



Cold Spring Harbor Laboratory
DNA LEARNING CENTER

Barcode Bioinformatics Part II

Jason Williams

Cold Spring Harbor Laboratory, DNA Learning Center

williams@cshl.edu



[@JasonWilliamsNY](https://twitter.com/JasonWilliamsNY)



Cold Spring Harbor Laboratory
DNA LEARNING CENTER

Barcoding Bioinformatics

Part II

(Sequence cleaning and BLAST)

Recap of the dataset

Steps to DNA Barcoding



Organism is sampled

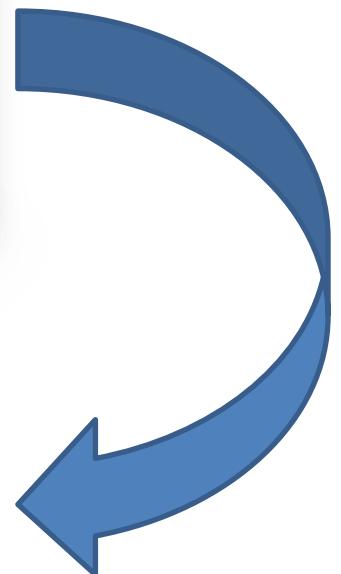


DNA is extracted



“Barcode” amplified

ACGAGTCGGTAGCTGCCCTTGACTGCATCGAA
TTGCTCCCTACTACGTGCTATATGCGCTTACGAT
CGTACGAAGATTATAGAACATGCTGCTACTGCTCC
CTTATTGATAACTAGCTGATTATAGCTACGATG



Sequenced DNA is compared with DNA in a barcode database



Example barcoding experiment



Mary Acheampong,
Bobby Glover, and Marisa
VanBrakle

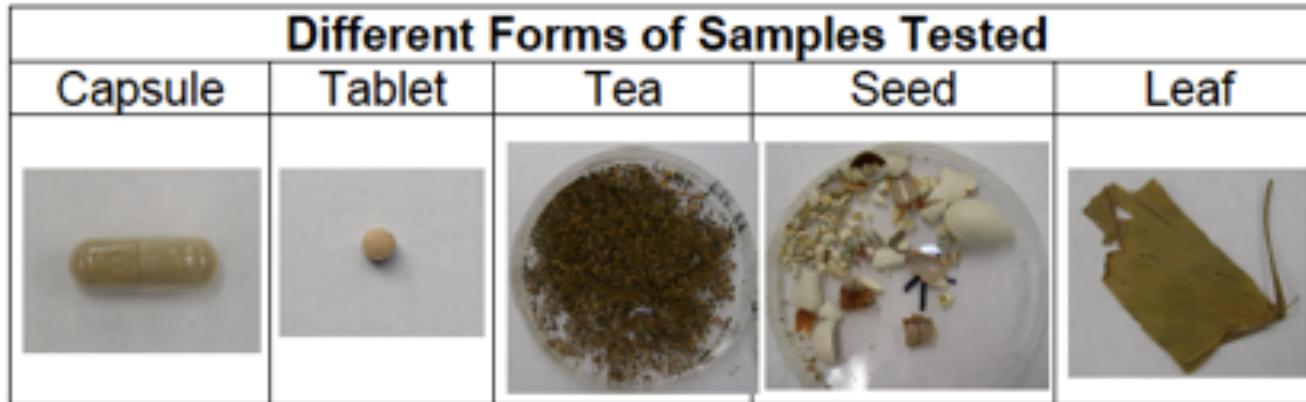
Mentor: Allison Granberry
Hostos-Lincoln Academy of
Science,
The Bronx

2012 UBP Grand Prize Winners



Cold Spring Harbor Laboratory
DNA LEARNING CENTER

Example barcoding experiment



Sample Letter	Form	DNA Expected	DNA Results
A	Capsule	Ginkgo biloba	Rice: <i>Oryza rufipogon</i>
B	Capsule	Ginkgo biloba	Rice: <i>Oryza rufipogon</i>
C	Capsule	Ginkgo biloba	Rice: <i>Oryza rufipogon</i>
D	Tablet	Ginkgo biloba	No sequence available.
E	Capsule	Ginkgo biloba	Rice: <i>Oryza rufipogon</i>
F	Liquid	Ginkgo biloba	No sequence available
G	Capsule	Ginkgo biloba	No sequence available
H	Tea	Ginkgo biloba	Other <i>rbcL</i> DNA present but not <i>Mentha piperita</i>
I	Capsule	Ginkgo	Rice: <i>Oryza</i>



Aedes adult



By Muhammad Mahdi Karim - Own work, GFDL 1.2, <https://commons.wikimedia.org/w/index.php?curid=11185617>

Anopheles adult



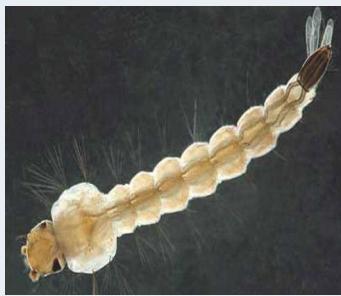
By Jim Gathany - (PHIL), ID #5814. <https://commons.wikimedia.org/w/index.php?curid=799284>

Culex adult



By Muhammad Mahdi Karim - Own work, GFDL 1.2, <https://commons.wikimedia.org/w/index.php?curid=7673048>

Aedes larva



Photograph by Michele M. Cutwa, University of Florida.

Anopheles larva

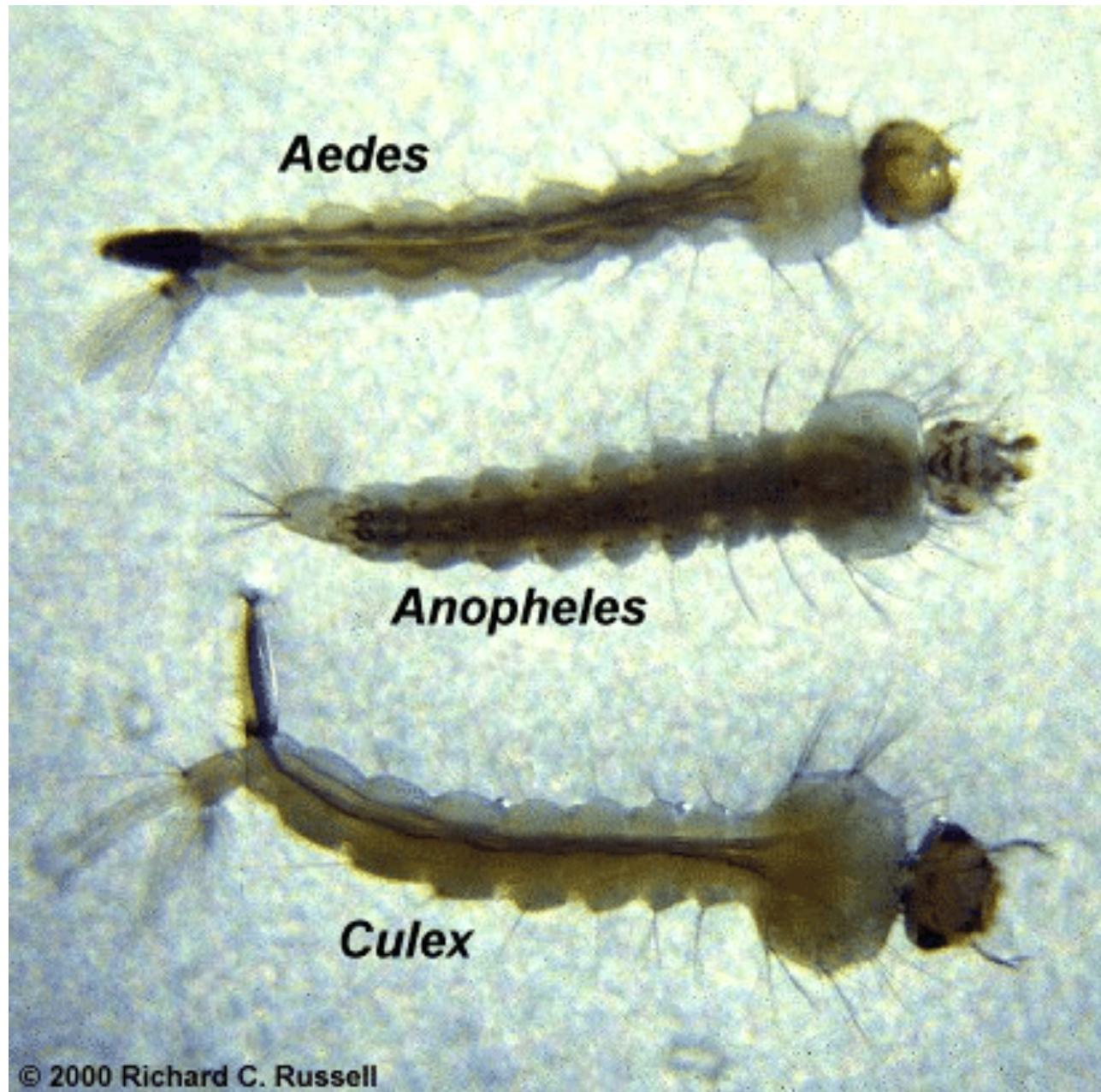


Culex larva



Photograph by Michelle Cutwa-Francis, University of Florida.





© 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)



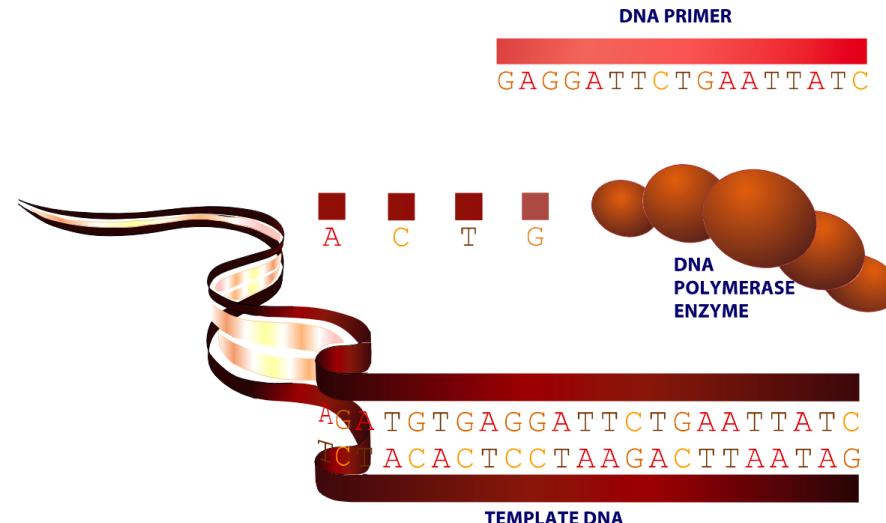
Review of sequencing and quality



DNA Sequencing

Cycle Sequencing

To this mix, we also add a second type of nucleotide; one that has a slightly different chemical formula. These "dideoxynucleotides (ddNTP)" can be recognized by a DNA sequencer.



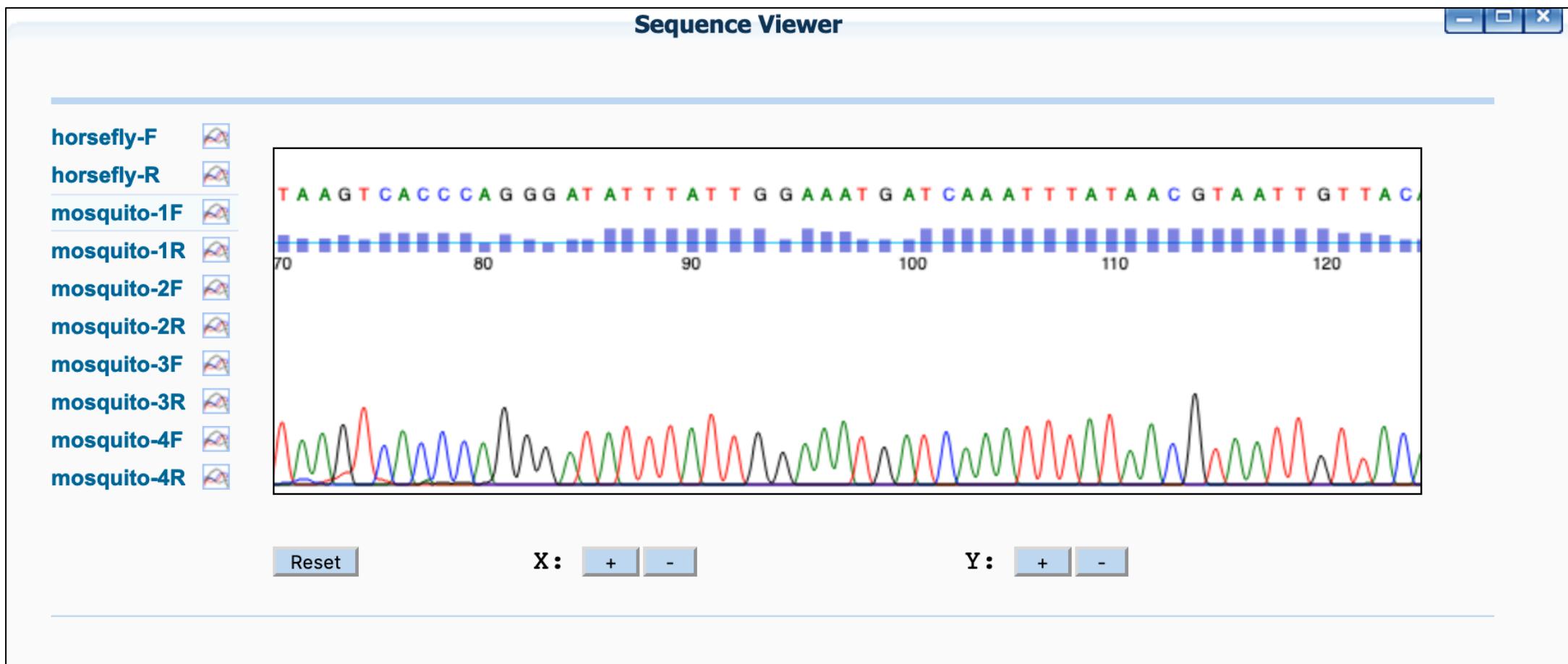
Jump to: 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15

<https://dnalc.cshl.edu/resources/animations/cycseq.html>



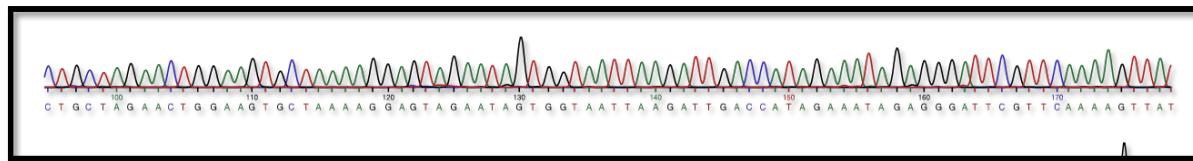
Cold Spring Harbor Laboratory
DNA LEARNING CENTER

Chromatogram/Electropherogram

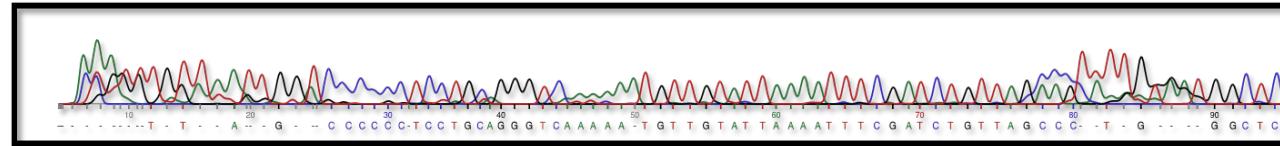


Some sequence examples...

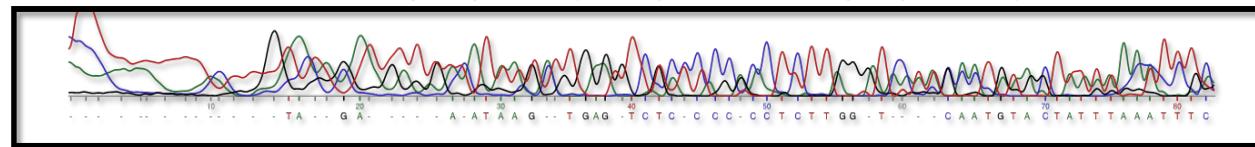
High Quality Sequence



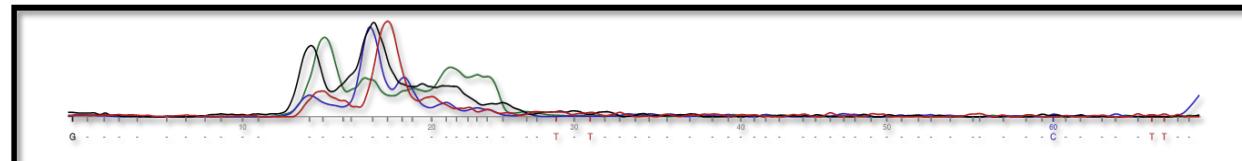
Acceptable Quality Sequence



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%

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.

Photo credit

<http://www.personal.psu.edu/sxt104/class/99percent.html>



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

Photo credits

https://commons.wikimedia.org/wiki/File:Water_in_a_beaker.JPG
<http://www.chem.uiuc.edu/webfunchem/temperature/Temp10.htm>
<https://www.dreamstime.com/photos-images/alcohol-thermometer.html>

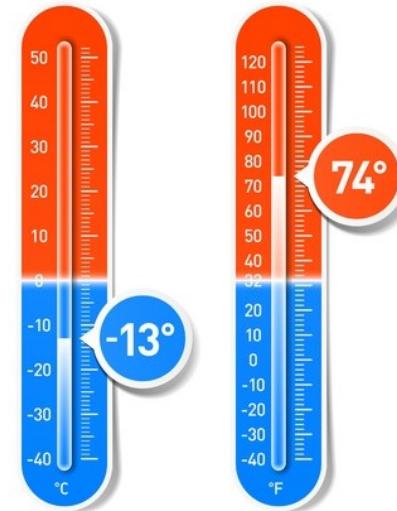
A note on controls



Positive control



Negative control

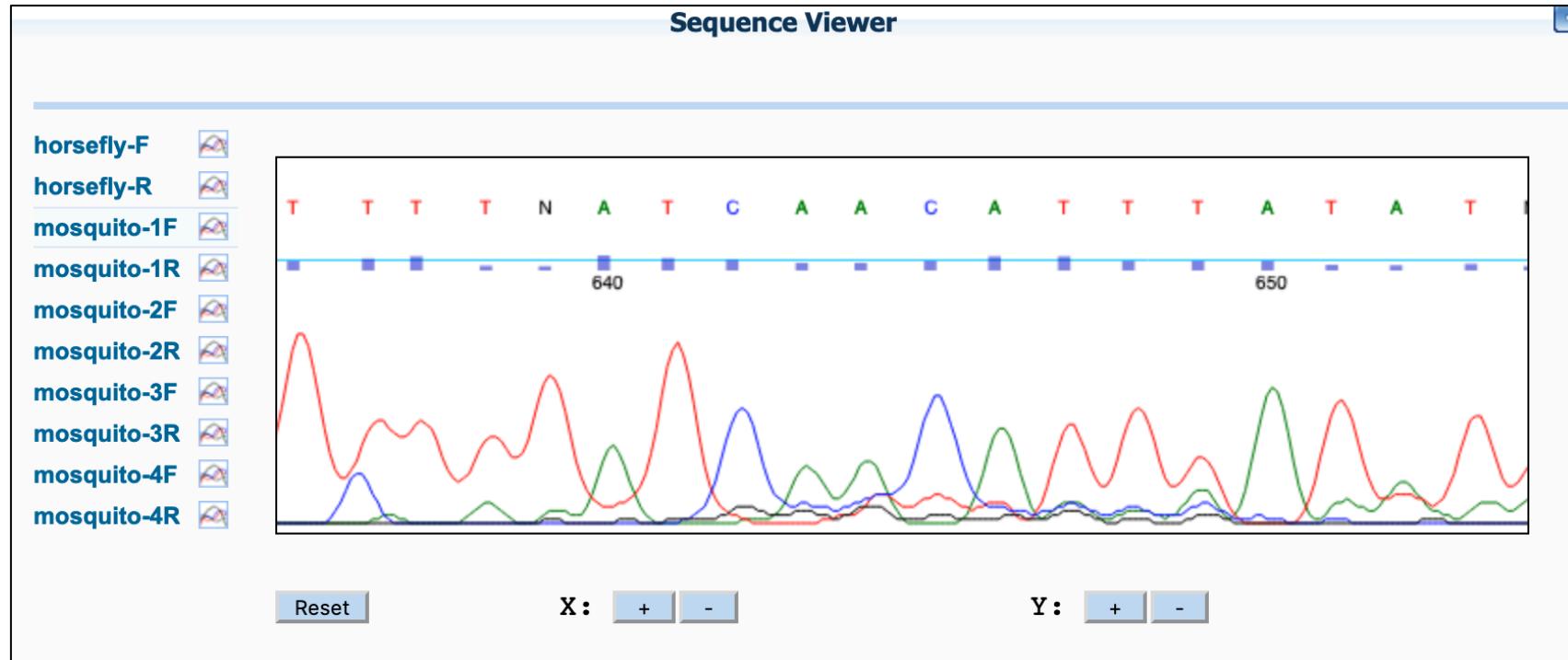


Sensitivity control

Photo credits

https://commons.wikimedia.org/wiki/File:Water_in_a_beaker.JPG
<http://www.chem.uiuc.edu/webfunchem/temperature/Temp10.htm>
<https://en.clipdealer.com/vector/media/A:17494508?>

Phred are our measure of quality (signal/noise)



Lower score = more noise than signal



Bi-directional sequencing

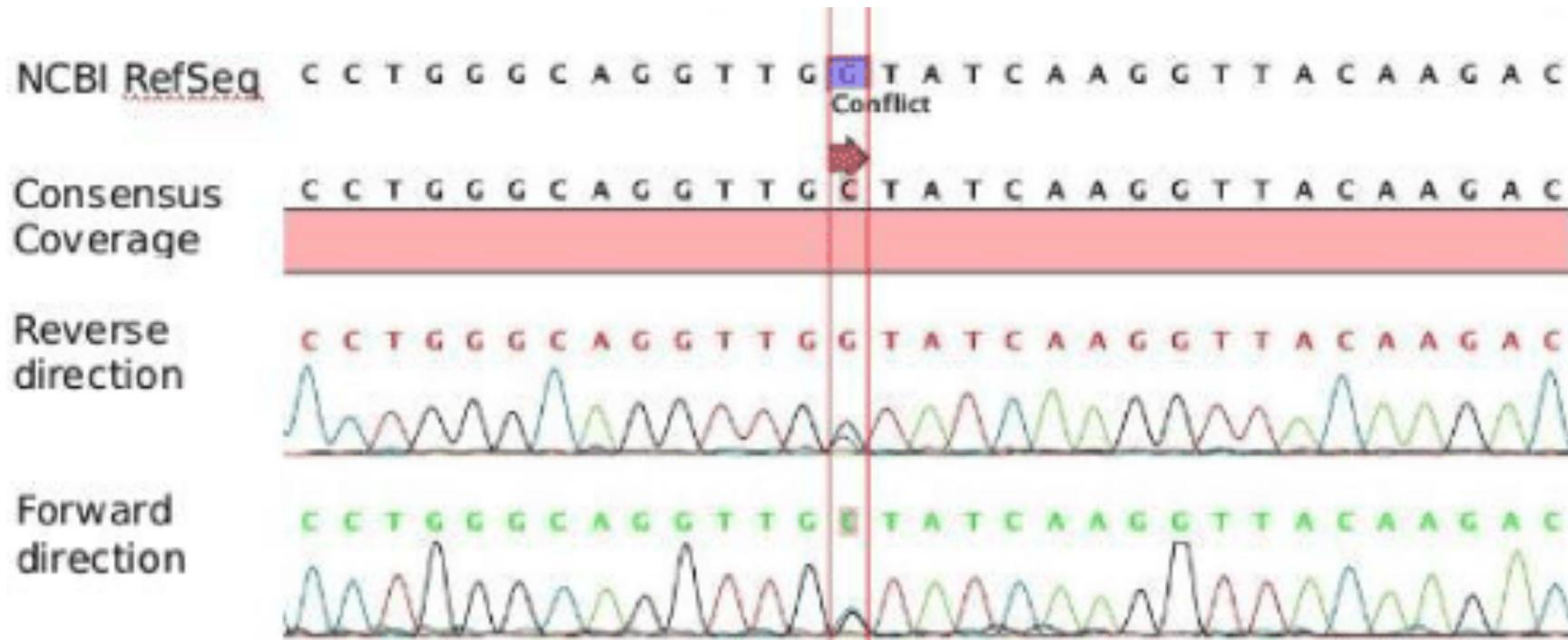


Photo credit

<https://www.omicsonline.org/articles-images/CMBO-2-108-g003.html>



Cold Spring Harbor Laboratory
DNA LEARNING CENTER

Reverse complementation

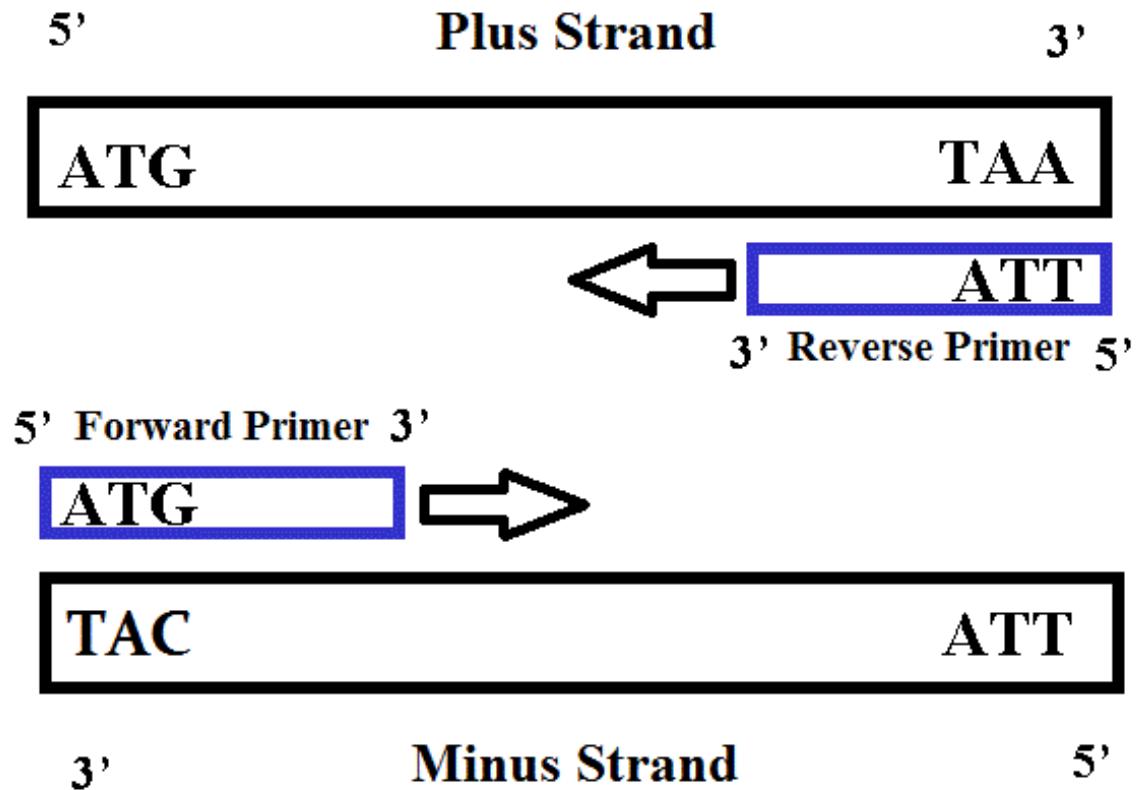


Photo credit

<https://biology.stackexchange.com/questions/56304/manual-primer-design-for-a-gene-on-the-reverse-strand>



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'

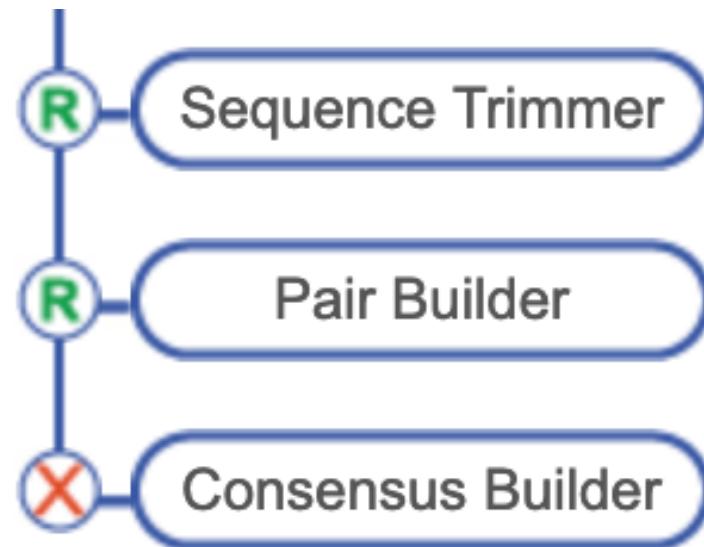
Complement: 5' CTGGAGCTTGGAGCTTGGAG 3'

Photo credit

<https://image.slidesharecdn.com/pcrprimerdesignenglishversion-160317161103/95/pcr-primer-design-english-version-10-638.jpg?cb=1458231192>



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 algorithm for searching a database of sequences
- “Google for DNA” (although works with any biological sequence, and started before Google ~1985)

Basic Local Alignment Search Tool

- An algorithm for searching a database 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



Cold Spring Harbor Laboratory
DNA LEARNING CENTER

BLAST algorithm analogy

Query sequence

ACTGACATCGGGGTGCTACG



Database



Cold Spring Harbor Laboratory
DNA LEARNING CENTER

BLAST algorithm analogy

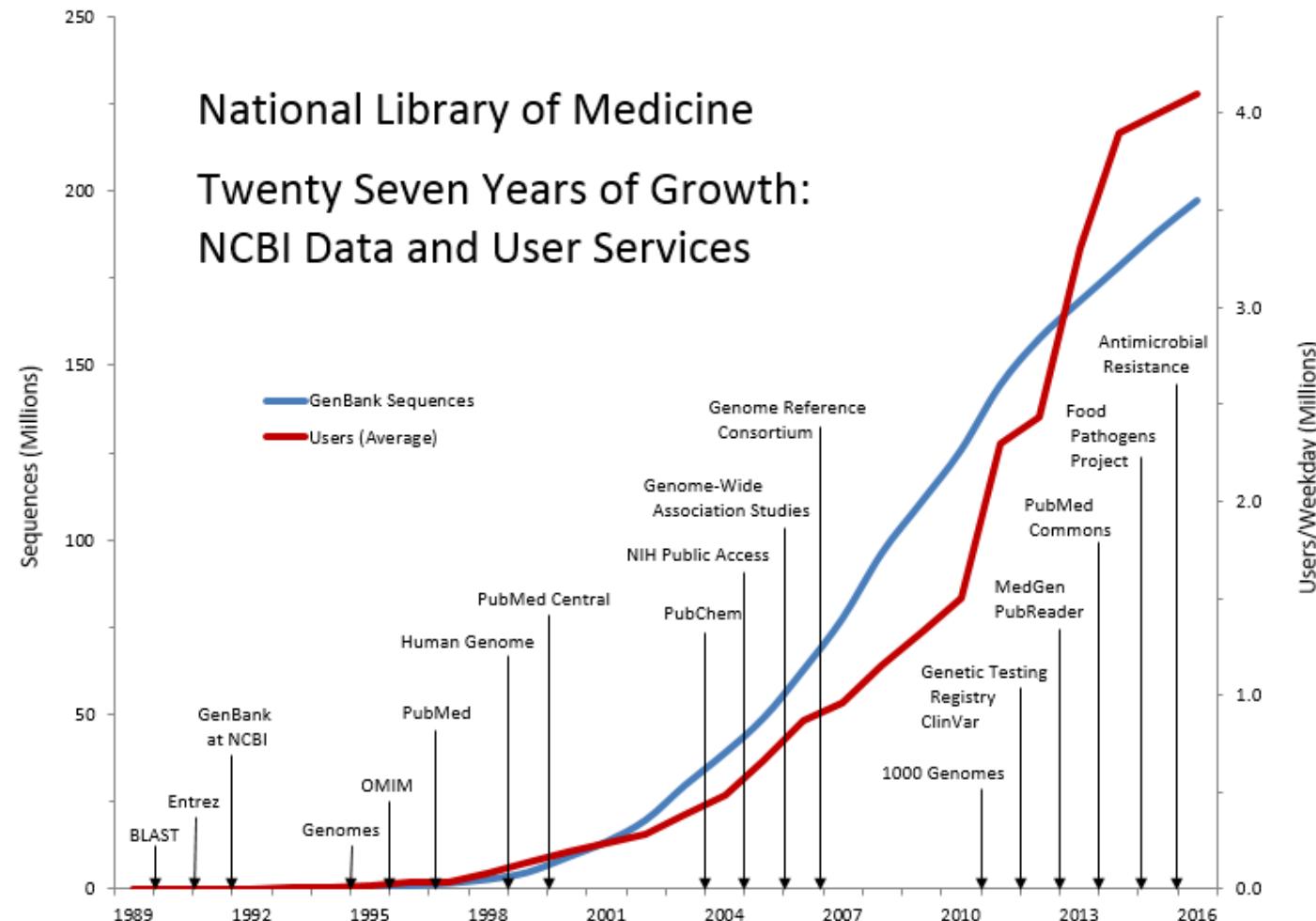


Photo credit
<https://www.nlm.nih.gov/about/2018CJ.html>



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



Database



Cold Spring Harbor Laboratory
DNA LEARNING CENTER

Let's BLAST a sequence

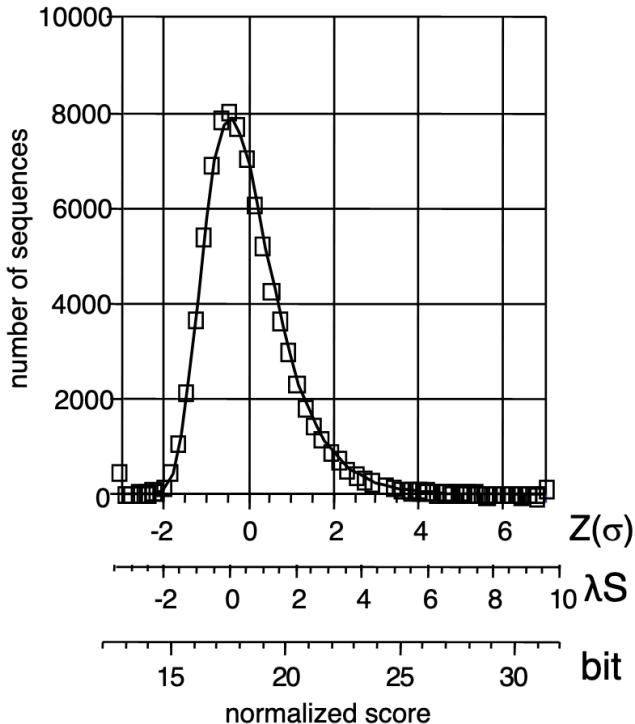
```
>mosquito-1F
CTTAAAGTATATTAAATTCTGCTGAATTAAAGTCACCCAGGGATTTAT
TGGAAATGATCAAATTATAACGTAATTGTTACAGCTCATGCATTATT
ATAATTTTTTTATAGTAATACCAATTATAATTGGAGGATTGGAAATT
GATTAGTTCCCTTAATATTAGGAGCTCCTGATATAGCATTCCCTCGAAT
AAATAATATAAGTTTGAAATATTACCTCCTTAACTCTACTACTTT
CTAGTTCAATAGTAGAAAATGGAGCAGGGACAGGATAACAGTTA
TCCTCCTCTTCATCAGGAACAGCACATGCTGGAGCTCTGTTGATT
AGCAATTTCCTCTTCATTAGCAGGGATTTCATCTATTAGGAGC
AGTAAATTATTACTACTGTTATTAAATACGATCATCTGGAATTACTT
TAGATCGATTACCTTATTGTTGATCTGTAGTAATTACTGCTATTAA
TTACTTTATCTCTCCTGTATTAGCTGGAGCTATTACTATTAACT
GATCGAAATTAAACTTCCTTGTACCCAATTGGAGGAGGAGA
```

<https://blast.ncbi.nlm.nih.gov/Blast.cgi>



BLAST and controls

Why smaller databases are better (more sensitive) – statistics



$$S' = \lambda S_{\text{raw}} - \ln K m n$$
$$S_{\text{bit}} = (\lambda S_{\text{raw}} - \ln K) / \ln(2)$$
$$P(S' > x) = 1 - \exp(-e^{-x})$$
$$P(S_{\text{bit}} > x) = 1 - \exp(-mn2^{-x})$$
$$E(S' > x | D) = P D$$

$$P(B \text{ bits}) = m n 2^{-B}$$
$$P(40 \text{ bits}) = 1.5 \times 10^{-7}$$
$$E(40 | D=4000) = 6 \times 10^{-4}$$
$$E(40 | D=80E6) = 12$$

Photo credit
https://fasta.bioch.virginia.edu/biol4230/lects/biol4230_4_blast2.pdf

BLAST algorithm analogy – searching by “word”

The *Query sequence*

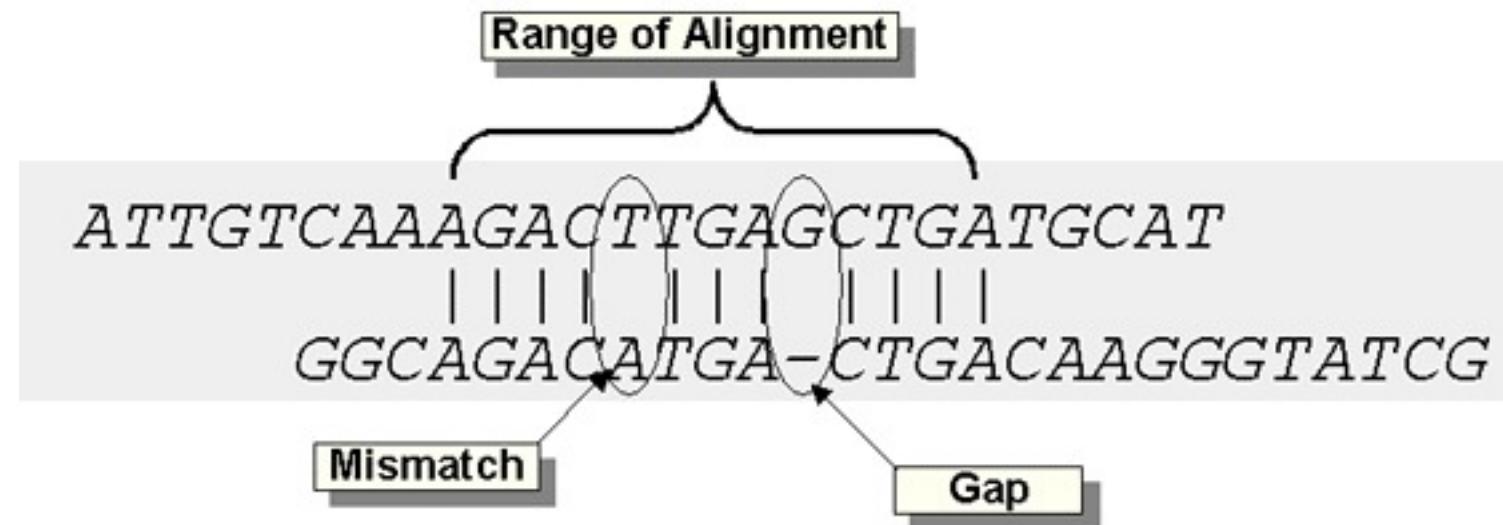
Is aligned to a *Subject* (a sequence in the database)

Q: **ACTGAC–ATCGGGGTGCTACG**

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

S: **ACTGACCATCGGAGTGCTACG**

BLAST algorithm analogy – alignment



$$S = \sum_{\text{identities, mismatches}} - \sum_{\text{gap penalties}}$$

Score = **Max(S)**

Photo credit

<https://www.ncbi.nlm.nih.gov/books/NBK62051/>



Cold Spring Harbor Laboratory
DNA LEARNING CENTER

Let's do a BLAST

Descriptions Graphic Summary Alignments Taxonomy

Sequences producing significant alignments Download ▾ Manage Columns ▾ Show 100 ▾ ?

select all 0 sequences selected GenBank Graphics Distance tree of results

	Description	Max Score	Total Score	Query Cover	E value	Per. Ident	Accession
<input type="checkbox"/>	Aedes vexans voucher BIOUG01574-F08 cytochrome oxidase subunit 1 (COI) gene, partial cds; mitochondrial	1053	1053	100%	0.0	99.83%	KR694809.1
<input type="checkbox"/>	Aedes vexans voucher BIOUG01519-A06 cytochrome oxidase subunit 1 (COI) gene, partial cds; mitochondrial	1053	1053	100%	0.0	99.83%	KT113440.1
<input type="checkbox"/>	Aedes vexans voucher BIOUG05112-D01 cytochrome oxidase subunit 1 (COI) gene, partial cds; mitochondrial	1053	1053	100%	0.0	99.83%	KM971547.1
<input type="checkbox"/>	Aedes sp. BOLD:AAA7067 voucher BIOUG08859-D04 cytochrome oxidase subunit 1 (COI) gene, partial cds; mitochondrial	1053	1053	100%	0.0	99.83%	KM910290.1
<input type="checkbox"/>	Culicinae sp. BOLD:AAA7067 voucher BIOUG03954-A01 cytochrome oxidase subunit 1 (COI) gene, partial cds; mitochondrial	1051	1051	99%	0.0	99.83%	KP039751.1
<input type="checkbox"/>	Aedes vexans voucher BIOUG24039-B11 cytochrome oxidase subunit 1 (COI) gene, partial cds; mitochondrial	1051	1051	99%	0.0	99.83%	KT707504.1
<input type="checkbox"/>	Aedes vexans voucher BIOUG27453-F12 cytochrome oxidase subunit 1 (COI) gene, partial cds; mitochondrial	1049	1049	100%	0.0	99.66%	MF820054.1



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.

Does BLAST tell me what species I have identified?

No*

(Some) Limitations to BLAST

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

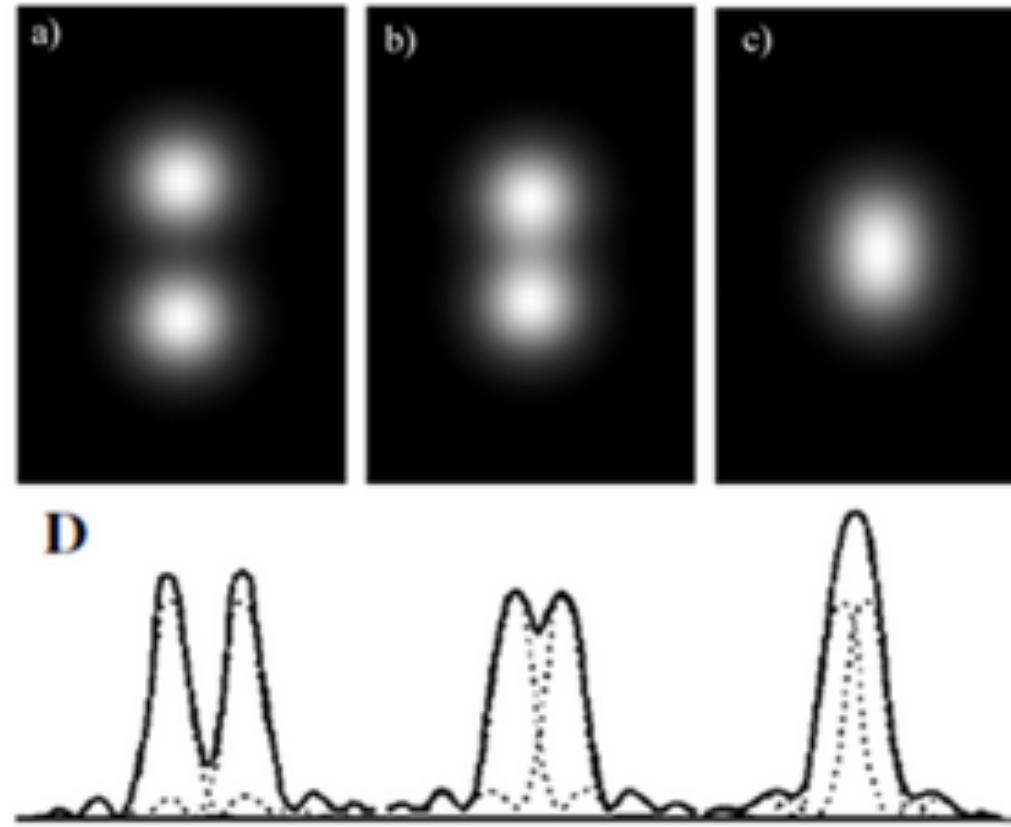
(Some) Limitations to BLAST

- **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)



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

Next time:

Multiple sequence alignments and
phylogenetics