```python
import argparse
import os
from datetime import datetime
from collections import Counter, defaultdict
import json
import pickle

import numpy as np
import torch
from torch import nn
from torchvision import datasets, transforms, models
import torchvision
import matplotlib.pyplot as plt
from torchsummary import summary
from tqdm import tqdm
import pandas as pd
from scipy.stats import pearsonr, spearmanr
import torch.nn.functional as F
from torch.nn.functional import softmax
import numpy

# ================================================================
# Parsing the arguments
# ================================================================
parser = argparse.ArgumentParser()
parser.add_argument("--test_run", type=int, required=True, help="Experiment run
    index")
parser.add_argument("--resnet_ckpt_path", type=str, required=True,
                    help="Path to saved ResNet checkpoint (.pth)")
args = parser.parse_args()

test_run = args.test_run
resnet_ckpt_path = args.resnet_ckpt_path

# ================================================================
# Device configuration
# ================================================================
device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
print("Device:", device)
print("test_run", test_run)
print("resnet_ckpt_path:", resnet_ckpt_path)

# ================================================================
# DATA LOADING
# ================================================================
transformations = transforms.Compose([
    transforms.Resize(255),
    transforms.CenterCrop(224),
    transforms.ToTensor(),
    transforms.Normalize(mean=[0.485, 0.456, 0.406],
                         std=[0.229, 0.224, 0.225])
])

base_data_dir =
    f"/leonardo_work/IscrC_ArtLLMs/alitariqnagi_work/training_data_split_run_{test_run}"

train_set_mix = datasets.ImageFolder(
    os.path.join(base_data_dir, "train"),
    transform=transformations
)
```

```python
59 val_set_mix = datasets.ImageFolder(
60     os.path.join(base_data_dir, "val"),
61     transform=transformations
62 )
63 test_set_mix = datasets.ImageFolder(
64     os.path.join(base_data_dir, "test"),
65     transform=transformations
66 )
67
68 print("Train size:", len(train_set_mix))
69 print("Val size:", len(val_set_mix))
70 print("Test size:", len(test_set_mix))
71
72 class_names = train_set_mix.classes
73 print("Classes:", class_names)
74
75 train_loader_mix = torch.utils.data.DataLoader(
76     train_set_mix, batch_size=32, shuffle=True, num_workers=8
77 )
78 val_loader_mix = torch.utils.data.DataLoader(
79     val_set_mix, batch_size=32, shuffle=False, num_workers=8
80 )
81 test_loader_mix = torch.utils.data.DataLoader(
82     test_set_mix, batch_size=32, shuffle=False, num_workers=8
83 )
84
85 # Optional: quick visualization
86 dataiter = iter(train_loader_mix)
87 images, labels = next(dataiter)
88
89
90 def imshow(img):
91     img = img / 2 + 0.5   # unnormalize
92     npimg = img.numpy()
93     plt.imshow(np.transpose(npimg, (1, 2, 0)))
94     plt.axis("off")
95
96
97 imshow(torchvision.utils.make_grid(images))
98 plt.show()
99
100 model = models.resnet50(weights='ResNet50_Weights.IMAGENET1K_V1')
101
102 for param in model.parameters():
103     param.requires_grad = True
104
105 model.fc = nn.Sequential(
106     nn.Linear(2048, 512),
107     nn.ReLU(),
108     nn.Dropout(0.2),
109     nn.Linear(512, len(class_names))
110 )
111
112 if torch.cuda.device_count() > 1:
113     print(f"Using DataParallel Mode with cuda count {torch.cuda.device_count()}")
114     model = torch.nn.DataParallel(model)
115
116 model.to(device)
117 print(summary(model, input_size=(3, 224, 224)))
118
```

```python
119  # ============================================================
120  # LOAD FINAL CHECKPOINT (NO TRAINING HERE)
121  # ============================================================
122  final_ckpt_path = resnet_ckpt_path
123  print("Loading final model checkpoint:", final_ckpt_path)
124
125  # Load on CPU first (safer across GPU setups)
126  state_dict = torch.load(final_ckpt_path, map_location="cpu")
127
128  # If checkpoint was saved using DataParallel, keys will start with "module."
129  if any(k.startswith("module.") for k in state_dict.keys()):
130      print("[INFO] Detected DataParallel checkpoint. Stripping 'module.' prefix...")
131      state_dict = {k.replace("module.", "", 1): v for k, v in state_dict.items()}
132
133  model.load_state_dict(state_dict, strict=True)
134
135  model.to(device)
136  model.eval()
137  print("Model loaded.\n")
138
139  # ============================================================
140  # CONFORMAL FUNCTIONS
141  # ============================================================
142
143  # different values will be overwritten in the loop
144  alpha = 0.1
145
146
147  #########################----------------------------------------
148  # Calculation of Calibration Scores
149  # ----------------------------------------------------------------
150  def calibration_scores(model, data_loader, device):
151      model.eval()
152
153      scores_calibration_set = []
154
155      with torch.no_grad():
156          for X, y in tqdm(data_loader,
157                          desc="#################Proceeding with the Calculation of
     Calibration Scores#################"):
158              X = X.to(device)
159              y = y.to(device)
160
161              y = y.cpu().numpy()
162
163              pred_logits = model(X)
164
165              cal_smx = softmax(pred_logits, dim=1).cpu().numpy()
166              scores_per_batch = 1 - cal_smx[np.arange(y.shape[0]), y]
167
168              scores_calibration_set.extend(scores_per_batch)
169
170      scores_calibration_set = numpy.array(scores_calibration_set)
171
172      return scores_calibration_set
173
174
175  ##############################################################################
     ###########################################################
176
```

```python
177
178  ########################---------------------------------------
179  # Calculation of Prediction Sets
180  # ----------------------------------------------------------------
181  def predict_with_conformal_sets(model, data_loader, calibration_scores_array, alpha,
     device):
182      model.eval()
183
184      prediction_sets = []
185      true_labels = []
186
187      n = len(calibration_scores_array)
188
189      q_level = np.ceil((n + 1) * (1 - alpha)) / n
190      qhat = np.quantile(calibration_scores_array, q_level, method='higher')
191
192      with torch.no_grad():
193          for X, y in tqdm(data_loader,
194                           desc="################Proceeding with the Calculation of
     Testing Scores################"):
195              X = X.to(device)
196              y = y.cpu().numpy()
197              pred_logits = model(X)
198
199              val_smx = softmax(pred_logits, dim=1)
200              val_smx = val_smx.cpu().numpy()
201
202              masks = val_smx >= (1 - qhat)
203
204              for _ in masks:
205                  prediction_sets.append(np.where(_)[0].tolist())
206
207              true_labels.extend(y.tolist())
208
209      return prediction_sets, true_labels
210
211
212  ####################################################################################
     ###########################################################
213
214
215  ########################---------------------------------------
216  # EVALUATION
217  # ----------------------------------------------------------------
218  def evaluate(prediction_sets, labels):
219      number_of_correct_predictions = 0
220      for pred_set, true_label in zip(prediction_sets, labels):
221          if true_label in pred_set:
222              number_of_correct_predictions += 1
223      empirical_test_coverage = number_of_correct_predictions / len(labels)
224      average_prediction_set_size = np.mean([len(s) for s in prediction_sets])
225      return empirical_test_coverage, average_prediction_set_size
226
227
228  ####################################################################################
     ###########################################################
229
230
231  def feature_stratified_coverage_by_class(prediction_sets, true_labels, alpha=0.1,
     class_names=None):
```

```python
        class_to_indices = defaultdict(list)
        for idx, label in enumerate(true_labels):
            class_to_indices[label].append(idx)

        class_coverages = {}
        for cls, indices in class_to_indices.items():
            covered = sum(true_labels[i] in prediction_sets[i] for i in indices)
            coverage = covered / len(indices)
            class_coverages[cls] = coverage

        print("\n Feature-Stratified Coverage (by True Class Label):")
        for cls, cov in sorted(class_coverages.items()):
            name = class_names[cls] if class_names else str(cls)
            print(f"  Class {name:<20}: Coverage = {cov:.3f} (Target ≥ {1 - alpha:.2f})")

        fsc_metric = min(class_coverages.values())
        print(f"\n FSC Metric (minimum class-wise coverage): {fsc_metric:.3f}")

        return fsc_metric, class_coverages


def plot_class_wise_coverage(class_coverages, alpha=0.1, class_names=None,
    title="Class-wise Coverage (FSC)"):
        classes = list(class_coverages.keys())
        coverages = [class_coverages[c] for c in classes]

        labels = [class_names[c] if class_names else str(c) for c in classes]

        plt.figure(figsize=(12, 6))
        bars = plt.bar(labels, coverages, color="pink", edgecolor="black")

        target = 1 - alpha
        plt.axhline(target, color='red', linestyle='--', label=f"Target Coverage
    ({target:.2f})")

        for bar, cov in zip(bars, coverages):
            if cov < target:
                bar.set_color('salmon')

        plt.xticks(rotation=45, ha="right")
        plt.ylim(0, 1.05)
        plt.ylabel("Coverage")
        plt.title(title)
        plt.legend()
        plt.tight_layout()
        filename = f"Class_wise_Coverage_(FSC)_{test_run}_{timestamp}.png"
        plt.savefig(os.path.join(results_directory, filename))
        plt.show()


def plot_coverage_vs_class_frequency(class_coverages, true_labels, alpha=0.1,
    class_names=None,
                                      title="Coverage vs Class Frequency"):
        class_indices = sorted(class_coverages.keys())
        coverages = [class_coverages[c] for c in class_indices]

        total = len(true_labels)
        counts = Counter(true_labels)
        frequencies = [counts[c] / total for c in class_indices]
```

```python
289        labels = [class_names[c] if class_names else str(c) for c in class_indices]
290        x = np.arange(len(class_indices))
291
292        fig, ax1 = plt.subplots(figsize=(12, 6))
293
294        bars1 = ax1.bar(x - 0.2, coverages, width=0.4, label="Coverage", color='pink',
      edgecolor='black')
295        ax1.axhline(1 - alpha, color='red', linestyle='--', label=f"Target Coverage ({1 -
      alpha:.2f})")
296        ax1.set_ylim(0, 1.05)
297        ax1.set_ylabel("Coverage", color='pink')
298        ax1.tick_params(axis='y', labelcolor='pink')
299
300        for bar, cov in zip(bars1, coverages):
301            if cov < (1 - alpha):
302                bar.set_color('salmon')
303
304        ax2 = ax1.twinx()
305        bars2 = ax2.bar(x + 0.2, frequencies, width=0.4, label="Class Frequency",
      color='gray', alpha=0.5)
306        ax2.set_ylabel("Class Frequency", color='gray')
307        ax2.tick_params(axis='y', labelcolor='gray')
308
309        plt.xticks(x, labels, rotation=45, ha="right")
310        plt.title(title)
311        fig.legend(loc='upper right', bbox_to_anchor=(0.85, 0.85))
312        plt.tight_layout()
313        filename = f"plot_coverage_vs_class_frequency_{test_run}_{timestamp}.png"
314        plt.savefig(os.path.join(results_directory, filename))
315        plt.show()
316
317
318 def plot_coverage_vs_class_frequency_with_correlation(class_coverages, true_labels,
      alpha=0.1, class_names=None,
319                                                        title="Coverage vs Class
      Frequency"):
320        class_indices = sorted(class_coverages.keys())
321        coverages = [class_coverages[c] for c in class_indices]
322
323        total = len(true_labels)
324        counts = Counter(true_labels)
325        frequencies = [counts[c] / total for c in class_indices]
326
327        labels = [class_names[c] if class_names else str(c) for c in class_indices]
328        x = np.arange(len(class_indices))
329
330        pearson_corr, pearson_p = pearsonr(frequencies, coverages)
331        spearman_corr, spearman_p = spearmanr(frequencies, coverages)
332
333        print(f"\n Correlation between class frequency and coverage:")
334        print(f"   - Pearson  r = {pearson_corr:.3f}, p = {pearson_p:.3e}")
335        print(f"   - Spearman ρ = {spearman_corr:.3f}, p = {spearman_p:.3e}")
336
337        fig, ax1 = plt.subplots(figsize=(12, 6))
338
339        bars1 = ax1.bar(x - 0.2, coverages, width=0.4, label="Coverage", color='pink',
      edgecolor='black')
340        ax1.axhline(1 - alpha, color='red', linestyle='--', label=f"Target Coverage ({1 -
      alpha:.2f})")
341        ax1.set_ylim(0, 1.05)
```

```python
342        ax1.set_ylabel("Coverage", color='pink')
343        ax1.tick_params(axis='y', labelcolor='pink')
344
345        for bar, cov in zip(bars1, coverages):
346            if cov < (1 - alpha):
347                bar.set_color('salmon')
348
349        ax2 = ax1.twinx()
350        bars2 = ax2.bar(x + 0.2, frequencies, width=0.4, label="Class Frequency",
       color='gray', alpha=0.5)
351        ax2.set_ylabel("Class Frequency", color='gray')
352        ax2.tick_params(axis='y', labelcolor='gray')
353
354        plt.xticks(x, labels, rotation=45, ha="right")
355        plt.title(title)
356        fig.legend(loc='upper right', bbox_to_anchor=(0.85, 0.85))
357        plt.tight_layout()
358        filename =
       f"plot_coverage_vs_class_frequency_with_correlation_{test_run}_{timestamp}.png"
359        plt.savefig(os.path.join(results_directory, filename))
360        plt.show()
361
362
363 def size_stratified_coverage(prediction_sets, true_labels, alpha=0.1,
364                              bins=[[1], [2], [3, 4, 5], list(range(6, 100))]):
365        bin_groups = defaultdict(list)
366
367        for i, pred_set in enumerate(prediction_sets):
368            set_size = len(pred_set)
369            for bin_range in bins:
370                if set_size in bin_range:
371                    bin_groups[str(bin_range)].append(i)
372                    break
373
374        bin_coverages = {}
375        for bin_key, indices in bin_groups.items():
376            covered = sum(true_labels[i] in prediction_sets[i] for i in indices)
377            coverage = covered / len(indices) if indices else 0
378            bin_coverages[bin_key] = coverage
379
380        print("\n Size-Stratified Coverage (SSC):")
381        for bin_key, cov in bin_coverages.items():
382            print(f"  Set Size {bin_key:10s}: Coverage = {cov:.3f} (Target ≥ {1 -
       alpha:.2f})")
383
384        ssc = min(bin_coverages.values())
385        print(f"\n SSC Metric (min bin-wise coverage): {ssc:.3f}")
386        return ssc, bin_coverages
387
388
389 def plot_ssc_coverage(bin_coverages, alpha=0.1, title="Size-Stratified Coverage"):
390        bin_labels = list(bin_coverages.keys())
391        coverages = [bin_coverages[k] for k in bin_labels]
392
393        plt.figure(figsize=(10, 5))
394        bars = plt.bar(bin_labels, coverages, color='pink', edgecolor='black')
395
396        target = 1 - alpha
397        plt.axhline(target, color='red', linestyle='--', label=f"Target Coverage
       ({target:.2f})")
```

```python
398
399        for bar, cov in zip(bars, coverages):
400            if cov < target:
401                bar.set_color('grey')
402
403        plt.xticks(rotation=45, ha='right')
404        plt.ylim(0, 1.05)
405        plt.ylabel("Coverage")
406        plt.title(title)
407        plt.legend()
408        plt.tight_layout()
409        filename = f"plot_ssc_coverage_{test_run}__{timestamp}.png"
410        plt.savefig(os.path.join(results_directory, filename))
411        plt.show()
412
413
414 def plot_prediction_set_size_by_class(prediction_sets, true_labels, class_names=None,
    results_dir="", test_run="",
415                                          timestamp=""):
416        size_per_class = defaultdict(list)
417        for pred_set, true_label in zip(prediction_sets, true_labels):
418            size_per_class[true_label].append(len(pred_set))
419
420        sorted_classes = sorted(size_per_class.keys())
421        means = [np.mean(size_per_class[c]) for c in sorted_classes]
422        stds = [np.std(size_per_class[c]) for c in sorted_classes]
423        labels = [class_names[c] if class_names else str(c) for c in sorted_classes]
424
425        plt.figure(figsize=(12, 6))
426        plt.bar(labels, means, yerr=stds, capsize=5, color="lightgreen",
    edgecolor="black")
427        plt.ylabel("Avg Prediction Set Size")
428        plt.xlabel("Class")
429        plt.title("Average Prediction Set Size by Class")
430        plt.xticks(rotation=45, ha="right")
431        plt.grid(axis="y", linestyle="--", alpha=0.5)
432        plt.tight_layout()
433
434        filename = f"prediction_set_size_by_class_{test_run}__{timestamp}.png"
435        plt.savefig(os.path.join(results_directory, filename))
436        plt.show()
437
438
439 def plot_calibration_score_distribution(calibration_scores, results_dir="",
    test_run="", timestamp=""):
440        plt.figure(figsize=(10, 5))
441        plt.hist(calibration_scores, bins=30, color="pink", edgecolor="black", alpha=0.8)
442        plt.xlabel("Nonconformity Score (1 - P(True Class))")
443        plt.ylabel("Frequency")
444        plt.title("Calibration Score Distribution")
445        plt.grid(axis='y', linestyle='--', alpha=0.6)
446        plt.tight_layout()
447
448        filename = f"calibration_score_distribution_{test_run}__{timestamp}.png"
449        plt.savefig(os.path.join(results_directory, filename))
450
451        plt.show()
452
453
454 # ============================================================
```

```python
# MAIN CONFORMAL EVALUATION FOR MULTIPLE ALPHA VALUES
# ==========================================================

coverage_levels = np.arange(0.65, 0.95 + 1e-9, 0.05)

timestamp = datetime.now().strftime("%Y-%m-%d_%H-%M-%S")

base_results_root =
    f"Results_normal_Conformal_Prediction_experiment_run_{test_run}_{timestamp}"
os.makedirs(base_results_root, exist_ok=True)
print("Base results directory:", base_results_root)

print("\n====== Computing calibration scores on validation set ======\n")
results_directory = base_results_root
scores_calibration = calibration_scores(model=model, data_loader=val_loader_mix,
    device=device)

np.save(os.path.join(base_results_root,

    f"calibration_scores_experiment_run_{test_run}_{timestamp}.npy"), scores_calibration)
with open(os.path.join(base_results_root,
                       f"scores_calibration_run_{test_run}__{timestamp}.pkl"), "wb")
    as f:
    pickle.dump(scores_calibration, f)

plot_calibration_score_distribution(
    calibration_scores=scores_calibration,
    results_dir=base_results_root,
    test_run=test_run,
    timestamp=timestamp
)

all_alpha_metrics = []

for coverage_target in coverage_levels:
    alpha = 1.0 - coverage_target

    print("\n==================================================")
    print(f" Target coverage: {coverage_target:.2f}  |  alpha = {alpha:.2f}")
    print("==================================================\n")

    cov_tag = int(round(coverage_target * 100))
    results_directory = os.path.join(base_results_root, f"coverage_{cov_tag}")
    os.makedirs(results_directory, exist_ok=True)

    prediction_sets, true_labels = predict_with_conformal_sets(
        model=model,
        data_loader=test_loader_mix,
        calibration_scores_array=scores_calibration,
        alpha=alpha,
        device=device
    )
    test_set_coverage, avg_size = evaluate(prediction_sets=prediction_sets,
    labels=true_labels)

    log_path = os.path.join(results_directory,

     f"logs_experiment_run_{test_run}_alpha_{alpha:.2f}_{timestamp}.txt")
    with open(log_path, "w") as f:
```

```python
508          f.write("######################### Results
    ##################################\n")
509          f.write(f"test_run {test_run}\n")
510          f.write(f"Target coverage: {coverage_target:.3f}\n")
511          f.write(f"Alpha: {alpha:.3f}\n")
512          f.write(f"Coverage: {test_set_coverage:.3f}\n")
513          f.write(f"Average Prediction Set Size for Test Set: {avg_size:.3f}\n")
514
     f.write("##################################################################\n")
515
516      print("########################Results:##################################")
517      print(f"test_run {test_run}")
518      print(f"Target coverage: {coverage_target:.3f}")
519      print(f"Alpha: {alpha:.3f}")
520      print(f"Coverage: {test_set_coverage:.3f}")
521      print(f"Average Prediction Set Size for Test Set: {avg_size:.3f}")
522      print("##################################################################")
523
524      set_sizes = [len(prediction_set) for prediction_set in prediction_sets]
525      counts = Counter(set_sizes)
526      total = sum(counts.values())
527      sizes = sorted(counts.keys())
528      frequencies = [counts[size] for size in sizes]
529      percentages = [100 * freq / total for freq in frequencies]
530
531      data = [{'size': len(prediction_set), 'correct': int(true in prediction_set)}
532              for prediction_set, true in zip(prediction_sets, true_labels)]
533      df = pd.DataFrame(data)
534      coverage_stats = df.groupby('size').agg({'correct': ['mean', 'sum', 'count']})
535      coverage_stats.columns = ['coverage', 'correct_count', 'total_count']
536      coverage_stats = coverage_stats.reindex(sizes, fill_value=0)
537
538      cumulative_correct = coverage_stats['correct_count'].cumsum()
539      cumulative_coverage = 100 * cumulative_correct / len(true_labels)
540      max_freq = max(frequencies)
541      cumulative_scaled = (cumulative_coverage / 100) * max_freq
542
543      fig, ax = plt.subplots(figsize=(13, 6))
544      bars = ax.bar(sizes, frequencies, color='pink', edgecolor='black',
    label='Frequency', alpha=0.7)
545      ax.plot(sizes, cumulative_scaled, color='red', marker='o', linestyle='-',
546              alpha=0.4, linewidth=2, label='Cumulative Coverage')
547
548      for bar, size, pct in zip(bars, sizes, percentages):
549          height = bar.get_height()
550          count = counts[size]
551          coverage_val = coverage_stats.loc[size, 'coverage'] * 100
552          ax.annotate(f'{count} ({pct:.1f}%)\nCov: {coverage_val:.1f}%',
553                      xy=(bar.get_x() + bar.get_width() / 2, height / 2),
554                      ha='center', va='center', fontsize=9)
555
556      for bar, cum_cov in zip(bars, cumulative_coverage):
557          height = bar.get_height()
558          ax.annotate(f'CumCov: {cum_cov:.1f}%',
559                      xy=(bar.get_x() + bar.get_width() / 2, height),
560                      xytext=(0, 8),
561                      textcoords="offset points",
562                      ha='center', va='bottom', fontsize=9, color='darkorange')
563
564      mean_size = np.mean(set_sizes)
```

```
565        median_size = np.median(set_sizes)
566        ax.axvline(mean_size, color='gray', linestyle='--', label=f'Mean =
      {mean_size:.2f}')
567        ax.axvline(median_size, color='orange', linestyle='--', label=f'Median =
      {median_size:.2f}')
568
569        ax.set_xlabel("Prediction Set Size")
570        ax.set_ylabel("Frequency")
571        ax.set_title(f"Prediction Set Size Distribution with Coverage (Test Set), alpha=
      {alpha:.2f}")
572        ax.set_xticks(sizes)
573        ax.grid(axis='y', linestyle='--', alpha=0.6)
574        ax.legend(loc='upper right')
575
576        plt.tight_layout()
577        plt.savefig(os.path.join(
578            results_directory,
579
       f"Prediction_Set_Size_Distribution_with_Coverage_(Test_Set)_{test_run}_alpha_{alpha:
      .2f}_{timestamp}.png"
580        ))
581        plt.show()
582
583        fsc_metric, per_class_cov = feature_stratified_coverage_by_class(
584            prediction_sets,
585            true_labels,
586            alpha=alpha,
587            class_names=class_names
588        )
589
590        fsc_metric, class_coverages = feature_stratified_coverage_by_class(
591            prediction_sets,
592            true_labels,
593            alpha=alpha,
594            class_names=class_names
595        )
596        plot_class_wise_coverage(class_coverages, alpha=alpha, class_names=class_names)
597
598        fsc_metric, class_coverages = feature_stratified_coverage_by_class(
599            prediction_sets,
600            true_labels,
601            alpha=alpha,
602            class_names=class_names
603        )
604        plot_coverage_vs_class_frequency(class_coverages, true_labels,
605                                         alpha=alpha, class_names=class_names)
606
607        plot_coverage_vs_class_frequency_with_correlation(
608            class_coverages=class_coverages,
609            true_labels=true_labels,
610            alpha=alpha,
611            class_names=class_names
612        )
613
614        ssc_metric, bin_coverages = size_stratified_coverage(
615            prediction_sets,
616            true_labels,
617            alpha=alpha,
618            bins=[[1], [2], [3], list(range(4, 100))]
619        )
```

```python
    ssc_metric, bin_coverages = size_stratified_coverage(
        prediction_sets,
        true_labels,
        alpha=alpha,
        bins=[[1], [2], [3], list(range(4, 100))]
    )
    plot_ssc_coverage(bin_coverages, alpha=alpha)

    plot_prediction_set_size_by_class(
        prediction_sets=prediction_sets,
        true_labels=true_labels,
        class_names=class_names,
        results_dir=results_directory,
        test_run=test_run,
        timestamp=timestamp
    )

    metrics = {
        "test_run": test_run,
        "timestamp": timestamp,
        "alpha": alpha,
        "target_coverage": coverage_target,
        "test_set_coverage": test_set_coverage,
        "average_prediction_set_size": avg_size,
        "FSC": fsc_metric,
        "SSC": ssc_metric,
        "class_coverages": {str(k): float(v) for k, v in class_coverages.items()},
        "bin_coverages": bin_coverages
    }

    with open(os.path.join(results_directory,

f"metrics_summary_{test_run}_alpha_{alpha:.2f}_{timestamp}.json"), "w") as f:
        json.dump(metrics, f, indent=4)

    all_alpha_metrics.append(metrics)

with open(os.path.join(base_results_root,
                       f"metrics_summary_all_alphas_{test_run}_{timestamp}.json"),
    "w") as f:
    json.dump(all_alpha_metrics, f, indent=4)

print("\nFinished multi-alpha conformal evaluation for alphas corresponding to
    coverage 0.65-0.95.")
print("Results saved under:", base_results_root)
```