

Welcome to the Plant Breeding and Genomics Webinar Series

Today's Presenter: Dr. Candice Hansey

Presentation: <http://www.extension.org/pages/60428>



Host: Heather Merk

Technical Production: John McQueen

PBG home page: www.eXtension.org/plant_breeding_genomics

Sign up for PBG News: <http://pbgworks.org>

Please fill out the survey
evaluation! (You will be contacted
via email)

Watch past webinars and sign up
for future webinars!

<http://www.extension.org/pages/60426>

How to Align Sequences

Presenter: Candice Hansey

hansey @msu.edu

Michigan State University



Overview

- Navigating NCBI to obtain sequences
- Using BLAST for sequence alignment
- Using other programs for specialized sequence alignment
- Next generation sequence alignment programs

Overview

- Navigating NCBI to obtain sequences
- Using BLAST for sequence alignment
- Using other programs for specialized sequence alignment
- Next generation sequence alignment programs

Goal

- Obtain sequence for the maize teosinte branched1 gene and then find organisms with orthologous genes

NCBI

National Center for Biotechnology Information

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

National Center for Biotechnology Information

NCBI Resources How To

Gene

All Databases
PubMed
Protein
Nucleotide
GSS
EST
Structure
Genome
BioProject (Genome Project)
BioSample
BioSystems
Books
CancerChromosomes
Conserved Domains
dbGaP
dbVar
Epigenomics
Gene
GENSAT
GEO DataSets

NCBI Home
Site Map (A-Z)
All Resources
Chemicals & Bioassays
Data & Software
DNA & RNA
Domains & Structures
Genes & Expression
Genetics & Medicine
Genomes & Maps
Homology
Literature
Proteins
Sequence Analysis
Taxonomy
Training & Tutorials
Variation

NCBI

National Center for Biotechnology Information advances science and access to biomedical and genomic information.

Session | Organization | Research | RSS Feeds

data using NCBI software
NCBI data or software
how to accomplish specific tasks at NCBI
submit data to GenBank or other NCBI databases

Genomic Structural Variation

dbVar archives large scale genomic variation data and associates defined variants with phenotypic information.

Popular Resources

BLAST
Bookshelf
Gene
Genome
Nucleotide
OMIM
Protein
PubChem
PubMed
PubMed Central
SNP

NCBI News

New NCBI News Issue 07 Sep 2011

New Feature Highlighter in the sequence databases and Simple Object Access Protocol

NCBI Discovery Workshop: A Practical Hands-On Course 02 Aug 2011

September 27-28, 2011 @ NLM: Space is still available in the 2-day Discovery Workshop

More...

NCBI

The screenshot shows a web browser window with the URL [http://www.ncbi.nlm.nih.gov/gene/?term='Zea Mays'\[Organism\]](http://www.ncbi.nlm.nih.gov/gene/?term='Zea Mays'[Organism]). The page title is "Zea Mays"[Organism] - Gene - NCBI. The search bar contains "Zea Mays"[Organism] and the search button is labeled "Search". Below the search bar, there are links for "Save search", "Limits", and "Advanced". The page shows the results for the search, with a summary of 26135 results. The first two results are listed:

- ☐ [adh1](#)
1. alcohol dehydrogenase1 [**Zea mays**]
Other Aliases: Adh1-1F, Adh1-1S
Other Designations: alcohol dehydrogenase 1
ID: 542363
- ☐ [o2](#)
2. **Official Symbol:** o2 and **Name:** opaque endosperm2 [**Zea mays**]
Other Aliases: Opaque-e
Other Designations: opaque 2; opaque-2 protein; p-O2; p-QDA2D8; p-pMM1a; p-pXho0.9; p-pcrO2; p-phi057; p-phi112; p-umc1066; regulatory protein opaque-2
ID: 542375

On the right side, there is a "Filter your results:" section with the following options:

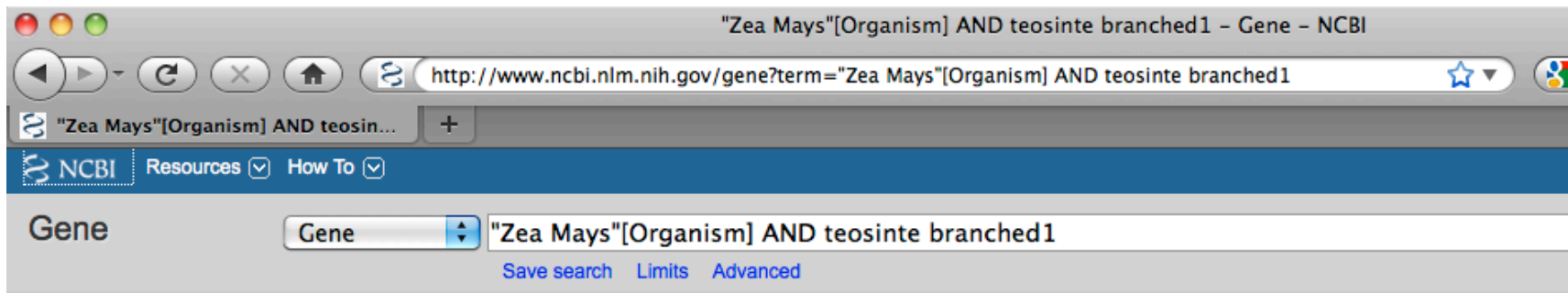
- All (26135)
- Current Only (26037)
- Genes Genomes (441)
- SNP GeneView (0)
- In Variation Viewer (0)

Below this, there is a "Top Organisms [Tree]" section with the following options:

- Zea mays (26135)
- Zea mays subsp. mays (218)
- Zea mays subsp. parviglumis (67)

Search specifying Zea Mays as the organism resulted in thousands of sequences

NCBI



[Display Settings:](#) ☐ Summary, Sorted by Relevance

[Send to:](#) ☐

Results: 3

☐ [tb1](#)

1. **Official Symbol:** tb1 and **Name:** teosinte branched1 [*Zea mays*]

Other Aliases: Z178A11.18

Other Designations: Transcription factor **TEOSINTE BRANCHED 1**; p-umc1082; **teosinte branched protein 1**; **teosinte branched1** protein; umc1082

ID: 542361

☐ [LOC100286105](#)

2. **teosinte branched1** protein [*Zea mays*]

ID: 100286105

☐ [LOC100282571](#)

3. **teosinte branched1** protein [*Zea mays*]

ID: 100282571

[Display Settings:](#) ☐ Summary, Sorted by Relevance

[Send to:](#) ☐

NCBI

http://www.ncbi.nlm.nih.gov/gene/100286105

LOC100286105 teosinte branched...

GeneRIFs: Gene References Into Functions [What's a GeneRIF?](#)

Submit: [New GeneRIF](#) [Correction](#)

General protein information

Preferred Names
teosinte branched1 protein

Names
teosinte branched1 protein

NCBI Reference Sequences (RefSeq)

RefSeqs maintained independently of Annotated Genomes

These reference sequences exist independently of genome builds. [Explain](#)

mRNA and Protein(s)

1. [NM_001158993.1](#) → [NP_001132405.1](#) teosinte branched1 protein

Status: **PROVISIONAL**

Source sequence(s) [EU975636](#)

UniProtKB/TrEMBL [P8UEU1](#)

Conserved Domains (1) [summary](#)

	pfam03634 Location:66 – 127 Blast Score: 275	TCP; TCP family transcription factor
--	--	--------------------------------------

Related Sequences

NCBI

http://www.ncbi.nlm.nih.gov/nuccore/EU975636

Zea mays clone 494952 teosinte b...

Display Settings: GenBank Send:

Zea mays clone 494952 teosinte branched1 protein mRNA, complete cds

GenBank: EU975636.1

[FASTA](#) [Graphics](#)

[Details](#)

LOCUS EU975636 1024 bp mRNA linear PLN 10-DEC-2008

DEFINITION Zea mays clone 494952 teosinte branched1 protein mRNA, complete cds.

ACCESSION EU975636

VERSION EU975636.1 GI:195656572

KEYWORDS FLI_CDNA.

SOURCE Zea mays

ORGANISM [Zea mays](#)

Eukaryota; Viridiplantae; Streptophyta; Embryophyta; Tracheophyta; Spermatophyta; Magnoliophyta; Liliopsida; Poales; Poaceae; PACMAD clade; Panicoideae; Andropogoneae; Zea.

REFERENCE 1 (bases 1 to 1024)

AUTHORS Alexandrov,N.N., Brover,V.V., Freidin,S., Troukhan,M.E., Tatarinova,T.V., Zhang,H., Swaller,T.J., Lu,Y.-P., Bouck,J., Flavell,R.B. and Feldmann,K.A.

TITLE Insights into corn genes derived from large-scale cDNA sequencing

JOURNAL Plant Mol. Biol. 69 (1-2), 179-194 (2009)

PUBMED [18937034](#)

REFERENCE 2 (bases 1 to 1024)

AUTHORS Alexandrov,N.N., Brover,V.V., Freidin,S., Troukhan,M.E., Tatarinova,T.V., Zhang,H., Swaller,T.J., Lu,Y.-P., Bouck,J., Flavell,R.B. and Feldmann,K.A.

TITLE Direct Submission

JOURNAL Submitted (04-AUG-2008) Ceres, Inc., 1535 Rancho Conejo Blvd., Thousand Oaks, CA 91320, USA

FEATURES

source Location/Qualifiers

1..1024

/organism="Zea mays"

/mol_type="mRNA"

/db_xref="taxon:[4577](#)"

/clone="494952"

CDS 48..854

/note="similar to NP_001063149.1"

/codon_start=1

/product="teosinte branched1 protein"

/protein_id="[ACG47754.1](#)"

/db_xref="GI:195656573"

/translation="MSSSDGYGGQIFPADMSSFRRHQDTLEAVFHHPPPTTTTGLLRNDGSSPVVDGGGHAAPRRPFRTDRHSKIRTAQGVDRDRRLSVGVAREFFALQDRLGFDKASKTVNWLTSQSKPAIDRLVDAADDPAAVAASGRRPTVVRGRGEGSSSSTCCCLTDSREAAEEATGNGRSRGPDGPPAALLEGHGGCGELGWIMSGAPTAAVATTTTTPQQPDGHEYYYQVCLQLEEMRCSNDEGETTPGDFLYGMQTRDRS"

ORIGIN

1 aaccctattg cagctgcttt tgtttcctta tgccctcgac cgcgcgatg tcgtgctctt

61 gggacgggta cggcgggcag atottcccg cagacatgtc gtcgttcac cgcaccagg

121 aacacctgga ggcggtgttc caaccaccgc gcgctgagac gacgacgac acggggctgc

181 tgaggaacga cgggtcatca ccggtggtgg acgacggcgg cggccacgcc gcaccgcgaa

241 ggcggccggt ccggacggac cgcacagca agatccgcac ggcgcagggc gtgcgcgacc

301 gcggtgagcg gctgtcggtc ggggtcgcgc gcgagttctt cgcgctcgac gaccgcctcg

361 gctgagcaga cggcgggcag aagcggcagt gctgctcgac cagtcgagac cggcggcgtg

NCBI

NCBI Resources ▾ How to ▾

Nucleotide Limits Advanced

[Display Settings:](#) ☒ FASTA

Zea mays clone 494952 teosinte branched1 protein mRNA, complete cds

GenBank: EU975636.1

[GenBank](#) [Graphics](#)

>gi|195656572|gb|EU975636.1| Zea mays clone 494952 teosinte branched1 protein mRNA, complete cds

```
AACCCCTATTGCAGCTGCTTTTGTTCCTTATGCCCTCGACCGCCGCGATGTCGTCGTCTTGGGACGGGTA
CGGCGGGCAGATCTTCCCGCAGACATGTCGTCGTTCCACCGCCACCAGGACACCCTGGAGGCGGTGTTT
CACCACCCGCCGCTGAGACGACGACGACGACGGGGCTGCTGAGGAACGACGGGTCATCACCAGGTGGTGG
ACGACGGCGGGCGGCCACGCCGACCCGGAAGGCGGGCGGTTCCGGACGGACCGCCACAGCAAGATCCGCAC
GGCGCAGGGCGTGCGCGACCGCCGGATGCGGCTGTCGGTTCGGGGTCGCGCGCGAGTTCTTCGCGCTGCAG
GACCGCCTCGGCTTCGACAAGGCCAGCAAGACGGTGAACCTGGCTCCTCACCAGTCCAAGCCGGCCATCG
ACCGCCTCGTCGACGCCGCGCCGCCGCCGACGACCCCGCGGCGTAGCAGCCTCAGGAGGCCGACGACC
GACGGTGGTGAGGGGCAGAGGCGAGGGCAGCTCCTCGAGCACTTGCTGCTGCTTGACGGACTCGAGAGAG
GCCGCCGAGGAGGCGACGGGGAACGGGAGAAGCAGAGGCGGGCCCTGACGACGGGCCACCGGCAGCGCTTC
TGGAAGGACACGGCGGCTGCGGCGAGCTGGGCTGGATCATGTCTGGGAGCGCCACAGCAGCGGTGGCAAC
GACGACGACGACGACGCCGACGAGCCGGACGGGACGAGTACTACTACCAGTATTGCCTGCAGCTCGAG
GAGATGATGCGATGCAGCAACGACGAAGGAGAAACAACGCCAGGTGATTTCTTGATGGTATGCAGACGC
GTGATAGGTCTTGAGCTCTCTAAACGCGCGTGAGAGGATTTCCATTACGTCTGAGATTATGCTGATCT
GCTGCCATGATGATCCATTACTACTGCATATCTATCTAGTACATATAAATCTCAACTGGTTCGATCTTTA
TCTCATCAATCAAGATCCAAAGGCCCAAAAAAAAAAAAAAAAAAAAAA
```

Overview

- Navigating NCBI to obtain sequences
- **Using BLAST for sequence alignment**
- Using other programs for specialized sequence alignment
- Next generation sequence alignment programs

BLAST

- BLAST = Basic Local Alignment Search Tool

BLAST

- BLAST = Basic Local Alignment Search Tool
- Used for searching a query sequence or sequences against a database or subject sequence

BLAST

- BLAST = Basic Local Alignment Search Tool
- Used for searching a query sequence or sequences against a database or subject sequence
- Uses a local alignment, which searches for regions in the query that locally align to subject sequences

BLAST

- BLAST = Basic Local Alignment Search Tool
- Used for searching a query sequence or sequences against a database or subject sequence
- Uses a local alignment, which searches for regions in the query that locally align to subject sequences
- BLAST uses word based heuristics to approximate the Smith-Waterman algorithm to find the near-optimal local alignments quickly, thus gaining speed over sensitivity.

BLAST

- BLAST = Basic Local Alignment Search Tool
- Used for searching a query sequence or sequences against a database or subject sequence
- Uses a local alignment, which searches for regions in the query that locally align to subject sequences
- BLAST uses word based heuristics to approximate the Smith-Waterman algorithm to find the near-optimal local alignments quickly, thus gaining speed over sensitivity.
- Two main versions: NCBI Blast and WU-BLAST

Local vs Global Alignments

Global FTFTALILLAVAV
 F--TAL-LLA-AV

Local FTFTALILL-AVAV
 --FTAL-LLAAV--

For sequences that are divergent, the optimal global alignment introduces gaps that can hide biologically relevant information such as motifs.

BLAST

- The process starts by searching the sequences for exact matches of small fixed length strings from the query called 'words'.

BLAST

- The process starts by searching the sequences for exact matches of small fixed length strings from the query called 'words'.
- These matches are the seeds for the local alignments.

BLAST

- The process starts by searching the sequences for exact matches of small fixed length strings from the query called 'words'.
- These matches are the seeds for the local alignments.
- The next step is to extend the seed by aligning the sequence until a gap is found – mismatches are inserted here.

BLAST

- The process starts by searching the sequences for exact matches of small fixed length strings from the query called 'words'.
- These matches are the seeds for the local alignments.
- The next step is to extend the seed by aligning the sequence until a gap is found – mismatches are inserted here.
- A gapped alignment is then performed using a modified Smith-Waterman algorithm – indels are added here.

BLAST

- The process starts by searching the sequences for exact matches of small fixed length strings from the query called 'words'.
- These matches are the seeds for the local alignments.
- The next step is to extend the seed by aligning the sequence until a gap is found – mismatches are inserted here.
- A gapped alignment is then performed using a modified Smith-Waterman algorithm – indels are added here.
- Only results scoring above a threshold (expect value or e value) are reported back to the user

BLAST Programs

- blastn - query DNA, subject DNA
- blastp- query Protein, subject Protein
- blastx - query DNA (6 frame translation)
subject Protein
- tblastn - query protein - subject DNA (6 frame translation) - slow
- tblastx - query DNA (6 frame translation),
subject DNA (6 frame translation) - very slow

Which Program Should You Use

What database do you have and how sensitive does your search has to be?

blastn, blastp – good for identifying sequences that are already in a database, finding local regions of similarity in closely related organisms.

Which Program Should You Use

What database do you have and how sensitive does your search has to be?

blastn, blastp – good for identifying sequences that are already in a database, finding local regions of similarity in closely related organisms.

blastx - Use when you have a nucleotide sequence with an unknown reading frame and/or sequencing errors that would lead to frame shifts or coding errors such as ESTs. More sensitive.

Which Program Should You Use

What database do you have and how sensitive does your search has to be?

blastn, blastp – good for identifying sequences that are already in a database, finding local regions of similarity in closely related organisms.

blastx - Use when you have a nucleotide sequence with an unknown reading frame and/or sequencing errors that would lead to frame shifts or coding errors such as ESTs. More sensitive.

tblastn- useful for finding homologs in sequence where the frame is unknown or sequencing errors are likely to be present such as ESTs and draft sequence.

Which Program Should You Use

What database do you have and how sensitive does your search has to be?

blastn, blastp – good for identifying sequences that are already in a database, finding local regions of similarity in closely related organisms.

blastx - Use when you have a nucleotide sequence with an unknown reading frame and/or sequencing errors that would lead to frame shifts or coding errors such as ESTs. More sensitive.

tblastn- useful for finding homologs in sequence where the frame is unknown or sequencing errors are likely to be present such as ESTs and draft sequence.

tblastx - used to detect novel ORFs/exons. Very Slow, use as the last resort.

Performing BLAST

The screenshot shows the NCBI BLAST website in a web browser. The address bar displays <http://blast.ncbi.nlm.nih.gov/Blast.cgi>. The page header includes the BLAST logo and navigation links: Home, Recent Results, Saved Strategies, and Help. Below the header, a banner states "BLAST finds regions of similarity between biological sequences." with a "more..." link. A "New" alert box promotes the COBALT Multiple Alignment Tool. The "BLAST Assembled RefSeq Genomes" section lists various species for selection, including Human, Mouse, Rat, Arabidopsis thaliana, Oryza sativa, Bos taurus, Danio rerio, Drosophila melanogaster, Gallus gallus, Pan troglodytes, Microbes, and Apis mellifera. The "Basic BLAST" section provides instructions on choosing a BLAST program to run. The "nucleotide blast" option is circled in red.

BLAST® Basic Local Alignment Search Tool

Home Recent Results Saved Strategies Help

NCBI/ BLAST Home

BLAST finds regions of similarity between biological sequences. [more...](#)

New Aligning Multiple Protein Sequences? Try the [COBALT Multiple Alignment Tool](#). [Go](#)

BLAST Assembled RefSeq Genomes

Choose a species genome to search, or [list all genomic BLAST databases](#).

- [Human](#)
- [Mouse](#)
- [Rat](#)
- [Arabidopsis thaliana](#)
- [Oryza sativa](#)
- [Bos taurus](#)
- [Danio rerio](#)
- [Drosophila melanogaster](#)
- [Gallus gallus](#)
- [Pan troglodytes](#)
- [Microbes](#)
- [Apis mellifera](#)

Basic BLAST

Choose a BLAST program to run.

- [nucleotide blast](#) Search a **nucleotide** database using a **nucleotide** query
Algorithms: blastn, megablast, discontinuous megablast
- [protein blast](#) Search **protein** database using a **protein** query
Algorithms: blastp, psi-blast, phi-blast
- [blastx](#) Search **protein** database using a **translated nucleotide** query
- [tblastn](#) Search **translated nucleotide** database using a **protein** query
- [tblastx](#) Search **translated nucleotide** database using a **translated nucleotide** query

Performing BLAST

http://blast.ncbi.nlm.nih.gov/Blast.cgi?PROGRAM=blastn&BLAST_PROGRAMS=megaBlast&PAGE_TYPE=BlastSearch

Nucleotide BLAST: Search nucleoti...

BLAST® Basic Local Alignment Search Tool

Home Recent Results Saved Strategies Help

NCBI/ BLAST/ blastn suite

blastn blastp blastx tblastn tblastx

BLASTN programs search nucleotide databases using a nucleotide query. [more...](#)

Enter Query Sequence

Enter accession number(s), gi(s), or FASTA sequence(s) [?](#) [Clear](#) [Query subrange](#) [?](#)

>gi|195656572|gb|EU975636.1| Zea mays clone 494952 teosinte
branched1 protein mRNA, complete cds
AACCTATTGCAGCTGCTTTTCCTTATGCGCTCGACCGCGGATGTCGTCGCTTGGGACGG
GTA
CGCGGGCAGATCTTCCCGCAGACATGTCGTCGTTCCACCGCCACGAGACACCTGGAGGCGGTG
TAA

From
To

Or, upload file [Browse...](#) [?](#)

Job Title
Enter a descriptive title for your BLAST search [?](#)

☐ Align two or more sequences [?](#)

Choose Search Set

Database ☐ Human genomic + transcript ☐ Mouse genomic + transcript ☒ Others (nr etc.):
☒ Nucleotide collection (nr/nt) [?](#)

Organism [Optional](#) ☐ Exclude [+](#)
Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown. [?](#)

Exclude [Optional](#) ☐ Models (XM/XP) ☐ Uncultured/environmental sample sequences

Entrez Query [Optional](#)
Enter an Entrez query to limit search [?](#)

Program Selection

Optimize for ☐ Highly similar sequences (megablast)
☐ More dissimilar sequences (discontiguous megablast)
☒ Somewhat similar sequences (blastn) [?](#)
Choose a BLAST algorithm [?](#)

Performing BLAST

http://blast.ncbi.nlm.nih.gov/Blast.cgi?PROGRAM=blastn&BLAST_PROGRAMS=megaBlast&PAGE_TYPE=BlastSearch&SHOW_I

Nucleotide BLAST: Search nucleoti...

BLAST® Basic Local Alignment Search Tool

Home Recent Results Saved Strategies Help

NCBI/ BLAST/ blastn suite

blastn blastp blastx tblastn tblastx

BLASTN programs search nucleotide databases using a nucleotide query. [more...](#)

Enter Query Sequence

Enter accession number(s), gi(s), or FASTA sequence(s)

>gi|195656572|gb|EU975636.1| Zea mays clone 494952 teosinte
branched1 protein mRNA, complete cds
AACCTTATTCAGCTGCTTTTGTACCTATGCCCTCGACCGCCGGGATGCTCGTCTGGGACGG
GTA
CGGCGGGCAGATCTTCCCGCAGACATGTCGTCGTTCCACCGCCACCAAGGACCCCTGGAGGCGGTG
TTC

Clear Query subrange

From

To

Or, upload file

Browse...

Job Title

gi|195656572|gb|EU975636.1| Zea mays clone...

Enter a descriptive title for your BLAST search

☐ Align two or more sequences

Choose Search Set

Database

☐ Human genomic + transcript ☐ Mouse genomic + transcript ☒ Others (nr etc.):

Nucleotide collection (nr/nt)

Organism

Optional

Exclude

Optional

Entrez Query

Optional

Program Selection

Optimize for

Exclude +

Exclude taxa will be shown.

Exclude sequences

Oryza s

Oryza sativa (taxid:4530)

Oryza sativa L. (taxid:4530)

Oryza sativa Japonica Group (taxid:39947)

Oryza sativa Indica Group (taxid:39946)

Oryza sativa subsp. indica Kato (taxid:39946)

Oryza sp. (taxid:52841)

Oryza glaberrima x Oryza sativa (taxid:441961)

Oryza subulata Nees (taxid:110453)

Oryza sativa x Oryza glaberrima (taxid:551113)

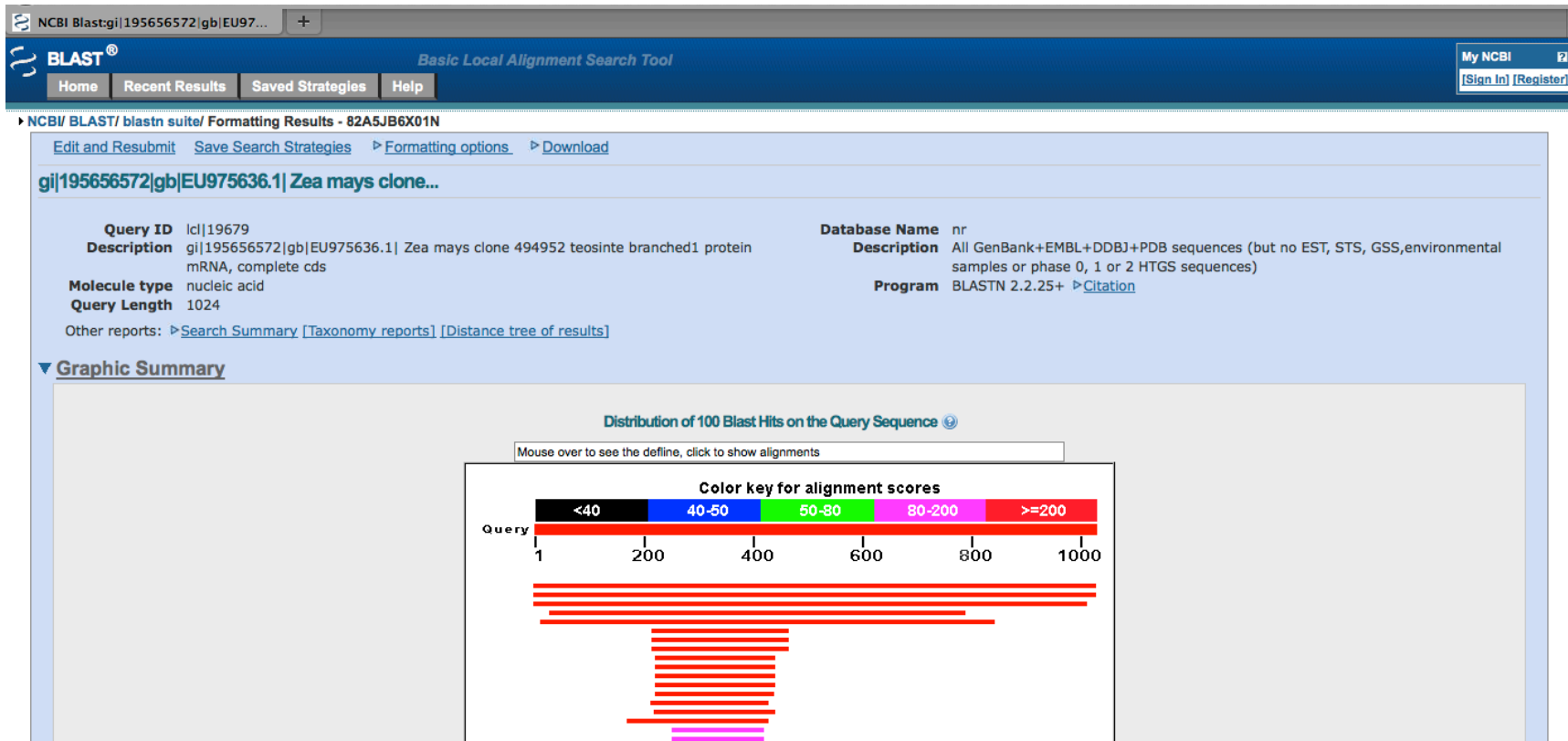
Oryza rufipogon x Oryza sativa (taxid:1077957)

more dissimilar sequences (discontiguous megablast)

☒ Somewhat similar sequences (blastn)

Choose a BLAST algorithm

Output from BLAST



Output from BLAST

▼ Descriptions							
Legend for links to other resources: U UniGene E GEO G Gene S Structure M Map Viewer P PubChem BioAssay							
Sequences producing significant alignments:							
Accession	Description	Max score	Total score	Query coverage	E value	Max ident	Links
NM_001158993.1	Zea mays teosinte branched1 protein (LOC100286105), mRNA >gb EU975636	1847	1847	100%	0.0	100%	UG
EU975768.1	Zea mays clone 499966 teosinte branched1 protein mRNA, complete cds	1833	1833	100%	0.0	99%	UG
BT063457.1	Zea mays full-length cDNA clone ZM_BFc0065H12 mRNA, complete cds	1793	1793	98%	0.0	99%	UG
XM_002444342.1	Sorghum bicolor hypothetical protein, mRNA	628	628	73%	2e-176	77%	G
NM_001155479.1	Zea mays teosinte branched1 protein (LOC100282571), mRNA >gb EU961496	600	600	80%	9e-168	75%	UG
NM_001069684.2	Oryza sativa Japonica Group Os09g0410500 (Os09g0410500) mRNA, complet	315	315	24%	5e-82	88%	UG
EU702407.1	Oryza sativa Japonica Group putative TCP transcription factor (DP1) mRNA, co	315	315	24%	5e-82	88%	U
AP005093.3	Oryza sativa Japonica Group genomic DNA, chromosome 9, BAC clone:OJ1294	315	315	24%	5e-82	88%	
HQ858764.1	Zea mays clone UT1680 TCP transcription factor mRNA, partial cds	289	289	21%	2e-74	89%	G
BT061081.1	Zea mays full-length cDNA clone ZM_BFb0095I10 mRNA, complete cds	289	289	21%	2e-74	89%	UG
NM_001158146.1	Zea mays teosinte-branched one (LOC100285252), mRNA >gb EU971416.1	289	289	21%	2e-74	89%	UG
NM_001136610.1	Zea mays TCP transcription factor (LOC100191175), mRNA >gb BT033200.1	289	289	21%	2e-74	89%	UG
XM_002460157.1	Sorghum bicolor hypothetical protein, mRNA	288	288	21%	7e-74	89%	G
AP003868.3	Oryza sativa Japonica Group genomic DNA, chromosome 8, BAC clone:OJ1111	280	280	20%	1e-71	89%	
NM_001196163.1	Zea mays hypothetical protein LOC100501452 (LOC100501452), mRNA >gb E	277	277	21%	1e-70	87%	G
AC157320.2	Zea mays clone ZMMBBb-7C14, complete sequence	259	259	25%	4e-65	81%	
NM_001057563.1	Oryza sativa Japonica Group Os03g0706500 (Os03g0706500) mRNA, complet	161	161	16%	7e-36	81%	UG
AC091775.10	Oryza sativa chromosome 3 BAC OSJNBa0004G17 genomic sequence, comple	161	161	16%	7e-36	81%	
AY286002.1	Oryza sativa (indica cultivar-group) teosinte-branching 1 (TB1) gene, complet	161	161	16%	7e-36	81%	
AY043215.1	Oryza sativa teosinte branched1 protein (tb1) mRNA, complete cds	161	161	16%	7e-36	81%	U E M
AK107083.1	Oryza sativa Japonica Group cDNA clone:002-121-G01, full insert sequence	161	161	16%	7e-36	81%	U E G
AF322143.1	Oryza sativa teosinte branched1 protein (tb1) gene, partial cds	161	161	16%	7e-36	81%	
AB088343.1	Oryza sativa Japonica Group OsTB1 gene for teosinte branched1 protein, com	161	161	16%	7e-36	81%	G
AC243260.1	Panicum virgatum clone PV_ABa103-K10, complete sequence	158	158	15%	9e-35	81%	
EF694162.2	Pennisetum glaucum isolate Tb103 teosinte-branched1 (Tb1) gene, complete	158	158	15%	9e-35	81%	

Output from BLAST

```
>\[ref|XM\_002444342.1\] \[G\] Sorghum bicolor hypothetical protein, mRNA
Length=774

GENE ID: 8075661 SORBIDRAFT 07g021140 | hypothetical protein [Sorghum bicolor]
(10 or fewer PubMed links)

Score = 628 bits (696), Expect = 2e-176
Identities = 627/808 (78%), Gaps = 105/808 (13%)
Strand=Plus/Plus

Query 30 ATGCCCTCGACCGCCGCGATGTCGTCGCTTTGGGACGGGTACGGCGGGCAGATCTTCCCC 89
Sbjct 1 ATGCCCTCGACCGC---GATGTC-----TTGGGACGGGTACGGCGGGCAGATCTTCCCC 51
Query 90 GCAGACATGTCGTCGTTCCACCGCCACAGGACACCTGGAGGCGGTGTTCACCAACC- 148
Sbjct 52 GCCGACATGTCGTCGTTCCACC---ACCAGGACACCTGGAGGCGGTGTTCGGCAGCCT 108
Query 149 -----GC-CGCCTG--AGACGACGACGACGACGGG-----CTGCTGAGG 185
Sbjct 109 GAGACGACGGCGCCCTGCAGGCGCCGCGCAGCAGCGGGGAGATGGAGCTGCTGTGAGG 168
Query 186 AACGACGGGTATC---ACCGGTGGTGGACGACGGCGGCG--GCCACG--CCGCACCGCGA 239
Sbjct 169 AACG---GGTCGCCCTACCGGTGGTGGATGCCGGCGTCCATGCCGCGCCGCGACCGCG 225
Query 240 AGGCGGCGGTTCCGGACGGACCGCCACAGCAAGATCCGCACGGCGCAGGGCGTGCAGC 299
Sbjct 226 AAGCGGCGGTTTCAGGACGGATCGGCACAGCAAGATCCGCACGGCGCAGGGCGTGCAGC 285
Query 300 CGCCGGATGCGGCTGTTCGGTCGGGGTCGCGCGCGAGTTCTTCGCGCTGCAGGACCGCTC 359
Sbjct 286 CGCCGGATGCGGCTGTTCGGTCGGGGTCGCGCGAGAGTTCTTCGCGCTGCAGGACCGCTC 345
Query 360 GGCTTCGACAAGGCCAGCAAGACGGTGAAGTGGCTCCTCACCAGTCCAAGCCGGCCATC 419
Sbjct 346 GGGTTCGACAAGGCCAGCAAGACGGTGAAGTGGCTCCTCACCAGTCCAAGCCGGCCATC 405
Query 420 GACCGCCTCGTCGACgcccgcgcgcgcgcgcgACGACCCCGCGGCCGTAGCAGCCTCAGGA 479
Sbjct 406 GACCGCCTCGTCGACGCCGCCG-AGCCGGCG-----GTGGCTCT---AGTCTCAGGA 453
Query 480 GGCCGACGACCGACGGTGGTGAAGGGGAGAGGGGAGGGGAGCTCCTCGAGCACTTGCTGC 539
Sbjct 454 G---GACCACCGACGGTGGTGAAGGGGAGAGGGGAGGGGAACCTCAAGCACT---TGC 507
Query 540 TGCTTGAC---GGACTCGAGAGAGGCCCGGAGGAGGCGACGGGGAACGGGAGAAGCAGA 596
Sbjct 508 TGTTCGACGGTGGACTCGAG-----GGAGGAGGCGACGGAGAAGGCAAGGAAGCAGA 558
Query 597 GCGGGC-----CC-----TGACGACGGGCCACCGGCAGCGCTTCTGGAAGGA 638
Sbjct 559 GCGGGCGGCGGCGCGGTACCGTGGTCTGATGGGCCACCG---GCGCTCATGGAAGAA 615
Query 639 CACGGCGGCTGCGCGAGCTGGGCTGGATCATGTCGGGAGCGCCACAGCAGCGGTGGCA 698
Sbjct 616 CA---CGCCCGCGGTGAGCTGGGCTGGATCATGACG---GAGGCCACAGCGGCAGCGG-- 667
Query 699 acgacgacgacgacgacgCCGACGACCGGACGGGACAGTACTACTACCACTATTGC 758
Sbjct 668 -CGGCAGCAACGGCGCAGCCGACGAGATGGACGGGTGGAGTACTACTACCACTATTGC 726
Query 759 CTGCAGCTCGAGGAGATGATGCGATGCA 786
Sbjct 727 CTGCAGCTCGAGGAGATGATGAGATGCA 754
```

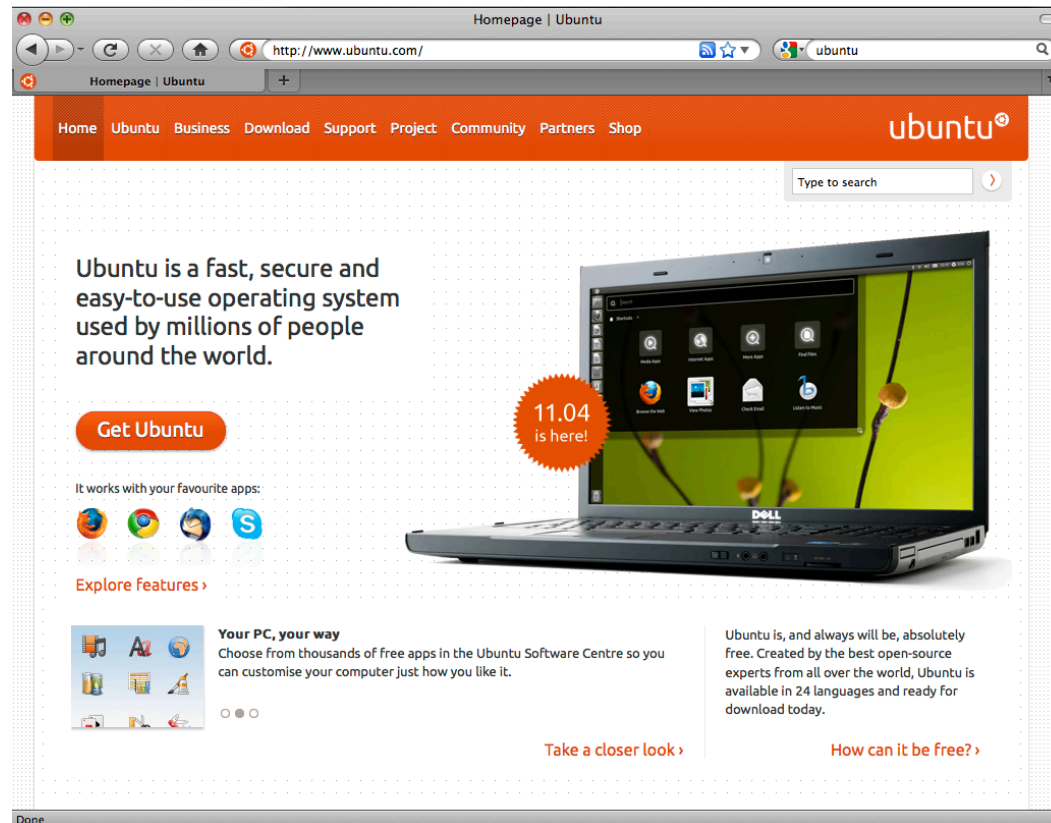
BLAST

- For additional information on BLAST go to http://blast.ncbi.nlm.nih.gov/Blast.cgi?CMD=Web&PAGE_TYPE=BlastDocs

BLAST

- For additional information on BLAST go to http://blast.ncbi.nlm.nih.gov/Blast.cgi?CMD=Web&PAGE_TYPE=BlastDocs
- Web based BLAST interfaces are designed for low throughput searching
- Personalized databases can be generated and subsequent BLASTs performed using the command line
- Using the command line BLAST, multiple sequences can be aligned simultaneously

Ubuntu



- For PC users I highly recommend getting Ubuntu
- This can be run as a virtual machine on your PC through programs such as VirtualBox <https://www.virtualbox.org/>

Overview

- Navigating NCBI to obtain sequences
- Using BLAST for sequence alignment
- Using other programs for specialized sequence alignment
- Next generation sequence alignment programs

EST and cDNA Alignment

EXONERATE

a generic tool for sequence alignment

- Exonerate is a generic tool for pairwise sequence comparison (<http://www.ebi.ac.uk/~guy/exonerate/>) and comes grid ready with the ability to chunk files directly through exonerate

EST and cDNA Alignment

EXONERATE

a generic tool for sequence alignment

- Exonerate is a generic tool for pairwise sequence comparison (<http://www.ebi.ac.uk/~guy/exonerate/>) and comes grid ready with the ability to chunk files directly through exonerate
- Exonerate can be run with many different models for gapped and ungapped alignments

EST and cDNA Alignment

EXONERATE

a generic tool for sequence alignment

- Exonerate is a generic tool for pairwise sequence comparison (<http://www.ebi.ac.uk/~guy/exonerate/>) and comes grid ready with the ability to chunk files directly through exonerate
- Exonerate can be run with many different models for gapped and ungapped alignments
- The est2genom model can be used for alignment of EST sequences to genomic sequence and the cdna2genome model can be used to align full length cDNAs that can be flanked by UTRs

EST and cDNA Alignment

EXONERATE

a generic tool for sequence alignment

- Exonerate is a generic tool for pairwise sequence comparison (<http://www.ebi.ac.uk/~guy/exonerate/>) and comes grid ready with the ability to chunk files directly through exonerate
- Exonerate can be run with many different models for gapped and ungapped alignments
- The est2genom model can be used for alignment of EST sequences to genomic sequence and the cdna2genome model can be used to align full length cDNAs that can be flanked by UTRs
- For additional models see the man page online at <http://www.ebi.ac.uk/~guy/exonerate/exonerate.man.html> or on the command line with the -h option
 - exonerate -h

EST and cDNA Alignment

EXONERATE

a generic tool for sequence alignment

- Command line example
 - `exonerate --query est_sequences.fasta --querytype DNA --target genome_sequence.fasta --targettype DNA --model est2genome --showalignment no --showvulgar no --showtargetgff no --ryo "%qi\t%ti\t%qab\t%qae\t%tab\t%tae\n" --minintron 10 --maxintron 1500 >exonerate_alignment.txt`

EST and cDNA Alignment

EXONERATE

a generic tool for sequence alignment

- Command line example
 - `exonerate --query est_sequences.fasta --querytype DNA --target genome_sequence.fasta --targettype DNA --model est2genome --showalignment no --showvulgar no --showtargetgff no --ryo "%qi\t%ti\t%qab\t%qae\t%tab\t%tae\n" --minintron 10 --maxintron 1500 >exonerate_alignment.txt`
 - The output from this command will be a tab delimited file with the following columns for each alignment: `query_id`, `target_id`, `query_start`, `query_end`, `target_start`, `target_end`
 - Note: using this output format the positions are in interbase coordinates

A C G T
0 1 2 3 4

- For complete list of options see the man page online at <http://www.ebi.ac.uk/~guy/exonerate/exonerate.man.html> or on the command line with the `-h` option

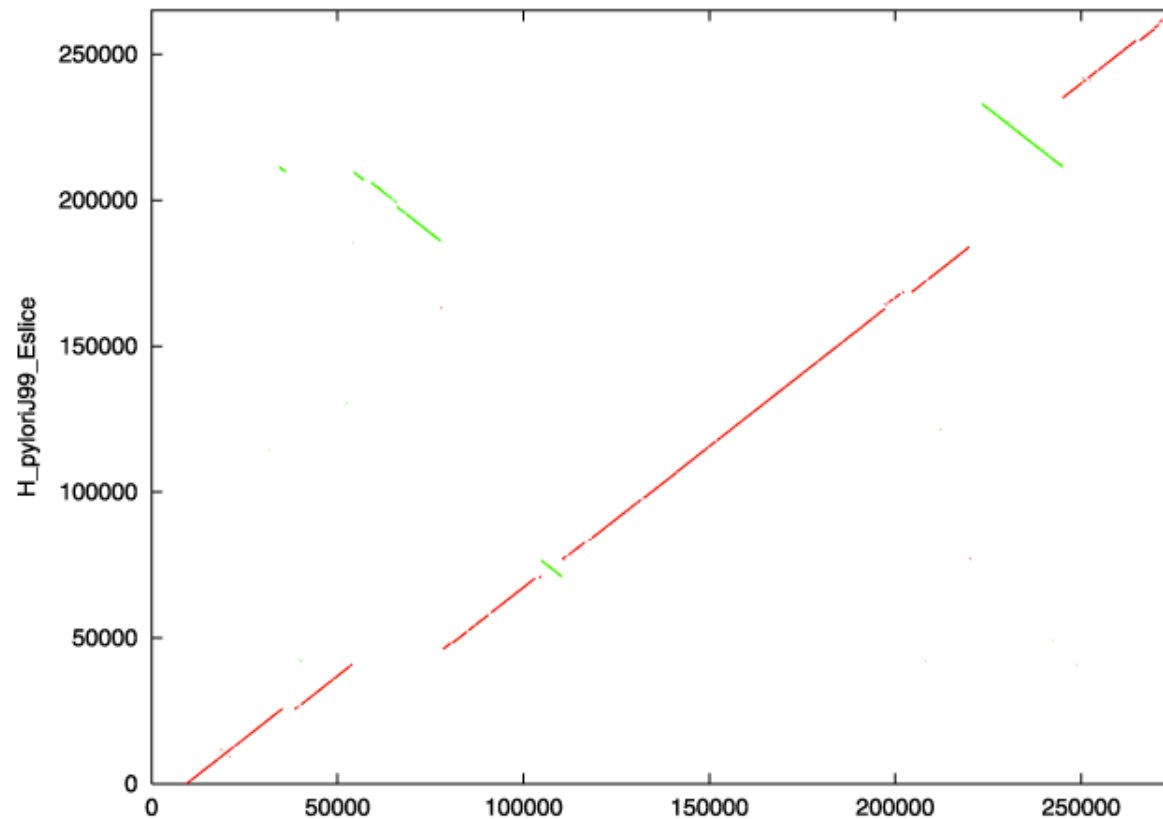
Whole Genome Alignment

- MUMmer – designed for rapid alignment of entire genomes
- NUCmer – designed for alignment of contigs to another set of contigs or a genome
- PROmer – designed for alignment of species too divergent for DNA alignment using six-frame translation of both input sequences
- <http://mummer.sourceforge.net/>

Whole Genome Alignment

- mummer can handle multiple reference and query sequences
- Command line usage
 - `mummer [options] <reference-file> <query-files>`
 - `mummerplot [options] <match file>`
- Example command line
 - `mummer -mum -b -c genotype1.fasta genotype2.fasta > output.mums`
 - -mum finds all maximal unique matches, -b will compute forward and reverse complement matches, -c reports all match positions relative to the forward strand
 - `mummerplot -x "[0,275287]" -y "[0,265111]" --png output.mums`
 - -x sets the x-axis range in the plot, -y sets the y-axis range in the plot, --png outputs the plot in .png format
- Additional options are available in the manual or by using the -h option for both mummer and mummerplot
- Examples for NUCmer and PROmer are available at <http://mummer.sourceforge.net/examples/>

Whole Genome Alignment



Forward MUMs are red and reverse MUMs are green. Dots on a line with a slope = 1 are unchanged, and dots on a line with a slope = -1 are inverted.

<http://mummer.sourceforge.net/examples/>

Short Sequence Alignment

- Vmatch is a tool that is ideal for aligning short sequences such as probes, primers, and SNPs with short context sequence
- This program requires a license
- <http://www.vmatch.de/>

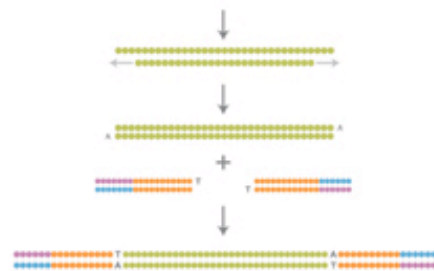
Overview

- Navigating NCBI to obtain sequences
- Using BLAST for sequence alignment
- Using other programs for specialized sequence alignment
- **Next generation sequence alignment programs**

High Throughput Sequencing Platforms

- Illumina HiSeq 1000 and HiSeq 2000
- Illumina Genome Analyzer Ix
- Life Sciences/Roche 454 pyrosequencing
- Pacific Biosciences
- Ion Torrent

High Throughput Sequencing



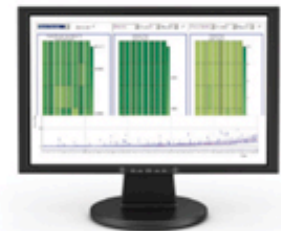
Library Preparation
~2 h [15 min hands-on (Nextera)]
< 6 h [< 3 h hands-on (TruSeq)]



Cluster Generation
~5 h (<10 min hands-on)



Sequencing by Synthesis
~1.5 to 8.5 days



CASAVA
2 days (30 min hands-on)

- HiSeq 2000
- Highly parallel sequencing by synthesis
- Single and paired-end reads between 50 bp and 100 bp
- 187 million single end or 374 million paired-end reads per lane
- High error rate in the 3' end

Sequence File Format

```
@HWUSI-EAS1789_0001:1:1:11120:1081#0/1 Read Name  
ACNACAGCTATGACCTCTAGGAATCTTTGTAAAGGCTTCGTAGTGAATCCCTGGCATTACCTTGGGAGTGAG Sequence  
+HWUSI-EAS1789_0001:1:1:11120:1081#0/1  
a_Baaeeeeefffaaaaaaaaaaaaaaaaaaaaaaaaedfdffffffffdffffffffffdfdfdd]fcfbf Quality
```

```
@HWUSI-EAS1789_0001:1:1:11154:1081#0/1  
GCNCTACAGCGGTCTTACTCGCGACTACAAATCTTGGGTTCTCCATAGATACTCTACAACCTTCGTTCTGAAATTA  
+HWUSI-EAS1789_0001:1:1:11154:1081#0/1  
aaB]_eeeeeHHHHHHHHHHHghhhhhhhhhhhhhhghgfhhehhhgfhggghghghhhhhhfhhgghfgghgg  
@HWUSI-EAS1789_0001:1:1:12819:1080#0/1  
TCNCCCTGCGCGAGCGGTACCAAATCGAGGCCAACTCTGAATACTAGATATGACCCEAAAATAACAGGGGTCAAG  
+HWUSI-EAS1789_0001:1:1:12819:1080#0/1  
bbBbbdeeeehhhhhhhghqhhhhhhghfhhhhhhhhhhhhhheqhchccqhghghhhhehh`bdf fdadgdq
```

Quality scores are in ASCII characters that are converted to Phred scores. These scores provide a likelihood that the base was called incorrectly.

- 10 – 1 in 10 chance the base call is incorrect
- 20 – 1 in 100 chance the base call is incorrect
- 30 – 1 in 1000 chance the base call is incorrect

Suggest checking the quality of the reads prior to performing sequence alignments

Read Quality with the FASTX-Toolkit



Introduction

The FASTX-Toolkit is a collection of command line tools for Short-Reads FASTA/FASTQ files preprocessing.

Next-Generation sequencing machines usually produce FASTA or FASTQ files, containing multiple short-reads sequences (possibly with quality information).

The main processing of such FASTA/FASTQ files is mapping (aka aligning) the sequences to reference genomes or other databases using specialized programs. Example of such mapping programs are: [Blat](#), [SHRiMP](#), [LastZ](#), [MAQ](#) and many many others.

However,

It is sometimes more productive to preprocess the FASTA/FASTQ files before mapping the sequences to the genome - manipulating the sequences to produce better mapping results.

The FASTX-Toolkit tools perform some of these preprocessing tasks.

http://hannonlab.cshl.edu/fastx_toolkit/

Read Quality with the FASTX-Toolkit

FASTX Statistics

>fastx_quality_stats -i sequence_file.fastq -o stats.txt

```
$ fastx_quality_stats -h
usage: fastx_quality_stats [-h] [-i INFILE] [-o OUTFILE]

version 0.0.6 (C) 2008 by Assaf Gordon (gordon@cshl.edu)
[-h] = This helpful help screen.
[-i INFILE] = FASTA/Q input file. default is STDIN.
               If FASTA file is given, only nucleotides
               distribution is calculated (there's no quality info).
[-o OUTFILE] = TEXT output file. default is STDOUT.

The output TEXT file will have the following fields (one row per column):
column = column number (1 to 36 for a 36-cycles read solexa file)
count  = number of bases found in this column.
min    = Lowest quality score value found in this column.
max    = Highest quality score value found in this column.
sum    = Sum of quality score values for this column.
mean   = Mean quality score value for this column.
Q1     = 1st quartile quality score.
med    = Median quality score.
Q3     = 3rd quartile quality score.
IQR    = Inter-Quartile range (Q3-Q1).
lW     = 'Left-Whisker' value (for boxplotting).
rW     = 'Right-Whisker' value (for boxplotting).
A_Count = Count of 'A' nucleotides found in this column.
C_Count = Count of 'C' nucleotides found in this column.
G_Count = Count of 'G' nucleotides found in this column.
T_Count = Count of 'T' nucleotides found in this column.
N_Count = Count of 'N' nucleotides found in this column.
max-count = max. number of bases (in all cycles)
```

FASTQ Quality Chart

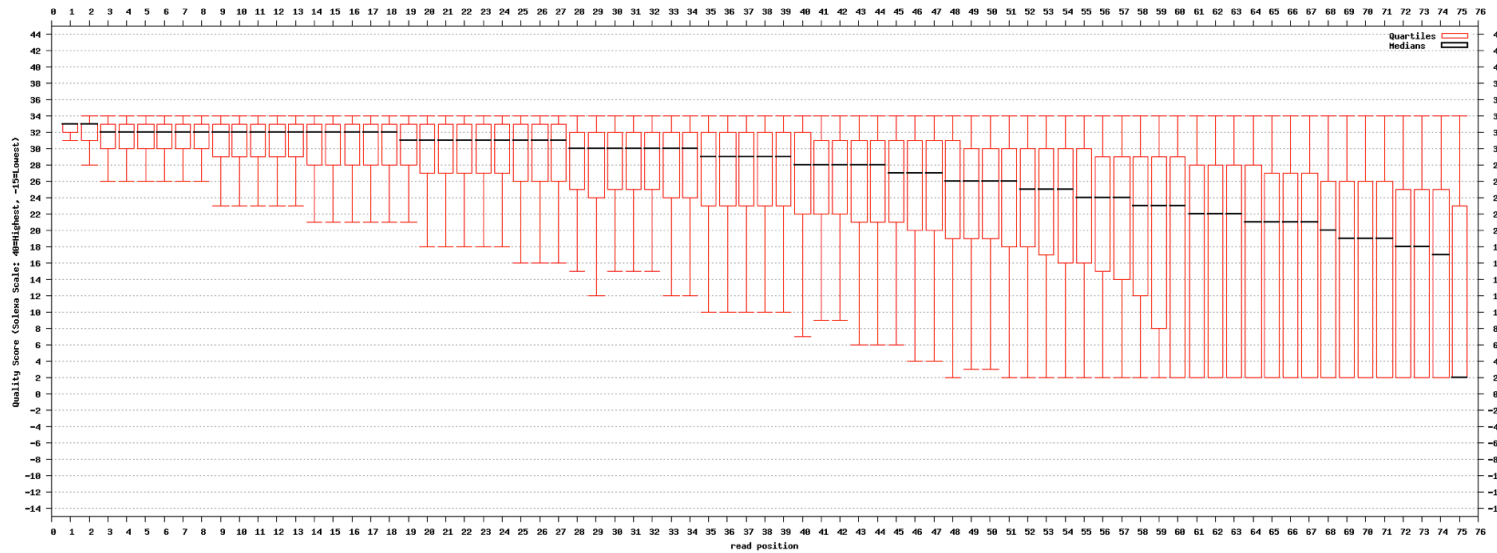
>fastx_quality_boxplot_graph.sh -i stats.txt -t sample1 -o quality.png

```
$ fastq_quality_boxplot_graph.sh -h
Solexa-Quality BoxPlot plotter
Generates a solexa quality score box-plot graph

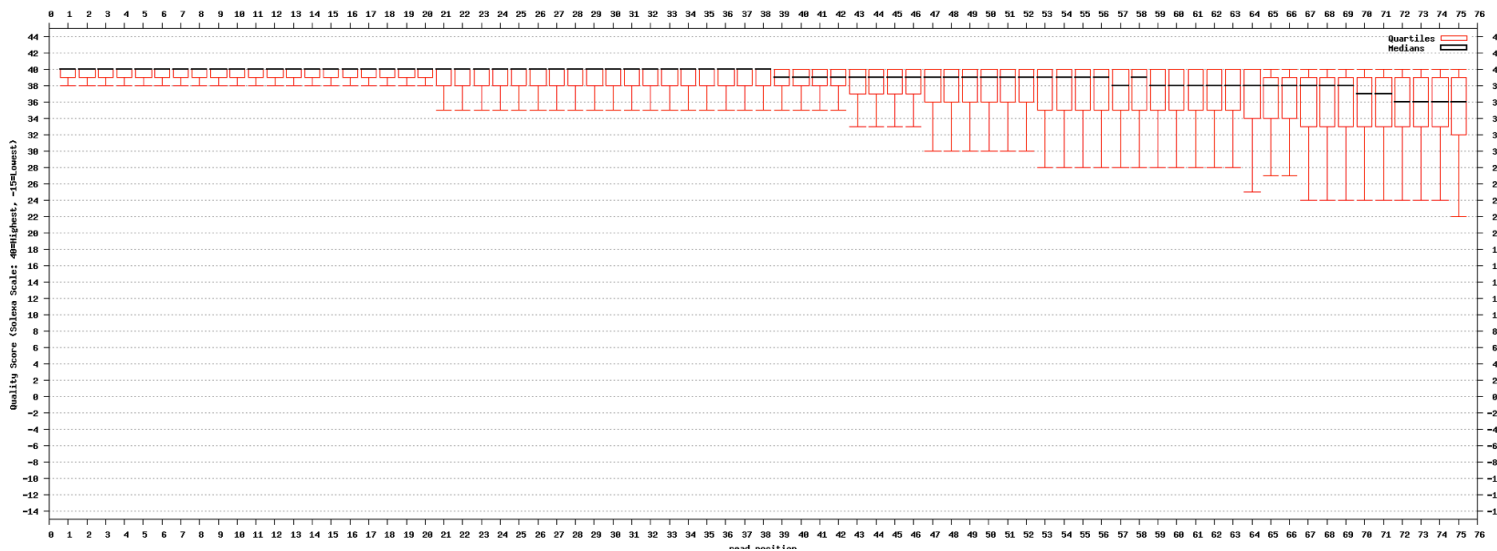
Usage: /usr/local/bin/fastq_quality_boxplot_graph.sh [-i INPUT.TXT] [-t TITLE] [-p] [-o OUTPUT]

[-p]          - Generate PostScript (.PS) file. Default is PNG image.
[-i INPUT.TXT] - Input file. Should be the output of "solexa_quality_statistics" program.
[-o OUTPUT]   - Output file name. default is STDOUT.
[-t TITLE]    - Title (usually the solexa file name) - will be plotted on the graph.
```

Read Quality with the FASTX-Toolkit



Bad
Sequence



Good
Sequence

High Throughput Sequence Alignment

- Traditional sequence alignment algorithms cannot be scaled to align millions of reads

High Throughput Sequence Alignment

- Traditional sequence alignment algorithms cannot be scaled to align millions of reads
- Utilize genome indexing such as Burrows-Wheeler for ultrafast and memory-efficient alignment programs

High Throughput Sequence Alignment

- Traditional sequence alignment algorithms cannot be scaled to align millions of reads
- Utilize genome indexing such as Burrows-Wheeler for ultrafast and memory-efficient alignment programs
- Next generation sequence alignment algorithms are rapidly evolving to accommodate the increasing sequence throughput

Tuxedo Suite

- Bowtie – fast and quality aware short read aligner for aligning DNA and RNA sequence reads (<http://bowtie-bio.sourceforge.net/index.shtml>)
- TopHat – fast, splice junction mapper for RNA-Seq reads built on the Bowtie aligner (<http://tophat.cbcb.umd.edu/>)
- Cufflinks – assembles transcripts, estimates their abundances, and test for differential expression and regulation using the alignments from Bowtie and TopHat

Bowtie

Bowtie

An ultrafast memory-efficient short read aligner



JOHNS HOPKINS
BLOOMBERG
SCHOOL of PUBLIC HEALTH

- Available for Windows, Mac OS X, Linux, and Solaris

Bowtie

Bowtie

An ultrafast memory-efficient short read aligner



JOHNS HOPKINS
BLOOMBERG
SCHOOL of PUBLIC HEALTH

- Available for Windows, Mac OS X, Linux, and Solaris
- Not a general-purpose alignment tool like MUMmer, BLAST, or Vmatch
- Ideal for aligning short reads to large genomes

Bowtie

Bowtie

An ultrafast memory-efficient short read aligner



JOHNS HOPKINS
BLOOMBERG
SCHOOL of PUBLIC HEALTH

- Available for Windows, Mac OS X, Linux, and Solaris
- Not a general-purpose alignment tool like MUMmer, BLAST, or Vmatch
- Ideal for aligning short reads to large genomes
- Forms the basis for TopHat, Cufflinks, Crossbow, and Myrna

Bowtie

Bowtie

An ultrafast memory-efficient short read aligner



JOHNS HOPKINS
BLOOMBERG
SCHOOL of PUBLIC HEALTH

- Available for Windows, Mac OS X, Linux, and Solaris
- Not a general-purpose alignment tool like MUMmer, BLAST, or Vmatch
- Ideal for aligning short reads to large genomes
- Forms the basis for TopHat, Cufflinks, Crossbow, and Myrna
- Online manual (<http://bowtie-bio.sourceforge.net/manual.shtml>) is very helpful to understand the options that are available

Bowtie

Bowtie

An ultrafast memory-efficient short read aligner



JOHNS HOPKINS
BLOOMBERG
SCHOOL of PUBLIC HEALTH

- Bowtie Index
 - To print a help screen with the optional parameters
 - `bowtie-build -h`
 - Command line usage
 - `bowtie-build [options]* <reference_in> <ebwt_base>`
 - Example command line
 - `bowtie-build chr1.fa,chr2.fa,chr3.fa,chr4.fa test_index`
 - Additional options can be used to improve performance

Bowtie

Bowtie

An ultrafast memory-efficient short read aligner



JOHNS HOPKINS
BLOOMBERG
SCHOOL of PUBLIC HEALTH

- Bowtie Index
 - To print a help screen with the optional parameters
 - `bowtie-build -h`
 - Command line usage
 - `bowtie-build [options]* <reference_in> <ebwt_base>`
 - Example command line
 - `bowtie-build chr1.fa,chr2.fa,chr3.fa,chr4.fa test_index`
 - Additional options can be used to improve performance
- Alignment with Bowtie
 - To print a help screen with the optional parameters
 - `bowtie -h`
 - Command line usage
 - `bowtie [options]* <ebwt> {-1 <m1> -2 <m2> | --12 <r> | <s>} [<hit>]`
 - Example command line for single end reads with our test_index
 - `Bowtie test_index --solexa1.3-quals test_index sequence_file.fastq > output_file`

Bowtie

Bowtie

An ultrafast memory-efficient short read aligner



JOHNS HOPKINS
BLOOMBERG
SCHOOL of PUBLIC HEALTH

- Additional Options:
- Quality aware MAQ-like mode is the default, Bowtie can also be run in a non-quality aware SOAP-like mode
 - `Bowtie test_index --solexa1.3-quals -v 2 test_index sequence_file.fastq > output_file`
 - `-v 2` is specifying that there can only be two mismatches in an alignment
- `-a` report all valid alignments
- `-k 2` report up to 2 valid alignments
- `--best` report the best alignment
- `-m 2` do not report any alignments for reads with greater than 2 reportable alignments
- `-S` print alignments in SAM format which is compatible with SAMtools for subsequent variant calling and manipulation of the alignments
- Many other options available

TopHat

TopHat

A spliced read mapper for RNA-Seq



- Available for Linux and OS X

TopHat

TopHat

A spliced read mapper for RNA-Seq



- Available for Linux and OS X
- Built on Bowtie and uses the same genome index

TopHat

TopHat

A spliced read mapper for RNA-Seq



- Available for Linux and OS X
- Built on Bowtie and uses the same genome index
- Used for alignment of RNA-Seq reads to a genome

TopHat

TopHat

A spliced read mapper for RNA-Seq



- Available for Linux and OS X
- Built on Bowtie and uses the same genome index
- Used for alignment of RNA-Seq reads to a genome
- Optimized for paired-end, Illumina sequence reads >70bp

TopHat

TopHat

A spliced read mapper for RNA-Seq

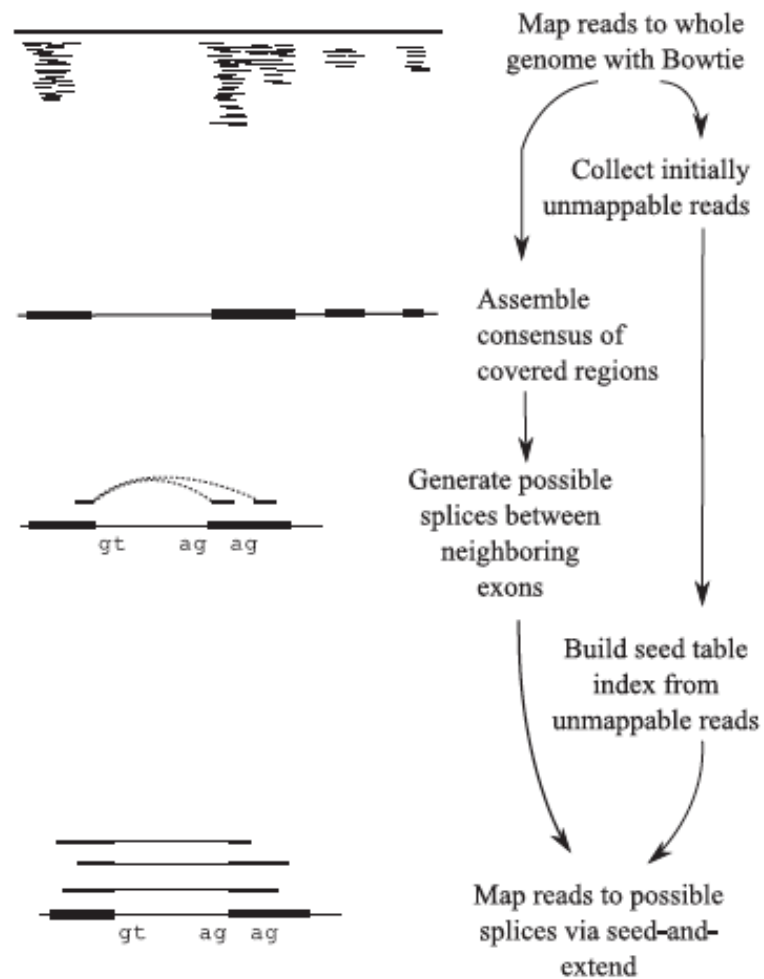


- Available for Linux and OS X
- Built on Bowtie and uses the same genome index
- Used for alignment of RNA-Seq reads to a genome
- Optimized for paired-end, Illumina sequence reads >70bp
- Online manual (<http://tophat.cbcb.umd.edu/manual.html>) is very helpful to understand the options that are available

TopHat

TopHat

A spliced read mapper for RNA-Seq



TopHat

TopHat

A spliced read mapper for RNA-Seq



- Alignment with TopHat

- To print a help screen with the optional parameters
 - tophat -h
- Command line usage
 - tophat [options]* <index_base> <reads1_1[,...,readsN_1]> [reads1_2,...readsN_2]
- Example command line for single end reads with our test_index from before
 - Tophat -o output_directory --solexa1.3-quals test_index sequence_file.fastq
- Additional options
 - --max-intron size (default is 500kb)
 - --min-intron size (default is 70bp)
 - --max-multihits (default is 40)
 - --mate-inner-dis (required for paired-end alignment mode)
 - --num_threads N (runs on N CPUs)

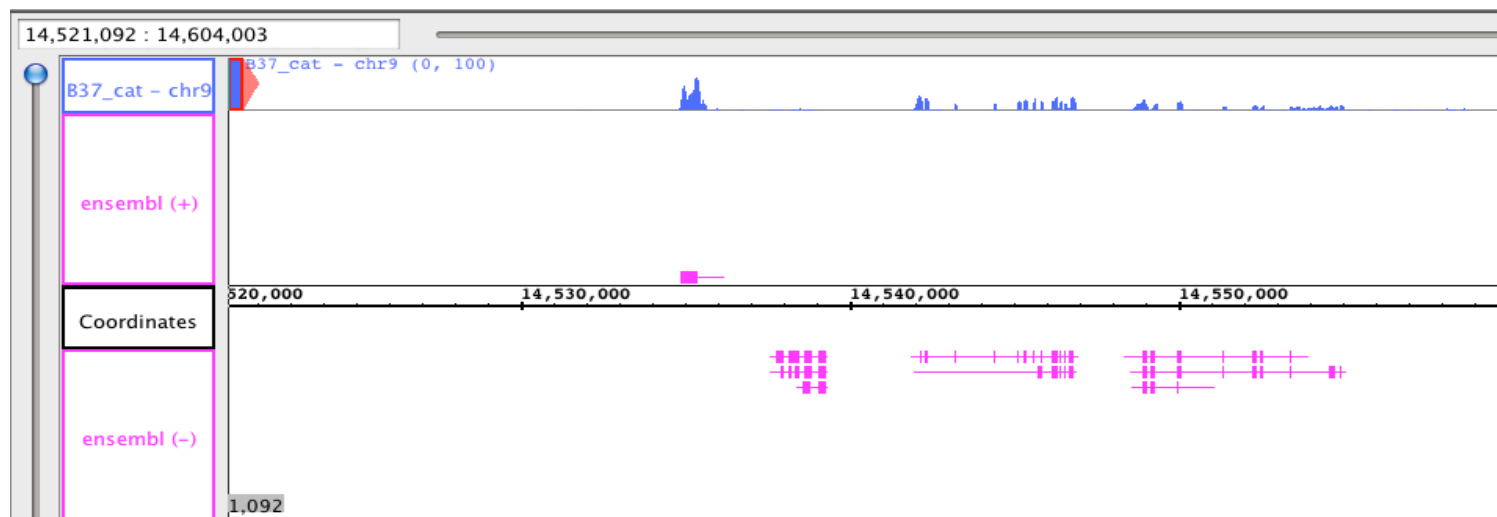
TopHat

TopHat

A spliced read mapper for RNA-Seq



- Wiggle tracks can be generated from the TopHat output file accepted_hits.bam
- Convert to SAM file
 - `samtools view -h -o accepted_hits.sam accepted_hits.bam`
- Generate wiggle file
 - `wiggles accepted_hits.sam coverage.wig`
- Wiggle tracks are viewable in programs such as Integrated Genome Browser (<http://bioviz.org/igb/>)



Additional Alignment Programs

How to map billions of short reads onto genomes

Cole Trapnell & Steven L Salzberg

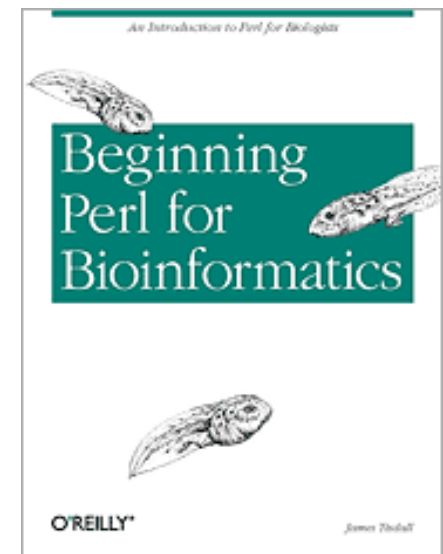
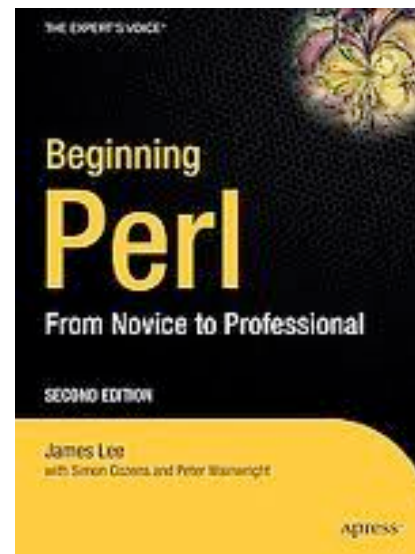
Mapping the vast quantities of short sequence fragments produced by next-generation sequencing platforms is a challenge. What programs are available and how do they work?

Nature Biotechnology **27**, 455-457 (2009)

Table 1 A selection of short-read analysis software

Program	Website	Open source?	Handles ABI color space?	Maximum read length
Bowtie	http://bowtie.cbcb.umd.edu	Yes	No	None
BWA	http://maq.sourceforge.net/bwa-man.shtml	Yes	Yes	None
Maq	http://maq.sourceforge.net	Yes	Yes	127
Mosaik	http://bioinformatics.bc.edu/marthlab/Mosaik	No	Yes	None
Novoalign	http://www.novocraft.com	No	No	None
SOAP2	http://soap.genomics.org.cn	No	No	60
ZOOM	http://www.bioinfor.com	No	Yes	240

Valuable Resources



Questions