# BLAST & FASTA

- FASTA

[Lipman, Pearson '85, '88]

- Basic Local Alignment Search Tool

[Altschul, Gish, Miller, Myers, Lipman '90]

# BLAST Overview

- Program(s) to search all sequence databases
- Tremendous Speed/Less Sensitive
- Statistical Significance reported
- WWWBLAST, QBLAST (send now, retrieve results later), Standalone BLAST, BLASTcl3 (Client version, TCP/IP connection to NCBI server), BLAST URLAPI (to access QBLAST, no local client)

# BLAST Strategy & Improvements

- Lipman et al.: speeded up finding "runs" of "hot spots".

- Eugene Myers '94: "Sublinear algorithm for approximate keyword matching".

- Karlin, Altschul, Dembo '90, '91: "Statistical Significance of Matches"

# BLAST Variants

- **Nucleotide BLAST**
  - **Standard**
  - **MEGABLAST** (Compare large sets, Near-exact searches)
  - **Short Sequences** (higher E-value threshold, smaller word size, no low-complexity filtering)
- **Protein BLAST**
  - **Standard**
  - **PSI-BLAST** (Position Specific Iterated BLAST)
  - **PHI-BLAST** (Pattern Hit Initiated BLAST; reg expr. Or Motif search)
  - **Short Sequences** (higher E-value threshold, smaller word size, no low-complexity filtering, PAM-30)
- **Translating BLAST**
  - **Blastx**: Search nucleotide sequence in protein database (6 reading frames)
  - **Tblastn**: Search protein sequence in nucleotide dB
  - **Tblastx**: Search nucleotide seq (6 frames) in nucleotide DB (6 frames)

# BLAST Cont'd

- **RPS BLAST**
  - Compare protein sequence against Conserved Domain DB; Helps in predicting rough structure and function
- **Pairwise BLAST**
  - blastp (2 Proteins), blastn (2 nucleotides), tblastn (protein-nucleotide w/ 6 frames), blastx (nucleotide-protein), tblastx (nucleotide w/6 frames-nucleotide w/ 6 frames)
- **Specialized BLAST**
  - Human & Other finished/unfinished genomes
  - P. falciparum: Search ESTs, STSs, GSSs, HTGs
  - VecScreen: screen for contamination while sequencing
  - IgBLAST: Immunoglobin sequence database

# BLAST Credits

- Stephen Altschul
- Jonathan Epstein
- David Lipman
- Tom Madden
- Scott McGinnis
- Jim Ostell
- Alex Schaffer
- Sergei Shavirin
- Heidi Sofia
- Jinghui Zhang

# Databases used by BLAST

- **Protein**
  - nr (everything), swissprot, pdb, alu, individual genomes

- **Nucleotide**
  - nr, dbest, dbsts, htgs (unfinished genomic sequences), gss, pdb, vector, mito, alu, epd

- **Misc**

# Rules of Thumb

- Most sequences with significant similarity over their entire lengths are homologous.

- Matches that are > 50% identical in a 20-40 aa region occur frequently by chance.

- Distantly related homologs may lack significant similarity. Homologous sequences may have few absolutely conserved residues.

- A homologous to B & B to C $\Rightarrow$ A homologous to C.

- Low complexity regions, transmembrane regions and coiled-coil regions frequently display significant similarity without homology.

- Greater evolutionary distance implies that length of a local alignment required to achieve a statistically significant score also increases.