

ゲノムデータベースとプログラミング

法政大学 生命科学部 大島研郎

本日の講義資料

kiso2



本日の講義で使用するWebページへのリンクが載せてあります。

blast.pl

result.txt

本日の講義では, Active perl を使います.

◆ コマンドプロンプトを立ち上げてください



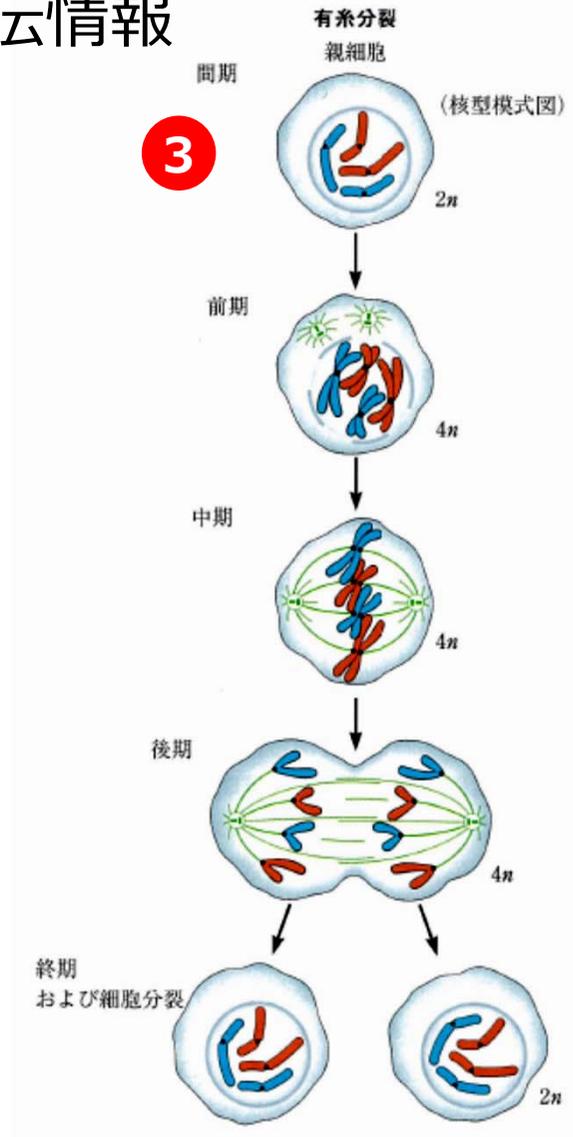
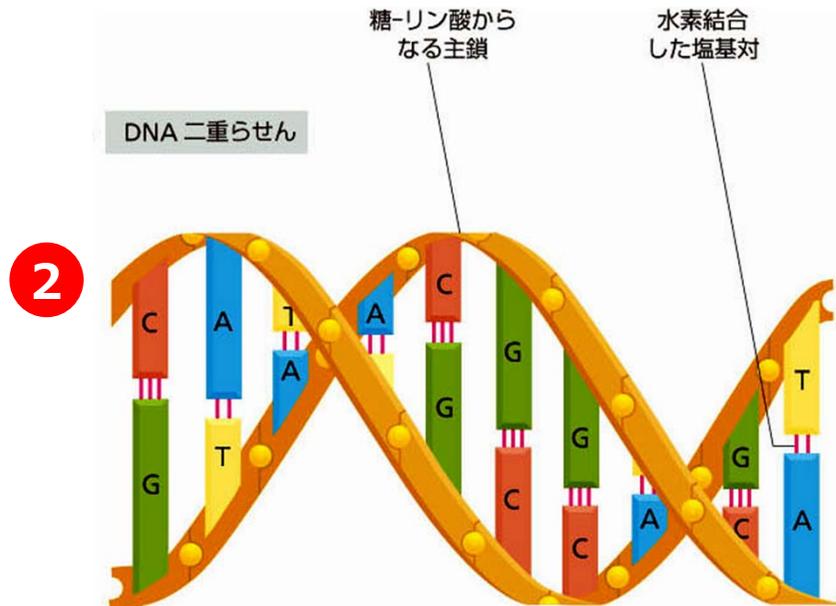
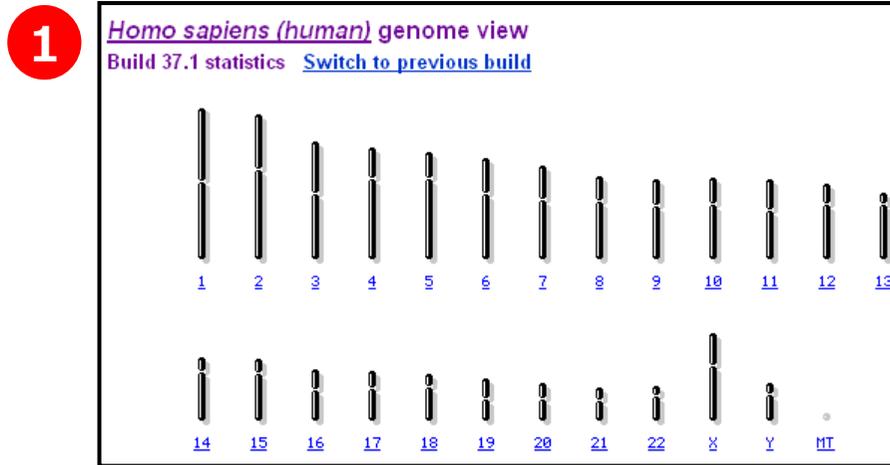
スタート → Windowsシステムツール → コマンドプロンプト

```
> perl -v
```

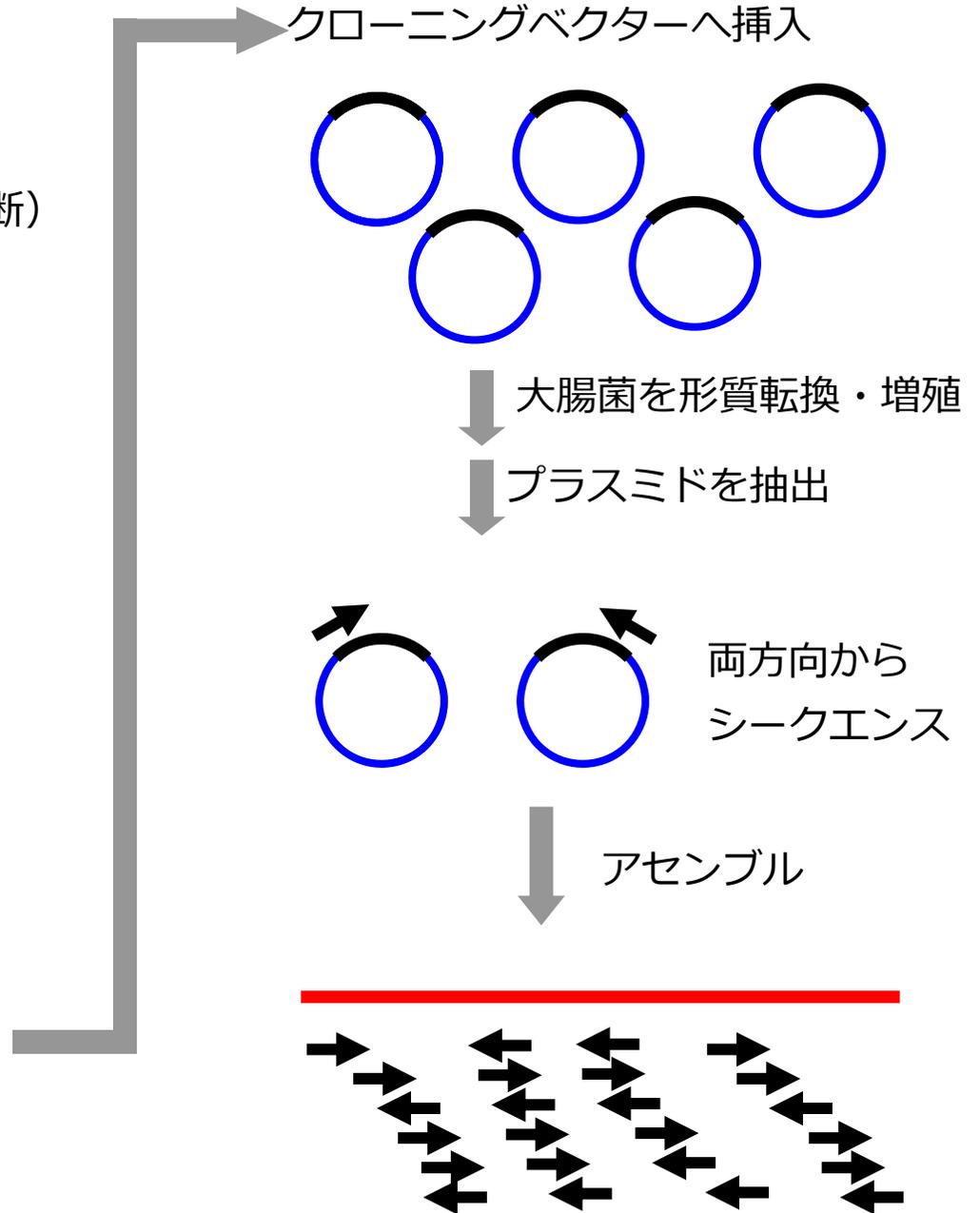
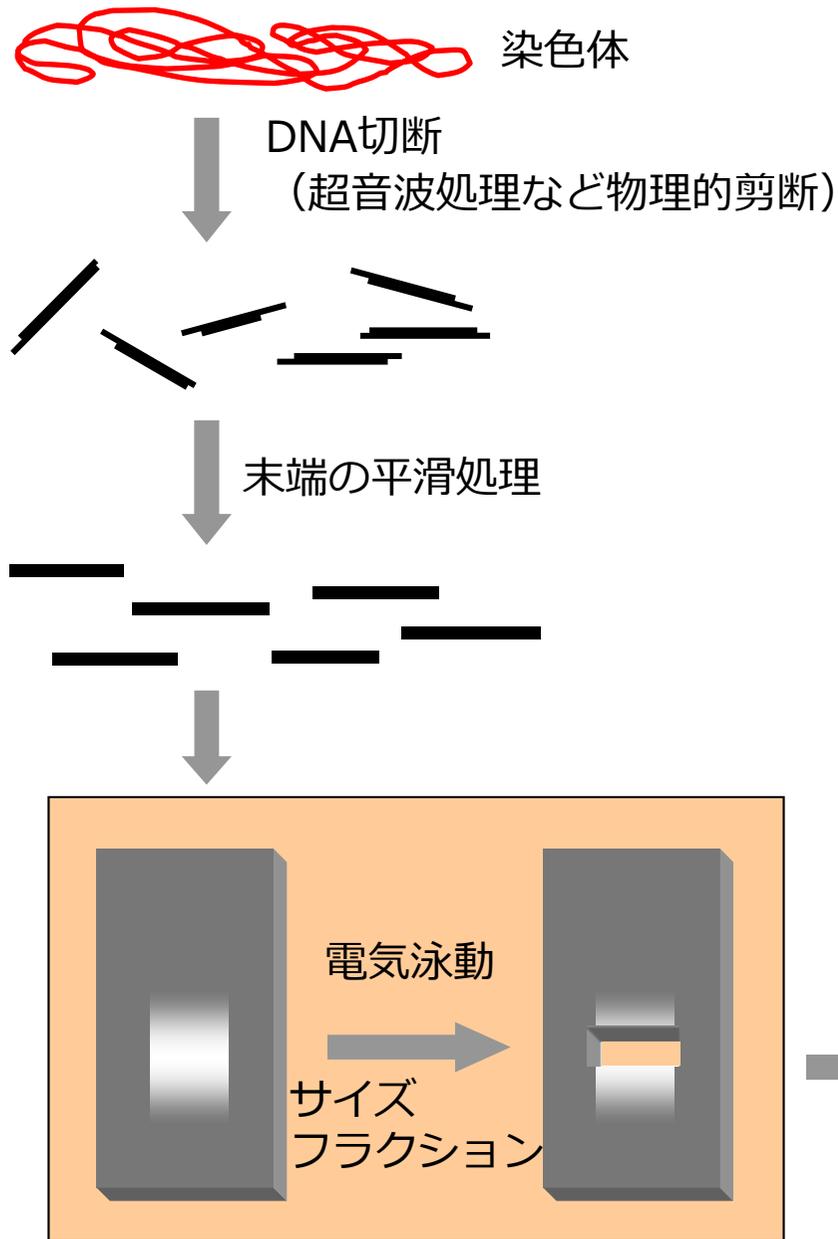
と入力して, エラーが出ないことを確認してください

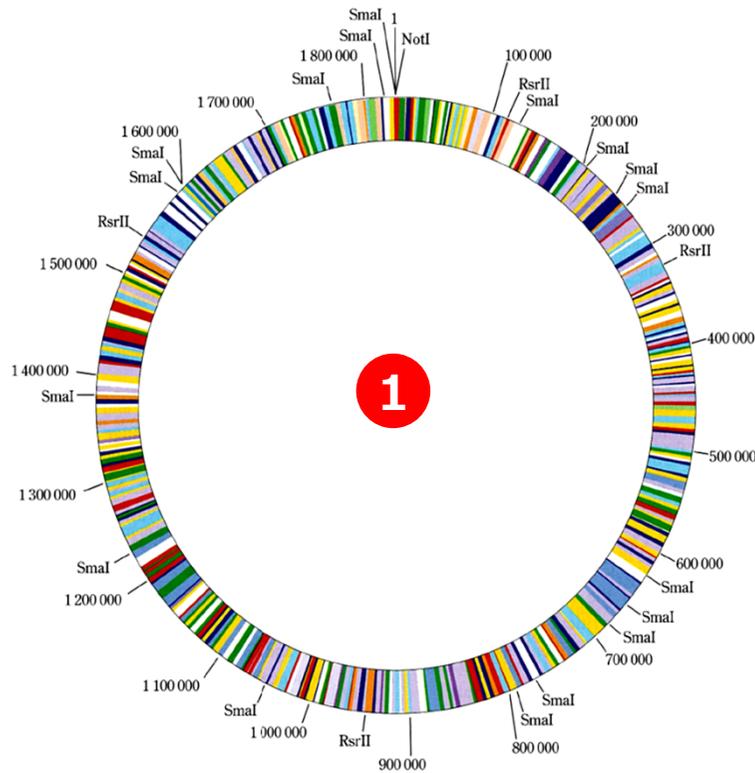
ゲノムとは gene (遺伝子) + -ome (総体)

ゲノム = ある生物のもつ全ての遺伝情報



ショットガン シーケンス法





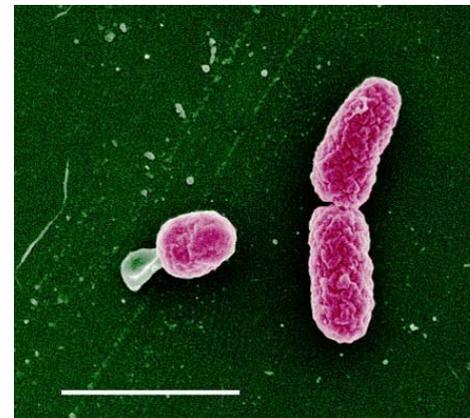
DNA は二本鎖なので、塩基対(base pair; bp)の数で分子の長さを表す。キロ塩基対(kilobase pair; kb)は 10^3 bp, メガ塩基対(megabase pair; Mb)は 10^6 bp, ギガ塩基対(gigabase pair; Gb)は 10^9 bp。まとめると,

1 kb = 1000 bp

2 1 Mb = 1000 kb = 1,000,000 bp

1 Gb = 1000 Mb = 1,000,000 kb = 1,000,000,000 bp

RNA 分子はたいてい一本鎖なので長さの単位に bp は使えず、ヌクレオチドの数で表す。



- 1995年, 生物として初めて*Haemophilus influenzae*の全ゲノムが解読された
- その後, 多くの生物の全ゲノムが解読され, 現在では1万種以上の生物のゲノム情報がデータベースに登録されている

▶ ヒトゲノムマップ を開く

<http://www.lif.kyoto-u.ac.jp/genomemap/>

Number. **09**
1億4000万bp
1076個

色別性染色体
染色体番号: Xp-A
決定要素: DR13C3
トリアニン有能タンパク質:
AIBP
決定要素: DR1B1
決定要素: DR1L4
A17決定要素:
アデニルキナーゼ1
細胞膜タンパク質:
スベクドリン
ABO血液型遺伝子
細胞膜タンパク質
タンパク質: Notch
グルタミル脱水素1

遺伝子名	通称名
ABO	ABO血液型遺伝子

◎赤血球に目印をつける酵素。
◎目印にはA型、B型の2種類があり、この組み合わせで血液型が決まる。
◎目印がつかない場合はO型になる。

同等の遺伝子(オノログ)を持つ生物

※2006年2月時点で公開されているデータベースに基づいて作成したものです。 CLOSE

Number. **07**
1億6300万bp
1378個

免疫系タンパク質:
インターロイキン5
細胞膜タンパク質:
シトクロームC
細胞膜タンパク質:
HDKA
細胞膜タンパク質:
HSP
コラーゲンIIα2
細胞膜タンパク質:
スフィンゴタン
アセチルコリン受容体:
ニコチン
細胞膜タンパク質:
FOXP2
細胞膜タンパク質:
シファン
タンパク質:
トリアニン

遺伝子名	通称名
FOXP2	発話と言語に関わる遺伝子

◎発話や言語に関わる脳の神経をつくるのに重要な役割を果たすタンパク質。
◎同:遺伝子モトとチンパンジーで比較すると、2塩基のみが異なることがわかっており、この差がヒトに優れた言語能力をもたらした可能性がある。

同等の遺伝子(オノログ)を持つ生物

※2006年2月時点で公開されているデータベースに基づいて作成したものです。 CLOSE

Number. **Y**
5100万bp
255個

性決定タンパク質:
SHOX
性決定遺伝子
性決定タンパク質:
DAZ
遺伝子砂漠

遺伝子名	通称名
SRY	性決定遺伝子

◎男性化に関わるタンパク質。
◎ヒトの体は元々女性型になっているが、このタンパク質が作用すると精果ができる。

Number. **12**
1億4200万bp
1268個

ヘルパーT細胞タンパク質: CD4
グルタミル脱水素2
乳糖分解酵素
コラーゲンIIα1
染色体タンパク質:
セパレーズ
細胞膜タンパク質:
HDXC
インスリン様増殖因子1
フェニルアラニン水酸化酵素
アルデヒド分解酵素2
一酸化窒素合成酵素1
神経細胞増殖因子:
タンパク質: Mサン
DNA複製酵素:
ポリメラーゼ

遺伝子名	通称名
ALDH2	アルデヒド分解酵素2

◎アルコールから生成される有毒なアセトアルデヒドを無毒な酢酸に変える酵素。
◎お酒に強い人は、この酵素のはたらきが弱い。

同等の遺伝子(オノログ)を持つ生物

Number. **X**
1億6300万bp
1141個

骨髄系タンパク質:
SHOX
インターロイキン5
細胞膜タンパク質:
HDKA
DNA複製酵素:
ポリメラーゼ
ポリストロフィーニ
細胞膜タンパク質:
シトクロームC
アンドロゲン受容体
インターロイキン5
細胞膜タンパク質:
BTK
細胞膜タンパク質:
タンパク質: CD44
細胞膜タンパク質:
OPN1LW
細胞膜タンパク質:
OPN1MW
細胞膜タンパク質

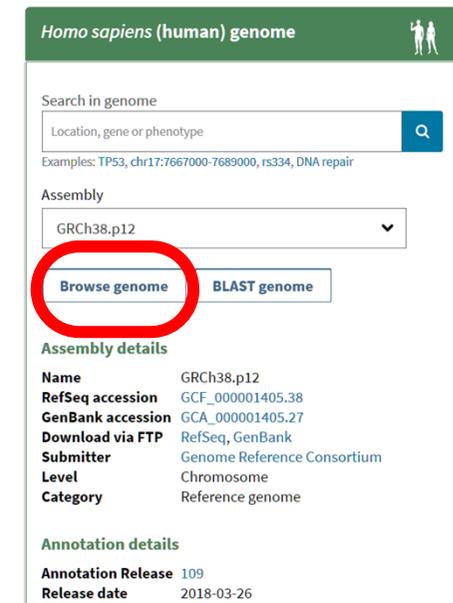
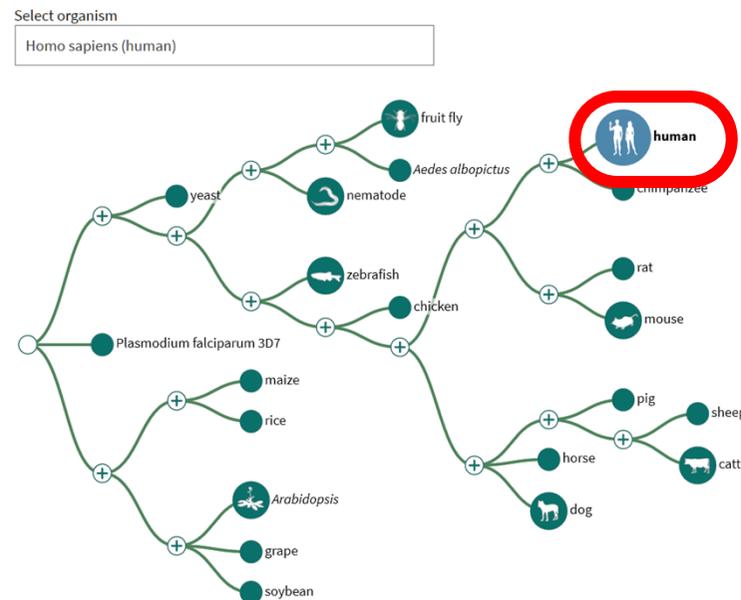
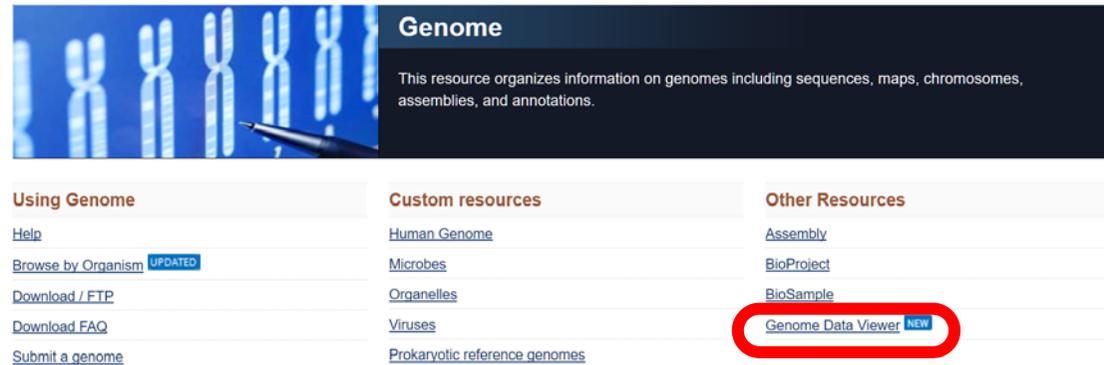
遺伝子名	通称名
OPN1LW	赤色識別遺伝子
OPN1MW	緑色識別遺伝子

◎OPN1LWは赤色を識別する際に、OPN1MWは緑色を識別する際に機能するタンパク質。
◎そのいずれかのタンパク質が変異すると、赤と緑が判別しにくい色覚を持つことになる。

遺伝子砂漠
<div style="text-align: center;"> </div> <p>◎非遺伝子領域が延々と続く不毛な地帯。 ◎このような領域はゲノム上の様々な場所に存在している。</p> <p style="text-align: center; font-size: 2em;">TTCCAの反復配列</p>

ゲノムブラウザを使ってヒトゲノムを見てみよう

NCBIトップページ右のリンクから「Genome」→「Genome Data Viewer」



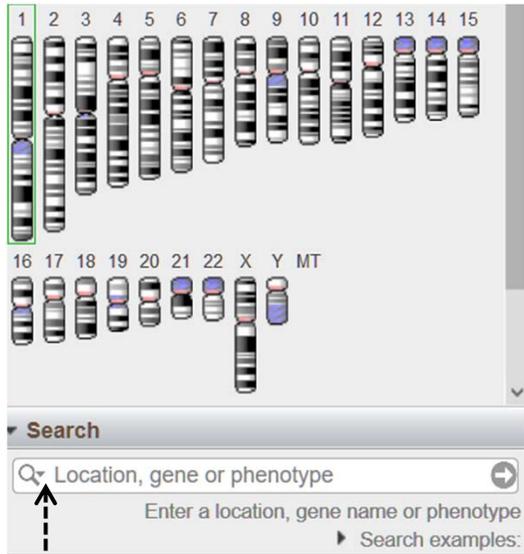
ゲノム解読された真核生物の
系統図が表示される

「human」



「Browse genome」
をクリック

クリックした染色体のゲノムマップが表示される



遺伝子マップ

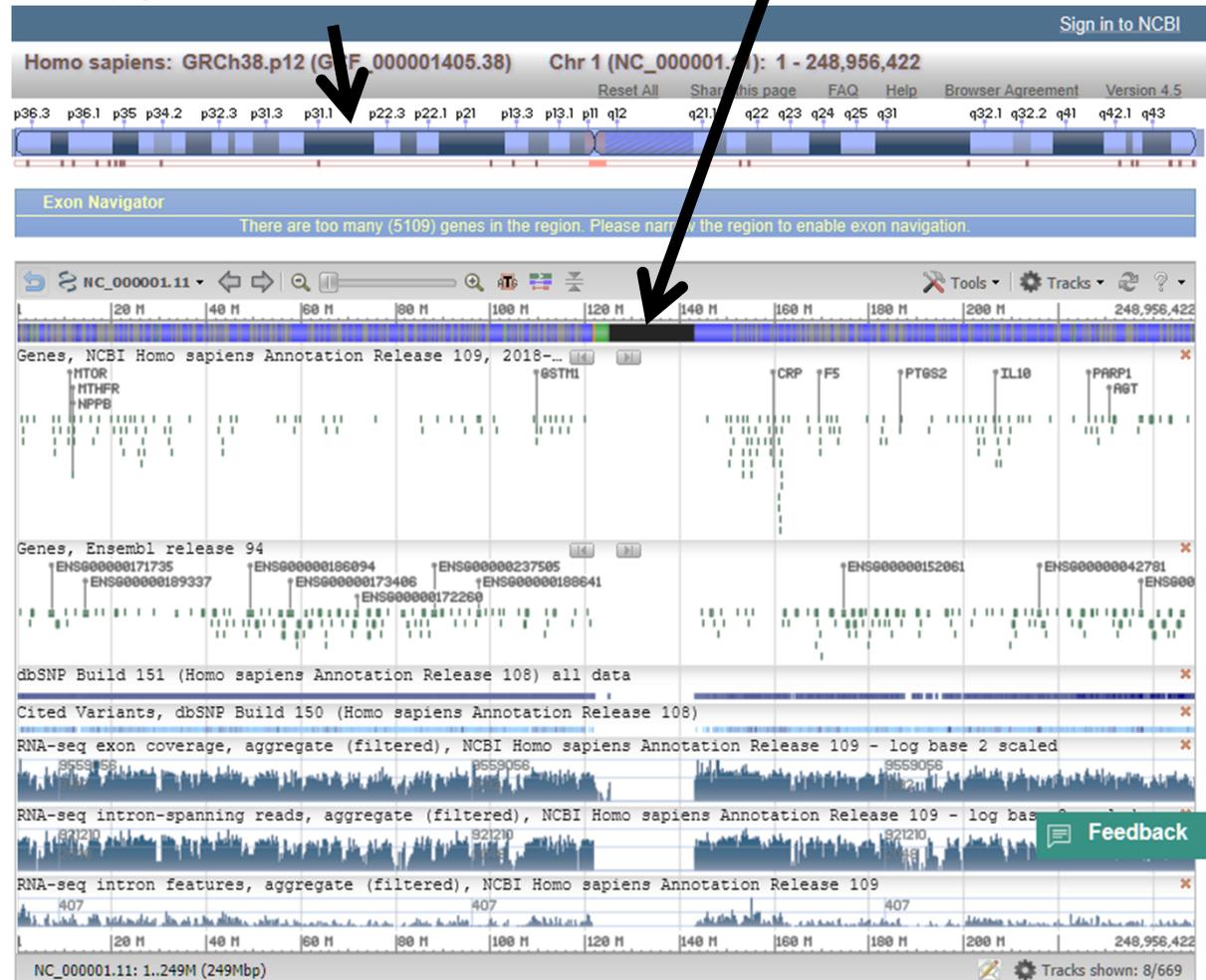
太線 : エキソン

細線 : イントロン

遺伝子発現量

1 青い領域にコードされる遺伝子が下に表示される

2 黒い領域は塩基配列が決まっていない

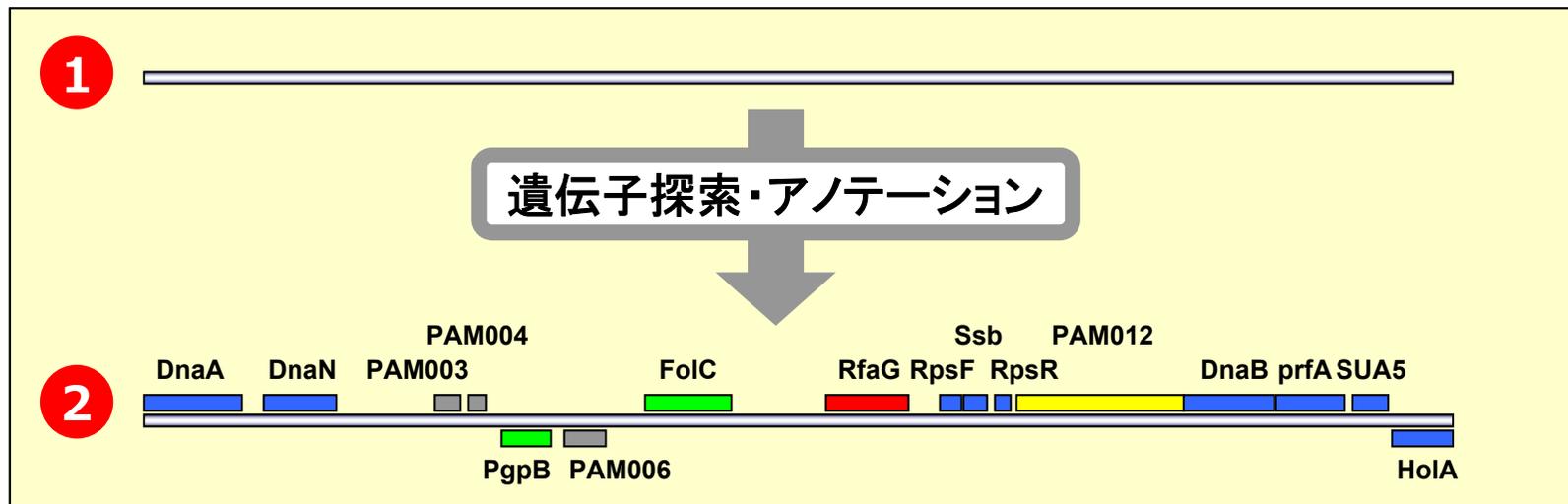


3 「ABO」と入力してみる

セントロメアには遺伝子がコードされていない

遺伝子探索

```
agttatTTTTTTTTcattttataattaaaagccaatttaaaaaaggagtattaattatgccat
atattgaaagtatTTTtagcgcgcgaagtgctagattccagaggaaatcctacagtagaagtaga
agttatacagaatcaggagcgtttggaagagctattgttccttcaggagcttctaccggacaa
tacgaagcagttgaattaagggatggggatgcccaaagatTTTTtaggtaaaggcgttttgcaag
ctgttaaaaatgttattgaagttattcaaccagaattagaaggttattctgtccttagaacaac
tttaattgataaattattaattaaacttgacggaactcctaacaatctaatttaggagctaac
gctatTTTtaggtgtttctttggcttgctaaagctgcagctaactacttaaatcttgagtttt
atcaatatgtaggagggcgttttacctaacaatgccagttcctatgatgaatggttatcaacgg
tgagc.....
```

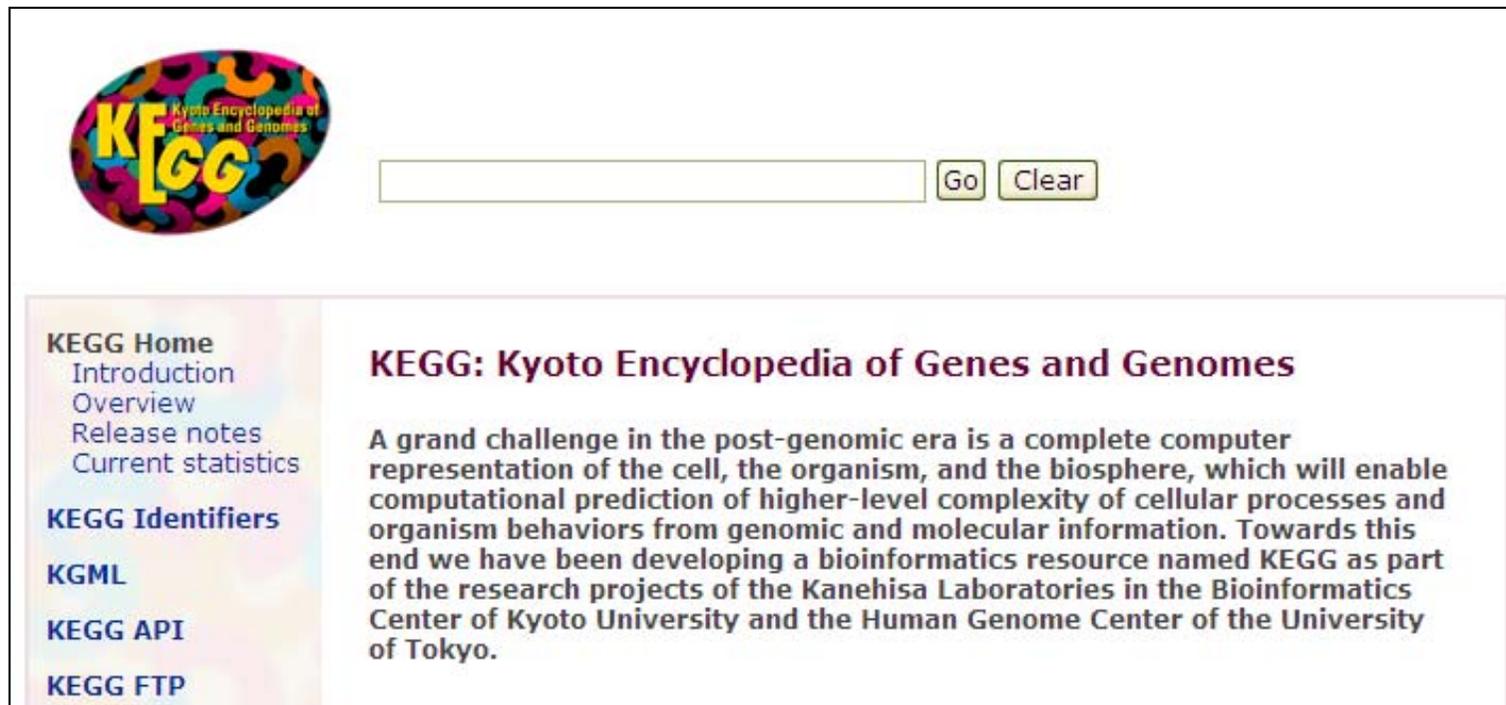


ゲノム配列から生命活動に関わる機能や分子進化に関する考察などを行うためには、タンパク質をコードしている遺伝子領域を同定することが重要となる。

代謝パスウェイデータベース

KEGG

<http://www.genome.jp/kegg/>



生命システム情報統合データベース。完全にゲノムが決まった生物種（一部、ドラフト配列も含む）の代謝系や一部の制御系（シグナル伝達や細胞周期など）をまとめ、そこから様々な物質データベースや酵素データベースを参照することができる。

次世代シーケンサー

Roche Diagnostics社

1

Genome Sequencer FLX System (454)

2005年発売



ライフテクノロジー社
Ion PGM



Applied Biosystems社

SOLiD 3

2007年発売



Solexa / illumina社

Genome Analyzer I/x

2005発売



イルミナ株式会社

MiSeq 2



PacBio RS II

次世代シーケンサーの比較

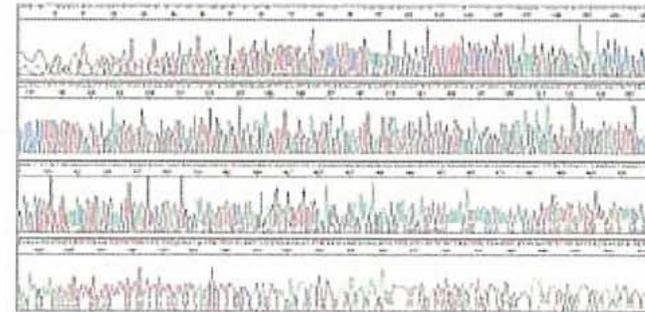
	Ion Protonシステム	Ion PGMシステム		MiSeq	HiSeq 2000/2500 (SBS v3試薬使用)	PacBio RS II
		Ion 318 chip				
1リード長	~200 base		150/250/300 base	100 base	約10,000 base	
リード数	約5,000万リード (1ランあたり)	約400万リード (1ランあたり)	約3,000万リード (1ランあたり) ※ペアエンド解析	1 約3億リード (1レーンあたり) ※ペアエンド解析	約5万リード (1セルあたり)	
データ量 (リード長 200 base の場合)	約7.5 Gb (平均150 bpの amplicon、 1チップあたり)	約800 Mb (1チップあたり)	約3~9 Gb (1ランあたり) ※ペアエンド解析	2 約30 Gb (1レーンあたり) ※ペアエンド解析	約500 Mb (1セルあたり)	
解析手法	Ion semiconductor sequencing法		Sequencing by Synthesis法	Sequencing by Synthesis法	SMRT(Silgle Molecule Real-Time) sequencing法	
アプリケー ション例	・癌遺伝子などの変異 解析 (409遺伝子をター ゲットとしたCancer Panelなど)	・癌遺伝子などの変異 解析 (50遺伝子をターゲッ トとしたCancer Panel など)	・微生物の新規ドラフ ト配列決定 ・癌遺伝子などの変異 解析 ・PCR産物のディープ シーケンス	・ゲノム変異解析 ・ChIP解析 ・small RNA解析 ・mRNA解析 ・cDNA配列解析	・ゲノムドラフト解析 ・cDNA配列解析	

イルミナ株式会社

次世代シーケンサー: Genome Analyzer

1 従来型キャピラリーシーケンサー

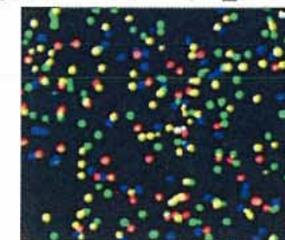
- 酵素反応+電気泳導+塩基読取 (384x600塩基)
- コスト、時間がかかる
- 例)「ヒトゲノムプロジェクト」
約13年、3000億円かかった



MiSeq

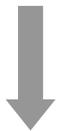
2 次世代シーケンサー

- 酵素反応+電気泳導+塩基読取 (100,000,000x50塩基)
- これまでの技術と比べて、「100分の1のコストで100倍のデータ」
- 例)現在ヒトゲノム1人読むのに 数週間、数千万円
→ 1週間 数百万円 ...



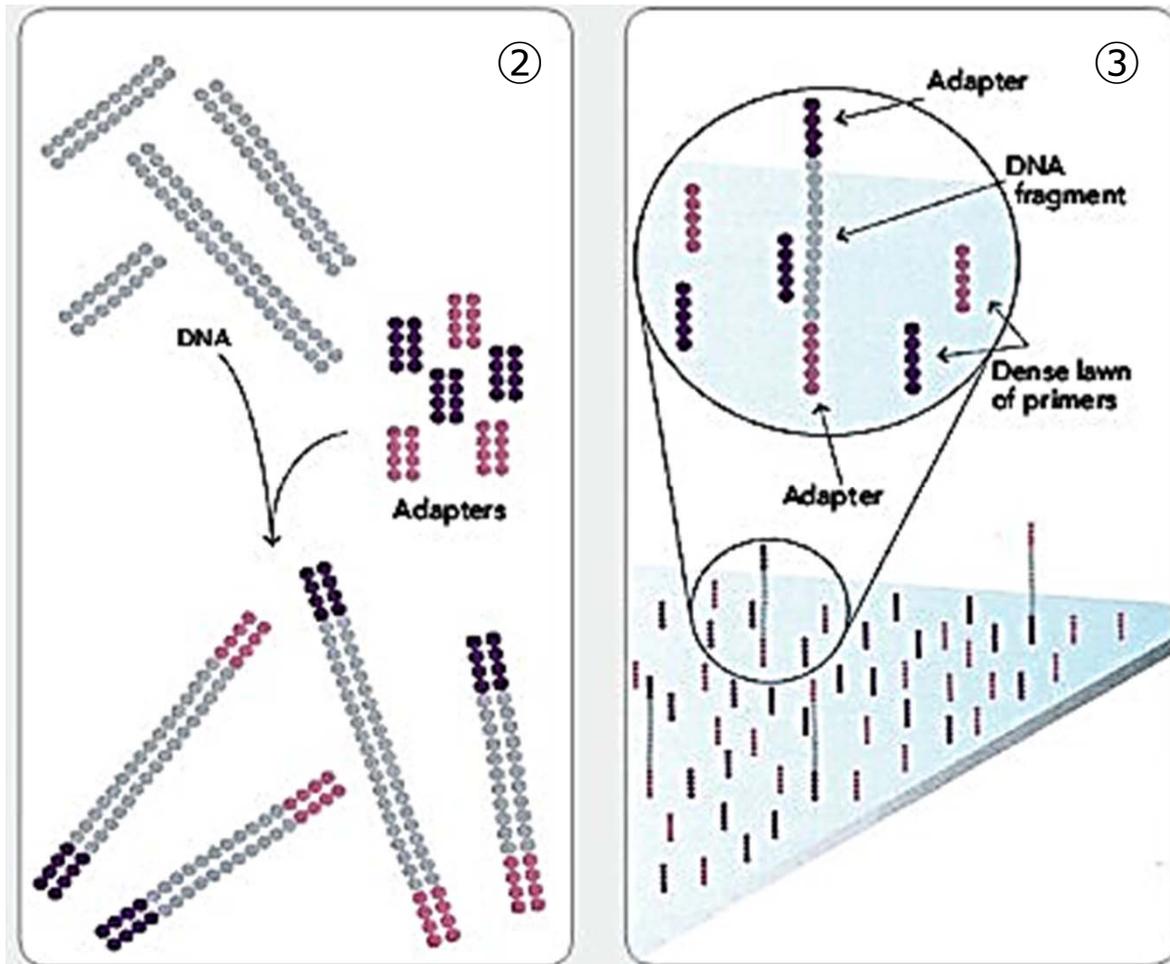
illumina®

ゲノムDNA



① 断片化

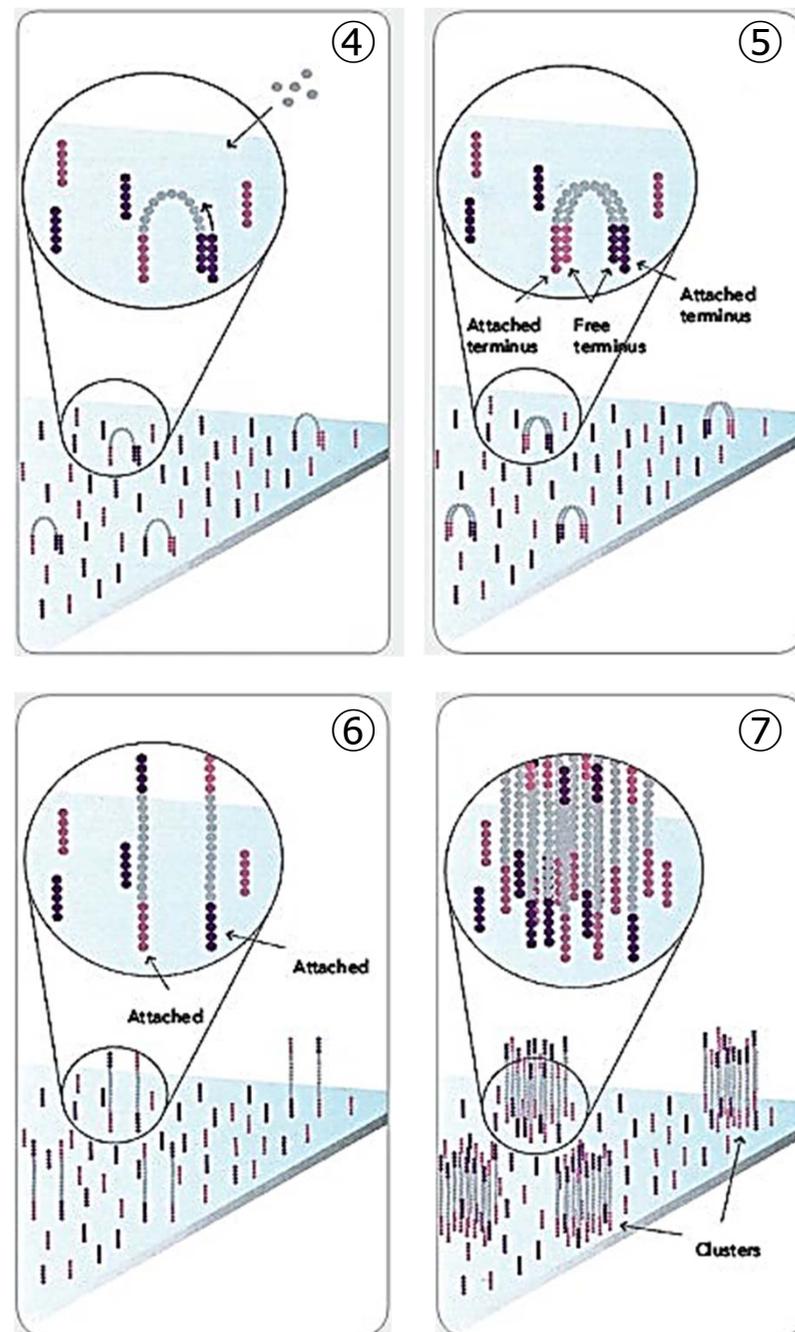
- ゲノムDNAを抽出し、断片化する
- DNA断片の両端に2種類の**アダプター**を連結させる



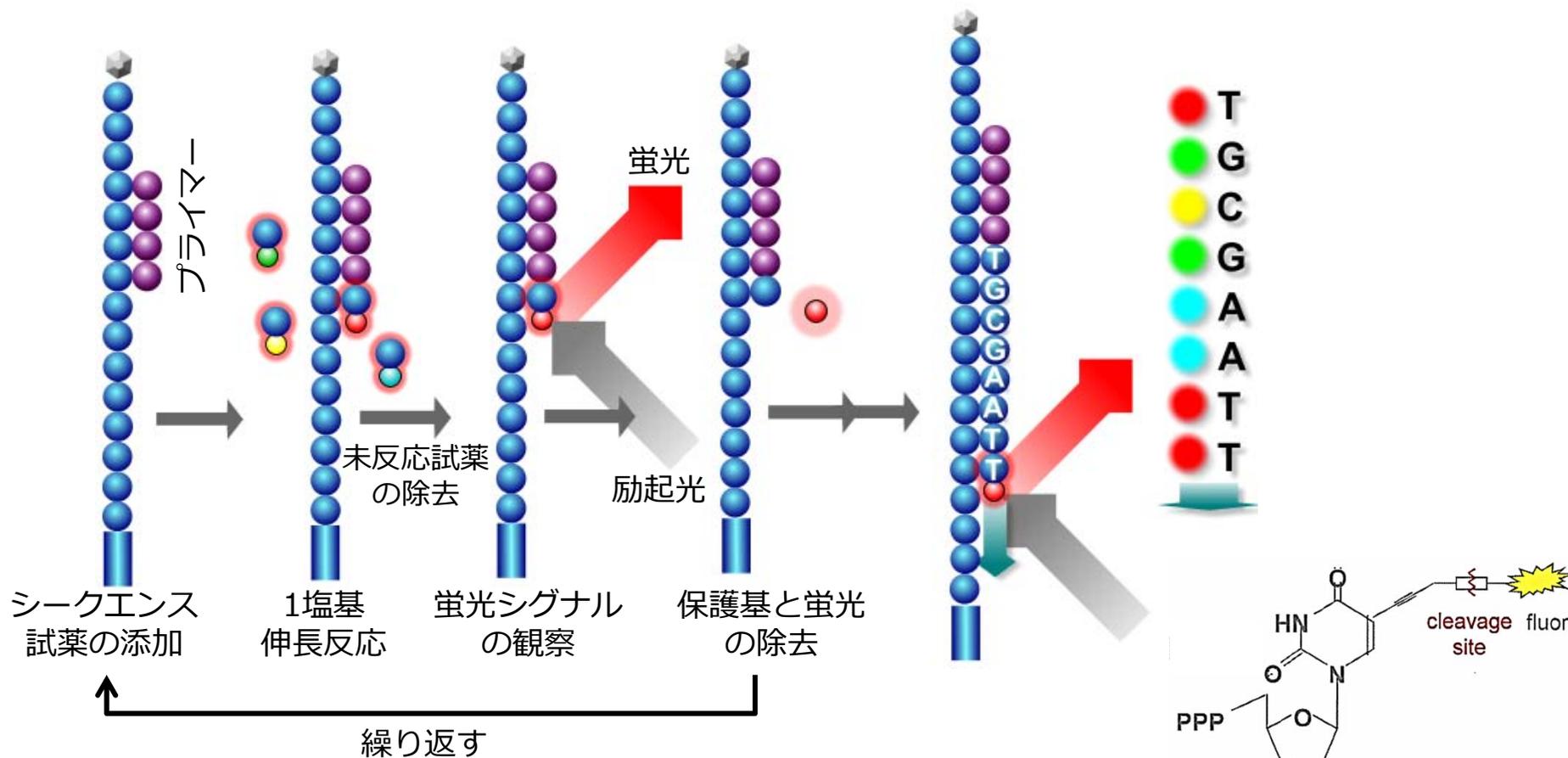
- 1本鎖にして、5'末端を**フローセル**上に固定する
- フローセル上には、あらかじめアダプターと相補的に結合するプライマーが高密度に配置されている

ブリッジPCR

- 固定された1本鎖DNAは、もう一方のアダプターの側でプライマーと結合する (橋がかかったような構造になる ④)
- DNAポリメラーゼによる伸長反応を行う ⑤
- 変性させると、フローセル上には根元がアダプター配列の1本鎖DNAが2本できあがる ⑥
- この反応を繰り返すことで、狭い面積の中でDNAを増幅することができる
 - フローセル上に多数のDNAの「束」ができる ⑦
- これらを鋳型として、配列解析を行う



Sequencing-by-synthesisによる塩基配列決定

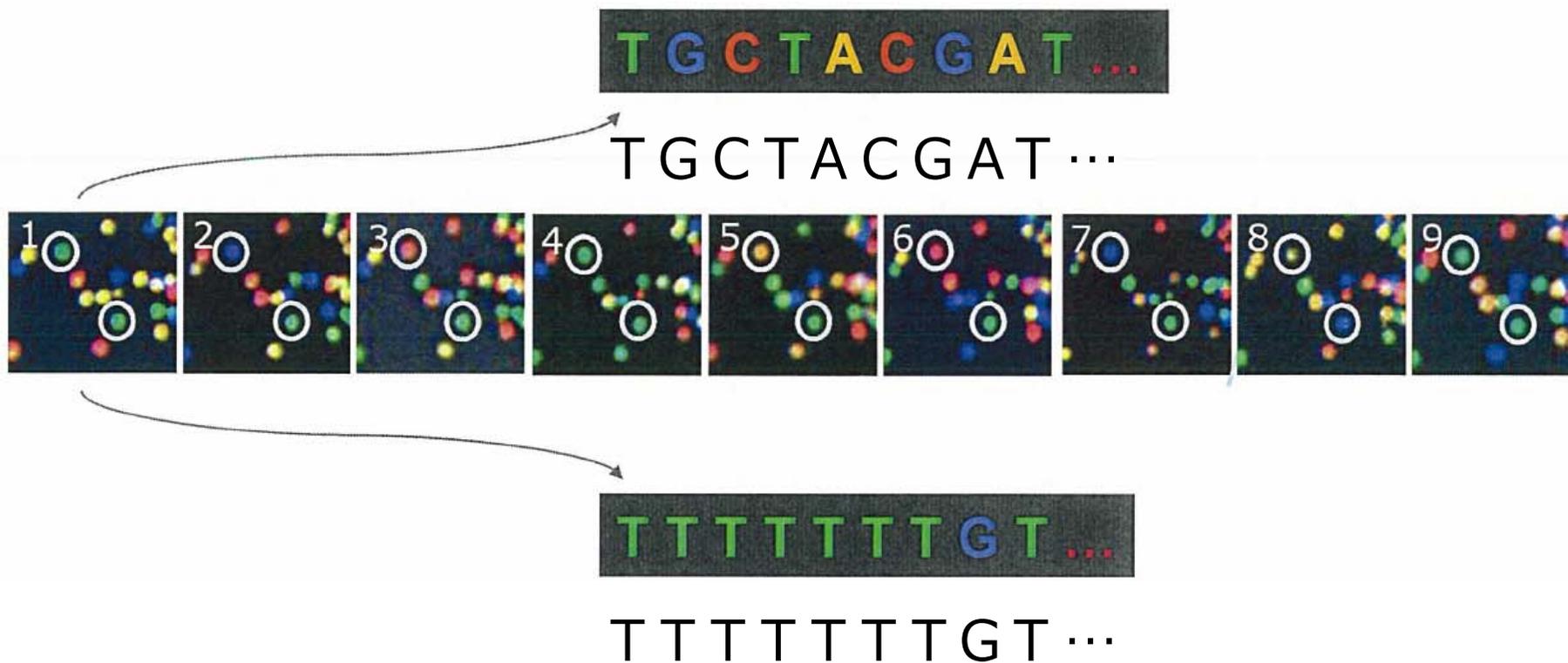
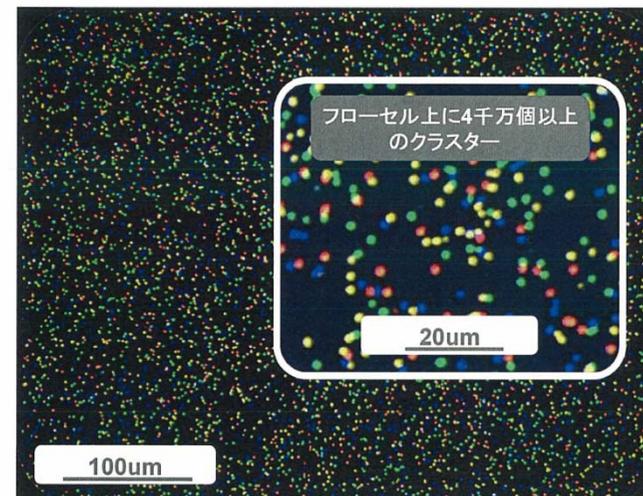


- 蛍光標識したdNTPの取り込みを**蛍光顕微鏡によって解析する**
- このdNTPは3'末端がブロックされており、1回の伸長反応で1塩基しか伸ばせない
- そのため、1塩基ごとにどのdNTPが取り込まれたかを観察し、蛍光物質とブロックを外して次の伸長反応を行うというステップで、解析を進めていく

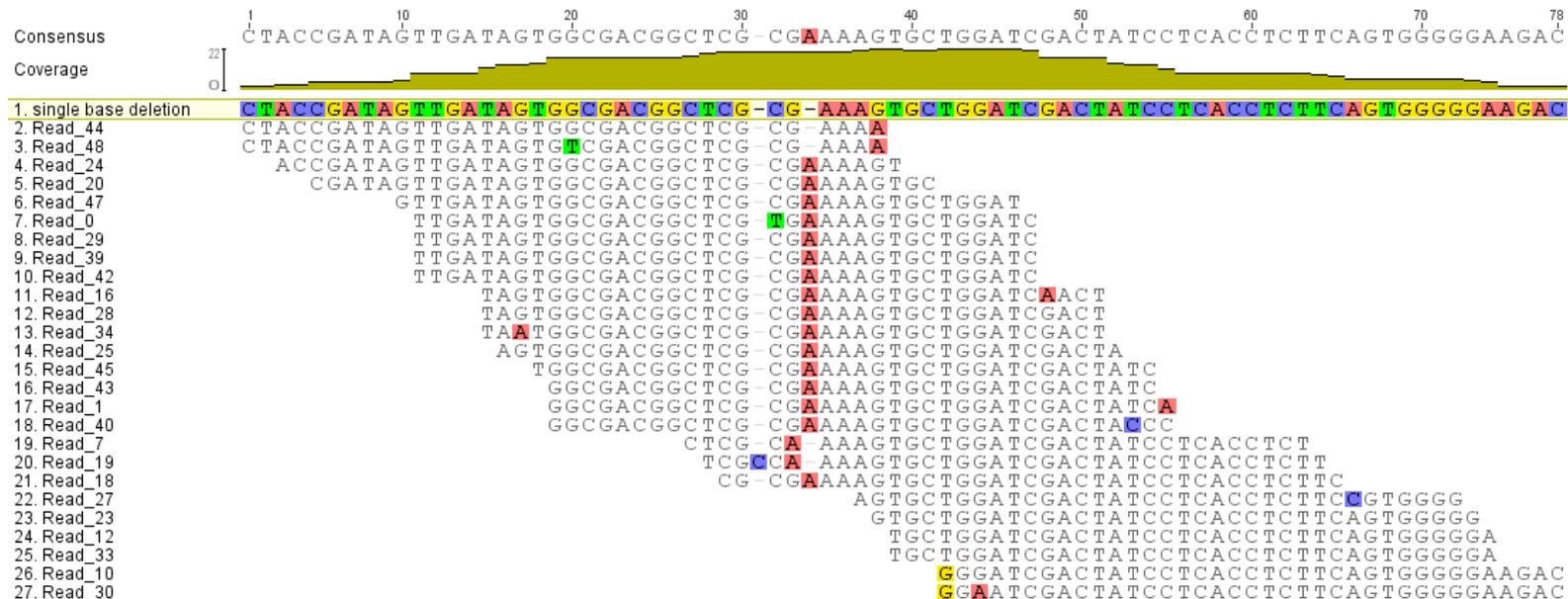
画像蛍光シグナルから塩基への返還

- **1塩基 伸長すること**に蛍光イメージを取得する
- それぞれのDNAの「束」の蛍光色の変化を調べることで、塩基配列を決定する
- 数千万～数億本の塩基配列が得られる

画像イメージの取り込み



- 一つ一つの断片の塩基配列が短いと、アセンブルするのが困難
- 次世代シーケンサーで読み取ることができる塩基配列長は短いので、既に全塩基が解読されているゲノム配列（リファレンス配列）を利用したリシークエンスや、リファレンス配列へのマッピングなどに用いられることが多い



マッピング : Bowtie, Bowtie2, BWA など
 アセンブル : Velvet, EDENA, Phrap など
 ビューア : Tablet, IGV など

有償ではあるが、CLC Genomics Workbenchなどの解析ソフトも良く使われる

Pectobacterium carotovorum
ssp. *carotovorum*

1



Pectobacterium carotovorum
PR1 strain (病原性強い)

2



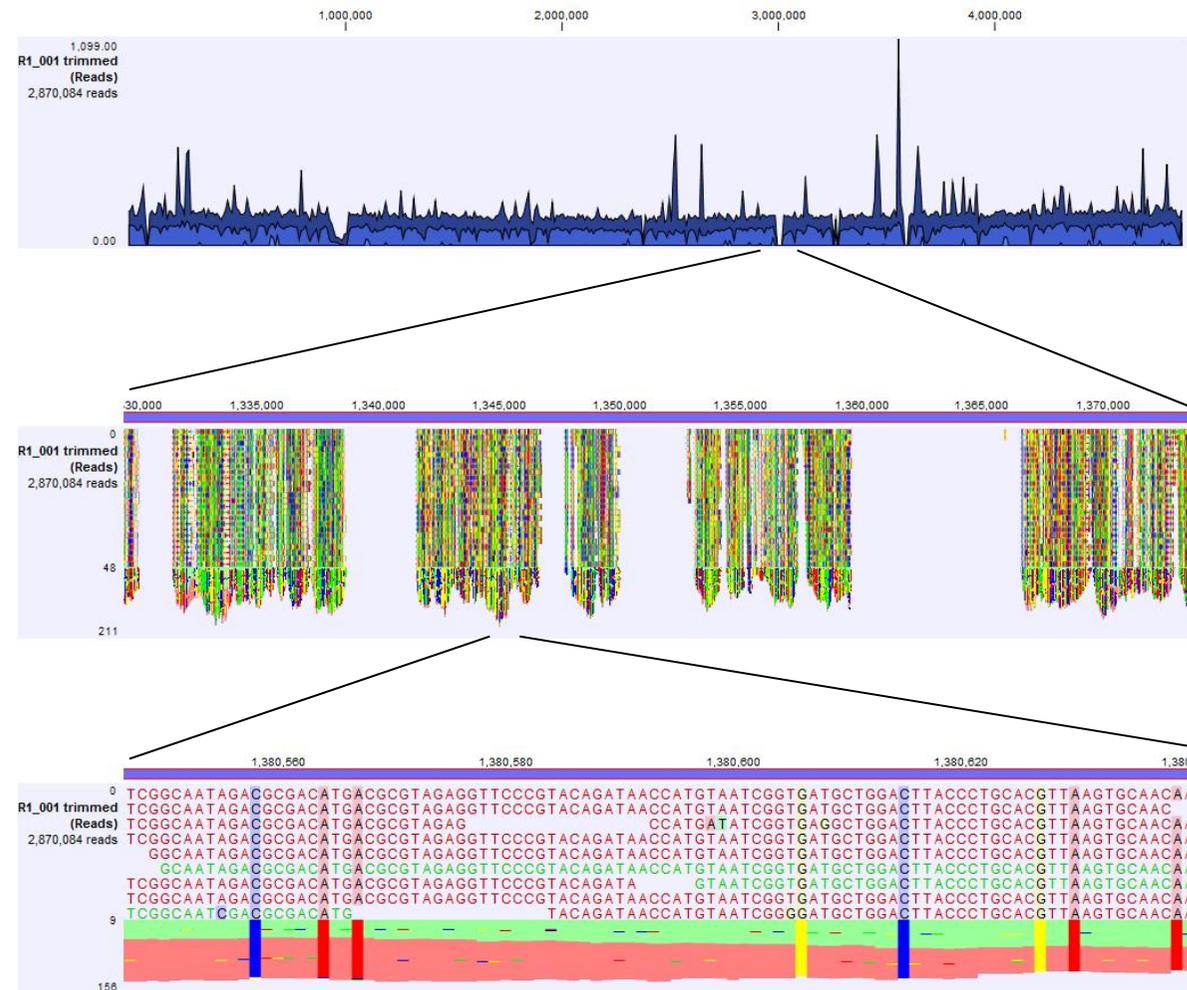
P. carotovorum PR1株
のゲノムを抽出



MiSeqを用いてシーケンス
(約300万リード)



P. carotovorum ssp.
carotovorum
のゲノム (リファレンス配列)
に対してマッピング



- 遺伝子の有無
- ゲノム構造の比較
- SNPの検出

等の比較ゲノム解析ができる

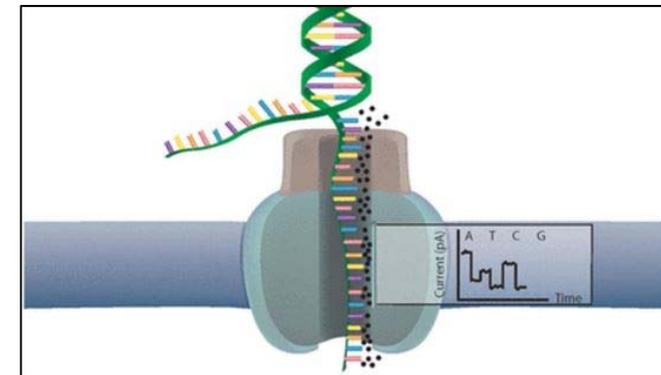
第3世代シーケンサー

Pacbio RS II DNA Sequencing System



- DNA **1分子**を鋳型としてDNAポリメラーゼによるDNA合成を行う
- 1分子レベルでリアルタイムに塩基を読み取る
- 長いリード(平均10,000bp)が出力される

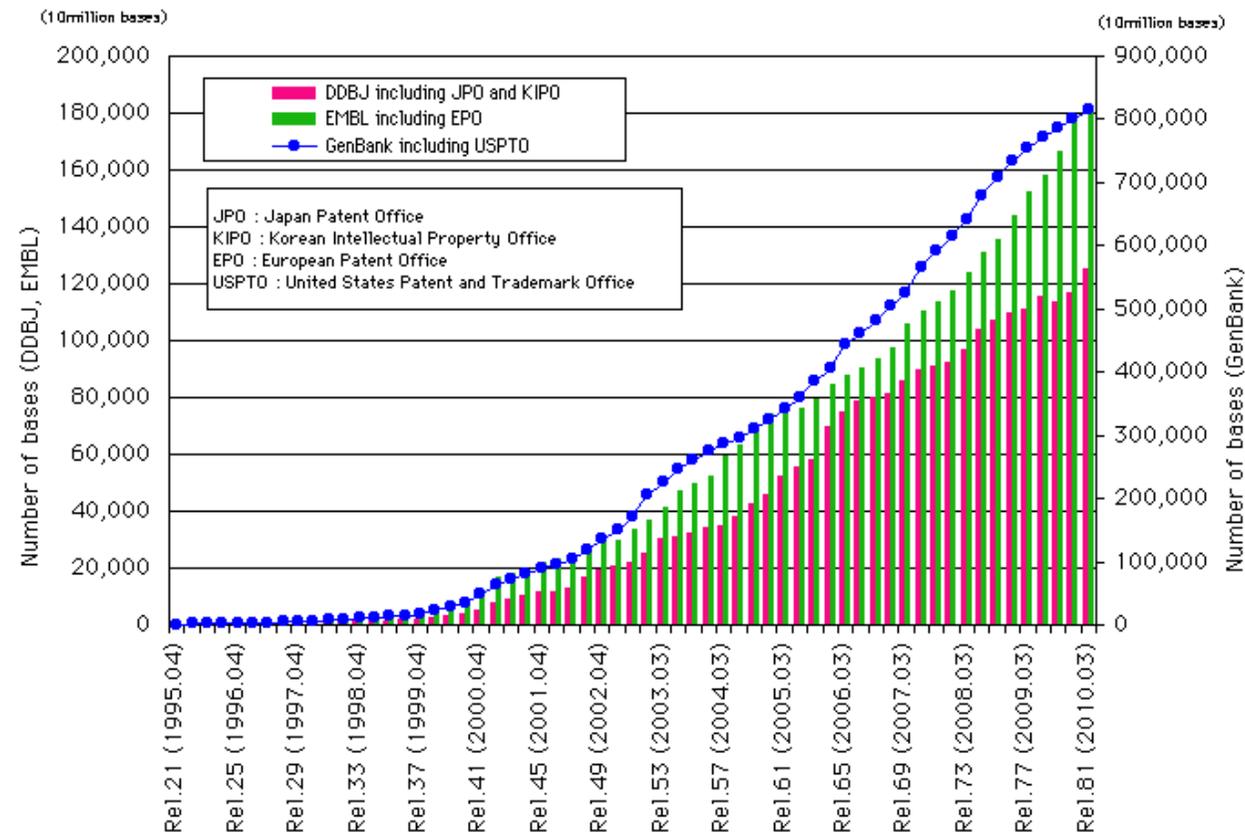
MinION



- USBメモリー用のシーケンサー
- DNAポリメラーゼを用いて1本鎖DNAに解きほぐす
- **ナノポア**を通過させる → 電流の変化を検知して配列を決定する

塩基データ登録数の推移

- シークエンス技術の進歩によって、塩基配列決定の速度はますます加速している



- 遺伝子の検出, アノテーション, 機能予測, 進化系統解析, 比較解析などを効率よく行い, 大量のシーケンスデータを有効に活用することが重要

ゲノムにコードされる遺伝子を網羅的に使用してホモロジー検索を行ったり，比較ゲノム解析を行いたい



大量のデータを処理するためのプログラミング技術が必要

バイオインフォマティクス分野では，Perl，C++，Java，Pythonなどが良く使われていますが，今日はPerlを用いて実習を行います



Perlの特徴

- テキスト処理が得意
- 歴史が長いのでライブラリーが豊富
- 掲示板やショッピングカートなど、CGIという仕組みの多くがPerlで書かれている
- LinuxやMacOSに標準でインストールされているほか、Windowsにもインストール可能

プログラミング言語の人気ランキング (2016)

- 【第1位】 Java
- 【第2位】 C言語
- 【第3位】 C++
- 【第4位】 Python
- 【第5位】 C#
- 【第6位】 Visual Basic
- 【第7位】 JavaScript
- 【第8位】 PHP
- 【第9位】 Perl
- 【第10位】 アセンブリ言語

perlを用いたデータ処理

- 大量のQueryに対して**BLAST検索**を行うと、結果が羅列した形で出力されます
- この中から、**必要な情報だけを取り出す**ためのプログラムをperlで組んでみましょう
- Queryのアクセッション番号と、検索の結果ヒットしたタンパク質のアクセッション番号とのリストを作成し、下に示したように検索結果を整理したいと思います

Query	1	2	3	4	5
gi 49176138 ref NP_416237.3	ref NP_009965.1	ref NP_010402.1	ref NP_012405.1	ref NP_012934.1	ref NP_015093.1
gi 16132212 ref NP_418812.1	ref NP_014926.1	ref NP_012380.1	ref NP_012969.1	ref NP_012770.1	ref NP_014585.1
gi 16131851 ref NP_418449.1	ref NP_009755.1	ref NP_011646.1	ref NP_013146.1	ref NP_013847.1	ref NP_013523.1
gi 16131757 ref NP_418354.1	ref NP_010335.1	ref NP_012586.1			
gi 16131754 ref NP_418351.1	ref NP_011756.1	ref NP_013932.1	ref NP_010104.1	ref NP_011671.1	
gi 16131018 ref NP_417595.1	ref NP_009362.1	ref NP_014992.1	ref NP_014792.1	ref NP_014227.1	ref NP_011015.1
gi 16130827 ref NP_417401.1	ref NP_009938.1	ref NP_011705.1	ref NP_011575.1	ref NP_011569.1	ref NP_012819.1
gi 16130826 ref NP_417400.1	ref NP_012863.1	ref NP_013282.1	ref NP_015308.1	ref NP_012835.1	ref NP_010022.1
gi 16130686 ref NP_417259.1	ref NP_011770.1	ref NP_012044.1	ref NP_014056.1	ref NP_015042.1	ref NP_015038.1
gi 16130106 ref NP_416673.1	ref NP_009965.1	ref NP_014276.1	ref NP_012639.1	ref NP_013060.1	ref NP_013066.1

BLASTP 2.2.5 [Nov-16-2002]

Reference: Altschul, Stephen F., Thomas L. Madden, Alejandro A. Schaffer, Jinghui Zhang, Zheng Zhang, Webb Miller, and David J. Lipman (1997), "Gapped BLAST and PSI-BLAST: a new generation of protein database search programs", Nucleic Acids Res. 25:3389-3402.

Query= gi|16131851|ref|NP_418449.1| glucosephosphate isomerase [Escherichia coli K12] (549 letters)

Database: yeast.aa
6298 sequences; 2,974,038 total letters

Sequences producing significant alignments:			Score	E
			(bits)	Value
ref NP_009755.1	Glucose-6-phosphate isomerase; Pgilp		641	0.0
ref NP_011646.1	Ygr130cp		30	0.98
ref NP_013146.1	spindle pole body component; Stu2p		29	1.7
ref NP_013847.1	(putative) involved in cell wall biogenesis; Ec...		28	3.7
ref NP_013523.1	Ylr419wp		28	3.7

>ref|NP_009755.1| Glucose-6-phosphate isomerase; Pgilp
Length = 554

Score = 641 bits (1654), Expect = 0.0
Identities = 326/549 (59%), Positives = 401/549 (73%), Gaps = 16/549 (2%)

```

Query: 7  TQTAAWQALQKHFDEM-KDVTIADLFAKDGDRFSKFSATFDD----QMLVDYSKNRITEE 61
          T+  AW  LQK ++   K +++   F KD  RF K + TF +   ++L DYSKN + +E
Sbjct: 13 TELPAWSK LQKIYESQGKTL SVKQEFQKDAKRFEKLNKTF TNYD GSKILFDYSKNLVNDE 72

Query: 62 TLAKLQDLAKECDLAGAIKSMFSGEKINRTENRAVLHVALRNRSNTPILVDGKDVMPEVN 121
          +A L +LAKE ++ G   +MF GE IN TE+RAV HVALRNR+N P+ VDG +V PEV+
Sbjct: 73 IIAALIELAKEANVTGLRDAMFKGEHINSTEDRAVYHVALRNRANKPMYVDGVNVAPEVD 132
    
```

- ◆ デスクトップ上に、「kiso」フォルダを作成してください



1. 生物配列解析基礎

授業の目標・概要

生命科学のためのデータベースの利用と基本的な解析手法について講義します。データベースの基礎、配列データベース、機能データベース、ホモロジー検索、モチーフ解析などの基本的な手法について解説します。

kiso2

blast.pl

result.txt

result.txt

blast.pl

の2つのファイルをダウンロードして、kisoフォルダに入れてください

- コマンドプロンプトを立ち上げてください

 スタート → Windowsシステムツール → コマンドプロンプト

まず, kisoフォルダに移動します

```
> cd ␣
```

「cd (スペース)」と入力した後 (まだEnterキーは押さない), kisoフォルダをコマンドプロンプト上にドラッグ&ドロップしてください

下記のように表示されますので, Enterキーを押してください

```
> cd C:¥Users¥iu¥Desktop¥kiso
```

- **blast.pl**をメモ帳やワードパッドなどを使って開いてください

```
#!/usr/local/bin/perl
```

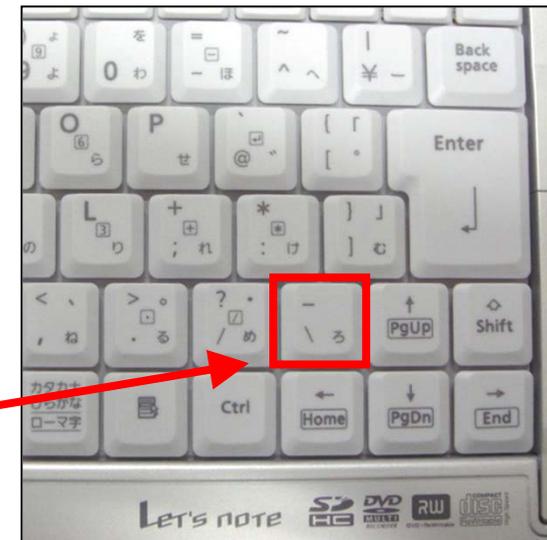
- 以下のように編集して、上書き保存してください

```
#!/usr/local/bin/perl  
print "Hello!¥n";
```

¥n は改行を表します

「¥」は、バックslash「\」を押してください

Windows上だと「¥」と表示されます



以下のコマンドを入力して、プログラムを実行してください

```
> perl blast.pl
```

変数

- 変数は, 「**\$文字列**」で表します
- 以下のように編集して保存してください

```
#!/usr/local/bin/perl  
$a = "Hello!¥n";  
print $a;
```

「;」を入力し忘れないように注意してください

Perlでは**行の終わり**に「;」をつける決まりになっています

以下のコマンドを入力して, プログラムを実行してください

```
> perl blast.pl
```

<STDIN>

- <STDIN>は, 1回呼び出すごとに**標準入力**から1行のデータを読み出す命令です. STDINとは, standard inputつまり標準入力の略です.
- 以下のプログラムに修正して保存してください.

blast.pl

```
#!/usr/local/bin/perl
$a = <STDIN>;      # 標準入力から1行のデータを読み出し, $aに代入する
print $a;          # $aを出力する
```

- 以下のコマンドを打ち込みblast.plを実行すると, **入力待ち**になります. キーボードで何か文字を入力すると, 入力した文字がそのまま出力されます.

```
> perl blast1.pl
```

リダイレクト

- キーボードから文字列を打ち込む代わりに、既存のテキストファイルから**データを読み込ませる**こともできます。
- 「result.txt」という**BLAST検索結果のファイル**を用意しておきました。中身を見てください。

```
BLASTP 2.2.19 [Nov-02-2008]
```

```
Reference: Altschul, Stephen F., Thomas L. Madden, Alejandro A. Schaffer,  
Jinghui Zhang, Zheng Zhang, Webb Miller, and David J. Lipman (1997),  
"Gapped BLAST and PSI-BLAST: a new generation of protein database search  
programs", Nucleic Acids Res. 25:3389-3402.
```

```
.  
.
```

リダイレクト

- 「result.txt」を読みこませるために、リダイレクトという機能を使います。

```
> perl blast.pl < result.txt
```

1

- 「result.txt」の一番最初の行だけが表示されます。
- 結果を、画面でなく、ファイルに出力することもできます。

```
> perl blast.pl < result.txt > output1.txt
```

2

3

whileループ

- 読み込むデータが何行であっても良いように、**whileループ**を利用してみましょう。

```
#!/usr/local/bin/perl
while ($a = <STDIN>) { # データを1行ずつ$aに
    print $a;          # 代入して, 出力する
}
```



読み込むデータ（行）がある限り「真」となり、{ }の中が繰り返されます。

- 実行してみましょう。全てのデータが出力されるはずです。

```
> perl blast.pl < result.txt > output1.txt
```

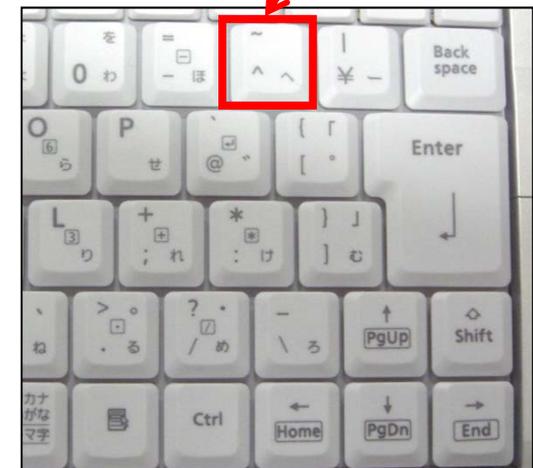
パターンに当てはまるかどうかを調べる「=~/」

- 文字列「DNA」を含む行を表示してみましょう。

```
#!/usr/local/bin/perl
while ($a = <STDIN>) {
    if ($a =~ /DNA/) {      # DNAを見つけたら
        print $a;          # 出力する
    }
}
```

「~」チルダ
Shiftを押しながら

- 「=～」をパターン結合演算子, 「/**文字列**/」をマッチ演算子といいます。
- 文字列の一部に一致すれば「**真**」を, 一致しなければ「**偽**」を返します
- `$a =~ /DNA/` は, `$a`にDNAという文字列が含まれていれば「**真**」となり, if文の{ }が実行されます。



- 文字列「Query=」を含む行を表示してみましょう。

```
#!/usr/local/bin/perl
while ($a = <STDIN>) {
    if ($a =~ /Query=/) { # Query=を見つけたら
        print $a;        # 出力する
    }
}
```

質問配列の行だけを抽出することができるようになりました

```
Query= gi|13507740|ref|NP_109689.1| DNA polymerase III beta subunit
Query= gi|13507741|ref|NP_109690.1| similar to j-domain of DnaJ
Query= gi|13507742|ref|NP_109691.1| DNA gyrase subunit B [Mycoplasma
Query= gi|13507743|ref|NP_109692.1| DNA gyrase subunit A [Mycoplasma
Query= gi|13507744|ref|NP_109693.1| seryl-tRNA synthetase [Mycoplasma
      .
      .
      .
```

正規表現による検索

- マッチ演算子の `/` と `/` の間には、文字列だけでなく、**パターン**と呼ばれるものを入れることができます
- パターンとは、「Mで始まる文字列」や「3文字の文字列」など、**文字列の特徴**を記述したものです
- このパターンの記述方法を**正規表現**といいます。

例えば、

DNA

DNNA

DNNNNNA

DNNNNNNNNNNNNNNNNNA

これらすべてを検索するには、**`/DN+A/`** と記述します



■文字クラス

以下のものは、次のような1文字にマッチします。

[abc]	aかbかcのどれか
[a-z]	任意の小文字
[^abc]	aでもbでもcでもない文字
3 \d	数字 (digit)
\D	数字以外
\w	英数字 (word)
\W	英数字以外
2 \s	空白文字 (space)
\S	空白文字以外
\b	単語境界 (word boundary)
1 .	任意の一文字

■位置指定

パターンの位置を指定します。

4 ^	先頭
\$	末尾

■エスケープ

/、^、\$などの、正規表現的に意味のある特殊記号自体を検索したい局面では、¥でエスケープします。

5 ^¥^	^という字で始まる行にマッチ
¥¥	¥自体にマッチ

■ 繰り返し

以下の記号を使って、文字または文字クラスの繰り返しとマッチします。ここでは文字または文字クラスをxと書きます。

x*	0回以上の繰り返し
6 x+	1回以上の繰り返し。xx*と同じ
x?	0回か1回
x{5}	5回繰り返し。xxxxxと同じ
x{3,}	3回以上繰り返し。xxx+と同じ
x{3,5}	3回以上5回以下繰り返し。xxxx?x?と同じ

■ グループと選択

文字列を繰り返すときは () を使ってグループ化します。

su(mo)+ sumo, sumomo, sumomomo などにマッチする

いくつかのパターンのどれかにマッチさせるときは | を使います。

love|kiss love か kiss にマッチする

stud(y|ies) study か studies にマッチする

su(mi|mo){2,3} sumimi, sumimo, sumomi, sumomo, sumimimi, sumimimo,
sumimomi, sumomimi, sumomomi, sumomimo, sumimomo,
sumomomo のいずれかにマッチする

正規表現による検索

- Gene indexを含む文字列を抽出してみましょう。

Query= gi|13507742|ref|NP_109691.1| DNA gyrase

Query= と ref ではさまれた連続した文字列を含む行を抽出するには

「.」（任意の文字） と 「+」（1文字以上の連続文字）

を使って以下のようにします

```
#!/usr/local/bin/perl
while ($a = <STDIN>) {
    if ($a =~ /Query= .+ref/) {
        print $a;
    }
}
```



カッコを使った記憶

- マッチ演算子のパターンの中で括弧 () を使うと、その括弧で囲まれた部分が、\$1, \$2,... という**特殊変数**に格納されます。

```
#!/usr/local/bin/perl
while ($a = <STDIN>) {
    if ($a =~ /Query= (.+)ref/) {
        print $1;
    }
}
```

改行

- 改行されるように, "¥n" を入れます

```
#!/usr/local/bin/perl
while ($a = <STDIN>) {
    if ($a =~ /Query= (.+)ref/) {
        print "¥n", $1;
    }
}
```

- BLAST検索の結果, ヒットしたタンパク質の情報

(例えば

>ref|NP_072866.1| topoisomerase IV, subunit A)

を含む行を抽出し, タブ区切りで表示してみましよう

```
#!/usr/local/bin/perl
while ($a = <STDIN>) {
    if ($a =~ /Query= (.+)ref/) {
        print "¥n", $1;
    }
    if ($a =~ />ref/) {
        print "¥t", $a;
    }
}
```

- ヒットしたタンパク質情報のref番号だけを抽出してみましょう
>ref|NP_072866.1| topoisomerase IV, subunit A)
- "|"ではさまれた文字列を取り出したいのですが、以下の表現ではうまくいきません

```
#!/usr/local/bin/perl
while ($a = <STDIN>) {
    if ($a =~ /Query= (.+)ref/) {
        print "¥n", $1;
    }
    if ($a =~ />ref|. +|/) {
        print "¥t", $a;
    }
}
}
```

「|」

Shiftを押しながら



- "|"は正規表現で使用する特殊な文字であるため、別の意味になってしまうからです
- ここで使う "|" が正規表現でないことを示すために、¥ (バックスラッシュ) を前につけます

```
#!/usr/local/bin/perl
while ($a = <STDIN>) {
    if ($a =~ /Query= (.+)ref/) {
        print "¥n", $1;
    }
    if ($a =~ />ref¥|. +¥|/) {
        print "¥t", $a;
    }
}
```

- 括弧を使って、ref番号だけを抽出してみましょう

>ref|NP_072866.1| topoisomerase IV, subunit A

```
#!/usr/local/bin/perl
while ($a = <STDIN>) {
    if ($a =~ /Query= (.+)ref/) {
        print "¥n", $1;
    }
    if ($a =~ />ref¥| (.+)¥| /) {
        print "¥t", $1;
    }
}
```

質問配列とヒットした配列のアクセシオン番号を抽出できるようになりました

QueryのGene Index

ヒットしたタンパク質のref番号

gi 13507740	NP_072661.1	
gi 13507741	NP_072662.1	
gi 13507742	NP_072663.1	NP_072865.1
gi 13507743	NP_072664.1	NP_072866.1
gi 13507744	NP_072665.1	
gi 13507745	NP_072666.1	
gi 13507746	NP_072667.1	
gi 13507747	NP_072668.1	NP_072998.1
gi 13507748	NP_072669.1	
.		
.		
.		

- **最も相同性の高い配列**の情報だけを表示するようにしてみましょう
- 新たに `$b` という変数を使い、これが **0** か **1** かを指標にします
- 質問配列 (Query= の行) を見つけたら **`$b = 1`** にします
- その後の最初に出てくるヒット配列 (>ref の行) を見つけたら、番号を抽出して **`$b = 0`** に戻します (次の質問配列を見つけるまで抽出しません)

```
#!/usr/local/bin/perl
while ($a = <STDIN>) {
    if ($a =~ /Query= (.+)ref/) {
        print "¥n", $1;
        $b = 1;
    }
    if ($a =~ />ref¥l(.+)¥l/ && $b == 1) {
        print "¥t", $1;
        $b = 0;
    }
}
```

質問配列と最も相同性の高いヒット配列のアクセション番号を抽出できるようになりました

QueryのGene Index

ヒットしたタンパク質のref番号

gi|13507740|

NP_072661.1

gi|13507741|

NP_072662.1

gi|13507742|

NP_072663.1

gi|13507743|

NP_072664.1

gi|13507744|

NP_072665.1

gi|13507745|

NP_072666.1

gi|13507749|

gi|13507750|

gi|13507751|

gi|13507752|

gi|13507753|

NP_072670.1

gi|13507754|

NP_072671.1

·
·
·

<課題 1>

- 質問配列のGene Indexのうち、**数字の部分だけ**を取り出して、以下のような出力結果になるプログラムを作成してください

(44ページ参照)

13507740	NP_072661.1
13507741	NP_072662.1
13507742	NP_072663.1
13507743	NP_072664.1
13507744	NP_072665.1
13507745	NP_072666.1
13507746	NP_072667.1
13507747	NP_072668.1
13507748	NP_072669.1
	.
	.

<課題 2>

- **E-valueの値**を抽出して、以下のような出力結果になるプログラムを作成してください

質問配列の Gene Index	ヒットしたタンパク質 のref番号	E-value
13507740	NP_072661.1	e-148
13507741	NP_072662.1	e-125
13507742	NP_072663.1	0.0
13507743	NP_072664.1	0.0
13507744	NP_072665.1	0.0
13507745	NP_072666.1	1e-077
	・	
	・	

課題 1 と同じファイル名にならないように、**output2.txt**という
ファイル名で結果を出力してください

- 出力したテキストファイル (**output1.txt** と **output2.txt**) を、メールに添付して提出してください
- 送付先は「kenro@hosei.ac.jp」です
- メールのはじめの件名は「**Perl課題**」にしてください
- メール本文に、以下のように「氏名」「所属」「学生証番号」「本日の講義の感想」を記載してください

氏名：○○ ○○

所属：××××専攻 △△△△研究室

学生証番号：□□□□□

講義の感想：