

# Análisis y Optimización de Algoritmos para Detección de Código Malicioso en Transmisiones de Datos

DIEGO VERGARA HERNÁNDEZ, Tecnológico de Monterrey

JOSÉ LEOBARDO NAVARRO MÁRQUEZ, Tecnológico de Monterrey

SANTIAGO ARISTA, Tecnológico de Monterrey

Este artículo presenta un análisis exhaustivo de algoritmos para la detección de código malicioso en transmisiones de datos hexadecimales. Se implementaron cuatro técnicas fundamentales: búsqueda de subsecuencias mediante KMP, detección de palíndromos con el algoritmo de Manacher, identificación del substring común más largo mediante programación dinámica optimizada, y detección de anomalías usando codificación de Huffman. Se analizan las complejidades computacionales de cada enfoque y se demuestra que la codificación de Huffman proporciona una mejora significativa en la detección de patrones anómalos basándose en análisis de frecuencias. Los resultados muestran que la optimización de espacio en el algoritmo de substring común (de  $O(n \times m)$  a  $O(m)$ ) permite el procesamiento eficiente de archivos grandes, mientras que el análisis estadístico con Huffman detecta código malicioso con alta precisión en distribuciones sesgadas.

CCS Concepts: • Theory of computation → Pattern matching; Design and analysis of algorithms; Data compression; • Security and privacy → Malware and its mitigation; Intrusion detection systems.

Additional Key Words and Phrases: Algoritmos de búsqueda, KMP, Manacher, Programación dinámica, Codificación de Huffman, Detección de anomalías, Complejidad computacional, Análisis de strings

## ACM Reference Format:

Diego Vergara Hernández, José Leobardo Navarro Márquez, and Santiago Arista. 2026. Análisis y Optimización de Algoritmos para Detección de Código Malicioso en Transmisiones de Datos. 1, 1 (February 2026), 9 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

## 1. Introducción

La detección de código malicioso en transmisiones de datos es un problema fundamental en ciberseguridad. Este trabajo aborda el desafío de identificar patrones sospechosos en flujos de datos hexadecimales mediante cuatro enfoques algorítmicos distintos, cada uno optimizado para un aspecto específico del problema.

### 1.1. Contexto del Problema

Las transmisiones de datos contienen secuencias de caracteres hexadecimales (0-9, A-F) que pueden incluir código malicioso. El objetivo es detectar:

1. Presencia exacta de patrones conocidos (subsecuencias)
2. Código “espejeado” mediante palíndromos

---

Authors' Contact Information: Diego Vergara Hernández, Tecnológico de Monterrey, diego.vergara@tec.mx; José Leobardo Navarro Márquez, Tecnológico de Monterrey, jose.navarro@tec.mx; Santiago Arista, Tecnológico de Monterrey, santiago.arista@tec.mx.

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2026 Copyright held by the owner/author(s). Publication rights licensed to ACM.

Manuscript submitted to ACM

3. Similitudes entre transmisiones
4. Patrones anómalos basados en frecuencias de caracteres

## 1.2. Contribuciones

Este trabajo presenta:

- Implementación eficiente de cuatro algoritmos clásicos adaptados al problema
- Optimización de espacio para el algoritmo de substring común
- Método estadístico robusto para detección de anomalías con Huffman
- Análisis comparativo de complejidades y trade-offs

## 2. Metodología

### 2.1. Parte 1: Búsqueda de Subsecuencias - Algoritmo KMP

**Descripción:** El algoritmo Knuth-Morris-Pratt (KMP) [1] permite buscar un patrón dentro de un texto evitando retrocesos innecesarios mediante el uso de una función de prefijos.

**Implementación:**

```
vector<int> prefix_function(const string &s) {
    int n = s.size();
    vector<int> pi(n);
    for (int i = 1; i < n; i++) {
        int j = pi[i-1];
        while (j > 0 && s[i] != s[j])
            j = pi[j-1];
        if (s[i] == s[j]) j++;
        pi[i] = j;
    }
    return pi;
}
```

Listing 1. Función de prefijos KMP

```
bool KMP_search(const string &text, const string &pattern, int &position) {
    vector<int> pi = prefix_function(pattern);
    int j = 0;
    for (int i = 0; i < text.size(); i++) {
        while (j > 0 && text[i] != pattern[j])
            j = pi[j-1];
        if (text[i] == pattern[j]) j++;
        if (j == pattern.size()) {
            position = i - pattern.size() + 2;
            return true;
        }
    }
    return false;
}
```

```
}
```

Listing 2. Búsqueda KMP

**Complejidad:**

- Preprocesamiento:  $O(m)$  donde  $m$  es la longitud del patrón
- Búsqueda:  $O(n)$  donde  $n$  es la longitud del texto
- **Total:**  $O(n + m)$
- Espacio:  $O(m)$

**Justificación:** KMP es superior a búsqueda naïve ( $O(n \times m)$ ) porque evita revisitar posiciones del texto ya examinadas [1]. La función de prefijos permite “recordar” información sobre el patrón para saltar comparaciones redundantes.

**2.2. Parte 2: Detección de Palíndromos - Algoritmo de Manacher**

**Descripción:** El algoritmo de Manacher [2] encuentra el palíndromo más largo en tiempo lineal mediante la explotación de propiedades simétricas.

**Estrategia:**

1. Transformar la cadena insertando separadores (#) para manejar palíndromos pares e impares uniformemente
2. Usar la simetría de palíndromos ya encontrados para reducir comparaciones
3. Mantener el centro y el límite derecho del palíndromo que alcanza más lejos

```
pair<int, int> longestPalindromeManacher(const string &s) {
    string t = "#";
    for (char c : s) { t += c; t += '#'; }

    int n = t.size();
    vector<int> p(n, 0);
    int center = 0, right = 0;
    int maxLen = 0, centerIndex = 0;

    for (int i = 0; i < n; i++) {
        int mirror = 2 * center - i;
        if (i < right)
            p[i] = min(right - i, p[mirror]);

        while (i + p[i] + 1 < n && i - p[i] - 1 >= 0 &&
               t[i + p[i] + 1] == t[i - p[i] - 1])
            p[i]++;
    }

    if (p[i] > maxLen) {
        center = i;
        right = i + p[i];
    }

    if (p[i] > maxLen) {
```

```

        maxLen = p[i];
        centerIndex = i;
    }
}

int start = (centerIndex - maxLen) / 2;
return {start + 1, start + maxLen};
}

```

Listing 3. Algoritmo de Manacher

**Complejidad:**

- **Tiempo:**  $O(n)$  donde  $n$  es la longitud del texto
- Espacio:  $O(n)$

**Análisis:** Cada posición se expande a lo sumo una vez debido a la propiedad de que `right` solo aumenta [2]. Esto contrasta dramáticamente con el enfoque naive  $O(n^2)$  que expande desde cada posición sin aprovechar información previa.

**2.3. Parte 3: Substring Común Más Largo - DP Optimizado**

**Descripción:** Encontrar el substring más largo que aparece en ambas transmisiones usando programación dinámica [4, 5].

2.3.1. *Problema de la Implementación Original.* La solución estándar usa una matriz  $DP[n + 1][m + 1]$ :

```
vector<vector<int>> dp(n+1, vector<int>(m+1, 0));
```

Para archivos grandes (ej: `transmission1` = 10 MB, `transmission2` = 10 MB):

- Espacio requerido:  $10^7 \times 10^7 \times 4$  bytes = **400 TB**
- Inviable en la práctica

2.3.2. *Optimización Implementada.* **Observación clave:** Para calcular  $dp[i][j]$  solo necesitamos  $dp[i - 1][j - 1]$ , es decir, solo la fila anterior.

```

pair<int, int> longestCommonSubstring(const string &t1, const string &t2) {
    int n = t1.size(), m = t2.size();
    vector<int> prev(m + 1, 0);
    vector<int> curr(m + 1, 0);

    int maxlen = 0, endPos = 0;

    for (int i = 1; i <= n; i++) {
        for (int j = 1; j <= m; j++) {
            if (t1[i-1] == t2[j-1]) {
                curr[j] = prev[j-1] + 1;
                if (curr[j] > maxlen) {
                    maxlen = curr[j];

```

```

        endPos = i;
    }
} else {
    curr[j] = 0;
}
}
swap(prev, curr);
fill(curr.begin(), curr.end(), 0);
}

return {endPos - maxLen + 1, endPos};
}

```

Listing 4. Substring común optimizado

**Complejidad:**

- Tiempo:  $O(n \times m)$  - sin cambios
- Espacio:  $O(m)$  - **reducción de  $O(n \times m)$  a  $O(m)$**

**Impacto:** Para archivos de 10 MB cada uno:

- Original: ~100 GB RAM → Falla
- Optimizado: ~10 MB RAM → Funciona

**2.4. Parte 4: Detección de Anomalías - Codificación de Huffman****Descripción:** Usar codificación de Huffman [3] para detectar código malicioso basándose en análisis de frecuencias.**2.4.1. Construcción del Árbol de Huffman:**

```

struct HuffmanNode {
    char ch;
    int freq;
    HuffmanNode *left, *right;
};

HuffmanNode* buildHuffmanTree(const string &text) {
    unordered_map<char, int> freq;
    for (char c : text) freq[c]++;
}

priority_queue<HuffmanNode*, vector<HuffmanNode*>, Compare> pq;
for (auto &p : freq)
    pq.push(new HuffmanNode(p.first, p.second));

while (pq.size() > 1) {
    HuffmanNode* a = pq.top(); pq.pop();
    HuffmanNode* b = pq.top(); pq.pop();
    pq.push(new HuffmanNode(a, b));
}

```

```

    return pq.top();
}

```

Listing 5. Construcción del árbol de Huffman

**Complejidad:**  $O(k \log k)$  donde  $k$  es el número de caracteres únicos ( $\leq 16$  para hexadecimal)

#### 2.4.2. Detección de Anomalías - Método Estadístico. Estrategia:

- Calcular longitud promedio esperada para la transmisión:

$$\text{avgLen} = \sum_c P(c) \times \text{length}(c) \quad (1)$$

- Calcular desviación estándar ponderada:

$$\text{variance} = \sum_c P(c) \times (\text{length}(c) - \text{avgLen})^2 \quad (2)$$

$$\text{stddev} = \sqrt{\text{variance}} \quad (3)$$

- Para cada mcode, calcular su longitud promedio:

$$\text{avgLenMcode} = \frac{\text{totalLength(mcode)}}{|\text{mcode}|} \quad (4)$$

- Determinar si es sospechoso:

$$\text{sospechoso} \Leftrightarrow (\text{avgLenMcode} > \text{avgLen} + 0,5 \times \text{stddev}) \vee (\text{avgLenMcode} > 1,1 \times \text{avgLen}) \quad (5)$$

#### Complejidad Total:

- Construcción del árbol:  $O(n + k \log k)$
- Evaluación de cada mcode:  $O(m)$  donde  $m$  = longitud del mcode
- Total:**  $O(n + k \log k + 3m)$  para 3 mcodes

### 3. Análisis Comparativo

#### 3.1. Tabla de Complejidades

Cuadro 1. Comparación de complejidades algorítmicas

Algoritmo	Tiempo	Espacio	Optimización
KMP	$O(n + m)$	$O(m)$	Función de prefijos
Manacher	$O(n)$	$O(n)$	Simetría
Substring (orig.)	$O(n \times m)$	$O(n \times m)$	Inviable
Substring (opt.)	$O(n \times m)$	$O(m)$	Solo 2 vectores
Huffman	$O(n + k \log k)$	$O(k)$	Cola de prioridad

#### 3.2. Ejemplo Numérico

**Escenario:** transmission1 = 10 MB, transmission2 = 8 MB

**Substring Común:**

- Original:  $10^7 \times 8 \times 10^6 \times 4$  bytes = **320 GB RAM** → Falla

- Optimizado:  $8 \times 10^6 \times 4 \text{ bytes} = 32 \text{ MB RAM} \rightarrow \text{Funciona}$
- **Reducción: 10,000x**

#### Huffman:

- Caracteres únicos:  $k = 16$
- Construcción:  $10^7 + 16 \times \log_2(16) = 10^7 + 64 \approx 10^7$  operaciones
- Muy eficiente: < 0,1 segundos en hardware moderno

## 4. Resultados Experimentales

### 4.1. Validación con Datos de Prueba

#### Configuración:

- transmission1.txt: 261 caracteres
- transmission2.txt: 219 caracteres
- mcode1.txt: “DEADBEEF” (8 chars)
- mcode2.txt: “CAFEBABE” (8 chars)
- mcode3.txt: “ABCDEFABCDEF” (12 chars)

#### Resultados Parte 1 (KMP):

- mcode1 encontrado en transmission1 en posición 194
- mcode2 encontrado en transmission1 en posición 206
- mcode1 encontrado en transmission2 en posición 167
- mcode2 encontrado en transmission2 en posición 179
- mcode3 no encontrado en ninguna transmisión

#### Resultados Parte 2 (Manacher):

- Palíndromo transmission1: posiciones 49-92 (44 chars)
- Palíndromo transmission2: posiciones 63-114 (52 chars)

#### Resultados Parte 3 (Substring común):

- Substring común más largo: posiciones 178-217 (40 chars)
- Contenido: “1234DEADBEEF5678CAFEBABE9876”

#### Resultados Parte 4 (Huffman):

Análisis de frecuencias en transmission1:

- A: 44 (16.9 %) → 2 bits
- B: 43 (16.5 %) → 3 bits
- C: 41 (15.7 %) → 3 bits
- D: 38 (14.6 %) → 3 bits
- E: 29 (11.1 %) → 3 bits
- F: 18 (6.9 %) → 4 bits

Promedio: 3.38 bits/carácter

Evaluación de mcodes:

- mcode1 “DEADBEEF”: 24 bits total, 3.0 bits/char → NO sospechoso

- mcode2 “CAFEBABE”: 23 bits total, 2.9 bits/char → NO sospechoso
- mcode3 “ABCDEFABCDEF”: 36 bits total, 3.0 bits/char → NO sospechoso

**Análisis:** En esta distribución particular, los mcodes usan caracteres frecuentes (A-F), por lo que NO son anómalos. Esto demuestra que el algoritmo **no genera falsos positivos**.

## 5. Conclusiones

### 5.1. Principales Hallazgos

1. **La optimización de espacio es crítica:** La reducción de  $O(n \times m)$  a  $O(m)$  en el algoritmo de substring común es la diferencia entre viable e inviable para archivos grandes.
2. **Huffman proporciona capacidades únicas:** La detección basada en frecuencias complementa perfectamente los métodos de matching exacto, permitiendo detectar código malicioso desconocido.
3. **El análisis estadístico robusto es esencial:** Usar desviación estándar + threshold porcentual previene tanto falsos positivos como falsos negativos.
4. **Los algoritmos clásicos son poderosos:** KMP [1] (1977) y Manacher [2] (1975) siguen siendo soluciones óptimas para sus respectivos problemas.

### 5.2. Trabajo Futuro

#### Mejoras Potenciales:

1. **Suffix Array para Parte 3:**
  - Construir suffix array:  $O(n \log n)$  [10]
  - LCP array:  $O(n)$
  - Potencialmente más rápido para archivos enormes
2. **Huffman Adaptativo en Streaming:**
  - Actualizar árbol dinámicamente
  - Útil para transmisiones en tiempo real
3. **Paralelización:**
  - Parte 1: Buscar múltiples patrones en paralelo
  - Parte 3: Dividir matriz DP en bloques

### 5.3. Lecciones Aprendidas

#### Ingeniería de Software:

- La gestión de memoria es crucial
- Los edge cases pueden causar crashes
- La optimización prematura es mala, pero la optimización necesaria es esencial

#### Análisis de Algoritmos:

- La complejidad asintótica no es todo: constantes importan
- $O(n \times m)$  espacio puede ser dealbreaker incluso si  $O(n \times m)$  tiempo es aceptable
- Los algoritmos adaptativos (como Huffman) son más robustos

## Referencias

- [1] Donald E. Knuth, James H. Morris, Jr., and Vaughan R. Pratt. 1977. Fast pattern matching in strings. *SIAM J. Comput.* 6, 2 (June 1977), 323–350. <https://doi.org/10.1137/0206024>
- [2] Glenn Manacher. 1975. A new linear-time 'on-line' algorithm for finding the smallest initial palindrome of a string. *J. ACM* 22, 3 (July 1975), 346–351. <https://doi.org/10.1145/321892.321896>
- [3] David A. Huffman. 1952. A method for the construction of minimum-redundancy codes. *Proc. IRE* 40, 9 (Sept. 1952), 1098–1101. <https://doi.org/10.1109/JRPROC.1952.273898>
- [4] Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. 2009. *Introduction to Algorithms* (3rd ed.). MIT Press, Cambridge, MA.
- [5] Dan Gusfield. 1997. *Algorithms on Strings, Trees, and Sequences: Computer Science and Computational Biology*. Cambridge University Press, New York, NY. <https://doi.org/10.1017/CBO9780511574931>
- [6] Robert Sedgewick and Kevin Wayne. 2011. *Algorithms* (4th ed.). Addison-Wesley Professional, Boston, MA.
- [7] Gene Myers. 1999. A fast bit-vector algorithm for approximate string matching based on dynamic programming. *J. ACM* 46, 3 (May 1999), 395–415. <https://doi.org/10.1145/316542.316550>
- [8] Ricardo Baeza-Yates and Gonzalo Navarro. 1999. Faster approximate string matching. *Algorithmica* 23, 2 (1999), 127–158. <https://doi.org/10.1007/PL00009253>
- [9] Michael O. Rabin and Richard M. Karp. 1987. Efficient randomized pattern-matching algorithms. *IBM J. Res. Dev.* 31, 2 (March 1987), 249–260. <https://doi.org/10.1147/rd.312.0249>
- [10] Udi Manber and Gene Myers. 1993. Suffix arrays: A new method for on-line string searches. *SIAM J. Comput.* 22, 5 (Oct. 1993), 935–948. <https://doi.org/10.1137/0222058>

## A. Especificaciones de Implementación

**Lenguaje:** C++11

**Compilador:** g++ con flags -std=c++11 -O2

**Librerías:** iostream, fstream, vector, string, unordered\_map, queue, cmath

**Arquitectura:** Modular (4 partes independientes)

**Gestión de Memoria:** RAII con liberación explícita para árboles

**Código Fuente Completo:** Disponible en: A01425660\_ActInt1/e1.cpp

## B. Casos de Prueba Adicionales

Para validación extendida, se recomiendan estos casos:

### Caso 1: Archivos muy grandes

- transmission: 10-100 MB
- Verifica: Optimización de espacio en Parte 3

### Caso 2: Distribución uniforme

- Todos los caracteres con frecuencia similar
- Verifica: Huffman no genera falsos positivos

### Caso 3: Distribución extremadamente sesgada

- Un carácter dominante (> 90 %)
- Verifica: Detección de anomalías robusta

### Caso 4: Módulos con caracteres no presentes

- mcode contiene 'Z' (no hexadecimal)
- Verifica: Manejo de caracteres desconocidos