# Bioinformatics for Biologists

Sequence Analysis: Part II. Pattern Searching and Gene Finding

Fran Lewitter, Ph.D.
Director
Bioinformatics & Research Computing
Whitehead Institute

---

## Topics to Cover

- Pattern searching
  – PSI-BLAST
  – PHI-BLAST
  – Finding patterns
- Gene finding

---

## PSI-BLAST

- **P**osition **S**pecific **I**terative BLAST uses a profile (or position specific scoring matrix, PSSM) that is constructed (automatically) from a multiple alignment of the highest scoring hits in an initial BLAST search.
- The PSSM is generated by calculating position-specific scores for each position in the alignment. Highly conserved positions receive high scores and weakly conserved positions receive scores near zero.
- The profile is used to perform a second (etc.) BLAST search and the results of each "iteration" is used to refine the profile. This iterative searching strategy results in increased sensitivity.

---

## Start with a BLASTP search

---

## PSI-BLAST - Iteration 1

---

## PSSM from PSI-BLAST

Amino acids

| | A | R | N | D | C | Q | E | G | H | I | L | K | M | F | P | S | T | W | Y | V |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 0 | 2 | 3 | 2 | 4 | 1 | 1 | 4 | 3 | 0 | 3 | 3 | 7 | 3 | 3 | 2 | 1 | 0 | 1 | 2 |
| 2 | 6 | 0 | 3 | 3 | 5 | 4 | 0 | 3 | 2 | 5 | 0 | 1 | 2 | 2 | 4 | 1 | 3 | 2 | 4 | 2 |
| 3 | 4 | 3 | 0 | 3 | 3 | 1 | 3 | 2 | 4 | 2 | 3 | 2 | 5 | 0 | 1 | 2 | 1 | 0 | 5 | 7 |
| 4 | 3 | 2 | 3 | 2 | 9 | 3 | 3 | 5 | 4 | 0 | 3 | 2 | 5 | 1 | 2 | 2 | 4 | 1 | 2 |   |
| 5 | 0 | 1 | 2 | 2 | 4 | 1 | 6 | 3 | 3 | 1 | 3 | 2 | 0 | 4 | 8 | 3 | 1 | 0 | 3 | 0 |
| 6 | 4 | 3 | 2 | … | | | | | | | | | | | | | | | | |
| • | … | | | | | | | | | | | | | | | | | | | |
| • | … | | | | | | | | | | | | | | | | | | | |
| N | | | | | | | | | | | | | | | | | | | | |

POSITIONS

## Pattern Hit Initiated (PHI)-BLAST

>HUMAN MSH2

MAVQPKETLQLESAAEVGFVRFFQGMPEKPTTTVRLFDRGDFYTAHGEDALLAAREVFKTQGVIKYMGPA
GAKNLQSVVLSKMNFESFVKDLLLVRQYRVEVYKNRAGNKASKENDWYLAYKASPGNLSQFEDILFGNND
MSASIGVVGVKMSAVDGQRQVGVGYVDSIQRKLGLCEFPDNDQFSNLEALLIQIGPKECVLPGGETAGDM
GKLRQIIQRGGILITERKKADFSTKDIYQDLNRLLKGKKGEQMNSAVLPEMENQVAVSSLSAVIKFLELL
SDDSNFGQFELTTFDFSQYMKLDIAAVRALNLFQGSVEDTTGSQSLAALLNKCKTPQGQRLVNQWIKQPL
MDKNRIEERLNLVEAFVEDAELRQTLQEDLLRRFPDLNRLAKKFQRQAANLQDCYRLYQGINQLPNVIQA
LEKHEGKHQKLLLAVFVTPLTDLRSDFSKFQEMIETTLDMDQVENHEFLVKPSFDPNLSELREIMNDLEK
KMQSTLISAARDLGLDPGKQIKLDSSAQFGYYFRVTCKEEKVLRNNKNFSTVDIQKNGVKFTNSKLTSLN
EEYTKNKTEYEEAQDAIVKEIVNISSGYVEPMQTLNDVLAQLDAVVSFAHVSNGAPVPYVRPAILEKGQG
RIILKASRHACVEVQDEIAFIPNDVYFEKDKQMFHIITGPNMGGKSTYIRQTGVIVLMAQIGCFVPCESA
EVSIVDCILARVGAGDSQLKGVSTFMAEMLETASILRSATKDSLIIIDELGRGTSTYDGFGLAWAISEYI
ATKIGAFCMFATHFHELTALANQIPTVNNLHVTALTTEETLTMLYQVKKGVCDQSFGIHVAELANFPKHV
                          FQYIGESQGYDIMEPAAKKCYLEREQQEKIIQEFLSKVKQMPFTEMSEENITIKLKQ
                     NEIISRIKVTT

DNA mismatch repair proteins mutS family signature

---

## PHI-BLAST

>gi|4099512|gb|AAD00647.1| (U87911) MutS homolog 2 [Arabidopsis thaliana]
        Length = 117

 Score =  136 bits (364), Expect = 1e-40
 Identities = 88/117 (75%), Positives = 98/117 (83%)

Query:  668  TGPNMGGKSTYIRQTGVIVLMAQIGCFVPCESAEVSIVDCILARVGAGDSQLKGVSTFMA  727
             TGPNMGGKST+IRQ GVIVLMAQ+G FVPC+ A +SI DCI ARVGAGD QL+GVSTFM
Sbjct:    1  TGPNMGGKSTFIRQVGVIVLMAQVGSFVPCDKASISIRDCIFARVGAGDCQLRGVSTFMQ   60

Query:  728  EMLETASILRSATKDSLIIIDELGRGTSTYDGFGLAWAISEYIATKIGAFCMFATHF  784
pattern 743          ****************
             EMLETASIL+ AT  SLIIIDELGRGTSTYDGFGLAW AI E++    A  +FATH+
Sbjct:   61  EMLETASILKSATDKSLIIIDELGRGTSTYDGFGLAWAICEHLVQVKRAPTLFATHY  117

---

## Pattern Searching

| RRRRYYYY or R(4)Y(4) | 4 purines followed by 4 pyrimidines |
| --- | --- |
| TATAA | Exact pattern: TATAA |
| [DE](2)HS{P}X(2)PX(2,4)C | [ ] Acceptable { } Not acceptable X(2)  2 of anything in a row X(2,4) from 2 to four of of anything |
| p1=6...8 GAGA ~p1 | hairpin with GAGA as the loop |
| p1=6...6 3...8 p1 | exact 6 character repeat separated by up to 8 |
| p1=6...6 3..8 p1[1,1,1] | allow one mismatch, deletion and insertion |

---

## Pattern Searching Programs

**Patscan**          scan_for_matches patfile < inputfile

**fuzznuc, fuzzprot, fuzztrans, dreg**          EMBOSS programs; web and Unix

---

## Interesting features in DNA

- ***Structure of genes*** – exons, introns, etc.
- ***Non-coding RNAs*** – including micro RNAs and other small RNAs
- ***Promoter sites***
- ***Alternative splice forms***

---

## Problem to Solve

# Types of Signals to Detect

- Transcriptional
  - TSS
  - TATA box
  - PolyA
- Translational
  - Kozak (CC A/G CCAUGG)
  - Termination codon (UAA, UAG, UGA)
- Splicing
  - Introns - GT……AG

# Gene Finding Strategies

- Content-based methods
  - codon usage, compositional complexity
- Site-based methods
  - presence or absence of specific pattern or sequence
- Comparative methods
  - determination based on homology

# RepeatMasker

# RepeatMasker

Summary:

Total length:   8750 bp

GC level:        35.61%

Bases masked: 6803 bp (77.75%)

| | number of elements* | length occupied | percentage of sequence |
|---|---|---|---|
| SINEs: | 4 | 1159 bp | 13.25 % |
| ALUs | 4 | 1159 bp | 13.25 % |
| MIRs | 0 | 0 bp | 0.00 % |
| LINEs: | 3 | 5605 bp | 64.06 % |
| LINE1 | 3 | 5605 bp | 64.06 % |
| LINE2 | 0 | 0 bp | 0.00 % |
| L3/CR1 | 0 | 0 bp | 0.00 % |
| LTR elements: | 0 | 0 bp | 0.00 % |
| MaLRs | 0 | 0 bp | 0.00 % |
| ERVL | 0 | 0 bp | 0.00 % |
| ERV_classI | 0 | 0 bp | 0.00 % |
| ERV_classII | 0 | 0 bp | 0.00 % |
| DNA elements: | 0 | 0 bp | 0.00 % |
| MER1_type | 0 | 0 bp | 0.00 % |
| MER2_type | 0 | 0 bp | 0.00 % |
| Unclassified: | 0 | 0 bp | 0.00 % |
| Total interspersed repeats: | | 6764 bp | 77.30 % |
| Small RNA: | 0 | 0 bp | 0.00 % |
| Satellites: | 0 | 0 bp | 0.00 % |
| Simple repeats: | 1 | 39 bp | 0.45 % |
| Low complexity: | 0 | 0 bp | 0.00 % |

# Coding Measures

- Look at frequencies of codons (e.g. redundancy of genetic code; Leucine = UUA, UUG, CUU, CUC, CUA, CUG)

- 6-tuple or hexamer approach
  ACCTCG  TACTCG  GCCCTC
  Thr   Ser   Tyr   Ser   Ala   Leu

# GENSCAN

- MM - prob for a given nuc to occur at position $p$ depends on nuc occupying previous $k$ positions
- Generalized Hidden Markov Model (GHMM)
- Optimize module performing signal recognition
- Incorporates influence of C+G content
- Considers gene models on both strands
- Can identify multiple genes

Burge and Karlin, JMB:268:78-94, 1997

## Gene Finding Programs

- FGENESH - Softberry
- GeneID - Barcelona
- Genscan - Stanford and MIT
- GenomeScan - MIT
- MZEF/First exon - Cold Spring Harbor
- Twinscan - WU
- Ecgene - Korea (EST clustering)
- SGP - Barcelona (human-mouse homology)

## HMR195  Test Set

- 103 human, 82 mouse, 10 rat sequences
- Sequence new since August, 1997
- Genomic sequences containing exactly one gene
- No mRNA sequences, pseudogenes or alternatively spliced genes
- The mean length of sequences is 7,096 bp

*Rogic, Mackworth and Ouellette, Genome Research 11:817-832, 2001*

## Definitions

- *Sensitivity:* the proportion of true sites(e.g., exons or donor splice sites) that are correctly predicted = TP/(TP + FN)

- ***Predictive value of positive results ("Specificity" in gene finding literature):*** the proportion of predicted sites that are correct = TP/(TP + FP)

- *Specificity:* the proportion of non-sites that are predicted to be non-sites = TN/(TN +FP)
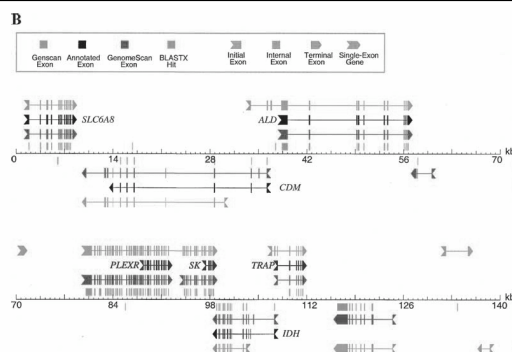
## Program Comparisons Results

- Genscan and HMMgene had reliable scores for exons
- Nucleotide Sn = .95 for Genscan and .93 for HMMgene.
- Sp = .90 and .93, respectively
- Accuracy dependent on G+C content

*Rogic, Mackworth and Ouellette, Genome Research 11:817-832, 2001*

**Table 2.**  Accuracy versus Signal Type

| Programs | start codon (195) | acceptor site (753) | donor site (753) | stop codon (195) |
|---|---|---|---|---|
| FGENES | 0.67 (0.63) | 0.80 (0.77) | 0.85 (0.82) | 0.75 (0.72) |
| GeneMark.hmm | 0.46 (0.60) | 0.81 (0.75) | 0.82 (0.78) | 0.57 (0.64) |
| Genie | 0.56 (0.57) | 0.77 (0.82) | 0.78 (0.83) | 0.72 (0.73) |
| Genscan | 0.61 (0.78) | 0.87 (0.80) | 0.90 (0.84) | 0.76 (0.86) |
| HMMgene | 0.75 (0.78) | 0.81 (0.85) | 0.83 (0.87) | 0.78 (0.81) |
| Morgan | 0.43 (0.43) | 0.66 (0.57) | 0.65 (0.56) | 0.39 (0.39) |
| MZEF | — | 0.59 (0.65) | 0.66 (0.73) | — |

For each program, the proportion of actual signals identified correctly (the upper number) and the proportion of predicted signals that are correct (the lower number) are averaged over all signals belonging to a particular type. The number in parenthesis in the header of each column represents the number of signals of each type in the HMR195 dataset.
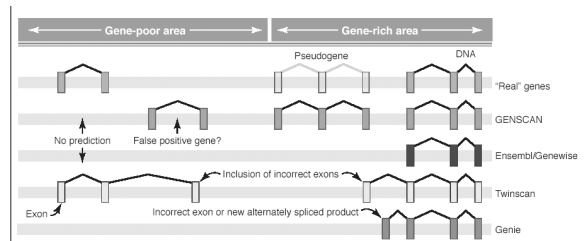
## GenomeScan



Yeh, Lim, and Burge, Genome Research 11:803-816, 2001.

## E. Pennisi, Science 301:1040, 2003



**Never perfect.** No program calls all genes correctly. Some see genes (shown here as coding regions, or exons, connected by bent lines) where there are none; some miss a gene altogether; and some don't put all the gene's parts in the right places.