# DNALC *Live*

This is an experiment; give us feedback on what you would like to see!

# DNALC *Live*

- Provide genetics, molecular biology, and bioinformatics learning resources

- Laboratory and computer demos, short online courses for middle school, high school, and the general public

- Interviews with scientists, help for teachers

- At-home activities, social media contests, and more

# DNALC Website and Social Media

## dnalc.cshl.edu



## dnalc.cshl.edu/dnalc-live

# Barcoding Bioinformatics
# Part II

# Who is this course for?

- Audience(s): US AP Biology (high school grades 10-12) AND Intro undergraduate biology

- Format: 3 sessions (1 per week); ~ 45 minutes each

- Exercises: Follow along with our online bioinformatics tool DNA Subway

- Learning resources: Slides and packet available (teachers can also request the teacher edition)

# Course Learning Goals

- Learn how DNA can be used to identify unknown organisms

- Understand how we obtain DNA Sequence and access its quality

- Use BLAST* to compare an unknown DNA Sequence to known sequences

- Compare DNA Sequences using phylogenetics

    *AP Bio (Lab 3 – Comparing DNA Sequences)

# Lab Setup

- We will be using DNA Subway – You can get a free account at **cyverse.org** (optional)

# Steps for today's session

- Recap on our experimental dataset

- Review of sequence quality

- Sequence cleaning and pairing

- Introduction to BLAST

# Recap of the dataset

# Steps to DNA Barcoding



**Organism is sampled**

**DNA is extracted**

**"Barcode" amplified**

ACGAGTCGGTAGCTGCCCTCTGACTGCATCGAA
TTGCTCCCCTACTACGTGCTATATGCGCTTACGAT
CGTACGAAGATTTATAGAATGCTGCTACTGCTCC
CTTATTCGATAACTAGCTCGATTATAGCTACGATG

Sequenced DNA is compared with DNA in a barcode database

# Example barcoding experiment



Mary Acheampong,
Bobby Glover, and Marisa
VanBrakle

Mentor: Allison Granberry
Hostos-Lincoln Academy of
Science,
The Bronx

2012 UBP Grand Prize Winners

# Example barcoding experiment



**Different Forms of Samples Tested**

| Capsule | Tablet | Tea | Seed | Leaf |
|---------|--------|-----|------|------|

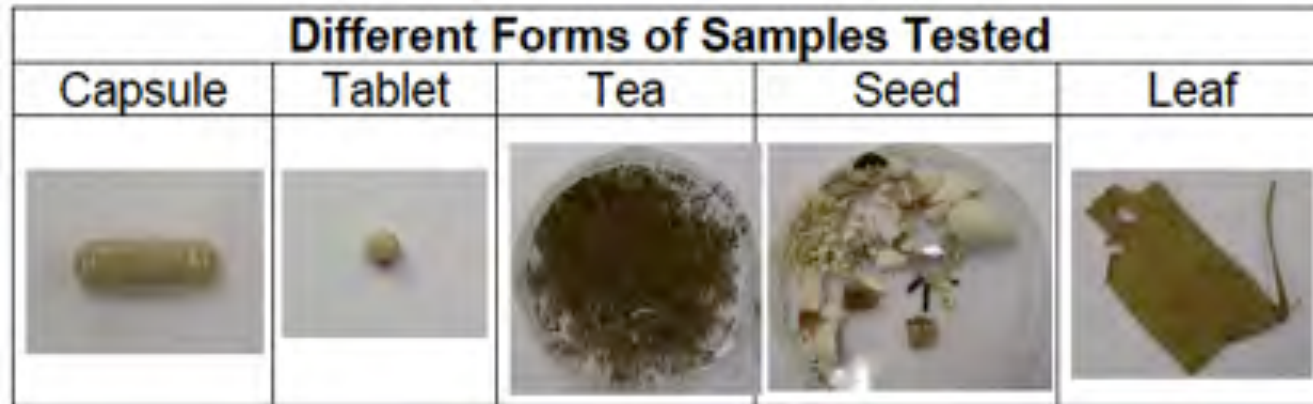| Sample Letter | Form | DNA Expected | DNA Results |
|---------------|------|--------------|-------------|
| A | Capsule | Ginkgo biloba | Rice: *Oryza rufipogon* |
| B | Capsule | Ginkgo biloba | Rice: *Oryza rufipogon* |
| C | Capsule | Ginkgo biloba | Rice: *Oryza rufipogon* |
| D | Tablet | Ginkgo biloba | No sequence available. |
| E | Capsule | Ginkgo biloba | Rice: *Oryza rufipogon* |
| F | Liquid | Ginkgo biloba | No sequence available |
| G | Capsule | Ginkgo biloba | No sequence available |
| H | Tea | Ginkgo biloba | Other rbcL DNA present but not *Mentha piperita* |
| L | Capsule | Ginkgo | Rice: *Oryza* |

| Aedes adult | Anopheles adult | Culex adult |
|---|---|---|
|  |  |  |
| By Muhammad Mahdi Karim - Own work, GFDL 1.2, https://commons.wikimedia.org/w/index.php?curid=11185617 | By Jim Gathany - (PHIL), ID #5814. https://commons.wikimedia.org/w/index.php?curid=799284 | By Muhammad Mahdi Karim - Own work, GFDL 1.2, https://commons.wikimedia.org/w/index.php?curid=7673048 |
| Aedes larva | Anopheles larva | Culex larva |
|  |  |  |
| Photograph by Michele M. Cutwa, University of Florida. | | Photograph by Michelle Cutwa-Francis, University of Florida. |

Aedes

Anopheles

Culex

© 2000 Richard C. Russell

Cold Spring Harbor Laboratory
DNA Learning Center

# Why does this matter?

**_Aedes_:**

- Chikungunya
- Dengue fever
- Lymphatic filariasis
- Rift Valley fever
- Yellow fever
- Zika

**_Anopheles_:**

- Malaria
- Lymphatic filariasis

**_Culex_:**

- Japanese encephalitis
- Lymphatic filariasis
- West Nile fever

# Experimental components/design

## Materials

- We have DNA from unknown mosquito samples
- We can obtain DNA from known samples

# Experimental components/design

## Materials

- We have DNA from unknown mosquito samples
- We can obtain DNA from known samples

## Hypothesis

- We can use computational methods (BLAST/phylogenetic analysis) to infer the species

# Experimental components/design

## Materials

- We have DNA from unknown mosquito samples
- We can obtain DNA from known samples

## Hypothesis

- We can use computational methods (BLAST/phylogenetic analysis) to infer the species
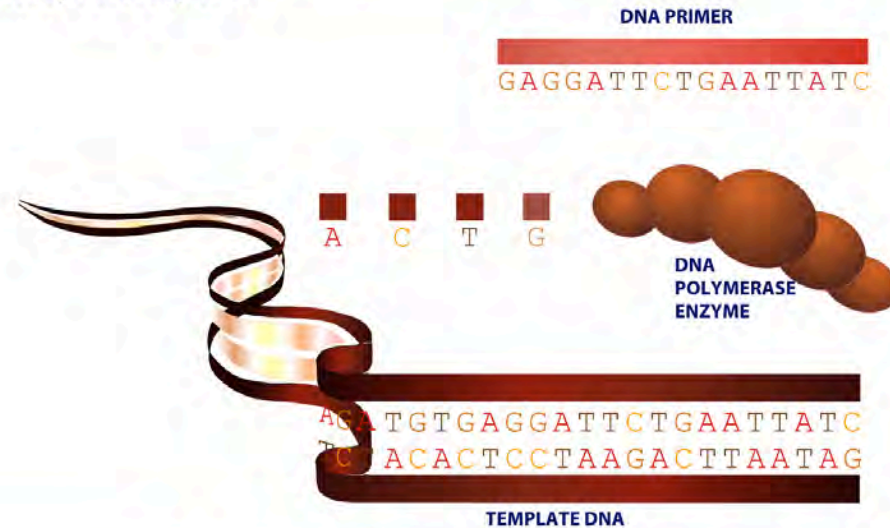
## Controls

- We have sensitivity controls (sequence quality, BLAST parameters)
- We have outgroup sequences (non-mosquito, negative controls) and known samples (positive controls)
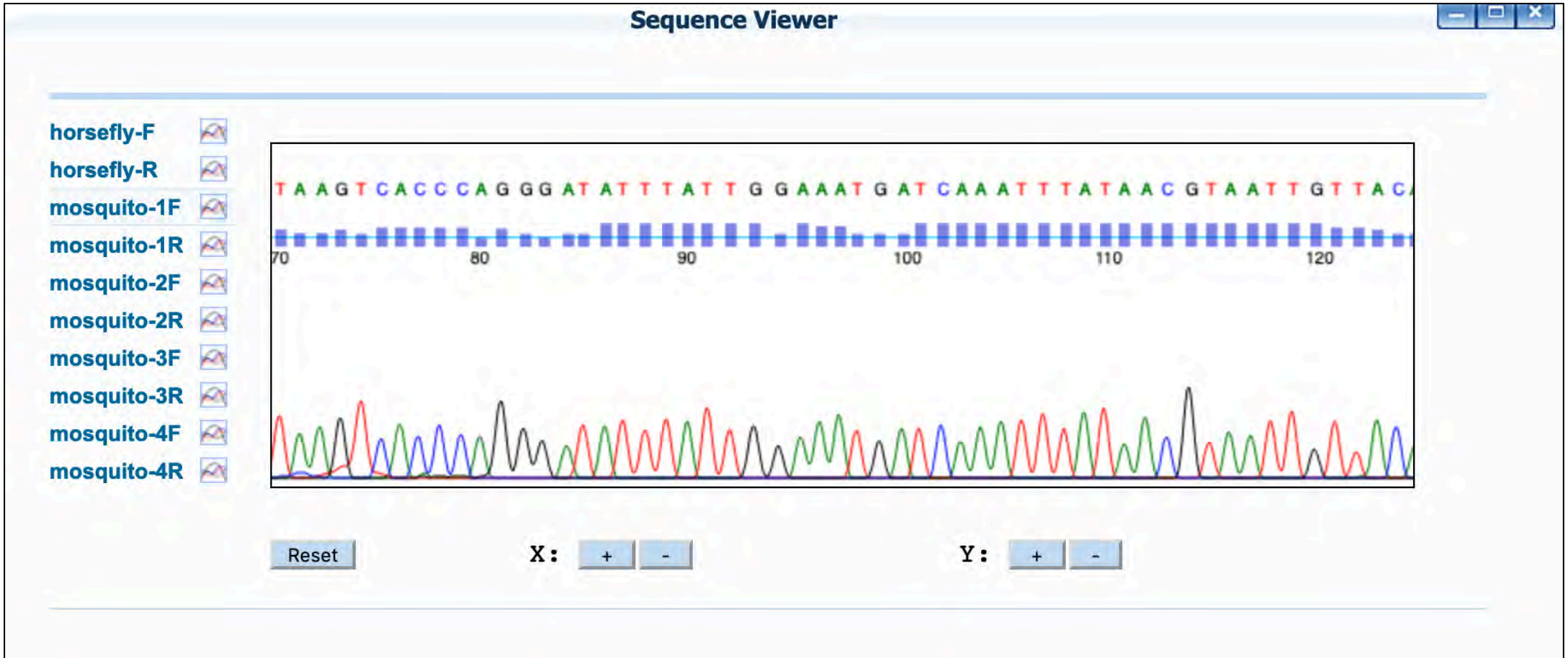
# Review of sequencing and quality
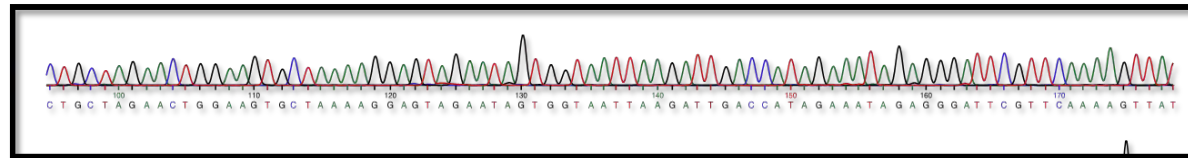
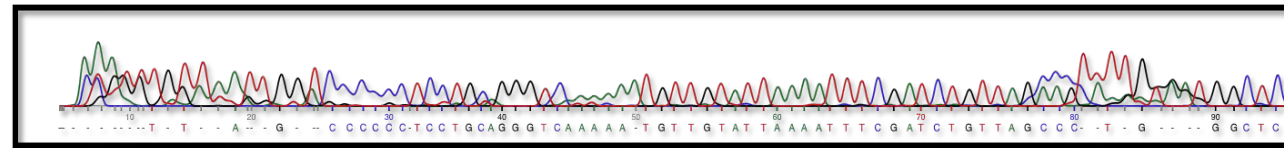# DNA Sequencing

# Chromatogram/Electropherogram
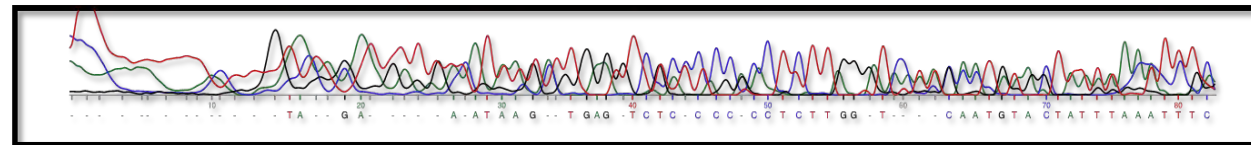
# Some sequence examples…

High Quality Sequence



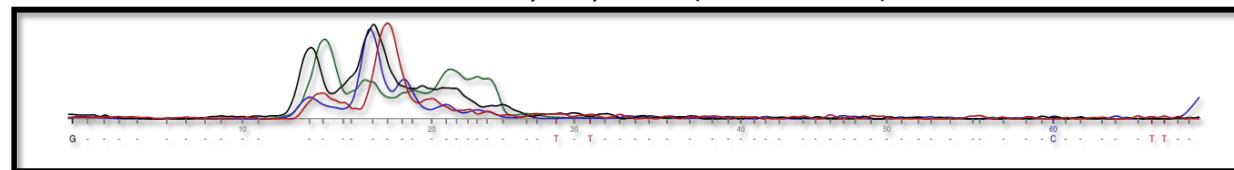Acceptable Quality Sequence



Low Quality Sequence (multiple base calls per position)



Low Quality Sequence (no base calls)

# Phred scores...

| Phred Score | Error (bases miscalled) | Accuracy |
|---|---|---|
| 10 | 1 in 10 | 90% |
| 20 | 1 in 100 | 99% |
| 30 | 1 in 1,000 | 99.9% |
| 40 | 1 in 10,000 | 99.99% |
| 50 | 1 in 100,000 | 99.999% |

Cold Spring Harbor Laboratory
DNA LEARNING CENTER

# If 99% was good enough

If things only work correctly 99.9% of the time...

- 12 newborns will be given to the wrong parents daily.
- 114,500 mismatched pairs of shoes will be shipped/year.
- 18,322 pieces of mail will be mishandled/hour.
- 2,000,000 documents will be lost by the IRS this year.
- 2.5 million books will be shipped with the wrong covers.
- Two planes landed at Chicago's O'Hare airport will be unsafe every day.
- 315 entries in Webster's Dictionary will be misspelled.
- 20,000 incorrect drug prescriptions will be written this year.
- 880,000 credit cards in circulation will turn out to have incorrect
- cardholder information on their magnetic strips.
- 103,260 income tax returns will be processed incorrectly during the year.
- 5.5 million cases of soft drinks produced will be flat.
- 291 pacemaker operations will be performed incorrectly.
- 3056 copies of tomorrow's Wall Street Journal will be missing one of the three sections.

CSH Cold Spring Harbor Laboratory
DNA LEARNING CENTER

# A note on controls

At what temperature does ice ($H_2O$) + Chemical "X" melt?

# A note on controls

**Positive control:** What does the effect look like if present?

# A note on controls

**Positive control:** What does the effect look like if present?

**Negative control:** What does the effect look like if absent?

# A note on controls

**Positive control:** What does the effect look like if present?

**Negative control:** What does the effect look like if absent?

**Sensitivity control:** Across what range of values can I measure the effect?

# A note on controls



**Positive control**     **Negative control**     **Sensitivity control**

CSH Cold Spring Harbor Laboratory
DNA LEARNING CENTER

# A note on controls



**Positive control**

**Negative control**

**Sensitivity control**

CSH Cold Spring Harbor Laboratory
DNA LEARNING CENTER

# Phred are our measure of quality (signal/noise)



**Lower score = more noise than signal**

# Bi-directional sequencing

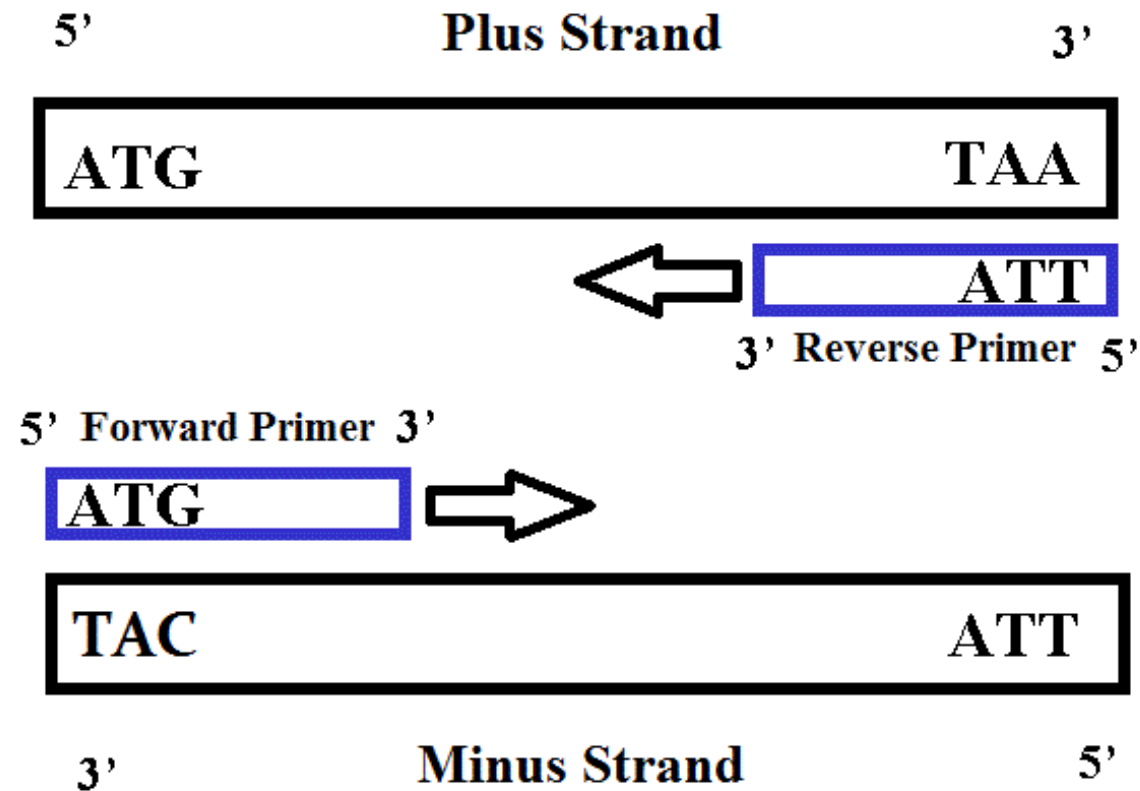Cold Spring Harbor Laboratory
DNA LEARNING CENTER

# Reverse complementation

CSH Cold Spring Harbor Laboratory
DNA LEARNING CENTER

# Reverse complementation



- Reverse: change nt. sequence from (5'→3') to (3'→5')

- Complement with the reversed sequence

```
              5' CTCCAAGCTCCAAGCTCCAG 3'
Reverse:      5' GACCTCGAACCTCGAACCTC 3'
Complement:   5' CTGGAGCTTGGAGCTTGGAG 3'
```
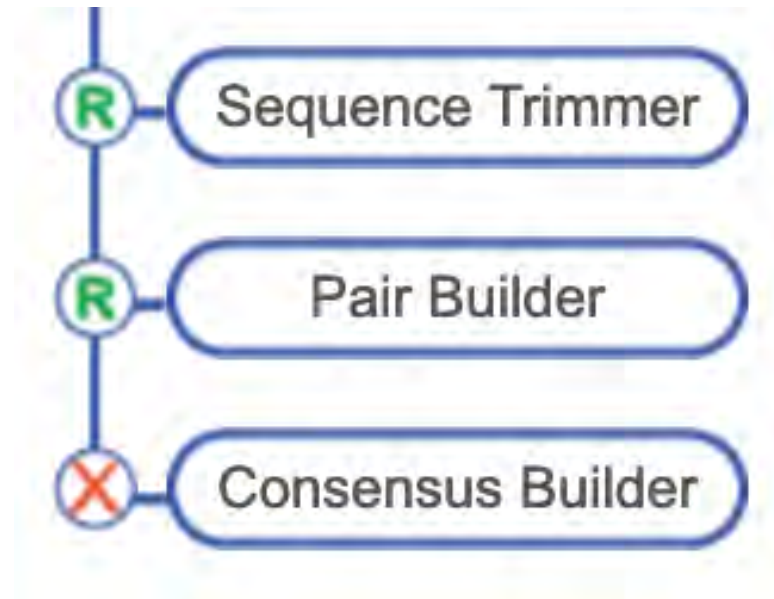
CSH Cold Spring Harbor Laboratory
DNA LEARNING CENTER

# Clean up and consensus

# Introduction to BLAST

# Basic Local Alignment Search Tool

- An algorithm for searching a database of sequences

# Basic Local Alignment Search Tool

- An <u>algorithm</u> for searching a <u>database</u> of sequences

- "Google for DNA" (although works with any biological sequence, and started before Google ~1985)

# Basic Local Alignment Search Tool

- An <u>algorithm</u> for searching a <u>database</u> of sequences

- "Google for DNA" (although works with any biological sequence, and started before Google ~1990 vs 1998)

- NCBI is the most popular interface, but this is software that can be run anywhere (including Subway)

# Warning: Analogy
(useful for discussion but not the whole picture)

# BLAST algorithm analogy

Query sequence
**ACTGACATCGGGGTGCTACG**



**Database**

# BLAST algorithm analogy

Query sequence
**ACTGACATCGGGGGTGCTACG**

**Database**

# BLAST algorithm analogy
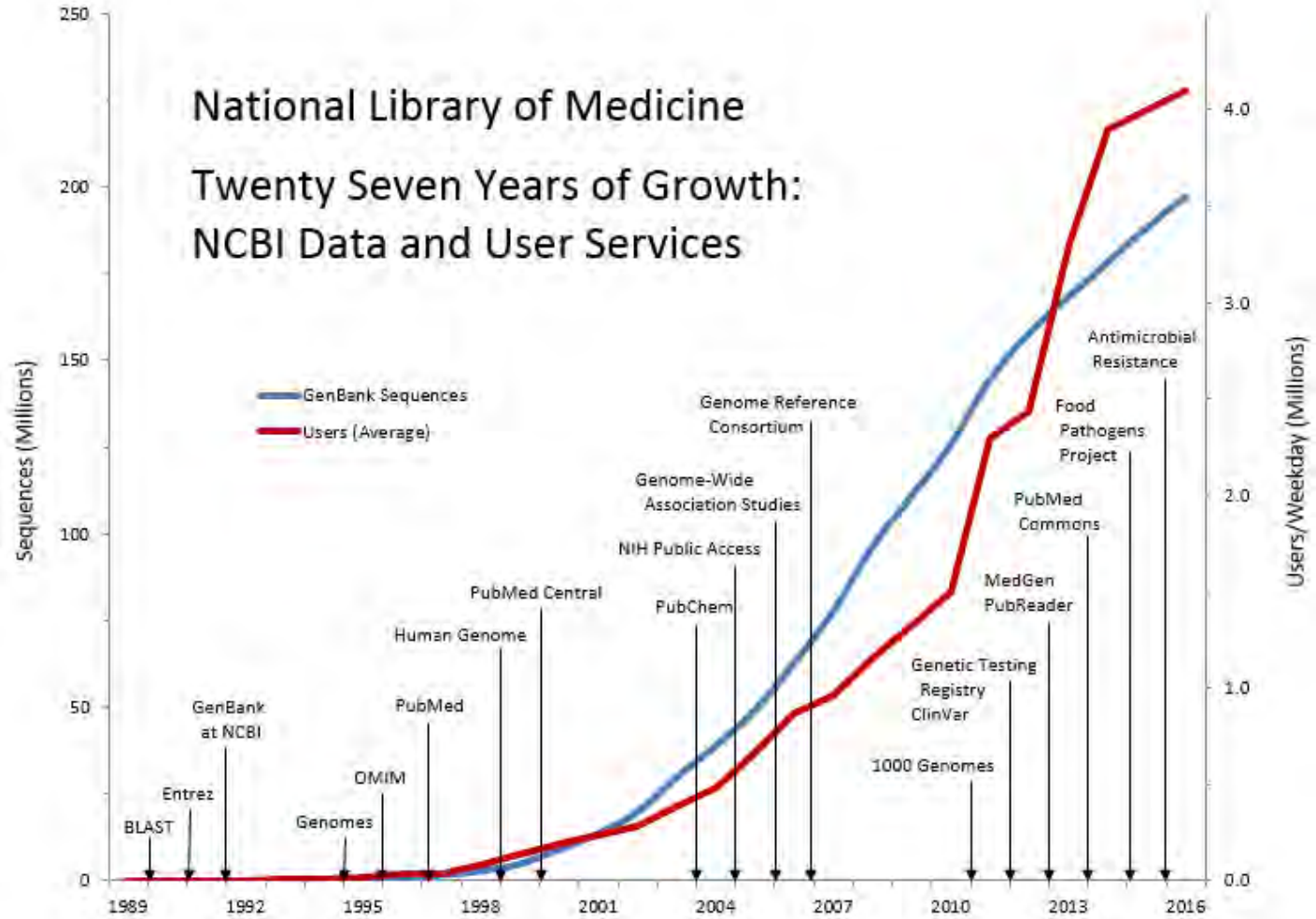
CSH Cold Spring Harbor Laboratory
DNA LEARNING CENTER

# BLAST algorithm analogy – searching by "word"

Break the *Query sequence*
Into "words" (k-mers)

**ACT GAC ATC GGG GTG CTA CG**



**Database**

Cold Spring Harbor Laboratory
DNA LEARNING CENTER

# BLAST algorithm analogy – searching by "word"

Break the *Query sequence*
Into "words" (k-mers)

ACT GAC ATC GGG GTG CTA CG

ACT... TCT ... GCT ...



**Database**

Cold Spring Harbor Laboratory
DNA LEARNING CENTER

# Let's BLAST a sequence
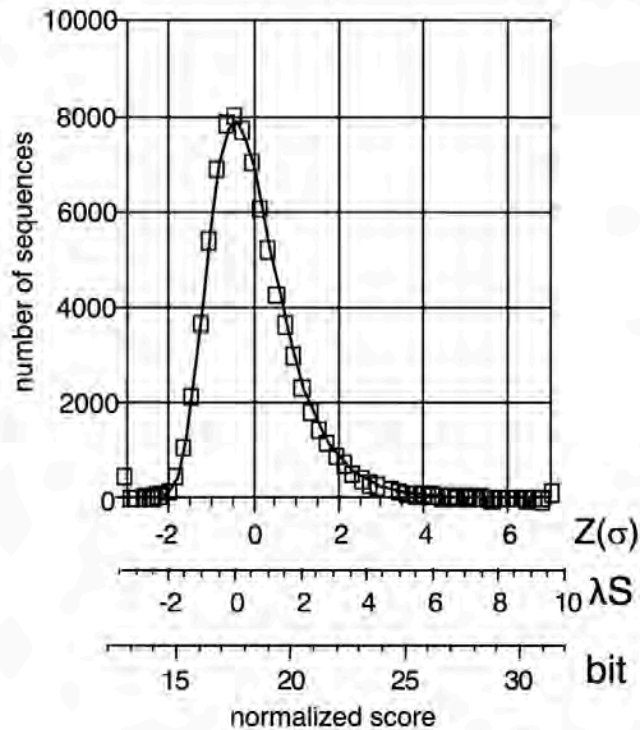
>mosquito-1F
CTTTAAGTATATTAATTCGTGCTGAATTAAGTCACCCAGGGATATTTAT
TGGAAATGATCAAATTTATAACGTAATTGTTACAGCTCATGCATTTATT
ATAATTTTTTTTATAGTAATACCAATTATAATTGGAGGATTTGGAAATT
GATTAGTTCCTTTAATATTAGGAGCTCCTGATATAGCATTTCCTCGAAT
AAATAATATAAGTTTTTGAATATTACCTCCTTCTTTAACTCTACTACTTT
CTAGTTCAATAGTAGAAAATGGAGCAGGGACAGGATGAACAGTTTA
TCCTCCTCTTTCATCAGGAACAGCACATGCTGGAGCTTCTGTTGATTT
AGCAATTTTCTCTCTTCATTTAGCAGGGATTTCATCTATTTTAGGAGC
AGTAAATTTTATTACTACTGTTATTAATATACGATCATCTGGAATTACTT
TAGATCGATTACCTTTATTTGTTTGATCTGTAGTAATTACTGCTATTTTA
TTACTTTTATCTCTTCCTGTATTAGCTGGAGCTATTACTATATTATTAACT
GATCGAAATTTAAATACTTCCTTCTTTGACCCAATTGGAGGAGGAGA

## https://blast.ncbi.nlm.nih.gov/Blast.cgi

# BLAST and controls

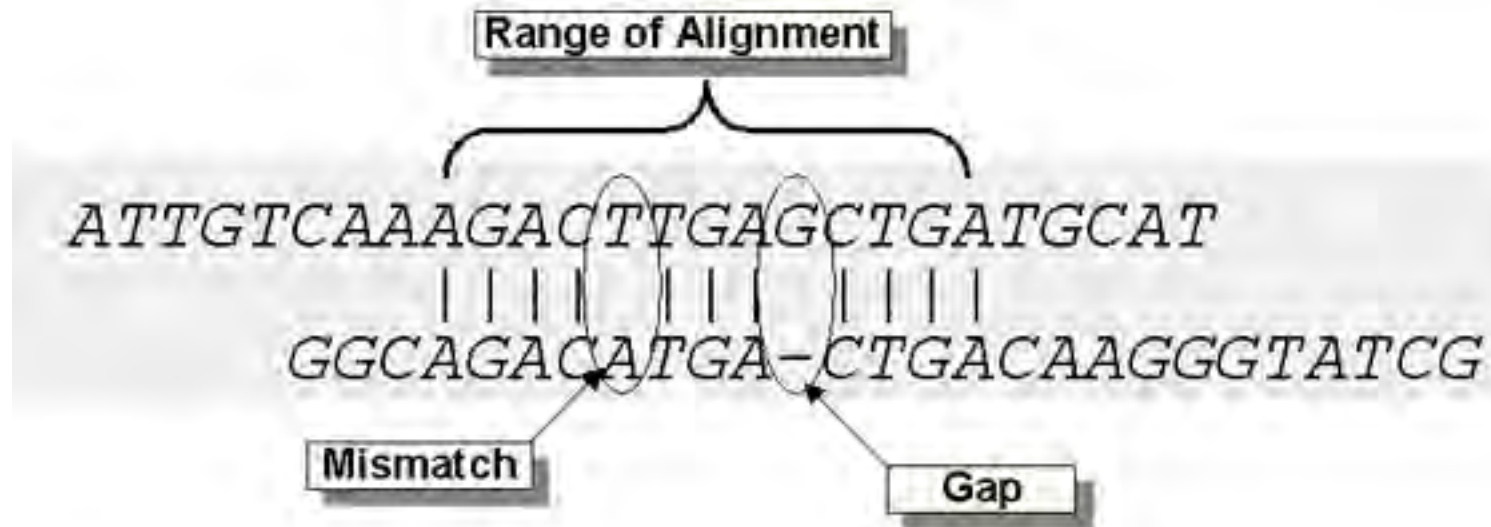# BLAST algorithm analogy – searching by "word"

The *Query sequence*
Is aligned to a *Subject* (a sequence in the database)

Q: ACTGAC–ATCGGGGTGCTACG

| | | | | | | | | | | | | | | | | | | |

S:  ACTGACCATCGGAGTGCTACG

# BLAST algorithm analogy – alignment



$$S = \sum(\text{identities, mismatches}) - \sum(\text{gap penalties})$$

$$\text{Score} = \text{Max}(S)$$

CSH Cold Spring Harbor Laboratory
DNA LEARNING CENTER

# Let's do a BLAST

# Some BLAST definitions

- **Max Score**: Highest alignment score (according to a formula)

# Some BLAST definitions

- **Max Score**: Highest alignment score (according to a formula)

- **Query Cover**: % of the query length included in aligned segment

# Some BLAST definitions

- **Max Score**: Highest alignment score (according to a formula)

- **Query Cover**: % of the query length included in aligned segment

- **E value**: The number of alignments expected by chance with the calculated score or better

# Some BLAST definitions

- **Max Score**: Highest alignment score (according to a formula)

- **Query Cover**: % of the query length included in aligned segment

- **E value**: The number of alignments expected by chance with the calculated score or better

- **Per. Identity**: Highest % identity for a set of aligned segments to the same subject sequence.

No*

# (Some) Limitations to BLAST

- **Homology**: BLAST is trying to indicate which homologous (related by ancestry) sequences are found in the database

# (Some) Limitations to BLAST

- **Homology**: BLAST is trying to indicate which homologous (related by ancestry) sequences are found in the database

- **Data base coverage**: BLAST returns its best result; that is not guaranteed to be the true result

# (Some) Limitations to BLAST

- **Homology**: BLAST is trying to indicate which homologous (related by ancestry) sequences are found in the database

- **Data base coverage**: BLAST returns its best result; that is not guaranteed to be the true result

- **Locus resolution**: Barcodes are often good for **genus**-level resolution
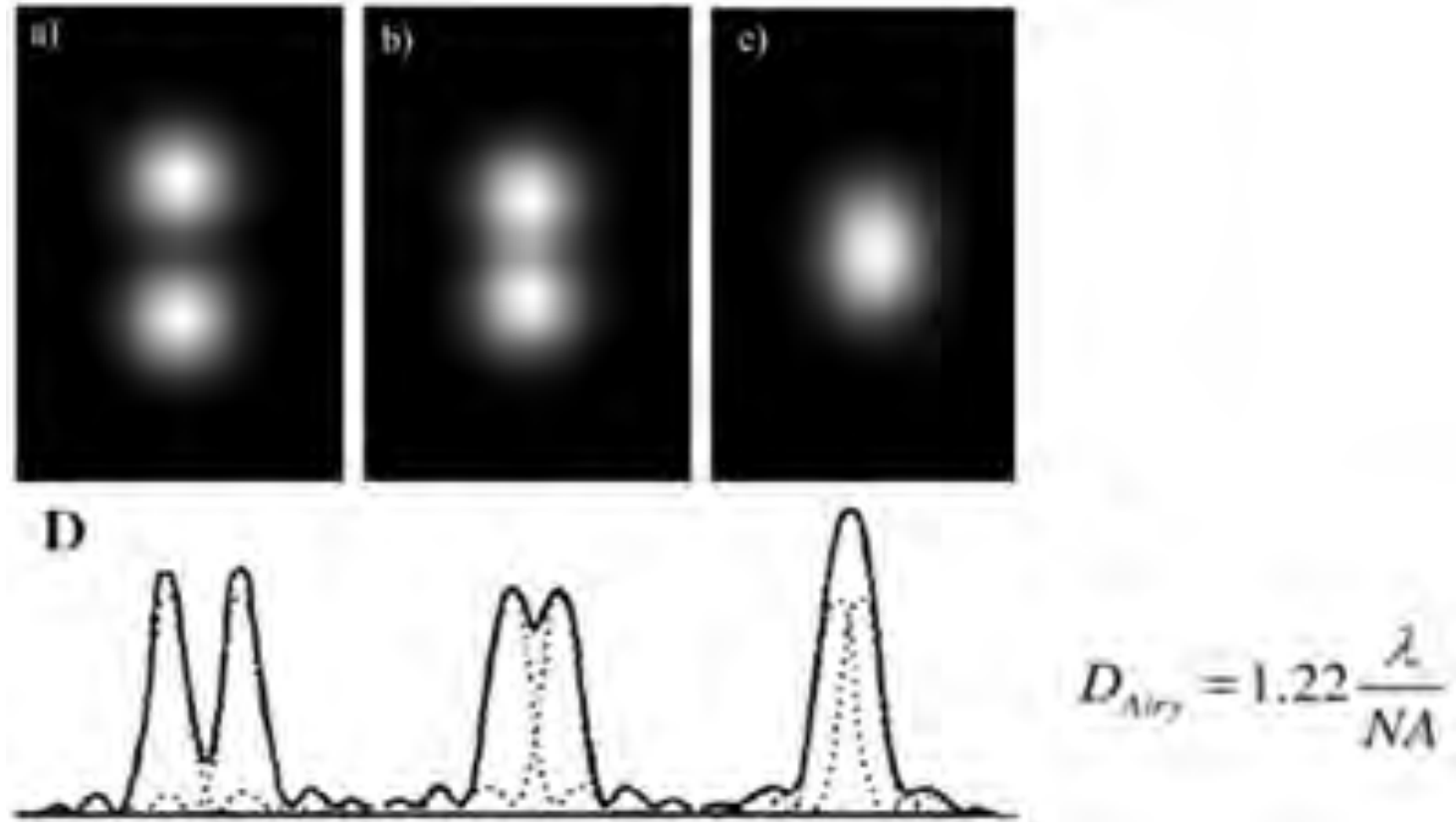
# A note on resolution (and controls)



$$D_{Airy} = 1.22 \frac{\lambda}{NA}$$

Cold Spring Harbor Laboratory
DNA LEARNING CENTER

# Next time:

# Multiple sequence alignments and phylogenetics

# DNALC Website and Social Media

## dnalc.cshl.edu



## dnalc.cshl.edu/dnalc-live