

# **ERROR CORRECTION**

**30 SEPT 15**

# ANNOUNCEMENTS

# FINAL PROJECT

- Groups of 3-5. No more than 1 grad per group.
- Can use data from published manuscript, downloaded from ENA, your own data.
- I approve topic (Oct14)
- Presentations last week of class
- Written paper due last day.

# FINAL PROJECT EXAMPLES

Genome

Transcriptome

Expression

Comparative Genomics

# EXAM REVIEW

type	query	database
blastn	nt	nt
blastx	Trans(nt)	prot
tblastx	Trans(nt)	Trans(nt)
blastp	prot	prot
tblastn	prot	Trans(nt)

# BLAST

## 1. Build Lookup table

**Preprocess:** Build a *lookup table* of size  $|\Sigma|^w$  for all  $w$ -length words in  $D$

$$\Sigma = \{A, C, G, T\}$$

$$w = 2$$

→  $4^2 (=16)$  entries in lookup table

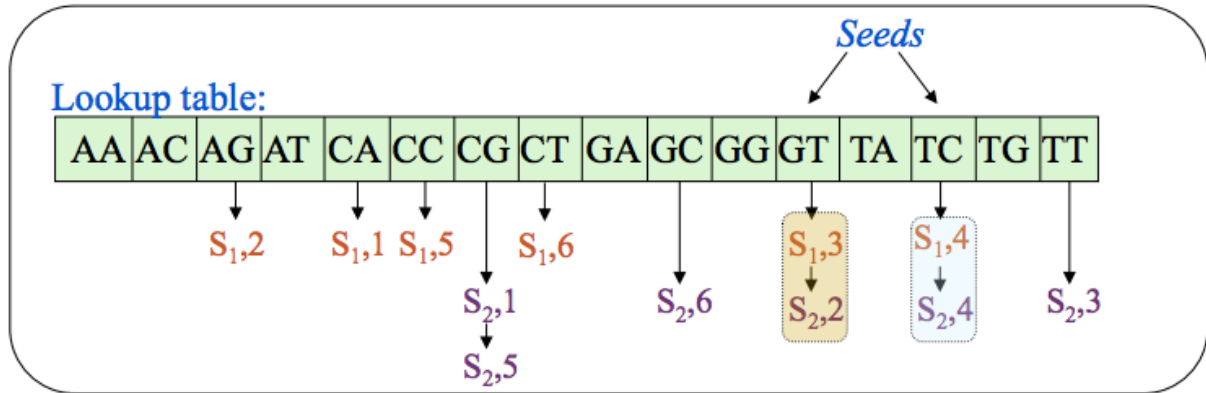
Lookup table:

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----

# BLAST

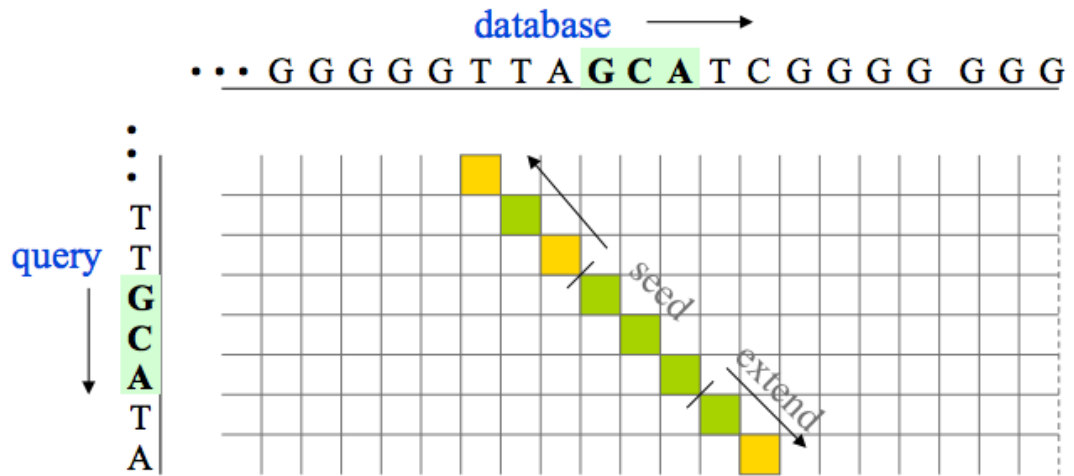
## 2. Filter low complexity and Identify Seeds

1 2 3 4 5 6 7  
S<sub>1</sub>: CAGTCTCT  
S<sub>2</sub>: CGTTCGC



# BLAST

## 3. Bidirectional extension – (Smith Waterman algorithm)





# BLAST

Stats

$$E = Kmne^{-\lambda S}$$

# BLAST

Stats

$$p = 1 - e^{-E}$$

# ERROR CORRECTION



# ERROR CORRECTION

TATACAAATTCGTTTTATGAAAACTCCTAAAAGCAAACAATTTACCAACAAATCCTTGCAACGAAATAACCGATTCATTTAAGCATTCGTCCTATTTATACAAATTCGTTTTATGAAAACTCCTAAAAGCAAACAATTTACCAACAAATTCCTTGCAACGAAATAACCGATTCATTTAAGCATTCG



ACTGTCATTCGGACTA  
ACTGTCATTCGGACTA  
ACTGTCATTCGGCCTA  
ACTGTCATTCGGACTA  
ACTGTCATTCGGACTA  
ACTGTCATTCGGACTA  
ACTGTCATTCGGACTA  
ACTGTCATTCGGACTA  
ACTGTCATTCGGACTA  
ACTGTCATTCGGACTA  
ACTGTCATTCGGACTA

Consensus= ACTGTCATTCGGACTA

# ERROR CORRECTION

TATACAAATTCGTTTTATGAAAACTCCTAAAAGCAAACATAATTACCAACAAATCCTTGCAACGAAATAACCGATTCTATTTAAGCATTCGTCCTATTTATACAAATTCGTTTTATGAAAACTCCTAAAAGCAAACATAATTACCAACAAATCCTTGCAACGAAATAACCGATTCTATTTAAGCATTCG



ACTGTCATTCGGACTA  
ACTGTCATTCGGACTA  
ACTGTCATTCGGCCTA  
ACTGTCATTCGGACTA  
ACTGTCATTCGGACTA  
ACTGTCATTCGGACTA  
ACTGTCATTCGGACTA  
ACTGTCATTCGGACTA  
ACTGACATTCGGACTA  
ACTGACATTCGGACTA  
ACTGTCATTCGGACTA

Consensus= ACTG{A,T}CATTCGGACTA

# ERROR CORRECTION

TATACAAATTCGTTTTATGAAAACTCCTAAAAGCAAACAATTTACCAACAAATCCTTGCAACGAAATAACCGATTCATTTAAGCATTCGTCCTATTTATACAAATTCGTTTTATGAAAACTCCTAAAAGCAAACAATTTACCAACAAATTCCTTGCAACGAAATAACCGATTCATTTAAGCATTCG



ACTGTCATTCGGACTA  
ACTGTCATTCGGACTA  
ACTGTCATTCGGCCTA  
ACTGTCATTCGGACTA  
ACTGTCATTCGGACTA  
ACTGTCATTCGGACTA  
ACTGTCATTCGGACTA  
ACTGTCATTCGGACTA  
**GCTGATAAC**CGGACTA  
ACTG**A**CATTCGGACTA  
ACTGTCATTCGGACTA

Consensus= ACTGTCATTCGGACTA

# ERROR CORRECTION

Hamming Distance:

[http://en.wikipedia.org/wiki/Hamming\\_distance](http://en.wikipedia.org/wiki/Hamming_distance)

# ERROR CORRECTION

3 different strategies



# ERROR CORRECTION

Kmer-spectra based

# ERROR CORRECTION

Suffix tree based

# ERROR CORRECTION

MSA based

# ERROR CORRECTION

Evaluation of Correction