Cold Spring Harbor Laboratory CSH DNA LEARNING CENTER

DNALCLIVE **Barcoding Bioinformatics** Part III

Jason Williams

Cold Spring Harbor Laboratory, DNA Learning Center

williams@cshl.edu



@JasonWilliamsNY



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



DNALC Website and Social Media



youtube.com/DNALearningCenter



facebook.com/cshldnalc



@dnalc



@dna_learning_center



Barcoding Bioinformatics Part III



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)





Barcoding Bioinformatics Part III

(Sequence alignment and phylogenetics)



Steps for today's session

- Recap on our experimental dataset
- Review of BLAST
- Introduction to multiple alignment
- Introduction to phylogenetic trees



Recap of the dataset





CSH Cold Spring Harbor Laboratory DNA LEARNING CENTER





Photograph by Michele M. Cutwa, University of Florida.

Anopheles larva



Culex larva



Photograph by Michelle Cutwa-Francis, University of Florida.

Cold Spring Harbor Laboratory DNA LEARNING CENTER

CSH S

Steps to DNA Barcoding



ACGAGTCGGTAGCTGCCCTCTGACTGCATCGAA TTGCTCCCCTACTACGTGCTATATGCGCTTACGAT CGTACGAAGATTTATAGAATGCTGCTACTGCTCC CTTATTCGATAACTAGCTCGATTATAGCTACGATG



Sequenced DNA is compared with DNA in a barcode database



Let's do a BLAST

| Descrip | ptions | Graphic Summary | Alignments | Taxonomy | | | | | | | |
|-----------|---|---------------------------|----------------------|-----------------------------------|----------------------------|--------------|----------------|----------------|------------|---------------|---------------|
| Seque | ences pro | oducing significant a | lignments | | Download ~ | Manage | e Colu | mns ~ | Sho | ow 10 | 00 🗸 🕜 |
| sele | lect all 0s | equences selected | | | | GenB | ank | Graphie | cs D | istance t | ree of result |
| | | | C | Description | | Max Score | Total Score | Query Cover | E value | Per. Ident | Accession |
| <u>Ae</u> | edes vexans | voucher BIOUG01574-F08 cy | tochrome oxidase sub | unit 1 (COI) gene, partial cds; n | nitochondrial | 1053 | 1053 | 100% | 0.0 | 99.83% | KR694809.1 |
| Ae | edes vexans | voucher BIOUG01519-A06 cy | tochrome oxidase sub | unit 1 (COI) gene, partial cds; n | nitochondrial | 1053 | 1053 | 100% | 0.0 | 99.83% | KT113440.1 |
| Ae | edes vexans | voucher BIOUG05112-D01 cy | tochrome oxidase sub | unit 1 (COI) gene, partial cds; n | nitochondrial | 1053 | 1053 | 100% | 0.0 | 99.83% | KM971547. |
| Ae | edes sp. BO | LD:AAA7067 voucher BIOUG0 | 8859-D04 cytochrome | oxidase subunit 1 (COI) gene, | partial cds; mitochondrial | 1053 | 1053 | 100% | 0.0 | 99.83% | KM910290. |
| | Culicinae sp. BOLD:AAA7067 voucher BIOUG03954-A01 cytochrome oxidase subunit 1 (COI) gene, partial cds; mitochondrial | | | | | | | | 0.0 | 99.83% | KP039751. |
| Ae | edes vexans | voucher BIOUG24039-B11 cy | tochrome oxidase sub | unit 1 (COI) gene, partial cds; n | nitochondrial | 1051 | 1051 | 99% | 0.0 | 99.83% | KT707504.1 |
| Ae | edes vexans | voucher BIOUG27453-F12 cy | tochrome oxidase sub | unit 1 (COI) gene, partial cds; n | nitochondrial | 1049 | 1049 | 100% | 0.0 | 99.66% | MF820054. |



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)



BLAST algorithm analogy

Query sequence ACTGACATCGGGGGTGCTACG



Database



Where in the DNA should we look for the "Barcode"?





Humans share greater than 99% of their DNA with other humans

Photo Credit: https://www.broadinstitute.org/files/styles/landing_page/public/genericpages/images/circle/MPG-mosaic_385x300.png?itok=B5zrVgYz





Photo credit: https://www.cnbc.com/2015/09/02/can-you-find-the-forgery.html





Photo credit: https://www.cnbc.com/2015/09/02/can-you-find-the-forgery.html



>Human Tubulin –Black?

>Human Tubulin –White?

>Human Tubulin –Asian?

Almost everywhere we look in our (human) genome, we all look the same



Cytochrome c oxidase I (COI)



Photo credit:

Structure: https://en.wikipedia.org/wiki/Cytochrome_c_oxidase_subunit_I#/m edia/File:PDB_locc_EBI.jpg Genome map: Emmanuel Douzery; https://en.wikipedia.org/wiki/Cytochrome_c_oxidase_subunit_I#/m

edia/File:Map_of_the_human_mitochondrial_genome.svg

CSH Cold Spring Harbor Laboratory DNA LEARNING CENTER



Choosing a barcoding locus



There are many criteria that go in to selecting an appropriate locus (location in the genome) that can serve as a barcode.

Three of them include:

- Universality
- Discrimination
- Robustness



Universality

Since barcoding protocols (typically) amplify a region of DNA by PCR, you need to choose a DNA sequence that every species has



Universality

Since barcoding protocols (typically) amplify a region of DNA by PCR, you need to choose a DNA sequence that every species has













Universality

Since barcoding protocols (typically) amplify a region of DNA by PCR, you need to choose a DNA sequence that every species has

3.0.3.3

| Carnivores | Omnivores | Herbivores |
|--|-------------------------------|----------------------------|
| Cat | • Pig | • Cow |
| 3.1.3.1 | 3.1.4.3 | 0.0.3.3 |
| 3.1.2.1 | 3.1.4.3 | 4.0.3.3 |
| • Dog | Human | Horse |
| 3.1.4.2 | 2.1.2.3 | 3.1.3/4.3 |
| 3.1.4.3 | 2123 | 3.1.3.3 |
| | 2.1.2.0 | Rabbit |
| | | 2.0.3.3 |
| | | 1.0.2.3 |
| Photo credit: | | Sheep |
| Dental formula https://slideplayer.com/slide/10972461/ | | 0.0.3.3 |
| Animal diagrams https://vet-science.blogspot.com/2012/01/dent | tition-in-sheep-and-goat.html | 3033 |

Cold Spring Harbor Laboratory

DNA LEARNING CENTER

CSH



Discrimination

Barcoding regions must be different for each species. Ideally you are looking for a single DNA locus which differs in each species



Discrimination

Barcoding regions must be different for each species. Ideally you are looking for a single DNA locus which differs in each species

HUMAN HAIR



Photo credit Human hair https://lewigs.com/human-hair-color-101/



Discrimination

Barcoding regions must be different for each species. Ideally you are looking for a single DNA locus which differs in each species



Photo credit http://loreal-dam-videos-corp-encdn.brainsonic.com/corpen/20160330pm/2016033 0-170533-c4bb4cb7/picture_photo_3eadc4.jpg



Discrimination (but not too much)

Fail: Sequence is completely conserved, good for PCR, but uninformative as barcode

| | • | | 1 | | | | | ••• | | | | | | | | | 1.1.1 | | | 1.1.1 | ••• | | | | |
|-----------|---|------|------|---------------------|---------------------|-----|----|-----|-----|------|----------------------|------|-----|-----|-----|-----|-------|--------------------|-----|-------|-----|-----|----|-----|---------------------|
| | • | | | 10 | | 2 | 20 | | | - 30 | 0 | | 4 | 10 | | | 50 | | | 60 | | | 70 | | |
| Species 1 | | atgt | ccgg | gc <mark>t</mark> a | ag <mark>e</mark> g | cga | ⊂g | tac | gat | cag | g <mark>et</mark> g | gtga | tct | ago | ga | ggt | acg | <mark>st</mark> ag | ttt | tgc | atg | cat | ga | cag | jaga <mark>t</mark> |
| Species 2 | 2 | atgt | ccgg | gcta | ag <mark>c</mark> g | cga | ⊂g | tac | gat | cag | g <mark>et</mark> g | g ga | tct | ago | gao | ggt | acg | gt ag | ttt | tgc | atg | cat | ga | cag | jagat |
| Species 3 | 3 | atgt | ccgg | gota | ageg | cga | cg | tac | gat | cag | g <mark>et</mark> g | g ga | tct | ago | ga | ggt | acg | gtag (| ttt | tge | atg | cat | ga | cag | jaga |
| Species 4 | Ł | atgt | ccgg | gota | ageg | cga | cg | tac | gat | cag | gete | g ga | tct | ago | ga | gg | acg | gtag | ttt | tge | atg | cat | ga | cag | jaga |
| Species 5 | 5 | atgt | ccgg | gota | ageg | cga | cg | tac | gat | cag | gete | gtga | tct | ago | ga | ggt | acg | gtag (| ttt | tge | atg | cat | ga | cag | jaga |
| Speices 6 | 5 | atgt | ccgg | gota | agcg | cga | cg | tac | ga | cac | g <mark>e t</mark> e | g ga | tct | ag | ga | gg | aco | <mark>rtag</mark> | ttt | tgca | a g | cat | ga | cac | jaga |



Discrimination (but not too much)

Fail: Sequence shows no conservation, impossible for PCR, but good as barcode

| | • | | | | | | | | |
|---------|---|---------|--|---------------------------|---|---|--------------------------|--|--------------------------------|
| | - | | 10 | 20 | 30 | 40 | 50 | 60 | 70 |
| Species | 1 | cttgaac | gaactcct | gcactgca <mark>c</mark> c | t <mark>c</mark> ta <u>a</u> ctgct | gggtta <mark>agct</mark> | cgttgctag | gaggtc <mark>atc</mark> a | atggggtgg <mark>cg</mark> aaac |
| Species | 2 | agactt | tt <u>a</u> ttggc | atgagagcac | a <mark>c</mark> gg <mark>t</mark> aggcga | a <mark>cttagt</mark> aa <mark>a</mark> g | ctcgcgagg | gttaag <u>g</u> a <mark>c</mark> o | ccaacggtt <mark>c</mark> ctgga |
| Species | 3 | caaagt | ta <mark>c</mark> taatc | tgtccgaa <mark>c</mark> c | g <mark>e</mark> gcacaatag | g <mark>t</mark> aaggg <mark>aagc</mark> | attgga <mark>a</mark> ag | tcaaga <mark>a</mark> aat | tgcgtagctagccta |
| Species | 4 | acttgc | aa <mark>c</mark> gc <mark>a</mark> aa | cgcttaaaat | catt <mark>t</mark> ggtag <mark>d</mark> | tcatgtacag | acgctaatc | tcagac <mark>a</mark> ggo | ccctaagg <mark>c</mark> cggat |
| Species | 5 | aaaacct | cc <mark>c</mark> tcaat | tcagcatc <mark>c</mark> t | cata <mark>t</mark> cagtc | t caaaagaag | tggactacg | gattat <mark>to</mark> a | agaaacgtg <mark>agggtt</mark> |
| Speices | 6 | gcatgtt | ct <mark>c</mark> aa <mark>a</mark> at | agactcgtcg | a <mark>c</mark> gc <mark>t</mark> ggctc | agacggacg | atgtccgtg | ataagg <mark>a</mark> a <mark>c</mark> o | caatattct <mark>cg</mark> cgcg |



Discrimination (but not too much)

Win: Sequence shows some (ideally ~70%) conservation, good for PCR, good as barcode

| | • 111 | | | | 111 | | | | | 111 | • • • | | 1 | 111 | | | | | 111 | ••• | | | ••• |
|-----------|-------------------|-----------------------|-----------------------|------|-----|------|--------------------|----------------------|--------------------|-----|--------------------|------|----------------------|------|---------------------|-----|--------------------|-----|-----|---------------------|--------------------|------|-----|
| | • | | 10 | | | 20 | | | 30 | | | 40 | | | 50 | | |) | 70 | | | | |
| Species 1 | . at | gge | <mark>jeget</mark> a | gege | gat | egt | c <mark>c</mark> g | a <mark>tc</mark> aa | actg | tga | t <mark>c</mark> t | ageg | <mark>ja</mark> ago | tac | g <mark>t</mark> ag | caa | a <mark>t</mark> g | cat | gct | t <mark>c</mark> ga | a <mark>t</mark> c | agag | gat |
| Species 2 | at at | g <mark>-cc</mark> gg | g <mark>eget</mark> a | gege | gat | cgt | a <mark>c</mark> g | a <mark>tc</mark> aa | actg | tga | tct | ageg | j <mark>ac</mark> go | tac | gtag | ctt | t <mark>t</mark> g | cat | gca | t <mark>e</mark> ga | atc | agag | gat |
| Species 3 | ; <mark>at</mark> | g <mark>tee</mark> gg | g <mark>eget</mark> a | gtgo | gat | cgt | c <mark>c</mark> g | a <mark>tc</mark> aa | actg | tga | tct | ageg | j <mark>ac</mark> go | taco | gtag | cta | a <mark>t</mark> g | cat | gca | t <mark>c</mark> ga | atc | agag | gat |
| Species 4 | a a t | g <mark>tee</mark> gg | g <mark>eget</mark> a | gtgo | gat | ⊂gt | a <mark>c</mark> g | a <mark>tc</mark> aa | actg | ca | tct | ageg | j <mark>ac</mark> go | too | gtag | ctt | t <mark>t</mark> g | cat | gca | tega | atc | agag | gat |
| Species 5 | 5 <u>e</u> t | g <mark>tee</mark> gg | g <mark>eget</mark> a | gtgo | gat | cg t | a <mark>c</mark> g | a <mark>tc</mark> aa | actg | ca | tct | ageg | j <mark>ac</mark> go | tac | gtag | ctt | t <mark>t</mark> g | cat | gca | tega | atc | agag | gat |
| Speices 6 | ; <mark>at</mark> | g <mark>te</mark> tge | gegeta | gege | gat | g | a <mark>c</mark> g | a <mark>tc</mark> a | g <mark>e g</mark> | ga | tct | ageo | j <mark>ac</mark> go | tac | <mark>j a</mark> g | cta | atg | cat | gca | t <mark>c</mark> ga | atc | agag | ja |



Robustness

Since barcoding protocols (typically) amplify a region of DNA by PCR, also need to select a locus that amplifies reliably and sequences well



Photo credit https://www.bio-rad.com/es-mx/applicationstechnologies/pcr-troubleshooting?ID=LUSO3HC4S



Choosing a barcoding locus

Cytochrome Oxidase C subunit 1 (COI):

- Universal
- Discriminating
- Robust



Photo credit: https://en.wikipedia.org/wiki/Cytochrome_c_oxidase#/media/File:Cmplx4.PNG


Reference data and phylogenetics





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

DNA LEARNING CENTER

- We have sensitivity controls (sequence quality, BLAST parameters)
- We have outgroup sequences (non-mosquito, negative controls) and known samples (positive controls)
 Cold Spring Harbor Laboratory

Intro to Multiple Sequence Alignment





To compare sequence features (nucleotides, amino acids, etc.) we need to line them up

Photo credit: https://www.pexels.com/photo/jigsaw-puzzle-1586950/



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



THEFISHISINTHEWATER TH'FESHISINTH'WATER DEVISISINHETWATER DIEVISISINDIEWATER DERFISCHISTIMWASSER FISKINERÍVATNI



ENGLISH SCOTTS DUTCH AFRIKAANS GERMAN ICELANDIC

. . . . | | | | 15 5 THEFISHISI NTHEWATER-TH-FESHISI NTH-WATER--D-EVISISI NHETWATER--DIEVISISI NDIEWATER--DERFISCHI STIMWASSER ----FIS-KI NERIVATNI-



ENGLISH SCOTTS DUTCH AFRIKAANS GERMAN ICELANDIC

| 5 | 15 ' |
|------------|-------------------|
| THEFISHISI | NTHEWATER- |
| TH-FESHISI | NTH-WATER- |
| -D-EVISISI | NHETWATER- |
| -DIEVISISI | NDIEWATER- |
| -DERFISCHI | STIMWASSER |
| FIS-KI | NERIVATNI- |

Colored at 60% identity



. . . . | | | | 15 5 THEFISHISI NTHEWATER-ENGLISH TH-FESHISI NTH-WATER-SCOTTS DUTCH -D-EVISISI NHETWATER--DIEVISISI NDIEWATER-AFRIKAANS -DERFISCHI STIMWASSER GERMAN ----FIS-KI NERIVATNI-ICELANDIC



ENGLISH SCOTTS DUTCH AFRIKAANS GERMAN ICELANDIC





ENGLISH SCOTTS DUTCH AFRIKAANS GERMAN ICELANDIC





ENGLISH SCOTTS DUTCH AFRIKAANS GERMAN ICELANDIC







Number of (global) alignments for 2 sequences of length n



Photo credit: https://www.pexels.com/photo/jigsaw-puzzle-1586950/





Number of (global) alignments for 2 sequences of length n



So, for 2 sequences (n=100) $\approx 10^{77}$ *

- We need software!

Photo credit: https://www.pexels.com/photo/jigsaw-puzzle-1586950/



*(some are trivial)

CLUSTAL alignment algorithm

ClustalW steps



Photo credit: https://slideplayer.com/slide/5155996/







Making phylogenetic trees



Tree relationships



CSH Cold Spring Harbor Laboratory DNA LEARNING CENTER

Tree relationships

This is not a phylogenetic tree





A phylogenetic tree is a <u>hypothesis</u> about how species may be related



Photo credit https://en.wikipedia.org/wiki/File:Phylogenetic_tree.svg



Trees can be created from (and therefore reflect) characteristics/traits



Photo credit https://kids.britannica.com/students/assembly/view/235364



Trees can be created from (and therefore reflect) characteristics/traits



Photo credit https://wildlifesnpits.wordpress.com/2014/03/22/understanding-phylogenies-terminology/





Photo credit https://wildlifesnpits.wordpress.com/2014/03/22/understanding-phylogenies-terminology/

&CSH &

DNA LEARNING CENTER

Tree Vocabulary

- Taxa: Individual species
- Tip: The endpoints of a tree
- Node: A point where branches split
- Branch: Any collection after a node



Photo credit https://wildlifesnpits.wordpress.com/2014/03/22/understanding-phylogenies-terminology/



Tree experiments

- Experiment 1: Unknown Mosquitos and outgroup
 - Method: neighbor joining
- Experiment 2: Unknown Mosquitos and outgroup
 - Method: maximum likelihood
- Experiment 3: Unknown Mosquitos + BLAST hits and outgroup
 - Method: neighbor joining
- Experiment 4: Unknown Mosquitos + BLAST hits and outgroup
 - Method: maximum likelihood



Trees are also computationally intensive

Number of possible rooted trees for n sequences

= (2n-3)! / (2ⁿ⁻² (n-2))!

| 2 sequences: | 1 |
|---------------|--|
| 3 sequences: | 3 |
| 4 sequences: | 15 |
| 5 sequences: | 105 |
| 6 sequences: | 954 |
| 7 sequences: | 10395 |
| 8 sequences: | 135135 |
| 9 sequences: | 2027025 |
| 10 sequences: | 34459425 |
| 51 sequences: | >10 [®] (nb of particles in the universe) |
| 51 sequences: | >10 [®] (nb of particles in the universe) |

Photo credit: https://www.slideshare.net/sebastiendelandtsheer/phylogenetics1



• There is no one "best" method



• There is no one "best" method

Neighbor joining:

"distance-based" method of tree building



• There is no one "best" method

Neighbor joining:

- "distance-based" method of tree building
- a matrix of the sequences is created based on distance between each pair of sequences in multiple alignment



• There is no one "best" method

Neighbor joining:

- "distance-based" method of tree building
- a matrix of the sequences is created based on distance between each pair of sequences in multiple alignment
- distance is related to the number of mismatches between the sequences



Bootstrap values (how we evaluate NJ trees):

 columns in the sequence alignment randomly resampled to make many new alignments (DNA subway does this 100x)



Bootstrap values (how we evaluate NJ trees):

- columns in the sequence alignment randomly resampled to make many new alignments (DNA subway does this 100x)
- a matrix of the sequences is created



Bootstrap values (how we evaluate NJ trees):

- columns in the sequence alignment randomly resampled to make many new alignments (DNA subway does this 100x)
- a matrix of the sequences is created
- each bootstrap value is the number of times that particular relationship appears in the 100 resampled trees



Bootstrap values (how we evaluate NJ trees):

- columns in the sequence alignment randomly resampled to make many new alignments (DNA subway does this 100x)
- a matrix of the sequences is created
- each bootstrap value is the number of times that particular relationship appears in the 100 resampled trees
- Values of 70 and above are "plausible"; above 95 considered highly supported



Maximum likelihood:

 Attempts to take into account observed patterns of how nucleotides and amino acids change over time. For instance, mutations from C to T are more common than mutations of C to A.



Maximum likelihood:

- Attempts to take into account observed patterns of how nucleotides and amino acids change over time. For instance, mutations from C to T are more common than mutations of C to A.
- The tree with highest overall likelihood score is accepted as the best estimate of the relationships between sequences.


Tree building methods

Maximum likelihood:

- Attempts to take into account observed patterns of how nucleotides and amino acids change over time. For instance, mutations from C to T are more common than mutations of C to A.
- The tree with highest overall likelihood score is accepted as the best estimate of the relationships between sequences.
- The length of the branches between nodes is also a measure of the confidence of the relationships. Very short branches separating sequences have lower confidence and the relationships are less certain, while longer branches are better supported.



Summary



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)



See the handouts for even more explanations and activities



DNALC Website and Social Media

dnalc.cshl.edu



dnalc.cshl.edu/dnalc-live

