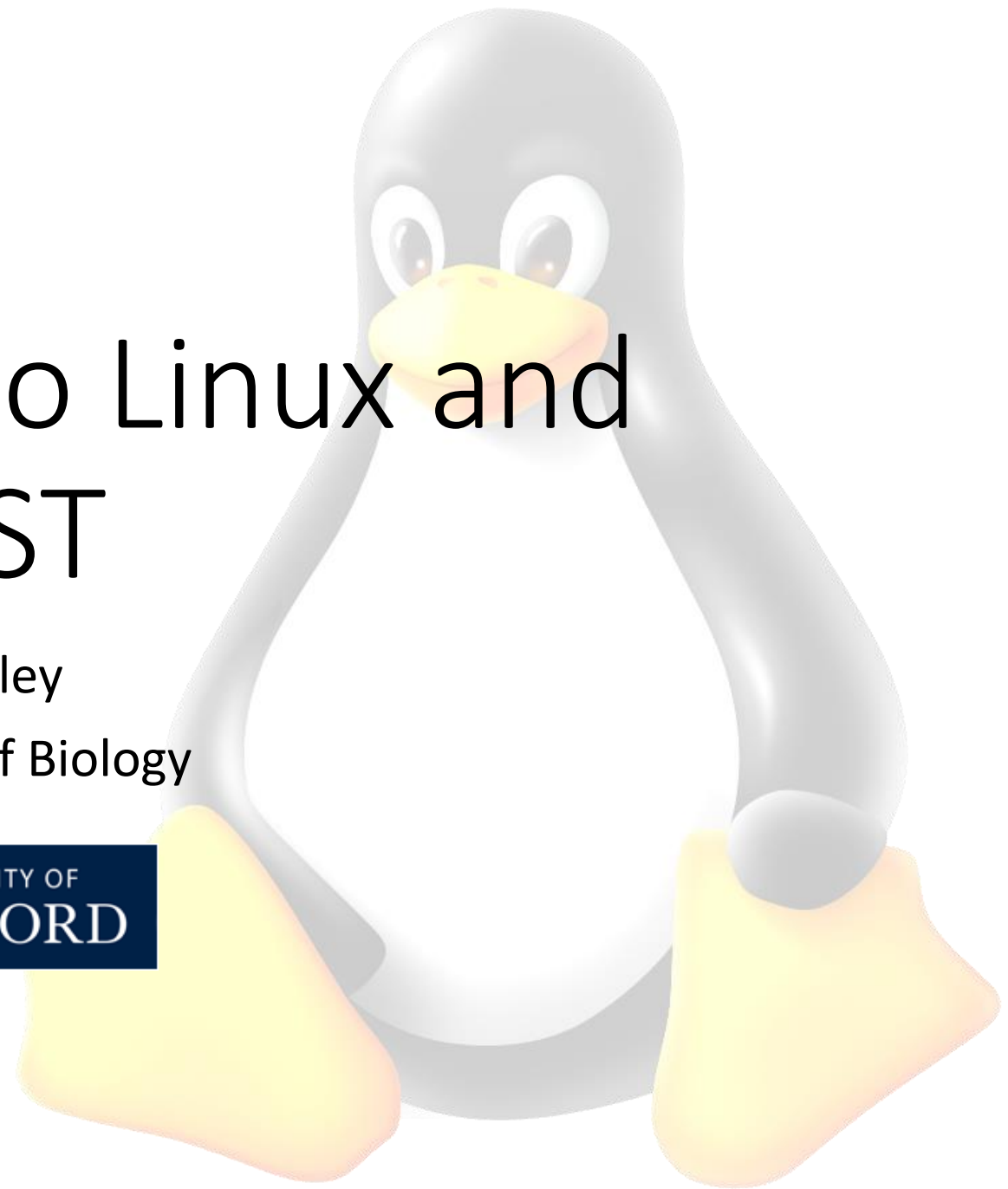


Introduction to Linux and BLAST

Keith Jolley
Department of Biology



UNIX

- Development dates back to 1960s.
- Philosophy is to use a large range of simple, dependable tools that each do one simple task.
- Combining tools facilitates complex analysis.
- Popular for high-performance computing.



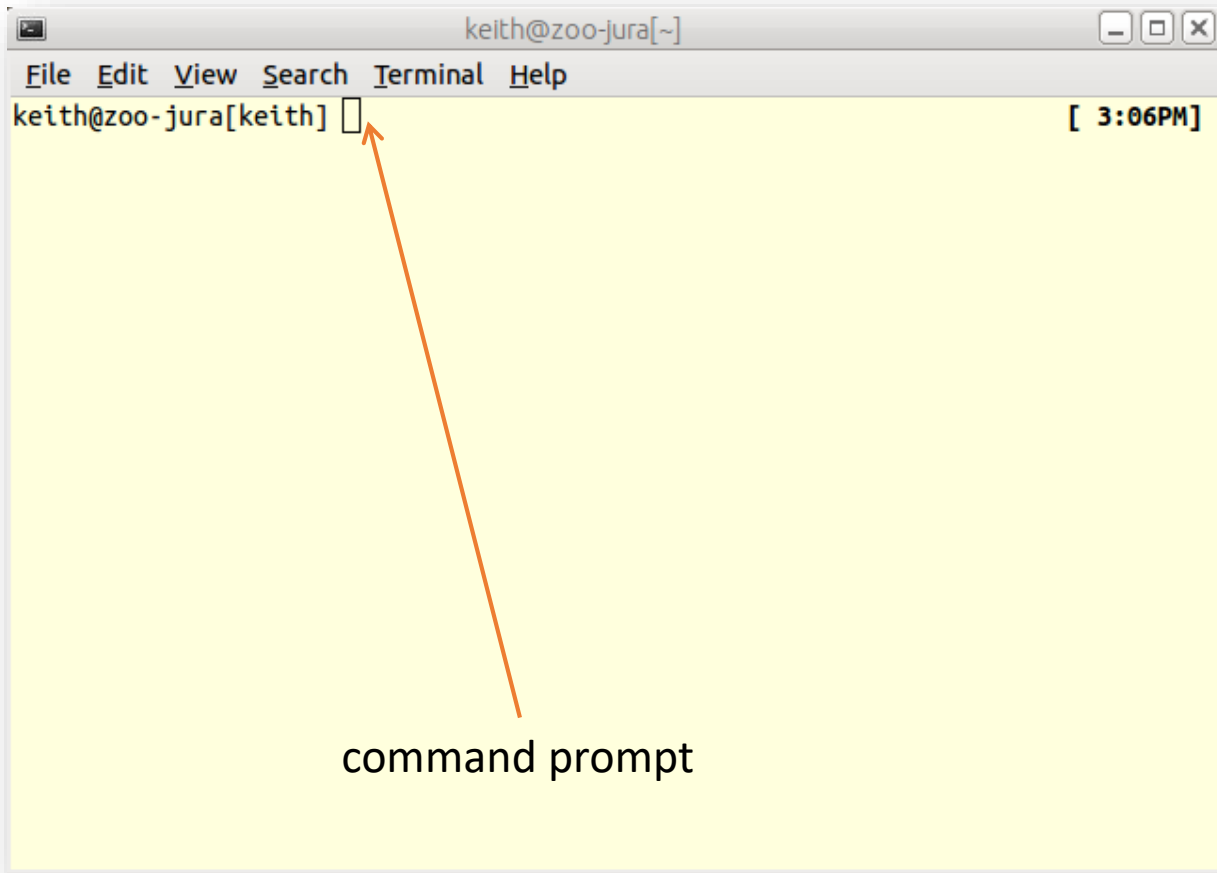
Ken Thompson (sitting) and Dennis Ritchie
at PDP-11 (photo: Peter Hamer) circa 1970

UNIX/Linux

- UNIX is the operating system of choice for engineering and scientific computing.
- Linux is a free Unix clone.
- Linux runs on most hardware.
 - Embedded systems
 - Mobile phone
 - Desktop computers
 - Super-computing clusters

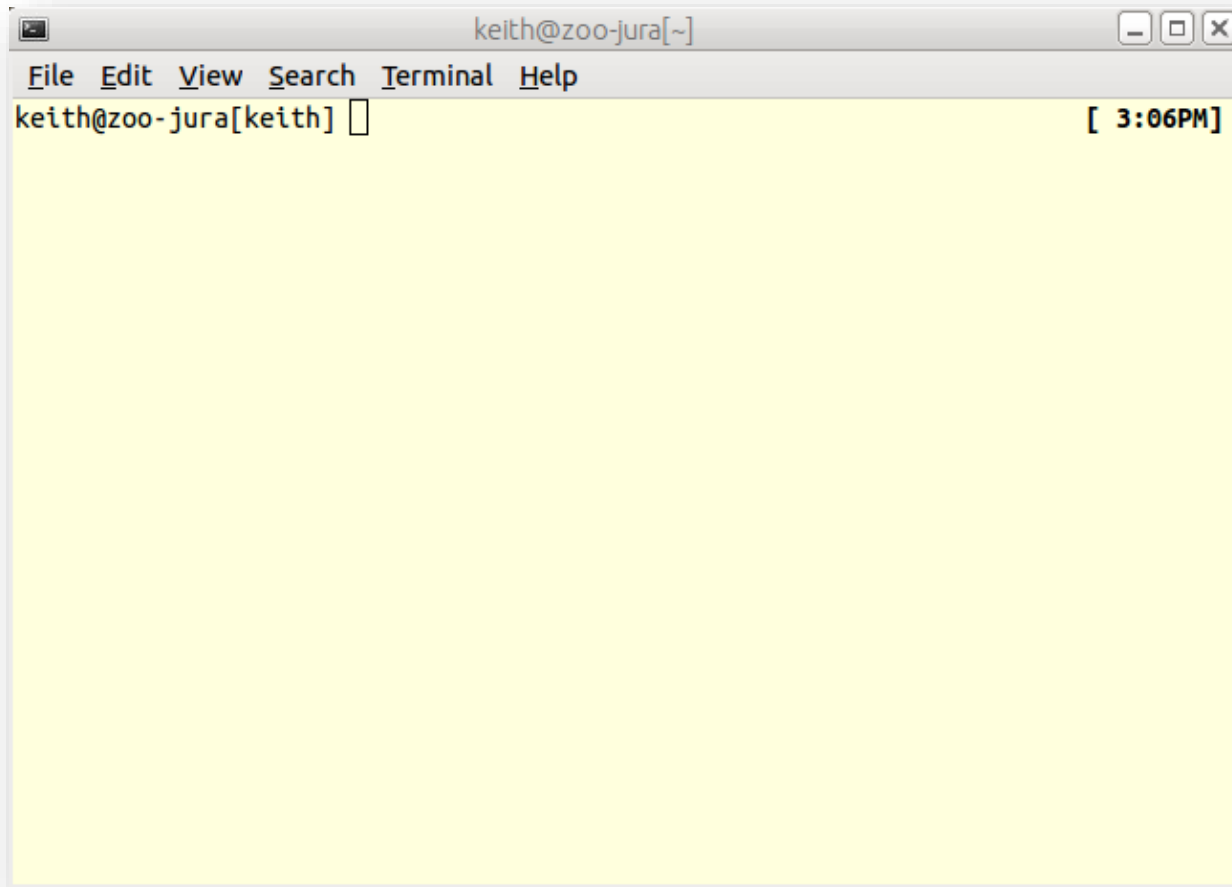


The command line



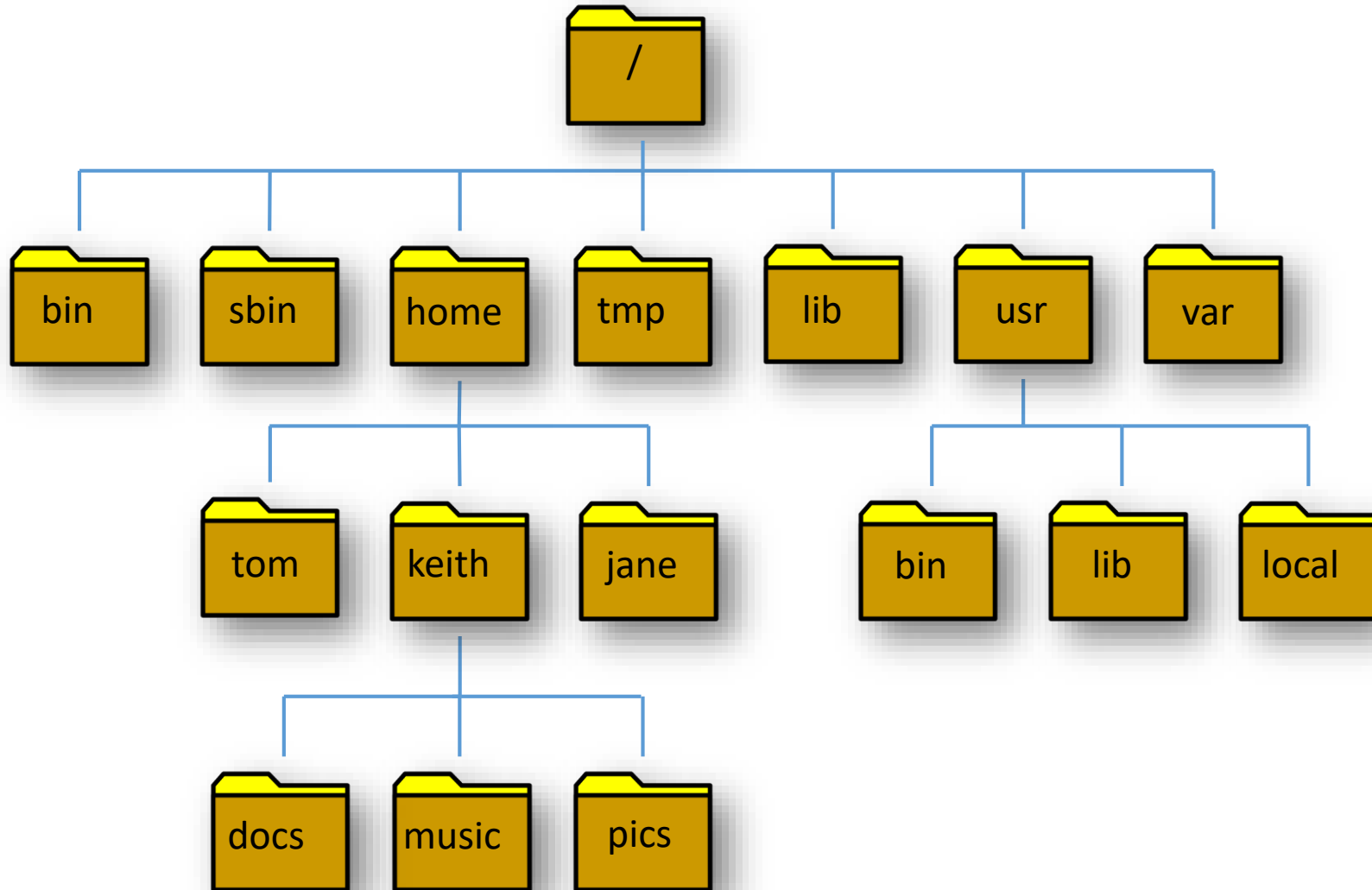
- A command line, or terminal, is a text-based interface to the system.
- You enter commands by typing them on the keyboard and feedback is given to you similarly as text.
- The command line usually presents you with a prompt.

The command line



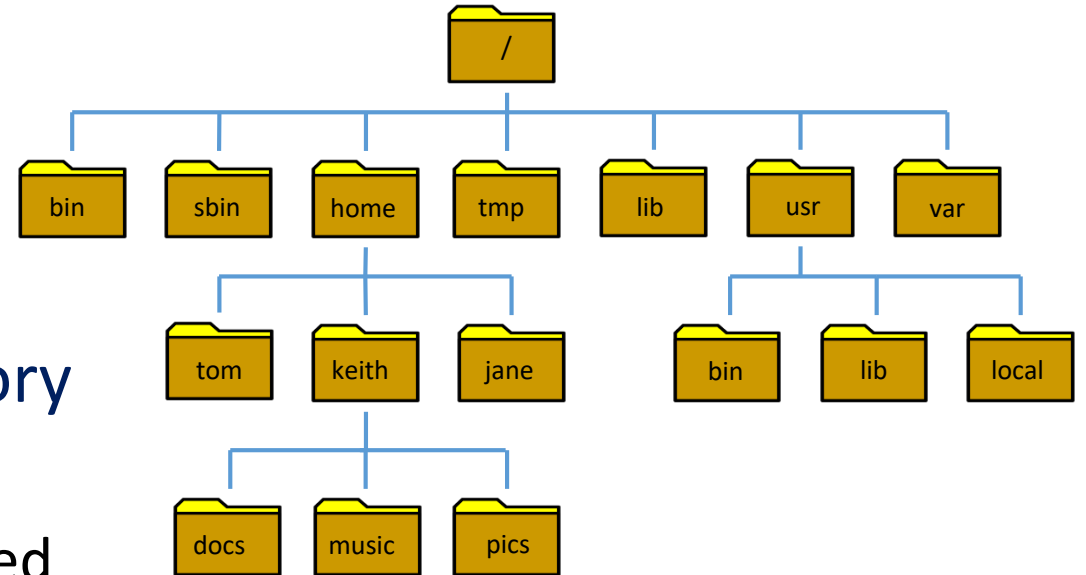
- Command syntax must be **exact!**
- Commands are **case-sensitive**.

The Linux file system is a hierarchy



Commands to navigate Linux filesystem

- **ls (list)**
 - Lists contents of current directory
- **cd (change directory)**
 - Change current focus to a new directory
 - Examples:
 - **cd documents** (enters directory called 'documents' in current directory)
 - **cd ..** (moves up a directory in the filesystem hierarchy)
 - **cd** (takes you to your home directory)
- **pwd (print working directory)**
 - Shows full path to current directory



Commands that operate on files

- mv (move file or directory)
 - Examples
 - `mv file1 dir1` (move 'file1' in to existing sub-directory 'dir1')
 - If destination does not exist, then file is renamed
 - `mv file1 file2` (renames 'file1' to 'file2')
- rm (remove file)
 - Examples
 - `rm file1` (remove 'file1')
 - `rm -fr dir1` (remove dir1 and all its contents)

BLAST

Basic Local Alignment Search Tool

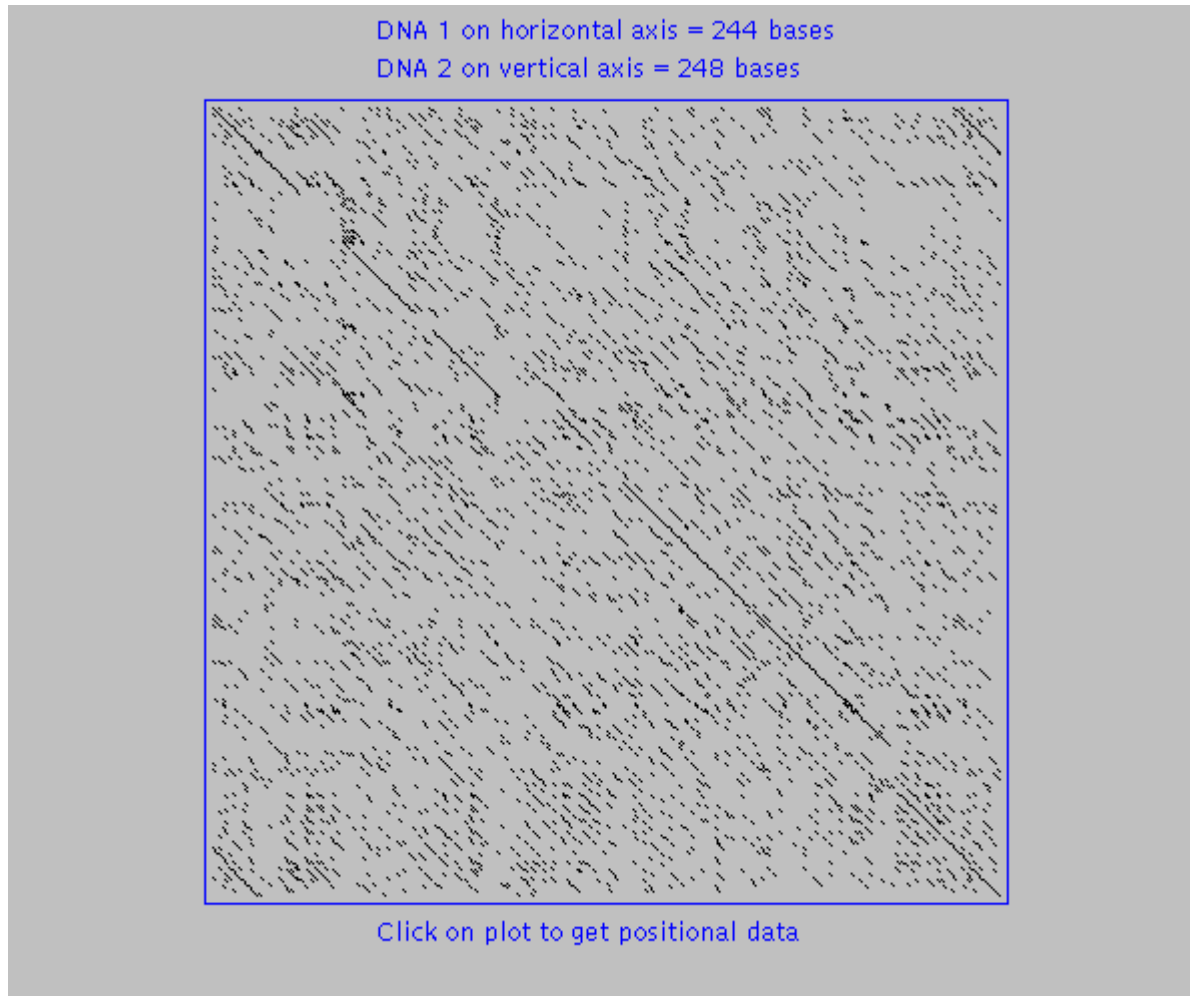
What is BLAST?

- Basic BLAST search
 - What is BLAST?
 - Different BLAST programs
 - BLAST databases you can search
 - Where can I run BLAST on the web?
- The BLAST algorithm
 - Stand-alone program is used behind-the-scenes in most sequence comparison applications

What is BLAST?

- BLAST stands for
Basic Local Alignment Search Tool
- Why BLAST is popular?
 - Good balance of sensitivity and speed
 - Reliable
 - Flexible
- Local alignments; search query divided into short sequences ('words') and exact matches identified. Once found, these matches are extended

BLAST extends local alignments



- Many nucleotide matches purely by chance (1:4)
- BLAST only considers significant 'seed' matches (default word size: 11)
- Scoring system takes account of matches, mismatches, gap formation and gap extension

BLAST output

1. List of sequences with scores

- Raw score
 - Higher is better
 - Depends on aligned length
- Expect Value (E-value)
 - Smaller is better
 - Independent of length and database size

2. List of alignments

Score = 248 bits (129), Expect = 1e-63
Identities = 213/263 (80%), Gaps = 34/263 (12%)
Strand = Plus / Plus

Query: 161 atatcaccacgtcaaagggtgactccaactcca---ccactccattttgttcagataaatgc 217
|||||||
Sbjct: 481 atatcaccacgtcaaagggtgactccaact-tattgatagtgttttatgttcagataaatgc 539

Query: 218 ccgatgatcatgtcatgcagctccaccgattgtgagaacgacagcgcacttccgtcccagc 277
|||||||
Sbjct: 540 ccgatgactttgtcatgcagctccaccgattttg-g-----ttccgtcccagc 586

Query: 278 c-gtgcc--aggtgctgcctcagattcagggttatgccgctcaattcgctgcgtatatcgc 334
| || | |||||
Sbjct: 587 caatgacgta-gtgctgcctcagattcagggttatgccgctcaattcgctgggtatatcgc 645

Query: 335 ttgctgattacgtgcagctttcccttcaggcgggga-----ccagccatccgtc 382
|||||||
Sbjct: 646 ttgctgattacgtgcagctttcccttcaggcggggattcatacagcggccagccatccgtc 705

Query: 383 ctccatatc-accacgtcaaagg 404
|||||||
Sbjct: 706 atccatataaccacgtcaaagg 728

BLAST Programs

Program	Database (Subject)	Query
BLAST ^N	Nucleotide	Nucleotide
BLAST ^P	Protein	Protein
BLAST ^X	Protein	Nt. → Protein
TBLAST ^N	Nt. → Protein	Protein
TBLAST ^X	Nt. → Protein	Nt. → Protein



Basic Local Alignment Search Tool

BLAST finds regions of similarity between biological sequences. The program compares nucleotide or protein sequences to sequence databases and calculates the statistical significance.

[Learn more](#)

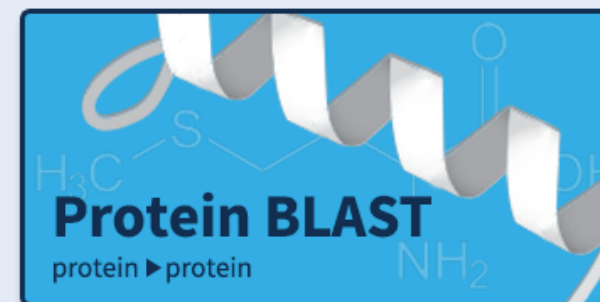
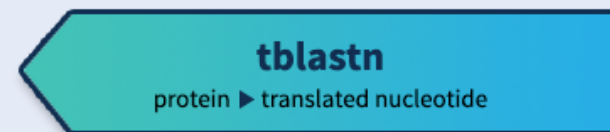
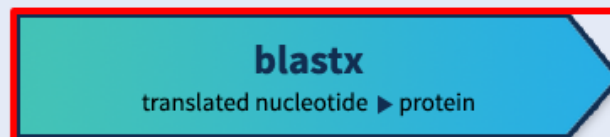
NEWS

BLAST+ 2.8.1 is released

New databases, better performance.
Wed, 19 Dec 2018 17:00:00 EST

[More BLAST news...](#)

Web BLAST



BLAST Genomes

[Search](#)[Human](#)[Mouse](#)[Rat](#)[Microbes](#)

Enter Query Sequence

BLASTX search protein databases using a translated nucleotide query. [more...](#)[Reset page](#)[Bookmark](#)

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

[Clear](#)

Query subrange

From

To

```
GGCGGCAAAA AACCTCAACA TACAAATGT CGCACCGCAC CCGCGCCAAA
CCTACGGGCT TTCCGGCGTA ACCACATTAA ATTGCGCCTA TGTCAICGAA
GACTATATTT TGGCGGAGAT TAAGAAAAAC CCGCATACGA GGTATGAAAT
TTATACCTTT TTACAGCGGCG CGGCGTTGAC GATGAAGGAT TTCCCAATG
TGCACGTTTA CGCATTGAAA CCGGCTTCCC TTCCGGAAGA TTATTGGCTC
AAGCCGGTGT ATGCCCTGTT TACCCAATCC GGCATCCCGA TTTTGACATT T
```

Or, upload file

Browse...

No file selected.

Genetic code

Standard (1)

Job Title

Enter a descriptive title for your BLAST search

☐ Align two or more sequences

Choose Search Set

Database

Non-redundant protein sequences (nr)

Organism

Optional

Enter organism name or id-completions will be suggested

☐ 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

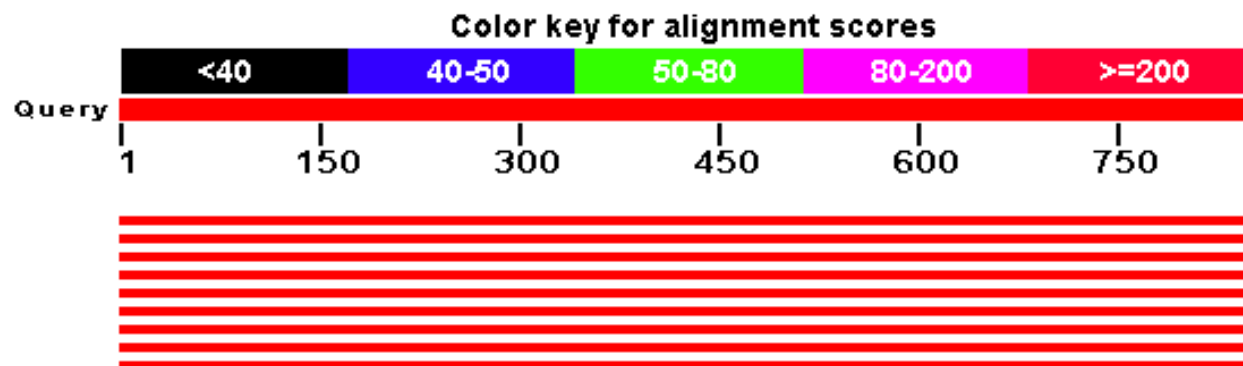
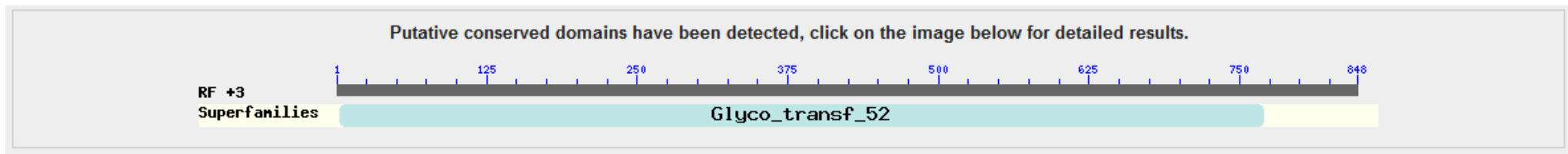
[YouTube](#) [Create custom database](#)

BLAST

Search database Non-redundant protein sequences (nr) using Blastx (search protein databases using a translated nucleotide query)

☐ Show results in a new window[+ Algorithm parameters](#)

NCBI BLAST output



Sequences producing significant alignments:

Select: [All](#) [None](#) Selected:0

Alignments Download GenPept Graphics							
	Description	Max score	Total score	Query cover	E value	Ident	Accession
<input type="checkbox"/>	alpha-2,3-sialyltransferase [Neisseria meningitidis]	576	576	99%	0.0	99%	AAC44543.1
<input type="checkbox"/>	CMP-N-acetylneuraminate-beta-galactosamide-alpha-2 3-sialyltransferase [Neisseria meningitidis]	576	576	99%	0.0	99%	WP_002234614.1
<input type="checkbox"/>	CMP-N-acetylneuraminate-beta-galactosamide-alpha-2 3-sialyltransferase [Neisseria meningitidis]	575	575	99%	0.0	99%	WP_014580539.1
<input type="checkbox"/>	CMP-N-acetylneuraminate-beta-galactosamide-alpha-2 3-sialyltransferase [Neisseria meningitidis]	574	574	99%	0.0	99%	WP_002236898.1
<input type="checkbox"/>	CMP-N-acetylneuraminate-beta-galactosamide-alpha-2 3-sialyltransferase [Neisseria meningitidis]	573	573	99%	0.0	98%	WP_002230437.1
<input type="checkbox"/>	CMP-N-acetylneuraminate-beta-galactosamide-alpha-2 3-sialyltransferase [Neisseria meningitidis]	573	573	99%	0.0	98%	WP_002239594.1

Sequence alignment output

CMP-N-acetylneuraminate-beta-galactosamide-alpha-2 3-sialyltransferase [Neisseria meningitidis]

Sequence ID: [ref|WP_002234614.1|](#) Length: 371 Number of Matches: 1

► [See 29 more title\(s\)](#)

Range 1: 86 to 367 [GenPept](#) [Graphics](#)

▼ Next Match ▲ Previous Match

Score	Expect	Method	Identities	Positives	Gaps	Frame
576 bits(1485)	0.0	Compositional matrix adjust.	279/282(99%)	279/282(98%)	0/282(0%)	+3
Query 3	NEKYDYYFKQIKDKAERAYFFHLPHYDLNKSFNVIPTMAELKVKSMLLPKVKRTYLASLEK				182	
	NEKYDYYFKQIKDKAERAYFFHLPY LNKSFN IPTMAELKVKSMLLPKVKR YLASLEK					
Sbjct 86	NEKYDYYFKQIKDKAERAYFFHLPHYGLNKSFNFIPTMAELKVKSMLLPKVKRIYLASLEK				145	
Query 183	VSIAAFLSTYPDAEIKTFDDGTGNLIQSSSYLGDEFSVNGTIKRNFARMMIGDWSIAKTR				362	
	VSIAAFLSTYPDAEIKTFDDGTGNLIQSSSYLGDEFSVNGTIKRNFARMMIGDWSIAKTR					
Sbjct 146	VSIAAFLSTYPDAEIKTFDDGTGNLIQSSSYLGDEFSVNGTIKRNFARMMIGDWSIAKTR				205	
Query 363	NASDEHYTIFKGLKNIMDDGRRKMTYLPLFDASELKAGDETGGTVRILLGSPDKEMKEIS				542	
	NASDEHYTIFKGLKNIMDDGRRKMTYLPLFDASELKAGDETGGTVRILLGSPDKEMKEIS					
Sbjct 206	NASDEHYTIFKGLKNIMDDGRRKMTYLPLFDASELKAGDETGGTVRILLGSPDKEMKEIS				265	
Query 543	EKAAKNFNIQYVAPHPRQTYGLSGVTTLNSPYVIEDYILREIKKNPHTRYEIYTFFSGAA				722	
	EKAAKNFNIQYVAPHPRQTYGLSGVTTLNSPYVIEDYILREIKKNPHTRYEIYTFFSGAA					
Sbjct 266	EKAAKNFNIQYVAPHPRQTYGLSGVTTLNSPYVIEDYILREIKKNPHTRYEIYTFFSGAA				325	
Query 723	LTMKDFPNVHVYALKPASLPEDYWLPVYALFTQSGIPILTF			848		
	LTMKDFPNVHVYALKPASLPEDYWLPVYALFTQSGIPILTF					
Sbjct 326	LTMKDFPNVHVYALKPASLPEDYWLPVYALFTQSGIPILTF			367		

Gotchas

- BLAST will usually return a result!
 - The top match may not be significant
 - E-values up to 10 returned by default
- Domains are often shared by different proteins
 - Significant match to parts of the sequence, but not to others