생명정보학에서 쓰이는 컴퓨터 Perl 언어의 기초 교육



서울시 강서구 화곡동 359-63 영주빌딩 201호 Tel 02 2698 1188 / Fax 02 6280 8821

www.infoboss.co.kr

2017/9/9 Jongsun Park, Ph. D.





Bioinformatics

From Wikipedia, the free encyclopedia

For the journal, see Bioinformatics (journal).

Bioinformatics /_bai.ou_infər'mætiks/ (listen) is an interdisciplinary field that develops methods and software tools for understanding biological data. As an interdisciplinary field of science, bioinformatics combines computer science, statistics, mathematics, and engineering to analyze and interpret biological data. Bioinformatics has been used for *in silico* analyses of biological queries using mathematical and statistical techniques.

Bioinformatics is both an umbrella term for the body of biological studies that use computer programming as part of their methodology, as well as a reference to specific analysis "pipelines" that are repeatedly used, particularly in the field of genomics. Common uses of bioinformatics include the identification of candidate genes and single nucleotide polymorphisms (SNPs). Often, such identification is made with the aim of better understanding the genetic basis of disease, unique adaptations, desirable properties (esp. in agricultural species), or differences between populations. In a less formal way, bioinformatics also tries to understand the organisational principles within nucleic acid and protein sequences, called proteomics.^[1]

From Wikipedia

- Bioinformatics is a field to **develop methods** and **software tools** for analyzing all kinds of biological data.
- Main way to analyze biological data (efferently or systematically) is writing a program.

- Bioinformatics started with needs to deal with huge amount of data (sequences) which human cannot handle.
- Human Genome Project (HGP), of which aim is unrevealing human genome sequences fully, generated
- a lot of data using Sanger Sequencing Method.
- Several programs, such as phred, phrap, and consed, had been developed for automation of sanger sequencing data.





Bioinformatics (3) History of Bioinformatics

- With the aid of program, in the late phase of human genome project, new sequencing method to uncover human genome sequences was addressed by Craig Venter.
- International group for human genome sequences keep going with traditional method and published paper in 2001.
- Craig Venter made his company (Celera Genomics) in 1999 and finished his genome sequences at the same time to international

group.

He published his genome sequence in nature and international group published in Science.

Venter, J.C. et al. The sequence of the human genome. Science 291, 1304-1351 (February 16, 2001). Lander, E.S. et al. The Genome International Sequencing Consortium. Intial sequencing and analysis of the human genome. Nature 409, 860-921 (February 15, 2001).



- Genomics is a starting point of bioinformatics in early phase.
- More technologies which can see global pictures at the diverse levels, DNA(=genome),
 - RNA(=transcriptome), proteins (=proteomics), and primary and secondary metabolites (=metabolomics),

have been developed with forming diverse bioinformatics area.



- Sequencing technology is for reading nucleic acid one by one.
- Sanger developed Sanger's Method and published in Nature, 1977, with uncovering sequence of bacteriophage.
- This technique has been automated with florescence dyes (ABI3730).

Journal home > Archive > Article > Full Text

Article

Nature 265, 687-695 (24 February 1977) | doi:10.1038/265687a0; Accepted 24 December 1976

Nucleotide sequence of bacteriophage ϕ X174 DNA

F. Sanger, G. M. Air^{*}, B. G. Barrell, N. L. Brown[†], A. R. Coulson, J. C. Fiddes, C. A. Hutchison, III[‡], P. M. Slocombe[§] & M. Smith[¶]

1. MRC Laboratory of Molecular Biology, Hills Road, Cambridge CB2 2QH, UK 2. Present addresses : *, John Curtin School of Medical Research, Microbiology Department, Canberra City ACT 2601, Australia, ^{+,} Department of Biochemistry, University of Bristol, Bristol BS8 1TD, UK, ^{‡,} Department of Bacteriology and Immunology, University of North Carolina, Chapel Hill, North Carolina 27514, ^{§,} Max-Planck-Institut für Molekulare Genetik, 1 Berlin 33, FRG, [¶]Department of Biochemistry, University of British Columbia, Vancouver BC, Canada V6T 1W5,

A DNA sequence for the genome of bacteriophage ϕ X174 of Top approximately 5,375 nucleotides has been determined using the rapid and simple 'plus and minus' method. The sequence identifies many of the features responsible for the production of the proteins of the nine known genes of the organism, including initiation and termination sites for the proteins and RNAs. Two pairs of genes are coded by the same region of DNA using different reading frames.

To read this story in full you will need to login or make a payment (see right).

ARTICLE TOOLS

- Send to a friend
- Export citation
- Export references
- Rights and permissions

SEARCH PUBMED FOR

- F. Sanger
- G. M. Air
- B. G. Barrell
- N. L. Brown
- A. R. Coulson
- J. C. Fiddes
- more authors of this article





- Sanger Sequencing technology allows us to read DNA sequences up to 1,000 bp per one reaction.
- However, chromosome or most of our targets are more than 1 kb, so that we need to make idea to get the data.
- In addition, through the automation, the cost of sequencing is not cheap due to the limitation of

technology.



Next Generation Sequencing Technology (1) Now, NGS

Next Generation Sequencing Technology was commercially launched in 2005: 454 Life Science.



-

-



Nextly, Solexa sequencer was launched in 2007.









Data format of NGS sequencer results is simple text-based one (fastq).

- It contains read name (=number), reads (nucleotide sequence), quality information (Q-33).
 - Based on ASCII code, quality data can be calculated.

Plant Genome / Structure of plant genomes





two organelle genomes: chloroplast and mitochondria.



Genome sequence (Nucleotide sequence)





- **1,381** plant genomes including red algae originated from **188** species are available.
- 418 Gbp genome sequences and 5.5 million plant genes (7.2 million transcripts) are ready to be analyzed.



Computer program

From Wikipedia, the free encyclopedia

For the TV programme, see The Computer Programme.

A **computer program** is a collection of instructions^[1] that performs a specific task when executed by a computer. A computer requires programs to function and typically executes the program's instructions in a central processing unit.^[2]

- Program is a collection of instructions to make computer works.



Programming Language /

A **programming language** is a formal language that specifies a set of instructions that can be used to produce various kinds of output. Programming languages generally consist of instructions for a computer. Programming languages can be used to create programs that implement specific algorithms.





- Perl is a family of high-level, general-purpose, interpreted, dynamic programming languages.
- Perl is OS-independent: it can be used under Windows and Linux(UNIX).
- Perl was originally developed by Larry Wall in 1987 as a general-purpose



Unix scripting language to make report processing easier. (Regular expression is the best component

in Perl.

- Perl syntax is similar to that of C, AW

```
Perl
#!/usr/bin/perl -w
use strict;
my $tot = 0;
for (my $i=0;$i<100;$i++) {
        $tot += $i;
}
print "Total : $tot\n";
~</pre>
```





Perl (2) Environment for using Perl



- In Windows, you have to install Perl Interpreter.
- ActiveState Perl (<u>http://www.activestate.com/activeperl/</u>) can be downloaded freely.

\rightarrow G	● 안전함 https://www.activestate.com/activeperl							G
	ActiveState The open source languages company		EDITIONS	INDUSTRY	SUPPORT	1.866.631.4581	Contact Sales	c Q
		SOLUTIONS					BLOG	
	Home » Solutions » Perl » ActivePerl							

ACTIVEPERL

ActivePerl Business and Enterprise Editions feature our precompiled, supported, quality-assured Perl distribution used by millions of developers around the world for easy Perl installation and quality-assured code. When you're using Perl on production servers or mission-critical applications, ActivePerl Business and Enterprise Editions offer significant time savings over open source Perl for installing, managing, and standardizing your Perl.

REDUCE RISK WITH COMMERCIALLY SUPPORTED PERL

- Comply with corporate policy requirements to have supported open source products
- Full license review including all precompiled third-party Perl modules with assurances to minimize risk (Enterprise Edition only)
- > Protect your organization from legal risk with indemnification coverage (Enterprise Edition only)

EXTENDED PLATFORM AND VERSION SUPPORT

Sale! ActivePerl Business Edition \$1200 \$999

Get a Quote Enterprise Edition or Business Edition

> **Download ActivePerl** Free Community Edition







- Here is simple Perl code.

print "Hello, perl.\n";

- Running this Perl script like below:
 - Make helloperl.pl on the Desktop.



Date modified: 9/9/2017 9:05 AM Size: 23 bytes

- Open 'cmd' via Windows+R key.



- Run the program

C:₩Users₩Jongsun Park>cd Desktop C:₩Users₩Jongsun Park₩Desktop>perl helloperl.pl Hello, perl. C:₩Users₩Jongsun Park₩Desktop>_





- Variable is a space for storing value in memory.
- In Perl, variable can be defined using 'my' command:

my \$a;

- '\$' means scalar variable in Perl, which is normal variable.
- After defining variable, you can 'assign' value on it because it is a space for storing value.



- Value can be numeric as well as string because Perl allow any types of value be stored in variable.



- Function is a set of commands with arguments and return values.



- Perl provide built-in functions in interpreter: for example, 'print'.

```
print "Hello perl!\n";
print $a;
```

- print function display string or variables on the screen.

```
print "a valie is ".$a."\n";
```

Perl (5)

Functions in Perl

Server Account / Let's make server account for further programming and running bioinformatic tools



- Server-client is important concept for bioinformatic analysis.
- Operating system is Linux (=UNIX) not Windows.



- 'Putty' is program to access server.



Thank you for your attention!

If you have question, please ask! =)

starflr@infoboss.co.kr

