

# Statistical Comparison of PoS Committee Selection Algorithms

Robert Jones

23 April 2025

## 1 Simulation Experiment

We describe a Monte Carlo simulation experiment designed to evaluate two algorithms (A and B) that select committee members in a Proof-of-Stake (PoS) blockchain to achieve Byzantine Fault Tolerance (BFT) consensus. Their performance is measured by their robustness against faults when maintaining the 2/3 majority required for consensus. In this experiment, we sample a set of system configuration parameters  $\alpha$  from a joint distribution. For each sampled  $\alpha$ , we run many epochs (committee selections) for both algorithms and estimate a *fault tolerance score* (FTS) for both algorithms. The results are then analyzed using statistical tests to determine if there is a significant difference in FTS between the two algorithms.

For a given algorithm, let

$$p_f = \text{Probability that the system tolerates } f \text{ faults}$$

for  $f = 1, 2, 3, \dots$ . The **Fault Tolerance Score (FTS)** for the algorithm is then defined as:

$$\text{FTS} = \sum_{f=1}^{\infty} p_f,$$

which can be interpreted as the expected number of faults the system can tolerate before losing consensus.

- For **Algorithm A**, we denote the score as  $\text{FTS}_A$ .
- For **Algorithm B**, we denote the score as  $\text{FTS}_B$ .

A higher Fault Tolerance Score indicates that an algorithm can, on average, tolerate more faults while still preserving its 2/3 majority, thereby demonstrating better resilience.

To compare the algorithms, we examine the difference in their scores:

$$\Delta\text{FTS} = \text{FTS}_B - \text{FTS}_A.$$

- If  $\Delta\text{FTS} > 0$ , then **Algorithm B** is more resilient (i.e., it tolerates a higher number of faults on average).
- If  $\Delta\text{FTS} < 0$ , then **Algorithm A** performs better in terms of fault tolerance.
- If  $\Delta\text{FTS} \approx 0$ , then there is no significant difference between the two in terms of their fault tolerance capability.

## 2 Paired Statistical Testing

A crucial aspect of this experiment is the variation introduced by sampling  $\alpha$  from a distribution. By pairing the outcomes for A and B within each  $\alpha$ , we have structured the test as a randomized block design, where each block is a fixed environment  $\alpha$  and we compare two treatments (A vs B) within that block. This design improves sensitivity by removing inter-block (inter- $\alpha$ ) variability. Essentially, each scenario acts as its own control. If we instead ignored  $\alpha$  and pooled all epochs together, random differences in scenario difficulty would obscure the comparison.

Given that we have a sample of paired differences  $\Delta\text{FTS}(\alpha)$  for many random scenarios, the natural approach uses a paired statistical test to assess whether the mean of  $\Delta\text{FTS}(\alpha)$  is more significant than zero. The paired test treats the set of differences as a single sample. If the number of sampled scenarios is large enough, one common choice is the paired Student's t-test, which tests whether the mean difference is significantly different from zero. The paired t-test assumes the differences  $\Delta\text{FTS}(\alpha)$  are approximately normally distributed (an assumption often reasonable by the Central Limit Theorem if the number of samples is moderately large). It also assumes the sample of scenarios is independent. Here we know each  $\alpha$  is drawn independently, but we cannot assess normality in the data distribution.

### 2.1 Hypothesis Test

To evaluate whether Algorithm B is statistically superior to Algorithm A, we conduct a *one-sided hypothesis test*:

- **Null Hypothesis**,  $H_0$ : Algorithm B offers no improvement over A.

$$\text{mean}(\Delta\text{FTS}) = 0$$

- **Alternative Hypothesis**,  $H_1$ : Algorithm B is more fault tolerant.

$$\text{mean}(\Delta\text{FTS}) > 0$$

### 2.2 Bootstrap Sampling

Traditional statistical tests assume normality in the data distribution, which, after some testing, we determine is not the case here—rather, the data distribution is heavy tailed. Therefore, we will use a non-parametric alternative called *bootstrap sampling*. This approach is particularly useful when the underlying distribution is unknown or when the sample size is small. The bootstrap method allows us to estimate the sampling distribution of a statistic—in this case,  $\Delta\text{FTS}(\alpha)$  for a given parameter set  $\alpha$ —by resampling the observed data *with replacement*.

Now, we can conduct bootstrap resampling as follows:

1. **Resample**: Draw a large number (e.g., 10,000) of bootstrap samples from the centered differences (with replacement).
2. **Statistic**: For each bootstrap sample, compute its mean. This collection of sample means forms an empirical null distribution of the mean difference.
3. **Construct Empirical Null Distribution**: Shift the observed differences by subtracting the empirical mean (centered data) to simulate what would be expected under  $H_0$ .

4. **Compute Bootstrap Means:** For each bootstrap sample, calculate the mean difference in FTS scores.
5. **Formulate the One-Sided p-Value:** The one-sided p-value is calculated as

$$p = \frac{\text{number of bootstrap replicates with mean} \geq \text{observed mean}}{\text{total number of replicates}}.$$

6. **Interpret the p-value:** If this p-value is small (say, less than 0.05), it implies that such a high mean difference is very unlikely under the null hypothesis, lending support to the claim that algorithm B's performance is significantly better than algorithm A's.
7. **Conclusion:** If the p-value is small, we reject the null hypothesis  $H_0$  in favor of the alternative hypothesis  $H_1$ . This suggests that Algorithm B is statistically significantly better than Algorithm A.

### 3 Experimental Results

Appendix 1 contains the Python code used to conduct the Monte Carlo simulation experiment. The code implements the algorithms, generates the parameter configurations, and computes the Fault Tolerance Scores (FTS) for both Algorithm A and Algorithm B. It also includes the bootstrap sampling procedure to assess the statistical significance of the differences in FTS scores.

Running the code produces the results of the Monte Carlo simulation experiment, including the histograms of FTS scores for both algorithms and the histogram of the differences in FTS scores ( $\Delta$ FTS). The code also performs the bootstrap hypothesis test to determine whether Algorithm B outperforms Algorithm A.

The output of the code:

```
simulation trials: 100% [1000/1000] [13:54<00:00, 1.20it/s]
Observed Mean Difference (B - A): 1.733
One-sided p-value: 0.0000
Conclusion: Algorithm B is significantly better than Algorithm A.
```

#### 3.1 Discussion

The simulation trials were run for 1000 iterations, and the observed mean difference in Fault Tolerance Scores (FTS) between Algorithm B and Algorithm A was found to be 1.733. The one-sided p-value was less than 0.0001, indicating strong evidence against the null hypothesis  $H_0$ . Therefore, we conclude that Algorithm B is significantly better than Algorithm A in terms of fault tolerance.

The code also generates visualizations of the results—figures 1, 2, and 3— including histograms of the FTS scores for both algorithms and the histogram of the differences in FTS scores ( $\Delta$ FTS). The histograms provide a visual representation of the distribution of FTS values for both algorithms across multiple parameter configurations.

The figures illustrate the results of our Monte Carlo simulation experiment. Figure 1 shows the comparison of Fault Tolerance Scores (FTS) between Algorithm A and Algorithm B across different parameter configurations. The overlaid histograms in figure 2 displays the distribution of FTS

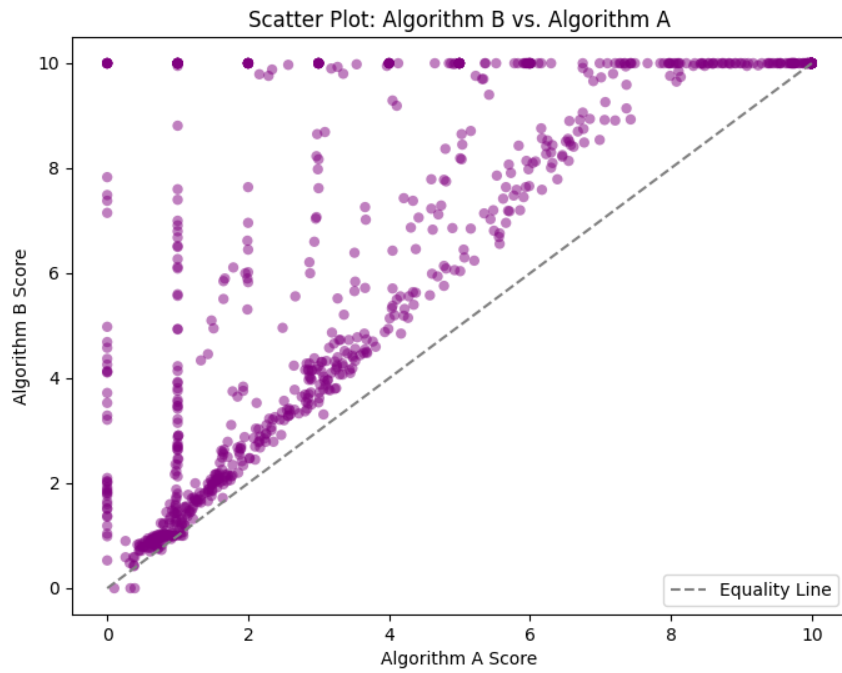


Figure 1: Comparison of Fault Tolerance Scores (FTS) between Algorithm A and Algorithm B across different parameter configurations. Higher scores indicate better fault tolerance.

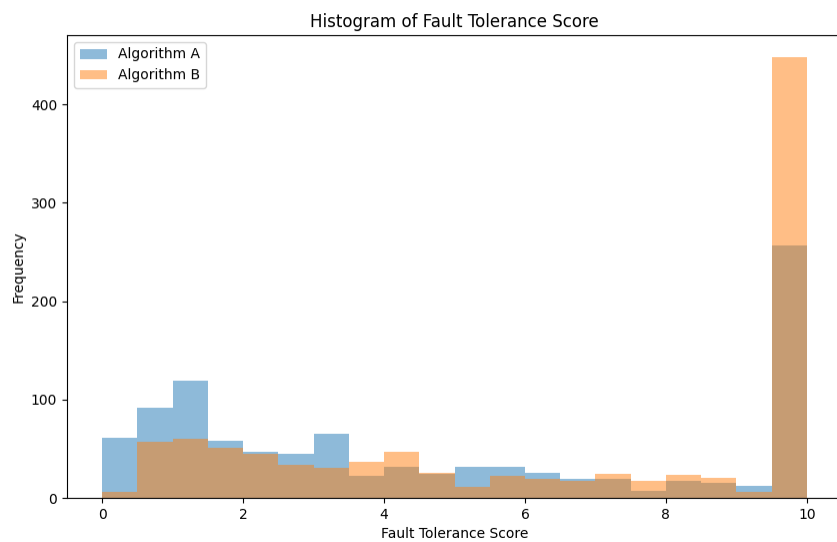


Figure 2: Comparison of FTS histograms between Algorithm A and Algorithm B. The histogram shows the distribution of FTS values for both algorithms across multiple parameter configurations.

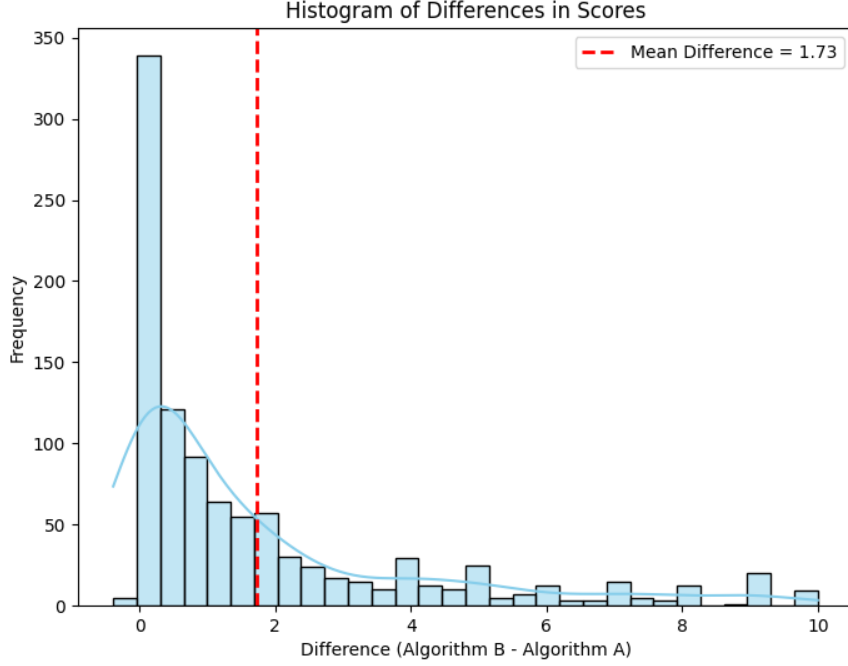


Figure 3: Histogram of the difference in FTS scores ( $\Delta\text{FTS}$ ) between Algorithm A and Algorithm B.

values for both algorithms. Finally, figure 3 shows the histogram of the difference in FTS scores ( $\Delta\text{FTS}$ ) between Algorithm A and Algorithm B. The mean of the  $\Delta\text{FTS}$  distribution is significantly greater than zero, with a p-value less than 0.05 from the bootstrap hypothesis test. This indicates that we can reject the null hypothesis  $H_0$  in favor of the alternative hypothesis  $H_1$ , indicating that Algorithm B consistently outperforms Algorithm A in terms of fault tolerance.

## 4 Conclusion

Our experimental results show the distribution of  $\Delta\text{FTS}$  values across multiple parameter configurations. As shown in figures, Algorithm B demonstrates consistently higher fault tolerance scores compared to Algorithm A. The bootstrap hypothesis test confirms that the observed differences are statistically significant, with a p-value less than 0.05. This suggests that Algorithm B is indeed more resilient in terms of fault tolerance, making it a preferable choice for PoS committee selection in Byzantine Fault Tolerance scenarios. The use of bootstrap sampling allowed us to overcome the limitations of traditional statistical tests, particularly in the presence of non-normal data distributions. By resampling the observed data and constructing an empirical null distribution, we were able to accurately assess the significance of the differences in FTS scores between the two algorithms. The results of this study provide valuable insights into the performance of PoS committee selection algorithms and highlight the importance of robust statistical methods in evaluating their effectiveness. Future work may involve further exploration of different parameter configurations, as well as the application of these methods to other PoS consensus algorithms.

## A Appendix: Code

```
1  # !/usr/bin/env python
2  # -*- coding: utf-8 -*-
3
4  import numpy as np
5  import pandas as pd
6  import scipy.stats as stats
7  import seaborn as sns
8  from matplotlib import pyplot as plt
9  from numpy import ceil, ndarray, ones, random, zeros
10 from tqdm import tqdm
11
12 from slot_alloc_sim import simulate_epoch_permissioned,
13 simulate_epoch_registered
14
15 def sample_multinomial(num_participants: int) -> ndarray:
16     """This function generates a probability vector from a flat Dirichlet
17     distribution
18     (with equal pseudo-counts for every participant)
19
20     Args:
21         num_participants: Number of participants (candidates from group)
22
23     Returns:
24         probabilities: A probability vector sampled from the Dirichlet
25     distribution
26     """
27     # Define a "flat" Dirichlet prior where every pseudo-count is 1.
28     alpha = ones(num_participants)
29
30     # Sample a probability vector from the Dirichlet distribution,
31     # which represents the multinomial probabilities for each participant.
32     return random.dirichlet(alpha)
33
34 def random_factor(n):
35     """
36     Computes all integer factors of a number and randomly selects one.
37
38     Args:
39         n (int): The number to find factors for
40
41     Returns:
42         int: A randomly selected factor of n
43     """
44     # Find all factors of n
45     factors = []
46     for i in range(1, int(n ** 0.5) + 1):
47         if n % i == 0:
48             factors.append(i)
49             if i != n // i: # Avoid adding the same factor twice for perfect
50                 squares
51                 factors.append(n // i)
52
53     # If no factors are found (should only happen if n=0), return 1
54     if not factors:
55         return 1
```

```

54
55     # Return a random factor
56     return random.choice(factors)
57
58
59 # Define the function to sample alpha from a joint distribution
60 def sample_alpha() -> dict:
61     """
62     Generates a sample allocation of seats among participants using a
multinomial
63     distribution.
64
65     The function determines the total number of seats and the number of
66     participants from predefined options. It then generates probabilities for
67     each participant and uses these probabilities to allocate seats
68     proportionally. The result is returned as a dictionary containing all the
69     generated values.
70
71     Returns:
72         dict: A dictionary containing the following keys and values:
73             - total_seats (int): Total number of available seats.
74             - num_participants (int): Number of participants.
75             - p (List[float]): Probabilities for each participant.
76             - seat_counts (List[int]): List containing the number of seats
77               allocated to each participant.
78
79     Raises:
80         AssertionError: If the internal logic fails, such as the sum of seat
81             counts exceeding the specified total seats due to an unexpected error.
82     """
83     # Consider a random number of registered seats over this range
84     num_registered = random.choice([10, 20, 40, 80, 160, 320])
85
86     # Let num_permitted seats be a random percentage of num_registered
87     num_permitted = int(ceil(random.choice([0.1, 0.3, 0.5, 0.7]) *
num_registered))
88
89     # Choose a random factor for seats_per_federated
90     seats_per_federated = random_factor(num_permitted)
91
92     # This is guaranteed now to be an integer number of federated participants
93     num_federated = num_permitted // seats_per_federated
94
95     total_seats = num_registered + num_permitted
96
97     registered_probabilities = sample_multinomial(num_registered)
98
99     max_faults = 10
100
101     sample = dict(
102         total_seats=total_seats,
103         num_registered=num_registered,
104         num_permitted=num_permitted,
105         num_federated=num_federated,
106         seats_per_federated=seats_per_federated,
107         registered_probabilities=registered_probabilities,
108         max_faults=max_faults,
109         num_epochs=100,
110     )

```

```

111     return sample
112
113
114 def faults_tolerated(committee_seats: pd.Series) -> int:
115     """
116     Compute the number of faults tolerated by the committee
117
118     Args:
119         committee_seats (pd.Series): Series of committee seats.
120
121     Returns:
122         int: The number of faults tolerated by the committee.
123     """
124     voting_strength = committee_seats.sort_values(ascending=False).divide(
125         committee_seats.sum(),
126     )
127     threshold = 1 / 3 # BFT finality risk threshold
128     faults = np.where(np.cumsum(voting_strength) > threshold)[0][0]
129     return faults
130
131
132 def simulate_epoch_federated(
133     registered_seat_counts: ndarray,
134     num_federated: int,
135     seats_per_federated: int,
136     verbose: bool = False,
137 ) -> pd.DataFrame:
138     """
139     Creates a committee with registered and federated seats.
140
141     The 'simulate_epoch_federated' function allocates committee seats between
142     two types of nodes: registered nodes
143     (whose seat allocation is proportional to their stake) and federated nodes
144     (which receive a fixed number of
145     seats). It calculates the voting strength of each group, ensures they sum
146     to 1.0, and returns a structured
147     DataFrame containing the committee composition with seat assignments for
148     both node types.
149
150     Args:
151         registered_seat_counts (ndarray): seat counts for registered
152         participants.
153         num_federated (int): Number of federated nodes.
154         seats_per_federated (int): Number of seats per federated node.
155         verbose (bool): Whether to print detailed information.
156         Default is False.
157
158     Returns:
159         pd.DataFrame: Committee seats information with kind and seat
160         assignments.
161     """
162     # Cast the seat counts to a pandas Series
163     registered_seat_counts = pd.Series(registered_seat_counts, name="
    registered seats")
164
165     # Get the number of registered participants (SPOs)
166     num_registered = registered_seat_counts.shape[0]
167
168     # Calculate the number of federated seats in the committee

```



```

163     federated_seats = seats_per_federated * num_federated
164
165     # Committee size is the total number of registered and federated seats
166     registered_seats = registered_seat_counts.sum()
167     committee_size = registered_seats + federated_seats
168
169     # Calculate the voting strength of the registered seats
170     registered_voting_strength = registered_seats / committee_size
171
172     # Calculate the voting strength of the federated seats
173     federated_voting_strength = federated_seats / committee_size
174
175     # Assert that the voting strengths sum to 1.0
176     assert (
177         registered_voting_strength + federated_voting_strength == 1.0
178     ), "Voting strength does not sum to 1.0"
179
180     # Create a series for the federated seats on the committee
181     federated_seat_counts = pd.Series(
182         ones(num_federated, dtype="int64") * seats_per_federated,
183         index=[str(i) for i in range(num_registered + 1, num_registered +
num_federated + 1)],
184         dtype="int64",
185         name="federated seats",
186     )
187
188     # Combine the federated and registered seats into a single DataFrame
189     committee_seats = (
190         pd.concat(
191             [federated_seat_counts, registered_seat_counts],
192             keys=["federated", "registered"],
193             names=["kind", "index"],
194             ignore_index=False,
195         )
196         .reset_index()
197         .rename(columns={0: "seats"})
198         .set_index("index")
199         .sort_values(by=["seats", "kind"], ascending=[False, True])
200     )
201     if verbose:
202         print(
203             f"Committee size .... = {committee_size} \n"
204             "-----\n"
205             "Registered:\n"
206             f"Number registered.. = {num_registered} \n"
207             f"Number of seats.... = {registered_seats} \n"
208             f"Voting strength.... = {registered_voting_strength:.2%} \n"
209             "-----\n"
210             "Federated:\n"
211             f"Number federated... = {num_federated} \n"
212             f"Seats per federated = {seats_per_federated} \n"
213             f"Number of seats.... = {federated_seats} \n"
214             f"Voting strength.... = {federated_voting_strength:.2%}\n",
215         )
216     return committee_seats
217
218
219 def simulate_committee_federated(
220     registered_probabilities: ndarray,

```

```

221     num_registered: int,
222     num_federated: int,
223     seats_per_federated: int,
224     num_epochs: int,
225     **kwargs,
226 ) -> pd.DataFrame:
227     """
228     Simulate the allocation of committee seats over multiple epochs.
229
230     This function simulates the process of assigning committee seats for a
231     specified
232     number of epochs using the simulate_epoch_federated function. It collects
233     seat
234     assignments across all epochs and returns them as a DataFrame.
235
236     Args:
237         registered_probabilities (ndarray): Probabilities for registered
238         participants.
239         num_registered (int): Number of registered participants.
240         num_federated (int): Number of federated nodes.
241         seats_per_federated (int): Number of seats per federated node.
242         num_epochs (int): Number of epochs to simulate.
243         kwargs: Additional keyword arguments.
244
245     Returns:
246         pd.DataFrame: DataFrame containing committee seat assignments for all
247         epochs,
248         with epochs as rows and committee members as columns.
249     """
250     committee_list = []
251     for epoch in range(num_epochs):
252         # Sample the committee seat configuration from the multinomial
253         # distribution
254         # parameterized by registered_probabilities.
255         registered_seat_counts = random.multinomial(num_registered,
256         registered_probabilities)
257
258         # Assign committee seats for the current epoch
259         committee = simulate_epoch_federated(
260             registered_seat_counts=registered_seat_counts,
261             num_federated=num_federated,
262             seats_per_federated=seats_per_federated,
263             )
264         committee_list.append(committee.seats)
265
266     assert len(committee_list) == num_epochs, "Number of epochs does not match
267     ."
268
269     committee_seats = pd.concat(committee_list, keys=range(num_epochs), axis
270     =1).T
271
272     return committee_seats
273
274 def simulate_proposed(
275     registered_probabilities: ndarray,
276     num_registered: int,
277     num_permitted: int,
278     num_epochs: int,

```

```

272     **kwargs,
273     ) -> pd.DataFrame:
274     """Runs the proposed algorithm simulation. For each epoch, simulates both
the
275     registered and permissioned candidate selection. Stores the number of
slots
276     each candidate gets and returns the cumulative results.
277
278     Args:
279         registered_probabilities (ndarray): Probabilities for registered
candidates.
280         num_registered (int): Number of registered candidates.
281         num_permitted (int): Number of permitted (federated) nodes.
282         num_epochs (int): Number of epochs to simulate.
283         kwargs: Additional keyword arguments.
284
285     Returns:
286         pd.DataFrame: DataFrame containing committee seat assignments for all
epochs,
287         with epochs as rows and committee members as columns.
288     """
289     permissioned_candidates = [f"P{i}" for i in range(num_permitted)]
290     permissioned_results = {name: [] for name in permissioned_candidates}
291     registered_results = {i: [] for i in range(num_registered)}
292
293     for _ in range(num_epochs):
294         registered_candidates = {i: registered_probabilities[i] for i in
registered_results.keys()}
295
296         # Simulate this epoch
297         reg_assign = simulate_epoch_registered(registered_candidates,
num_registered)
298         perm_assign = simulate_epoch_permissioned(permissioned_candidates,
num_permitted)
299
300         # Collect results
301         for name in registered_candidates:
302             registered_results[name].append(reg_assign[name])
303         for name in permissioned_candidates:
304             permissioned_results[name].append(perm_assign[name])
305
306         # Combine the two dictionaries
307         combined = registered_results.copy()
308         combined.update(permissioned_results)
309
310         # Convert to DataFrame
311         committee_seats = (
312             pd.DataFrame.from_dict(combined, orient="index")
313             .astype(int)
314             .fillna(0)
315             .transpose()
316         )
317         return committee_seats
318
319
320 def calculate_fault_tolerance_probability(
321     committee_seats: pd.DataFrame,
322     fault_tolerance: int = 1,
323     ) -> float:

```

```

324     """
325     Calculate the probability of tolerating a given number of faults
326     in a committee of a given size.
327     The function uses Monte Carlo simulation to estimate the probability.
328     Args:
329         committee_seats (pd.DataFrame): DataFrame of committee seat
330             assignments of both registered and permissioned members
331             for each epoch.
332         fault_tolerance (int): Number of faults to tolerate.
333
334     Returns:
335         float: Estimated probability of tolerating the given number of faults.
336     """
337     if fault_tolerance == 0:
338         probability = 1.0
339     else:
340         # Calculate success rate through Monte Carlo simulation
341         probability = (
342             committee_seats.apply(faults_tolerated, axis=1) >= fault_tolerance
343         ).mean()
344     return probability
345
346
347 # Stub functions for Algorithm A and B
348 def algorithm(function, **params) -> float:
349     """
350     Computes fault tolerance probabilities using an algorithm that processes a
351     defined function and its
352     parameters to determine committee seats and calculate probabilities for a
353     fault tolerance metric.
354
355     Args:
356         function (Callable): A function that determines committee seats based
357         on provided parameters.
358         **params: A set of keyword arguments to be passed into the function '
359         function'. Must include
360         'max_faults' which defines the maximum faults for calculating
361         probabilities; defaults to 5
362         if not specified.
363
364     Returns:
365         mean_fault_tolerance: (float) The mean fault tolerance probability for
366         the algorithm.
367     """
368     committee_seats = function(**params)
369
370     # Calculate fault tolerance probabilities for the algorithm
371     faults = np.arange(1, params.get("max_faults", 5) + 1)
372     ftprob = pd.Series(0.0, index=faults, name="probability")
373     for f in faults:
374         # p is the probability of tolerating f faults
375         p = calculate_fault_tolerance_probability(committee_seats,
376 fault_tolerance=f)
377         ftprob.loc[f] = p
378         if p == 0: # since the rest of the series will be zero as well
379             break
380
381     # Compute the performance score as the sum of probabilities
382     score = ftprob.sum()

```

```

376
377     return score
378
379
380 def hypothesis_test(
381     results_a: ndarray,
382     results_b: ndarray,
383     kind: str = "bootstrap",
384     n_bootstrap: int = 10000,
385 ):
386     """
387     Conducts a hypothesis test to compare two result sets, 'results_a' and '
388     results_b', based on the specified
389     kind of statistical testing approach. By default, a bootstrap approach is
390     used to calculate the one-sided
391     p-value for the null hypothesis. Alternatively, a paired t-test can be
392     performed. The results help determine
393     whether algorithm B (corresponding to 'results_b') is significantly better
394     than algorithm A (corresponding
395     to 'results_a').
396
397     Args:
398         results_a: An array of numerical results for algorithm A.
399         results_b: An array of numerical results for algorithm B.
400         kind: The kind of statistical test to perform. Options are "bootstrap"
401         (default) or "paired_t".
402         n_bootstrap: The number of bootstrap resampling iterations to perform.
403         Only applicable
404         if the 'kind' is "bootstrap".
405
406     Raises:
407         ValueError: If 'kind' is not "bootstrap" or "paired_t".
408         RuntimeError: If inconsistencies in data prevent calculations, such as
409         mismatched array lengths
410         or insufficient data.
411     """
412     if kind == "bootstrap":
413         # Calculate the observed mean difference.
414         differences = results_b - results_a
415         obs_diff = np.mean(differences)
416
417         # For a one-sided test (B > A), if the observed mean difference is non
418         -positive,
419         # the p-value would be 1 (or you might decide not to proceed).
420         if obs_diff <= 0:
421             p_value = 1.0
422         else:
423             # Center the differences: adjust the data so that its mean becomes
424             0
425             # (the null hypothesis condition)
426             differences_centered = differences - obs_diff
427
428             # Bootstrap resampling from the centered differences.
429             bootstrap_means = np.empty(n_bootstrap)
430             for i in range(n_bootstrap):
431                 sample = np.random.choice(differences_centered, size=len(
432                 differences_centered), replace=True)
433                 bootstrap_means[i] = np.mean(sample)

```

```

425         # One-sided p-value: proportion of bootstrap means greater than or
         equal to the observed difference.
426         p_value = np.mean(bootstrap_means >= obs_diff)
427
428         print(f"Observed Mean Difference (B - A): {obs_diff:.3f}")
429         print(f"One-sided p-value: {p_value:.4f}")
430
431     else:
432         # Perform a paired t-test to determine if B is significantly better
than A
433         t_stat, p_value = stats.ttest_rel(results_a, results_b, alternative='
less')
434
435         print(f"Mean mA: {results_a.mean():.3f}, Mean mB: {results_b.mean():.3
f}")
436         print(f"Paired t-test result: t-statistic = {t_stat:.3f}, p-value = {
p_value:.3f}")
437
438         if p_value < 0.05:
439             print("Conclusion: Algorithm B is significantly better than Algorithm
A.")
440         else:
441             print("Conclusion: No statistically significant difference between A
and B.")
442
443
444     def monte_carlo_simulation(num_trials: int = 100) -> tuple[np.ndarray, np.
ndarray]:
445         """
446         Run the Monte Carlo simulation to compare Algorithms A and B across random
scenarios.
447
448         Args:
449             num_trials: (int) The number of trials to run the simulation for.
Defaults to 50.
450
451         Returns:
452             tuple[ndarray[Any, float], ndarray[Any, float]]: results_a, results_b
453         """
454         a = simulate_committee_federated
455         b = simulate_proposed
456
457         results_a = zeros(num_trials)
458         results_b = zeros(num_trials)
459
460         for i in tqdm(range(num_trials), desc="simulation trials"):
461             alpha = sample_alpha()
462
463             results_a[i] = algorithm(a, **alpha)
464             results_b[i] = algorithm(b, **alpha)
465
466         hypothesis_test(results_a, results_b, kind="bootstrap")
467
468         return results_a, results_b
469
470
471     def visualize_algorithm_comparison(results_a, results_b, show_plots=True):
472         """
473         Visualize the comparison between two algorithms through various plots.

```

```

474
475     This function creates three visualization plots:
476     1. Histogram of scores for both algorithms
477     2. Histogram of differences between algorithm scores
478     3. Scatter plot of paired algorithm results
479
480     Args:
481         results_a (ndarray): Array of scores from Algorithm A
482         results_b (ndarray): Array of scores from Algorithm B
483         show_plots (bool): Whether to display the plots immediately. Default
484         is True.
485
486     Returns:
487         tuple: Three matplotlib figure objects (hist_fig, diff_fig,
488         scatter_fig)
489         """
490         # Calculate differences
491         differences = results_b - results_a # difference: B - A
492
493         # -----
494         # 1. Scatter Plot of Paired Data
495         # -----
496         scatter_fig = plt.figure(figsize=(8, 6))
497         plt.scatter(results_a, results_b, alpha=0.5, color='purple', edgecolor='
498         none')
499         plt.xlabel("Algorithm A Score")
500         plt.ylabel("Algorithm B Score")
501         plt.title("Scatter Plot: Algorithm B vs. Algorithm A")
502         plt.plot(
503             [min(results_a), max(results_a)], [min(results_a), max(results_a)],
504             linestyle='--', color='gray', label='Equality Line',
505             )
506         plt.legend()
507
508         # -----
509         # 2. Histogram each Score
510         # -----
511         hist_fig = plt.figure(figsize=(10, 6))
512         plt.hist(results_a, bins=20, label="Algorithm A", alpha=0.5)
513         plt.hist(results_b, bins=20, label="Algorithm B", alpha=0.5)
514         plt.xlabel("Fault Tolerance Score")
515         plt.ylabel("Frequency")
516         plt.title("Histogram of Fault Tolerance Score")
517         plt.legend()
518
519         # -----
520         # 3. Histogram of Differences
521         # -----
522         diff_fig = plt.figure(figsize=(8, 6))
523         sns.histplot(differences, bins=30, kde=True, color='skyblue', edgecolor='
524         black')
525         plt.xlabel("Difference (Algorithm B - Algorithm A)")
526         plt.ylabel("Frequency")
527         plt.title("Histogram of Differences in Scores")
528         plt.axvline(
529             np.mean(differences), color='red', linestyle='dashed', linewidth=2,
530             label=f'Mean Difference = {np.mean(differences):.2f}',
531             )
532         plt.legend()

```

```
529     if show_plots:
530         plt.show()
531
532     return hist_fig, diff_fig, scatter_fig
533
534
535 if __name__ == "__main__":
536     # Run the simulation
537     results_a, results_b = monte_carlo_simulation(num_trials=1000)
538
539     visualize_algorithm_comparison(results_a, results_b)
```

Listing 1: Python code for the Monte Carlo simulation experiment