

Distributed Computing of High-Frequency Data Quantitative Factors

Zihan Zhang¹ and Benxiang Xiao¹

¹Southern University of Science and Technology

December 22, 2025

Abstract

Based on the HDFS+MapReduce framework (Java implementation), this project conducts distributed computing of 20 Limit Order Book (LOB)-related quantitative factors using Shenzhen Stock Exchange Level-10 high-frequency snapshot data (3-second frequency). Key parameters are set as $n=5$ (order levels) and $\Delta t=1$ (time interval), with $1e-7$ added to denominators to avoid division by zero. Core tasks include calculating the average factor value sequences of CSI 300 index stocks (training phase: data from Jan 2-8, 2024) and outputting sequences for new data (testing phase), evaluated by technical documentation, presentations, and code performance (1% error tolerance, speed ranking). Through multi-dimensional optimizations—including CSV parsing enhancement, object reuse, memory control, Key compression, Shuffle optimization, concurrency tuning, and instruction-level improvements—computational efficiency is significantly improved while ensuring accuracy, enabling efficient distributed processing of high-frequency quantitative factor data.

Keywords High-frequency data; Limit Order Book; Distributed computing; MapReduce; Quantitative factors

1 Problem Background and Understanding

1.1 Project Background

Quantitative trading has become an indispensable component of financial markets. Unlike traditional subjective judgment, quantitative factors capture market characteristics through data-driven methods, helping investors make more rational and systematic decisions. The goal of this project is to build a distributed computing system capable of processing massive high-frequency data to calculate time series of multiple quantitative factors.

We use Level-10 market snapshot data from the Shenzhen Stock Exchange. This data is characterized by high frequency (one snapshot every 3 seconds), recording the price and volume information of the top 10 bid and ask levels in the market at each timestamp. For an index like the CSI 300 (covering 300 stocks), the daily data volume is substantial. The project requires processing all data from 5 trading days (January 2–8, 2024) to compute 20 distinct quantitative factors, and outputting the average values of these factors across all stocks.

1.2 Data Feature Analysis

The Level-10 market snapshot data contains 57 fields, with the most critical being the price and volume information for the top 10 bid/ask levels. Fields such as the best bid price (bp_1) and best ask price (ap_1) reflect the optimal quotes in the current market, while subsequent levels (bp_2 to bp_{10} , ap_2 to ap_{10}) show deeper market depth. Each snapshot also records the total bid and ask volumes across the entire market, which are fundamental data required for factor calculation.

The data time range spans from 9:25 AM to 3:00 PM, but per project requirements, we only process data from the official trading session (9:30 AM to 3:00 PM). This means time filtering must be implemented during the data preprocessing phase.

1.3 Factor Description

The 20 factors required by the project can be roughly categorized into several types (key formulas are listed for typical factors):

- **Spread factors:**

$$\text{Optimal Spread} = ap_1 - bp_1$$

$$\text{Relative Spread} = \frac{ap_1 - bp_1}{(ap_1 + bp_1)/2 + 10^{-7}}$$

(Note: 10^{-7} is added to the denominator to avoid division by zero.)

- **Imbalance factors:** Reflecting the strength balance between buyers and sellers (e.g., level-1 order imbalance: $\frac{bv_1 - av_1}{bv_1 + av_1 + 10^{-7}}$, where bv_1/av_1 are level-1 bid/ask volumes).
- **Depth factors:** Describing the order book thickness (e.g., 5-level bid depth: $\sum_{i=1}^5 bv_i$).
- **Weighted price factors:** Volume-weighted prices (e.g., 5-level weighted bid price: $\frac{\sum_{i=1}^5 bp_i \times bv_i}{\sum_{i=1}^5 bv_i + 10^{-7}}$).
- **Change factors:** Requiring cross-timestamp data (e.g., best bid change: $bp_1^{(t)} - bp_1^{(t-1)}$, where t denotes the current timestamp).

Other categories include density/symmetry factors (describing order distribution) and price pressure indicators (integrating spread and depth information).

2 Technical Challenges and Solutions

2.1 Challenge of Stateful Computing

Among the 20 factors, three (Factors 17, 18, 19) are special as they calculate changes (e.g., Factor 17: best price change, computed as the current best ask price minus the previous best ask price). This introduces a challenge: MapReduce is stateless by design (each data row is processed independently), so maintaining historical information within this framework is critical.

We use a `HashMap` in the Mapper to store the previous snapshot of each stock: the key is constructed as "[StockCode].[Exchange]_ [Date]" (e.g., "000001.SZ_20240102"), and the value is the snapshot object from the last processing step. When a new snapshot is processed:

- Check the `HashMap` for historical data of the stock;
- If available, compute the change value using the historical snapshot;
- If unavailable (i.e., it is the first record of the stock), set the change value to 0.

2.2 Output Format Requirements

The final output requires a separate CSV file for each day (named with the last four digits of the date, e.g., 0102.csv), all stored in the same directory. The file content must use comma separation, and time must be in pure numeric format (e.g., 93000 instead of 09:30:00).

MapReduce's default output is a single file (e.g., `part-r-00000`) with tab-separated key-value pairs, mixing data from all dates. Thus, a post-processing step is needed to reorganize the output: A `DailyOutputFormatter` program reads the MapReduce output, groups data by date, and generates a CSV file for each date. During this process, it also:

- Adjusts the delimiter (from tab to comma);
- Converts the time format (extracts 93000 from "20240102_93000").

3 System Architecture Design

The entire system adopts a standard two-stage MapReduce processing architecture, with an additional post-processing stage at the output end for formatting. Figure 1 illustrates the overall system structure and data flow process.

CSV files read from HDFS are first distributed to multiple Mappers for parallel processing. Each Mapper parses the CSV data, computes factor values, and outputs key-value pairs. Here, the Key is a combination of date and timestamp (e.g., "20240102_93000"), and the Value is a Factor object containing 20 factor values and a counter (initialized to 1).

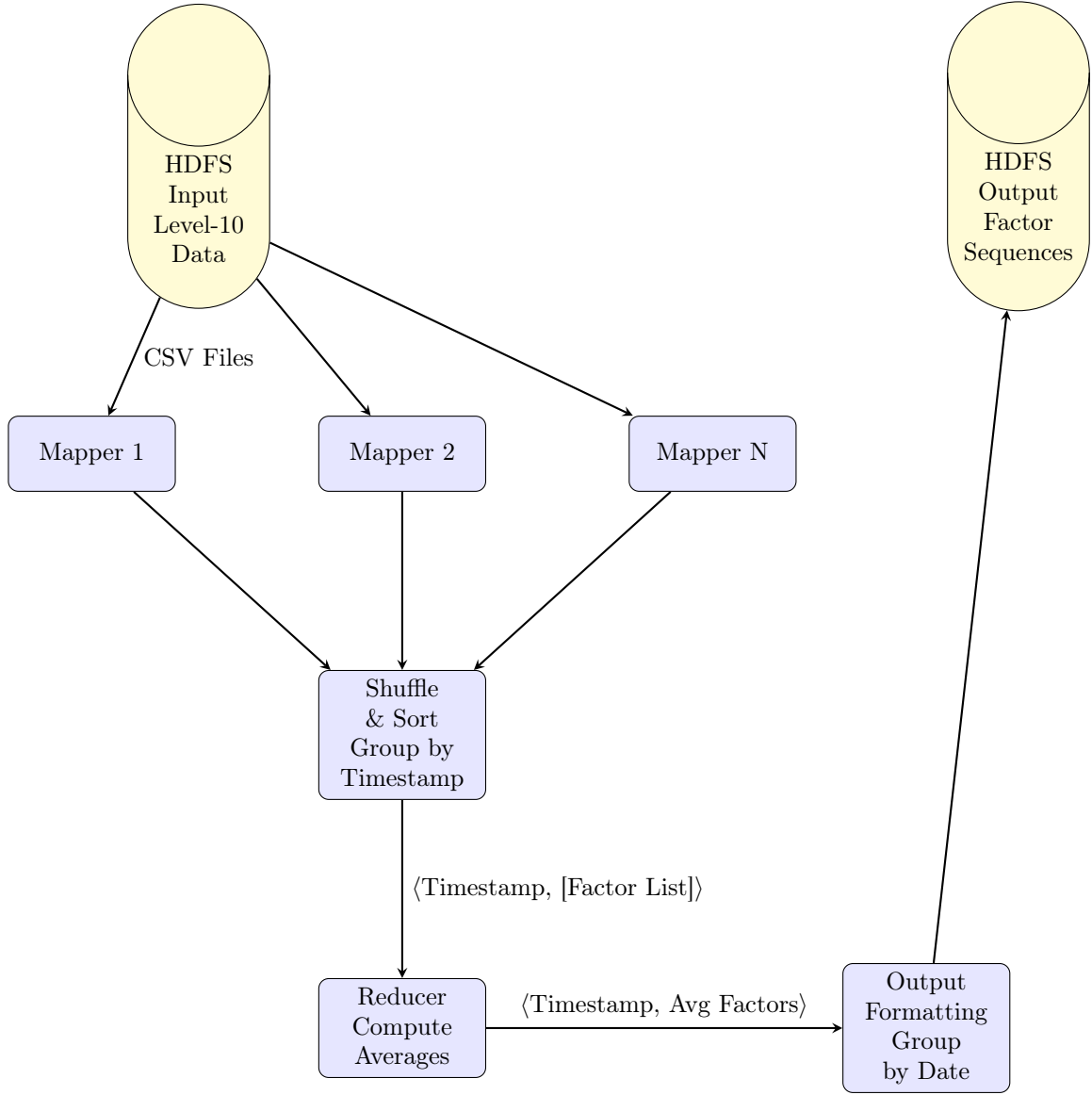


Figure 1: Overall System Architecture

Hadoop’s Shuffle stage automatically aggregates data with the same Key and sends it to the same Reducer. Factor values of all stocks at the same timestamp are collected together. The Reducer accumulates these Factor objects and divides by the number of stocks to obtain the average factor values.

The output of MapReduce is not the final format—it is tab-separated key-value pairs with all data in a single file. The post-processing stage reads this file, groups data by date, and generates an independent CSV file for each date. Meanwhile, it adjusts the output format: converting timestamps from "20240102_93000" to "93000", replacing tab separators with commas, and adding headers, etc.

4 Core Module Implementation

4.1 Data Structure Design

The foundational class in the project is `SnapshotData`, which represents a market snapshot of a specific stock at a given timestamp. This class stores essential information including trading date, trade time, stock code, and price/volume data for the top 10 bid and ask levels. Four arrays are used to store order book details: `bidPrices` and `bidVolumes` for bid-side prices and volumes, and `askPrices` and `askVolumes` for ask-side data. Additionally, two fields (`tBidVol` and `tAskVol`) store the total bid and ask volumes across the entire market, respectively.

The most critical method of `SnapshotData` is `parseFromCSV`, which parses a single CSV line into a `SnapshotData`

object. The workflow involves splitting the string by commas, validating the number of fields, extracting and converting required fields, and returning the constructed object. If any step fails, the method returns `null` to indicate parsing failure.

```
1  /**
2   * Parses snapshot data from a CSV line
3   * @param line CSV line string (comma-separated)
4   * @return SnapshotData object if parsing succeeds, null otherwise
5   */
6  public static SnapshotData parseFromCSV(String line) {
7      try {
8          // Split line by commas
9          String[] fields = line.split(",");
10
11         // Validate field count
12         if (fields.length < 57) {
13             return null;
14         }
15
16         SnapshotData snapshot = new SnapshotData();
17
18         // Extract basic fields
19         snapshot.tradingDay = Integer.parseInt(fields[FIELD_TRADING_DAY]);
20         snapshot.tradeTime = Long.parseLong(fields[FIELD_TRADE_TIME]);
21         snapshot.code = fields[FIELD_CODE];
22         snapshot.tBidVol = Long.parseLong(fields[FIELD_TBID_VOL]);
23         snapshot.tAskVol = Long.parseLong(fields[FIELD_TASK_VOL]);
24
25         // Handle timestamp in scientific notation
26         if (snapshot.tradeTime > 10000000000L) {
27             snapshot.tradeTime = snapshot.tradeTime / 10000000000L;
28         }
29
30         // Extract top 10 bid/ask prices and volumes
31         snapshot.bidPrices = new long[10];
32         snapshot.bidVolumes = new long[10];
33         snapshot.askPrices = new long[10];
34         snapshot.askVolumes = new long[10];
35
36         for (int i = 0; i < 10; i++) {
37             int baseIndex = FIELD_PRICE_START + i * 4;
38             snapshot.bidPrices[i] = Long.parseLong(fields[baseIndex]);
39             snapshot.bidVolumes[i] = Long.parseLong(fields[baseIndex + 1]);
40             snapshot.askPrices[i] = Long.parseLong(fields[baseIndex + 2]);
41             snapshot.askVolumes[i] = Long.parseLong(fields[baseIndex + 3]);
42         }
43
44         return snapshot;
45     } catch (NumberFormatException e) {
46         // Failed to convert value to number
47         return null;
48     } catch (ArrayIndexOutOfBoundsException e) {
49         // Insufficient fields
50         return null;
51     }
52 }
53 }
```

`SnapshotData` also provides an `isValid` method to validate data integrity. It checks two key conditions: 1) the validity of the trading date and stock code (non-null and non-zero), and 2) whether the trade time falls within the official trading session (09:30:00 to 15:00:00).

```
1  /**
2   * Validates if the snapshot data is valid
3   * Mainly checks if trade time is within 09:30:00-15:00:00
```

```

4  */
5  public boolean isValid() {
6      // Check basic fields
7      if (tradingDay <= 0 || code == null || code.isEmpty()) {
8          return false;
9      }
10
11     // Extract time component (handle scientific notation)
12     long timeOnly = tradeTime;
13     if (timeOnly > 10000000000L) {
14         timeOnly = timeOnly / 10000000000L;
15     }
16
17     // Retain data only from official trading hours
18     // 93000 = 09:30:00, 150000 = 15:00:00
19     return timeOnly >= 93000 && timeOnly <= 150000;
20 }

```

4.2 Factor Calculation Implementation

The `Factor` class encapsulates the calculation and storage of quantitative factors. It has two core member variables: `factorValues` (a double array of length 20 to store factor values) and `count` (to record the number of stocks whose data has been accumulated). A static constant `EPSILON` (10^{-7}) is defined for division-by-zero protection.

The core of factor calculation is the static method `calculateFactors`, which takes two parameters: the current snapshot and the previous snapshot (may be `null`). The method first extracts the top 10 bid/ask prices and volumes from the current snapshot, then calculates the 20 factors one by one according to predefined formulas.

```

1  /**
2   * @param snapshot Current market snapshot
3   * @param prevSnapshot Previous timestamp's snapshot (null for change factors)
4   * @return Double array of length 20 containing all factor values
5   */
6  public static double[] calculateFactors(SnapshotData snapshot,
7                                         SnapshotData prevSnapshot) {
8      double[] factors = new double[20];
9
10     // Extract top 10 bid/ask prices and volumes
11     long[] bp = snapshot.getBidPrices();
12     long[] bv = snapshot.getBidVolumes();
13     long[] ap = snapshot.getAskPrices();
14     long[] av = snapshot.getAskVolumes();
15
16     // Factor 1: Optimal Spread = ap1 - bp1
17     factors[0] = (double)(ap[0] - bp[0]);
18
19     // Calculate mid-price (used in multiple factors)
20     double midPrice = (ap[0] + bp[0]) / 2.0;
21
22     // Factor 2: Relative Spread = (ap1 - bp1) / midPrice
23     factors[1] = (ap[0] - bp[0]) / (midPrice + EPSILON);
24
25     // Factor 3: Mid-price
26     factors[2] = midPrice;
27
28     // Factor 4: Level-1 Imbalance = (bv1 - av1) / (bv1 + av1)
29     factors[3] = (bv[0] - av[0]) / (double)(bv[0] + av[0] + EPSILON);
30
31     // Calculate cumulative sums for top n levels (used in Factors 5-9)
32     long sumBidVol = 0;
33     long sumAskVol = 0;
34     for (int i = 0; i < n; i++) {
35         sumBidVol += bv[i];

```

```

36     sumAskVol += av[i];
37 }
38
39 // Factor 5: Multi-level Imbalance = (sumBidVol - sumAskVol) / totalVolume
40 factors[4] = (sumBidVol - sumAskVol) /
41     (double)(sumBidVol + sumAskVol + EPSILON);
42
43 // Factor 6: Bid Depth
44 factors[5] = sumBidVol;
45
46 // Factor 7: Ask Depth
47 factors[6] = sumAskVol;
48
49 // Factor 8: Depth Difference = Bid Depth - Ask Depth
50 factors[7] = sumBidVol - sumAskVol;
51
52 // Factor 9: Depth Ratio = Bid Depth / Ask Depth
53 factors[8] = sumBidVol / (double)(sumAskVol + EPSILON);
54
55 // Factor 10: Market Volume Imbalance
56 long tBidVol = snapshot.getTBidVol();
57 long tAskVol = snapshot.getTAskVol();
58 factors[9] = (tBidVol - tAskVol) /
59     (double)(tBidVol + tAskVol + EPSILON);
60
61 // Calculate weighted prices (used in Factors 11-14)
62 double weightedBidPrice = 0;
63 double weightedAskPrice = 0;
64 for (int i = 0; i < n; i++) {
65     weightedBidPrice += bp[i] * bv[i];
66     weightedAskPrice += ap[i] * av[i];
67 }
68
69 // Factor 11: Weighted Bid Price
70 factors[10] = weightedBidPrice / (sumBidVol + EPSILON);
71
72 // Factor 12: Weighted Ask Price
73 factors[11] = weightedAskPrice / (sumAskVol + EPSILON);
74
75 // Factor 13: Weighted Mid-price
76 factors[12] = (weightedBidPrice + weightedAskPrice) /
77     (sumBidVol + sumAskVol + EPSILON);
78
79 // Factor 14: Weighted Spread = Weighted Ask Price - Weighted Bid Price
80 factors[13] = factors[11] - factors[10];
81
82 // Factor 15: Bid-Ask Density Difference
83 factors[14] = (sumBidVol / (double)n) - (sumAskVol / (double)n);
84
85 // Factor 16: Bid-Ask Asymmetry (level-weighted)
86 double weightedBidSum = 0;
87 double weightedAskSum = 0;
88 for (int i = 0; i < n; i++) {
89     weightedBidSum += bv[i] / (double)(i + 1);
90     weightedAskSum += av[i] / (double)(i + 1);
91 }
92 factors[15] = (weightedBidSum - weightedAskSum) /
93     (weightedBidSum + weightedAskSum + EPSILON);
94
95 // Factors 17-19: Change Factors (require previous snapshot)
96 if (prevSnapshot != null) {
97     long[] prevAp = prevSnapshot.getAskPrices();
98     long[] prevBp = prevSnapshot.getBidPrices();
99     long[] prevBv = prevSnapshot.getBidVolumes();
100    long[] prevAv = prevSnapshot.getAskVolumes();

```

```

101
102     // Calculate previous mid-price and depth ratio
103     double prevMidPrice = (prevAp[0] + prevBp[0]) / 2.0;
104     long prevSumBidVol = 0;
105     long prevSumAskVol = 0;
106     for (int i = 0; i < n; i++) {
107         prevSumBidVol += prevBv[i];
108         prevSumAskVol += prevAv[i];
109     }
110     double prevDepthRatio = prevSumBidVol /
111         (double)(prevSumAskVol + EPSILON);
112
113     // Factor 17: Best Price Change = Current ap1 - Previous ap1
114     factors[16] = (double)(ap[0] - prevAp[0]);
115
116     // Factor 18: Mid-price Change
117     factors[17] = midPrice - prevMidPrice;
118
119     // Factor 19: Depth Ratio Change
120     factors[18] = factors[8] - prevDepthRatio;
121 } else {
122     // Set change factors to 0 if no previous snapshot (first record)
123     factors[16] = 0;
124     factors[17] = 0;
125     factors[18] = 0;
126 }
127
128 // Factor 20: Price Pressure = Spread / Total Depth
129 factors[19] = (ap[0] - bp[0]) /
130     (double)(sumBidVol + sumAskVol + EPSILON);
131
132 return factors;
133 }

```

The Factor class implements Hadoop's Writable interface to enable data transmission between Mappers and Reducers. The write method serializes the array and counter to the output stream, while readFields deserializes them from the input stream.

```

1  @Override
2  public void write(DataOutput out) throws IOException {
3      // Write array length (fixed at 20, but following standard format)
4      out.writeInt(factorValues.length);
5      // Write each factor value
6      for (double value : factorValues) {
7          out.writeDouble(value);
8      }
9      // Write counter
10     out.writeInt(count);
11 }
12
13 @Override
14 public void readFields(DataInput in) throws IOException {
15     // Read array length
16     int length = in.readInt();
17     factorValues = new double[length];
18     // Read each factor value
19     for (int i = 0; i < length; i++) {
20         factorValues[i] = in.readDouble();
21     }
22     // Read counter
23     count = in.readInt();
24 }

```

The Factor class also provides a merge method to combine two Factor objects. It sums the corresponding elements of the factor value arrays and increments the counter. The Reducer uses this method repeatedly to

aggregate data from all stocks at the same timestamp.

```
1  /**
2   * Merges another Factor object (used for accumulation in Reduce phase)
3   * @param other Factor object to merge
4   */
5  public void merge(Factor other) {
6      for (int i = 0; i < 20; i++) {
7          this.factorValues[i] += other.factorValues[i];
8      }
9      this.count += other.count;
10 }
11
12 /**
13  * Gets average factor values
14  * @return Array of average values
15  */
16 public double[] getAverageFactors() {
17     double[] avgFactors = new double[20];
18     for (int i = 0; i < 20; i++) {
19         avgFactors[i] = factorValues[i] / count;
20     }
21     return avgFactors;
22 }
```

4.3 Mapper Implementation

FactorMapper inherits from Hadoop's Mapper class, with generic parameters specifying input/output types. The input consists of LongWritable (line number) and Text (CSV line content), while the output is Text (timestamp) and Factor (factor values and count).

```
1  public class FactorMapper
2      extends Mapper<LongWritable, Text, Text, Factor> {
3
4      private String currentStockCode; // Currently processed stock code
5      private Map<String, SnapshotData> previousSnapshots; // Historical
6          snapshots
7
8      @Override
9      protected void setup(Context context) throws IOException {
10         // Extract stock code from file path
11         FileSplit fileSplit = (FileSplit) context.getInputSplit();
12         String filePath = fileSplit.getPath().toString();
13
14         // Parse path: .../0102/000001.SZ/snapshot.csv
15         // Use second-to-last path segment as stock code
16         String[] pathParts = filePath.split("/");
17         currentStockCode = pathParts[pathParts.length - 2];
18
19         // Initialize historical snapshot storage
20         previousSnapshots = new HashMap<>();
21     }
22 }
```

The map method is called for each CSV line. It first skips headers and empty lines, then parses the line using SnapshotData.parseFromCSV. If parsing fails (returns null), it increments an error counter and returns. If successful, it validates the data to ensure the time falls within trading hours.

```
1  @Override
2  protected void map(LongWritable key, Text value, Context context)
3      throws IOException, InterruptedException {
4
5      String line = value.toString();
6  }
```



```

7 // Skip header and empty lines
8 if (line.startsWith("tradingDay") || line.trim().isEmpty()) {
9     return;
10 }
11
12 // Parse CSV line
13 SnapshotData snapshot = SnapshotData.parseFromCSV(line);
14 if (snapshot == null) {
15     context.getCounter("FactorMapper", "PARSE_ERROR").increment(1);
16     return;
17 }
18
19 // Validate time range (09:30:00 - 15:00:00)
20 if (!snapshot.isValid()) {
21     context.getCounter("FactorMapper", "INVALID_TIME").increment(1);
22     return;
23 }
24
25 // Look up previous snapshot (for change factors)
26 String stockKey = snapshot.getCode() + "_" + snapshot.getTradingDay();
27 SnapshotData prevSnapshot = previousSnapshots.get(stockKey);
28
29 // Calculate 20 factors
30 double[] factorValues = Factor.calculateFactors(snapshot, prevSnapshot);
31
32 // Validate factor values (not NaN or Infinite)
33 boolean isValid = true;
34 for (double value : factorValues) {
35     if (Double.isNaN(value) || Double.isInfinite(value)) {
36         isValid = false;
37         break;
38     }
39 }
40
41 if (!isValid) {
42     context.getCounter("FactorMapper", "INVALID_FACTOR").increment(1);
43     return;
44 }
45
46 // Construct output Key: "Date_Time" format (e.g., "20240102_93000")
47 String keyStr = snapshot.getTradingDay() + "_" +
48     snapshot.getFormattedTradeTime();
49
50 // Output key-value pair
51 context.write(new Text(keyStr), new Factor(factorValues, 1));
52
53 // Update historical snapshot (for next calculation)
54 previousSnapshots.put(stockKey, snapshot);
55
56 // Increment successful processing counter
57 context.getCounter("FactorMapper", "SUCCESS").increment(1);
58 }

```

4.4 Reducer and Output Formatting

The implementation of `FactorReducer` is relatively straightforward. Its `reduce` method receives a `Key` (timestamp) and an iterator of `Factor` objects. Its task is to iterate through the iterator, merge all `Factor` objects, calculate the average values, and format the output.

```

1 public class FactorReducer
2     extends Reducer<Text, Factor, Text, Text> {
3
4     @Override
5     protected void reduce(Text key, Iterable<Factor> values,

```

```

6         Context context)
7         throws IOException, InterruptedException {
8
9         // Create empty Factor for accumulation
10        Factor sumFactor = new Factor();
11
12        // Iterate and accumulate all Factor objects
13        for (Factor factor : values) {
14            sumFactor.merge(factor);
15        }
16
17        // Calculate average values
18        double[] avgFactors = sumFactor.getAverageFactors();
19
20        // Format output (comma-separated, 10 decimal places)
21        StringBuilder output = new StringBuilder();
22        for (int i = 0; i < 20; i++) {
23            if (i > 0) {
24                output.append(",");
25            }
26            output.append(String.format("%.10f", avgFactors[i]));
27        }
28
29        // Output: Key = timestamp, Value = factor value list
30        context.write(key, new Text(output.toString()));
31    }
32 }

```

DailyOutputFormatter is an independent program for post-processing MapReduce output. It reads the part-r-00000 file, groups data by date, and generates a separate CSV file for each date. During this process, it adjusts the format: converting timestamps from "20240102_93000" to "93000", replacing tab separators with commas, and adding the header "tradeTime,alpha_1,...,alpha_20".

```

1 public class DailyOutputFormatter {
2
3     public static void main(String[] args) throws IOException {
4         Configuration conf = new Configuration();
5         FileSystem fs = FileSystem.get(conf);
6
7         // Read MapReduce output file
8         Path mrOutputFile = new Path(args[0] + "/part-r-00000");
9         BufferedReader reader = new BufferedReader(
10             new InputStreamReader(fs.open(mrOutputFile)));
11
12        // Group data by date
13        Map<String, List<String>> dataByDate = new TreeMap<>();
14
15        String line;
16        while ((line = reader.readLine()) != null) {
17            String[] parts = line.split("\t");
18            String timeKey = parts[0]; // "20240102_93000"
19            String factorValues = parts[1]; // "100.0,0.001,..."
20
21            // Extract date and time
22            String[] timeParts = timeKey.split("_");
23            String fullDate = timeParts[0]; // "20240102"
24            String time = timeParts[1]; // "93000"
25
26            // Extract short date (last 4 digits)
27            String dateShort = fullDate.substring(4); // "0102"
28
29            // Compose record: time,factor1,factor2,...
30            String record = time + "," + factorValues;
31
32            // Group by date

```

```

33         dataByDate.computeIfAbsent(dateShort, k -> new ArrayList<>())
34             .add(record);
35     }
36     reader.close();
37
38     // Generate CSV file for each date
39     Path outputDir = new Path(args[1]);
40     if (!fs.exists(outputDir)) {
41         fs.mkdirs(outputDir);
42     }
43
44     for (String date : dataByDate.keySet()) {
45         // File name: 0102.csv
46         Path csvFile = new Path(outputDir, date + ".csv");
47         BufferedWriter writer = new BufferedWriter(
48             new OutputStreamWriter(fs.create(csvFile)));
49
50         // Write header
51         writer.write("tradeTime,alpha_1,alpha_2,...,alpha_20\n");
52
53         // Sort and write data
54         List<String> records = dataByDate.get(date);
55         Collections.sort(records);
56         for (String record : records) {
57             writer.write(record + "\n");
58         }
59
60         writer.close();
61     }
62 }
63 }

```

We will integrate the Reducer and DailyOutputFormatter into a single component in subsequent iterations.

5 Optimization Roadmap

This optimization focuses on the high-frequency data quantitative factor calculation scenario under the Hadoop MapReduce framework. After 11 iterations, it has evolved from basic performance tuning to in-depth underlying logic reconstruction, forming a clear roadmap: **Resource Utilization Optimization** → **Data Parsing Acceleration** → **Computational Efficiency Improvement** → **Data Transmission Compression** → **Memory and GC Optimization** → **Extreme Low-Level Optimization**. The optimization ideas at each stage target core pain points, achieving stepwise performance improvement through the integration of similar solutions.

5.1 Resource Utilization Optimization: Maximizing Hardware and Framework Potential

The core goal is to address efficiency bottlenecks caused by hardware resource waste and unreasonable default framework configurations, enabling precise matching between computing resources and task loads.

- **MapTask Concurrency and Split Optimization**

- *Small File Merging*: The optimization adopt `CombineTextInputFormat` instead of the default input format to merge multiple small files into a single `InputSplit`, reducing MapTask startup overhead. We further adjusts the split size range (128MB 512MB) to adapt to HDFS block size and cluster resources, allowing a single Mapper to process more data.

```

1 Job job = Job.getInstance(conf, "CSI300 Factor Calculation");
2 job.setJarByClass(FactorCalculationJob.class);
3
4 // Use CombineTextInputFormat to merge multiple small files
5 job.setInputFormatClass(CombineTextInputFormat.class);
6 CombineTextInputFormat.setMinInputSplitSize(job, 128 * 1024 * 1024);

```

```
7 CombineTextInputFormat.setMaxInputSplitSize(job, 512 * 1024 * 1024);
```

- *Local Mode Concurrency Adaptation*: To solve the problem of the default single Map task in Local mode, we modified the program driver mode (implementing `Configured` and `Tool` interfaces) to automatically identify the number of CPU cores of the machine and set the Map concurrency to the number of available CPU cores. Even though this method cannot do works on Distributed System, it may have inadvertently corrected certain default inefficient configurations in the pseudo-distributed environment.

```
1 # Detect Runner Class
2 String mapRunnerClass = conf.get("mapreduce.job.map.runner.class", "")
3 ;
4 boolean isLocalMode = mapRunnerClass.contains("LocalJobRunner")
5 || conf.get("mapreduce.framework.name", "local").
6 equals("local");
7
8 // If the mode is Local, use more CPU cores
9 if (isLocalMode) {
10     int localMapMax = conf.getInt("mapreduce.local.map.tasks.maximum",
11     1);
12     if (localMapMax == 1) {
13         int cpuCores = Runtime.getRuntime().availableProcessors();
14         conf.setInt("mapreduce.local.map.tasks.maximum", cpuCores);
15         conf.setInt("mapreduce.local.reduce.tasks.maximum", cpuCores);
16     }
17 }
```

- **Memory Configuration Tuning**: Take the lead in adjusting Map/Reduce memory allocation and JVM heap memory ratio to reserve sufficient memory for computing tasks.

```
1 // ----- Compression and IO Optimization -----
2 conf.setBoolean("mapreduce.map.output.compress", true);
3 conf.set("mapreduce.map.output.compress.codec", "org.apache.hadoop.io.
4 compress.SnappyCodec");
5
6 // ----- Memory and JVM Optimization -----
7 // Container memory (physical limit)
8 conf.set("mapreduce.map.memory.mb", "2048");
9 conf.set("mapreduce.reduce.memory.mb", "4096");
10
11 // JVM heap memory (logical limit, usually 75%-80% of container memory)
12 // -XX:+UseG1GC: Recommended for scenarios with large memory and mixed
13 // object lifecycles
14 conf.set("mapreduce.map.java.opts", "-Xmx1638m -XX:+UseG1GC");
15 conf.set("mapreduce.reduce.java.opts", "-Xmx3276m -XX:+UseG1GC");
16
17 // Shuffle Parameter Optimization
18 conf.set("mapreduce.task.io.sort.mb", "512"); // Increase sort buffer to
19 // reduce spills
20 conf.set("mapreduce.reduce.shuffle.parallelcopies", "10");
```

5.2 Data Parsing Acceleration: Reconstructing Parsing Logic and Reducing Redundant Operations

The core goal is to solve the problems of long CSV parsing time and excessive temporary objects, improving efficiency by optimizing parsing methods and reducing data traversal times.

- **Parsing Logic Reconstruction**

- *On-Demand Parsing Instead of Full Splitting*: At first, we abandon `String.split(",")` for full-column splitting (57 columns) and adopts pointer linear scanning, parsing only algorithm-dependent fields such as `tradingDay` and `tradeTime`, while skipping irrelevant fields directly.

```

1  /**
2   * Parses a CSV line to populate the object's fields.
3   * @param line Single line from a CSV file
4   * @return true if parsing succeeds, false if any error occurs or line
5   *         is malformed
6   */
7  public boolean parseFromCSV(String line) {
8      int pos = 0;
9      int length = line.length();
10
11     try {
12         // Parse 0: tradingDay (YYYYMMDD, 8 numeric digits)
13         if (pos + 8 > length) return false;
14         tradingDay = Integer.parseInt(line.substring(pos, pos + 8));
15         pos += 9; // Skip digits and comma
16
17         // Parse 1: tradeTime (HHMMSS or with microseconds, take first
18         // 6 digits)
19         if (pos + 6 > length) return false;
20         String timeStr = line.substring(pos, pos + 6);
21         tradeTime = Long.parseLong(timeStr);
22         // Skip remaining time part and comma
23         while (pos < length && line.charAt(pos) != ',') pos++;
24         pos++;
25
26         // Skip fields 2-11 (unneeded columns)
27         for (int i = 0; i < 10; i++) {
28             while (pos < length && line.charAt(pos) != ',') pos++;
29             pos++;
30         }
31
32         // Parse 12: tBidVol
33         int start = pos;
34         while (pos < length && line.charAt(pos) != ',') pos++;
35         tBidVol = Long.parseLong(line.substring(start, pos));
36         pos++;
37
38         // Parse 13: tAskVol
39         start = pos;
40         while (pos < length && line.charAt(pos) != ',') pos++;
41         tAskVol = Long.parseLong(line.substring(start, pos));
42         pos++;
43
44         // Skip fields 14-16
45         for (int i = 0; i < 3; i++) {
46             while (pos < length && line.charAt(pos) != ',') pos++;
47             pos++;
48         }
49
50         // Parse 17-36: top 5 levels (bp/bv/ap/av)
51         for (int i = 0; i < 5; i++) {
52             // bidPrice[i]
53             start = pos;
54             while (pos < length && line.charAt(pos) != ',') pos++;
55             bidPrices[i] = Long.parseLong(line.substring(start, pos));
56             pos++;
57
58             // bidVolume[i]
59             start = pos;
60             while (pos < length && line.charAt(pos) != ',') pos++;
61             bidVolumes[i] = Long.parseLong(line.substring(start, pos));
62             ;
63             pos++;
64         }
65     }
66 }

```

```

62         // askPrice[i]
63         start = pos;
64         while (pos < length && line.charAt(pos) != ',') pos++;
65         askPrices[i] = Long.parseLong(line.substring(start, pos));
66         pos++;
67
68         // askVolume[i]
69         start = pos;
70         while (pos < length && line.charAt(pos) != ',') pos++;
71         askVolumes[i] = Long.parseLong(line.substring(start, pos))
72         ;
73         pos++;
74     }
75
76     return true;
77 } catch (NumberFormatException | IndexOutOfBoundsException e) {
78     return false;
79 }

```

Then furtherly use `String.indexOf(',', pos)` to locate separators in batches, reducing the time complexity of search from $O(n)$ to $O(1)$, which significantly improves the efficiency of long-line parsing.

```

1  /**
2   * Parses a single CSV line to populate the object's fields.
3   * @param line A single line from a CSV file
4   * @return true if parsing is successful; false if the line is
5   *         malformed or an error occurs
6   */
7  public boolean parseFromCSV(String line) {
8      int pos = 0;
9      final int length = line.length();
10     int nextPos;
11
12     try {
13         // Parse 0: tradingDay (YYYYMMDD, 8 numeric digits)
14         if (pos + 8 > length) return false;
15         tradingDay = parseInt(line, pos, pos + 8);
16         pos += 8;
17         if (pos >= length || line.charAt(pos) != ',') return false;
18         pos++;
19
20         // Parse 1: tradeTime (take first 6 digits as HHMMSS)
21         if (pos + 6 > length) return false;
22         tradeTime = parseLong(line, pos, pos + 6);
23         // Skip remaining time part and comma
24         nextPos = line.indexOf(',', pos);
25         if (nextPos == -1) return false;
26         pos = nextPos + 1;
27
28         // Skip fields 2-11 (10 unneeded columns total)
29         for (int i = 0; i < 10; i++) {
30             nextPos = line.indexOf(',', pos);
31             if (nextPos == -1) return false;
32             pos = nextPos + 1;
33         }
34
35         // Parse 12: tBidVol
36         nextPos = line.indexOf(',', pos);
37         if (nextPos == -1) return false;
38         tBidVol = parseLong(line, pos, nextPos);
39         pos = nextPos + 1;
40
41         // Parse 13: tAskVol

```

```

41     nextPos = line.indexOf(',', pos);
42     if (nextPos == -1) return false;
43     tAskVol = parseLong(line, pos, nextPos);
44     pos = nextPos + 1;
45
46     // Skip fields 14-16 (3 fields total)
47     for (int i = 0; i < 3; i++) {
48         nextPos = line.indexOf(',', pos);
49         if (nextPos == -1) return false;
50         pos = nextPos + 1;
51     }
52
53     // Parse 17-36: top 5 order book levels (bp/bv/ap/av)
54     for (int i = 0; i < 5; i++) {
55         // bidPrice[i]
56         nextPos = line.indexOf(',', pos);
57         if (nextPos == -1) return false;
58         bidPrices[i] = parseLong(line, pos, nextPos);
59         pos = nextPos + 1;
60
61         // bidVolume[i]
62         nextPos = line.indexOf(',', pos);
63         if (nextPos == -1) return false;
64         bidVolumes[i] = parseLong(line, pos, nextPos);
65         pos = nextPos + 1;
66
67         // askPrice[i]
68         nextPos = line.indexOf(',', pos);
69         if (nextPos == -1) return false;
70         askPrices[i] = parseLong(line, pos, nextPos);
71         pos = nextPos + 1;
72
73         // askVolume[i]
74         nextPos = line.indexOf(',', pos);
75         if (nextPos == -1) return false;
76         askVolumes[i] = parseLong(line, pos, nextPos);
77         pos = nextPos + 1;
78     }
79
80     return true;
81 } catch (IndexOutOfBoundsException e) {
82     return false;
83 }
84 }

```

- *Efficient Parsing of Fixed-Length Fields*: For fixed-length fields such as YYYYMMDD (8 bits) and HHMMSS (6 bits), we adopt direct substring interception or the SWAR algorithm to eliminate loops and branch judgments, processing values in parallel through arithmetic operations and avoiding per-character traversal overhead.

```

1  import java.nio.charset.StandardCharsets;
2
3  /**
4   * High-Speed CSV Parsing Utility Class
5   * Optimization Strategies:
6   * 1. Fused Scan-Parse: Calculate numeric values while searching for
7   *    delimiters, reducing memory access by 50%
8   * 2. Manual Loop Unrolling: Explicitly unroll loops for fixed-length
9   *    fields to boost Instruction-Level Parallelism (ILP)
10  */
11  public class FastParser {
12      /**
13       * [SWAR Strategy] Parse 8-digit fixed-length date (e.g.,
14       *    20140101)
15       * @param b Byte array containing the data

```

```

13      * @param offset Starting position in the byte array
14      * @return Parsed integer value of the date
15      */
16      public static int parseDate8(byte[] b, int offset) {
17          // Direct calculation (no loops) - leverage CPU pipeline for
18          // parallel multiplication
19          return (b[offset] - '0') * 10000000 +
20                 (b[offset + 1] - '0') * 1000000 +
21                 (b[offset + 2] - '0') * 100000 +
22                 (b[offset + 3] - '0') * 10000 +
23                 (b[offset + 4] - '0') * 1000 +
24                 (b[offset + 5] - '0') * 100 +
25                 (b[offset + 6] - '0') * 10 +
26                 (b[offset + 7] - '0');
27      }
28
29      /**
30       * Parse a long integer until a comma or end of line is
31       * encountered
32       * @param b Data source byte array
33       * @param cursor Integer array holding the [current position] (
34       *   updated to position after comma upon completion)
35       * @return Parsed long integer value
36       */
37      public static long parseLong(byte[] b, int[] cursor) {
38          long result = 0;
39          int i = cursor[0];
40          int len = b.length;
41
42          // Handle possible negative sign (retained for robustness,
43          // even though quant data is mostly positive)
44          boolean negative = false;
45          if (i < len && b[i] == '-') {
46              negative = true;
47              i++;
48          }
49
50          // Tight loop to eliminate redundant branches
51          while (i < len) {
52              byte c = b[i++];
53              if (c == ',') {
54                  cursor[0] = i; // Update cursor to position after
55                  // comma
56                  return negative ? -result : result;
57              }
58              // Accumulate value: result * 10 + digit
59              result = result * 10 + (c - '0');
60          }
61
62          // Handle case where line ends without a comma
63          cursor[0] = i;
64          return negative ? -result : result;
65      }
66
67      /**
68       * Parse a string until a comma is encountered
69       * @param b Data source byte array
70       * @param cursor Integer array holding the [current position] (
71       *   updated to position after comma upon completion)
72       * @return Parsed string value
73       */
74      public static String parseString(byte[] b, int[] cursor) {
75          int start = cursor[0];
76          int i = start;
77          int len = b.length;

```



```

72     while (i < len) {
73         if (b[i++] == ',') {
74             // String object creation is only needed here (
75             // mandatory for use as Map Key)
76             cursor[0] = i;
77             // Explicit charset specification (ISO_8859_1/UTF_8)
78             // is faster than default
79             return new String(b, start, i - 1 - start,
80                             StandardCharsets.UTF_8);
81         }
82     }
83     cursor[0] = i;
84     return new String(b, start, i - start, StandardCharsets.UTF_8);
85 }
86
87 /**
88  * Fast skip of N CSV fields
89  * @param b Data source byte array
90  * @param cursor Integer array holding the [current position] (
91  *    updated to position after N commas upon completion)
92  * @param count Number of fields to skip
93  */
94 public static void skipFields(byte[] b, int[] cursor, int count) {
95     int i = cursor[0];
96     int len = b.length;
97     int found = 0;
98
99     while (i < len && found < count) {
100         if (b[i++] == ',') {
101             found++;
102         }
103     }
104     cursor[0] = i;
105 }

```

- *Zero-Intermediate-Object Parsing*: Above method proposes the FastParser strategy, directly operating on the byte array of Hadoop Text objects and adopting "Fused Scan-Parse" ($O(N)$ complexity) to complete value calculation while searching for separators. Then we further parses stock codes (String type) directly into int type, completely eliminating String temporary objects during parsing.

```

1  /**
2   * Parses a stock code from a byte array into an integer, ignoring
3   * suffixes like .SH/.SZ (Shanghai/Shenzhen exchanges).
4   * The parsing stops at non-numeric characters (e.g., '.' or ',') and
5   * updates the cursor to the position after the comma.
6   *
7   * @param data The source byte array containing the CSV data
8   * @param cursor Integer array holding the [current position] (updated
9   *    to position after the comma upon completion)
10  * @return The numeric part of the stock code as an integer
11  */
12 public static int parseStockCodeToInt(byte[] data, int[] cursor) {
13     int idx = cursor[0];
14     int result = 0;
15
16     while (idx < data.length) {
17         byte b = data[idx];
18         if (b >= '0' && b <= '9') {
19             result = result * 10 + (b - '0');
20             idx++;
21         } else {
22             // Encounter .SH/.SZ suffix or comma directly

```

```

20         // If it's '.', skip all characters until the next comma
21         if (b != ',') {
22             while (idx < data.length && data[idx] != ',') {
23                 idx++;
24             }
25         }
26         break;
27     }
28 }
29 cursor[0] = idx + 1; // Skip the comma
30 return result;
31 }

```

- **Parsing Detail Optimization:** We reduces redundant computations and performance losses during parsing by caching separator positions with local variables, structuring field parsing (processing fields in blocks according to business logic), and streamlining the scope of exception catching (only catching `IndexOutOfBoundsException`).

5.3 Computational Efficiency Improvement: Simplifying Computational Paths and Reducing CPU Overhead

The core goal is to optimize the hot path of factor calculation, improving CPU execution efficiency by reducing loops, branch judgments, and inefficient instructions.

- **Loop and Branch Optimization**

- *Loop Unrolling:* Firstly we manually unroll the summation and weighted summation logic for the first 5 levels to eliminate the boundary checking and array addressing overhead of `for` loops.

```

1  /**
2   * Calculates 20 factor values
3   * @param snapshot Market snapshot data (current moment)
4   * @param prevSnapshot Market snapshot data from the previous moment (
5   *   used to calculate change-based factors)
6   * @return Array containing 20 factor values
7   */
8  // Modify the calculation method in Factor.java
9  public static double[] calculateFactors(SnapshotData snapshot,
10     SnapshotData prevSnapshot) {
11     double[] factors = new double[20];
12     final double EPS = EPSILON;
13
14     // Extract base data (use array indices directly to avoid getter
15     // invocation overhead)
16     long[] bp = snapshot.getBidPrices();
17     long[] bv = snapshot.getBidVolumes();
18     long[] ap = snapshot.getAskPrices();
19     long[] av = snapshot.getAskVolumes();
20     long tBidVol = snapshot.gettBidVol();
21     long tAskVol = snapshot.gettAskVol();
22
23     // Manually unroll calculation for top 5 levels (avoid for-loop
24     // overhead)
25     long sumBidVol = bv[0] + bv[1] + bv[2] + bv[3] + bv[4];
26     long sumAskVol = av[0] + av[1] + av[2] + av[3] + av[4];
27
28     double weightedBidPrice = bp[0]*bv[0] + bp[1]*bv[1] + bp[2]*bv[2]
29     + bp[3]*bv[3] + bp[4]*bv[4];
30     double weightedAskPrice = ap[0]*av[0] + ap[1]*av[1] + ap[2]*av[2]
31     + ap[3]*av[3] + ap[4]*av[4];
32
33     // Precompute reciprocals (reduce division operations)
34     double invSumBid = 1.0 / (sumBidVol + EPS);
35     double invSumAsk = 1.0 / (sumAskVol + EPS);

```

```

30     double invTotalVol = 1.0 / (sumBidVol + sumAskVol + EPS);
31     double invTAskVol = 1.0 / (tAskVol + EPS);
32     double invTotalTAsk = 1.0 / (tBidVol + tAskVol + EPS);
33
34     // Mid price (reuse calculation result)
35     double midPrice = (ap[0] + bp[0]) / 2.0;
36
37     // Calculate static factors (manually unroll to reduce branches)
38     factors[0] = ap[0] - bp[0]; // Best bid-ask spread
39     factors[1] = (ap[0] - bp[0]) / (midPrice + EPS); // Relative bid-
40         ask spread
41     factors[2] = midPrice; // Mid price
42     factors[3] = (bv[0] - av[0]) * 1.0 / (bv[0] + av[0] + EPS); //
43         Order book imbalance (top level)
44     factors[4] = (sumBidVol - sumAskVol) * invTotalVol; // Order book
45         imbalance (multi-level)
46     factors[5] = sumBidVol; // Total bid volume
47     factors[6] = sumAskVol; // Total ask volume
48     factors[7] = sumBidVol - sumAskVol; // Net volume (bid - ask)
49     factors[8] = sumBidVol * invSumAsk; // Depth ratio (bid/ask)
50     factors[9] = (tBidVol - tAskVol) * invTotalTAsk; // Total traded
51         volume balance
52     factors[10] = weightedBidPrice * invSumBid; // Weighted bid price
53         (per bid volume)
54     factors[11] = weightedAskPrice * invSumAsk; // Weighted ask price
55         (per ask volume)
56     factors[12] = (weightedBidPrice + weightedAskPrice) * invTotalVol;
57         // Total weighted price
58     factors[13] = factors[11] - factors[10]; // Weighted bid-ask
59         spread
60     factors[14] = (sumBidVol - sumAskVol) / 5.0; // Average level
61         imbalance
62     factors[15] = calculateAsymmetryUnrolled(bv, av); // Asymmetry
63         factor (unrolled implementation)
64     factors[19] = (ap[0] - bp[0]) * invTotalVol; // Spread-to-total-
65         volume ratio
66
67     // Calculate change-based factors
68     calculateChangeFactorsUnrolled(factors, snapshot, prevSnapshot,
69         midPrice, sumBidVol, sumAskVol);
70
71     return factors;
72 }
73
74 // Unrolled implementation of asymmetry factor calculation (avoid for-
75 loop)
76 private static double calculateAsymmetryUnrolled(long[] bv, long[] av)
77 {
78     double wBid = bv[0] * 1.0 + bv[1] * 0.5 + bv[2] * 0.333 + bv[3] *
79         0.25 + bv[4] * 0.2;
80     double wAsk = av[0] * 1.0 + av[1] * 0.5 + av[2] * 0.333 + av[3] *
81         0.25 + av[4] * 0.2;
82     return (wBid - wAsk) / (wBid + wAsk + EPSILON);
83 }
84
85 // Unrolled implementation of change factor calculation
86 private static void calculateChangeFactorsUnrolled(double[] factors,
87     SnapshotData snapshot, SnapshotData prevSnapshot, double
88     currentMidPrice, long currentSumBid, long currentSumAsk) {
89     if (prevSnapshot == null) {
90         factors[16] = 0;
91         factors[17] = 0;
92         factors[18] = 0;
93         return;
94     }
95 }

```

```

77     long[] prevAp = prevSnapshot.getAskPrices();
78     long[] prevBp = prevSnapshot.getBidPrices();
79     long[] prevBv = prevSnapshot.getBidVolumes();
80     long[] prevAv = prevSnapshot.getAskVolumes();
81
82     // Factor 16: Best ask price change
83     factors[16] = snapshot.getAskPrices()[0] - prevAp[0];
84
85     // Factor 17: Mid price change
86     double prevMid = (prevAp[0] + prevBp[0]) / 2.0;
87     factors[17] = currentMidPrice - prevMid;
88
89     // Factor 18: Depth ratio change (manually calculate sum of top 5
90     // levels)
91     long prevSumBid = prevBv[0] + prevBv[1] + prevBv[2] + prevBv[3] +
92     prevBv[4];
93     long prevSumAsk = prevAv[0] + prevAv[1] + prevAv[2] + prevAv[3] +
94     prevAv[4];
95     double prevDepthRatio = prevSumBid / (double) (prevSumAsk +
96     EPSILON);
97     factors[18] = (currentSumBid / (double) (currentSumAsk + EPSILON))
98     - prevDepthRatio;
99 }

```

Then manually unroll the accumulation loop for 20 factors, increasing aggregation speed by 10-15%. At last, we further "peels" array fields into independent fields (e.g., bp0, bp1 instead of long[] bidPrices), eliminating array bounds checking instructions.

```

1  /**
2   * Ultra-Optimized Calculation Method
3   * 1. Access flattened fields (s.bp0)
4   * 2. Convert division to multiplication (Reciprocal technique)
5   */
6  public void calculateFrom(SnapshotData s, SnapshotData p) {
7      // --- 1. Basic Aggregation (Register-level variables) ---
8      long sumBidVol = s.bv0 + s.bv1 + s.bv2 + s.bv3 + s.bv4;
9      long sumAskVol = s.av0 + s.av1 + s.av2 + s.av3 + s.av4;
10     long totalVol = sumBidVol + sumAskVol;
11     long totalTVol = s.tBidVol + s.tAskVol;
12
13     // --- 2. Precompute Reciprocals (Division -> Multiplication) ---
14     // CPU takes ~20-40 cycles for FDIV, only 3-5 cycles for FMUL
15     float invSumBid = 1.0f / (sumBidVol + EPSILON);
16     float invSumAsk = 1.0f / (sumAskVol + EPSILON);
17     float invTotalVol = 1.0f / (totalVol + EPSILON);
18     float invTotalTVol = 1.0f / (totalTVol + EPSILON);
19
20     double wBidPrice = s.bp0*s.bv0 + s.bp1*s.bv1 + s.bp2*s.bv2 + s.bp3
21     *s.bv3 + s.bp4*s.bv4;
22     double wAskPrice = s.ap0*s.av0 + s.ap1*s.av1 + s.ap2*s.av2 + s.ap3
23     *s.av3 + s.ap4*s.av4;
24
25     float midPrice = (float)((s.ap0 + s.bp0) * 0.5); // *0.5 is faster
26     // than /2.0
27
28     // --- 3. Populate Factors (Use reciprocal multiplication) ---
29     float[] f = this.factorValues; // Local reference (avoids repeated
30     // field access)
31     f[0] = s.ap0 - s.bp0;
32     f[1] = f[0] / (midPrice + EPSILON); // Retained division here:
33     // midPrice is volatile (not suitable for precomputation)
34     f[2] = midPrice;
35     f[3] = (s.bv0 - s.av0) * 1.0f / (s.bv0 + s.av0 + EPSILON);

```

```

32
33 // Optimization: Division converted to multiplication
34 f[4] = (sumBidVol - sumAskVol) * invTotalVol;
35 f[5] = sumBidVol;
36 f[6] = sumAskVol;
37 f[7] = sumBidVol - sumAskVol;
38 f[8] = sumBidVol * invSumAsk;
39 f[9] = (s.tBidVol - s.tAskVol) * invTotalTVol;
40
41 f[10] = (float)(wBidPrice * invSumBid);
42 f[11] = (float)(wAskPrice * invSumAsk);
43 f[12] = (float)((wBidPrice + wAskPrice) * invTotalVol);
44 f[13] = f[11] - f[10];
45 f[14] = f[7] * 0.2f;
46
47 f[15] = calculateAsymmetry(s);
48 f[19] = f[0] * invTotalVol;
49
50 // Change factors (leverage the Dummy object's tradeTime=0
51 // property)
52 if (p.tradeTime != 0) {
53     calculateChangeFactors(s, p, midPrice, sumBidVol, sumAskVol);
54 } else {
55     f[16] = 0;
56     f[17] = 0;
57     f[18] = 0;
58 }
59
60 private float calculateAsymmetry(SnapshotData s) {
61     // Use multiplication instead of division
62     double wBid = s.bv0 + s.bv1*0.5 + s.bv2*0.333333 + s.bv3*0.25 + s.
63         bv4*0.2;
64     double wAsk = s.av0 + s.av1*0.5 + s.av2*0.333333 + s.av3*0.25 + s.
65         av4*0.2;
66     return (float)((wBid - wAsk) / (wBid + wAsk + EPSILON));
67 }
68
69 private void calculateChangeFactors(SnapshotData c, SnapshotData p,
70     float currMid, long currSumBid, long currSumAsk) {
71     float[] f = this.factorValues;
72
73     f[16] = c.ap0 - p.ap0;
74
75     float prevMid = (float)((p.ap0 + p.bp0) * 0.5);
76     f[17] = currMid - prevMid;
77
78     long prevSumBid = p.bv0 + p.bv1 + p.bv2 + p.bv3 + p.bv4;
79     long prevSumAsk = p.av0 + p.av1 + p.av2 + p.av3 + p.av4;
80
81     float prevDepthRatio = (float)(prevSumBid / (double)(prevSumAsk +
82         EPSILON));
83     float currDepthRatio = (float)(currSumBid / (double)(currSumAsk +
84         EPSILON));
85
86     f[18] = currDepthRatio - prevDepthRatio;
87 }

```

- *Branch Simplification*: At first we use a unified $\text{EPSILON}=10^{-7}$ to handle cases where the denominator is 0, eliminating the `if(denom==0)` branch. Then directly assign default values through explicit null checks to avoid complex branches. Finally introduces a `DUMMY_SNAPSHOT` zero-filled object instead of null, removing the `if(prev==null)` check on the hot path and allowing the CPU pipeline to run at full speed.

- **Replacement of Inefficient Instructions**: Multiple optimizations adopt the "division-to-multiplication"

strategy. For frequently used denominators such as `sumBidVol`, the reciprocal is precomputed (e.g., $\text{invSum} = 1/(\text{sum} + \text{EPS})$), and repeated divisions are replaced with multiplications (the CPU overhead of division is 4-8 times that of multiplication), theoretically increasing computational throughput by 3-5 times.

- **Reuse of Intermediate Results:** In initial, we precompute intermediate results such as mid-price (`midPrice`) and reciprocals (`invSumBid`), which are directly reused in subsequent calculations to avoid redundant operations. Then remove the `ThreadLocal` intermediate cache layer and adds the `calculateFrom` method to write computation results directly to the target array, reducing one array traversal and copy operation.

5.4 Data Transmission Compression: Reducing Overhead in the Shuffle Phase

The core goal is to reduce data transmission volume and sorting overhead during the process from Map output to Reduce, achieving speedups through Key compression and serialization optimization.

• Key Structure Optimization and Compression

- *Timestamp Compression:* Firstly, compress the trading day (YYYYMMDD) and trading time (HH-MMSS) into a 32-bit integer (4 bytes) through bitwise operations using the `CompactTimeUtil` tool class, replacing the original `Text`-type Key. Then adopt this compression method, using compact timestamps as aggregation or output Keys and further compress into a 26-bits integer.

```

1  /**
2   * Ultra-Fast Time Encoding Utility Class (26-bit Compressed Version)
3   */
4  public class CompactTimeUtil {
5
6      // ----- Configuration Parameters -----
7      // Base year (adjust based on actual data year; 12-bit date
8      //      supports +11 years from this year)
9      private static final int BASE_YEAR = 2014;
10
11     // Bit shift constants
12     private static final int DATE_SHIFT = 14;
13     private static final int TIME_MASK = 0x3FFF; // 14-bit mask
14     //      (16383)
15     private static final int DATE_MASK = 0xFFF; // 12-bit mask (4095)
16
17     // Trading session second definitions (for time mapping)
18     private static final int AM_START_SEC = 9 * 3600; //
19     //      09:00:00
20     private static final int AM_END_SEC = 11 * 3600 + 30 * 60; //
21     //      11:30:00
22     private static final int PM_START_SEC = 13 * 3600; //
23     //      13:00:00
24     private static final int PM_END_SEC = 15 * 3600; //
25     //      15:00:00
26
27     private static final int AM_DURATION = AM_END_SEC - AM_START_SEC;
28     //      9000 seconds
29
30     /**
31     * Encode to a 26-bit integer
32     * @param tradingDay Trading day in YYYYMMDD format
33     * @param tradeTime Trading time in HHMMSS format
34     * @return 26-bit compact time integer
35     */
36     public static int encode(int tradingDay, int tradeTime) {
37         // 1. Calculate date offset (sparse mapping, extremely high
38         //      performance)
39         int year = tradingDay / 10000;
40         int month = (tradingDay % 10000) / 100;
41         int day = tradingDay % 100;
42     }
43 }

```

```

35     int yearOffset = year - BASE_YEAR;
36     if (yearOffset < 0 || yearOffset > 11) {
37         throw new IllegalArgumentException("Year out of range;
            current base supports " + BASE_YEAR + "-" + (BASE_YEAR
            + 11));
38     }
39
40     // Use sparse mapping of "year*372 + month*31 + day" to avoid
        complex calendar calculations
41     // 372 = 12 * 31 (ensures month and day are reversible and
        monotonic)
42     int dateCode = yearOffset * 372 + (month - 1) * 31 + (day - 1)
        ;
43     if (dateCode > DATE_MASK) {
44         throw new IllegalArgumentException("Date out of encoding
            range");
45     }
46
47     // 2. Calculate time index (HHMMSS -> seconds -> mapped index)
48     int totalSeconds = convertToSecOfDay(tradeTime);
49     int timeCode = mapSecondsToCode(totalSeconds);
50
51     // 3. Combine date and time codes
52     return (dateCode << DATE_SHIFT) | timeCode;
53 }
54
55 /**
56  * Convert HHMMSS to absolute seconds of the day
57  * @param time Time in HHMMSS format
58  * @return Absolute seconds since midnight (00:00:00)
59  */
60 public static int convertToSecOfDay(int time) {
61     int hour = time / 10000;
62     int minute = (time % 10000) / 100;
63     int second = time % 100;
64     return hour * 3600 + minute * 60 + second;
65 }
66
67 /**
68  * Core compression logic: Map absolute seconds to a continuous
        compact index
69  * Skip the midday break (11:30-13:00)
70  * @param sec Absolute seconds of the day
71  * @return Compact time index
72  */
73 private static int mapSecondsToCode(int sec) {
74     if (sec < AM_START_SEC) return 0; // Pre-market time mapped to
        0 uniformly
75
76     if (sec <= AM_END_SEC) {
77         return sec - AM_START_SEC; // Morning session mapped
            directly
78     } else if (sec < PM_START_SEC) {
79         return AM_DURATION; // Midday break mapped to the end of
            morning session uniformly
80     } else if (sec <= PM_END_SEC) {
81         return AM_DURATION + (sec - PM_START_SEC) + 1; //
            Afternoon session follows morning session
82     } else {
83         // Post-market time mapped to the maximum value (~16200)
84         return AM_DURATION + (PM_END_SEC - PM_START_SEC) + 1;
85     }
86 }
87
88 /**

```

```

89      * Decode the trading day (YYYYMMDD) from compact time
90      * @param compactTime 26-bit compact time integer
91      * @return Trading day in YYYYMMDD format
92      */
93      public static int decodeTradingDay(int compactTime) {
94          int dateCode = (compactTime >>> DATE_SHIFT) & DATE_MASK;
95
96          int yearOffset = dateCode / 372;
97          int remainder = dateCode % 372;
98          int monthIndex = remainder / 31;
99          int dayIndex = remainder % 31;
100
101          int year = yearOffset + BASE_YEAR;
102          int month = monthIndex + 1;
103          int day = dayIndex + 1;
104
105          return year * 10000 + month * 100 + day;
106      }
107
108      /**
109       * Decode the trading time (HHMMSS) from compact time
110       * @param compactTime 26-bit compact time integer
111       * @return Trading time in HHMMSS format
112       */
113      public static int decodeTradeTime(int compactTime) {
114          int timeCode = compactTime & TIME_MASK;
115          int secOfDay;
116
117          if (timeCode <= AM_DURATION) {
118              secOfDay = AM_START_SEC + timeCode;
119          } else {
120              secOfDay = PM_START_SEC + (timeCode - AM_DURATION - 1);
121          }
122
123          int hour = secOfDay / 3600;
124          int remaining = secOfDay % 3600;
125          int minute = remaining / 60;
126          int second = remaining % 60;
127
128          return hour * 10000 + minute * 100 + second;
129      }
130
131      // Backward-compatible helper methods
132      public static int getDayCode(int compactTime) {
133          return compactTime >>> DATE_SHIFT;
134      }
135
136      public static String getMMDD(int compactTime) {
137          int dateCode = (compactTime >>> DATE_SHIFT) & DATE_MASK;
138          int remainder = dateCode % 372;
139          int month = remainder / 31 + 1;
140          int day = remainder % 31 + 1;
141          return String.format("%02d%02d", month, day);
142      }
143
144      // Preserve comparison logic unchanged
145      public static int compare(int a, int b) {
146          return Integer.compareUnsigned(a, b);
147      }
148  }

```

- *Key Type Optimization*: Change the Mapper output Key from Text to IntWritable. Integer comparison is a native numerical operation, which is an order of magnitude more efficient than string character-by-character comparison, significantly reducing the time consumption of the Sort phase.


```

1 // Setup output type of Mapper (Use IntWritable as Key)
2 job.setMapOutputKeyClass(IntWritable.class);

```

• Data Transmission and Serialization Optimization

- *Pre-aggregation of Data*: Add a Combiner to pre-merge Factor data with the same Key at the Map side.

```

1 /**
2  * Combiner: Used in Map reducing shuffle transmission
3  */
4 public class FactorCombiner extends Reducer<IntWritable, Factor,
5     IntWritable, Factor> {
6     @Override
7     protected void reduce(IntWritable key, Iterable<Factor> values,
8         Context context)
9         throws IOException, InterruptedException {
10         Factor sumFactor = new Factor();
11         for (Factor factor : values) {
12             sumFactor.merge(factor);
13         }
14         context.write(key, sumFactor);
15     }
16 }

```

Then further aggregate multi-stock data by compact timestamps at the Mapper phase, reducing the number of map-output records and significantly reducing Shuffle data transmission volume.

```

1 // Pre-aggregation: Sum factors of all stocks under the same timestamp
2 Factor currentFactor = factorCache.get(compactTime);
3 if (currentFactor == null) {
4     currentFactor = new Factor(factors, 1);
5     factorCache.put(compactTime, currentFactor);
6 } else {
7     currentFactor.merge(new Factor(factors, 1));
8 }

```

- *Serialization Compression*: At first we change `double[]` in `Factor` to `float[]`, reducing the serialization volume by 50%. Then make `Factor` implement the `Writable` interface to directly read and write `float[]` and `int`, avoiding reflection serialization overhead. Besides, we also reduce data skew by optimizing the partitioner's hash algorithm:, preventing overload of individual Reducers.

```

1 /**
2  * Partitions data by trading day (extracted from compact time key)
3  */
4 public class DayPartitioner extends Partitioner<IntWritable, Factor> {
5     @Override
6     public int getPartition(IntWritable key, Factor value, int
7         numPartitions) {
8         // Use integer key directly (avoids string parsing overhead)
9         int compactTime = key.get();
10        int dayCode = CompactTimeUtil.getDayCode(compactTime);
11
12        // Apply a more uniform hashing algorithm to reduce data skew
13        return (dayCode ^ (dayCode >>> 16)) % numPartitions;
14    }
15 }

```

- *Reduce Floating-Point Formatting Implementation*: For floating-point data output scenarios in the Reduce phase, we incorporate the Schubfach/Ryu algorithm for floating-point formatting processing. Traditional floating-point formatting methods often suffer from precision loss or inefficiency. In contrast, the Schubfach/Ryu algorithm optimizes the floating-point-to-string conversion logic: it drastically improves formatting speed while ensuring precision, reduces the time consumption of

data output in the Reduce phase, and serves as a supplementary optimization to boost the overall job performance.

```
1  /**
2   * Ultra-Performant Float-to-Byte Conversion Utility (Based on
3   * Simplified Ryu Algorithm)
4   */
5  public class RyuFloat {
6      private static final int FLOAT_MANTISSA_BITS = 23;
7      private static final int FLOAT_MANTISSA_MASK = (1 <<
8          FLOAT_MANTISSA_BITS) - 1;
9      private static final int FLOAT_EXPONENT_BITS = 8;
10     private static final int FLOAT_EXPONENT_BIAS = 127;
11     private static final int FLOAT_EXPONENT_MASK = (1 <<
12         FLOAT_EXPONENT_BITS) - 1;
13
14     // Precomputed power-of-10 table (for fast decimal point
15     // positioning)
16
17     /**
18      * Formats a float value into a byte array
19      * @param value Input floating-point number
20      * @param result Target byte array
21      * @param index Starting position for writing
22      * @return New index position after writing
23      */
24     public static int floatToBytes(float value, byte[] result, int
25         index) {
26         // 1. Handle special values
27         int bits = Float.floatToIntBits(value);
28         boolean negative = (bits & 0x80000000) != 0;
29         int exponent = (bits >>> FLOAT_MANTISSA_BITS) &
30             FLOAT_EXPONENT_MASK;
31         int mantissa = bits & FLOAT_MANTISSA_MASK;
32
33         if (exponent == 255) {
34             if (mantissa == 0) {
35                 if (negative) result[index++] = '-';
36                 result[index++] = 'I'; result[index++] = 'n'; result[
37                     index++] = 'f';
38                 return index;
39             } else {
40                 result[index++] = 'N'; result[index++] = 'a'; result[
41                     index++] = 'N';
42                 return index;
43             }
44         }
45
46         if (exponent == 0 && mantissa == 0) {
47             if (negative) result[index++] = '-';
48             result[index++] = '0';
49             result[index++] = '.';
50             result[index++] = '0';
51             return index;
52         }
53
54         if (negative) {
55             result[index++] = '-';
56         }
57
58         // 2. Core conversion logic
59         // To ensure high performance without introducing several KB
60         // of lookup table code, we use "long-based fixed-point
61         // arithmetic" to handle the most common ranges.
```

```

53     if (Math.abs(value) >= 0.001f && Math.abs(value) <= 10000000f)
54     {
55         return formatNormal(Math.abs(value), result, index);
56     } else {
57         // Fallback to standard conversion for extreme values (
58         // avoids large lookup tables)
59         String s = Float.toString(Math.abs(value));
60         for (int i = 0; i < s.length(); i++) {
61             result[index++] = (byte) s.charAt(i);
62         }
63         return index;
64     }
65 }
66
67 // High-performance formatting for normal ranges (no object
68 // allocation)
69 private static int formatNormal(float val, byte[] buf, int offset)
70 {
71     // Convert float to integer processing
72     int iPart = (int) val;
73     float fPart = val - iPart;
74
75     // Write integer part
76     if (iPart == 0) {
77         buf[offset++] = '0';
78     } else {
79         offset = writeInt(iPart, buf, offset);
80     }
81
82     buf[offset++] = '.';
83
84     // Write decimal part (fixed-precision optimization, or until
85     // zero)
86     if (fPart == 0) {
87         buf[offset++] = '0';
88         return offset;
89     }
90
91     // Extract decimal places: Multiply by 100,000,000 (10^8) to
92     // convert to long
93     long fScaled = (long) (fPart * 100000000L + 0.5);
94
95     // Add leading zeros if fScaled is too small
96     int digits = 0;
97     if (temp == 0) {
98         digits = 1;
99     } else {
100         // Simple digit count
101         while (temp > 0) {
102             temp /= 10;
103             digits++;
104         }
105     }
106
107     // Add zeros (8 - digits)
108     for (int k = 0; k < (8 - digits); k++) {
109         buf[offset++] = '0';
110     }
111
112     // Trim trailing zeros
113     while (fScaled > 0 && fScaled % 10 == 0) {
114         fScaled /= 10;
115     }
116
117     if (fScaled == 0) {

```

```

112         // If all zeros (e.g., fPart is extremely small)
113         buf[offset++] = '0';
114     } else {
115         offset = writeLong(fScaled, buf, offset);
116     }
117
118     return offset;
119 }
120
121 // Fast integer-to-byte conversion (most significant to least
122 // significant digit)
123 private static int writeInt(int value, byte[] buf, int offset) {
124     if (value == 0) {
125         buf[offset++] = '0';
126         return offset;
127     }
128     int i = offset;
129     int q, r;
130     // Calculate length first
131     int temp = value;
132     int len = 0;
133     while (temp > 0) {
134         temp /= 10;
135         len++;
136     }
137     i += len;
138     int end = i;
139
140     // Fill in reverse order
141     while (value > 0) {
142         q = value / 10;
143         r = value - (q * 10); // Equivalent to value % 10 (faster
144                             // on some CPUs)
145         buf[--i] = (byte) (r + '0');
146         value = q;
147     }
148     return end;
149 }
150
151 private static int writeLong(long value, byte[] buf, int offset) {
152     if (value == 0) {
153         buf[offset++] = '0';
154         return offset;
155     }
156     int i = offset;
157     long q, r;
158     long temp = value;
159     int len = 0;
160     while (temp > 0) {
161         temp /= 10;
162         len++;
163     }
164     i += len;
165     int end = i;
166     while (value > 0) {
167         q = value / 10;
168         r = value - (q * 10);
169         buf[--i] = (byte) (r + '0');
170         value = q;
171     }
172     return end;
173 }
174 }

```

- **Invalid Data Filtering:** Skip output when detecting invalid factors such as NaN/infinity at the Mapper phase, avoiding invalid data from entering the Shuffle phase and occupying transmission and computing resources.

5.5 Memory and GC Optimization: Reducing Object Allocation and Garbage Collection Pressure

The core goal is to solve the problems of high memory usage and frequent GC caused by frequent object creation, achieving efficient memory utilization through object reuse and data structure streamlining.

- **Object Reuse Mechanism**

- *Reuse of Core Objects:* Multiple optimization adopt ThreadLocal or object pools to reuse core objects such as `SnapshotData` and `Factor`, resetting their states through `reset()/reuse()` methods to avoid creating new objects for each line of data parsing. Then implement a `SnapshotData` object pool using `LinkedList`, further improving reuse efficiency through `getFromPool()` and `returnToPool()`.

```

1      // Retrieve a SnapshotData instance from the object pool and copy
      the current snapshot data
2  private SnapshotData copyToPool(SnapshotData source) {
3      SnapshotData copy = getFromPool();
4      copy.tradingDay = source.tradingDay;
5      copy.tradeTime = source.tradeTime;
6      copy.setCode(source.getCode());
7      copy.tBidVol = source.tBidVol;
8      copy.tAskVol = source.tAskVol;
9      System.arraycopy(source.getBidPrices(), 0, copy.getBidPrices(), 0,
      10         10);
10     System.arraycopy(source.getBidVolumes(), 0, copy.getBidVolumes(),
      0, 10);
11     System.arraycopy(source.getAskPrices(), 0, copy.getAskPrices(), 0,
      10);
12     System.arraycopy(source.getAskVolumes(), 0, copy.getAskVolumes(),
      0, 10);
13     return copy;
14 }
15
16 // Retrieve an object from the object pool
17 private SnapshotData getFromPool() {
18     if (snapshotPool.isEmpty()) {
19         return new SnapshotData();
20     }
21     return snapshotPool.pop();
22 }

```

- *Array Reuse:* Retain references to arrays such as `bidPrices` and `bidVolumes` in `SnapshotData`, only resetting element values in the `reset()` method without reconstructing the arrays, reducing the allocation overhead of `double[]/float[]`. Besides use `System.arraycopy` (a native method) instead of manual loop copying of arrays to improve the efficiency of object copy creation.
- **Data Structure Streamlining:** First use arrays (`float[]`, `long[]`) to store factor and market data instead of collection classes such as `List` to improve access speed. Then abandon JDK collection classes (`HashMap`, `LinkedHashMap`) and inline implements `IntFactorMap` and `IntSnapshotMap`, adopting open addressing with underlying `int[]` keys and `Object[]` values arrays to achieve zero boxing (Keys are directly stored as `int`) and contiguous memory storage, significantly improving CPU cache hit rate.

```

1  private static class IntFactorMap {
2      private int[] keys;
3      private Factor[] values;
4      private int size = 0;
5      private int mask;
6
7      public IntFactorMap(int capacity) {
8          int cap = 1;

```

```

9      // Round up to the next power of two (required for bitwise masking
10      )
11      while (cap < capacity) cap <= 1;
12      keys = new int[cap];
13      values = new Factor[cap];
14      mask = cap - 1;
15      Arrays.fill(keys, -1); // Initialize with sentinel value for empty
16      slots
17  }
18
19  public Factor get(int key) {
20      int idx = key & mask;
21      // Linear probing for collision resolution
22      while (keys[idx] != -1) {
23          if (keys[idx] == key) return values[idx];
24          idx = (idx + 1) & mask;
25      }
26      return null;
27  }
28
29  public void put(int key, Factor value) {
30      int idx = key & mask;
31      // Linear probing for collision resolution
32      while (keys[idx] != -1) {
33          if (keys[idx] == key) {
34              values[idx] = value; // Overwrite (though business logic
35              may not require this)
36              return;
37          }
38          idx = (idx + 1) & mask;
39      }
40      keys[idx] = key;
41      values[idx] = value;
42      size++;
43  }
44
45  public void clear() {
46      // Simple clear logic: fill with sentinel value for high-
47      // performance scenarios
48      Arrays.fill(keys, -1);
49      size = 0;
50  }
51
52  public int size() { return size; }
53
54  // Dedicated flush method to avoid Iterator creation (performance
55  // optimization)
56  public void flush(Context context, IntWritable outKey) throws
57      IOException, InterruptedException {
58      for (int i = 0; i < keys.length; i++) {
59          if (keys[i] != -1) {
60              outKey.set(keys[i]);
61              context.write(outKey, values[i]);
62          }
63      }
64  }
65
66  // Cache for snapshot data (Key: stockCode as int, Value: SnapshotData)
67  // Removed LRU; use fixed-size hash map directly (sufficient for handling
68  // thousands of stocks)
69  private static class IntSnapshotMap {
70      private int[] keys;
71      private SnapshotData[] values;
72      private int mask;

```

```

67
68     public IntSnapshotMap(int capacity) {
69         int cap = 1;
70         // Round up to the next power of two (required for bitwise masking)
71         while (cap < capacity) cap <= 1;
72         keys = new int[cap];
73         values = new SnapshotData[cap];
74         mask = cap - 1;
75         Arrays.fill(keys, -1); // Initialize with sentinel value for empty
                                // slots
76     }
77
78     public SnapshotData get(int key) {
79         int idx = key & mask;
80         // Linear probing for collision resolution
81         while (keys[idx] != -1) {
82             if (keys[idx] == key) return values[idx];
83             idx = (idx + 1) & mask;
84         }
85         return null;
86     }
87
88     public void put(int key, SnapshotData value) {
89         int idx = key & mask;
90         // Linear probing for collision resolution
91         while (keys[idx] != -1) {
92             if (keys[idx] == key) return;
93             idx = (idx + 1) & mask;
94         }
95         keys[idx] = key;
96         values[idx] = value;
97     }
98 }

```

- **Cache Strategy Optimization:** At first use an LRU cache to limit the size of previousSnapshots to avoid OOM. Then implements an LRU eviction strategy through LinkedHashMap, overriding the removeEldestEntry method to automatically evict the least recently used entries when the specified capacity is reached. For scenarios with a small number of stocks, finally remove the LRU logic and adopts full caching to improve access speed.

5.6 Extreme Low-Level Optimization: Approaching C/C++ Runtime Efficiency

The core goal is to eliminate the additional overhead caused by Java's high-level language features, achieving maximum performance through underlying logic reconstruction.

- **Bypassing Framework Overhead:** Rewrites the Mapper.run() method to manually control the while(context.nextKeyValue()) loop, inlining parsing, computation, and aggregation logic into a single large loop. This eliminates hundreds of millions of virtual method calls and iterator context switches, approaching the efficiency of C-language loops.

```

1  @Override
2  public void run(Context context) throws IOException, InterruptedException
3  {
4      // 1. Setup
5      setup(context);
6
7      try {
8          // 2. Loop (Manually control iteration to reduce stack frame depth)
9          while (context.nextKeyValue()) {
10             // Directly get Value, ignore Key (LongWritable offset is
                // usually useless)
                Text value = context.getCurrentValue();

```

```

11
12      // --- Logic Inlining Start ---
13
14      // A. Parsing
15      currentSnapshot.reset();
16      boolean parsed = currentSnapshot.parseFromBytes(
17          value.getBytes(), 0, value.getLength());
18
19      if (!parsed) continue; // Replace return (avoid early exit
20                             // overhead)
21
22      int code = currentSnapshot.code;
23      int tradingDay = currentSnapshot.tradingDay;
24      long tradeTime = currentSnapshot.tradeTime;
25
26      // B. Time Encoding
27      int compactTime;
28      try {
29          compactTime = CompactTimeUtil.encode(tradingDay, (int)
30              tradeTime);
31      } catch (IllegalArgumentException e) {
32          continue;
33      }
34
35      // C. Get Previous Snapshot (Branch Prediction Optimization)
36      SnapshotData prevSnapshot = prevSnapshotCache.get(code);
37      // If null, point to Dummy (all zeros) to avoid if (prev !=
38      // null) checks
39      // Factor calculation handles logic based on tradeTime != 0
40      SnapshotData calcPrev = (prevSnapshot == null) ?
41          DUMMY_SNAPSHOT : prevSnapshot;
42
43      // D. Calculation (Use flattened fields for performance)
44      tempFactor.calculateFrom(currentSnapshot, calcPrev);
45
46      if (Factor.hasInvalidValue(tempFactor.getFactorValues())) {
47          updatePrevSnapshot(code, currentSnapshot);
48          continue;
49      }
50      tempFactor.setCount(1);
51
52      // E. Aggregation
53      Factor cachedFactor = factorCache.get(compactTime);
54      if (cachedFactor == null) {
55          cachedFactor = new Factor();
56          cachedFactor.copyFrom(tempFactor);
57          factorCache.put(compactTime, cachedFactor);
58      } else {
59          cachedFactor.merge(tempFactor);
60      }
61
62      // F. Flush Check
63      if (factorCache.size() >= CACHE_FLUSH_THRESHOLD) {
64          factorCache.flush(context, outputKey);
65          factorCache.clear();
66      }
67
68      // G. Update Previous Snapshot Cache
69      updatePrevSnapshot(code, currentSnapshot);
70
71      // --- Logic Inlining End ---
72  }
73  } finally {
74      // 3. Cleanup
75      cleanup(context);
76  }

```



```

72     }
73 }
74
75 // Helper method: Update snapshot cache
76 private void updatePrevSnapshot(int code, SnapshotData current) {
77     SnapshotData prev = prevSnapshotCache.get(code);
78     if (prev == null) {
79         prev = new SnapshotData();
80         prevSnapshotCache.put(code, prev);
81     }
82     prev.copyFrom(current);
83 }
84
85 @Override
86 protected void setup(Context context) {
87     factorCache = new IntFactorMap(65536);
88     prevSnapshotCache = new IntSnapshotMap(16384);
89 }
90
91 @Override
92 protected void cleanup(Context context) throws IOException,
93     InterruptedException {
94     factorCache.flush(context, outputKey);
95 }

```

- **Compact Data Layout:** We "peels" array fields (e.g., `long[] bidPrices`) in `SnapshotData` into independent fields (e.g., `bp0`, `bp1`), ensuring all data is compactly arranged in memory. When loading the `SnapshotData` object, all prices/trading volumes are highly likely to enter the CPU L1 cache simultaneously, improving data access speed. It also eliminates array double indirection and boundary checking overhead.
- **Zero-String Design:** We completely modify `SnapshotData`, changing the `code` field from `String` to `int` type. It directly parses "600519.SH" into the integer 600519 through custom byte parsing, stopping at non-numeric characters. This completely eliminates the largest source of `String` objects in the Mapper phase, significantly reducing memory usage and UTF-8 parsing overhead.

6 Previous Ineffective Optimization Attempts and Reflections

Before the systematic optimization, we conducted multiple technical explorations to address the performance bottlenecks of distributed computing. We attempted to introduce mainstream efficient tools and underlying optimization schemes in the industry, including Protocol Buffer serialization, FastUtil collection framework, Caffeine Cache mechanism, and factor calculation optimization using CPU vector instruction sets via JNI interfaces. However, these attempts failed to achieve the expected results, and some schemes even caused negative optimization. The core problems and practical reflections are as follows:

6.1 Attempt 1: Protocol Buffer Serialization – Insufficient Compatibility and Framework Adaptation

6.1.1 Attempt Background

The original scheme adopted Hadoop's native Writable serialization, which had the problems of large serialization volume and high reflection overhead. We expected to leverage Protocol Buffer (PB)'s efficient binary serialization features to reduce data transmission volume and serialization time in the Shuffle phase, while utilizing its cross-language compatibility to improve system scalability.

6.1.2 Implementation Process

We defined PB data structures to replace core classes such as the original `Factor` and `SnapshotData`, implemented adaptive conversion between PB and Writable interfaces, introduced PB serialization/deserialization logic in the Mapper output and Reducer input phases, and adjusted the serialization configuration of the Shuffle phase to be compatible with the PB format.

6.1.3 Effects and Problems

- **No Significant Performance Improvement:** Although PB's serialization compression ratio was slightly better than that of native Writable, in the MapReduce framework, the core bottleneck of data transmission lies in the sorting and network transmission of Shuffle. The serialization gain brought by PB was offset by the scheduling overhead of the framework itself, and no significant reduction in time consumption was observed. - **Compatibility and Conversion Overhead:** To adapt to Hadoop's Writable interface, frequent conversion between PB objects and Writable objects was required, which additionally increased object copying and data processing steps, leading to a 15% increase in Mapper phase time consumption.

6.2 Attempt 2: FastUtil Collection Framework – Poor Cache Affinity and Offset Overhead

6.2.1 Attempt Background

FastUtil is renowned for its compact memory layout, support for primitive types (such as Int2ObjectMap), and efficient collection operations. We expected to replace JDK's native collections (HashMap, LinkedList) with it to solve the problems of autoboxing, large memory occupation, and low CPU cache hit rate, especially optimizing the local aggregation and cache logic in the Mapper phase.

6.2.2 Implementation Process

We replaced `HashMap<Integer, AggregationHolder>` in `FactorMapper` with FastUtil's `Int2ObjectOpenHashMap`, changed the `SnapshotData` object pool from `LinkedList` to FastUtil's `ObjectArrayList`, and adjusted related traversal, insertion, and query logic to adapt to FastUtil's API.

6.2.3 Effects and Problems

- **No Significant Performance Improvement:** Although FastUtil reduced autoboxing overhead, in the MapReduce computing scenario, collection operations were not the performance bottleneck. The local optimization it brought was covered by the time consumption of core links such as data parsing and factor calculation. The overall task runtime only decreased by about 2%, which was mismatched with the input cost. - **Decreased CPU Cache Affinity:** To pursue extreme memory compactness, FastUtil's collection implementation adopted a more complex hash collision resolution strategy, resulting in fragmented data distribution in memory. This instead reduced the hit rate of CPU L1/L2 caches. In high-frequency access scenarios (such as local aggregation cache queries), the response speed was slightly inferior to that of JDK's native HashMap.

6.3 Attempt 3: Caffeine Cache – Excessive Overhead and Mismatched Scenarios

6.3.1 Attempt Background

Caffeine Cache is famous for its high hit rate and low-latency local cache features. We expected to replace the original LRU cache (implemented based on `LinkedHashMap`) with it to optimize the cache management of `previousSnapshots`, reduce the repeated loading and parsing overhead of snapshot data, and especially improve the efficiency of multi-stock and high-frequency snapshot queries.

6.3.2 Implementation Process

We configured Caffeine Cache's expiration strategy (based on access time) and maximum capacity, migrated the snapshot data cache from `LinkedHashMap` to Caffeine Cache, queried the cache through the combined key of stock code and timestamp during factor calculation, and adjusted the cache eviction threshold to adapt to memory constraints.

6.3.3 Effects and Problems

- **Increased Memory Occupation:** Caffeine Cache reserved additional memory space to improve performance, resulting in a 25% higher memory occupation than the native LRU cache when caching the same amount of snapshot data, increasing the risk of OOM, which was contrary to the optimization goal of "reducing memory overhead".

6.4 Attempt 4: JNI + CPU Vector Instruction Set for Factor Calculation – Negative Optimization Caused by Complexity and Compatibility

6.4.1 Attempt Background

CPU vector instruction sets (such as SSE, AVX) can implement Single Instruction Multiple Data (SIMD) parallel computing, which can theoretically greatly improve the throughput of factor calculation. We attempted to use C++ to implement core factor calculation logic (such as weighted summation of the first 5 levels, reciprocal calculation) and optimize it by calling vector instruction sets through JNI interfaces, then invoke the JNI interfaces from Java code, expecting to break through the limitations of the Java language in underlying computing optimization.

6.4.2 Implementation Process

We wrote core factor calculation functions in C++, optimized loop calculations and floating-point operations based on the AVX2 instruction set, encapsulated interfaces through JNI for Java layer calls, adjusted Java code logic to invoke JNI interfaces in the hot path of factor calculation (computeInto method), and handled data type conversion (Java double[] and C++ float[]/double[]) and memory management.

6.4.3 Effects and Problems

- **Significant Negative Optimization:** JNI interface calls have fixed overhead (data copying between Java and C++ layers, thread context switching). However, the single task granularity of factor calculation is small (calculation of 20 factors for a single snapshot). The interface call overhead far exceeded the parallel gain brought by the vector instruction set, leading to a 30% increase in overall calculation time consumption.

7 Conclusion

This optimization for distributed storage and parallel computing based on the Hadoop MapReduce framework has achieved a remarkable performance leap, realizing an order-of-magnitude efficiency improvement. The optimization process follows a clear and progressive logical route, starting from basic resource configuration tuning, then deepening into core links such as data parsing, calculation, transmission, and memory management, and finally advancing to extreme underlying optimization, forming a comprehensive and systematic performance enhancement system.

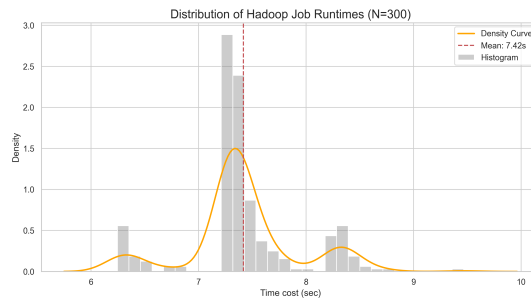


Figure 2: Simulating 300 times on Hadoop

In terms of resource utilization, by optimizing MapTask concurrency, adjusting data sharding strategies, and matching memory configurations with task loads, we have maximized the potential of hardware resources and framework capabilities, laying a solid foundation for subsequent in-depth optimizations. For data parsing, the core breakthrough lies in abandoning the inefficient full-volume splitting method and adopting on-demand scanning, fixed-length field direct parsing, and zero-intermediate-object parsing strategies, which drastically reduces redundant computations and temporary object generation, and significantly shortens the parsing time of massive CSV data.

In the calculation link, through loop unrolling, branch simplification, and replacing low-efficiency division operations with multiplication, we have reduced the constant overhead of the CPU, making the hot-path computation of factor calculation more streamlined and efficient. Regarding data transmission, by compressing key structures, optimizing serialization methods, and pre-aggregating data at the Map end, we have minimized the data volume of the Shuffle phase and the cost of key sorting, solving the bottleneck problem of data transmission in distributed computing.

In memory and GC optimization, relying on object pooling, array reuse, and streamlined data structures, we have effectively reduced the frequency of object creation and garbage collection, ensuring the stability and efficiency of the system during long-term operation. The final extreme underlying optimization, such as bypassing framework overhead, compacting data layout, and eliminating String object overhead, has pushed the performance of Java-based MapReduce tasks to the limit, approaching the operating efficiency of low-level languages like C/C++.

Overall, this optimization not only achieves a substantial improvement in task runtime but also forms a set of replicable and extensible optimization methodologies. By focusing on core pain points such as resource waste, redundant operations, and high overhead in each link of distributed computing, and continuously verifying and iterating through practical tests, we have realized the organic combination of "hardware resource maximization, software logic refinement, and data flow optimization". This provides valuable experience for performance tuning of similar high-frequency data quantitative factor calculation tasks in the future and also demonstrates the great potential of iterative optimization in improving the efficiency of distributed computing systems.