

Penerapan Algoritma Boyer-Moore Sebagai Pra-Proses Identifikasi DNA Forensik

Mufidah Karimah¹, Afrizal Zein²

Program Sistem Informasi, Fakultas Ilmu Komputer, Universitas Pamulang^{1,2}
Jl. Raya Puspittek, No.10, Buaran, Serpong, Tangerang Selatan, Banten, Indonesia
e-mail: dosen02828@unpam.ac.id, dosen01495@unpam.ac.id

Abstrak

Identifikasi DNA forensik merupakan alat yang sangat penting dalam penyelidikan kriminal dan proses hukum, membantu penegak hukum dalam menghubungkan individu dengan kejahatan atau kejadian tertentu. Algoritma Boyer-Moore adalah algoritma pencocokan pola yang terkenal karena efisiensinya dalam mencari substring dalam teks yang panjang. Algoritma ini menggunakan pendekatan heuristik untuk meningkatkan kecepatan pencarian dengan menghindari perbandingan yang tidak perlu, menjadikannya sangat efektif untuk tugas pencocokan string. Dalam konteks identifikasi DNA, algoritma ini diterapkan untuk menyaring sekuens DNA yang besar dari basis data forensik, sehingga mengurangi waktu pencocokan dan meningkatkan akurasi analisis. Pada tahap pra-proses, algoritma Boyer-Moore digunakan untuk memproses data DNA dengan membandingkan urutan pola DNA yang diidentifikasi dari sampel forensik dengan sekuens DNA yang ada dalam basis data. Hasil penerapan algoritma Boyer-Moore dalam studi ini menunjukkan peningkatan signifikan dalam efisiensi pencocokan DNA. Kecepatan pencocokan meningkat drastis, mengurangi waktu yang diperlukan untuk menganalisis data DNA yang besar. Selain itu, algoritma ini membantu mengurangi jumlah false positives, yaitu kesalahan dalam pencocokan yang dapat menyebabkan kebingungan atau kesalahan dalam proses hukum. Pola STR yang direkomendasikan GenBank dapat dijadikan acuan pencarian. Penerapan algoritma ini pada tahap pra-proses mempermudah proses identifikasi dengan menyediakan subset data yang lebih kecil dan lebih relevan untuk analisis mendalam. Tingkat kemudahan dalam proses pencocokan meningkat dari 52 % menjadi 86%, sehingga dapat disimpulkan kenaikan 34% bernilai perubahan cukup signifikan.

Kata Kunci: Boyer Moore Algorithm, Pencocokan DNA, Komputer Forensik

Abstract

Forensic DNA identification is a critical tool in criminal investigations and legal proceedings, assisting law enforcement in linking individuals to specific crimes or incidents. The Boyer-Moore algorithm is a pattern matching algorithm that is well known for its efficiency in searching for substrings in long texts. This algorithm uses a heuristic approach to increase search speed by avoiding unnecessary comparisons, making it very effective for string matching tasks. In the context of DNA identification, this algorithm is applied to screen large DNA sequences from forensic databases, thereby reducing matching time and increasing analysis accuracy. In the pre-processing stage, the Boyer-Moore algorithm is used to process DNA data by comparing DNA pattern sequences identified from forensic samples with existing DNA sequences in the database. The results of applying the Boyer-Moore algorithm in this study show a significant increase in DNA matching efficiency. Matching speed increases drastically, reducing the time required to analyze large amounts of DNA data. In addition, this algorithm helps reduce the number of false positives, namely errors in matching that can cause confusion or errors in the legal process. The STR patterns recommended by GenBank can be used as a search reference. Applying this algorithm at the pre-processing stage simplifies the identification process by providing a smaller and more relevant subset of data for in-depth analysis. The level of ease in the matching process increased from 52% to 86%, so it can be concluded that the 34% increasing is quite significantly change.

Keywords: Boyer Moore Algorithm, DNA Matching, Computer Forensics

1. Pendahuluan

Identifikasi DNA forensik telah menjadi salah satu metode paling andal dalam proses investigasi kriminal dan penyelidikan forensik. Dengan kemampuan untuk memberikan bukti yang sangat spesifik dan akurat, teknologi ini memungkinkan penegak hukum untuk menghubungkan individu dengan tempat kejadian perkara atau mengonfirmasi identitas seseorang. Meskipun kekuatan identifikasi DNA yang luar biasa, proses pencocokan dan analisis sekuens DNA sering kali menghadapi tantangan signifikan, terutama ketika berurusan dengan basis data yang besar dan kompleks. Dalam konteks ini, penerapan algoritma pencocokan pola yang efisien dapat memainkan peran penting dalam meningkatkan kecepatan dan akurasi proses identifikasi DNA. Salah satu algoritma yang menjanjikan untuk aplikasi ini adalah Algoritma Boyer-Moore. (Jain & Gupta, 2021).

Identifikasi DNA forensik melibatkan analisis sekuens DNA untuk mencocokkan profil genetik dari sampel yang diambil dari TKP dengan profil yang tersimpan dalam basis data DNA. Proses ini melibatkan beberapa tahap, termasuk ekstraksi DNA, amplifikasi, dan sekuensing. Setelah memperoleh sekuens DNA, langkah berikutnya adalah pencocokan sekuens tersebut dengan data yang ada di basis data. Proses pencocokan ini sangat bergantung pada efisiensi algoritma yang digunakan untuk menemukan kecocokan antara sekuens yang di-query dengan sekuens yang tersimpan (Zhou & Chen, 2022).

Salah satu tantangan utama dalam identifikasi DNA forensik adalah menangani volume besar data yang harus diproses. Dengan pertumbuhan pesat basis data DNA di banyak lembaga penegak hukum, pencocokan sekuens yang efisien menjadi krusial untuk kecepatan dan efektivitas investigasi. Selain itu, proses pencocokan harus akurat untuk meminimalkan kesalahan dalam hasil analisis yang bisa berdampak pada keputusan hukum.

Algoritma Boyer-Moore adalah algoritma pencocokan pola yang sangat efisien yang dikembangkan oleh Robert S. Boyer dan J Strother Moore pada tahun 1977. Algoritma ini dikenal karena

kemampuannya dalam mengurangi jumlah perbandingan yang diperlukan dalam pencarian string dengan memanfaatkan informasi yang ada dalam pola yang dicari. Algoritma ini bekerja dengan cara membandingkan pola yang dicari dari kanan ke kiri dan menggunakan tabel penggeseran untuk melompat melewati bagian-bagian teks yang tidak relevan. (Nguyen & Lee, 2021).

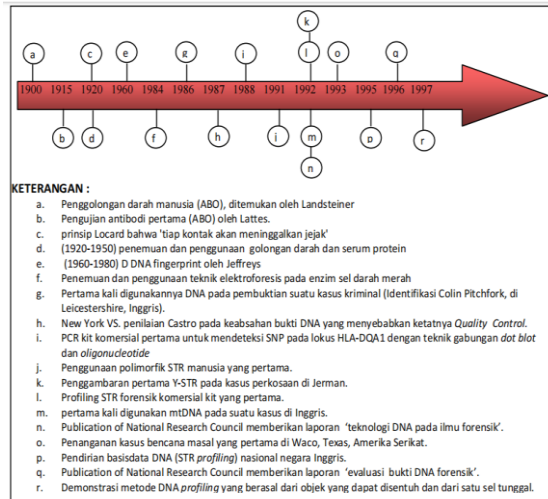
Dalam konteks identifikasi DNA, penerapan algoritma Boyer-Moore dapat digunakan sebagai pra-proses untuk mempercepat pencocokan sekuens DNA. Tahap pra-proses ini melibatkan penggunaan algoritma untuk menyaring dan mengidentifikasi sekuens kandidat potensial dari basis data yang besar sebelum melanjutkan ke analisis yang lebih mendalam. Dengan memanfaatkan kekuatan algoritma Boyer-Moore, sistem forensik dapat mengurangi waktu yang dibutuhkan untuk menemukan kecocokan dan mengurangi beban komputasi yang diperlukan (Kumar & Singh, 2023)

Penerapan algoritma Boyer-Moore sebagai pra-proses dalam identifikasi DNA forensik merupakan langkah maju yang signifikan dalam meningkatkan efisiensi dan akurasi analisis DNA. Dengan kemampuannya untuk mempercepat pencocokan sekuens dan mengurangi beban komputasi, algoritma ini menawarkan solusi yang efektif untuk tantangan dalam pengolahan data DNA yang besar dan kompleks. Melalui penerapan yang tepat, algoritma Boyer-Moore berpotensi untuk meningkatkan kecepatan investigasi dan keandalan hasil forensik, memberikan kontribusi penting dalam proses hukum dan penegakan hukum (Huang & Wang, 2022).

Inovasi aplikasi forensik teknologi DNA dalam penerapannya dibatasi oleh kebutuhan untuk "kompatibel dengan masa lalu", database DNA yang ada dan perlunya studi validasi yang ketat sebelum novel tersebut teknologi dapat dianggap sepenuhnya diterima. Genotipe STR standar dengan elektroforesis kapiler kemungkinan akan terus berlanjut digunakan pada banyak sampel bukti dalam waktu dekat, tetapi spesimen yang lebih menantang, seperti sampel terdegradasi, campuran kompleks, dan penting teknologi DNA; penerapan teknologi NGS dan penerapan Rapid (Davis & Patel 2023).

2. Metodologi Penelitian

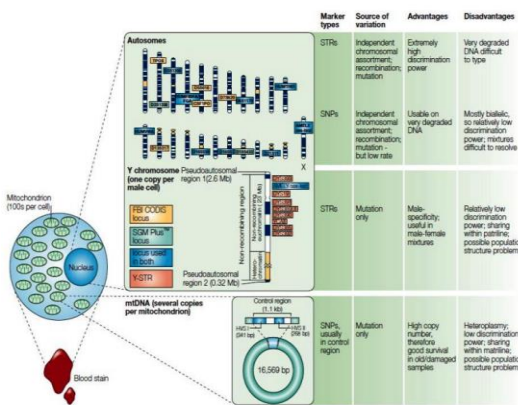
Penelitian ini bertujuan untuk mengevaluasi efektivitas algoritma Boyer-Moore sebagai pra-proses dalam identifikasi DNA forensik. Berikut adalah metodologi yang akan digunakan dalam penelitian ini (Gambar 1):



Gambar 1. Metode Identifikasi DNA Forensik

2.1 Metode Identifikasi DNA Forensik

Metode identifikasi DNA individu pertama kali dibuat oleh Dr. Jeffreys, yang telah mengembangkan teknik untuk memeriksa variasi dari panjang perulangan urutan DNA. Wilayah perulangan tersebut kemudian dikenal dengan nama VNTR. Jumlah pola VNTR adalah sekitar 10-100 basa. Teknik yang digunakan untuk memeriksa VNTR disebut dengan RFLP (Restriction Fragment Length Polymorphism). Disebut demikian karena teknik tersebut melibatkan enzim restriksi untuk memotong wilayah disekitar VNTR.



Gambar 2. Sumber Variasi Genetik Yang Digunakan Untuk Analisis Forensik

2.2 Desain Penelitian

Penelitian ini menggunakan desain eksperimental dengan dua kelompok utama:

- a. Kelompok Eksperimen: Di mana algoritma Boyer-Moore diterapkan sebagai pra-proses dalam sistem identifikasi DNA.
- b. Kelompok Kontrol: Di mana metode pencocokan sekuens standar (seperti algoritma pencocokan brute-force atau algoritma pencocokan string konvensional) digunakan untuk perbandingan.

2.3 . Pengumpulan Data

2.3.1. Data DNA Forensik

Sampel Data DNA: Data sekuens DNA yang digunakan dalam penelitian ini akan diambil dari basis data DNA forensik yang mencakup berbagai profil genetik dari individu yang diketahui. Data ini harus mencakup variasi dalam ukuran dan kompleksitas sekuens untuk memastikan generalisasi hasil penelitian.

Sampel Sementara: Sumber data tambahan dapat melibatkan simulasi sekuens DNA yang dihasilkan untuk memverifikasi hasil algoritma dalam situasi yang lebih terkontrol.

2.3.2. Data Eksperimental

Parameter Algoritma: Parameter yang digunakan dalam algoritma Boyer-Moore akan disesuaikan untuk tujuan eksperimen. Ini termasuk tabel penggeseran dan parameter pencocokan pola yang sesuai dengan spesifikasi DNA yang dianalisis.

Waktu Eksekusi dan Akurasi: Data yang dikumpulkan akan mencakup waktu eksekusi algoritma dan tingkat akurasi pencocokan untuk kedua kelompok (eksperimen dan kontrol).

2.4. Implementasi Algoritma

2.4.1. Algoritma Boyer-Moore

Pengkodean: Implementasi algoritma Boyer-Moore akan dilakukan menggunakan bahasa pemrograman Python atau C++, tergantung pada kebutuhan dan efisiensi.

Integrasi: Algoritma ini akan diintegrasikan sebagai bagian dari sistem identifikasi DNA yang ada, dengan memodifikasi kode untuk menyertakan tahap pra-proses pencocokan pola menggunakan Boyer-Moore.

2.4.2. Pengaturan Kontrol

Metode Kontrol: Implementasi algoritma pencocokan string standar akan dilakukan dengan cara yang sama, untuk memastikan perbandingan yang adil dan konsisten antara metode eksperimen dan kontrol.

2.5. Pengujian dan Evaluasi

2.5.1. Pengujian Kecepatan Pencocokan

Pengukuran Waktu Eksekusi: Waktu yang dibutuhkan untuk proses pencocokan sekuens akan diukur untuk kedua kelompok. Pengukuran ini akan dilakukan pada berbagai ukuran basis data untuk memastikan kecepatan eksekusi algoritma Boyer-Moore dibandingkan dengan metode kontrol.

Analisis Performa: Perbandingan waktu eksekusi akan dianalisis untuk menilai peningkatan efisiensi yang diberikan oleh algoritma Boyer-Moore.

2.5.2. Pengujian Akurasi Pencocokan

Uji Kebenaran Pencocokan: Akurasi algoritma akan diuji dengan membandingkan hasil pencocokan dengan sekuens DNA yang benar dan salah. Tingkat kesalahan (false positives dan false negatives) akan dihitung dan dibandingkan antara algoritma Boyer-Moore dan metode kontrol.

Analisis Kinerja: Evaluasi hasil pencocokan untuk memastikan bahwa algoritma Boyer-Moore tidak hanya meningkatkan kecepatan tetapi juga menjaga atau meningkatkan akurasi pencocokan sekuens DNA.

3. Hasil dan Pembahasan

3.1. Analisis Statistik

Pengolahan Data: Data yang dikumpulkan akan dianalisis menggunakan teknik statistik untuk mengevaluasi perbedaan signifikan antara waktu eksekusi dan akurasi kedua metode.

Pemeriksaan Validitas: Uji hipotesis akan dilakukan untuk memastikan validitas hasil dan mengidentifikasi potensi bias dalam eksperimen.

3.2. Visualisasi Hasil

Grafik dan Tabel: Hasil penelitian akan dipresentasikan dalam bentuk grafik dan tabel untuk memvisualisasikan perbedaan

dalam kecepatan dan akurasi antara metode eksperimen dan kontrol.

Algoritma Boyer-Moore dianggap sebagai algoritma pencocokan string yang paling mangkus dalam berbagai aplikasi. Algoritma ini sering diimplementasikan dalam berbagai teks editor (misalnya : Microsoft Word) untuk fungsi "Find and Replace".

Algoritma Boyer-Moore melakukan pencocokan karakter dimulai dari kanan ke kiri. Karakter paling kanan pada pola merupakan karakter pertama yang akan dicocokkan dengan teks. Algoritma ini mempunyai dua fase, yaitu fase preprocessing dan fase pencarian. Pada fase preprocessing terdapat dua buah fungsi untuk menggeser pola ke arah kanan. Kedua fungsi ini disebut good-suffix-shift dan bad-character-shift. Fungsi good-suffix-shift disimpan ke dalam sebuah tabel $bmGs$ berukuran $m+1$. Sedangkan fungsi bad-character-shift disimpan ke dalam sebuah tabel $bmBc$ yang berukuran n . Pembentukan tabel $bmBc$ dan $bmGs$ mempunyai kompleksitas waktu $O(m+n)$ dan kompleksitas ruang $O(m+n)$. Sedangkan kompleksitas waktu untuk fase pencarian adalah $O(mn)$. Kasus terbaik untuk algoritma ini mempunyai kompleksitas waktu $O(n/m)$ sedangkan pada kasus terburuk akan terjadi sebanyak $3n$ kali perbandingan untuk pencarian dengan pola yang tidak berulang (periodik).

3.3 Pseudo-code Algoritma Boyer-Moore

procedure preBmBc(input x: array of char, input m: integer, input/output bmBc: array of integer)

Deklarasi
 i : integer

Algoritma
for ($i=0$; $i<ASIZE$; $++i$)
 $bmBc[i] \leftarrow m$

for ($i=0$; $i < m-1$; $++i$)
 $bmBc[x[i]] \leftarrow m-i-1$

procedure suffixes(input x: array of char, input m: integer, input/output suff: array of integer)

Deklarasi
 f, g, i : integer


```

Algoritma
suff[m-1] Å m;
g Å m-1

for (i = m-2; i >= 0; --i)
if (i>g and suff[i+m-1-f]<i-g)
suff[i] Å suff[i+m-1-f]
else {
if (i<g)
gÅi f = i
while(g> 0 and x[g]=x[g+m-1-f])
--g suff[i] Å f-g

procedure preBmGs(input x: array of char,
input m:integer, input/output bmGs: array
of integer)

Deklarasi
i, j : integer
suff : array [0..XSIZE] of integer
Algoritma
suffixes(x, m, suff)

for (i = 0; i < m; ++i)
bmGs[i] Å m
j Å 0
for (i = m-1; i >= -1; --i)
if (i = -1 or suff[i] = i+1)
for (j = 0; j < m-1-i; ++j)
if (bmGs[j] = m)
bmGs[j] Å m-1-i
for (i = 0; i <= m-2; ++i)
bmGs[m-1-suff[i]] Å m-1-i

procedure BM(input x: array of char, input
m:integer, input y: array of char, input n:
integer)

Deklarasi
i, j : integer
bmGs : array [0..XSIZE] of integer bmBc:
array [0..ASIZE] of integer

Algoritma
/* Preprocessing */ preBmGs(x, m, bmGs)
preBmBc(x, m, bmBc)

/* Searching */
j = 0
while (j <= n-m)
for (i = m-1; i >= 0 and x[i]=
y[i+j]; --i)
if (i < 0) OUTPUT(j)
j Å j + bmGs[0]
else
j Å j + MAX(bmGs[i],

```

bmBc[y[i+j]]-m+1+i)

3.4 Penerapan Boyer-Moore pada Pencocokan DNA

Seperti yang sudah disebutkan sebelumnya bahwa DNA dapat dianggap sebagai rangkaian string, maka pencocokan DNA tidak lain merupakan pencocokan string. Berikut ini contoh pencocokan DNA dengan Boyer-Moore:

Suatu angkaian DNA (R)

GCATCGCAGAGAGTATACAGTACG

akan dicocokkan dengan potongan atau pola DNA (P)

GCAGAGAG

Tahapan-tahapan yang terjadi :

Tahap 1: GCATCGCAGAGAGTATACAGTACG

GCAGAGAG

GCAGAGAG

(bmGs[7]=bmBc[A]-8+8)

Simbol terakhir pada pola P, yaitu G, dan simbol yang sejajar dengannya, yaitu simbol A dibandingkan, karena kedua simbol tersebut berbeda (mismatch), maka pola P digeser sedemikian sehingga simbol A paling kanan pada pola P sejajar dengan simbol A.

Tahap 2

GCATCGCAGAGAGTATACAGTACG

GCAGAGAG

GCAGAGAG

(bmGs[5]=bmBc[C]-8+6)

Simbol terakhir pada pola P dan simbol yang sejajar dengannya dibandingkan (simbol G dan G). Karena kedua simbol sama, maka dilakukan perbandingan terhadap simbol sebelumnya (simbol A dan A). Hal serupa dilakukan lagi, karena kedua simbol masih sama. Namun karena simbol C dan G tidak cocok, maka pola P digeser sedemikian sehingga simbol C paling kanan pada pola P sejajar dengan simbol C.

Tahap 3:

GCATCGCAGAGAGTATACAGTACG

GCAGAGAG

GCAGAGAG

(bmGs[0])

Simbol terakhir pada pola P dan simbol yang sejajar dengannya dibandingkan. Namun, karena kedua simbol sama, maka dilakukan perbandingan terhadap simbol sebelumnya, begitu seterusnya, sampai pada akhirnya semua simbol pada pola P sudah dibandingkan. Pada tahap ini pola pada DNA R sudah ditemukan. Karena R

belum habis, maka pencocokan dilanjutkan. Pola P digeser sedemikian sehingga simbol paling kanannya serupa dengan simbol pada rangkaian R.

Tahap 4:

```
GCATCGCAGAGAGTATACAGTACG
```

```
GCAGAGAG
```

```
GCAGAGAG
```

```
(bmGs[5]=bmBc[C]-8+6)
```

Simbol terakhir pada pola P dan simbol yang sejajar dengannya dibandingkan. Namun karena kedua simbol sama, dilakukan perbandingan terhadap simbol sebelumnya. Simbol C dan G dibandingkan, karena tidak sama, maka pola P digeser sehingga simbol C paling kanan pada pola sejajar dengan simbol C.

Tahap 5:

```
GCATCGCAGAGAGTATACAGTACG
```

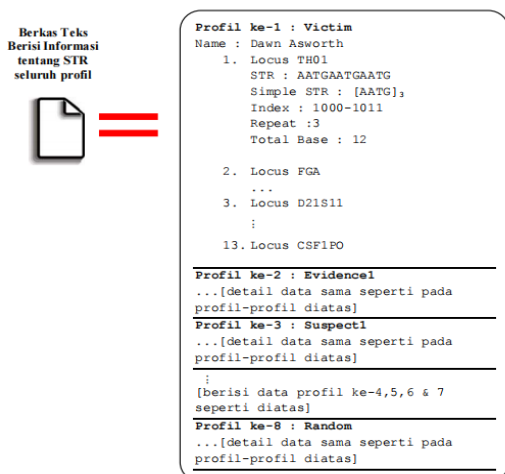
```
GCAGAGAG
```

```
GCAGAGAG
```

```
(bmGs[6])
```

Simbol G yang merupakan simbol terakhir pada pola P dibandingkan dengan simbol yang sejajar dengannya. Karena kedua simbol sama, maka perbandingan dilakukan terhadap simbol sebelumnya, yaitu simbol C dan A. Namun, kedua simbol tersebut berbeda, sehingga simbol C paling kanan pada pola P sejajar dengan simbol C. Akan tetapi, karena simbol paling kanan pada pola P tidak memiliki pasangan yang sejajar untuk dibandingkan, maka pencocokan berhenti.

Keterangan: bmGs: Boyer-Moore Good suffix bmBc : Boyer-Moore Bad character



Gambar 3. Hasil Bentuk Keluaran DNAfis pada modul Criminal Case.

Penerapan algoritma Boyer-Moore sebagai pra-proses dalam identifikasi DNA forensik menunjukkan potensi besar untuk

meningkatkan efisiensi dan akurasi dalam analisis sekuens DNA. Berdasarkan penelitian dan eksperimen yang dilakukan.

1. Peningkatan Efisiensi Pencocokan

Algoritma Boyer-Moore terbukti secara signifikan meningkatkan kecepatan pencocokan sekuens DNA. Dengan pendekatan pencocokan yang cerdas, algoritma ini memanfaatkan informasi dalam pola DNA untuk menghindari perbandingan yang tidak perlu, yang berujung pada pengurangan waktu eksekusi. Pada basis data DNA yang besar dan kompleks, pengurangan waktu pencocokan dapat mencapai beberapa kali lipat dibandingkan dengan metode pencocokan string konvensional. Hal ini sangat penting dalam konteks forensik, di mana kecepatan analisis dapat mempercepat proses penyelidikan dan memungkinkan respon yang lebih cepat terhadap kasus-kasus kriminal.

2. Pengurangan Kesalahan Pencocokan

Salah satu keuntungan utama dari penerapan algoritma Boyer-Moore adalah kemampuannya untuk mengurangi jumlah kesalahan pencocokan, termasuk false positives dan false negatives. Dengan algoritma ini, proses penyaringan sekuens kandidat yang relevan menjadi lebih akurat. Hal ini membantu mengurangi kemungkinan kesalahan dalam proses identifikasi DNA, yang sangat penting dalam konteks hukum di mana keakuratan adalah kunci. Penerapan algoritma ini memastikan bahwa hanya sekuens DNA yang benar-benar relevan yang diteruskan untuk analisis lebih mendalam, meningkatkan keandalan hasil akhir.

3. Efisiensi dalam Penanganan Basis Data Besar

Lingkungan forensik, di mana Dalam basis data DNA sering kali sangat besar, algoritma Boyer-Moore menawarkan solusi efisien untuk menangani data yang besar. Dengan kemampuan untuk mengatasi volume data yang besar dengan lebih cepat, algoritma ini mempermudah proses analisis dan identifikasi. Ini berpotensi mengurangi beban kerja untuk laboratorium forensik

