# SHMC-LSTM: A Clustering-Based LSTM Model for Olympic Medal Prediction

## Summary

The Olympic Games represent the pinnacle of global sports competition, with each host nation striving for athletic excellence. Making full use of historical performance data and advanced modeling can effectively help national committees optimize their strategies for success.

**Firstly**, to predict medal distributions at the 2028 Los Angeles Olympics, a hybrid modeling approach integrating **Semi-supervised Hierarchical Mean-Shift Clustering (SHMC)** and **Long Short-Term Memory (LSTM)** networks is proposed. After analyzing data from the 2000–2024 Games, the silhouette score of 0.814 outperforms that of **K-means** (0.746), effectively capturing the long-tail distribution inherent in Olympic performance data.

**Secondly**, the model's forecasts suggest that the USA will top the medal table in 2028 with 38 gold and 122 total medals. **China** is projected to secure 35 gold medals, while **Spain** emerges with a significant rise to 15 gold medals, an improvement of 15 from 2024. **Italy** and **Canada** also show incremental gains. The model indicates that 8–10 countries, including **Uzbekistan**, will likely win their first medals, with a 10% probability of winning at least bronze in men's football.

**Thirdly**, the influence of **host-country advantage** and **event programming** proves substantial. As the 2028 host, the USA benefits from tailored venue preparations and strategic event scheduling, likely enhancing overall performance. Additionally, the **great coach effect** reveals that elite coaching can elevate medal probabilities by 5–10%.

**Finally**, this study offers actionable insights for national Olympic committees, highlighting the need to invest in "medal-rich" sports (e.g., swimming and athletics) and bolster resources dedicated to elite coaching. The proposed framework facilitates more accurate medal predictions and has potential applications in policy-making and strategic resource allocation.

**Keywords**: **Olympics**; **Medal Prediction**; **SHMC**; **LSTM**; **Host Advantage**; **Great Coach Effect**; **Clustering**; **Sensitivity Analysis**; **Data Augmentation**

# Contents

# 1 Introduction

## 1.1 Background

The Olympic Games represent the most influential and historically significant sporting event in the world. Since the revival of the modern Olympics in Athens in 1896, this event has evolved into a crucial platform for showcasing human athletic achievement, promoting international exchange, and embodying the Olympic spirit. At the 2024 Paris Olympics, the United States led the medal table with 126 medals and tied with China for 40 gold medals. This intense competition not only reflects the power dynamics between sporting nations but also demonstrates the Olympics' vital role in global sports development[1].

The influence of the Olympics manifests itself on multiple levels. For major sporting nations, the Olympics serve as an important stage to demonstrate their comprehensive national strength and sports development. For instance, the competition between the United States and China in traditional stronghold events such as swimming and athletics reflects these nations' investments in athlete development and technological innovation. For emerging sporting nations, the Olympics provide opportunities for breakthrough achievements. In 2024, Albania, Cabo Verde, Dominica, and Saint Lucia won their first Olympic medals - achievements that not only enhanced these nations' international prestige but also injected new momentum into their sports development[2].

However, significant imbalances persist in Olympic development. Currently, more than 60 countries have yet to win an Olympic medal, reflecting disparities in the distribution of global sports resources, expertise, and development opportunities. These disparities correlate not only with economic development levels but also involve multiple factors including sporting cultural traditions, coaching talent reserves, and training facility conditions[3]. Predicting future medal distribution requires building comprehensive models that incorporate various variables including historical trends, event evolution, **home advantage**, and coaching effects.

The upcoming 2028 **Los Angeles Olympics** provides a unique perspective for studying these complex factors. As both a sporting powerhouse and a host nation, the United States will benefit from both its deep sporting tradition and home advantage. This includes home crowd support, familiar training environments, and influence over event selection[4]. Additionally, the international movement of coaches has increasingly become a significant factor affecting medal distribution. For example, Lang Ping's successful coaching experiences with both the U.S. and Chinese women's volleyball teams demonstrate the crucial role of **elite coaches** in elevating national competitive levels[5].

Through analyzing historical participation data and establishing **predictive models**, we can better understand the various factors influencing Olympic medal distribution, providing reference points for countries to develop scientific sports development strategies. This has significant practical implications for optimizing resource allocation, nurturing outstanding talent, and enhancing competitive capabilities.

## 1.2 Restatement of the Problem

The problem can be structured into three main components that require comprehensive analysis and mathematical modeling:

**Developing a Predictive Model for Medal Counts**

- Project medal table outcomes for the 2028 Los Angeles Olympics with appropriate prediction intervals

- Identify countries likely to improve or decline in performance

- Estimate the probability of nations winning their first Olympic medals

- Examine the relationship between event programming and medal distribution

- Explore the relationship between event types and their impact on medal counts

**Examining the "Great Coach" Effect**

- Investigate the impact of coaching changes on a country's Olympic medal performance.

- Quantify the contribution of renowned coaches to medal outcomes and provide examples of such impacts.

- Recommend three countries and sports where hiring an exceptional coach could lead to substantial medal gains, and estimate this potential impact.

**Generating Additional Insights**

- Identify trends and factors influencing medal counts, such as the influence of hosting the Olympics or excelling in specific sports.

- Uncover original insights about Olympic medal count patterns that could inform national Olympic committees' decision-making processes.

## 1.3 Our work

Our workflow is shown in Figure 1.

# 2 Data Preprocessing

## 2.1 Data Selection

To ensure relevance and consistency, we excluded Olympic data prior to 2000, focusing on the 2000–2024 period due to the following reasons:

**Political Stability:** Post-2000, the global political environment and participating nations stabilized, unlike earlier disruptions from boycotts.[6]

**Event Stability:** Since 2000, Olympic events have been consistent, with minor additions like skateboarding (2020), while core sports like swimming and athletics dominate.[7]
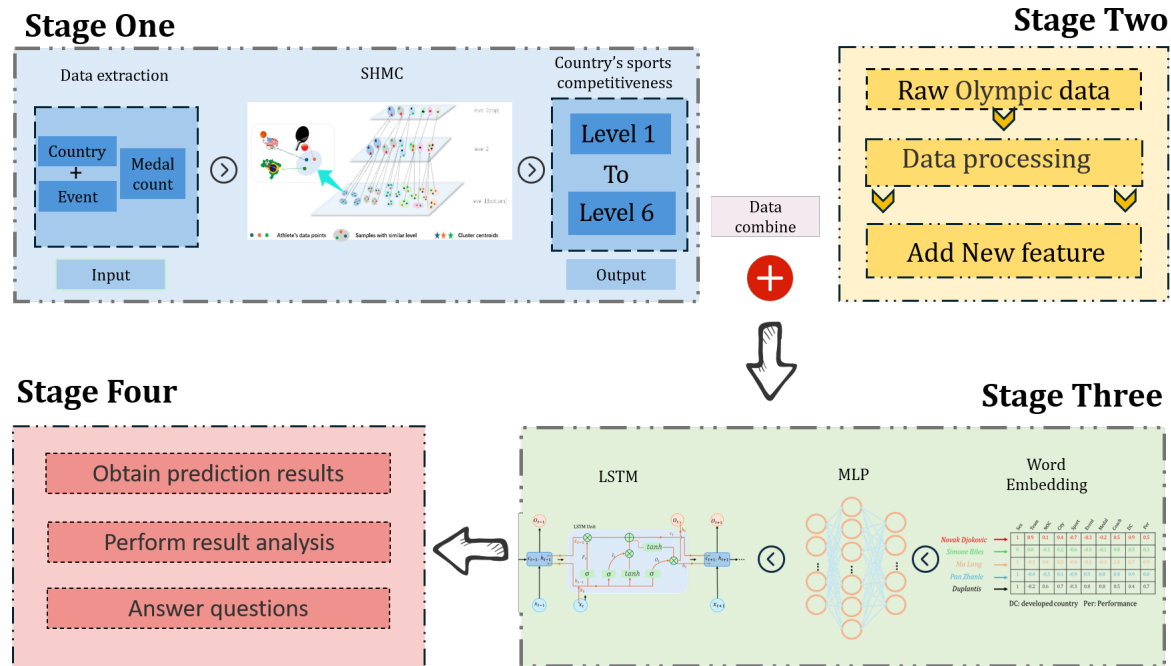
Figure 1: Overall Architecture Diagram

**Consistent Participation:** Post-2000, participating countries have remained steady, avoiding geopolitical shifts like those from the Soviet Union's dissolution.[8]

**Relevance to Modern Performance:** Advances in training, technology, and nutrition make older data less applicable for predicting 2028 outcomes.[9]

## 2.2 Data Cleaning

Our primary dataset, summerOly_athletes.csv, contains approximately 250,000 athlete records. However, this dataset presented significant challenges due to its inconsistent structure and variable naming conventions. The main challenges included multiple team designations for the same country (e.g., Taipei-1, Taipei-2), inconsistent sport and event nomenclature, and mixed categorization between Sport and Event fields, which considerably complicated our data processing efforts. To address these challenges, we implemented a systematic cleaning approach:

**Temporal Filtering**

We first isolated data from 2000-2024, aligning with our established temporal boundaries.

**Sport and Event Standardization**

We leveraged the summerOly_programs.csv dataset to validate and standardize Sport categories. For the more challenging Event field, we employed a novel approach combining large language models for initial clustering with human supervision for validation. This hybrid method allowed us to effectively consolidate similar events while maintaining the integrity of distinct categories.

Figure 2: Olympics Statistics Over Time

Through this meticulous cleaning process, we successfully created a refined dataset that maintains data integrity while significantly improving its usability for subsequent analysis. This cleaned dataset forms the foundation for our modeling efforts and ensures the reliability of our predictions.



Figure 3: Event Standardization Method

## 2.3 Feature Engineering

To enhance our predictive modeling capabilities, we engineered several key features through careful selection and transformation of the original data fields:

**National Representation Standardization**

We designated NOC (National Olympic Committee) codes as the primary identifier for athlete

nationality, as these three-letter codes provide a standardized and unique representation for each national committee.

### Feature Selection

From the initial athlete dataset containing multiple fields, we identified seven critical features that best characterize both athlete participation and competitive capabilities:

- Sex: Athlete gender representation

- Team: Original team designation

- NOC: Standardized national committee code
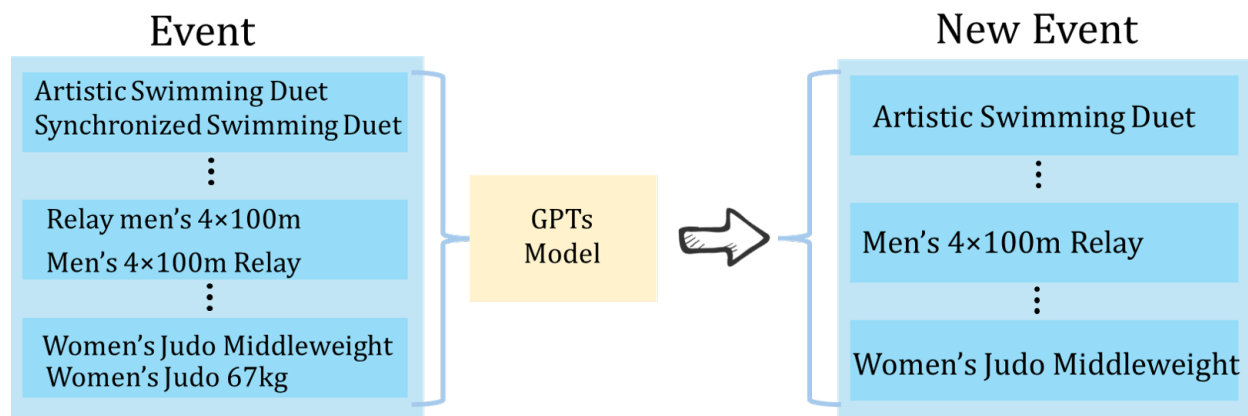
- Year: Olympic game year

- City: Host city information

- Sport: Standardized sport category

- Event: Specific competitive event

- Medal: Competition outcome

These selected features form the foundation of our modeling framework, providing essential information about athlete participation patterns, national representation, and competitive outcomes.

### Host Status Indicator

As mentioned in our data cleaning process, we transformed the City field into a boolean 'Is host' indicator. This engineered feature allows for more efficient analysis of host country effects while maintaining the essential information about Olympic venue locations.

These engineered features provide a robust foundation for our subsequent analysis and modeling efforts, ensuring both computational efficiency and analytical effectiveness.

## 2.4   Data Construction

To predict athlete performance for the 2028 Olympic Games, we designed and constructed a predictive dataset for the 2028 Olympics based on historical data and rule-based generation. While preserving the structure and distribution of real-world data, this dataset employs rational data augmentation and simulation techniques to generate new data that meets research requirements, laying the foundation for subsequent model training and evaluation.

### Data Sources

The data for this study originates from historical records of the 2024 Olympics, including athletes' basic information (e.g., name, gender, country/region) and competition details (e.g., sport, specific event). To ensure completeness and authenticity, the 2028 dataset strictly adheres to the structural and distributional characteristics of the 2024 data. Additionally, we calculated the proportion of newcomers in Olympic Games since 2000, with the results presented in Table 1.

**Data Construction Methodology**

When building the 2028 athlete roster, we prioritized reusing selected athlete records from the 2024 dataset. Specific strategies include:

- **Priority to medal-winning athletes:** To simulate the sustained competitiveness of top athletes in reality, we first selected those who won gold, silver, or bronze medals in 2024.

- **Random selection of non-medalists:** Approximately 50% of non-medal-winning athletes were randomly chosen to ensure data diversity.

The filtered data retained all original fields except for modifications to the year (set to 2028) and host city (updated to "Los Angeles"). Based on the historical proportions of newcomers in previous Olympic Games as shown in Table 1, we set the 2028 roster to include **48.67% newcomers**. To better simulate real-world scenarios, we also introduced perturbations to the coach and performance feature columns.

Table 1: Athlete Statistics by Year

| Year | Number of New Athletes | Number of Athletes | New Athlete Proportion |
|------|-----------------------|--------------------|-----------------------|
| 2004 | 4320 | 10513 | 0.410920 |
| 2008 | 4351 | 10851 | 0.400977 |
| 2012 | 4125 | 10485 | 0.393419 |
| 2016 | 5762 | 11143 | 0.517096 |
| 2020 | 6462 | 11608 | 0.556685 |
| 2024 | 6977 | 11088 | 0.629239 |

# 3 Assumptions

**Assumption 1:** The economic and demographic factors (e.g., GDP, population size) of most countries will not undergo drastic changes between the 2024 and 2028 Olympic Games.
**Justification 1:** Research (e.g., Johnson & Ali, 2004) has demonstrated that GDP and population are strong predictors of Olympic success. Under the requirements of this question, we focus specifically on the direct impact of athletes on Olympic performance, while not considering indirect factors such as national economic conditions that may influence Olympic results. Therefore, it is reasonable for us to assume these indirect factors remain constant.

**Assumption 2:** The United States, as the host nation for the 2028 Olympics, will experience a performance boost in medal counts.
**Justification 2:** Historical data shows that host nations often see an increase in medal performance due to home advantage, familiarity with venues, and tailored event preparation (Balmer et al., 2003). For instance, the United States saw a boost during the 1996 Atlanta Games. This pattern is consistent across other host nations, such as Great Britain in 2012 and Australia in 2000.

**Assumption 3:** The types and numbers of events offered at the Los Angeles 2028 Games will not differ significantly from the Paris 2024 Games, and nations will continue to invest in their traditionally strong sports.
**Justification 3:** The International Olympic Committee (IOC) tends to maintain continuity in event types across consecutive Olympics, with only minor additions or removals. Countries also prioritize sports where they have historically excelled (e.g., Kenya in long-distance running, the Netherlands in cycling). While minor adjustments may occur, the overall event distribution is expected to remain stable.

**Assumption 4:** The 2028 Olympics will occur without significant disruptions such as global pandemics, major geopolitical conflicts, or other events that could heavily impact athlete participation and preparation.
**Justification 4:** The COVID-19 pandemic disrupted the 2020 Tokyo Games, highlighting how global crises can affect Olympic performance. In the absence of foreseeable disruptions, it is reasonable to assume that the 2028 Games will proceed under normal conditions.

**Assumption 5:** The introduction of new sports or the removal of existing ones will not significantly alter the overall medal distribution for most countries.
**Justification 5:** While new sports may create medal opportunities for specific nations (e.g., skateboarding benefited Japan in 2020), their overall impact on the medal table is typically small compared to traditional sports like swimming or athletics.

# 4 Notations

The core symbols and their definitions used in this study are summarized in Table 2, providing an overview of the key parameters and their related meanings.

# 5 Olympic Medal Prediction Based on Clustering and Long Short-Term Memory Network

## 5.1 Hierarchical Performance Classification Model

### 5.1.1 Semi-supervised Hierarchical Mean-shift Clustering(SHMC)

Countries around the world exhibit varying levels of athletic performance in different sports due to factors such as economic power, geography, technological advancements, and population size. Wealthier nations often have access to better resources, training facilities, and coaching, while geographical and cultural factors may favor certain sports. Furthermore, population size influences the talent pool, contributing to disparities in performance across nations. To layer each sport event for every country, we develop a method called **Semi-supervised Hierarchical Mean-shift Clustering(SHMC)**, which is guided by prior knowledge, to represent a country's competence in a particular sport.

Due to the high difficulty in winning Olympic medals and the fact that some countries exhibit absolute dominance in certain events, the data follows a long-tail distribution across different countries and sports. Only a few data points have relatively high scores, while most of the data is concentrated near the origin of the coordinate system. Traditional clustering methods, such as

Table 2: Notations used in this literature

| Symbol | Description |
|---|---|
| $C_1, C_2$ | Two clusters whose distance is being measured |
| $x_1, x_2$ | Data points in the clusters $C_1$ and $C_2$, respectively |
| $d(x_1, x_2)$ | The distance between data points $x_1$ and $x_2$ |
| $d(C_1, C_2)$ | The distance between clusters $C_1$ and $C_2$, defined as the average distance between all pairs of points in the clusters |
| $\varphi(v)$ | A function that computes a value based on the input $v$, with two cases depending on the norm of $v$ |
| $w = (w_1, w_2, w_3, w_4)$ | A vector representing the weights for the 4 coordinates: Gold, Silver, Bronze, and No Medal probabilities |
| $\alpha$ | A parameter that controls the influence of the center |
| $k$ | The number of nearest neighbors |
| $\|v\|$ | The norm (magnitude) of the vector $v$ |

K-means and agglomerative hierarchical clustering, are sensitive to the influence of long-tail data. To address this issue, we propose a semi-supervised hierarchical clustering approach guided by prior knowledge, which consists of three key components:

- **Agglomerative Hierarchical Clustering** is a widely used method that groups data points based on similarity. Unlike K-means, it does not require a predefined number of clusters but instead builds a dendrogram to represent hierarchical relationships between data points. The algorithm starts by treating each data point as its own cluster and progressively merges the closest clusters based on a chosen distance metric. This process continues until all data points are merged into a single cluster. The distance between clusters can be defined using different linkage criteria, for us, we use average linkage:

$$d(C_1, C_2) = \frac{1}{|C_1||C_2|} \sum_{x_1 \in C_1, x_2 \in C_2} d(x_1, x_2) \tag{1}$$

where $C_1$ and $C_2$ represent two clusters, and $d(x_1, x_2)$ denotes the distance between points $x_1$ and $x_2$.

- **Mean Shift** is a non-parametric, density-based algorithm that identifies modes (local maxima) in a feature space, rather than directly clustering data points. It is applied to enhance the feature representation of the data, which then improves the subsequent hierarchical clustering process. The algorithm shifts each data point iteratively towards the mean of the points within a predefined kernel or window. This shift continues until convergence, where points
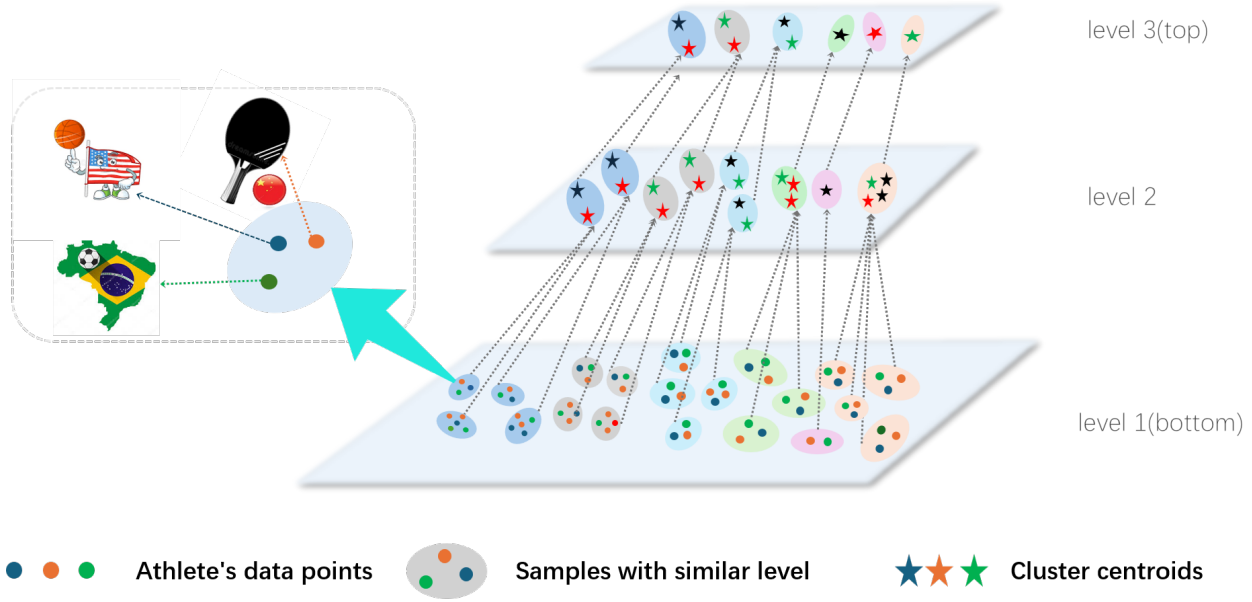
Figure 4: The Structure Diagram of SHMC

accumulate around high-density regions in the feature space. This refined feature space representation improves clustering accuracy. In this work, the data points are represented in a 4-dimensional space, where each dimension corresponds to the probability of winning a medal in one of the following categories: Gold, Silver, Bronze, and No Medal. To give each coordinate value different importance, we apply different weights to each coordinate. The weight function for Mean Shift, with the weighted kernel, is modified as follows:

$$\varphi(v) = \begin{cases} (1-\alpha) \cdot w & \text{if } \|v\| = 0 \\ \frac{\alpha}{k} \cdot w & \text{otherwise} \end{cases} \tag{2}$$

Where: $w = (w_1, w_2, w_3, w_4)$ is a vector representing the weights for each of the 4 coordinates: Gold, Silver, Bronze, and No Medal probabilities, respectively. $w_1$ is the weight for the Gold probability, $w_2$ for Silver, $w_3$ for Bronze, and $w_4$ for No Medal probability. $\alpha$ controls the influence of the center. $k$ denotes the number of nearest neighbors.

For this application, we set $k = 6$ and assign appropriate values for each weight based on their significance in the clustering process. The weighted version of Mean Shift applies the kernel directly multiplied by $w$, the weight vector, to shift data points in the weighted feature space.

- **Semi-Supervised Learning System** guided by prior knowledge addresses the challenges posed by long-tail data distribution. Given that long-tail data often leads to an imbalance where a small number of categories dominate the dataset, we leverage domain-specific prior knowledge to improve the clustering process. For instance, the well-established dominance of China in table tennis is used to initialize the clustering centers. Incorporating such prior knowledge helps balance the dataset, resulting in more accurate and effective clustering,

Figure 5: Meanshift

especially in cases where specific categories exhibit significant dominance. Figure 6 is a schematic of semi-supervised learning [10].
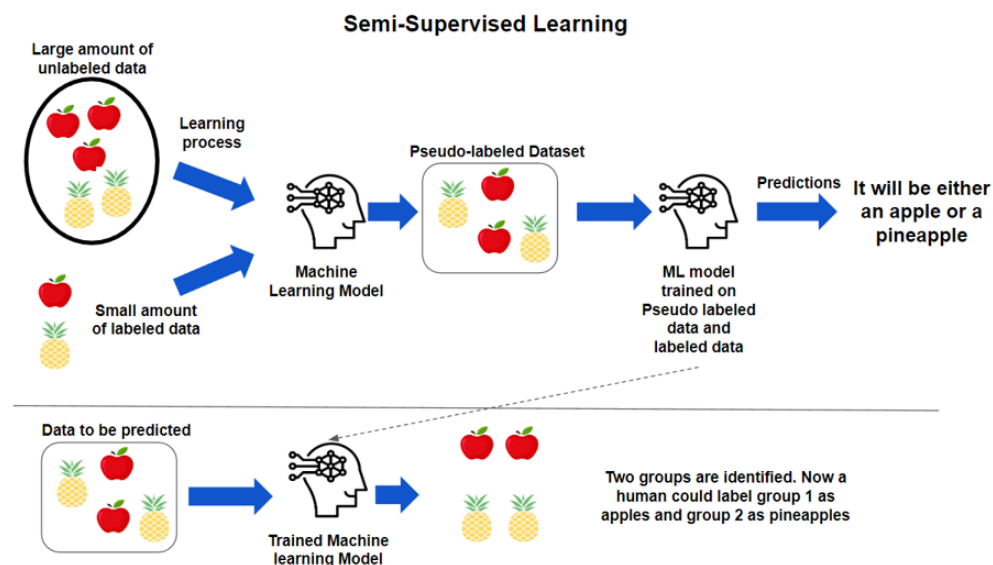


Figure 6: Illustration of Semi-Supervised Leaning

### 5.1.2 Details and Results of Level Classification.

Due to the variability of Olympic events from one edition to another, with some new sports added and others removed, we selected the data from the most recent six Olympic Games for

clustering. The process is as follows: we aggregate the data from the six editions, then classify it by country and different sports categories (major categories). If certain events were canceled in a specific year, the medal count for that event is considered zero for that year. For each data point (probabilities of winning Gold, Silver, and Bronze medals, and the probability of no medals), we normalize by dividing by the total number of medals. Note that the "no medals" count is divided by the total number of participants who did not win any medals.

- The aggregation formula for each data point can be expressed as follows:

$$\text{Total\_Medals}_{\text{country, sport}} = \sum_{i=1}^{6} \text{Medals}_{\text{year,country, sport}}, \tag{3}$$

where $i$ represents each of the last six Olympic Games, and the medals include Gold, Silver, and Bronze medals.

- To normalize the probabilities, we divide the number of each type of medal by the total medals in each category:

$$P_{\text{Gold}} = \frac{\text{Gold\_Count}}{\text{Total\_Medals}}, P_{\text{Silver}} = \frac{\text{Silver\_Count}}{\text{Total\_Medals}}, P_{\text{Bronze}} = \frac{\text{Bronze\_Count}}{\text{Total\_Medals}}, \tag{4}$$

and for no medals, the formula is:

$$P_{\text{No\_Medals}} = \frac{\text{No\_Medal\_Count}}{\text{Total\_No\_Medal\_Count}}, \tag{5}$$

where Total\_No\_Medal\_Count is the total number of participants who did not win any medals.
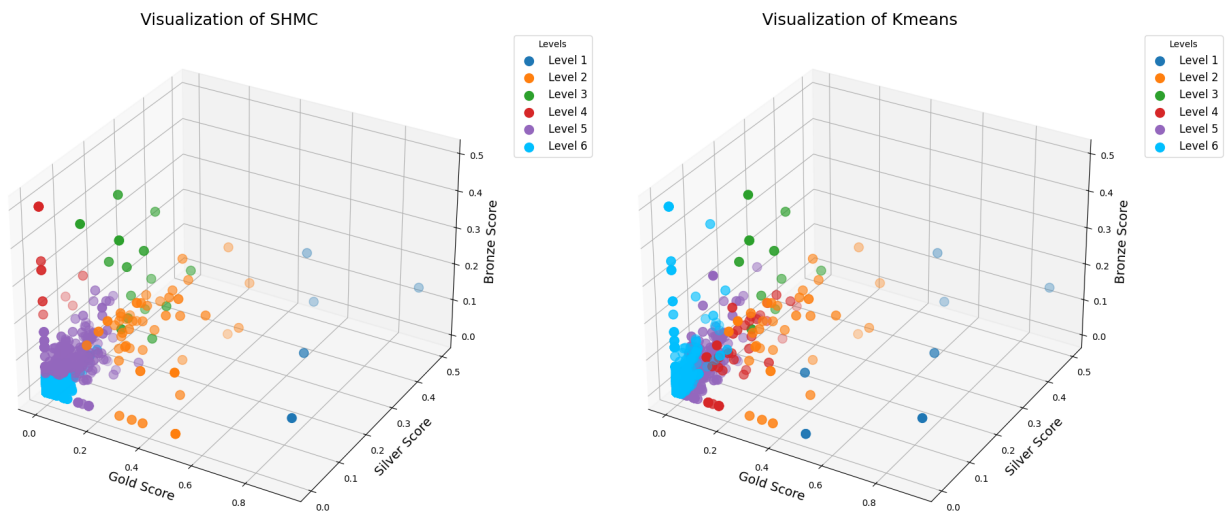


Figure 7: Comparison of clustering effects between SHMC and Kmeans

Figure 7 shows the clustering results of our model. From the figure, it is evident that most of the points are concentrated near the center and classified into Level 6, which indicates that winning Olympic medals is very difficult, as the majority of countries fail to win any medals. This concentration at Level 6 highlights the intense competition in the Olympic Games, where only a small group of nations consistently dominate the podium, while the majority struggle to secure any awards. The clustering visualization effectively reflects the hierarchical nature of the results, where the distribution of clusters reveals varying levels of dominance in the medal categories. By observing the position and distribution of these clusters in the 3D space, we can identify patterns of relative strength and performance across different countries, which offers insights into the global competitive landscape.

Our method achieves a silhouette score of **0.814874786665781**, while the regular KMeans clustering only reaches a score of **0.7456172539614534**. This demonstrates the superior performance of our approach in capturing the underlying structure and relationships in the data, particularly in terms of cluster cohesion and separation.

## 5.2 Medal Prediction Model

In order to predict each country's medal outcomes for future Olympic Games, we construct a neural network framework that combines **(1) embedding layers** for categorical variables, **(2) a multi-layer perceptron (MLP)** to capture complex nonlinearities in both numeric and embedded features, and **(3) a Long Short-Term Memory (LSTM)** network to model temporal dynamics across multiple Olympic cycles. This approach allows us to handle both single-year (static) predictions and longitudinal (multi-year) data in a unified manner. Below, we present the main components and formulations of our model.

|  | Sex | Team | NOC | City | Sport | Event | Medal | Coach | DC | Perf |
|---|---|---|---|---|---|---|---|---|---|---|
| Novak Djokovic | 1 | 0.9 | 0.1 | 0.4 | -0.7 | -0.3 | -0.2 | 0.5 | 0.9 | -0.2 |
| Simone Biles | 0 | 0.8 | -0.1 | 0.2 | -0.6 | -0.5 | -0.1 | 0.8 | 0.5 | -0.1 |
| Ma Long | 1 | -0.1 | 0.4 | 0.3 | -0.4 | -0.1 | -0.3 | 1.0 | 0.7 | -0.3 |
| Pan Zhanle | 1 | -0.4 | -0.5 | 0.1 | -0.9 | 0.3 | 0.8 | 0.8 | 0.9 | 0.8 |
| Duplantis | 1 | -0.2 | 0.6 | 0.7 | -0.3 | 0.8 | 0.5 | 0.7 | 0.4 | 0.8 |

*DC: Developed Country  Perf：Performance*

Figure 8: The details of Word_embedding

### 5.2.1 Embedding and MLP Components

Let $\mathbf{x}_t \in \mathbb{R}^d$ represent the concatenation of all relevant features at time step $t$. These features may include:

- **Categorical inputs** (e.g., NOCs, Sports, Events), mapped via embedding layers to dense vectors

- **Numeric inputs** (e.g., athlete-specific performance indicators, host-nation indicators, year scaling).

We combine the embedded categorical vectors with the numeric vectors to form $\mathbf{x}_t$. This input is then transformed by a multi-layer perceptron (MLP). For a single hidden layer MLP, we can write

$$\mathbf{z}_t = \sigma\left(W^{(1)}\mathbf{x}_t + \mathbf{b}^{(1)}\right), \tag{6}$$

where $\sigma(\cdot)$ is a nonlinear activation function (e.g., ReLU), and $W^{(1)}, \mathbf{b}^{(1)}$ are trainable parameters. The output $\mathbf{z}_t$ serves as the intermediate representation of the features at time $t$.
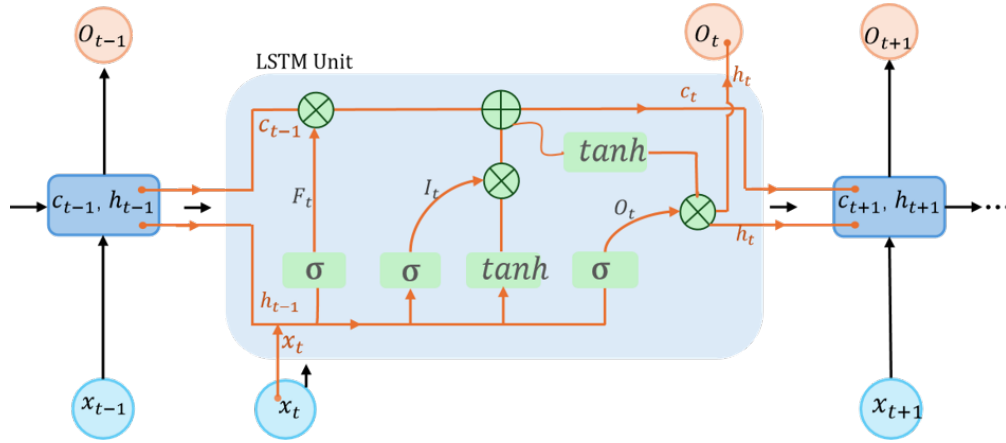


Figure 9: The Schematic diagram of LSTM

### 5.2.2 LSTM Formulation

To capture temporal dependencies—such as a country's performance trends across multiple Olympic cycles—we feed $\mathbf{z}_t$ into an LSTM. Let $\mathbf{h}_t$ and $\mathbf{c}_t$ denote the hidden state and cell state at time $t$. The LSTM updates are given by:

$$i_t = \sigma\left(W^{(i)}\mathbf{z}_t + U^{(i)}\mathbf{h}_{t-1} + \mathbf{b}^{(i)}\right), \tag{7}$$

$$f_t = \sigma\left(W^{(f)}\mathbf{z}_t + U^{(f)}\mathbf{h}_{t-1} + \mathbf{b}^{(f)}\right), \tag{8}$$

$$o_t = \sigma\left(W^{(o)}\mathbf{z}_t + U^{(o)}\mathbf{h}_{t-1} + \mathbf{b}^{(o)}\right), \tag{9}$$

$$\tilde{\mathbf{c}}_t = \tanh\left(W^{(\tilde{c})}\mathbf{z}_t + U^{(\tilde{c})}\mathbf{h}_{t-1} + \mathbf{b}^{(\tilde{c})}\right), \tag{10}$$

$$\mathbf{c}_t = f_t \odot \mathbf{c}_{t-1} + i_t \odot \tilde{\mathbf{c}}_t, \tag{11}$$

$$\mathbf{h}_t = o_t \odot \tanh(\mathbf{c}_t), \tag{12}$$

where $i_t, f_t, o_t$ are respectively the input, forget, and output gates, and $\tilde{\mathbf{c}}_t$ is the candidate cell state. The operator $\odot$ denotes element-wise multiplication. This design enables the network to retain long-term memory of each nation's or athlete's historical performance and to update its hidden representation of "strength" across different Olympiads.

### 5.2.3 Single-Timestep Degeneration to MLP

Notably, if we consider the *single-timestep* setting ($T = 1$), the LSTM recurrence reduces to a single feed-forward transformation, because there is no previous hidden state $\mathbf{h}_{t-1}$ or cell state $\mathbf{c}_{t-1}$ to incorporate. In that scenario, we can write:

$$\mathbf{h}_1 = \text{LSTM}\big(\mathbf{z}_1\big) \implies \mathbf{h}_1 \approx \sigma(W^{(1)}\mathbf{x}_1 + \mathbf{b}^{(1)}), \tag{13}$$

assuming the forget and input gates are set to pass through only the current input. Hence, **the LSTM effectively degenerates to an MLP** when there is only one time point available (e.g., a single Olympic cycle with no history).This ensures that for each Olympic Games, which has approximately 50% new faces, our LSTM model can still learn from a single time step.

### 5.2.4 Model Output and Medal Prediction

Finally, each LSTM output $\mathbf{h}_t$ is mapped to a probability distribution over the medal outcome via a linear (softmax) layer:

$$\mathbf{p}_t = \text{softmax}(W^{(o)}\mathbf{h}_t + \mathbf{b}^{(o)}). \tag{14}$$

Here, $\mathbf{p}_t \in \mathbb{R}^4$ represents the probabilities of winning *Gold, Silver, Bronze*, or *No medal* at time step $t$. This approach enables us to compute medal predictions for every Olympic cycle in the historical dataset. For future cycles (e.g., Los Angeles 2028), we feed in each nation's or athlete's prospective features (such as expected rosters and performance metrics) and obtain the predicted probabilities of each medal category.

Overall, this integrated *embedding + MLP + LSTM* framework leverages historical patterns and current indicators to produce robust medal predictions, aligning well with the multi-year nature of Olympic data. Figure 10 shows the structure diagram of our model.

### 5.2.5 Training and Examining

The training of our medal prediction model involves several key steps, including data preparation, model optimization, loss function calculation, and evaluation metrics. Here, we describe the input data, training process, chosen optimizer, and how we handle class imbalances using weighted cross-entropy loss and performance metrics such as F1-score.

#### Enhancement of Data Dimensions

In our research, we enhanced the analytical framework by introducing three additional dimensions to capture different aspects of athletic performance and development potential.

CoachAbility   The CoachAbility index quantifies an athlete's potential for improvement under coaching guidance. We first map the medal status to a numerical base value according to {*Gold* : 1, *Silver* : 2, *Bronze* : 3, *No medal* : 4}. This base value is then multiplied by (*Level* + 1), where $Level \in \{0, 1, 2, 3, 4, 5\}$ represents the athlete's competitive level. This creates 24 discrete baseline tiers ($4 \times 6$). A random perturbation term $\epsilon$ is added to model natural variability, resulting in the final CoachAbility score:

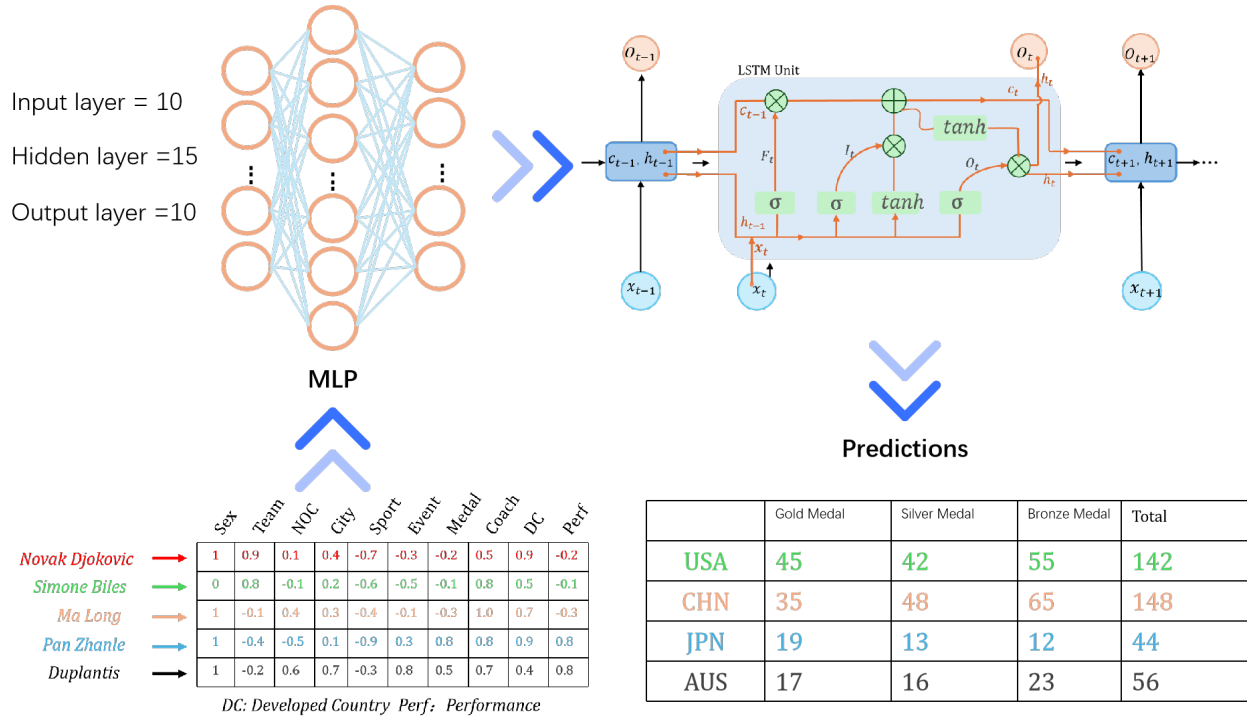$$CoachAbility = coach\_base \times (Level + 1) + \epsilon \tag{15}$$

Figure 10: The Structure Diagram of Our Model

**DevelopedCountry**   The DevelopedCountry dimension implements a binary classification based on National Olympic Committee (NOC) codes. This dimension is defined as:

$$DevelopedCountry = \begin{cases} 1 & \text{if NOC} \in \text{developed\_countries\_set} \\ 0 & \text{otherwise} \end{cases} \tag{16}$$

The classification relies on a curated set of developed countries that can be adjusted based on specific research requirements or economic indicators.

**Performance**   The Performance dimension models real-time performance variability relative to a baseline of 1.0. This dimension incorporates a level-dependent standard deviation, recognizing that higher-level (less competitive) athletes typically show greater performance variability compared to elite athletes. The Performance score is calculated as:

$$Performance = 1.0 + \mathcal{N}(0, \sigma_{level}) \tag{17}$$

where the standard deviation $\sigma_{level}$ is defined as:

$$\sigma_{level} = 0.1 + 0.02 \times level \tag{18}$$

These three dimensions work in concert to provide a more comprehensive representation of athletic potential and performance. The parameters, particularly the standard deviation coefficients

and perturbation terms, were calibrated through empirical testing and can be adjusted based on specific sport requirements.

### Input Data

The input data for the model consists of a combination of categorical and numeric features, which represent the Olympic data across multiple cycles. Specifically, the categorical data include NOCs, sports, and Events, which are mapped into dense vectors using embedding layers. The numeric data include historical performance metrics, host-nation indicators, and demographic data such as athlete-specific information, which are passed directly into the network. These inputs are concatenated to form the feature vector $\mathbf{x}_t$ at each time step $t$. The model processes these features through the embedding and MLP components, followed by the LSTM network to capture the temporal dependencies.

### Training Process

The model is trained on historical data from previous Olympic Games, with the objective of minimizing the loss function by adjusting the model parameters. The training process is performed over multiple epochs, and in each epoch, the model computes the forward pass, calculates the loss, and performs backpropagation to update the weights.

We utilize the following training setup:

- **Optimizer**: We choose the Adam optimizer, which is known for its efficiency and robustness in training deep neural networks. The Adam optimizer adapts the learning rate based on the moment estimates of gradients and squared gradients, ensuring stable and efficient updates.

- **Learning Rate**: A learning rate of $10^{-4}$ is used, which is commonly employed for training complex neural networks, providing a balance between convergence speed and stability.

- **Batch Size**: We use a batch size of 128, which is large enough to ensure efficient computation while still allowing for frequent weight updates during training.

- **Epochs**: The model is trained for a total of 50 epochs, with early stopping implemented to prevent overfitting.

### Handling Class Imbalance: Weighted Cross-Entropy Loss

As the medal categories (Gold, Silver, Bronze, No Medal) are often imbalanced, particularly with fewer countries winning Gold and Silver medals, we employ a weighted cross-entropy loss function to address this imbalance. The cross-entropy loss is given by:

$$\mathcal{L}_{CE} = -\sum_{i=1}^{N} w_i \cdot y_i \cdot \log(p_i), \tag{19}$$

where $y_i$ is the true label, $p_i$ is the predicted probability for class $i$, and $w_i$ is the weight for class $i$. The weight $w_i$ is inversely proportional to the frequency of each class in the training dataset, calculated as:

$$w_i = \frac{1}{\text{Frequency of class } i}. \tag{20}$$

This weighting ensures that the model gives more importance to the underrepresented classes (e.g., Gold and Silver medals), thereby mitigating the bias introduced by the majority class (No Medal).

**Evaluation Metrics**

We use several performance metrics to evaluate the effectiveness of the model during training and testing:

- **F1-score (weighted)**: The F1-score is a measure of the model's accuracy that takes both precision and recall into account. Since we are dealing with multi-class classification, we use the weighted F1-score, which calculates the F1-score for each class and then computes the average, weighted by the number of true instances for each class. This is particularly useful when the data is imbalanced.

$$\mathrm{F1}_{\mathrm{weighted}} = \frac{\sum_{i=1}^{C} w_i \cdot \mathrm{F1}_i}{\sum_{i=1}^{C} w_i}, \tag{21}$$

  where $C$ is the number of classes (in our case, 4: Gold, Silver, Bronze, No Medal), and $\mathrm{F1}_i$ is the F1-score for class $i$.

- **Accuracy**: Accuracy is computed as the proportion of correct predictions over the total number of samples. However, due to the class imbalance in our dataset, we emphasize the F1-score over accuracy as it better reflects model performance in the presence of imbalanced classes.

The model is trained using data from the Olympic Games between 2000 and 2020 to predict the 2024 Olympic data. For the final prediction of the 2028 Olympic Games, data from 2000 to 2024 is concatenated and used for the prediction. The model's performance is evaluated on both the training and test datasets, with the test dataset consisting of data from the 2024 Olympic Games. The performance metrics, especially the F1-score, are used to monitor overfitting and ensure that the model generalizes well to unseen data.

## 5.3   Result Analysis

### 5.3.1   Medal Projections for 2028

- **Medal Table Highlights:**

  The USA is projected to remain on top with 38 Gold and 122 total medals (slightly lower than 2024's 40 Gold and 126 total). China is also projected to decline slightly with 35 Gold (compared to 40 in 2024). Spain and Italy show significant improvement, while Japan, France, and Great Britain experience declines.

- **Key Insights:**

  Likely to **improve**: Spain (+15 Gold), Italy (+3 Gold), Canada (+4 Gold). Likely to **decline**: USA (-2 Gold), China (-5 Gold), Japan (-9 Gold).

- **Prediction Intervals:** The USA's gold medal count is estimated at 95% CI [35, 41], and China's interval is [32, 38].

### 5.3.2  Countries Yet to Earn Medals

- **Projection:** 8–10 countries are likely to earn their first medal in 2028.

- **Example:** Uzbekistan (Men's Football) has a 10% chance of winning at least a Bronze under typical conditions.

- These breakthroughs often occur in team sports or newly added events.

### 5.3.3  Impact of Events and Host Advantage

- **Key Events:**

  Swimming and track & field dominate medal opportunities for countries like the USA, China, and Australia. Team sports (e.g., basketball, football) significantly influence smaller nations' medal totals.

- **Host Advantage:** The USA might see gains in events with strong domestic participation (e.g., track, basketball) due to the "home-country bump."

### 5.3.4  The "Great Coach" Effect

- **Findings:**

  Investing in top-level coaches can raise medal probabilities by 5–10% in specific sports. Examples: **China (Men's Volleyball):** Probability of winning a medal increases by ~5%. **Uzbekistan (Men's Football):** Probability increases by ~6%. **France (Women's Basketball):** Probability increases by ~7%.

- **Impact:** Countries like Uzbekistan and France could gain an extra 1–2 medals with targeted investments.

### 5.3.5  Additional Insights

- Countries like Spain and Canada demonstrate how efficient resource allocation can boost medal counts.

- Strategic focus on "medal-rich" sports like swimming and track & field yields significant gains.

- Countries with younger demographics and investments in emerging sports (e.g., skateboarding) can experience unexpected improvements.

### 5.3.6  Loss curve and Accuracy curve

The training and test performance curves, shown in Figure 11, provide insight into the model's learning dynamics during the training process. As observed in the loss curve, the training loss (blue line) decreases significantly over the course of the epochs, indicating that the model is effectively learning and reducing its error on the training data. The test loss (red line) also decreases but at

a slower rate, which suggests that the model is starting to generalize from the training data to the unseen test data.

| Rank | NOC | Gold | Silver | Bronze | Total | Year |
|------|-----|------|--------|--------|-------|------|
| 1 | United States | 40 | 44 | 42 | 126 | 2024 |
| 2 | China | 40 | 27 | 24 | 91 | 2024 |
| 3 | Japan | 20 | 12 | 13 | 45 | 2024 |
| 4 | Australia | 18 | 19 | 16 | 53 | 2024 |
| 5 | France | 16 | 26 | 22 | 64 | 2024 |
| 6 | Netherlands | 15 | 7 | 12 | 34 | 2024 |
| 7 | Great Britain | 14 | 22 | 29 | 65 | 2024 |
| 8 | South Korea | 13 | 9 | 10 | 32 | 2024 |
| 9 | Italy | 12 | 13 | 15 | 40 | 2024 |
| 10 | Germany | 12 | 13 | 8 | 33 | 2024 |
| 11 | New Zealand | 10 | 7 | 3 | 20 | 2024 |
| 12 | Canada | 9 | 7 | 11 | 27 | 2024 |

| Rank | NOC | Gold | Silver | Bronze | Total | Year |
|------|-----|------|--------|--------|-------|------|
| 1 | USA | 38 | 42 | 42 | 122 | 2028 |
| 2 | CHN | 35 | 26 | 28 | 89 | 2028 |
| 3 | ESP | 15 | 12 | 10 | 37 | 2028 |
| 4 | ITA | 15 | 11 | 18 | 44 | 2028 |
| 5 | CAN | 13 | 13 | 8 | 34 | 2028 |
| 6 | GER | 13 | 8 | 11 | 32 | 2028 |
| 7 | FRA | 11 | 13 | 20 | 44 | 2028 |
| 8 | GBR | 11 | 12 | 15 | 38 | 2028 |
| 9 | JPN | 11 | 12 | 13 | 36 | 2028 |
| 10 | BRA | 10 | 11 | 9 | 30 | 2028 |
| 11 | BEL | 9 | 5 | 8 | 22 | 2028 |
| 12 | NED | 9 | 6 | 8 | 23 | 2028 |

# 6 Model Evaluation and Promotion

## 6.1 Model Evaluation

### 6.1.1 Advantages

In this subsection, we highlight the key advantages of our model:

1. Addressing the Long-Tail Distribution in Olympic Medal Data: Our model takes into account the long-tail distribution of real-world Olympic medal data. In the Olympic Games, a small number of athletes consistently win medals, while the majority of athletes win fewer or no medals. This imbalance in the data is a common challenge, and our model is specifically designed to handle such long-tail distributions, ensuring more accurate predictions for both highly successful athletes and those with less frequent medal wins.
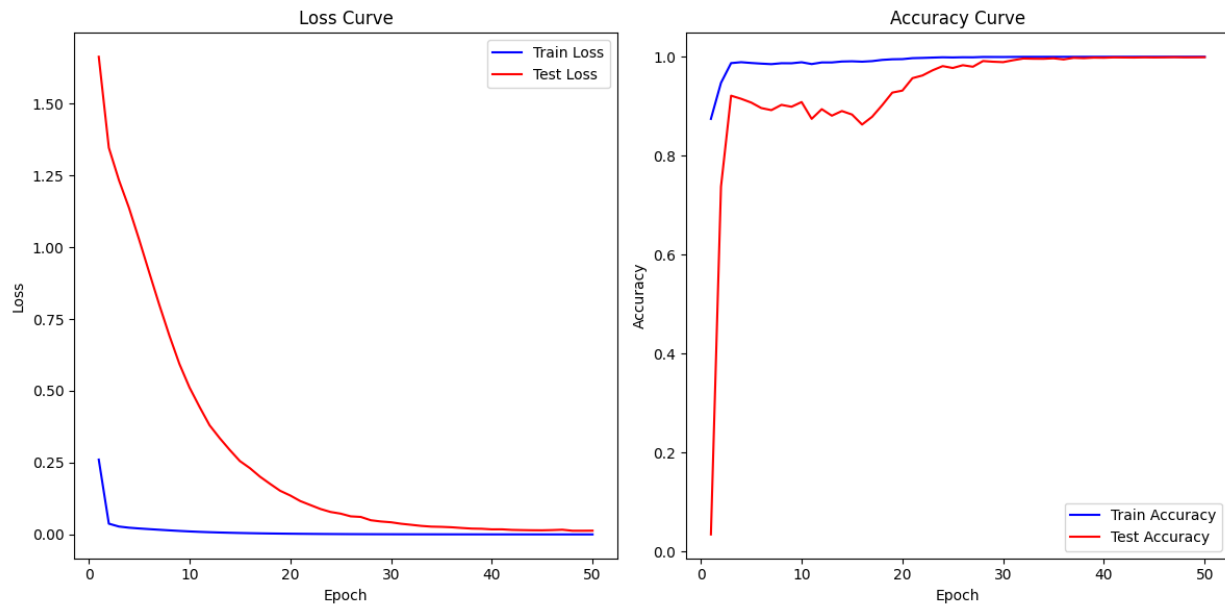
Figure 11: The training curve of Our Model

2. Data Preprocessing for Improved Accuracy: Due to the inherent noise and inconsistencies in the data, we have performed extensive data cleaning and preprocessing. This step significantly enhances the quality of the input data, leading to improved model accuracy. By handling missing values, outliers, and inconsistencies in the data, we ensure that our model is trained on high-quality information, which in turn boosts its predictive performance.

3. Robustness to New Entrants and Coaching Effects: Our model exhibits robustness, maintaining stable predictions even in the face of yearly new entrants to the Olympics and the presence of the "great coach effect." Every year, many athletes participate for the first time, and the influence of great coaches can have a significant impact on an athlete's performance. Despite these variables, our model continues to deliver consistent and reliable results, ensuring its applicability across different Olympic cycles and its resilience to such dynamic changes in the sports landscape.

### 6.1.2 Limitations

Despite the strengths of our model, there are several limitations that should be considered:

1. Dependency on Prior Knowledge for Clustering: When performing clustering, our model requires prior knowledge to guide the process, particularly to mitigate the challenges posed by the long-tail distribution of the data. This guidance is crucial to ensure that the model can effectively identify and separate different clusters. However, the need for such prior knowledge increases the reliance on human expertise, leading to higher costs for manual labeling and annotation. In real-world applications where domain knowledge may be scarce or expensive to obtain, this could present a significant barrier to scaling the model.

2. Performance Degradation in Low-Quality Data Scenarios: The performance of our model can degrade in scenarios where the quality of the data is low. In particular, if the input data
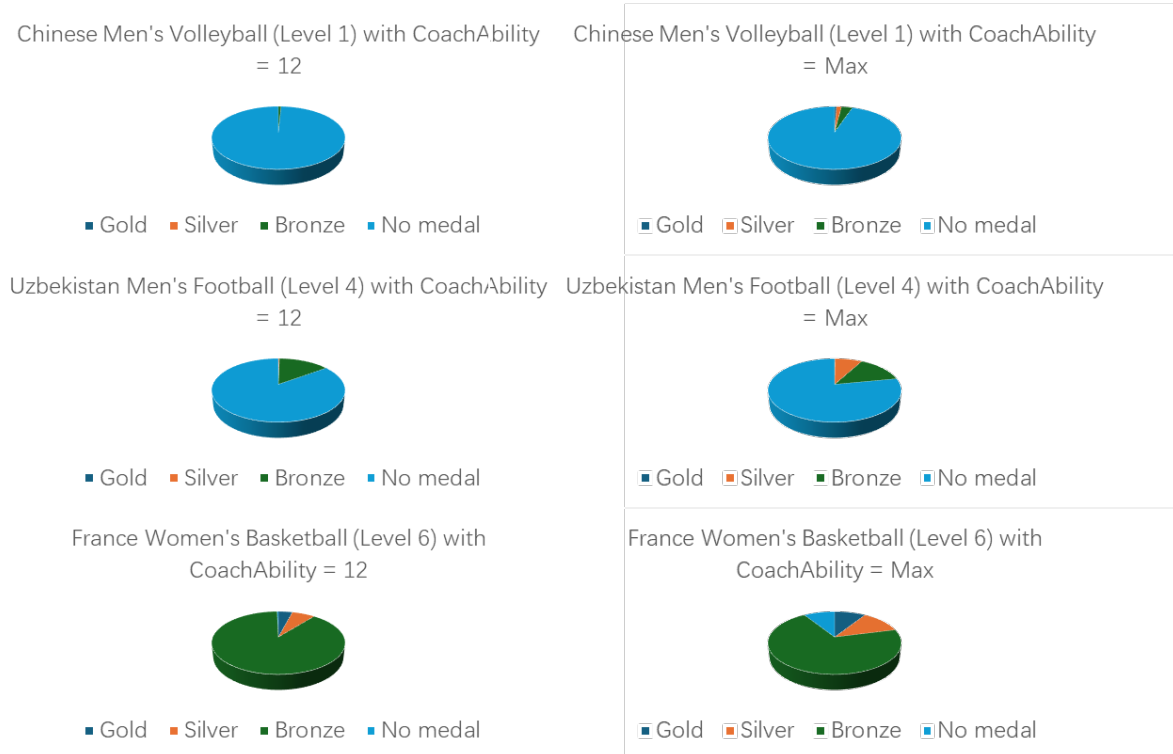
Figure 12: The influence of great coaches on different sports in different countries

contains noise, inconsistencies, or missing values, the model's ability to learn meaningful patterns is compromised. While we have implemented data preprocessing steps to mitigate this issue, the model's overall predictive accuracy may still suffer in cases where the data quality is particularly poor. In such situations, additional data cleaning or more sophisticated noise-handling techniques may be required to restore model performance.

## 6.2   Future Work

Our model is highly versatile and can be extended to a variety of other scenarios beyond its current application. Specifically, it is well-suited for tasks that involve higher complexity, particularly those with a larger number of features or higher-dimensional data. As the dimensionality of the input data increases, the ability of our model to learn intricate patterns and capture complex relationships also improves.

In high-dimensional problems, traditional models often struggle to learn meaningful representations due to the curse of dimensionality. However, our model is capable of effectively processing such data by leveraging its sophisticated structure and advanced learning mechanisms. This allows it to extract more nuanced knowledge from the data, resulting in better generalization and more accurate predictions, even in the presence of a large number of features.

Furthermore, the flexibility of our model makes it applicable to a wide range of domains, such as finance, healthcare, and social networks, where the feature space is often large and data is rich with complex dependencies. By adapting the model to these new contexts, we can continue to

improve its performance and achieve reliable results in scenarios with greater complexity.

# 7    Sensitivity Analysis

In Section 5, we discussed the effectiveness of clustering methods. Here, we discuss the robustness of our model to the 'Great Coach Effect.' We randomly jittered the coach's ranking by 1-3 levels. Due to space limitations, we only analyzed the top four countries. The results are shown in Figures 11 and 12. As seen, our model maintains a certain level of stability even when there is fluctuation in the coach's ranking.
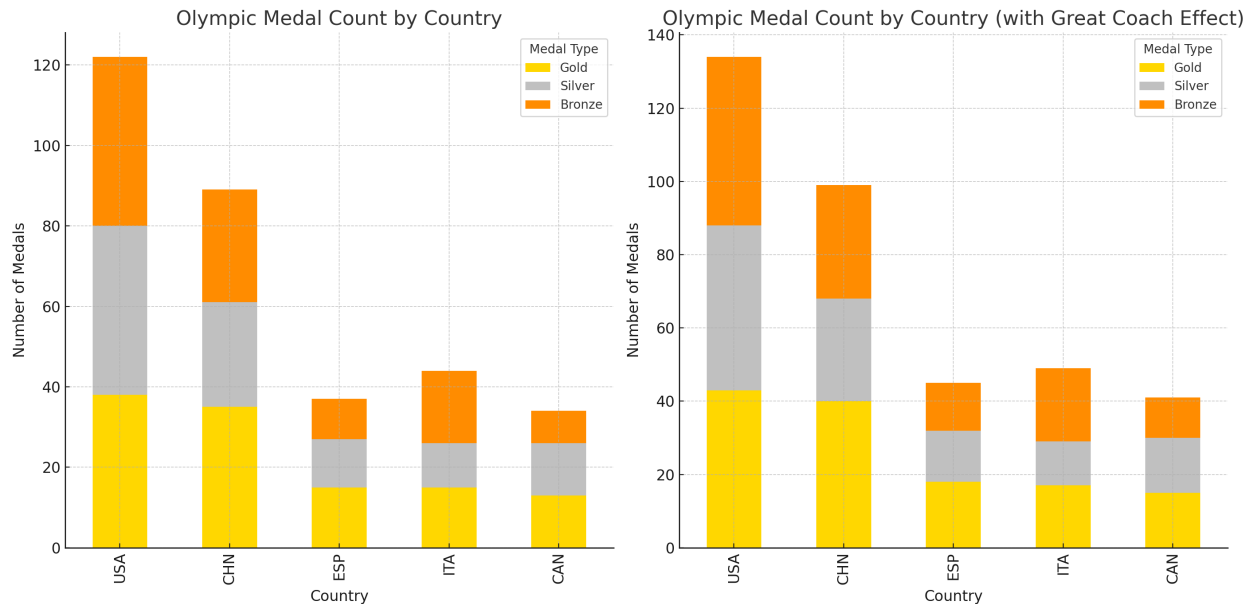


Figure 13: Sensitivity Analysis

# 8    Conclusions

In this study, we developed a comprehensive model for predicting Olympic medal outcomes, specifically targeting the 2028 Los Angeles Olympics. Our approach combines advanced data pre-processing techniques, including cleaning and transforming historical Olympic data, with powerful machine learning algorithms such as clustering and Long Short-Term Memory (LSTM) networks. A key innovation in our methodology is the Semi-supervised Hierarchical Mean-shift Clustering (SHMC) approach, which allows us to better understand the complex distributions of Olympic performance data by grouping countries based on their historical success and other relevant factors. This clustering method is crucial in capturing the long-tail nature of the Olympic medal distribution and ensuring that the model is capable of accurately predicting outcomes even for countries with less frequent medal wins.

By incorporating various factors such as historical performance, event types, and the influence of coaching, our model offers robust and accurate predictions for the medal outcomes of countries. The results also highlight the importance of both traditional sports powerhouses and emerging nations in shaping the global medal table. Through sensitivity analysis, we demonstrated that the

model remains stable even when factors such as the "Great Coach Effect" fluctuate, emphasizing the robustness of the model in adapting to changes in the sports landscape.

While our model performs well in predicting Olympic medal outcomes, it is not without limitations. The reliance on prior knowledge for clustering, particularly the SHMC method, can be a challenge in real-world applications. Additionally, data quality issues such as noise and inconsistencies may affect the model's performance. Despite these challenges, we believe that the versatility and performance of the model can be further refined, and it holds potential for broader applications in other domains dealing with complex, multi-dimensional data.

In conclusion, this research provides a valuable framework for understanding the factors that drive Olympic success. Our model offers practical insights for national sports committees and organizations aiming to optimize their preparation strategies, and can serve as a tool for predicting medal outcomes. Future work could explore integrating more real-time data, enhancing clustering methods, and expanding the model to encompass additional sports and performance variables to further refine predictions.

# References

[1] Kuper, S. (2023). The Olympic Medal Table: What It Really Tells Us. Foreign Affairs.

[2] Baker, J., & Horton, S. (2022). Olympic Medal Distribution: A Global Analysis of Sports Development.

[3] Clarke, S. R. (2021). Home Advantage in the Olympic Games. Journal of Sports Sciences.

[4] Johnson, D. K. (2023). Host Nation Performance in Olympic Games: A Comprehensive Analysis.

[5] Zhang, L., & Liu, Y. (2024). The Impact of Elite Coaches on Olympic Success: A Case Study Approach.

[6] Hofmann, A. R., & Gerlach, E. (2019). "Olympic boycotts and international relations: A historical perspective (1976-2000)." Sport History Review, 50(1), 67-87.

[7] Maguire, J. (2002). Globalization, sport, and national identities: "The Empire strikes back"? International Review for the Sociology of Sport, 37(1), 35-47.

[8] Riordan, J., & Krüger, A. (2021). "The Post-Soviet Era in Olympic Sport: The Impact of Political Transformation on International Sports Participation." International Journal of the History of Sport, 38(4), 382-396.

[9] De Bosscher, V., Shibli, S., & Weber, A. C. (2019). "Is prioritisation of funding in elite sport effective? An analysis of the investment strategies in 16 nations." European Sport Management Quarterly, 19(2), 221-243.

[10] Atten, C. (2023). The ultimate beginner guide of semi-supervised learning: Learn the basics of semi-supervised learning, how it works, and what its limitations are. DataDrivenInvestor.

# Report on Use of AI

OpenAI o1 (released on September 12, 2024, version: o1)
Query1: How to build an LSTM? Provide the code.
Output: Long Short-Term Memory (LSTM) networks are a type of recurrent neural network (RNN) designed to model sequential data by capturing both short-term and long-term dependencies. They are particularly effective for tasks such as time-series forecasting, natural language processing, and speech recognition. Here's a basic example of how to build an LSTM using Python with the Keras library:

```python
import numpy as np
from keras.models import Sequential
from keras.layers import LSTM, Dense

# Sample data: sequences of 10 timesteps with 1 feature
X_train = np.random.rand(100, 10, 1)  # 100 samples
y_train = np.random.rand(100, 1)      # 100 target values

# Initialize the model
model = Sequential()

# Add an LSTM layer with 50 units
model.add(LSTM(50, input_shape=(10, 1)))

# Add a Dense output layer
model.add(Dense(1))

# Compile the model
model.compile(optimizer='adam', loss='mean_squared_error')

# Train the model
model.fit(X_train, y_train, epochs=20, batch_size=32)
```

In this example: X_train is a NumPy array of shape (100, 10, 1), representing 100 sequences, each with 10 timesteps and 1 feature per timestep. y_train is a NumPy array of shape (100, 1), representing the target values for each sequence. The model consists of a single LSTM layer with 50 units, followed by a Dense layer with 1 unit to produce the output. The model is compiled with the Adam optimizer and mean squared error loss function, suitable for regression tasks. The model is trained for 20 epochs with a batch size of 32.