```
1    import pandas as pd
2    import numpy as np
3    import matplotlib.pyplot as plt
4    import seaborn as sns
5    from sklearn.model_selection import train_test_split
6    from sklearn.linear_model import LinearRegression
7    from sklearn.ensemble import RandomForestRegressor
8    from sklearn.metrics import mean_squared_error
9
10   file_path = '/content/Large language models (2024).csv'
11   df = pd.read_csv(file_path, encoding='latin1')
12
13   df.columns = [col.strip().replace(' ', '_').lower() for col in df.columns]
14
15   df['parameters'] = pd.to_numeric(df['parameters'], errors='coerce')
16   df['tokens'] = pd.to_numeric(df['tokens'], errors='coerce')
17   df['alscore'] = (df['parameters'] * df['tokens']) ** 0.5
18
19   df.head()
```

| | model | comapany | arch | parameters | tokens | ratio | alscore | training_dataset | release_date | notes |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Olympus | Amazon | TBA | 2000.0 | 40000.0 | 20:01 | 8944.27191 | TBA | TBA | New related Titan details: '$65m training run.... |
| 1 | GPT-5 | OpenAI | TBA | 2000.0 | NaN | TBA | NaN | TBA | TBA | Due 2024. |
| 2 | GPT-6 | OpenAI | TBA | NaN | NaN | TBA | NaN | TBA | TBA | Due 2025. |
| 3 | AuroraGPT (ScienceGPT) | ANL | TBA | 1000.0 | NaN | TBA | NaN | TBA | TBA | https://tpc.dev/2023/11/10/tpc-announced-with-... h |
| 4 | Grok-2 | xAI | TBA | NaN | NaN | TBA | NaN | TBA | TBA | Due 2025. h |

Next steps:  [ Generate code with `df` ]   [ ⬤ View recommended plots ]

```
 1 data = df[['parameters', 'tokens']].dropna()
 2
 3 X = data[['parameters']]
 4 y = data['tokens']
 5
 6 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
 7
 8 # RandomForestRegressor model
 9 model = RandomForestRegressor(n_estimators=100, random_state=42)
10 model.fit(X_train, y_train)
11
12 y_pred = model.predict(X_test)
13 rmse = mean_squared_error(y_test, y_pred, squared=False)
14
15 rmse
```

6460.983122100423

```
 1 data = df[['parameters', 'tokens', 'alscore']].dropna()
 2
 3 X = data[['parameters', 'tokens']]
 4 y = data['alscore']
 5
 6 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
 7
 8 # LinearRegression model
 9 model = LinearRegression()
10 model.fit(X_train, y_train)
11
12 y_pred = model.predict(X_test)
13 rmse = mean_squared_error(y_test, y_pred, squared=False)
14
15 rmse
```
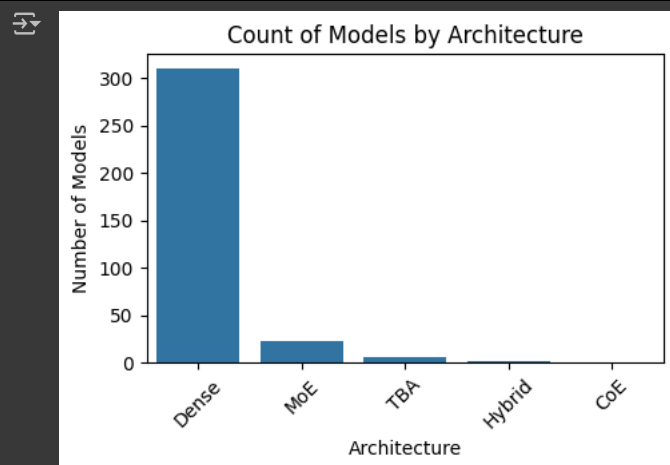
344.0360842023981

```
 1 # Count models by architecture
 2 arch_count = df['arch'].value_counts()
```

```
 1
 2
 3 # Plotting model count by architecture
 4 plt.figure(figsize=(5, 3))
 5 sns.barplot(x=arch_count.index, y=arch_count.values)
 6 plt.title('Count of Models by Architecture')
 7 plt.xlabel('Architecture')
 8 plt.ylabel('Number of Models')
 9 plt.xticks(rotation=45)
10 plt.show()
```
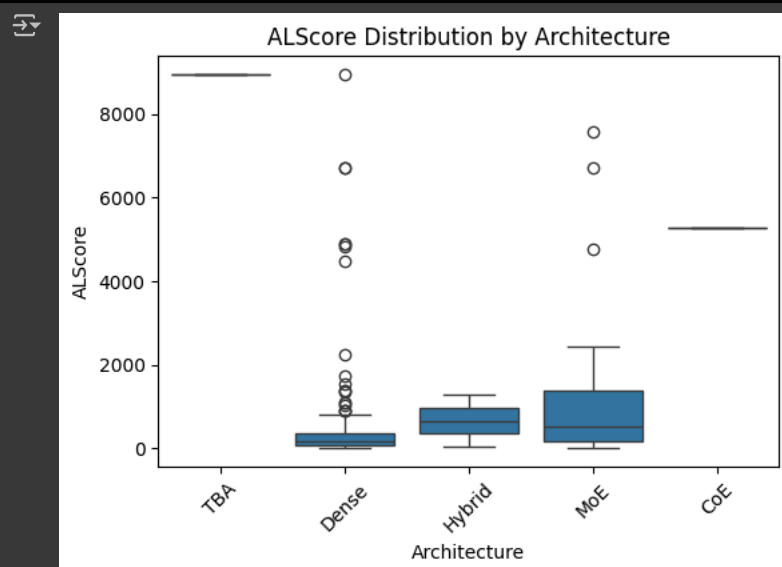


```
1 # Boxplot of ALScore by architecture
2 plt.figure(figsize=(6, 4))
3 sns.boxplot(x='arch', y='alscore', data=df)
4 plt.title('ALScore Distribution by Architecture')
5 plt.xlabel('Architecture')
6 plt.ylabel('ALScore')
7 plt.xticks(rotation=45)
8 plt.show()
```
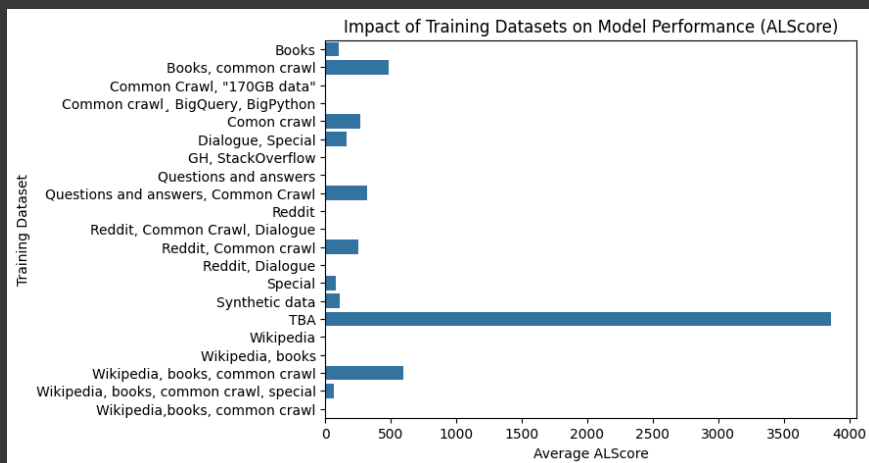


```
 1 # Analyzing the impact of training datasets
 2 # Grouping by training dataset and calculate mean ALScore
 3 dataset_impact = df.groupby('training_dataset')['alscore'].mean().reset_index()
 4
 5 # Plotting the impact of training datasets
 6 plt.figure(figsize=(7, 5))
 7 sns.barplot(x='alscore', y='training_dataset', data=dataset_impact)
 8 plt.title('Impact of Training Datasets on Model Performance (ALScore)')
 9 plt.xlabel('Average ALScore')
10 plt.ylabel('Training Dataset')
11 plt.show()
```

Impact of Training Datasets on Model Performance (ALScore)

```
# Scatter plot of parameters vs tokens
plt.figure(figsize=(4, 3))
sns.scatterplot(x='parameters', y='tokens', data=df)
plt.title('Relationship between Parameters and Tokens')
plt.xlabel('Parameters (Billions)')
plt.ylabel('Tokens (Billions)')
plt.show()
```



Relationship between Parameters and Tokens