

# Penggunaan Pattern Matching dalam Penentuan Gen pada BLAST

Rama Febriyan 13511067  
Program Studi Teknik Informatika  
Sekolah Teknik Elektro dan Informatika  
Institut Teknologi Bandung, Jl. Ganesha 10 Bandung 40132, Indonesia  
13511067@std.stei.itb.ac.id

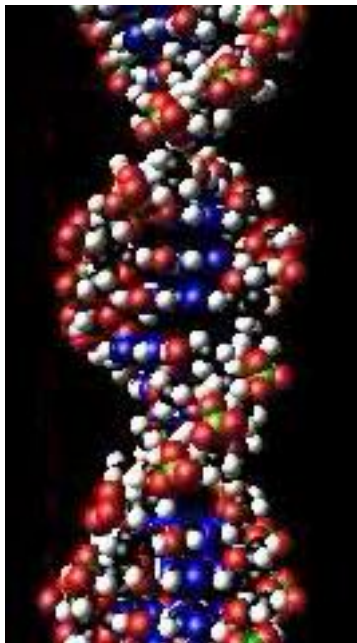
**Abstract**—Makalah ini membahas tentang penerapan strategi *pattern matching* dalam pencocokan struktur DNA yang digunakan pada kaskas bernama Basic Local Alignment Search Tool (BLAST).

**Index Terms**—gen, BLAST, DNA, FASTA, Smith-Waterman.

## I. PENDAHULUAN

Makhluk hidup, termasuk manusia memiliki ciri khas masing-masing. Semua ciri khas tersebut berasal dari sebuah komponen kecil yang di sebut DNA.

Secara garis besar, DNA di dalam sebuah sel adalah sebagai materi genetik, dengan kata lain, DNA merupakan *blueprint* bagi setiap aktivitas sel. Hal ini berlaku bagi seluruh organisme.



Gambar 1 DNA

Setiap makhluk hidup memiliki ciri khas tertentu pada DNA mereka. Dengan mengetahui perbedaan dan ciri khas tersebut, DNA dapat digunakan dalam berbagai bidang. Misalnya dalam forensik, DNA yang terdapat dalam darah, sperma, kulit, liur, dapat digunakan untuk

melakukan identifikasi pelaku kejahatan. Proses ini di sebut fingerprinting genetika. Metode ini merupakan tekni yang paling terpercaya meskipun tidak selalu sempurna, misal jika tempat kejadian merupakan tempat yang terkontaminasi oleh DNA dari banyak orang.

Selain bidang forensik, keunikan DNA juga dapat digunakan di bidang bioinformatika. Info yang terdapat di dalam DNA dapat di gunakan untuk mengidentifikasi sebuah DNA tak dikenal yang ditemukan di sebuah lokasi penelitian. Dengan melakukan pencocokan terhadap DNA yang terdapat pada basis data dengan DNA temuan, dapat diketahui apakah DNA temuan tersebut berasal dari organisme yang sudah dikenal atau sebuah organisme baru.

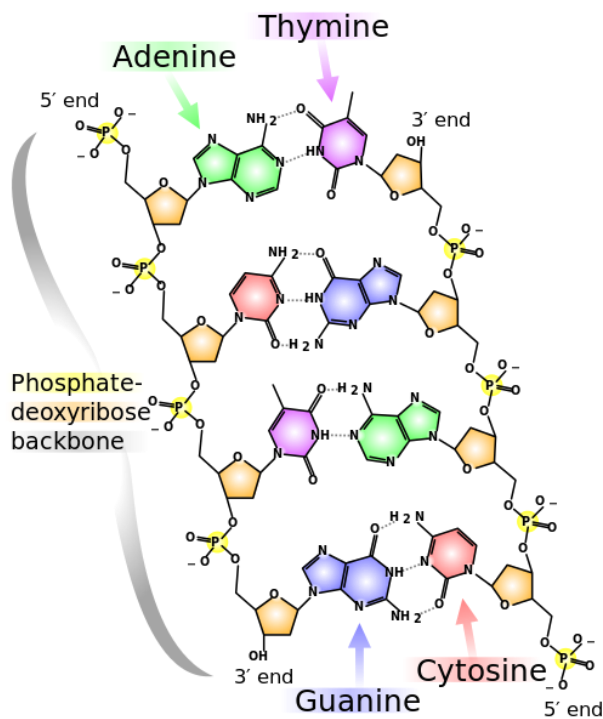
## II. DNA

DNA (deoxyribonucleic acid) atau asam deoksiribonucleat adalah sejenis asam nukleat yang tergolong biomolekul utama dari organisme. Pada umumnya, DNA terletak di dalam inti sel

DNA terdiri dari tiga komponen utama, yaitu:

- Gugus fosfat
- Gula deoksiribosa
- Basa nitrogen, yang terdiri dari:
  - ◆ Adenin (A)
  - ◆ Guanin (G)
  - ◆ Sitosin (S)
  - ◆ Timin (T)

Satu unit monomer DNA dinamakan nukleotida. Basa nitrogen pada rantai DNA saling berkaitan. Adenin berkaitan dengan Timin, sedangkan Guanin berkaitan dengan Sitosin.



Gambar 2 Struktur kimia DNA

### III. STRING MATCHING

String matching merupakan sebuah teknik yang digunakan untuk melakukan untuk mencari sebuah pola tertentu dalam teks

Untuk melakukan proses ini dibutuhkan 2 jenis data yaitu:

1. Teks, yaitu string sepanjang n karakter sebagai sumber pencarian
2. Pattern, yaitu string sepanjang m dengan  $m < n$  yang akan dicari di dalam teks

#### Contoh:

Pattern: rumah  
 Teks: kami pulang ke **rumah** sore hari

Bagian bercetak tebal di atas adalah string yang ditemukan sesuai pola.

String matching dapat di temukan dalam kehidupan sehari-hari. Google Search yang sering kita gunakan, memanfaatkan string matching dalam melakukan pencarian. Pertama pengguna memasukkan kata kunci yang akan di cari. Kata kunci ini digunakan sebagai *pattern*. Kemudian dengan menggunakan algoritma tertentu, kata kunci di cocokkan dengan data yang dimiliki oleh Google. Semua dokumen yang memiliki kecocokan dengan *pattern* yang telah diberikan, ditampilkan sebagai hasil pencarian dengan pola pengurutan tertentu.

String matching juga digunakan dalam kriminalitas. Pada tempat kejadian perkara, terkadang tersisa sidik jari dari pelaku kejahatan. Sidik jari tersebut diambil dan disimpan dalam format tertentu. Kemudian hasil simpanan diubah dalam bentuk string dan dibandingkan dengan

string yang terdapat dalam basis data sidik jari yang dimiliki pihak kepolisian. Dengan ini tersangka kejahatan dapat dibatasi ruang pencariannya.

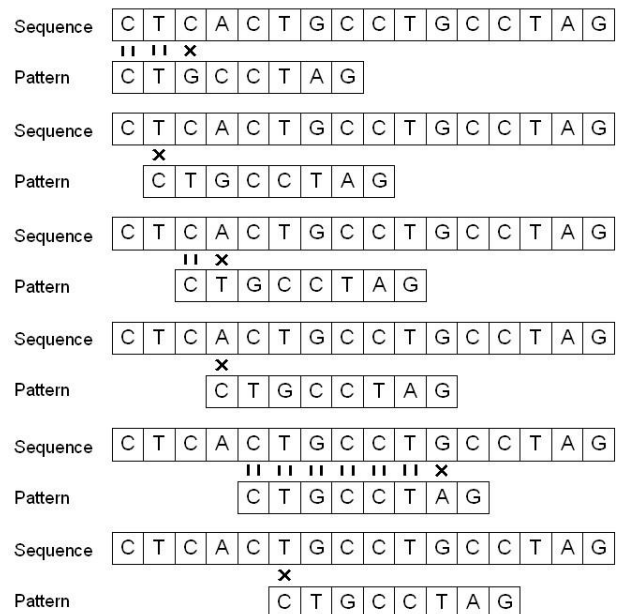
### IV. ALGORITMA STRING MATCHING

Terdapat banyak algoritma yang digunakan dalam string matching. Seperti *Brute-force*, *Knuth-Morris-Pratt*, *Boyer-Moore*. Algoritma-algoritma ini memiliki karakteristik tersendiri dalam melakukan pencarian.

#### A. Brute-force

Algoritma Brute-force merupakan algoritma paling sederhana dalam melakukan string matching. Algoritma ini berkerja dengan cara berikut:

1. Mula-mula, pattern di cocokkan pada awal teks
2. Dengan bergerak dari kiri ke kanan, setiap karakter pada pattern dibandingkan dengan karakter yang bersesuaian dengan teks, sampai semua karakter sama, atau ditemukan ketidakcocokan.
3. Apabila ditemukan ketidakcocokan dan teks belum habis, maka pattern digeser satu karakter ke kanan dan langkah 2 diulang.



Gambar 3 Contoh Bruteforce String Matching

```

BruteForceStringMatch(T[0...n-1],
P[0...m-1])
  for i ← 0 to n-m do
    j ← 0
    while j < m and P[j]
      = T[i+j] do
        j++
    if j = m then return
  i
  return -1

```

Algoritma *Brute-force* dalam string matching memiliki kelemahan. Pattern digeser satu karakter ke kanan setiap ditemukan ketidakcocokan. Ini akan memakan waktu yang lama hingga ditemukan ketidakcocokan atau teks habis.

### B. Knuth-Morris-Pratt

Algoritma KMP memelihara informasi yang digunakan saat melakukan pergeseran. Informasi ini digunakan untuk melakukan pergeseran yang lebih jauh, tidak seperti *brute-force* yang melakukan pergeseran per karakter.

Pergeseran dilakukan berdasarkan suffix kesamaan suffix dan prefix dalam pattern dan yang ditemukan di dalam teks.

```

m: 01234567890123456789012
S: ABC ABCDAB ABCDABCDABDE
W: ABCDABD
i: 0123456

m: 01234567890123456789012
S: ABC ABCDAB ABCDABCDABDE
W:      ABCDABD
i:      0123456

m: 01234567890123456789012
S: ABC ABCDAB ABCDABCDABDE
W:          ABCDABD
i:          0123456

m: 01234567890123456789012
S: ABC ABCDAB ABCDABCDABDE
W:              ABCDABD
i:              0123456

m: 01234567890123456789012
S: ABC ABCDAB ABCDABCDABDE
W:                  ABCDABD
i:                  0123456

```

Gambar 4 Pergeseran pada KMP

### KMP(P, T)

```

f ← compute failure function of
Pattern P
i ← 0
j ← 0
while i < length[T] do
    if j ← m-1 then
        return i- m+1 // we have a
match
    i ← i +1
    j ← j +1
else if j > 0
    j ← f(j -1)
else
    i ← i +1

```

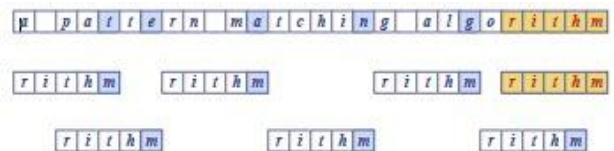
Seperti terlihat pada *pseudocode* diatas, algoritma KMP menggunakan fungsi pingiran yang mengindikasikan pergeseran terbesar yang mungkin dilakukan dalam melakukan pencarian. Fungsi pingiran ini membuat pencarian lebih efektif dengan menghindari pergeseran tidak berguna seperti halnya algoritma *Brute-force*.

### C. Boyer-Moore

Seperti Algoritma KMP, algoritma BM juga memanfaatkan informasi yang dimiliki untuk melakukan pergeseran. Hanya saja, Algoritma BM melakukan pencocokan string dari belakan pattern, tidak seperti algoritma KMP yang melakukannya dari depan.

Secara sistematis, algoritma BM bekerja sesuai langkah berikut:

1. Pattern dicocokkan dari awal teks
2. Dari kanan ke kiri, karakter per karakter dicocokkan
3. Algoritma menggeser pattern dengan memaksimalkan pergeseran *good-suffix* dan pergeseran *bad-character*, kemudian mengulangi langkah 2 sampai ujung teks atau kecocokan di temukan.



Gambar 5 Ilustrasi Algoritma Boyer-Moore

### Fungsi pingiran (P)

```

i ← 1
j ← 0
f(0) ← 0
while i < m do
    if P[j] = P[i]
        f(i) ← j +1
        i ← i +1
        j ← j + 1
    else if
        j ← f(j - 1)
    else
        f(i) ← 0
        i ← i +1

```

```

void preBmBc(char *x, int m, int
bmBc[]) {
    int i;

    for (i = 0; i < ASIZE; ++i)
        bmBc[i] = m;
    for (i = 0; i < m - 1; ++i)
        bmBc[x[i]] = m - i - 1;
}

```

```

void suffixes(char *x, int m, int
*suff) {
    int f, g, i;

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

```

```

void preBmGs(char *x, int m, int
bmGs[]) {
    int i, j, suff[XSIZE];

    suffixes(x, m, suff);

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

```

```

void BM(char *x, int m, char *y, int
n) {
    int i, j, bmGs[XSIZE], bmBc[ASIZE];

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

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

```

Pada *Boyer-Moore* dikenal istilah *bad-character*. *Bad-character* adalah karakter pada teks yang menjadikan perbandingan gagal. Pada pattern, dilihat karakter yang sama yang terdekat dengan karakter tempat terjadi kegagalan. Kemudian pattern digeser sehingga karakter pada pattern berada pada posisi yang sama dengan karakter pada teks.

### V. ALGORITMA SMITH-WATERMAN

Algoritma *Smith-Waterman* adalah algoritma pencarian yang dikembangkan T.F.Smith dan M.S.Waterman. algoritma ini merupakan algoritma pencarian yang mengimplementasikan program dinamis yang mengambil *alignment* dari panjang dari lokasi, dalam sebuah *sequence* dan menentukan apakah susunan yang optimal dapat ditemukan.

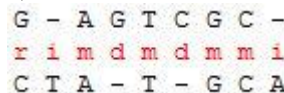
Pada algoritma ini, program dinamis bertugas mencari solusi terkecil dari masalah dan menggabungkannya menjadi solusi final dari keseluruhan masalah.

Algoritma SW membandingkan segmen dengan panjang berbeda, melihat apakah segmen memaksimalkan hasil pengukuran. Bentuk Rekursif dari algoritma ini adalah:

$$H_{ij} = \max\{H_{i-1,j-1} + s(a_i, b_j); H_{i,k,j} - W_k; H_{i,j-1} - W_1; 0\}$$

Gambar 6 Bentuk Rekursif dari SW

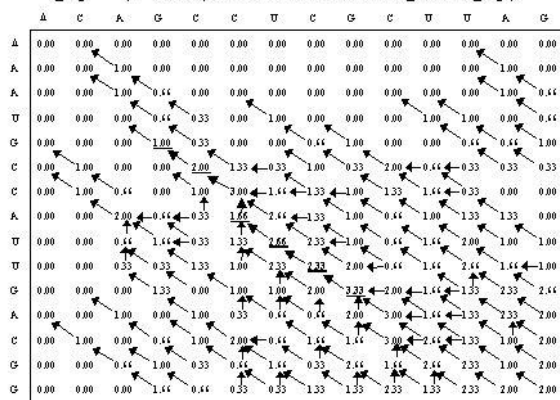
Misal kita memiliki *sequence* GAGTCGC dan CTATGCA. Pertama, bariskan *sequence* tersebut dan membandingkan jarak. Hitung r dan d (dalam kasus ini bernilai 5) sementara nilai 4 adalah nilai yang sama (dalam kasus ini 4)



Gambar 7 contoh penggunaan

Pada gambar diatas, r=replacement; m=match; i,d=insertion/deletion.

Smith Waterman Scoring Matrix  
(matches=1; mismatches=-.33; single gap = -(1 + k/3); where k is the length of gap)



Gambar 8 Scoring dalam Matrics

## VI. BASIC LOCAL ALIGNMENT SEARCH TOOL

Basic Local Alignment Search Tool atau BLAST, adalah sebuah kakas yang di gunakan dalam melakukan perbandingan *sequence* pada info biologis, seperti asam amino, protein atau nukleotida.

BLAST merupakan algoritma yang lebih cepat dibandingkan dengan algoritma Smith-Waterman, meskipun tidak menjamin alignment yang seoptimal Smith-Waterman. BLAST menggunakan format FASTA sebagai input. FASTA merupakan sebuah format berbasis teks sebagai representasi dari nukleotida atau peptida.

Keluaran dari BLAST dapat berupa HTML, plain-text, atau XML. Jka dijalankan langsung pada halaman web NCBI, keluaran default adalah HTML. Hasil juga ditampilkan berupa grafis jika dijalankan melalui NCBI, juga tabel berisi pengidentifikasi *sequence* terhadap data yang berkaitan sehingga hasil dari pencarian lebih mudah dibaca.

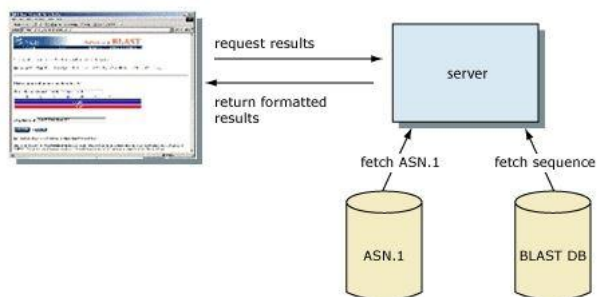
### A. Cara Kerja

Kebanyakan protein modular di alam, dengan domain fungsional seringkali berulang dalam protein yang sama begitu juga diantara protein yang berbeda dari spesies yang berbeda. Agoritma BLAST menggunakan informasi ini untuk mempercepat pencarian kesamaan *sequence*. mRNA yang ditemukan dapat di *align* dengan kepingan DNA, karena dibutuhkan dalam analisis sebuah genom.

Ketika sebuah query di submit, *sequence* dan beberapa info lainnya dimasukkan kedalam algoritma BLAST di dalam server. BLAST pertama kali dengan membuat sebuah tabel *lookup* dari semua “kata”(sebuah sub-*sequence* dimana protein biasanya terdiri dari tiga huruf) dan “kata terdekat”. Ketika kecocokan ditemukan, kecocokan ini digunakan untuk menginisiasi *gapped extensions* dari “kata” tersebut.

BLAST tidak mencari di dalam GenBank seara langug. Setiap entri dibagi dua, salah satu mengandung header dan yang lainnya berisi informasi *sequence*. Inilah data yang digunakan dalam algoritma.

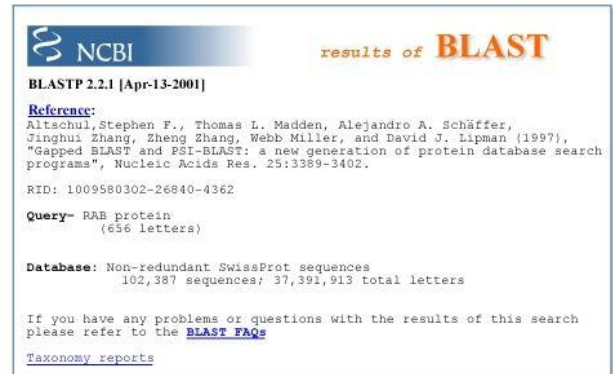
Setelah algoritma selesai melakukan pencarian terhadap semua “kata” yang mungkin dan memperluas secara maksimal, setiap *query-sequence* disusun kembali dan menuliskan informasinya ke sebuah SeqAlign.



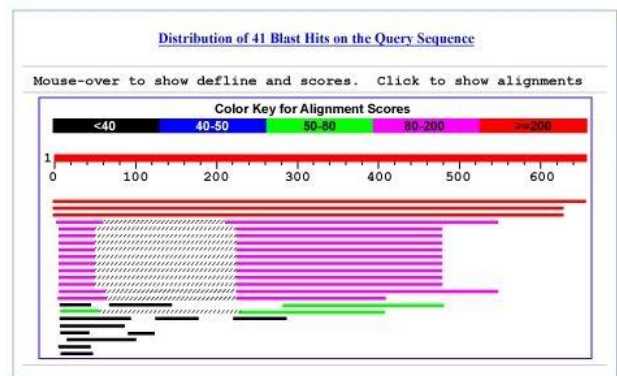
Gambar 9 Bagaimana hasil pada web dibentuk

### B. Keluaran

Keluaran dari BLAST dapat berupa informasi “tradisional”. Keluaran ini terdiri dari (1) header, yang berisi informasi dari *sequence*, dan grafis menyeluruh; (2) deskripsi dari *sequence* pada database yang sesuai dengan query; (3) alignment dari setiap database yang sesuai



Gambar 10 Header BLAST



Gambar 11 Hasil dalam grafik

Sequences producing significant alignments:	Score	E Value	
(a)	(b)	(c) (d)	
gi 116365 sp P26374 RAR2 HUMAN	Rab proteins geranylgeranyl...	1216	0.0
gi 21431807 sp P24386 RAR1 HUMAN	Rab proteins geranylgeranyl...	879	0.0
gi 585775 sp P37727 RAR1 RAT	Rab proteins geranylgeranyltra...	846	0.0
gi 13626886 sp Q61598 GDIC MOUSE	RAB GDP dissociation inhib...	127	5e-29
gi 7295661 sp F39958 GD11 YEAST	SECRETORY PATHWAY GDP DISSOC...	127	5e-29
gi 13626813 sp O97556 GDIB CANFA	Rab GDP dissociation inhib...	126	1e-28
gi 13638228 sp P50397 GDIB MOUSE	RAB GDP dissociation inhib...	125	3e-28
gi 1707888 sp P50398 GDIA RAT	RAB GDP dissociation inhibito...	124	7e-28
gi 121108 sp P21856 GDIA BOVIN	Rab GDP dissociation inhibiti...	124	7e-28
gi 21903424 sp P50396 GDIA MOUSE	Rab GDP dissociation inhib...	124	7e-28
gi 13626812 sp O97555 GDIA CANFA	RAB GDP dissociation inhib...	124	8e-28
gi 1707886 sp P31150 GDIA HUMAN	Rab GDP dissociation inhibi...	123	9e-28
gi 13638228 sp P50395 GDIB HUMAN	Rab GDP dissociation inhib...	122	2e-27
gi 1707891 sp P50399 GDIB RAT	RAB GDP DISSOCIATION INHIBITO...	121	5e-27
gi 1723461 sp Q10305 YD4C SCHPO	Putative secretory pathway ...	120	8e-27
gi 585776 sp F32864 RAEP YEAST	RAB proteins geranylgeranyl...	97	7e-20
gi 10720243 sp O93831 RAEP CANAL	RAB proteins geranylgeranyl...	74	9e-13
gi 2498411 sp Q49398 GLF MYCGE	UDP-galactopyranose mutase	35	0.63
gi 11135401 sp Q9X8Q2 STHA AZCVI	Soluble pyridine nucleotid...	34	1.0
gi 11135075 sp O05139 STHA PSEFL	Soluble pyridine nucleotid...	33	1.3
gi 11135195 sp P57112 STHA PSEAE	Soluble pyridine nucleotid...	33	1.8
gi 22257022 sp Q8TJ38 RLA0 PYRFU	Acidic ribosomal protein P...	33	2.1
gi 3915516 sp P94488 YNAJ BACSU	Hypothetical symporter ynaJ	32	3.4
gi 231788 sp P30599 CHS2 USTMA	CHITIN SYNTHASE 2 (CHITIN-UD...	32	3.7
gi 2498412 sp P75499 GLF MYCPN	UDP-galactopyranose mutase	32	4.2
gi 547891 sp P36225 MAP4 BOVIN	Microtubule-associated prote...	32	4.2
gi 586602 sp P37747 GLF ECOLI	UDP-galactopyranose mutase	32	4.6

Gambar 12 on-line description

