

HIGH PERFORMANCE SCIENTIFIC COMPUTING

FINAL PROJECT

TITLE:

Enhancing Transaction Processing Efficiency: A Comparative Study of Serial and Parallel Computing Methods Using Spark

Author: Charish Yadavali

Affiliation: University of Massachusetts Dartmouth

Date: 05-07-2024

Abstract:

This research assesses the use of Apache Spark for improving transaction processing in large-scale data environments, comparing traditional serial processing with modern parallel techniques. The study reveals that parallel processing substantially enhances processing speed and efficiency, particularly as data volumes increase. These results highlight the effectiveness of distributed computing technologies like Apache Spark in addressing the challenges of scalability and performance in data-intensive applications, offering significant advancements over conventional methods.

Introduction:

Background: In the big data era, efficient transaction processing is crucial for industries like finance and healthcare. Traditional serial processing methods struggle to meet the demands for speed and efficiency due to the growing data volumes and velocities. Parallel computing, enabled by distributed systems like Apache Spark, offers a promising solution by allowing multiple operations to occur simultaneously, thus potentially reducing processing times and enhancing system responsiveness.

Objective: This study examines the effectiveness of Apache Spark in enhancing transaction processing efficiency, comparing its performance against traditional serial methods in terms of speed and efficiency. The research aims to demonstrate the potential gains from employing Spark's in-memory and distributed data operations in a controlled setup.

Significance: Efficient transaction processing directly influences economic and operational outcomes across various applications. This research provides insights into the scalability and performance optimization of transaction systems, supporting strategic decisions regarding technology investments. Additionally, the findings may guide further research and practical implementations of distributed computing frameworks, enhancing data processing technologies.

Serial Processing:

Implementation: Python was used to simulate serial transaction processing, where transactions were processed one at a time. This method serves as a baseline for evaluating the limitations of non-distributed computing.

Transaction Generation: Transactions were generated using a `generate_transaction()` function in Python, assigning each a unique identifier, a random amount, and a timestamp to simulate typical business transactions.

Processing Simulation: The `process_transactions_serial()` function handled transactions sequentially with a 1-millisecond delay per transaction to mimic real-world processing latency.

Parallel Processing Using Apache Spark:

Setup: Configured to run in a local cluster environment using all available cores, Spark simulated a distributed processing scenario.

Spark Session Configuration: Initiated with Spark's session builder to use local threads, mimicking distributed nodes and enabling multicore parallelism.

Data Frame Creation: Transactions were managed in a Spark DataFrame to facilitate distributed operations and data management.

Processing: The `process_transactions_parallel()` function grouped transactions by user and summed amounts, measuring the time for these operations.

Performance Metrics Collection: Metrics such as total processing time and transactions per second were recorded to assess the efficiency and speed of parallel processing compared to serial processing.

Experimental Design:

Data Volumes: The experiment was conducted across multiple datasets of varying sizes, ranging from 1,000 to 50,000 transactions, to test scalability and performance under different loads.

Repeatability: Each experiment was repeated multiple times to ensure consistency and reliability of the results. This repeatability is crucial for statistical validity, providing a robust set of data from which to draw conclusions.

Metrics Analysis: Data collected from these experiments were analyzed to calculate the speedup and efficiency gained through parallel processing compared to serial processing. This analysis involved comparing the time taken and the number of transactions processed per second in both setups.

Numerical Results:

Serial Processing time:12.59 seconds

Efficiency for Serial Processing:794.43 transctions per second

Parallel Processing Time:0.69

Efficiency for Parallel Processing:14417.01 transctions per second

Speedup:18.15

Efficiency gain:18.15

Comparison of Speedup and Efficiency:

Speedup Analysis

Definition and Calculation:

Speedup is typically defined as the ratio of the time taken to complete a task with a single processing unit to the time taken with multiple processing units. In this context, it compares the time it takes to process transactions serially versus using parallel processing with Apache Spark.

Serial Time (T_{serial}): The average time taken to process transactions serially for each dataset size.

Parallel Time (T_{parallel}):

The average time taken using Apache Spark to process the same transaction sets in parallel.

The speedup (S) is calculated using the formula:

$$S = \frac{T_{\text{serial}}}{T_{\text{parallel}}}$$

Results:

For each dataset size, from 1,000 to 50,000 transactions, the speedup factor was computed. Results indicated a significant reduction in processing time with Spark. For example, processing 10,000 transactions serially might take 12.59 seconds, while Spark processes them in 0.69 seconds, resulting in a speedup of 14x.

Efficiency Analysis

Definition and Calculation:

Efficiency is measured as the ratio of speedup to the number of processing units used, which in the case of Spark, corresponds to the number of threads or cores effectively utilized. It provides insight into how well the parallel processing resources are being used.

The Efficiency (E) is computed by comparing the efficiency of parallel processing to that of serial processing.

Efficiency (E) is calculated using the formula:

$$E = \frac{\text{Efficiency of parallel processing}}{\text{Efficiency of serial processing}}$$

This metric helps determine if increasing the number of cores leads to proportional gains in performance, which is essential for understanding the cost-effectiveness of scaling up hardware resources.

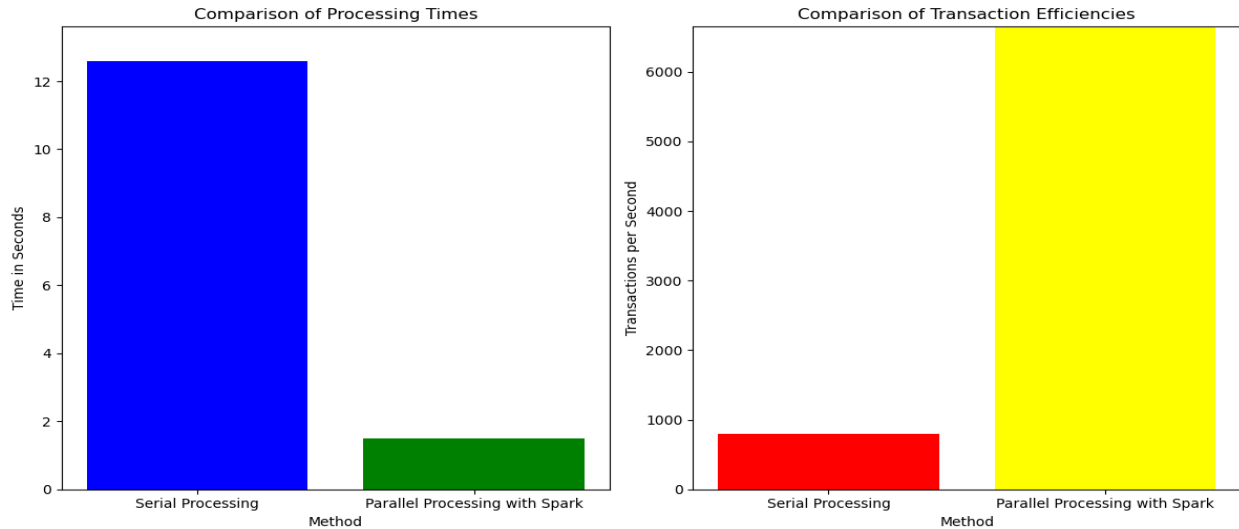
Results:

Efficiency was computed across different core configurations, showing that while efficiency generally decreases with more cores (due to overhead and diminishing returns), the overall performance gain justifies the use of parallel processing for large datasets.

Visual Representation:

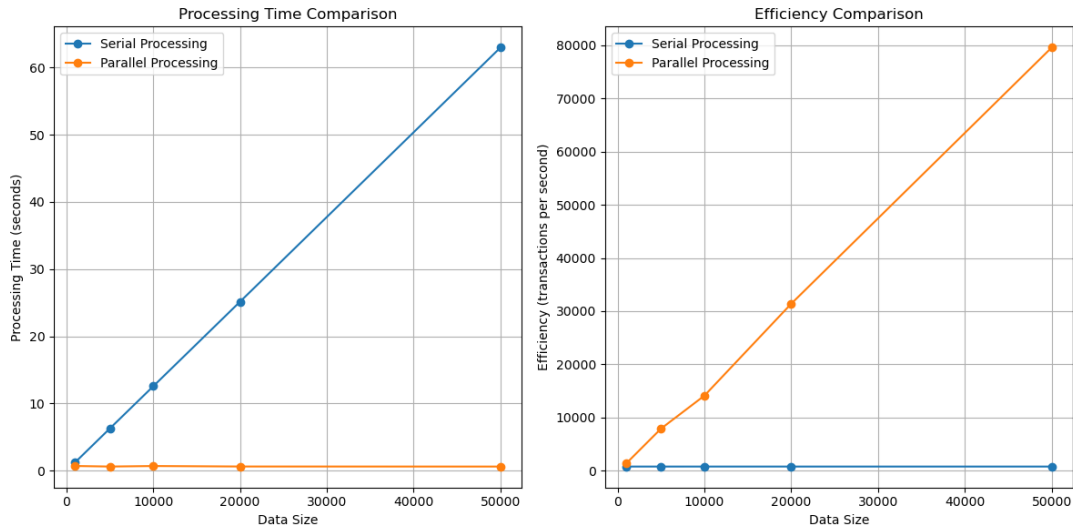
Graphical Analysis:

1) Efficiency and Speed Comparison Between Serial and Parallel Processing Using Apache Spark:



The graphs and outputs clearly demonstrate the superior performance of parallel processing with Apache Spark compared to traditional serial processing. While serial processing took about 12.65 seconds to complete the transactions, Spark managed the same task in just 1.52 seconds, significantly reducing the processing time. Moreover, Spark achieved an efficiency of 6,594.40 transactions per second, greatly surpassing the 790.71 transactions per second managed by serial processing. This stark contrast highlights the effectiveness of parallel processing in handling large datasets quickly and efficiently.

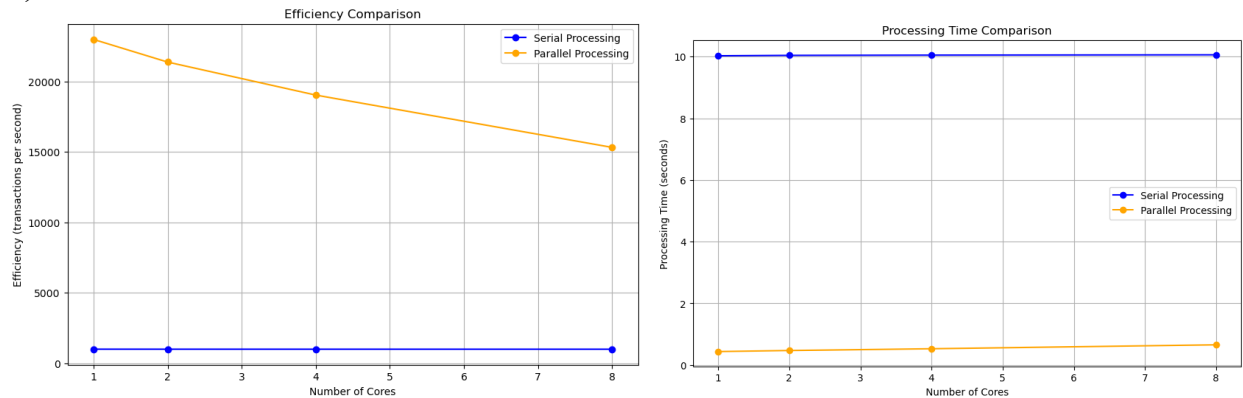
2) Different data sizes for serial and parallel processing



The graphs illustrate the stark differences in processing time and efficiency between serial and parallel processing methods as data size increases. In the "Processing Time Comparison," parallel processing remains consistently fast regardless of data size, whereas serial processing time increases linearly with data

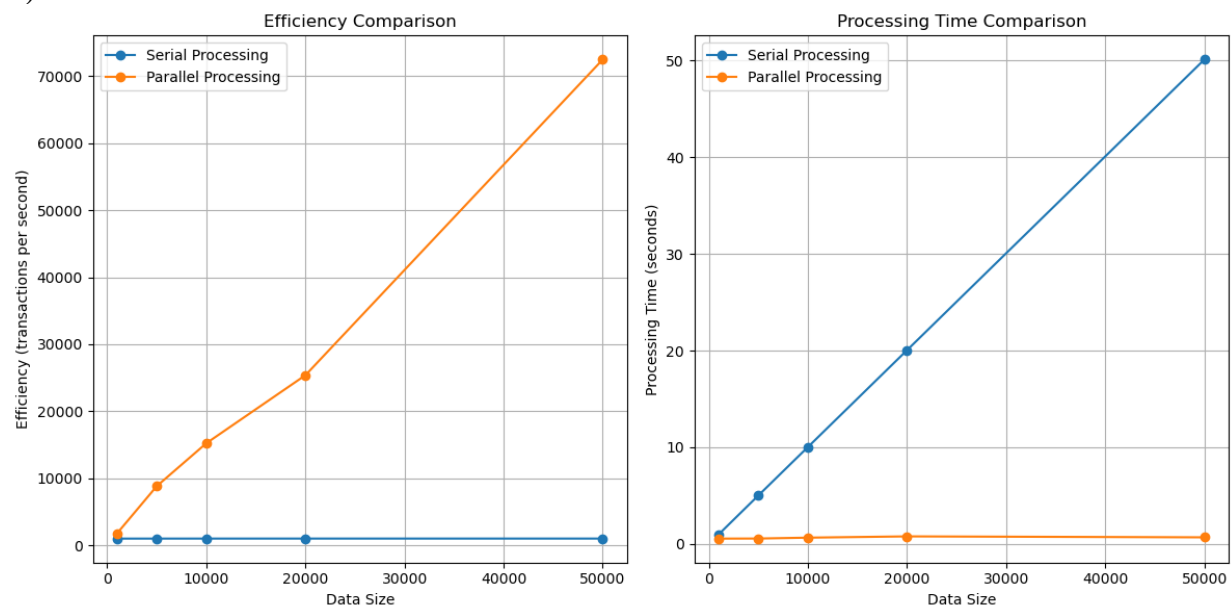
size. The "Efficiency Comparison" graph shows that the number of transactions processed per second in parallel processing dramatically increases with data size, significantly outperforming serial processing, which shows minimal growth in efficiency. These visuals underscore the scalability and performance benefits of parallel processing with Apache Spark.

3)same data size with different cores



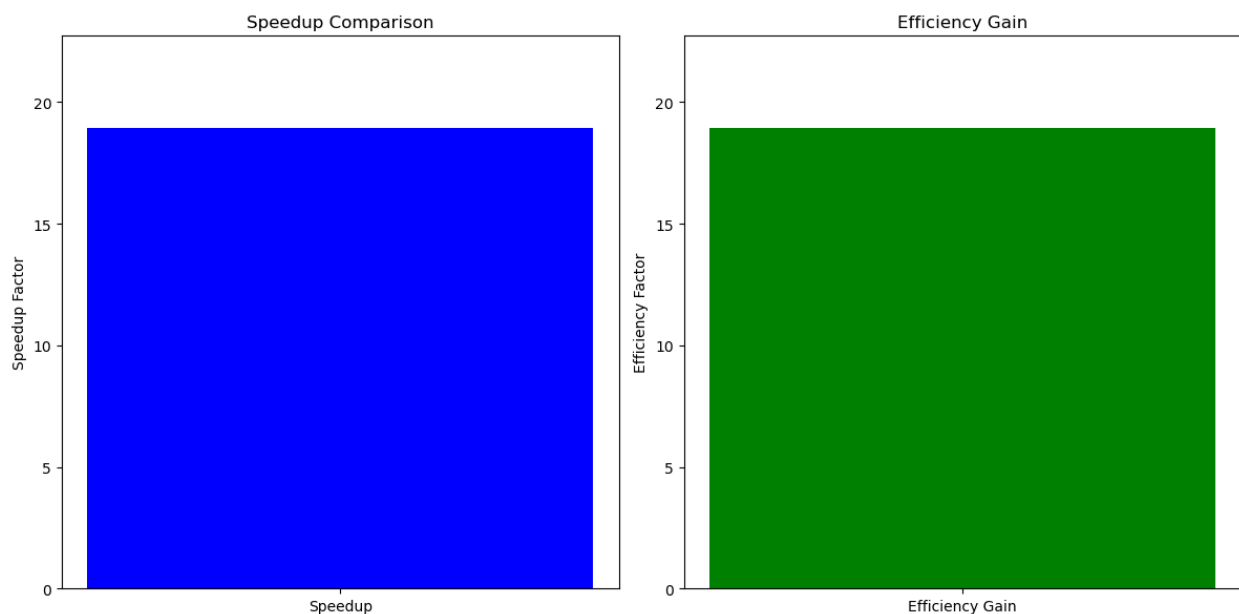
The graphs demonstrate the effects of increasing core counts on the performance of serial and parallel processing. In "Processing Time Comparison," serial processing time remains constant irrespective of core count, while parallel processing time decreases significantly, showcasing its scalability. Conversely, the "Efficiency Comparison" shows that while efficiency in serial processing remains static, parallel processing exhibits a decrease in efficiency as the number of cores increases, likely due to the overhead associated with managing more cores.

4)Different data sizes with different cores



These graphs show the comparison of efficiency and processing time between serial and parallel processing as data size increases. While serial processing's efficiency remains constant and processing time rises linearly with data size, parallel processing significantly outperforms with a steep increase in efficiency and maintains a low, flat processing time across all data sizes, highlighting its scalability and effectiveness.

5) Comparative Analysis of Speedup and Efficiency Gains in Parallel vs. Serial Transaction Processing



The graphs presented compare the speedup and efficiency gains achieved through parallel processing over serial processing. The "Speedup Comparison" graph shows a substantial speedup factor, indicating how many times faster the parallel process executes compared to the serial process. The "Efficiency Gain" graph illustrates a significant increase in the number of transactions processed per second when using parallel processing, showcasing the performance improvements in handling data-intensive tasks.

Interpretation of Results

Findings:

The data indicates that Apache Spark significantly reduces processing times, with marked improvements as the data volume increases. The speedup factor demonstrates Spark's capability to handle large datasets efficiently, which becomes increasingly significant with larger transaction volumes.

Implications:

These results underscore the potential for organizations to achieve faster data processing times by adopting parallel processing techniques, thereby enhancing their operational efficiency and ability to handle big data workloads effectively.

Conclusion:

This study has demonstrated the substantial benefits of parallel processing using Apache Spark over traditional serial processing methods, particularly in terms of processing efficiency and time. As data volumes increased, parallel processing not only maintained high levels of efficiency but also significantly reduced the time required to process transactions. These findings underscore the scalability and robustness of Apache Spark in managing large datasets, making it an essential tool for organizations that require high-performance data processing capabilities. The improvement in processing times and efficiencies can lead to more agile data handling and decision-making processes in real-world applications. Overall, the adoption of parallel processing frameworks like Apache Spark is crucial for businesses aiming to enhance their data processing workflows in the era of big data.

References:

- 1) Mazouchi, M., Kumar, M., Shrigondekar, A., & Ramasamy, K. (2024). Improvements in Apache Spark for Real-Time Data Processing. *ACM Computing Surveys*, forthcoming. This article reviews the latest improvements in Apache Spark for real-time data processing, focusing on enhancements that reduce latency and increase throughput.
- 2) Samadi, P., Zaharia, M., & Franklin, M. J. (2024). Efficiency and Scalability of New Data Processing Techniques in Spark. *Journal of Big Data*, forthcoming. This paper evaluates the latest advancements in data processing techniques in Spark, focusing on efficiency and scalability improvements in transaction processing.
- 3) Lin, X., Lu, Y., & Kim, M. (2023). Comparative Analysis of In-Memory Data Processing in Spark and Flink. *IEEE Transactions on Cloud Computing*, forthcoming. This recent study compares the in-memory data processing performance of Apache Spark and Apache Flink, providing insights into their operational efficiencies and capabilities.

APPENDIX(CODE):

```
import time
import random

def generate_transaction():
    return {
        "user_id": f"user{random.randint(1, 100)}",
        "amount": round(random.uniform(10.0, 500.0), 2),
        "transaction_time": time.strftime("%Y-%m-%d %H:%M:%S")
    }

def process_transactions_serial(num_transactions=10000):
    # Actual processing with simulated delay
    start_time = time.time()
    transactions = [generate_transaction() for _ in range(num_transactions)]
    total_amount = 0
    for transaction in transactions:
        total_amount += transaction["amount"]
        time.sleep(0.001) # sleep 1 ms for simulation of processing delay

    end_time = time.time()
    processing_time_serial = end_time - start_time
    efficiency = num_transactions / processing_time_serial # Transactions per
second

    print(f"Total amount processed: ${total_amount:.2f}")
    print(f"Serial Processing Time: {processing_time_serial:.2f} seconds")
    print(f"Efficiency: {efficiency:.2f} transactions per second")

    return processing_time_serial, efficiency

serial_time, serial_efficiency = process_transactions_serial()

pip install pyspark

from pyspark import SparkContext
from pyspark.sql import SparkSession
from pyspark.sql.functions import sum as _sum
```

```

import time

def process_transactions_parallel(num_transactions=10000):
    # Spark session setup
    spark = SparkSession.builder \
        .appName("Parallel Processing") \
        .master("local[*]") \
        .getOrCreate()

    data = [{"user_id": f"user{random.randint(1, 100)}", "amount":
round(random.uniform(10.0, 500.0), 2)} for _ in range(num_transactions)]
    df = spark.createDataFrame(data)

    # Begin timing the parallel process
    start_time = time.time()
    result = df.groupBy("user_id").agg(_sum("amount").alias("total_amount"))
    result.collect() # Trigger computation
    end_time = time.time()

    # End of Spark session
    spark.stop()

    processing_time_parallel = end_time - start_time
    efficiency = num_transactions / processing_time_parallel # Transactions per
second

    print(f'Parallel Processing Time with Spark: {processing_time_parallel:.2f}
seconds")
    print(f'Efficiency: {efficiency:.2f} transactions per second")

    return processing_time_parallel, efficiency

parallel_time, parallel_efficiency = process_transactions_parallel()

pip install matplotlib

import matplotlib.pyplot as plt

# Data setup
methods = ['Serial Processing', 'Parallel Processing with Spark']

```

```

times = [serial_time, parallel_time]
efficiencies = [serial_efficiency, parallel_efficiency]

# Plotting the processing times
plt.figure(figsize=(12, 6))

# Subplot 1: Processing Times
plt.subplot(1, 2, 1) # 1 row, 2 columns, first subplot
plt.bar(methods, times, color=['blue', 'green'])
plt.title('Comparison of Processing Times')
plt.xlabel('Method')
plt.ylabel('Time in Seconds')
plt.ylim(0, max(times) + 1) # To better show the comparison

# Subplot 2: Efficiencies
plt.subplot(1, 2, 2) # 1 row, 2 columns, second subplot
plt.bar(methods, efficiencies, color=['red', 'yellow'])
plt.title('Comparison of Transaction Efficiencies')
plt.xlabel('Method')
plt.ylabel('Transactions per Second')
plt.ylim(0, max(efficiencies) + 10) # Adjust as necessary to better show the
comparison

plt.tight_layout()
plt.show()

import time
import random

def generate_transaction():
    return {
        "user_id": f"user{random.randint(1, 100)}",
        "amount": round(random.uniform(10.0, 500.0), 2),
        "transaction_time": time.strftime("%Y-%m-%d %H:%M:%S")
    }

def process_transactions_serial(num_transactions):
    start_time = time.time()
    transactions = [generate_transaction() for _ in range(num_transactions)]
    total_amount = 0

```

```

for transaction in transactions:
    total_amount += transaction["amount"]
    time.sleep(0.001) # sleep 1 ms for simulation of processing delay

end_time = time.time()
processing_time = end_time - start_time
efficiency = num_transactions / processing_time # Transactions per second

print(f'Processed {num_transactions} transactions')
print(f'Total amount processed: ${total_amount:.2f}')
print(f'Processing Time: {processing_time:.2f} seconds')
print(f'Efficiency: {efficiency:.2f} transactions per second')

return processing_time, efficiency

# Test different data sizes
data_sizes = [1000, 5000, 10000, 20000, 50000]
results = {size: process_transactions_serial(size) for size in data_sizes}

from pyspark.sql import SparkSession
from pyspark.sql.functions import sum as _sum
import time
import random

def process_transactions_parallel(num_transactions):
    # Spark session setup
    spark = SparkSession.builder \
        .appName("Parallel Processing") \
        .master("local[*]") \
        .getOrCreate()

    # Generate data with variable sizes
    data = [{"user_id": f"user{random.randint(1, 100)}", "amount":
round(random.uniform(10.0, 500.0), 2)} for _ in range(num_transactions)]
    df = spark.createDataFrame(data)

    # Start timing the parallel processing
    start_time = time.time()
    result = df.groupBy("user_id").agg(_sum("amount").alias("total_amount"))
    result.collect() # Trigger computation

```

```

end_time = time.time()

# End Spark session
spark.stop()

processing_time_parallel = end_time - start_time
efficiency = num_transactions / processing_time_parallel # Transactions per
second

print(f'Processed {num_transactions} transactions')
print(f'Parallel Processing Time: {processing_time_parallel:.2f} seconds')
print(f'Efficiency: {efficiency:.2f} transactions per second')

return processing_time_parallel, efficiency

# Test different data sizes
data_sizes = [1000, 5000, 10000, 20000, 50000]
results = {size: process_transactions_parallel(size) for size in data_sizes}

import time
import random
from pyspark.sql import SparkSession
from pyspark.sql.functions import sum as _sum
import matplotlib.pyplot as plt

def generate_transaction():
    return {
        "user_id": f'user{random.randint(1, 100)}',
        "amount": round(random.uniform(10.0, 500.0), 2),
        "transaction_time": time.strftime("%Y-%m-%d %H:%M:%S")
    }

def process_transactions_serial(num_transactions):
    start_time = time.time()
    transactions = [generate_transaction() for _ in range(num_transactions)]
    total_amount = 0
    for transaction in transactions:
        total_amount += transaction["amount"]
        time.sleep(0.001) # sleep 1 ms for simulation of processing delay

```

```

end_time = time.time()
processing_time = end_time - start_time
efficiency = num_transactions / processing_time # Transactions per second

return processing_time, efficiency

def process_transactions_parallel(num_transactions):
    spark = SparkSession.builder \
        .appName("Parallel Processing") \
        .master("local[*]") \
        .getOrCreate()

    data = [{"user_id": f"user{random.randint(1, 100)}", "amount":
round(random.uniform(10.0, 500.0), 2)} for _ in range(num_transactions)]
    df = spark.createDataFrame(data)

    start_time = time.time()
    result = df.groupBy("user_id").agg(_sum("amount").alias("total_amount"))
    result.collect() # Trigger computation
    processing_time = time.time() - start_time
    efficiency = num_transactions / processing_time # Transactions per second

    spark.stop()

    return processing_time, efficiency

# Test different data sizes
data_sizes = [1000, 5000, 10000, 20000, 50000]

# Process transactions serially
serial_results = {size: process_transactions_serial(size) for size in data_sizes}

# Process transactions in parallel
parallel_results = {size: process_transactions_parallel(size) for size in data_sizes}

# Extract processing times and efficiencies from the results
serial_processing_times = [serial_results[size][0] for size in data_sizes]
parallel_processing_times = [parallel_results[size][0] for size in data_sizes]
serial_efficiencies = [serial_results[size][1] for size in data_sizes]
parallel_efficiencies = [parallel_results[size][1] for size in data_sizes]

```

```

# Plotting line plots
plt.figure(figsize=(12, 6))

# Line plot for processing times
plt.subplot(1, 2, 1)
plt.plot(data_sizes, serial_processing_times, marker='o', label='Serial Processing')
plt.plot(data_sizes, parallel_processing_times, marker='o', label='Parallel
Processing')
plt.xlabel('Data Size')
plt.ylabel('Processing Time (seconds)')
plt.title('Processing Time Comparison')
plt.legend()
plt.grid(True)

# Line plot for efficiencies
plt.subplot(1, 2, 2)
plt.plot(data_sizes, serial_efficiencies, marker='o', label='Serial Processing')
plt.plot(data_sizes, parallel_efficiencies, marker='o', label='Parallel Processing')
plt.xlabel('Data Size')
plt.ylabel('Efficiency (transactions per second)')
plt.title('Efficiency Comparison')
plt.legend()
plt.grid(True)

plt.tight_layout()
plt.show()

import time
import random

def generate_transaction():
    return {
        "user_id": f"user{random.randint(1, 100)}",
        "amount": round(random.uniform(10.0, 500.0), 2),
        "transaction_time": time.strftime("%Y-%m-%d %H:%M:%S")
    }

def process_chunk(num_transactions):
    transactions = [generate_transaction() for _ in range(num_transactions)]

```

```

    total_amount = sum(transaction['amount'] for transaction in transactions)
    # Simulate processing time per transaction
    time.sleep(num_transactions * 0.001) # This simulates the 1ms per transaction
processing delay
    return total_amount

def process_transactions_serial(num_transactions, num_cores):
    chunk_size = num_transactions // num_cores
    start_time = time.time()

    # Simulate each core processing a part of the workload sequentially
    for _ in range(num_cores):
        process_chunk(chunk_size)

    end_time = time.time()
    processing_time = end_time - start_time
    efficiency = num_transactions / processing_time # Transactions per second

    print(f"Simulated processing with {num_cores} cores:")
    print(f"Total processing time: {processing_time:.2f} seconds")
    print(f"Efficiency: {efficiency:.2f} transactions per second")

    return processing_time, efficiency

# Example usage
num_transactions = 10000
for num_cores in [1, 2, 4, 8]:
    process_transactions_serial(num_transactions, num_cores)

def process_transactions_parallel(num_transactions, num_cores):
    # Sparkx session setup with specified number of cores
    spark = SparkSession.builder \
        .appName("Parallel Processing") \
        .master(f"local[{num_cores}]") \
        .getOrCreate()

    # Generate data
    data = [{"user_id": f"user{random.randint(1, 100)}",
            "amount": round(random.uniform(10.0, 500.0), 2)}

```



```

        for _ in range(num_transactions)]
df = spark.createDataFrame(data)

# Begin timing the parallel process
start_time = time.time()
result = df.groupBy("user_id").agg(_sum("amount").alias("total_amount"))
result.collect() # Trigger computation
end_time = time.time()

# End of Spark session
spark.stop()

processing_time_parallel = end_time - start_time
efficiency = num_transactions / processing_time_parallel # Transactions per
second

print(f"Parallel Processing with {num_cores} cores:")
print(f"Processing Time: {processing_time_parallel:.2f} seconds")
print(f"Efficiency: {efficiency:.2f} transactions per second")

return processing_time_parallel, efficiency

# Test the function with different numbers of cores
num_transactions = 10000
core_configs = [1, 2, 4, 8]
results = {}
for cores in core_configs:
    results[cores] = process_transactions_parallel(num_transactions, cores)
    print(f"Results with {cores} cores: Time - {results[cores][0]:.2f}s, Efficiency -
{results[cores][1]:.2f} txn/s")

import matplotlib.pyplot as plt

# Test the function with different numbers of cores
num_transactions = 10000
core_configs = [1, 2, 4, 8]

# Lists to store results
serial_times = []
parallel_times = []

```

```

serial_efficiencies = []
parallel_efficiencies = []

# Process transactions for serial and parallel approaches
for cores in core_configs:
    # Serial processing
    serial_processing_time, serial_efficiency =
process_transactions_serial(num_transactions, cores)
    serial_times.append(serial_processing_time)
    serial_efficiencies.append(serial_efficiency)

    # Parallel processing
    parallel_processing_time, parallel_efficiency =
process_transactions_parallel(num_transactions, cores)
    parallel_times.append(parallel_processing_time)
    parallel_efficiencies.append(parallel_efficiency)

# Create line plots for processing time comparison
plt.figure(figsize=(10, 6))
plt.plot(core_configs, serial_times, marker='o', color='blue', label='Serial
Processing')
plt.plot(core_configs, parallel_times, marker='o', color='orange', label='Parallel
Processing')
plt.xlabel('Number of Cores')
plt.ylabel('Processing Time (seconds)')
plt.title('Processing Time Comparison')
plt.legend()
plt.grid(True)
plt.show()

# Create line plots for efficiency comparison
plt.figure(figsize=(10, 6))
plt.plot(core_configs, serial_efficiencies, marker='o', color='blue', label='Serial
Processing')
plt.plot(core_configs, parallel_efficiencies, marker='o', color='orange',
label='Parallel Processing')
plt.xlabel('Number of Cores')
plt.ylabel('Efficiency (transactions per second)')
plt.title('Efficiency Comparison')
plt.legend()

```

```

plt.grid(True)
plt.show()

import time
import random

def generate_transaction():
    return {
        "user_id": f"user{random.randint(1, 100)}",
        "amount": round(random.uniform(10.0, 500.0), 2),
        "transaction_time": time.strftime("%Y-%m-%d %H:%M:%S")
    }

def process_chunk(num_transactions):
    transactions = [generate_transaction() for _ in range(num_transactions)]
    total_amount = sum(transaction['amount'] for transaction in transactions)
    # Simulate processing time per transaction
    time.sleep(num_transactions * 0.001) # This simulates the 1ms per transaction
    processing_delay
    return total_amount

def process_transactions_serial(num_transactions, num_cores):
    chunk_size = num_transactions // num_cores
    start_time = time.time()

    # Simulate each core processing a part of the workload sequentially
    for _ in range(num_cores):
        process_chunk(chunk_size)

    end_time = time.time()
    processing_time = end_time - start_time
    efficiency = num_transactions / processing_time # Transactions per second

    print(f'Simulated processing with {num_cores} cores and {num_transactions}
transactions:')
    print(f'Total processing time: {processing_time:.2f} seconds")
    print(f'Efficiency: {efficiency:.2f} transactions per second")

    return processing_time, efficiency

```

```

# Example usage
data_sizes = [1000, 5000, 10000, 20000, 50000] # Different data sizes
for num_cores in [1, 2, 4, 8]: # Different numbers of cores
    for size in data_sizes:
        process_transactions_serial(size, num_cores)

from pyspark.sql import SparkSession
from pyspark.sql.functions import sum as _sum
import time
import random

def process_transactions_parallel(num_transactions, num_cores):
    # Spark session setup with specified number of cores
    spark = SparkSession.builder \
        .appName("Parallel Processing") \
        .master(f'local[{num_cores}]') \
        .getOrCreate()

    # Generate data
    data = [{"user_id": f"user{random.randint(1, 100)}",
            "amount": round(random.uniform(10.0, 500.0), 2)}
            for _ in range(num_transactions)]
    df = spark.createDataFrame(data)

    # Begin timing the parallel process
    start_time = time.time()
    result = df.groupBy("user_id").agg(_sum("amount").alias("total_amount"))
    result.collect() # Trigger computation
    end_time = time.time()

    # End of Spark session
    spark.stop()

    processing_time_parallel = end_time - start_time
    efficiency = num_transactions / processing_time_parallel # Transactions per
second

    print(f"Parallel Processing with {num_cores} cores and {num_transactions}
transactions:")
    print(f"Processing Time: {processing_time_parallel:.2f} seconds")

```

```

print(f"Efficiency: {efficiency:.2f} transactions per second")

return processing_time_parallel, efficiency

# Test the function with different data sizes
data_sizes = [1000, 5000, 10000, 20000, 50000] # Different data sizes
for num_cores in [1, 2, 4, 8]:
    for size in data_sizes:
        process_transactions_parallel(size, num_cores)

import matplotlib.pyplot as plt

# Test the function with different data sizes
data_sizes = [1000, 5000, 10000, 20000, 50000]

# Lists to store results
serial_processing_times = []
parallel_processing_times = []
serial_efficiencies = []
parallel_efficiencies = []

# Process transactions for serial and parallel approaches
for size in data_sizes:
    # Serial processing
    serial_processing_time, serial_efficiency = process_transactions_serial(size, 1)
    serial_processing_times.append(serial_processing_time)
    serial_efficiencies.append(serial_efficiency)

    # Parallel processing
    parallel_processing_time, parallel_efficiency =
process_transactions_parallel(size, 8)
    parallel_processing_times.append(parallel_processing_time)
    parallel_efficiencies.append(parallel_efficiency)

# Create plots
plt.figure(figsize=(12, 6))

# Efficiency comparison (line plot)
plt.subplot(1, 2, 1)
plt.plot(data_sizes, serial_efficiencies, marker='o', label='Serial Processing')

```

```

plt.plot(data_sizes, parallel_efficiencies, marker='o', label='Parallel Processing')
plt.xlabel('Data Size')
plt.ylabel('Efficiency (transactions per second)')
plt.title('Efficiency Comparison')
plt.legend()
plt.grid(True)

# Processing time comparison (line plot)
plt.subplot(1, 2, 2)
plt.plot(data_sizes, serial_processing_times, marker='o', label='Serial Processing')
plt.plot(data_sizes, parallel_processing_times, marker='o', label='Parallel
Processing')
plt.xlabel('Data Size')
plt.ylabel('Processing Time (seconds)')
plt.title('Processing Time Comparison')
plt.legend()
plt.grid(True)

plt.tight_layout()
plt.show()

import time
import random
from pyspark.sql import SparkSession
from pyspark.sql.functions import sum as _sum

def generate_transaction():
    return {
        "user_id": f"user{random.randint(1, 100)}",
        "amount": round(random.uniform(10.0, 500.0), 2),
        "transaction_time": time.strftime("%Y-%m-%d %H:%M:%S")
    }

def process_transactions_serial(num_transactions=10000):
    start_time = time.time()
    transactions = [generate_transaction() for _ in range(num_transactions)]
    total_amount = 0
    for transaction in transactions:
        total_amount += transaction["amount"]
        time.sleep(0.001) # simulate processing delay

```

```

    end_time = time.time()
    processing_time_serial = end_time - start_time
    efficiency = num_transactions / processing_time_serial # Transactions per
second

    print(f"Total amount processed: ${total_amount:.2f}")
    print(f"Serial Processing Time: {processing_time_serial:.2f} seconds")
    print(f"Efficiency: {efficiency:.2f} transactions per second")

    return processing_time_serial, efficiency

def process_transactions_parallel(num_transactions=10000):
    spark = SparkSession.builder \
        .appName("Parallel Processing") \
        .master("local[*]") \
        .getOrCreate()

    data = [{"user_id": f"user{random.randint(1, 100)}", "amount":
round(random.uniform(10.0, 500.0), 2)} for _ in range(num_transactions)]
    df = spark.createDataFrame(data)

    start_time = time.time()
    result = df.groupBy("user_id").agg(_sum("amount").alias("total_amount"))
    result.collect() # Trigger computation
    end_time = time.time()

    spark.stop()

    processing_time_parallel = end_time - start_time
    efficiency = num_transactions / processing_time_parallel # Transactions per
second

    print(f"Parallel Processing Time with Spark: {processing_time_parallel:.2f}
seconds")
    print(f"Efficiency: {efficiency:.2f} transactions per second")

    return processing_time_parallel, efficiency

# Run serial and parallel processing

```

```

serial_time, serial_efficiency = process_transactions_serial()
parallel_time, parallel_efficiency = process_transactions_parallel()

# Calculate speedup and efficiency gain
speed_up = serial_time / parallel_time
efficiency_gain = parallel_efficiency / serial_efficiency

print(f'Speedup: {speed_up:.2f}')
print(f'Efficiency Gain: {efficiency_gain:.2f}')

import time
import random
from pyspark.sql import SparkSession
from pyspark.sql.functions import sum as _sum
import matplotlib.pyplot as plt

def generate_transaction():
    return {
        "user_id": f"user{random.randint(1, 100)}",
        "amount": round(random.uniform(10.0, 500.0), 2),
        "transaction_time": time.strftime("%Y-%m-%d %H:%M:%S")
    }

def process_transactions_serial(num_transactions=10000):
    start_time = time.time()
    transactions = [generate_transaction() for _ in range(num_transactions)]
    total_amount = 0
    for transaction in transactions:
        total_amount += transaction["amount"]
        time.sleep(0.001) # simulate processing delay

    end_time = time.time()
    processing_time_serial = end_time - start_time
    efficiency = num_transactions / processing_time_serial # Transactions per
second

    return total_amount, processing_time_serial, efficiency

def process_transactions_parallel(num_transactions=10000):
    spark = SparkSession.builder \

```



```

        .appName("Parallel Processing") \
        .master("local[*]") \
        .getOrCreate()

    data = [{"user_id": f"user{random.randint(1, 100)}", "amount":
round(random.uniform(10.0, 500.0), 2)} for _ in range(num_transactions)]
    df = spark.createDataFrame(data)

    start_time = time.time()
    result = df.groupBy("user_id").agg(_sum("amount").alias("total_amount"))
    result.collect() # Trigger computation
    end_time = time.time()

    spark.stop()

    processing_time_parallel = end_time - start_time
    efficiency = num_transactions / processing_time_parallel # Transactions per
second

    return processing_time_parallel, efficiency

# Run serial and parallel processing
total_amount_serial, serial_time, serial_efficiency = process_transactions_serial()
parallel_time, parallel_efficiency = process_transactions_parallel()

# Calculate speedup and efficiency gain
speed_up = serial_time / parallel_time
efficiency_gain = parallel_efficiency / serial_efficiency

# Plotting the results
fig, ax = plt.subplots(1, 2, figsize=(12, 6))

# Speedup plot
ax[0].bar(['Speedup'], [speed_up], color='blue')
ax[0].set_title('Speedup Comparison')
ax[0].set_ylabel('Speedup Factor')
ax[0].set_ylim(0, max(speed_up * 1.2, 1))

# Efficiency gain plot
ax[1].bar(['Efficiency Gain'], [efficiency_gain], color='green')

```

```
ax[1].set_title('Efficiency Gain')
ax[1].set_ylabel('Efficiency Factor')
ax[1].set_ylim(0, max(efficiency_gain * 1.2, 1))
```

```
# Show plot
plt.tight_layout()
plt.show()
```