# EDA

For VQA TextMining's Project

# Data info
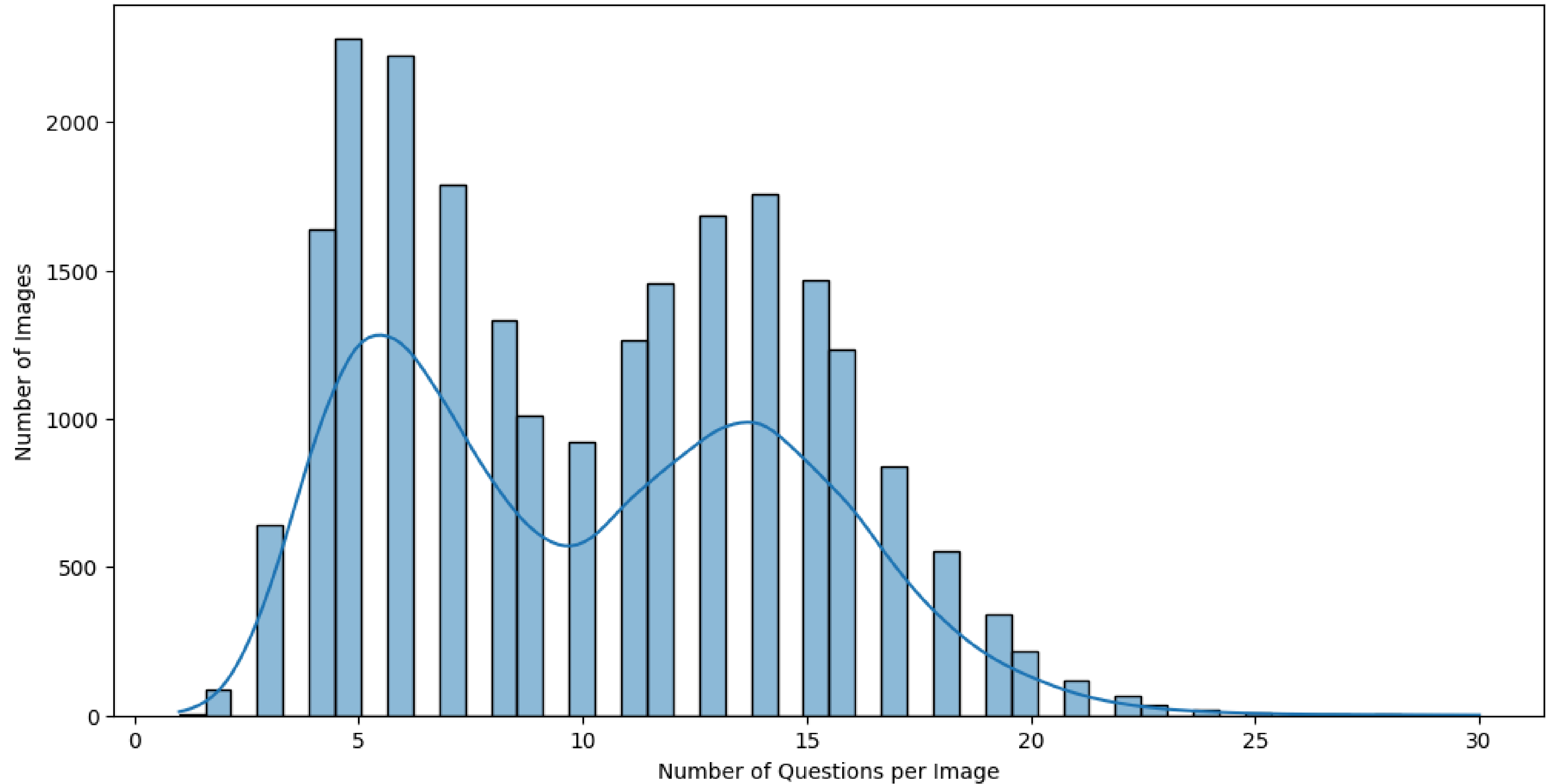
Images: 23008
QA-pairs: 235421
Questions: 57300
Answers: 108

Distribution of Questions per Image

# Top 10 question

**01**

Top 10 Questions Leading to Yes/No Answers:

Are there any person in image? 12300

Are there any cup in image? 2152

Are there any bowl in image? 2080

Are there any vase in image? 1566

Are there any car in image? 1452

Are there any wine glass in image? 1418

Are there any dining table in image? 1371

Are there any dog in image? 1354

Are there any bird in image? 1338

Are there any bottle in image? 1331

**02**

Top 10 Questions Leading to Numeric Answers:

How many person in image? 12207

How many cup in image? 1950

How many bowl in image? 1879

How many vase in image? 1368

How many car in image? 1284

How many dog in image? 1182

How many wine glass in image? 1175

How many dining table in image? 1151

How many bottle in image? 1129

How many bird in image? 1114

# Top 10 question

**01**

Top 10 Questions Leading to Color Answers:
What color is the sky? 562
What color is the background? 353
 What color are the leaves? 263
What color are the flowers? 248
What color is the shirt? 228
What color is the grass? 217
What is the color of the sky? 205
What color is the jacket? 177
 What color is the table? 167
What color is the water? 149

**02**

Top 10 Questions Leading to Position Answers:
Where is the vase? 73
 Where is the teddy bear? 65
Where is the fork? 51
Where is the clock? 44
Where is the book? 35
are the flowers? 34
Where is the knife? 34
Where is the glass? 33
Where is the plant? 31
Where is the wine glass? 31

# Top 10 Questions



question

Contribution of Different Answer Types

# Distribution of Yes/No Answers

yes

60.9%

39.1%

no

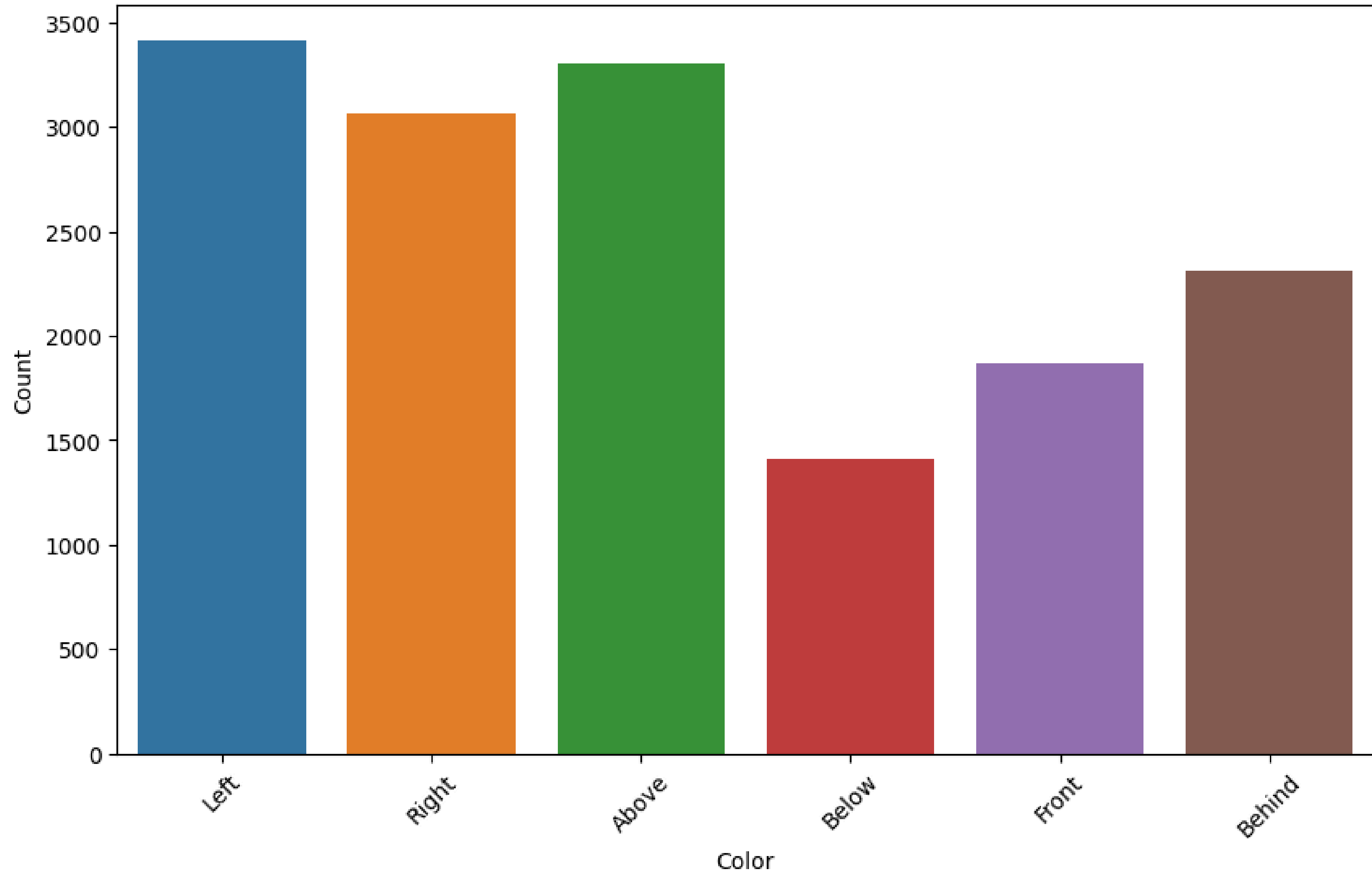Distribution of Object Answers

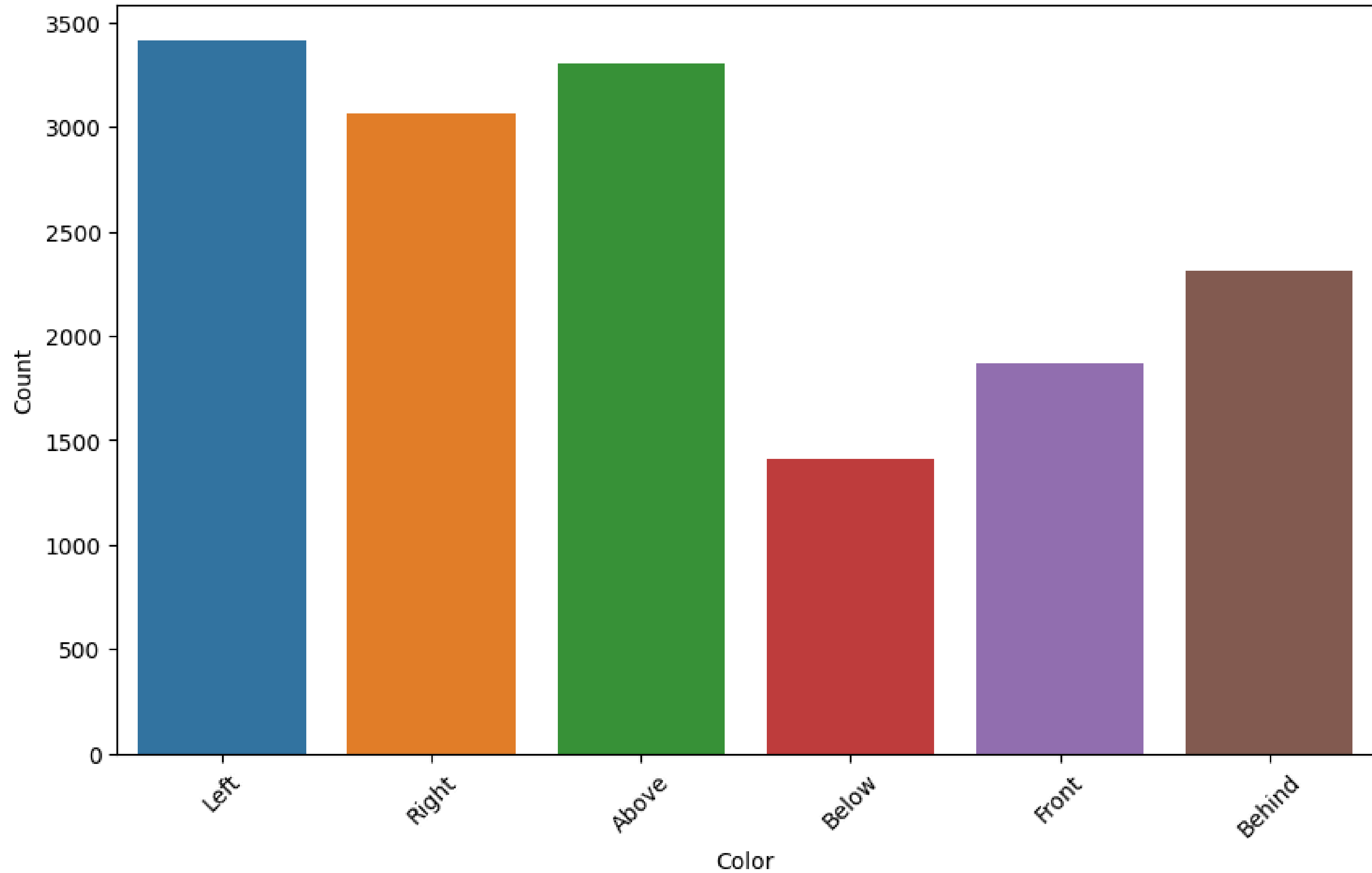Distribution of Numberic Answers

Distribution of Color Answers

Distribution of Position Answers

Distribution of Position Answers
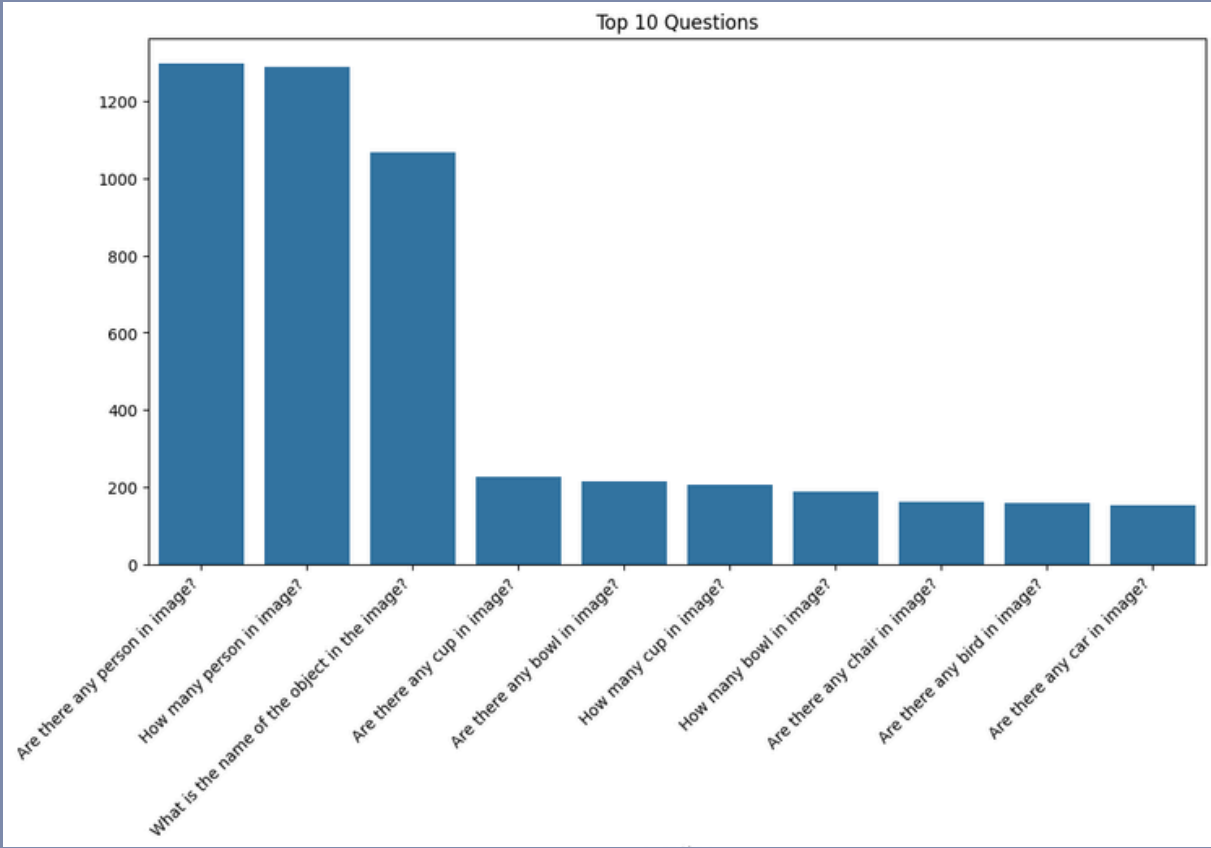
# Test info

# Eval info

Images: 2417
Questions: 7961
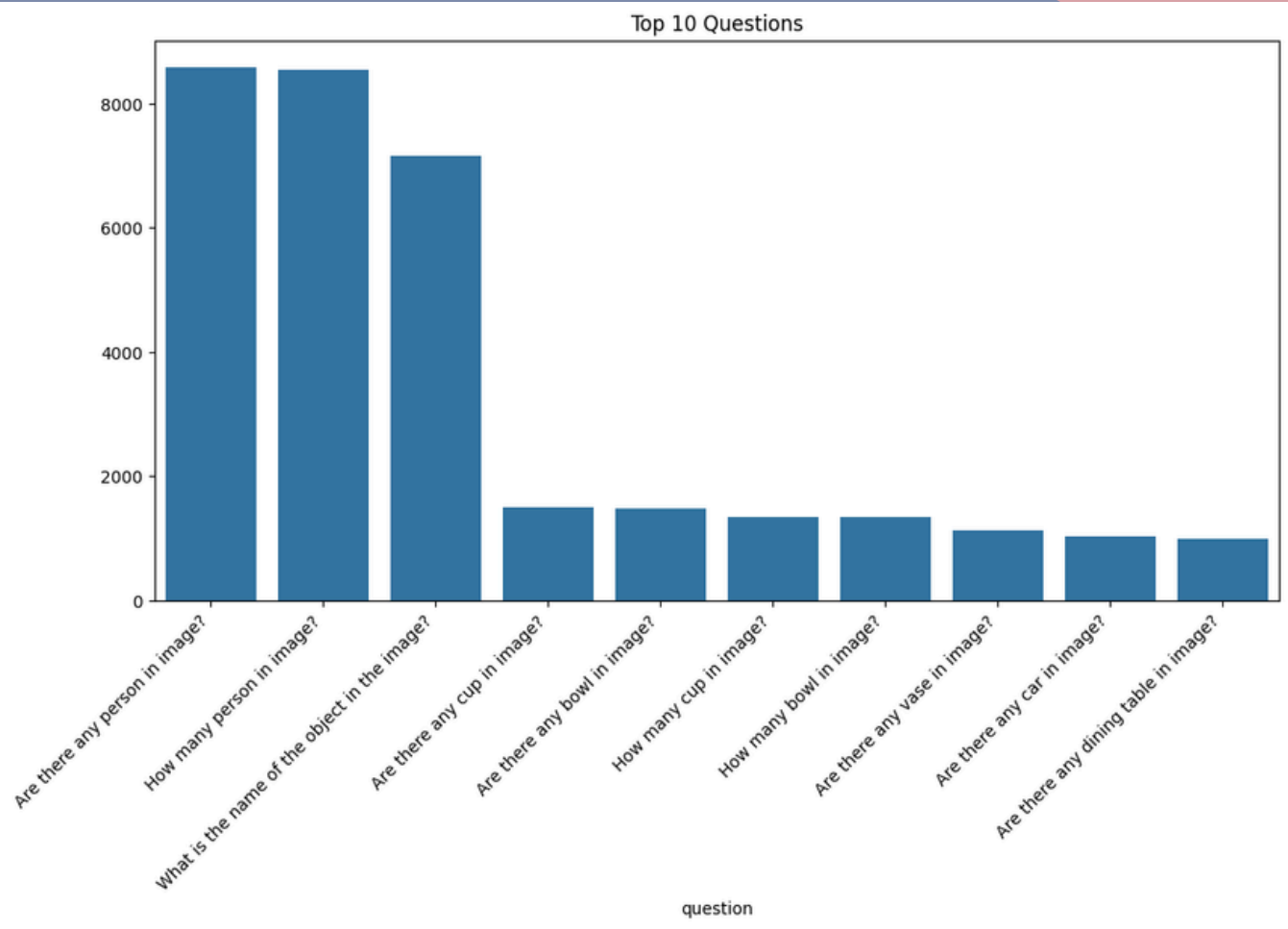Answers: 87

Images: 4486
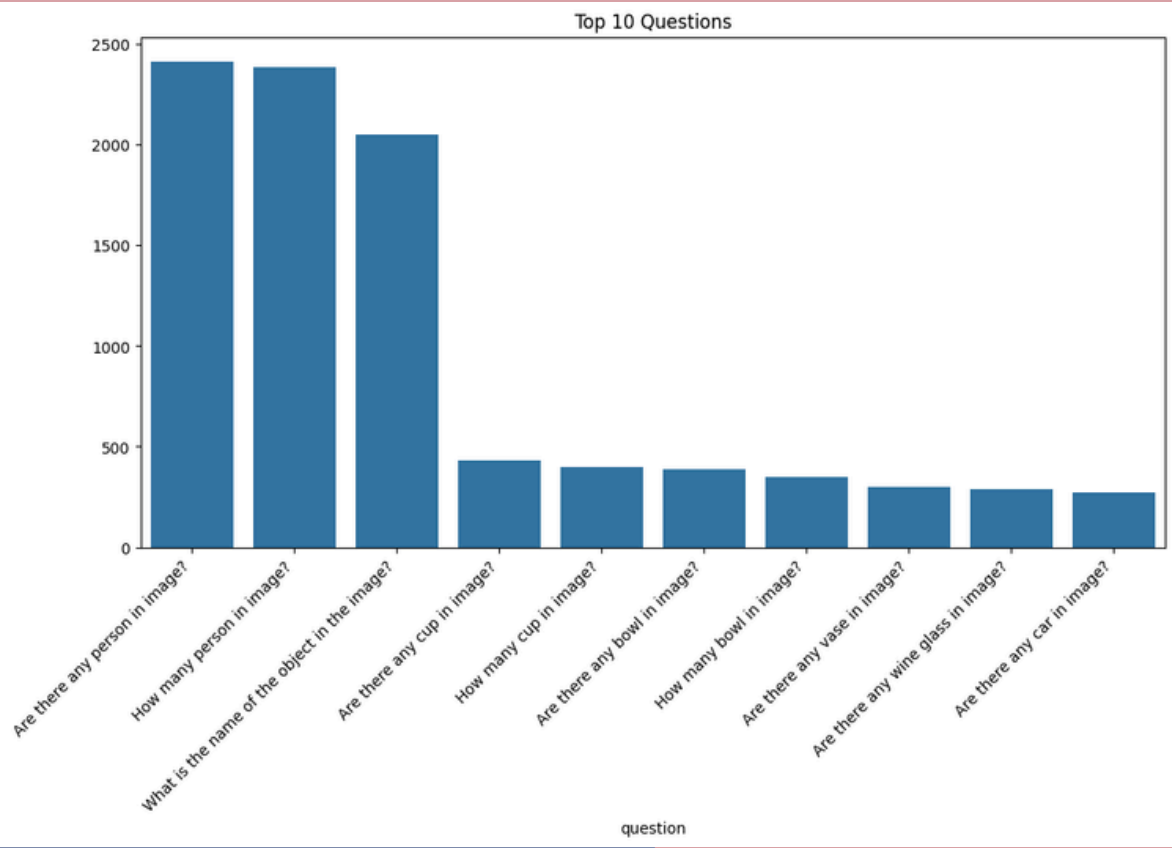Questions: 13460
Answers: 91
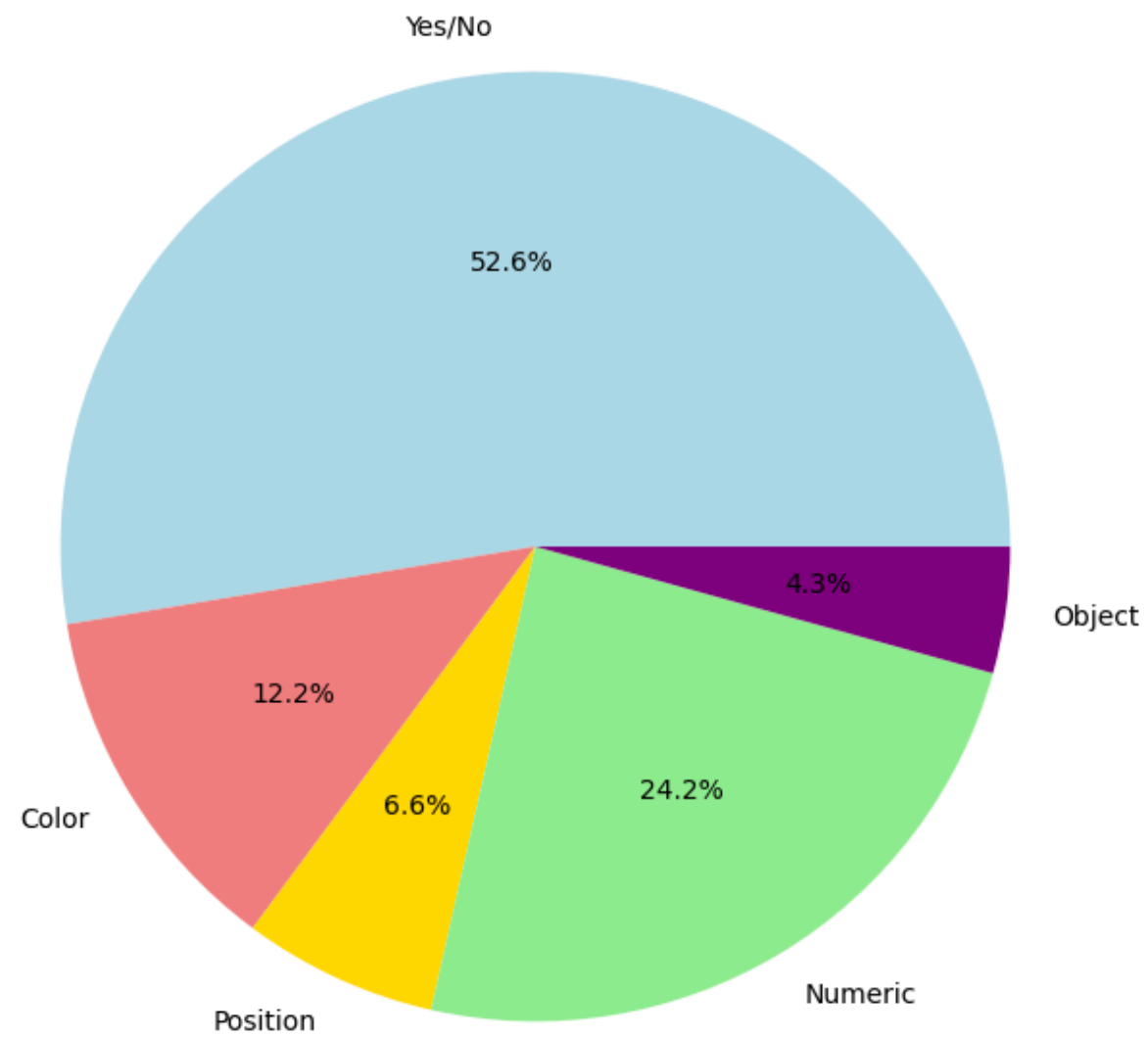
# Train info

Images: 16105
Questions: 42073
Answers: 105

Top 10 Questions

Test sets


Top 10 Questions

Train set


Top 10 Questions

Eval set

Contribution of Different Answer Types

**Test set**

Contribution of Different Answer Types

**Train set**

Contribution of Different Answer Types

**Eval set**

Distribution of Yes/No Answers

yes — 60.4%
no — 39.6%

Test set

Distribution of Yes/No Answers

yes — 60.9%
no — 39.1%

Train set

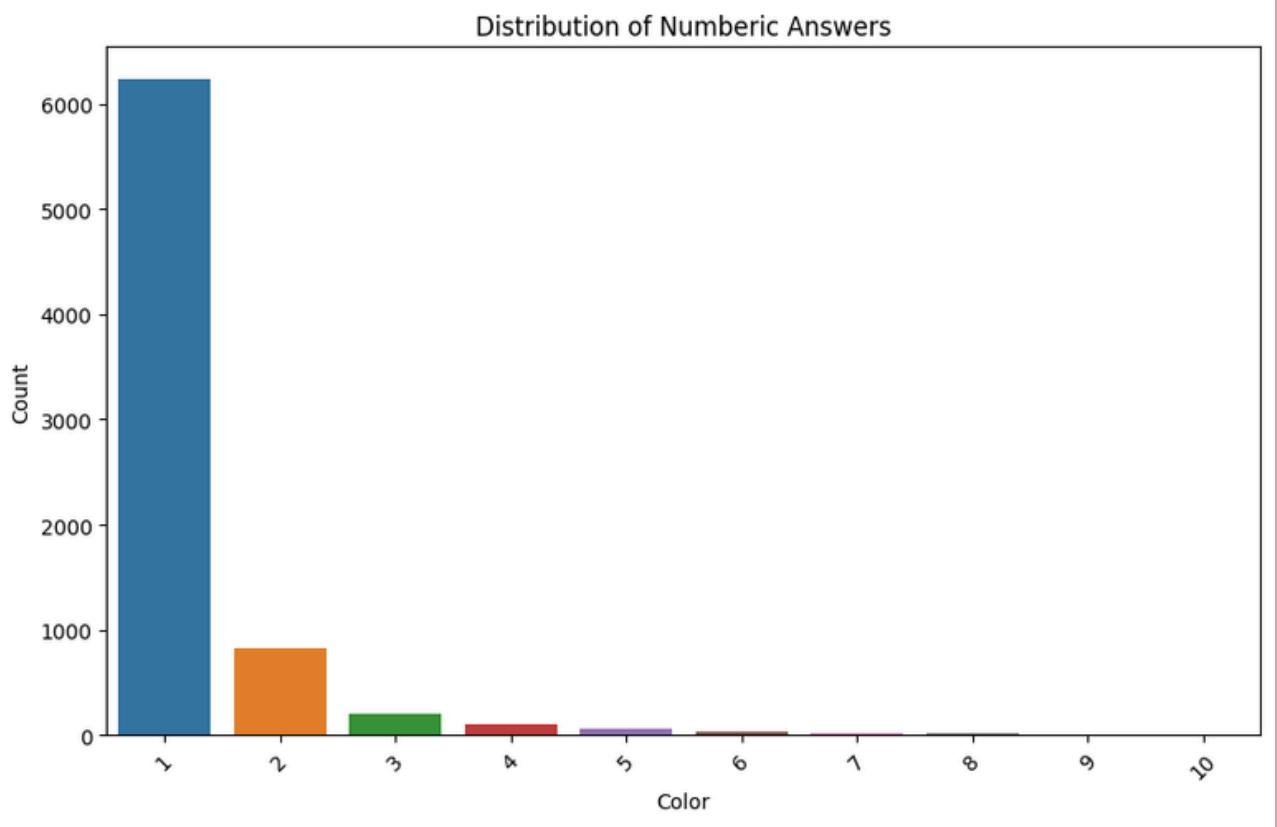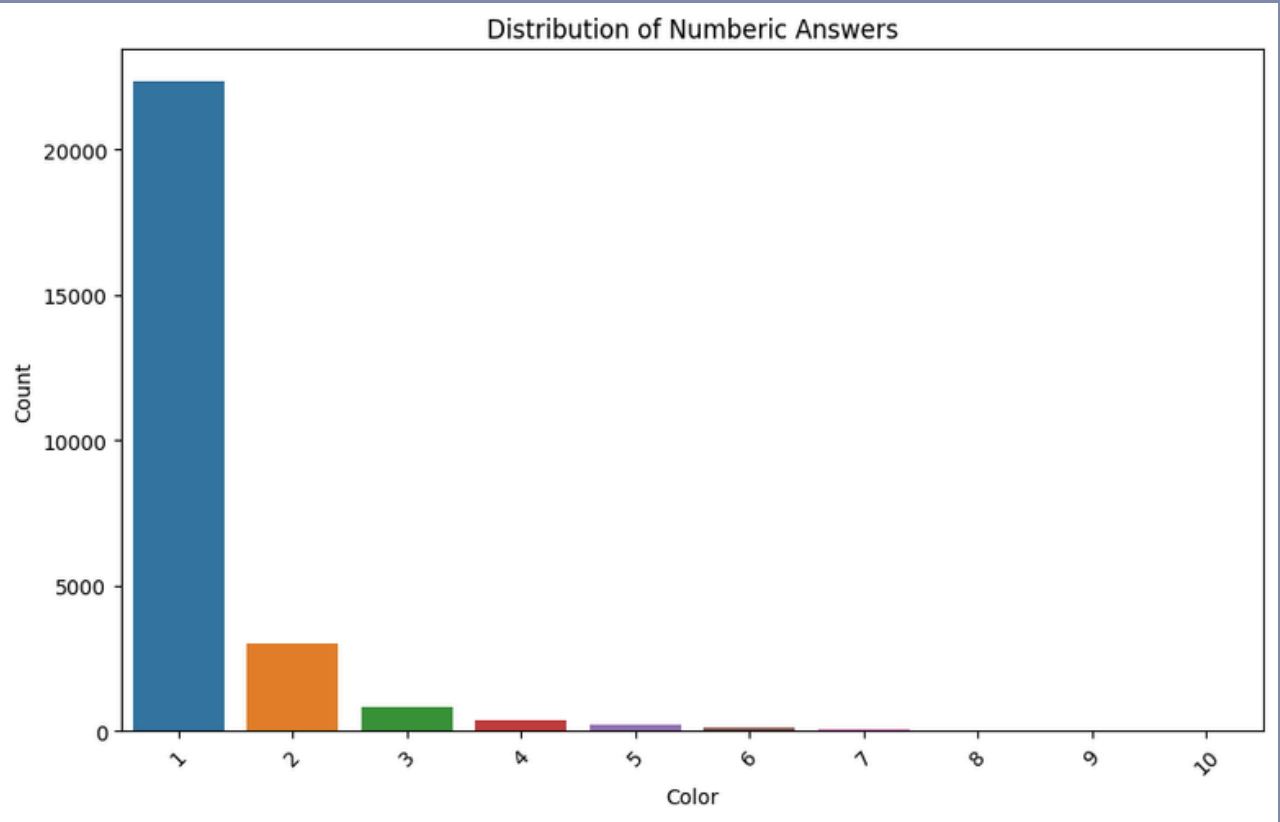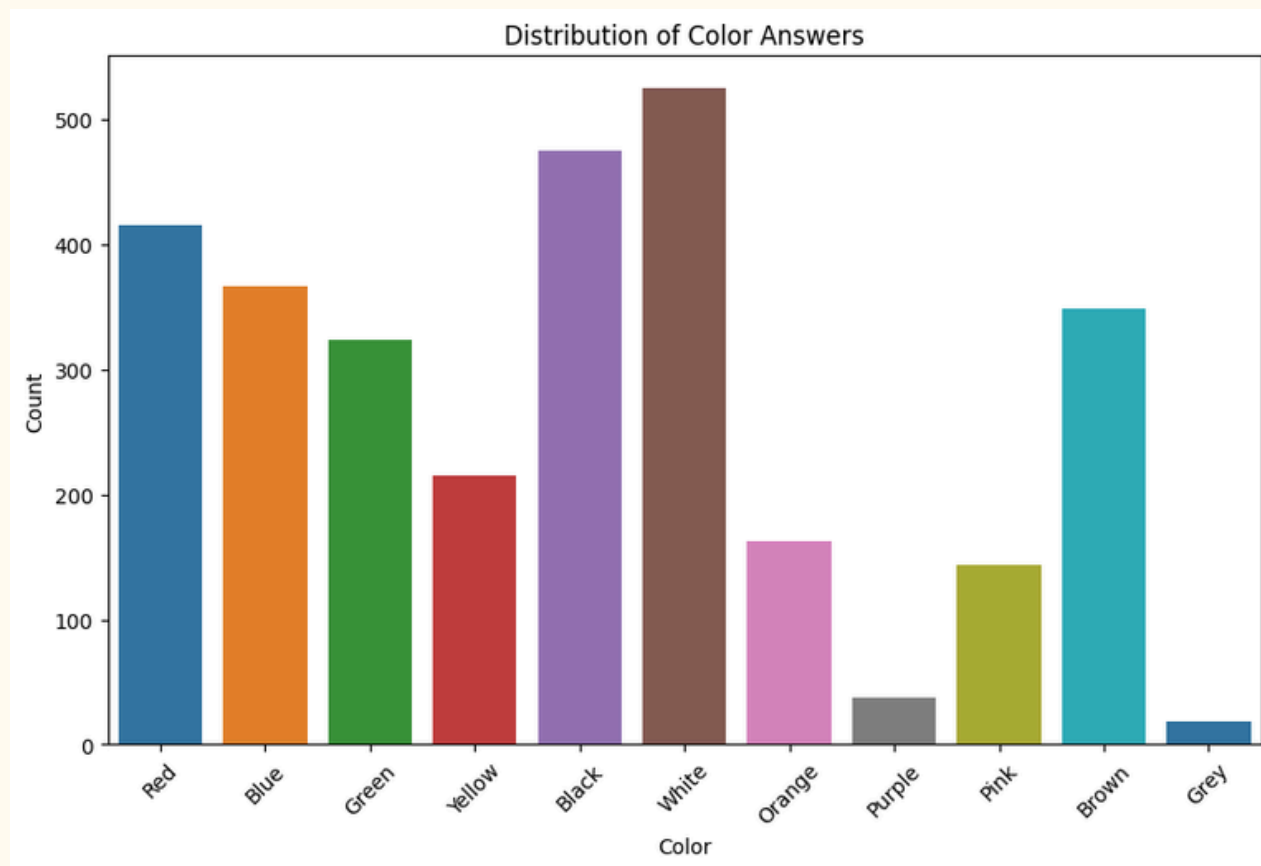Distribution of Yes/No Answers

yes — 61.1%
no — 38.9%
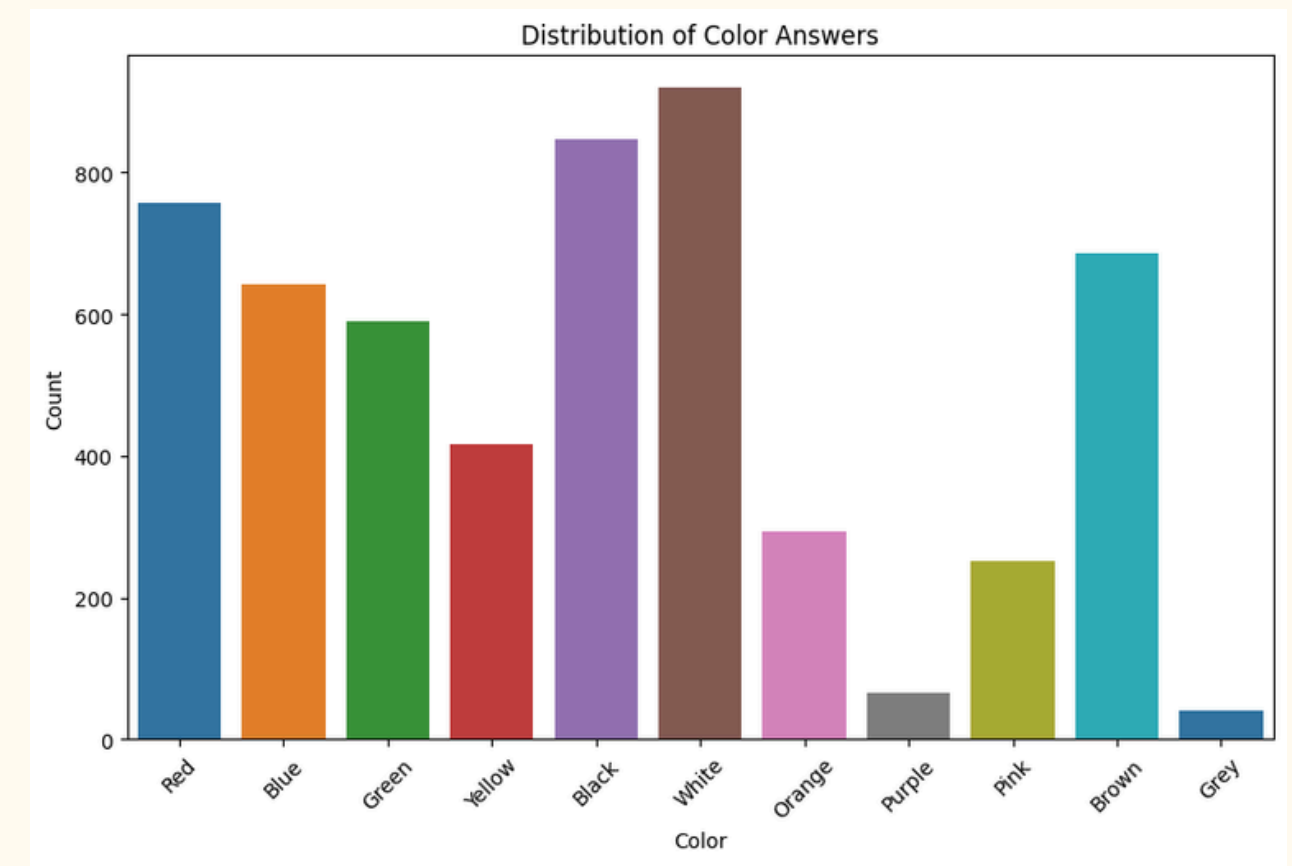
Eval set

Test set

Eval set

Train set

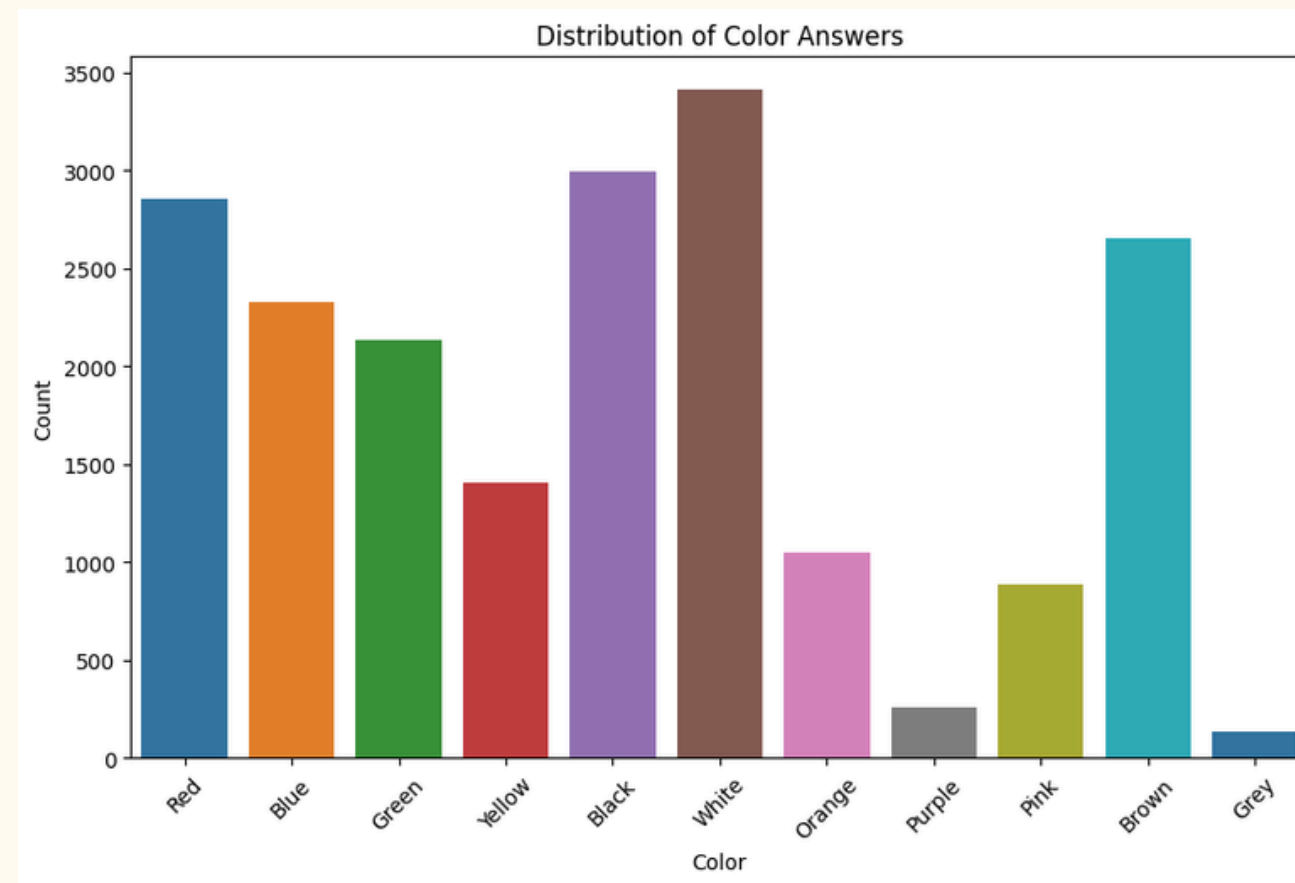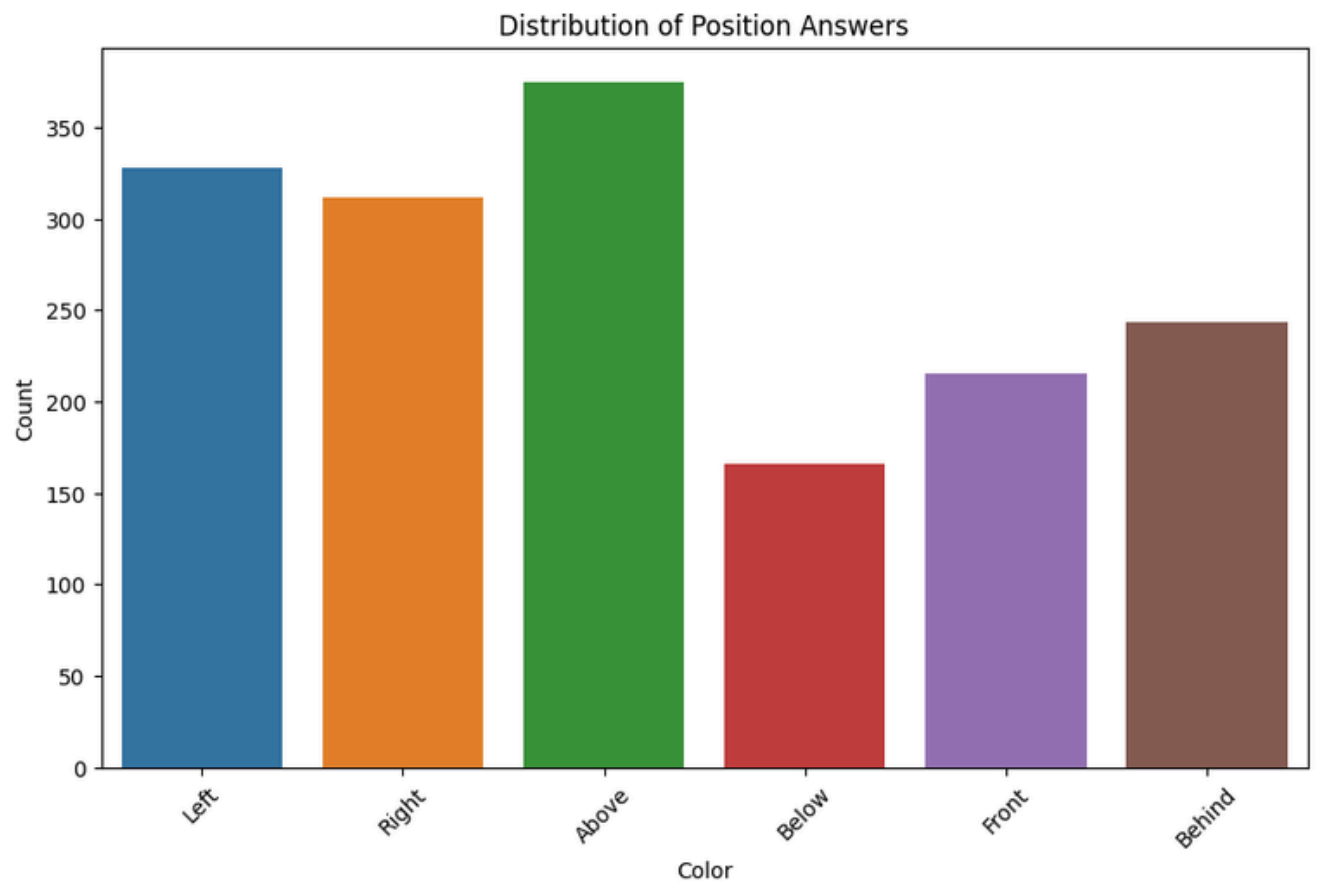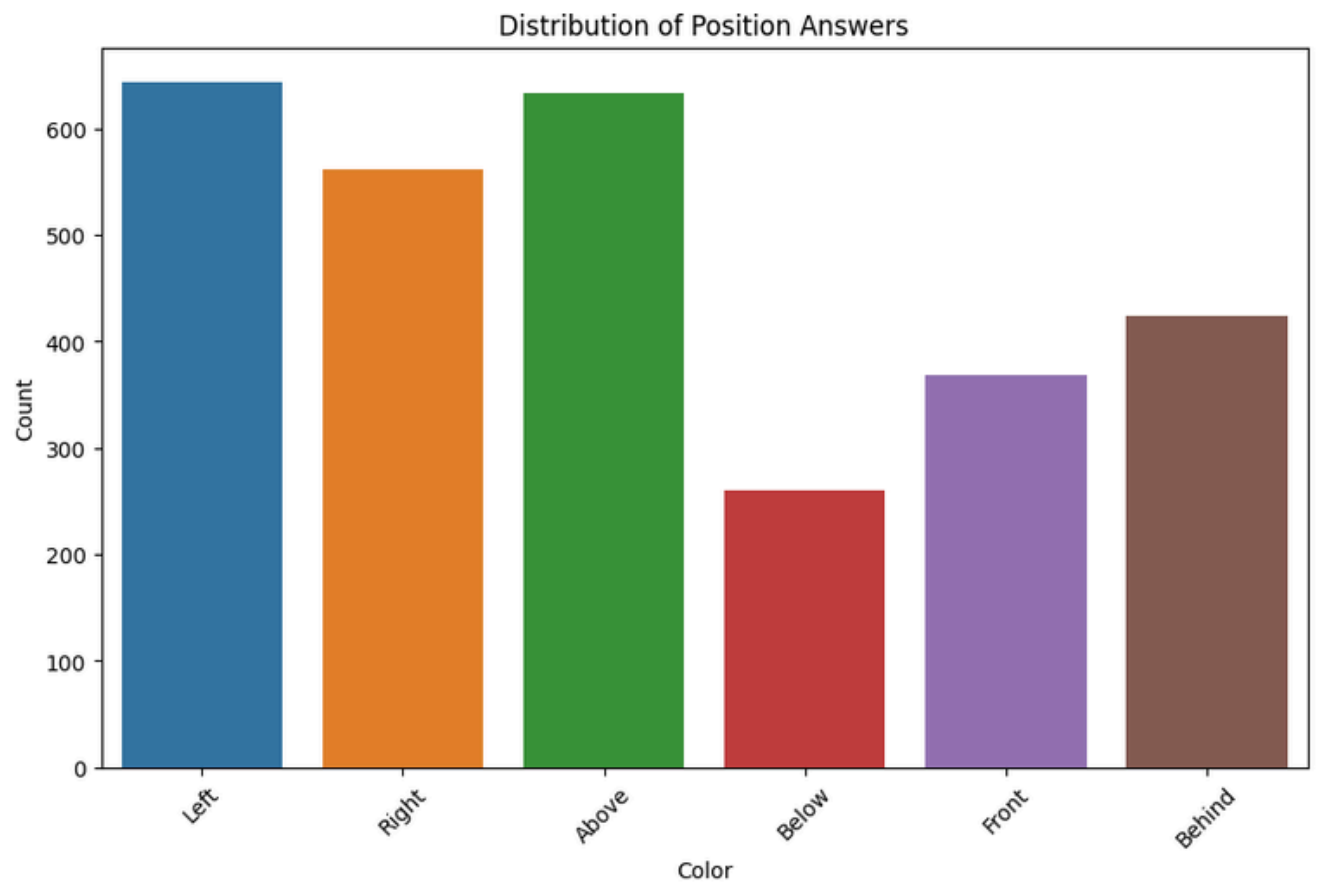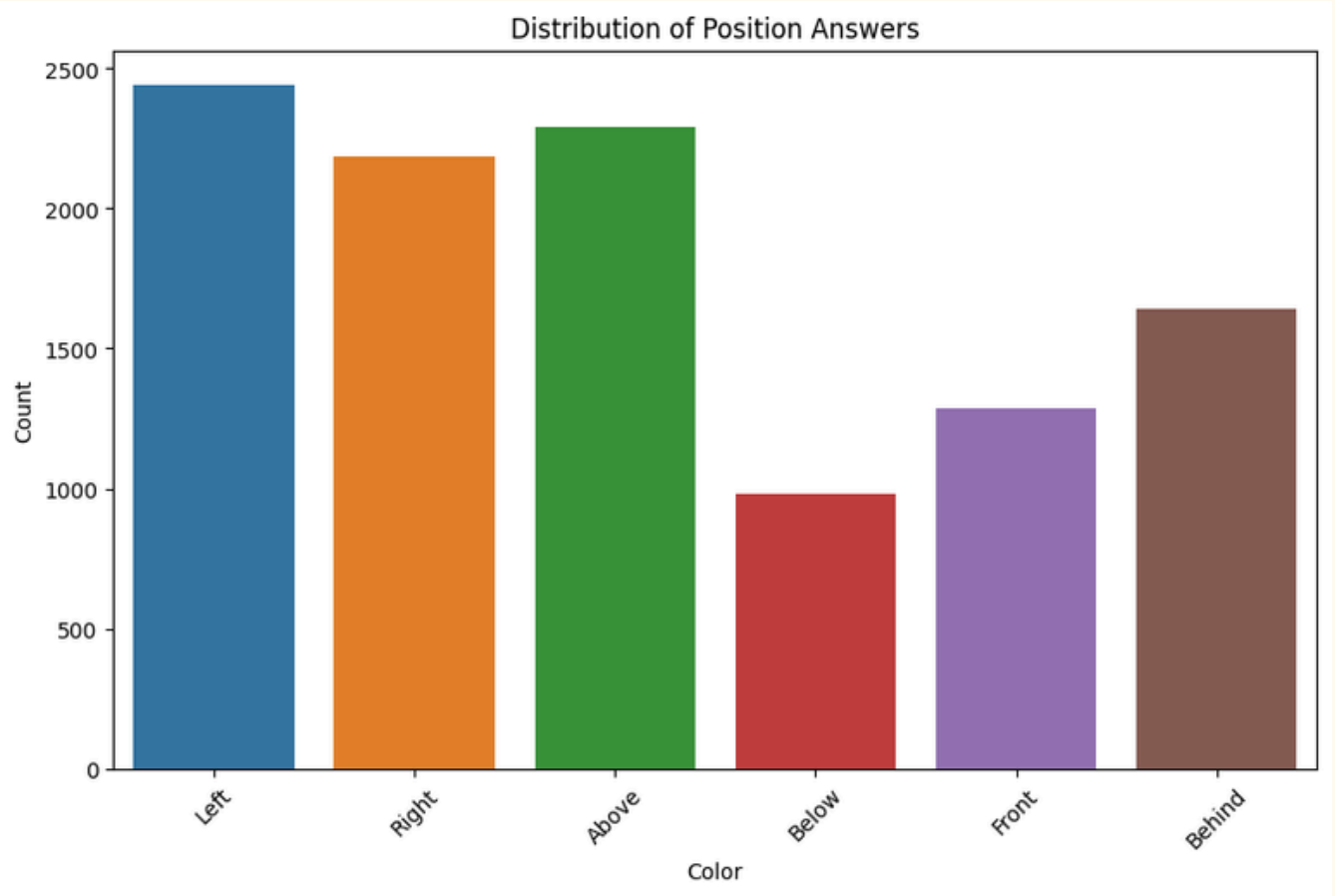Test set

Train set

Eval set

Test set

Eval set
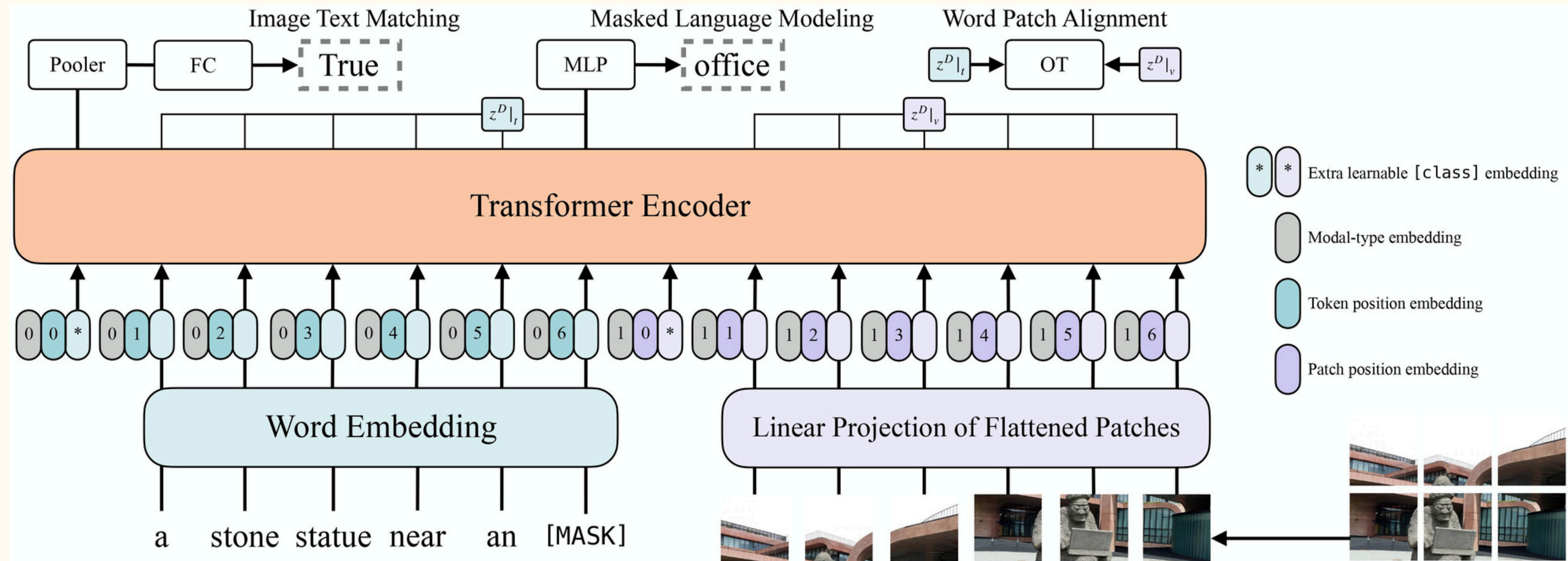
Train set

Test set

Eval set

Train set
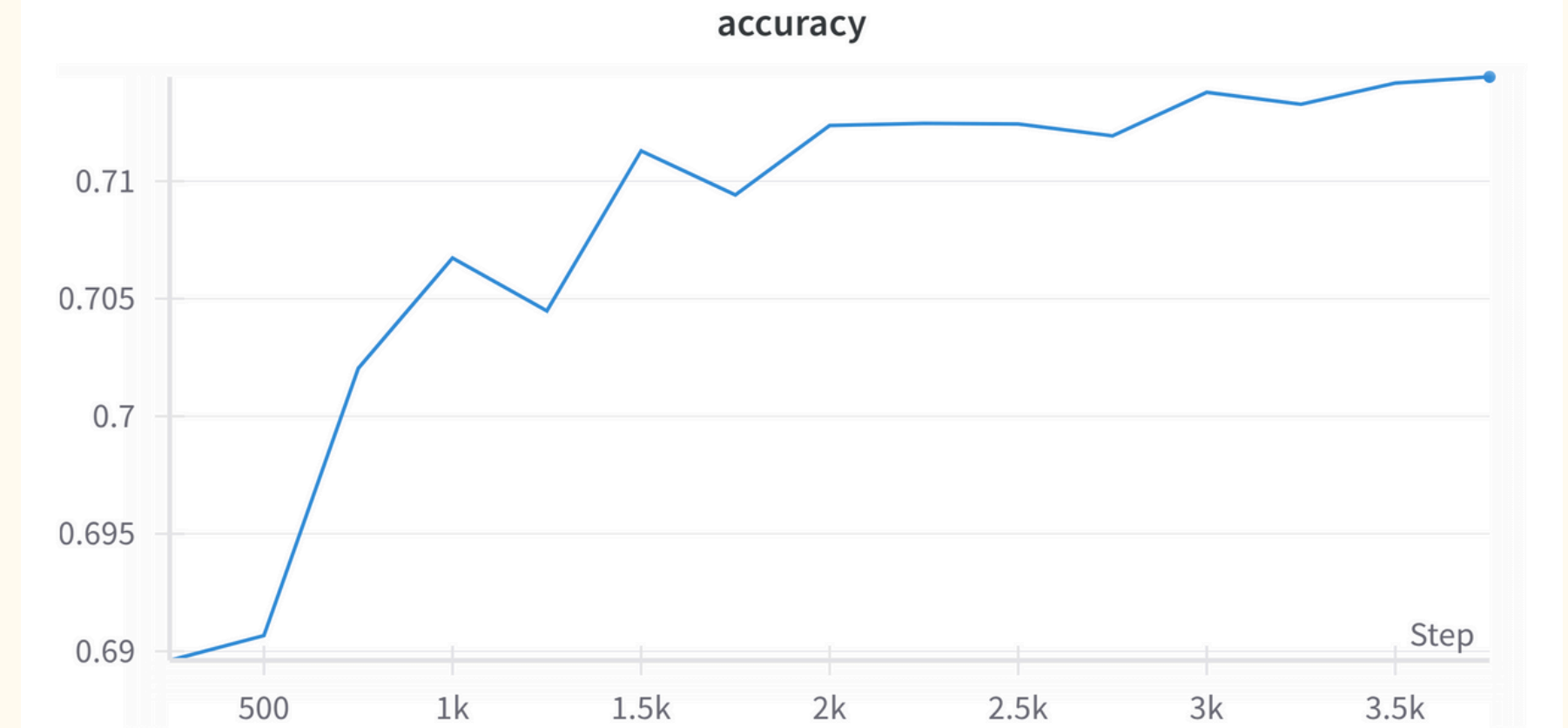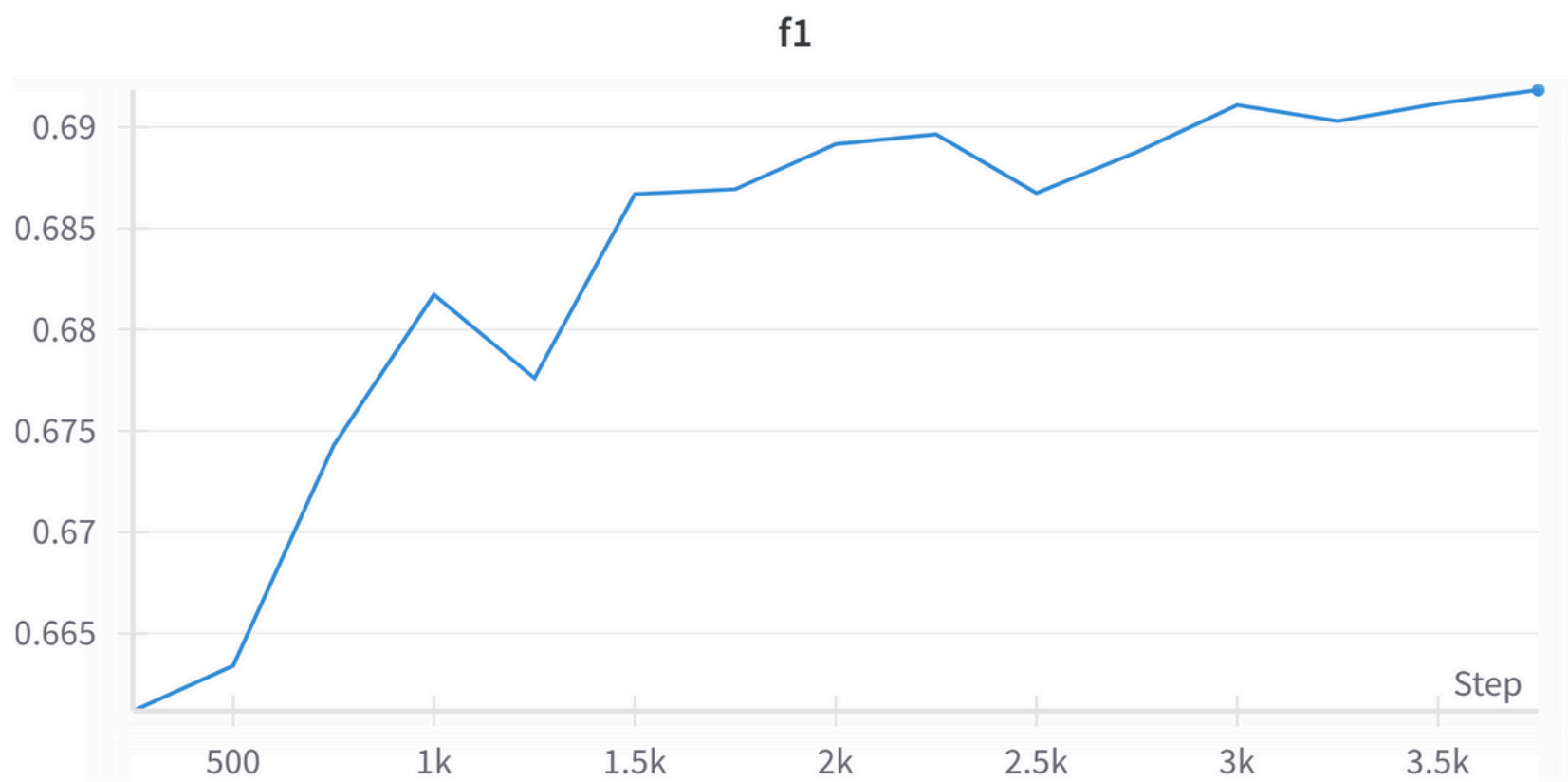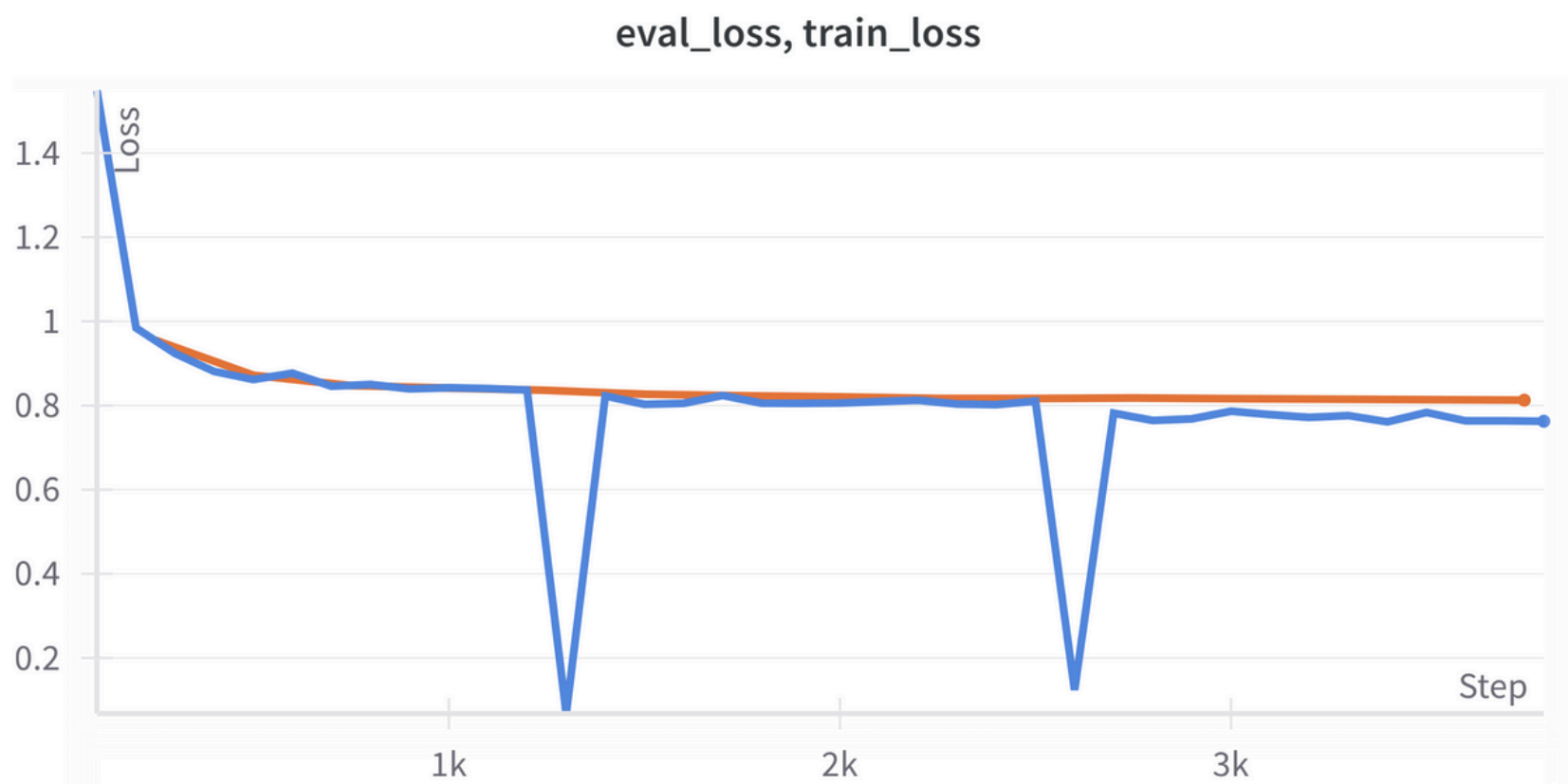
Test set

Eval set

Train set

# Models

ViLt
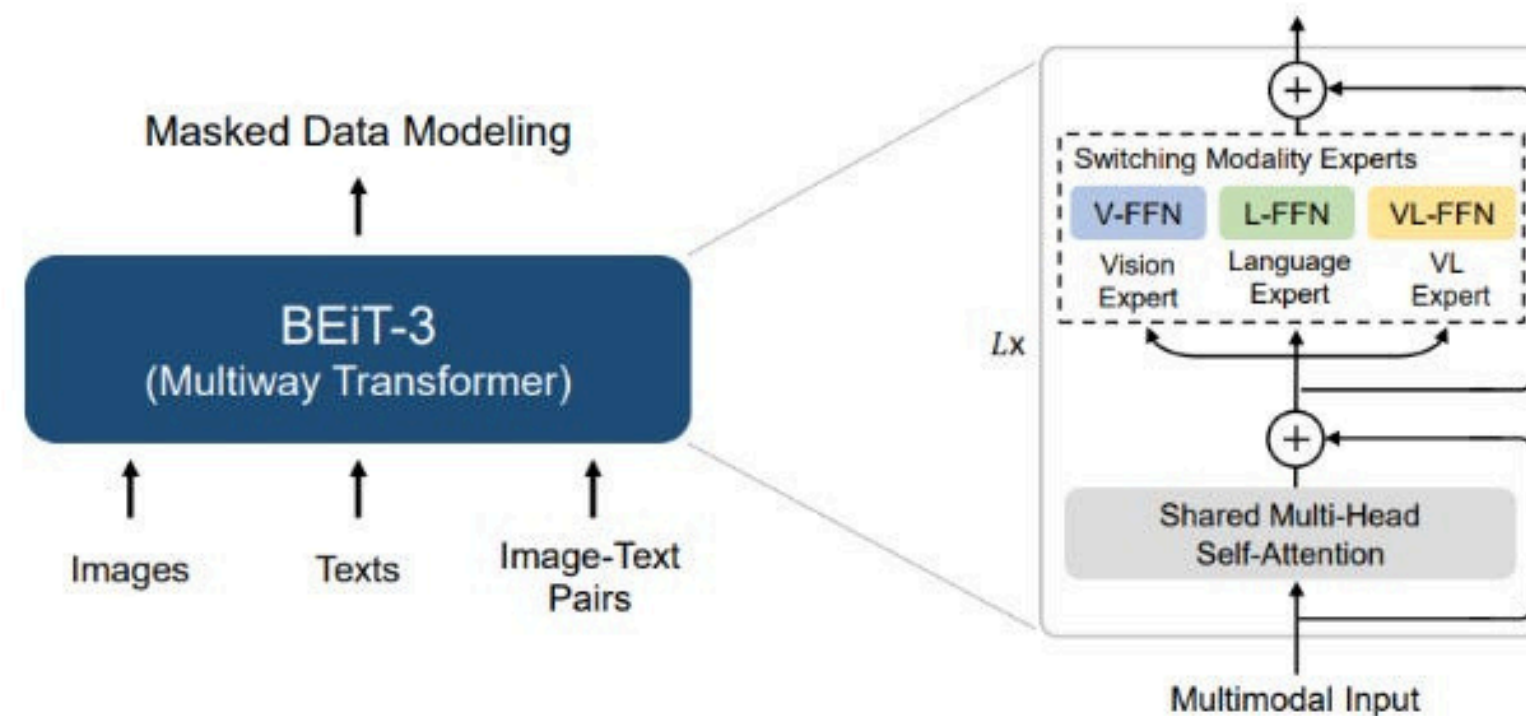
Beit-3

# ViLT

Finetune from
dandelin/vilt-b32-finetuned-vqa

# Beit-3



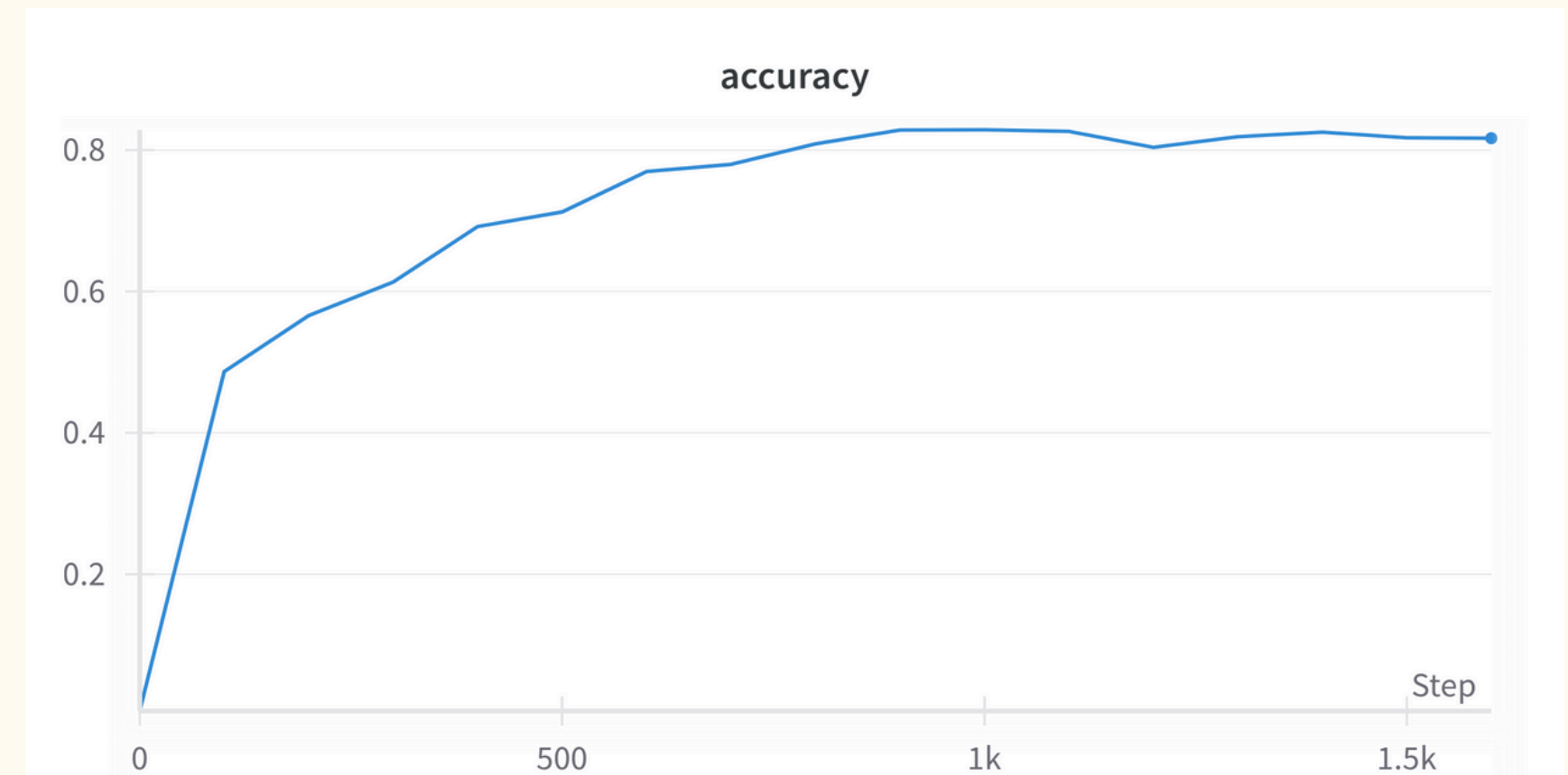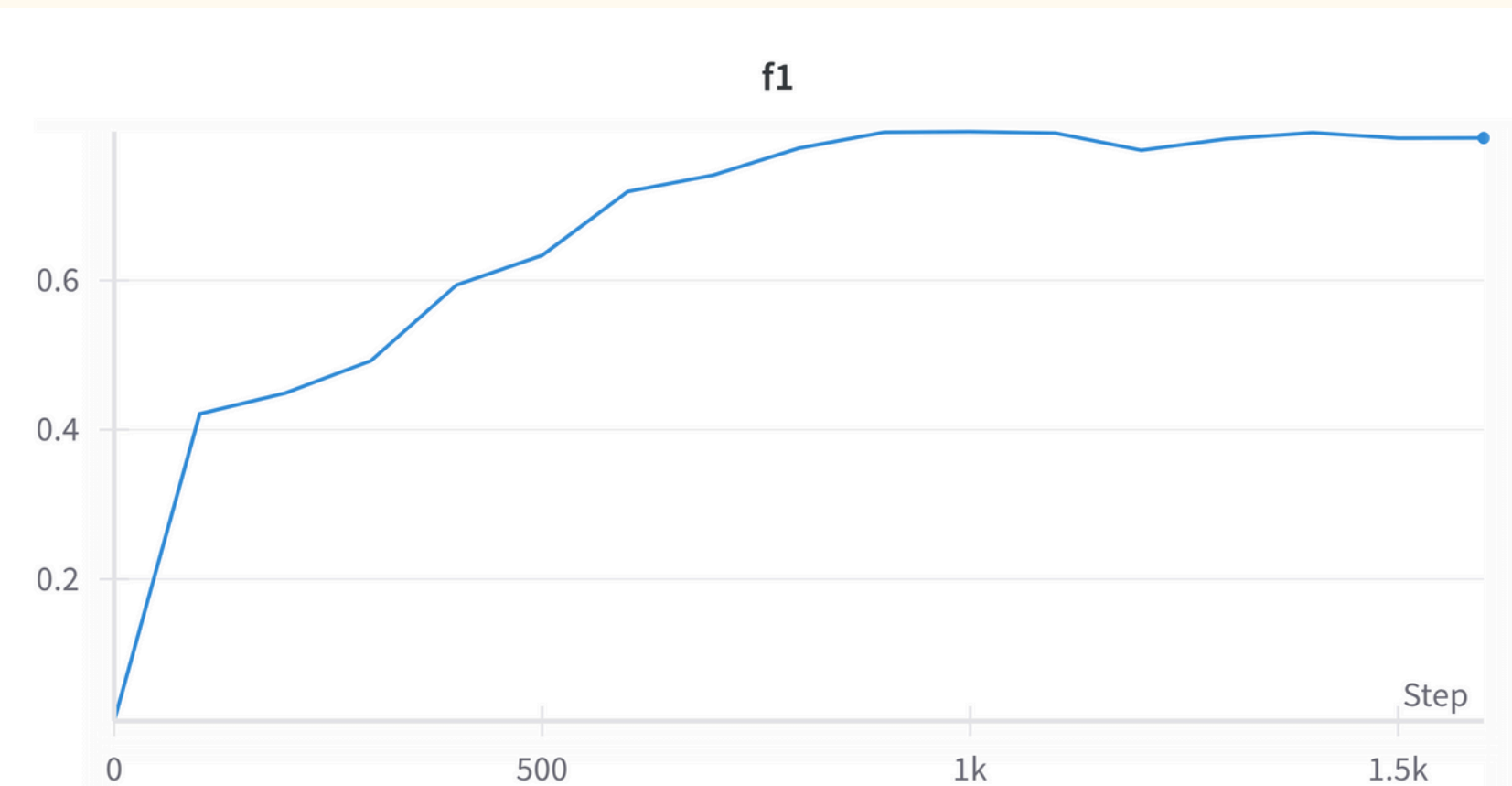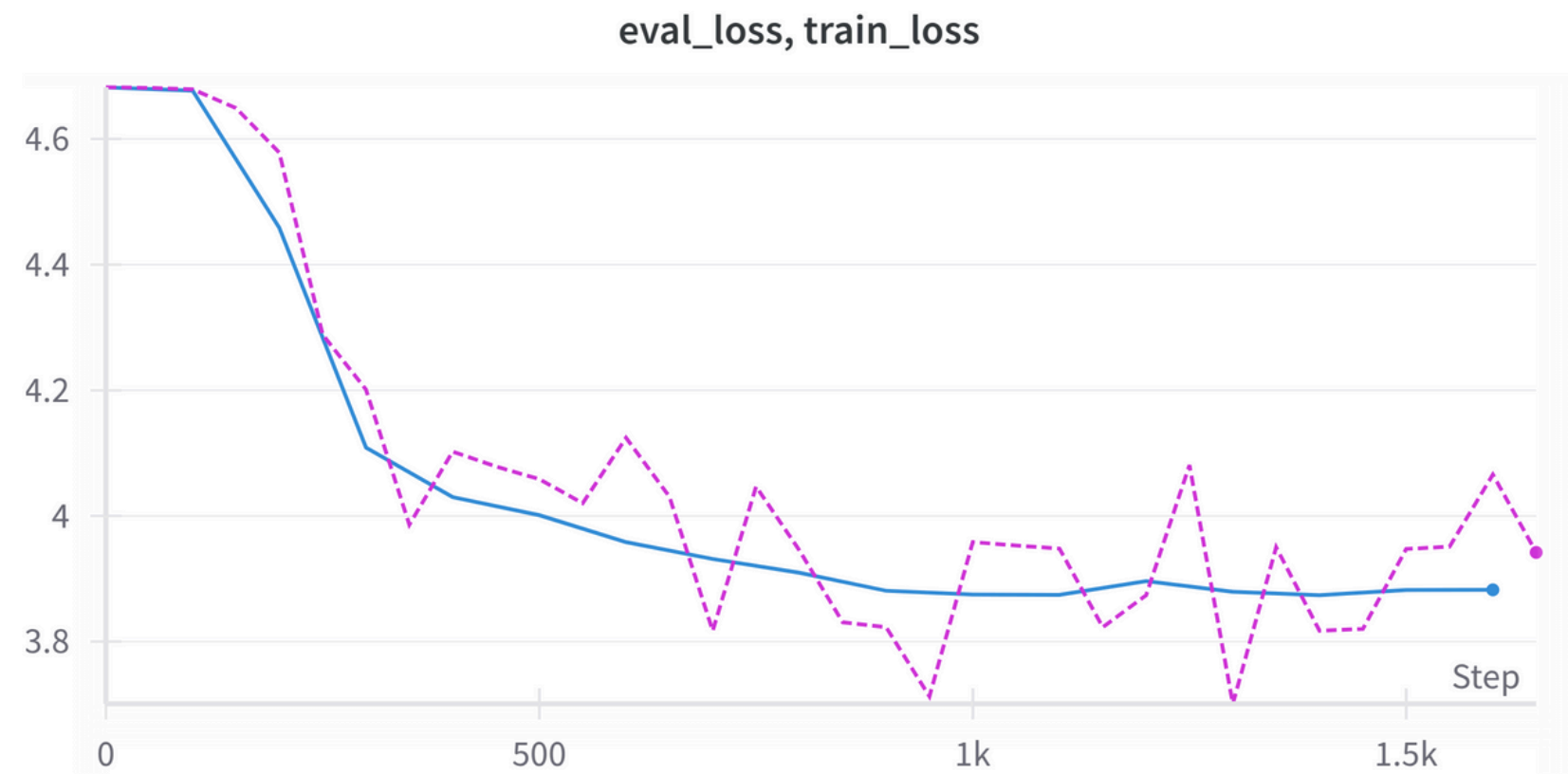Image as a Foreign Language: BEiT Pretraining for All Vision and Vision-Language Tasks

Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, Zhiliang Peng, Qiang Liu
Kriti Aggarwal, Owais Khan Mohammed, Saksham Singhal, Subhojit Som, Furu Wei[†]
Microsoft Corporation
https://aka.ms/beit-3

beit3_base_patch16_224
(VQAv2 finetuned)

# Thanks for watching