

United States Department of Agriculture National Institute of Food and Agriculture





Bioinformatics 101 Part I: the tools

David Francis The Ohio State University, OARDC SolCAP workshop

This module Introduces some basic tools used for bioinformatics. After following this module, you should be able to:

Describe the purpose of BLAST, perl and BioPerl in building pipelines for marker discovery

Find and install BLAST

Format a FASTA file as a database for BLAST searches

Perform BLAST searches

Use PERL and BioPerl to parse BLAST searches



Basic Local Alignment Search Tool

National Center for Biotechnology Information (NCBI) – 1988 **BLAST – Basic Local Alignment Search Tool**

Blast finds regions of similarity between biological sequences. The program compares nucleotide or protein sequences to sequence databases

NCBI maintains databases that are freely available to the public for download. YOU MAY WANT TO USE YOUR OWN OR **CREATE ONE TO ADDRESS SPECIFIC QUESTIONS**



NC











Basic Local Alignment Search Tool

SETUP TUTORIAL

STANDALONE

The BLAST Family

blastp protein/protein

blastn

nucleotide/nucleotide

tblastx

Translated nucleotide vs Translated nucleotide

blastn

blastx

nucleotide/protein

Similarity search programs

Resources () How To () S NCBI

NCBI Search All Databases ¥ National Center for Search Clear Biotechnology Information Resources Welcome to NCBI Popular Resources NCBI Home All Resources (A-Z) The National Center for Biotechnology Information advances science and BLAST health by providing access to biomedical and genomic information. Bookshelf Chemicals & Bioassays Gene Data & Software More about the NCBI | Mission | Organization | Research | RSS Genome Nucleotide DNA & RNA OMIM Domains & Structures Protein Genome Genes & Expression

Genetics & Medicine

Genomes & Maps

Homology

Literature

Proteins

Sequence Analysis

Taxonomy

Training & Tutorials

Variation

1000 prokaryotic genomes are now completed and available in the Genome database.



How To...

- Determine conserved synteny between the genomes of two organisms
- Find a homolog for a gene in another organism
- Obtain the full text of an article
- Design PCB primers and check them for specificity



- PubChem
- PubMed
- PubMed Central
- SNP

NCBI News

NCBI Workshop: A Practical Guide to Genome-Scale Data

03 Nov 2010

Presented at the ASHG Meeting on Thursday, November 4, 2010,

New Find-in-sequence feature for Nucleotide and Protein databases



Training & Tutorials Variation	How To Determine co Find a homolo Obtain the full Design PCR p See all	nserved synteny between the gen og for a gene in another organism text of an article orimers and check them for specifi	omes of two organisms city	03 Nov 2010 Presented at the ASHG Meeting on Thursday, November 4, 2010, New Find-in-sequence feature for Nucleotide and Protein databases 18 Oct 2010 See the featured article in the September issue of the NCBI More
You are here: NCBI		8		Write to the Help Desk
GETTING STARTED	RESOURCES	POPULAR	FEATURED	NCBLINFORMATION
NCBI Help Manual	Literature	PubMed	GenBank	About NCBI
NCBI Handbook	DNA & BNA	Nucleotide	Reference Sequence	ces Besearch at NCBI
Training & Tutorials	Proteins	BLAST	Man Viewer	NCBI Newsletter
Training a Fatoriaio	Sequence Analysis	PubMed Central	Genome Projects	NCBI FTP Site
	Genes & Expression	Gene	Human Genome	
	Genomes & Maps	Bookshelf	Mouse Genome	
	Domains & Structures	Protein	Influenza Virus	
	Genetics & Medicine	OMIM	Primer-BLAST	
	Taxonomy	Genome	Sequence Read Arc	chive
	Data & Software	SNP		
	Training & Tutorials	Structure		
	Homology			
	Small Molecules			
	Variation			-
				L
'http	ו.www.ו	ncbi.nlm	.nih.go	V/'



'http://www.ncbi.nlm.nih.gov/'

~

^

🔶 👻 🍥 😣	🕋 😒 http://www.ncbi.nlm.nih.gov/Ftp/
S NCBI	FTP site
PubMed Entr Search All Databa	ez BLAST OMIM Books TaxBrowser Structure
NCBI SITE MAP Guide to NCBI	Major resources available by ftp (ftp.ncbi.nih.gov): BLAST Basic Local Alignment Search Tool
resources	Download the BLAST database and stand-alone sequence comparison software.
About NCBI The science behind our resources. An introduction for researchers,	<u>CDD Data</u> Download data from the Conserved Domain Database. <u>CD-Tree</u>
educators and the public.	Download the protein domain hierarchy viewer and editor.
GenBank sequence submission support and software	<u>Cn3D</u> Download the stand-alone software for viewing 3-dimensional structures.
Molecular databases	Data Repository Download collections of contributed molecular biology data.
sequences, structures and taxonomy	dbGaP Download open access Genotype and Phenotype data.
ł	nttp://www.ncbi.nlm.nih.gov/Ftp/

Example 2.2.20 (April 2009)

	📷 Most Visited 🗸 🌘 Getting Started 🔋	🔂 Latest Headlines 🗸				
	Index of ftp://ftp.r	ncbi.nih.gov/blast	/executables/			
	🕆 Up to higher level dire	ectory				
	Name	\$		Size	Last M	odified
	LATEST			0	3/24/2010	12:00:00
				0	8/23/2010	01:54:00
	🚔 snapshot			0	9/13/2010	02:46:00
	Done					
	Done	atSh) 📄 [Questions.doc	.] 📄 [Untitled 1 - O] 🗑	BioInformat1	0	Blastinstru
	Done	atSh) 📄 [Questions.doc	. [[Untitled 1 - 0] @ 03/08/2006	BioInformat1	o 📄 [Blastinstru
	Done	atSh) 📄 [Questions.doc	. [Untitled 1 - 0] @ 03/08/2006 05/26/2006	BioInformat1 12:00:00 / 12:00:00 /	0) 🗇 [AM AM	BlastInstru
	Done Index of ftp://ft	tSh) 📄 [Questions.doc	. [Iuntitled 1 - 0] () 03/08/2006 05/26/2006 10/20/2006	BioInformat1 12:00:00 / 12:00:00 / 12:00:00 /	0) () (AM AM	BlastInstru
Do =	ne Index of ftp://ft 📔 🗈 [LinuxChea	atSh) 📄 [Questions.doc	. [Untitled 1 - 0) 03/08/2006 05/26/2006 10/20/2006 04/18/2007	BioInformat1 12:00:00 / 12:00:00 / 12:00:00 /	0) () (AM AM AM	Blastinstru
	Done	atSh) 📄 [Questions.doc	.) [Untitled 1 - 0] () 03/08/2006 05/26/2006 10/20/2006 04/18/2007	BioInformat1 12:00:00 / 12:00:00 / 12:00:00 / 12:00:00 /	0) (C) (AM AM AM AM	BlastInstru
	Done	atSh) 📄 [Questions.doc	. Untitled 1 - 0 03/08/2006 05/26/2006 10/20/2006 04/18/2007 11/01/2007	BioInformat1 12:00:00 / 12:00:00 / 12:00:00 / 12:00:00 / 12:00:00 /	0) [] AM AM AM AM AM	BlastInstru
	Done	atSh)	. [Untitled 1 - 0] 03/08/2006 05/26/2006 10/20/2006 04/18/2007 11/01/2007 05/02/2008	BioInformat1 12:00:00 / 12:00:00 / 12:00:00 / 12:00:00 / 12:00:00 /	0) () (AM AM AM AM AM AM	BlastInstru

?

🔝 🕋 📮 ftp://ftp.ncbi.nih.gov/blast/executables/

<u>H</u>elp

Index of ftp://ftp.ncbi.nih.gov/blast/executables/ - Mozilla Firefox

📶 📋 🕸 Sat Nov 6, 9:32 AM 🛛 Da

☆ 🗸 🖌 Google

<u>F</u>ile

Applications Places System

Ĉ.

Edit View History Bookmarks Tools

✤ Up to higher level directory

Name

Index of ftp://ftp.ncbi.nih.gov/blast/executables/release/2.2.20/

🕆 Up to higher level directory

Name

last-2.2.20-ia32-freebsd.tar.gz blast-2.2.20-ia32-linux.tar.gz 📔 blast-2.2.20-ia32-solaris.tar.gz Interpretation blast-2.2.20-ia32-win32.exe blast-2.2.20-ia64-linux.tar.gz blast-2.2.20-sparc64-solaris.tar.gz blast-2.2.20-universal-macosx.tar.gz 当 blast-2.2.20-x64-linux.tar.gz blast-2.2.20-x64-solaris.tar.gz blast-2.2.20-x64-win64.exe ncbi.tar.gz ncbiz.exe netblast-2.2.20-ia32-freebsd.tar.gz netblast-2.2.20-ia32-linux.tar.gz netblast-2.2.20-ia32-solaris.tar.gz netblast-2.2.20-ia32-win32.exe netblast-2.2.20-ia64-linux.tar.gz netblast-2.2.20-sparc64-solaris.tar.gz netblast-2.2.20-universal-macosx.tar.gz netblast-2.2.20-x64-linux.tar.gz netblast-2.2.20-x64-solaris.tar.gz netblast-2.2.20-x64-win64.exe

Size	Last Modified						
25514 KB	03/01/2009	12:00:00 AM					
31357 KB	03/01/2009	12:00:00 AM					
29056 KB	03/01/2009	12:00:00 AM					
14364 KB	03/01/2009	12:00:00 AM					
77572 KB	03/04/2009	12:00:00 AM					
43674 KB	03/01/2009	12:00:00 AM					

Select version for your OS; Generally Win32 for Windows; ia32-linux for Linux or MAC (32 or 64 bit depends on system's processor)

My notebook processor is 32 bit

-🛟 Арр	olications	Place	es Sy	stem	82
*					
<u>M</u> onitor	<u>E</u> dit	<u>v</u> iew	<u>H</u> elp		
System	Processe	s Reso	urces	File S	ystems



Ubuntu

Release 9.04 (jaunty) Kernel Linux 2.6.28-11-generic **GNOME 2.26.1**

Hardware

Memory: 2.0 GiB Processor 0: Intel(R) Atom(TM) CPU N270 @ 1.60GHz Processor 1: Intel(R) Atom(TM) CPU N270 @ 1.60GHz

0

System Status

Available disk space: 121.6 GiB

Essentials	
Status	Launched
Launch Date	Q2'08
Processor Number	N270
# of Cores	1
# of Threads	2
Clock Speed	1.6 GHz
L2 Cache	512 KB
Bus/Core Ratio	12
FSB Speed	523 MHz
FSB Parity	No
Instruction Set	32-bit
Instruction Set Extensions	SSE2, SSE3, SSC +
Embedded Options Available	Tes .
Supplemental SKU	No
Lithography	45 nm
Max TDP	2.5 W
VID V Lige Range	0.9V-1.1625V
ay 1ku Budgetary Price	\$44.00
Package Specifications	



- 🚯 [PLAST] 🖹 [LipuxC] 🖹 [Questi] 🖹 [Liptitle] 📄 [Pielef] 📄 [Plast] 🔲 [devid

All Essentials

Package Specifications

Advanced Technologies COMPATIBLE PRODUCTS

BLOCK DIAGRAMS

2

ORDERING / SSPECS / STEPPINGS

[20040

Custom

2

<u>F</u> ile <u>E</u> dit <u>V</u> iew Hi <u>s</u> tory <u>B</u> ookmarks <u>T</u> o	ols <u>H</u> elp		5 ¹ 0 010
츶 🗼 🗸 🍪 🛞 💼 🖳 ftp://ftp.nc	bi.nih.gov/blast/executables/release/2.2.20/	☆ 🗸 🔽 Google	0
📷 Most Visited 🗸 🌸 Getting Started 🔝 Late	Opening blast-2.2.20-ia32-linux.tar.gz		
🕈 Up to higher level directory	You have chosen to open		^
Name blast-2.2.20-ia32-freebsd. blast-2.2.20-ia32-linux.tar blast-2.2.20-ia32-solaris.t blast-2.2.20-ia32-win32.es blast-2.2.20-ia64-linux.tar blast-2.2.20-sparc64-solari blast-2.2.20-x64-linux.tar blast-2.2.20-x64-solaris.tar blast-2.2.20-x64-solaris.tar blast-2.2.20-x64-win64.ex ncbi.tar.gz ncbi.tar.gz ncbiz.exe netblast-2.2.20-ia32-freek netblast-2.2.20-ia32-freek netblast-2.2.20-ia32-solar netblast-2.2.20-ia32-solar netblast-2.2.20-ia32-solar netblast-2.2.20-ia32-win3	blast-2.2.20-ia32-linux.tar.gz which is a: Gzip archive from: ftp://ftp.ncbi.nih.gov What should Firefox do with this file? Open with Archive Manager (default) • Save File • Do this automatically for files like this from now on. • Cancel • Save File • Do this automatically for files like this from now on. • Save File • Do this automatically for files like this from now on. • Save File • Do this automatically for files like this from now on. • Save File • Do this automatically for files like this from now on. • Save File • Do this automatically for files like this from now on. • Save File • Do this automatically for files like this from now on. • Save File • Do this automatically for files like this from now on.	Modified 009 12:00:00 AM 009 12:00:00 AM	
Done		0000 10000111	
🕑 Index of ftp 📄 [LinuxCheat 📄	[Questions 📄 [Untitled 1 📄 BioInformat 📄 [BlastInstr	u 国 david@dell [ງ 🔂

Create a folder called 'BLAST' (preferably on your main drive to simplify PATH statements e.g. C:\) and save the file to that folder.

'Go to the Blast directory: cd c: /Blast
 'Windows Install: ./blast-2.2.20-ia32.win32.exe
 'Unix (incl. MAC) Install: ./tar zxf blast-2.2.20-ia32-linux.tar.gz

(MAC: you may find that an additional folder was created "blast-2.2.20". In there you will find the folders below)

In the Blast directory there will now be 3 new folders



This ends our discussion of BLAST Next up: BioPerl



BioPerl



Windows users will first need to install perl Perl comes installed with MAC and Linux operating systems

Perl is: <u>Practical Extraction and Report</u> <u>Language</u>; a programing language for easily manipulating text, files, and processes

BioPerl



BioPerl is an open source project that develops modules for biological data in Perl.

A Perl module is a reusable package defined in a library file.

BioPerl modules are stable and "easy" to use.

Modules include objects for sequence files, alignment files and database searching. These objects can interact: the objects provide a coordinated and extensible framework for computational biology.

BioPerl



BioPerl module names minimize 'namespace' collisions by separating parts of a name by a double colon (::). For example:.

- The module 'Bio::DB::GenBank'; instructs Perl to go to the <u>Database GenBank</u>
- This module can automate retrieval of a set of sequences
- The 'Bio::SearchIO' module is used for parsing an input file and creating an output file with the specified information.
- This module can be used to create tables that summarize results from BLAST searches

Historial web | Configuración de búsqueda | Acceder

~6

		Buscar
Aproximadamente	266.000 resultados (0,87 segundos)	Búsqueda avanzada
Source Source </th <th>e Bioperi Project is an international associa s for bioinformatics, genomics and life science g/ - En caché Perl Browse Modules</th> <th>tion of developers of open</th>	e Bioperi Project is an international associa s for bioinformatics, genomics and life science g/ - En caché Perl Browse Modules	tion of developers of open
Granadero Baigorria, Santa Fe ▼ Cambiar ubicación BioPerl Tuto HOWTOs Downloads Más resultad	ial Installing Bioperl for Unix Getting Started SeqIO los de bioperl.org »	
La Web Páginas en español Páginas de Argentina ▼ Más herramientas BioPerl Tuto 25 May 2010 from various Bi www.bioperl.or	rial - BioPerl - [Traducir esta página] This tutorial is an overview of Bioperl , it inclue operl documents including module documentat g/wiki/ BioPerl _Tutorial - En caché - Similares	des snippets of code and text ion
BioPerl - Wi BioPerl es una para aplicacione es.wikipedia.org	kipedia, la enciclopedia libre colección de módulos de Perl que facilitan el d es de bioinformática y/wiki/ BioPerl - En caché - Similares	lesarrollo de scripts en Perl
Bioperl cour Introduction to l examples and e www.pasteur.fr/	<u>Se</u> - [Traducir esta página] Bioperl (www.bioperl.org). This course introduc exercises. The course content has been upgrade recherche/unites/sis//bioperl/ - En caché - S	ces to the bioperl modules with ed imilares
(PDF) Taller de Formato de arc de B Contreras-	e (bio)Perl nivo: PDF/Adobe Acrobat - Vista rápida Moreira - 2010 que necesitamos, con dos subrutinas de BioPe	ri podemos obtenerlas:



http://www.bioperl.org/wiki/Main_Page

http://www.bioperl.org/wiki/Getting_Started



Ouick avample



United States Department of Agriculture National Institute of Food and Agriculture





Bioinformatics 101 Part II: using the tools to for marker discovery

David Francis The Ohio State University, OARDC SolCAP workshop



olcap tomato PI128216 19659 olcap tomato PI128216 37082

00% + B C 129216 649

637 scaffold06004,

669 scaffold06004

2761110 2761110

genes



🖉 NCBI Sequence Viewer v2.0 - Windows Internet Explorer	
🕒 💽 👻 🕞 http://www.ncbi.nlm.nih.gov/sviewer/viewer.fcgi?too	il=portal&db=nucest&t 🔽 👉 🗙 Google
<u>File E</u> dit <u>V</u> iew F <u>a</u> vorites <u>T</u> ools <u>H</u> elp	
Google 🚽 🚼 Search	n 🕶 🧭 T 🖶 T 🔯 T 😭 Bookmarks T 🔉 🔌 🔹 🌑 franci T
😪 🏟 😂 NCBI Sequence Viewer v2.0	
SNCBI	GGATCCCCGGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
PubMed Nucleotide Protein Genome Search EST Image: for phytophthora [ORG]	Structure PMC Taxonomy OMIM Books N] AND Judelson [AUTH] AND T Go Clear Image: Clear
Display FASTA Show 20 Send to	
Item 1 - 20 of 3921	page 1 of 197 Previous <u>Next</u>
□ 1: <u>CV969339</u> . Reports PI010H11 infected[gi:5815]	9088] Links
>gi 58159088 gb CV969339.1 CV969339 PI010H11 lesion 3 dpi Phytophthora infestans cDNA, mF CGGGTCGACCCACGCGTCCGATTCATTTAATTGTTTATACACTT AAAGGTTCTTTCTAGAGTGGATCGTGAAATTTTTCGCTCGAAGO CTTCGGCGTCCGTATTCGTAGCCACCTTCAACTCAAGCTTGTAG ATATCATGAATTTATTCGAAGTAGCCAATATGGAAGAAGATATG CAATTCACATTCAAGATCGGTATTATTCTTAAATTCACGATGCA ATGAAACTTGCGGAAGGGTACTGAGTACAGAAGATCCATCTTAT AGTACATGTGGGCTTGATCCCACTTTAGAATCAGTGTTGGGCTA TTTTATGTTCTAGCCCTTATTAGTTTAACTAGTGTCTGGTGTCT C	Next sequence I infected tomato, outside of RNA sequence TAAAATGGATAACGATTGGAACAGGAA CCTGGCGGACATTGAAATCTGTGCCAT GATGCCATCAAAGAAGCTGGAAATTAA GTGGTTCAATGTCCCTTCAACACAAAC ACTAGTCAGATATGAAGAAAGTAAAGA TAAAAAAAAATATTGATTTGTTAAAA ACTTAAGACGGTATGTGACGAAAATAA TAAGCTATTTTTTTAAAAATCATTAGAA Internet Internet Internet Internet
🏄 Start 😰 🧊 🏉 🔹 👋 E. 🥭 2- 💹 2- 🕅 B. 🗁 H	1. 🕂 My Documents 🛅 Family » Desktop » « 🜖 뢧 12:28 PM

USE the Unix grep command to verify that you downloaded the entire file

NOTE: a 'cheet sheet of UNIX commands' is available on the PBGWorks wiki at:

http://pbgworks.org/node/901

Example: \$ grep -c '>' will count the number of times the '>' occurs, and therefore the number of sequences in a FASTA file.

For the file downloaded using following the NCBI EST database search using ENTREZ: Lycopersicum [ORGN] AND TA496 grep returns 116711, which matches our expectations

Lycopersicum [ORGN] AND Rio Grande

21973

Lycopersicum [ORGN] AND Rio Fuego

171

Lycopersicum [ORGN] AND MicroTom

120462

Lycopersicum [ORGN] AND TA496

116711

Lycopersicum [ORGN] AND Moneymaker

833

Phytophthora [ORGN] AND Judelson [AU] AND Tomato 3921 Format your database for BLAST

1) Use the formatdb command for this task.

2) formatdb.exe is located in the bin folder that was created when you unpacked BLAST. So, if you saved BLAST to a folder named Blast, bin will be located within Blast.

3) You need to tell the computer where to look for formatdb.exe, and where to look for the file that you want to format. This means specifying a PATH. Use cd to navigate to the bin folder (\$ cd c:/Blast/bin). The pwd and Is commands can be used to verify that you are in the proper path and that formatdb.exe is in the folder.



SYNTAX of the command:

\$/formatdb -i ./DatabaseName -p F

Note:

the ./ implies the database input file is in the bin folder with the formatdb.exe

If the database is in another folder, you must specify the path (C:/BLAST/DF/)

-p asks if the file contains protein data. Our answer is False, because the file contains DNA sequence. Use Is to list the files in the folder containing the database.

You should now see three new files: DatabaseName.nhr DatabaseName.nin DatabaseName.nsq Now we're ready to run a stand alone BLAST. SYNTAX of the command:

\$/blastall -p blastn -d ./DatabaseName -i ./QueryFile.txt -o Output

Note:

the ./ implies all files are in the bin folder

-p asks which program. We are using blastn
-d asks for the database (must be formatted)
-i asks for the input or query file (FASTA format)
-o Tell BLAST what you want to name the output file

Viewing results of a BLAST search will depend on the search:

A simple search may be viewed by opening the output file in a text editor

Some BLAST searches will return very large files. These are best examined with some basic UNIX commands (grep, head, tail, and less), and then parsed to organize the data.

Viewing the output file

<u>File Edit View T</u>erminal <u>H</u>elp

BLASTN 2.2.18 [Mar-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.

Query=

Database: ./TA496Seq1.txt 116,711 sequences; 60,705,274 total letters

Searching......done

Use the UNIX "less" command (followed by q to quit)

Sequences producing significant alignments:	Score (bits)	E Value
ail162298291ab187925062 1187925062 EST544951 tomato flower b	uds 15	30 0.0
gi 18263929 gb BM412299.1 BM412299 EST586626 tomato breaker f	rui 15	16 0.0
gi 13778955 gb BG643730.1 BG643730 EST511924 tomato shoot/mer	ist 15	07 0.0
gi 18262786 gb BM411156.1 BM411156 EST585483 tomato breaker f	rui 14	71 0.0
gi 18262599 gb BM410969.1 BM410969 EST585296 tomato breaker f	rui 14	53 0.0
gi 16232161 gb BI925777.1 BI925777 EST545666 tomato flower, b	uds 14	49 0.0
gi 18262889 gb BM411259.1 BM411259 EST585586 tomato breaker f	[.] rui 14	47 0.0
gi 16248431 gb BI933959.1 BI933959 EST553848 tomato flower, a	nth 14	33 0.0
gi 16236252 gb BI927083.1 BI927083 EST546972 tomato flower, 3	14	23 0.0
gi 14683291 gb BI205567.1 BI205567 EST523607 cTOS Solanum lyc	ope 14	03 0.0
gi 12627674 gb BG127486.1 BG127486 EST473132 tomato shoot/mer	ist 13	90 0.0

Parsing the output file

usage:

perl program <BLAST-report-file> # to extract <output_name >
#

use strict; use warnings; use lib "/home/users/David/lib/perl5"; use Bio::SearchIO;

Parsing the output file

```
# Usage information
die "Usage: $0 <BLAST-report-file> <number-of-top-hits> <output-file>\n", if (@ARGV != 3);
```

```
my ($infile,$numHits,$outfile) = @ARGV;
print "Parsing the BLAST result ...";
my $in = Bio::SearchIO->new(-format => 'blast', -file => $infile);
open (OUT,">$outfile") or die "Cannot open $outfile: $!";
```

```
# print the header info for tab-deliminated columns
```

```
print OUT "query_name\tquery_length\taccession_number\tlength\tdescription\tE value\tbit score\tframe\tquery_start\t";
print OUT "query_end\thit_start\thit_end\tpositives\tidentical\n";
```

```
# extraction of information for each result recursively
while ( my $result = $in->next result ) {
       # the name of the guery sequence
       print OUT $result->query name . "\t";
       # the length of the query sequence
       print OUT $result->query length;
       # output "no hits found" if there is no hits
       if ( $result->num hits == 0 ) {
                print OUT "\tNo hits found\n";
       } else {
                my scount = 0;
                # process each hit recursively
               while (my $hit = $result->next hit) {
                        print OUT "\t" if ($count > 0);
                        # get the accession numbers of the hits
                        print OUT "\t" . $hit->accession . "\t";
                        # get the lengths of the hit sequences
                        print OUT $hit->length . "\t";
                        # get the description of the hit sequences
                        print OUT $hit->description . "\t";
```

The Perl script checks for the expected three arguments (input file, number of hits to extract, output file); then it uses Bio::SearchIO to pull information from the blast report and put that information into a tab delimited file

3

Key Commands

perl blast_parsing_pl1.pl out_Ch11_blast 100 Ch11Parse

Results (open in EXCEL)

ces Sys	tem 🥹 📄 🕢							= 🐼	📶 🕸 Sເ	in Nov 7	7, 10:44 AI	M David	I
		Ch11F	Parse -	OpenO	ffice.or	g Calc	:					_ 0][
rt F <u>o</u> rma	t <u>T</u> ools <u>D</u> ata <u>W</u> i	indow <u>H</u> el	р										
à I 🗾 🔝 🚔 💩 I 🥙 👺 I 🔏 🗊 🗊 × 🍰 I 🥱 × 🛷 × I 🚳 🖧 Ž _n I 💣 📝 🖗 📼 🗃													
;	· 10 · A	<u> </u>					1 14 🗎		%	000. F +0 000			
				÷ ''									-
fee S		me											-
J(4) Z		inc											
В	С	D	E	F	G	н	1	J	К	L	М	N	
ry_length	accession_number	length	descrip	E value	bit score	frame	query_start	query_end	hit_start	hit_end	positives	identical	_
3203940	BI925062	805	EST54	0	1530	0	456189	456993	1	805	98.80%	98.80%	
	BM412299	765	EST58	0	1516	0	634586	635350	1	765	100.00%	100.00%	
	2					0	1223235	1223983	13	761	89.10%	89.10%	
	BG643730	781	EST51?	0	1507	0	971799	972577	2	781	99.50%	99.50%	
	BM411156	778	EST58	0	1471	0	1968468	1969245	1	. 778	98.80%	98.80%	
	BM410969	751	EST58	0	1453	0	1339681	1340433	1	751	99.60%	99.60%	
	BI925777	745	EST549	0	1449	0	456530	457272	1	. 742	99.70%	99.70%	
	BM411259	746	EST58	0	1447	0	634910	635653	1	. 746	99.70%	99.70%	
						0	1223550	1223992	299	743	88.50%	88.50%	
						0	623730	623926	1	. 197	99.50%	99.50%	
	BI933959	739	EST55	0	1433	0	634519	635257	1	739	99.50%	99.50%	
						0	1223157	1223893	2	738	88.30%	88.30%	
	BI927083	750	EST54	0	1423	0	146056	146799	1	. 745	99.50%	99.50%	
	BI205567	770	EST523	0	1404	0	419550	420265	17	732	99.70%	99.70%	
	BG127486	733	EST47?	0	1390	0	597213	597944	2	733	99.20%	99.20%	

Next Step:

Retrieve the desired sequence:

- a) Directly from the FASTA file
- b) from GenBank (using BioPerl)
- For (b) create a text file with the gb id of the sequences you want to retrieve:

BG643730

AF536200

AF536199...

(Chr11EST.txt = file containing a single column of gb id numbers)

Retrieving Sequence Data

#This script will extract sequences from Genbank# Only 2 arguments are required

an input file (with the accession numbers) and output file

use strict;

- use warnings;
- use lib "/home/users/David/lib/perl5";
- use Bio::DB::GenBank;
- use Bio::SearchIO;

Key Commands

perl GenbankSearch2.pl Chr11EST.txt Chr11_FASTA

Workshop Resources:

http://pbgworks.org/node/901

Perl script to test if BioPerl has been properly installed (returns "it works!)

Perl script that will parse a BLAST search

Perl scripts that allow user defined criteria to be considered during the BLAST parsing

Perl script that will retrieve specified sequences from GenBank (NCBI)



Questions?

