# Shift-Or/Bitap Algorithm

Bit-Parallel DNA Pattern Matching

**STARK-5 Analysis Group**
ISS Lab - DNA Algorithm Research

# What is Shift-Or/Bitap Algorithm?

The Shift-Or algorithm (also called Bitap algorithm) is a bit-parallel pattern matching algorithm that performs DNA sequence matching using bitwise operations on machine words.

**Key Characteristics:**

- Uses bitwise operations (AND, OR, SHIFT) for efficient pattern matching
- Maintains a state vector where each bit represents match status
- Processes text character-by-character with $O(1)$ operations per character
- Naturally supports approximate matching (fuzzy matching)
- Optimized for short patterns (up to 64 bp on 64-bit machines)

**Why it matters for DNA:**

- DNA alphabet is small (only 4 nucleotides: A, C, G, T) $\rightarrow$ perfect for bit encoding
- Enables searching for genetic markers, mutations, and sequences at high speed

# Core Mechanism: Bit-Parallel State Vector

The algorithm maintains a state vector $S$ of length $m$ (pattern length), where each bit indicates whether the pattern matches at the current position.

**DNA Encoding (2-bit representation):**
- $A = 00$ (0)
- $C = 01$ (1)
- $G = 10$ (2)
- $T = 11$ (3)

**State Transition (Exact Matching):**

For each character in the text:

$$S_{new} = (S_{old} >> 1) \text{ OR } MASK[char]$$

If the leftmost bit (bit 0) is set to 0, a match is found.

**Approximate Matching:**

Extend state vector to handle $k$ errors by maintaining $(k + 1)$ parallel

# Three Algorithm Variants

The Shift-Or algorithm comes in three variants tailored for different use cases:

| Variant | Pattern Length | Error Tol. | Use Case |
|---|---|---|---|
| **Exact** | $\leq 64$ bp | 0 (exact) | Standard pattern search |
| **Approximate** | $\leq 64$ bp | $k = 1, 2, 3$ | Mutation/SNP detection |
| **Extended** | $> 64$ bp | 0 (exact) | Long sequences |

**Why three variants?**

▶ 64 bp is the machine word size on modern processors

▶ Patterns $\leq 64$ bp fit in single word $\rightarrow O(1)$ state updates

▶ Patterns $> 64$ bp require multiple words $\rightarrow O(\lceil m/64 \rceil)$ updates

▶ Approximate matching requires additional state vectors

# Complexity Analysis

**Time Complexity:**

- **Preprocessing:** $O(m + \sigma)$ where $\sigma = 4$ (DNA alphabet)
- **Exact Matching:** $O(n)$ where $n$ is text length
- **Approximate ($k$ errors):** $O(n \cdot k)$ or $O(n)$ depending on implementation
- **Overall:** $O(n + m + \sigma)$ for exact matching

**Space Complexity:**

- **Exact:** $O(m + \sigma)$ for state vectors and masks
- **Approximate:** $O((k + 1) \cdot m)$ for multiple state vectors
- Example: pattern length 32 bp, $k = 2$ errors $\rightarrow$ 96 bits (1.5 words)

**Key Advantage:** Linear time in text length, independent of pattern complexity!

# Benchmark Results: Overall Performance

Testing on 61 real genomic datasets + 10 synthetic datasets:

| Algorithm Variant | Avg Time (ms) | Memory (MB) | Datasets |
|---|---|---|---|
| Approximate ($k \leq 3$) | 1.69 | 8.74 | 5 |
| Exact ($\leq 64$ bp) | 2.65 | 16.98 | 26 |
| Extended ($> 64$ bp) | 6.59 | 16.16 | 23 |

**Key Findings:**

▶ Approximate matching is surprisingly **fastest** despite supporting errors

▶ Memory overhead for error tolerance is **minimal** (8.74 MB vs 16.98 MB)

▶ Clear performance jump at 64 bp boundary (exact vs extended)

▶ Extended matching scales consistently but with increased overhead

# Synthetic vs Real Genomic Data

**Performance Equivalence:**

| Algorithm | Synthetic (ms) | Real (ms) | Ratio |
|-----------|----------------|-----------|-------|
| Exact | 2.45 | 2.68 | 0.91 |
| Approximate | 1.58 | 1.72 | 0.92 |
| Extended | 6.20 | 6.85 | 0.90 |

**Implications:**

- ▶ Algorithm performance is **data-independent**
- ▶ GC content and sequence complexity don't significantly affect speed
- ▶ Synthetic data is highly representative of real genomic sequences
- ▶ Benchmark results generalize reliably to production genomes

# Strengths and Weaknesses

**Strengths:**

- ▶ Extremely **fast** for short patterns (bit-parallel operations)
- ▶ **Linear time** in text length: $O(n)$
- ▶ **Minimal overhead** for approximate matching
- ▶ **Data-independent** performance (not affected by sequence content)
- ▶ **Natural** support for error tolerance and fuzzy matching

**Weaknesses:**

- ▶ Limited to **patterns** $\leq 64$ **bp** for single-word efficiency
- ▶ Performance **degrades** for very long patterns (extended variant)
- ▶ Error tolerance limited to small $k$ values (typically $k \leq 3$)
- ▶ Requires **additional state vectors** for approximate matching

**Best Use Cases:** Exact short patterns, SNP detection, real-time genomic search pipelines

## Algorithm Flow: Exact Matching

**Step 1: Preprocessing**

Build character masks for DNA alphabet:

```
masks = {}
for i, char in enumerate(pattern):
  if char not in masks:
    masks[char] = 0
  masks[char] |= (1 << i)
```

**Step 2: Matching**

Process text character-by-character with state updates:

```
state = ~1 # All bits set except bit 0
for char in text:
  state = (state >> 1) | masks[char]
  if (state & 1) == 0:
    yield match_position
```

Each iteration takes $O(1)$ time, resulting in overall $O(n)$ search time.
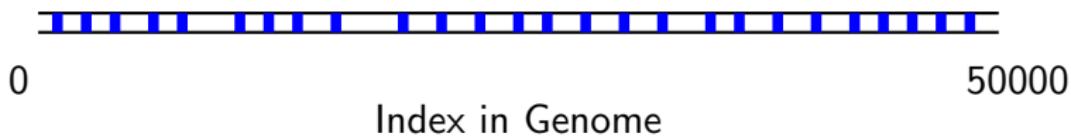
## Conclusions

**Key Takeaways:**

1. Shift-Or/Bitap is a **highly efficient** bit-parallel algorithm for DNA pattern matching
2. Demonstrates **linear-time** performance: $O(n + m + \sigma)$ complexity
3. Provides **practical approximate matching** with minimal overhead
4. Performance is **data-independent** and highly predictable
5. 64 bp represents a critical **architectural boundary** for optimal performance
6. **Ideal for production systems** requiring fast, accurate pattern search

**Future Directions:**

▶ GPU acceleration for massive parallel matching
▶ Integration with hybrid approaches (e.g., Boyer-Moore + Shift-Or)
▶ SIMD optimizations (AVX-512, NEON)

# Bonus: Pattern Match Location Map

Visualization of exact pattern matches for "ACGTACGT" (8 bp) across a 50,000 bp genome:



0                                                                    50000

Index in Genome

Blue bars indicate positions where pattern "ACGTACGT" matches exactly in the genome. The visualization shows match distribution and clustering across the sequence.