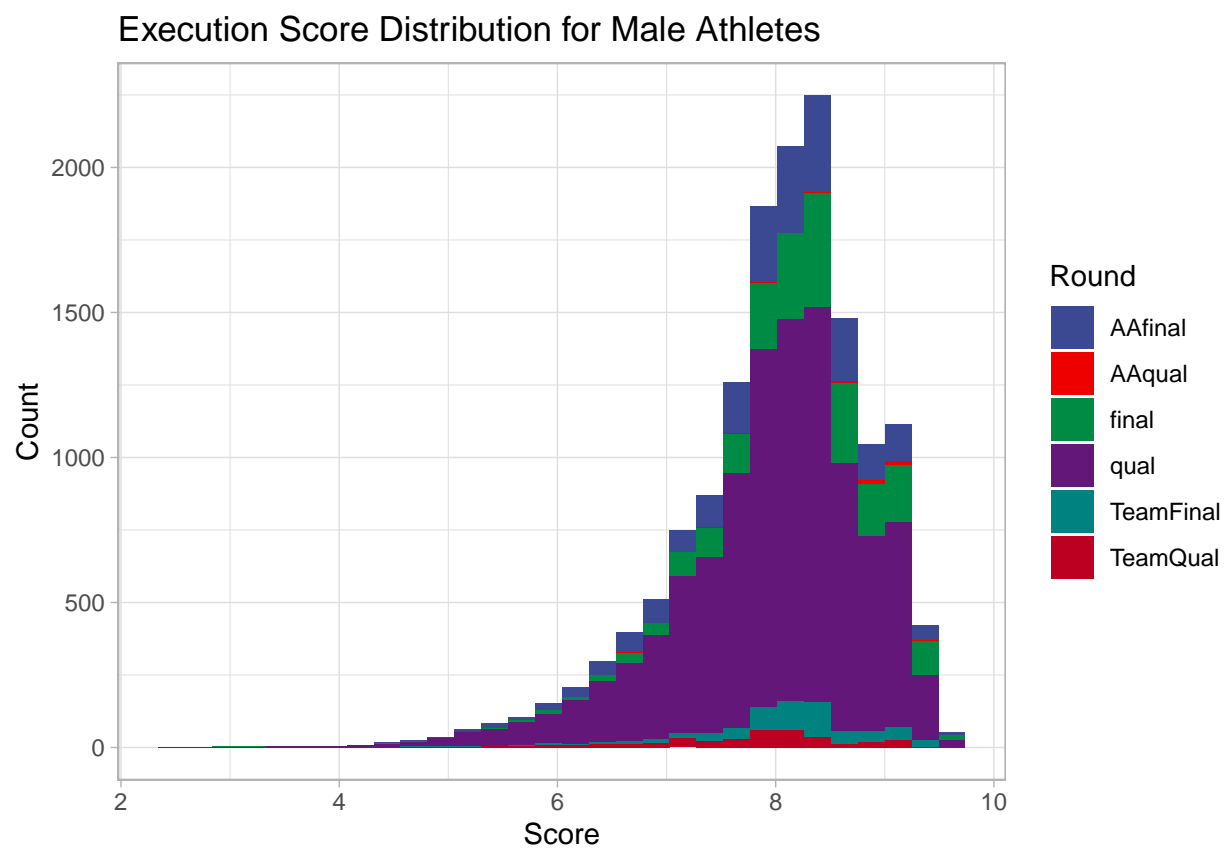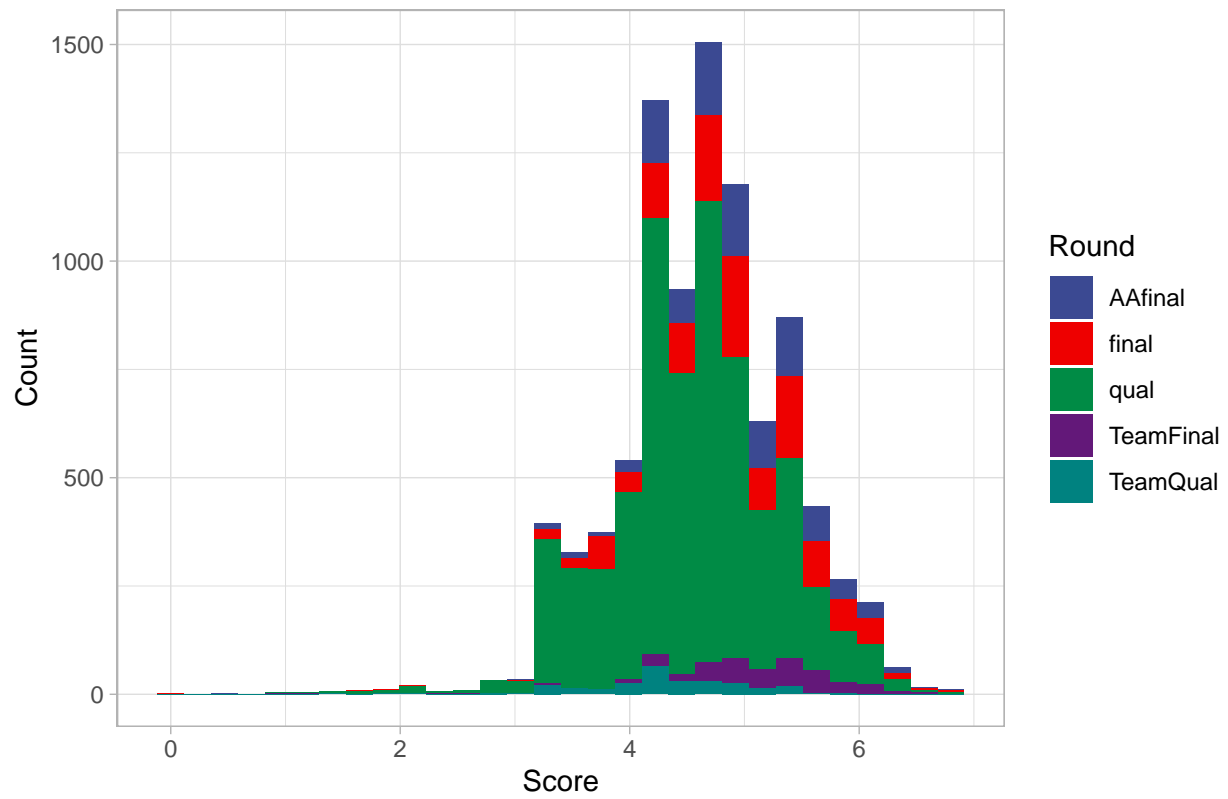
```
## [1] 12.4
```

```
## [1] 13.04569
```

```
## [1] 13.233
```

Execution scores have a higher range than difficulty scores, for both men and women athletes (colored by round but doesn't seem to be much difference based on round):

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```
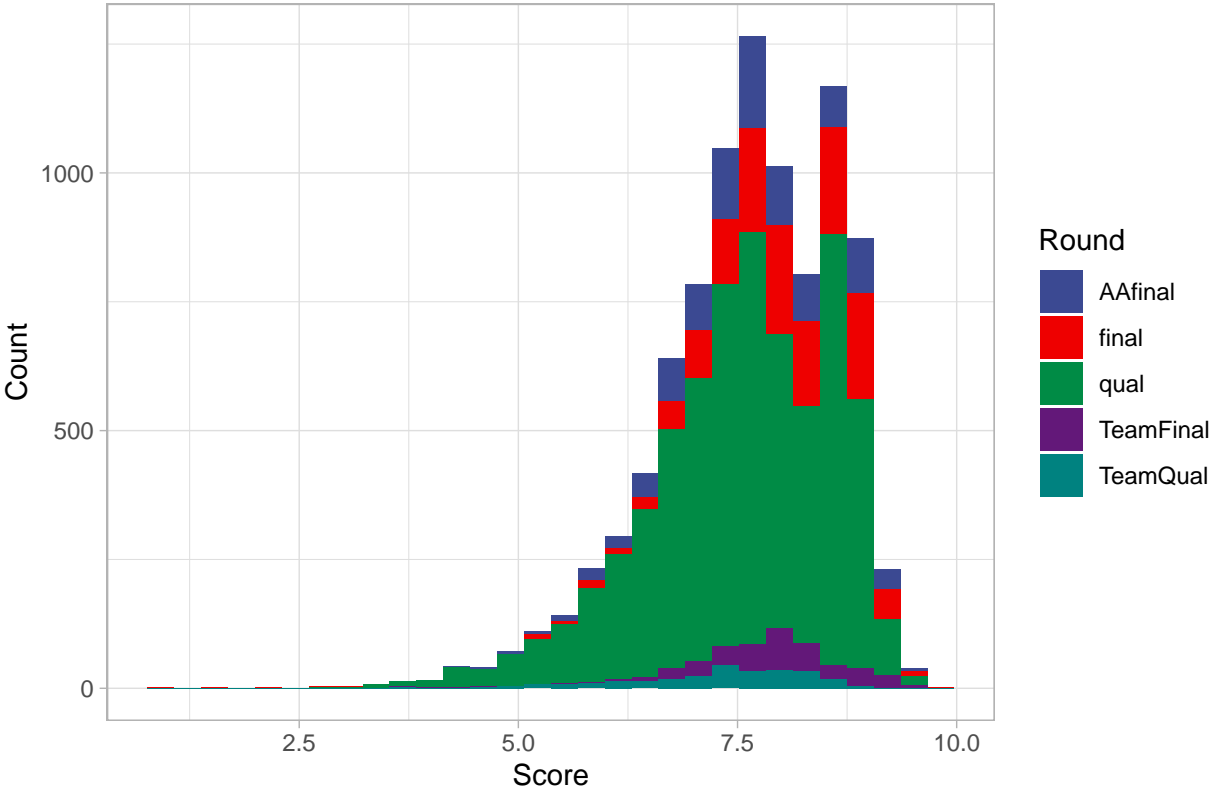
Figure 8: Difficulty Score Distribution for Male Athletes

# Execution Score Distribution for Male Athletes



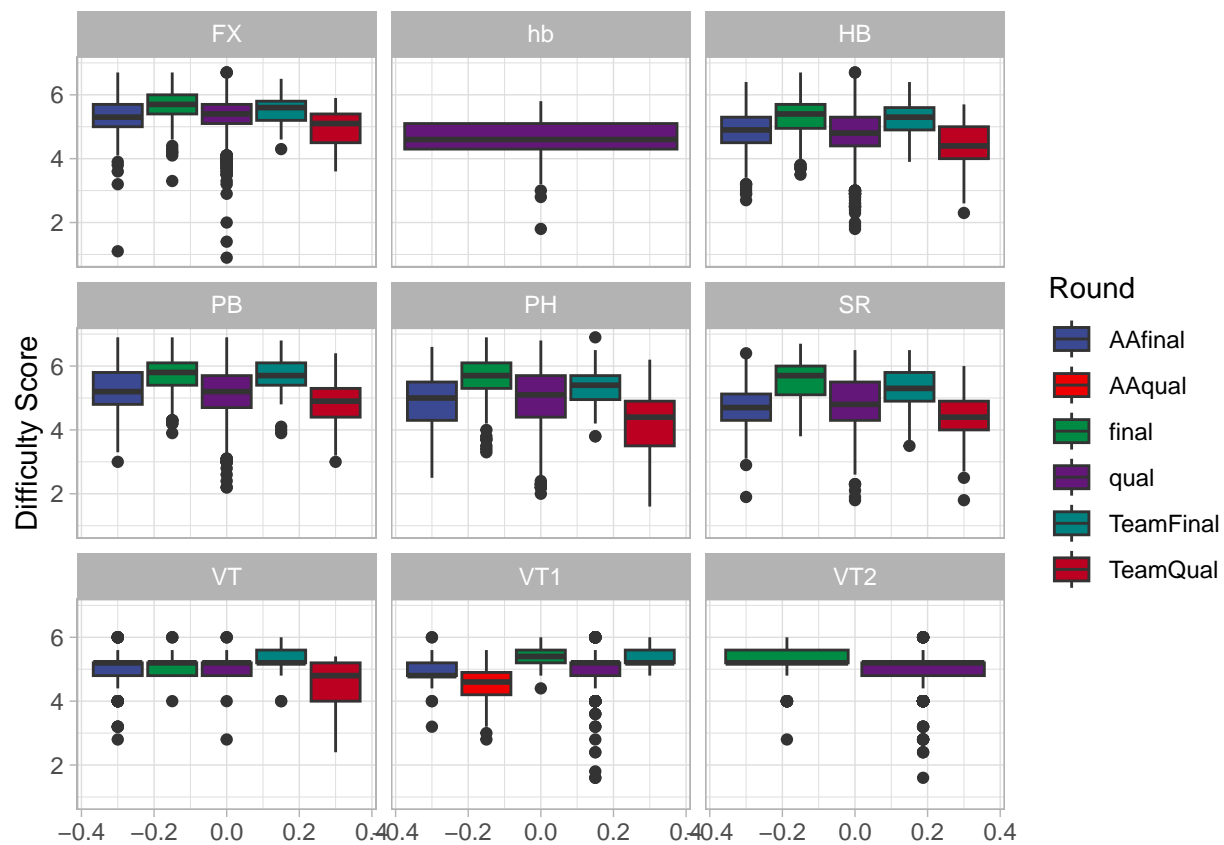Figure 9: Execution Score Distribution for Male Athletes

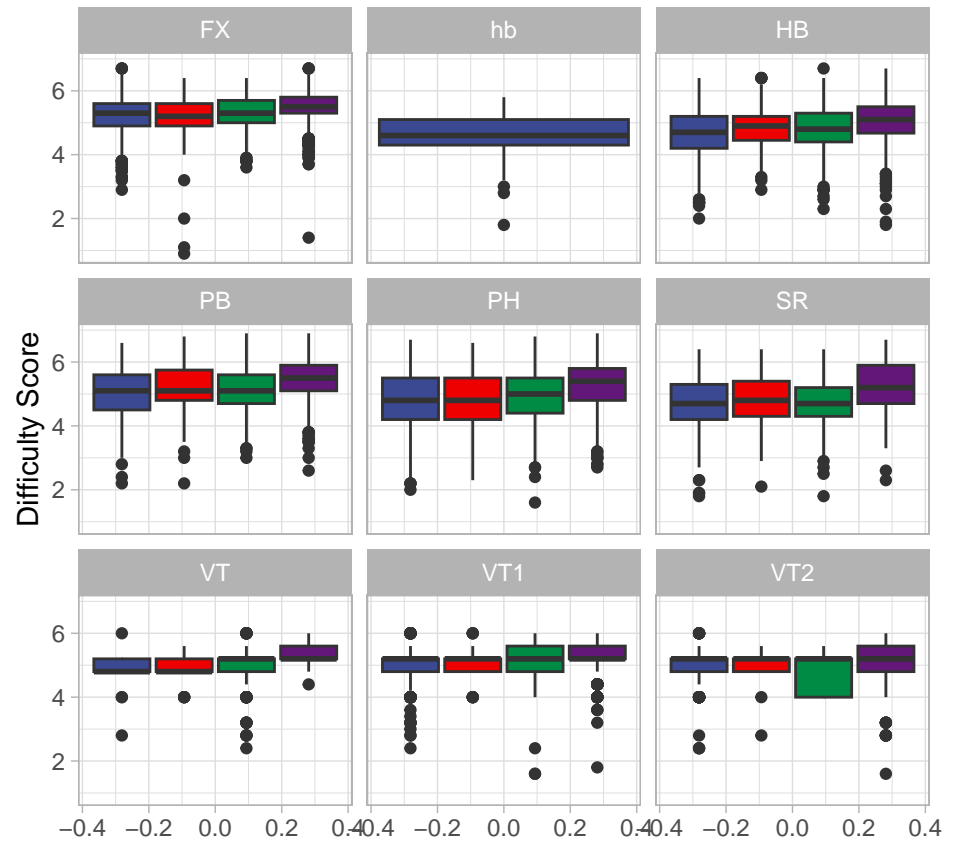# Difficulty Score Distribution for Women Athletes



```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

Execution Score Distribution for Women Athletes
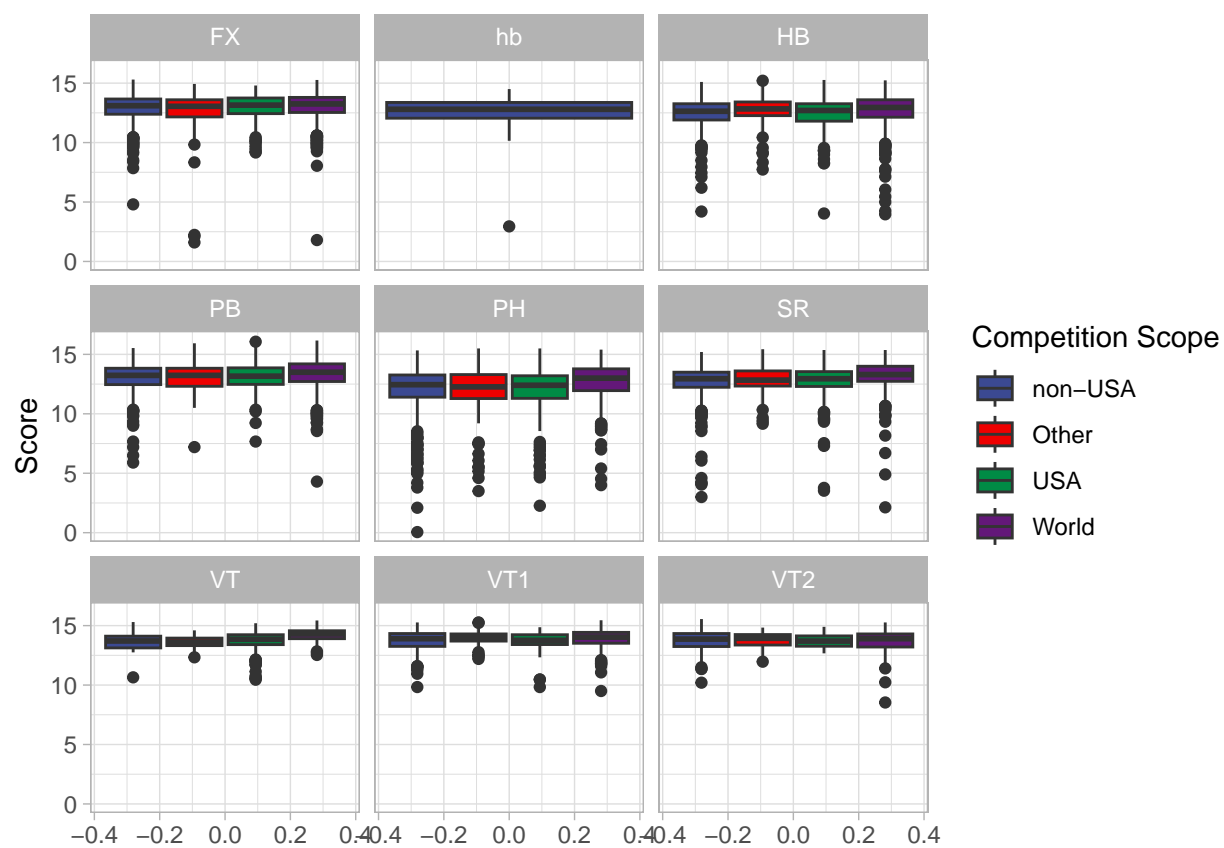
Waiting for competition scope for men

Figure 10: Distribution of Scores by Scope of Competition and Apparatus for Men