# Bit Packing Compression for Optimized Network Transmission

## Software Engineering Project 2025

Youness RADI
22311111

November 2025

**Abstract**

This report presents the design, implementation, and analysis of bit packing compression algorithms for integer array transmission. I implemented three distinct strategies: overlapping, non-overlapping, and overflow compression. each optimized for different use cases. the benchmarks demonstrates compression ratios of 2-8x with millisecond-scale computational overhead, making compression beneficial for virtually all network transmission scenarios. The implementation provides O(1) random access without full decompression, a critical feature for practical applications.

## 1 Introduction

Integer array transmission is a fundamental problem in distributed systems, databases, and network protocols. Standard approaches waste significant bandwidth:transmitting 10,000 integers requires 320,000 bits, even when values only need 12 bits each.

**The Core Problem**: How to compress integer arrays while maintaining:

- **Lossless compression** - exact reconstruction

- **O(1) random access** - retrieve any element without full decompression

- **Computational efficiency** - compression overhead must be less than transmission time saved

**The proposed solution**: Bit packing: Use only the minimum bits needed per value, not the full 32 bit representation. We had to implement three algorithms, each with distinct trade-offs.

## 2 Problem Analysis

### 2.1 The bandwidth Waste

Consider an array of 10,000 values where $max = 4095$ (12-bit values):

- **Standard transmission**: $10,000 \times 32 = 320,000$ bits

- **Optimal transmission**: $10,000 \times 12 = 120,000$ bits

- **Waste**: $200,000$ bits ($62.5\%$)

## 2.2 The Trade-Off Space

Bit packing presents fundamental trade-offs:

1. **Space vs Time**: Maximum compression vs. fastest access

2. **Complexity vs Efficiency**: Simple algorithms vs. optimal results

3. **General vs Specialized**: Works for all data vs. optimal for specific patterns

## 2.3 Requirements

From the project specification:

- Implement two compression variants (overlapping and non-overlapping)

- Provide `compress()`, `decompress()`, and `get(index)` functions

- Measure execution time with rigorous methodology

- Calculate transmission break-even point

- **Required**: Implement overflow compression for sparse data with outliers

# 3 Solution Design

## 3.1 Core Principle

All three algorithms share a fundamental principle:

$$\text{bitsNeeded} = \lceil \log_2(\max(\text{data})) \rceil + 1$$

For an array with $n$ elements:

$$\text{compressedSize} = \left\lceil \frac{n \times \text{bitsNeeded}}{32} \right\rceil$$

## 3.2 Architecture

We employ three design patterns:

- **Factory Pattern**: `BitPackingFactory` creates compressor instances

- **Strategy Pattern**: `BitPacking` interface with multiple implementations

- **Template Method**: `AbstractBitPacking` defines workflow, subclasses implement specifics

# 4 Core Implementation Details

## 4.1 Bit Utilities: The Foundation

All compression algorithms rely on low-level bit manipulation utilities in `BitUtils.java`.

### 4.1.1  Calculating Required Bits

```java
public static int bitsNeeded(int value) {
    if (value == 0) return 1;
    return 32 - Integer.numberOfLeadingZeros(value);
}
```

Listing 1: Determining minimum bits needed

**How it works**: `Integer.numberOfLeadingZeros()` counts zero bits from the most significant bit. For value 15 (binary: 1111), leading zeros = 28, so bits needed = $32 - 28 = 4$.

### 4.1.2  Creating Bit Masks

```java
public static int createMask(int bits) {
    if (bits == 32) return -1;
    return (1 << bits) - 1;
}
```

Listing 2: Dynamic mask generation

**Usage**: Masks isolate specific bit ranges. For 4 bits: $(1 \ll 4) - 1 = 16 - 1 = 15 = $ 0b00001111.

### 4.1.3  Overlapping Write: The Critical Operation

```java
public static void writeBitsOverlapping(int[] data, int bitPosition, int value,
     int bitsPerValue) {
    int intIndex = bitPosition / 32;
    int bitOffset = bitPosition % 32;
    int bitsAvailable = 32 - bitOffset;
    int mask = createMask(bitsPerValue);
    value &= mask;

    if (bitsAvailable >= bitsPerValue) {
        // all bits in one integer
        int clearMask = ~(mask << bitOffset);
        data[intIndex] = (data[intIndex] & clearMask) | (value << bitOffset);
    } else {
        // spans two integers
        int lowMask = createMask(bitsAvailable);
        int highMask = createMask(bitsPerValue - bitsAvailable);

        int clearLowMask = ~(lowMask << bitOffset);
        data[intIndex] = (data[intIndex] & clearLowMask) | ((value & lowMask)
            << bitOffset);

        int highBits = value >>> bitsAvailable;
        data[intIndex + 1] = (data[intIndex + 1] & ~highMask) | (highBits &
            highMask);
    }
}
```

Listing 3: Writing bits that span integer boundaries

**Detailed breakdown**:

1. **Position calculation**: convert absolute bit position to (integer index, bit offset) pair

2. **Space check**: if remaining space $\geq$ bits needed $\rightarrow$ single write, else split

3. **Mask operations**:

3

- clearMask: Zeros target bits while preserving others
- value << bitOffset: Shifts value to correct position
- data & clearMask | shifted_value: merge without affecting unrelated bits

4. **Split write**: Low bits fill current integer, high bits start next integer

**Example**: Writing 12-bit value 0xABC at bit position 28:

- intIndex = 0, bitOffset = 28, bitsAvailable = 4
- Low 4 bits (0xC) → bits 28-31 of data[0]
- High 8 bits (0xAB) → bits 0-7 of data[1]

### 4.1.4 Overlapping Read: The Extraction

```java
public static int extractBitsOverlapping(int[] data, int bitPosition, int
    bitsPerValue) {
    int intIndex = bitPosition / 32;
    int bitOffset = bitPosition % 32;
    int bitsAvailable = 32 - bitOffset;

    if (bitsAvailable >= bitsPerValue) {
        // all bits in one integer
        return (data[intIndex] >>> bitOffset) & createMask(bitsPerValue);
    } else {

        // spans two integers
        int lowBits = (data[intIndex] >>> bitOffset) & createMask(bitsAvailable
            );
        int highBits = data[intIndex + 1] & createMask(bitsPerValue -
            bitsAvailable);
        return lowBits | (highBits << bitsAvailable);
    }
}
```

Listing 4: Extracting bits across boundaries

**Key operations**:

- >>> bitOffset: Unsigned right shift to align bits to position 0
- & createMask(): Isolate only the required bits
- lowBits | (highBits << bitsAvailable): reconstruct original value

## 4.2 Overlapping compression

### 4.2.1 Compression Algorithm

```java
protected int[] performCompression(int[] data, int bitsPerValue) {
    int totalBits = data.length * bitsPerValue;
    int compressedSize = (totalBits + 31) / 32;  // ceiling division
    int[] compressed = new int[compressedSize];

    for (int i = 0; i < data.length; i++) {
        int bitPosition = i * bitsPerValue;
        BitUtils.writeBitsOverlapping(compressed, bitPosition,
                                      data[i], bitsPerValue);
    }

    return compressed;
```

Listing 5: Overlapping compression implementation

**Complexity**: $O(n)$ time, $O(\lceil \frac{n \cdot k}{32} \rceil)$ space where $k =$ bits per value.

**Why ceiling division**: $(totalBits + 31)/32$ ensures partial integers are allocated. Example: 125 bits requires $\lceil 125/32 \rceil = 4$ integers.

### 4.2.2 Decompression Algorithm

```java
protected int[] performDecompression(int[] compressed, int originalSize) {
    int[] decompressed = new int[originalSize];

    for (int i = 0; i < originalSize; i++) {
        int bitPosition = i * this.bitsPerValue;
        decompressed[i] = BitUtils.extractBitsOverlapping(
            compressed, bitPosition, this.bitsPerValue);
    }
    return decompressed;
}
```

Listing 6: Overlapping decompression

**Symmetry**: Compression and decompression are exact inverses. Bit position calculation is identical, enabling O(1) random access.

### 4.2.3 Random Access

```java
protected int performGet(int index) {
    int bitPosition = index * bitsPerValue;
    return BitUtils.extractBitsOverlapping(compressedData,
                                           bitPosition, bitsPerValue);
}
```

Listing 7: O(1) random access

**Efficiency**: No iteration required. Direct arithmetic: $bitPosition = index \times bitsPerValue$.

## 4.3 Non-Overlapping Compression

### 4.3.1 Slot Calculation

```java
protected int[] performCompression(int[] data, int bitsPerValue) {
    this.valuesPerInt = 32 / bitsPerValue;  // floor division
    int compressedSize = (data.length + valuesPerInt - 1) / valuesPerInt;
    int[] compressed = new int[compressedSize];

    for (int i = 0; i < data.length; i++) {
        int intIndex = i / valuesPerInt;        // which integer
        int slotIndex = i % valuesPerInt;       // which slot [0..n]
        int bitOffset = slotIndex * bitsPerValue;

        BitUtils.writeBitsNonOverlapping(compressed, intIndex,
                                         bitOffset, data[i], bitsPerValue);
    }
    return compressed;
}
```

Listing 8: Non-overlapping compression with slots

**Key difference**: Each value stays within a single integer. No boundary crossing.

**Wasted space**: For 12-bit values, $valuesPerInt = \lfloor 32/12 \rfloor = 2$. Each integer holds 2 values (24 bits), wasting 8 bits.

5

### 4.3.2 Non-Overlapping Write

```java
public static void writeBitsNonOverlapping(int[] data, int intIndex,
                                            int bitOffset, int value,
                                            int bitsPerValue) {
    int mask = createMask(bitsPerValue);
    value &= mask;
    int clearMask = ~(mask << bitOffset);        // zero target bits
    data[intIndex] = (data[intIndex] & clearMask) | (value << bitOffset);
}
```

Listing 9: Single-integer write operation

**Simpler logic**: No split case. Always single integer operation. Faster but less space-efficient.

## 4.4 Overflow Compression

### 4.4.1 Two-Tier Storage Design

**Encoding scheme**: [1-bit flag][k bits payload]

- flag=0: payload = actual value (direct storage)

- flag=1: payload = index into overflow array (indirect storage)

### 4.4.2 Threshold Optimization

```java
private OverflowStats analyzeOverflow(int[] data) {
    int max = findMax(data);
    int maxBits = BitUtils.bitsNeeded(max);
    int bestTotalBits = data.length * maxBits;  // baseline: no overflow

    // Try thresholds: maxBits-8 to maxBits-1
    for (int thresholdBits = max(1, maxBits - 8);
         thresholdBits < maxBits; thresholdBits++) {

        int threshold = (1 << thresholdBits);
        int overflowCount = countAbove(data, threshold);

        // Constraint: index bits must fit in main storage
        int indexBits = BitUtils.bitsNeeded(overflowCount);
        if (indexBits > thresholdBits) continue;  // skip invalid

        // Calculate total storage
        int mainBits = data.length * (thresholdBits + 1);  // +1 flag
        int overflowBits = overflowCount * 32;
        int totalBits = mainBits + overflowBits;

        if (totalBits < bestTotalBits) {
            bestTotalBits = totalBits;
            // update best threshold...
        }
    }
    return bestStats;
}
```

Listing 10: Finding optimal overflow threshold

**Algorithm analysis**:

1. **Search space**: Try thresholds 8 bits below maximum down to 1 bit below

2. **Constraint check**: Overflow index must fit in payload bits

3. **Cost function**: $totalBits = n \times (k' + 1) + overflowCount \times 32$

4. **Optimization**: Select threshold minimizing total storage

**Example**: For data $[1, 2, 3, 1024, 4, 5, 2048]$, maxBits=11:

- Try threshold=3 bits: 2 overflow values, needs 1-bit index $\rightarrow$ valid

- Cost: $7 \times 4 + 2 \times 32 = 28 + 64 = 92$ bits

- vs no overflow: $7 \times 11 = 77$ bits

- Result: Overflow not beneficial for this case

### 4.4.3 Overflow Compression

```java
protected int[] performCompression(int[] data, int bitsPerValue) {
    OverflowStats stats = analyzeOverflow(data);
    this.mainBits = stats.mainBits;
    this.totalMainBits = mainBits + 1;  // +1 for flag bit

    int[] overflowArea = new int[stats.overflowCount];
    int overflowIndex = 0;

    int mainStorageSize = (data.length * totalMainBits + 31) / 32;
    int[] compressed = new int[mainStorageSize + overflowArea.length];

    for (int i = 0; i < data.length; i++) {
        int bitPosition = i * totalMainBits;

        if (data[i] >= overflowThreshold) {
            // Overflow: encode as [1][index]
            int encoded = (1 << mainBits) | overflowIndex;
            BitUtils.writeBitsOverlapping(compressed, bitPosition,
                                          encoded, totalMainBits);
            overflowArea[overflowIndex++] = data[i];
        } else {
            // Direct: encode as [0][value]
            int encoded = (0 << mainBits) | data[i];
            BitUtils.writeBitsOverlapping(compressed, bitPosition,
                                          encoded, totalMainBits);
        }
    }

    // Append overflow area
    System.arraycopy(overflowArea, 0, compressed,
                     mainStorageSize, overflowArea.length);
    return compressed;
}
```

Listing 11: Two-tier compression logic

**Encoding details**:

- `(1 << mainBits)`: Sets flag bit to 1

- `| overflowIndex`: OR combines flag and payload

Example: mainBits=3, index=2:

- Binary result: `0b1010` (4 bits total)

- Structure: `[flag=1][payload=010]`

- Interpretation : "This is overflow value, stored at index 2"

### 4.4.4 Overflow Random Access

```java
protected int performGet(int index) {
    int bitPosition = index * totalMainBits;
    int encoded = BitUtils.extractBitsOverlapping(
        compressedData, bitPosition, totalMainBits);

    int flag = (encoded >>> mainBits) & 1;        // extract flag bit
    int payload = encoded & createMask(mainBits);  // extract payload

    if (flag == 1) {
        // Overflow: lookup in overflow area
        int mainStorageSize = (originalSize * totalMainBits + 31) / 32;
        return compressedData[mainStorageSize + payload];
    } else {
        return payload;  // direct value
    }
}
```

Listing 12: Random access with overflow lookup

# 5 Benchmarking Methodology

## 5.1 Protocol

Must follow rigorous timing methodology:

1. **Separate measurements**: Time compression, decompression, and random access independently

2. **Nanosecond precision**: Use `System.nanoTime()` for accurate measurements

3. **Multiple scenarios**: Test random data, sparse data, and transmission analysis

4. **Statistical validity**: Average 100 random accesses for stable measurements

## 5.2 Test Scenarios

- **Basic demo**: 8 elements (max=15, 4-bit)   validates correctness

- **Random data**: 10,000 elements (max=4095, 12-bit)   typical compression case

- **Sparse data**: 10,000 elements (90% in [0,15], 10% in [1000,10000])   tests overflow

- **Transmission**: 100,000 elements at 100 Mbps   calculates break-even

## 5.3  Transmission Break-Even Analysis

The critical question: *When is compression worth it?*

**Without compression**:

$$T_{\text{uncompressed}} = \frac{n \times 32}{\text{bandwidth}} \tag{1}$$

**With compression**:

$$T_{\text{compressed}} = T_{\text{compress}} + \frac{c \times 32}{\text{bandwidth}} + T_{\text{decompress}} \tag{2}$$

— $n$ = original size  $c$ = compressed size.

**Break-even point** (when times are equal):

$$T_{\text{compress}} + T_{\text{decompress}} = \frac{(n - c) \times 32}{\text{bandwidth}} \tag{3}$$

**Interpretation**: If transmission time saved, exceeds compression overhead, compression is beneficial.

# 6  Results and Analysis

Note: Benchmark execution times vary based on system conditions, but compression ratios and relative performance ordering remain consistent across runs

## 6.1  Basic Demo Validation

Testing with 8 elements [1, 5, 12, 7, 3, 9, 15, 2]:

- Original: 256 bits ($8 \times 32$)

- Compressed: 1 integer (32 bits)

- Bits per value: 4

- Compression ratio: **8.00x**

- All algorithms: Correctness verified

## 6.2  Random Data Performance

Testing with 10,000 random 12-bit values (max=4095):

| Algorithm | Compress | Decompress | Access | Ratio |
|---|---|---|---|---|
| Overlapping | 1.33 ms | 0.86 ms | 1224 ns | 2.67x |
| Non-Overlapping | 1.07 ms | 0.88 ms | 531 ns | 2.00x |
| Overflow | 2.62 ms | 0.49 ms | 324 ns | 2.67x |

Table 1: Performance on random 12-bit data (10K elements)

**Key observations**:

- **Compression speed**: Non-overlapping 20% faster (1.07 vs 1.33 ms)—simpler bit operations

- **Decompression**: Overflow 43% faster (0.49 vs 0.86 ms)—efficient extraction with branch prediction

- **Random access**: Overflow fastest (324 ns), Overlapping slowest (1224 ns)—boundary-spanning penalty visible

- **Compression ratio**: Overlapping and Overflow tie at 2.67x (optimal), Non-overlapping 25% worse due to wasted bits

**Surprising result**: Overflow has fastest random access despite conditional logic. Branch predictor handles the flag=0 common case efficiently.

## 6.3   Sparse Data Analysis

Testing with 10,000 values: 90% in [0, 15], 10% in [1000, 10000]:

| Algorithm | Compress | Decompress | Access | Ratio |
|---|---|---|---|---|
| Overlapping | 0.88 ms | 0.63 ms | 502 ns | 2.29x |
| Non-Overlapping | 1.02 ms | 0.80 ms | 486 ns | 2.00x |
| Overflow | 1.80 ms | 0.66 ms | 264 ns | 2.29x |

Table 2: Performance on sparse data (14-bit required)

**Critical insight**: Overflow does *not* achieve better compression here. Why?

- Data requires 14 bits (max  10,000)

- 10% overflow (1,000 values) requires 10-bit index

- Constraint: `indexBits > mainBits` check fails for low thresholds

- Result: Algorithm falls back to same compression as overlapping

- **Lesson**: Overflow excels when outliers are *rare* (¡1%), not 10%

## 6.4   Transmission Analysis

For 100,000 elements at 100 Mbps:

| Algorithm | Overhead | Time Saved | Break-Even |
|---|---|---|---|
| Overlapping | 7.75 ms | 20.00 ms | +12.25 ms |
| Non-Overlapping | 9.84 ms | 16.00 ms | +6.16 ms |
| Overflow | 16.39 ms | 20.00 ms | +3.61 ms |

Table 3: Transmission break-even analysis at 100 Mbps

**Practical implications**:

- All algorithms show **positive break-even**—compression saves time

- Typical network latency (50-200 ms) far exceeds overhead (8-16 ms)

- **Overlapping**: Best net savings (12.25 ms) due to maximum compression

- **Overflow**: Highest overhead (16.39 ms) from threshold analysis, but still worthwhile

# 7 Design Choices and Justification

## 7.1 Why These Strategies?

Each addresses different constraints:

- **Overlapping**: When bandwidth is the bottleneck (mobile networks, cloud storage)

- **Non-overlapping**: When simplicity and balanced performance matter (real-time systems)

- **Overflow**: When data has rare outliers (¡1% of values with large magnitude)

## 7.2 Why Factory Pattern?

```
BitPacking packer = BitPackingFactory.create(CompressionType.OVERLAPPING);
int[] compressed = packer.compress(data);
```

Listing 13: Clean API for algorithm selection

Benefits:

- Client code independent of concrete classes

- Easy to add new algorithms

- Centralized configuration

- Follows Open-Closed Principle

## 7.3 Why Template Method Pattern?

`AbstractBitPacking` handles common operations:

- Input validation (null checks, empty arrays)

- Bits-per-value calculation via `BitUtils`

- Compression ratio computation

- Bounds checking for random access

Subclasses implement only algorithm-specific operations:

- `performCompression()`

- `performDecompression()`

- `performGet()`

**Result**: 40% code reuse, eliminates duplication, ensures consistent behavior.

# 8    Conclusion

This project successfully implements three bit packing algorithms with distinct trade-offs:

- **Overlapping**: 2.67x compression, optimal space efficiency

- **Non-overlapping**: 2.00x compression, simpler implementation, 20% faster compression

- **Overflow**: 2.67x compression, fastest random access (324 ns vs 1224 ns)

Key achievements:

1. **O(1) random access**: Direct bit position calculation eliminates sequential access

2. **Practical performance**: Compression overhead (1-3 ms) far smaller than transmission savings (12-20 ms at 100 Mbps)

3. **Real-world applicability**: Techniques used in Apache Parquet, Protocol Buffers, time-series databases

4. **Design patterns**: Factory, strategy, and template method ensure extensibility and maintainability

5. **Comprehensive testing**: Basic validation, performance benchmarks, and transmission analysis

The implementation demonstrates that *appropriate complexity*—not maximum complexity—gives the best engineering solutions. Each algorithm serves specific use cases:

- Use **overlapping** for maximum compression when bandwidth is constrained

- Use **non-overlapping** for simpler code and balanced performance

- Use **overflow** for data with rare large outliers (¡1% of values)

**Closer**: Bit packing compression is essential for efficient integer array transmission. With positive break-even times (3-12 ms) far below typical network latencies (50-200 ms), compression is virtually always beneficial for networked applications.