

James Fogelman  
Silvana Gaudieri  
Phillip Harris  
Eleanor Steinberg  
Anne Yoder

mail

- mailroom
- ev0/5
- passwords are mole94

computers

- talk to Brenden
- don't use Swafford
- Vax α ... ev0 3
- ev0 5 ... vax ... normal
- ev0 1 ... vax ... normal
- ev0/4
- two Indigos

two disks ... vms?  
unix.?

## MACS

- WORKSHOP1-15

- SERVER

- HAS MISC. UTILITIES

- MOUNTED BY RESTARTED

- TRANSFER FILES

- CHOOSER -- APPLESARE

- WORKSHOP COMPUTERS - name = Workshop-#

- PASSWORD = WORK#

## EVOL1

- doubleclick

- go to open to get to other computers

- dir

## MAIL

- mail enter

- send ... enter

- to ? --- "name"

- subject --

quit when done  
exit

INSIDE DOMAIN

OUTSIDE

- send enter

- to type 2 on keypad get SMTP "address"

## SUBDIRECTORIES



# LOGOFF

## FTP

### TELNET PROGRAM

- ① make sure FTP is enabled
- ② set transfer directory
  - unselect drive

③ type FTP ... go to network ... send IP #.

④ put "file name"

⑤

## PUTTING TO VAX

- ① set directory
- ② enable

} DON'T ~~DO~~ VIEW SEQ FILES

## FETCH

host

user ID

passwd

automatically does binaries

MANUALS - BLACK BOOKS

mput  
mget

WORKSHOP

SEQUENCE

- load master (for master seg)
- C
- move cursor to right

- font

BKG Cyclic --- color

ctrl Z axis

W80  
W132

D32

width

sign

round at ...

... ..



ERLDIR IN MANUAL  
 PFI 4 ... GOLD BOTTOM

SEQUENCE NOTES

STARTING

- ① TYPE SEQUENCE
- ② ENTER # SEQs, LENGTH

COMMANDS

- CONSENSUS      LOAD filename ... opens file
- BASE-COMPOSITION      WRITE filename ... save file
- NWS              PRINT filename ... saves file in print format
- INVERT              # \_      1      .. executes vms command
- PLOT-PLLOT              @ \_      .. executes SEQ commands from file
- CREATE              .. creates new seq. line
- DIRECTORY              .. lists directory
- EXIT              ... leaves + save
- QUIT              ... leaves w/o save
- CHANGE              .. enter edit mode

GOLD PFI GOLD	Help PF2 Help	Find PF3 Find Next	undel. char PF4 del. Char
command		replace	
7	8	9 append	—
top	bottom	paste	undel. char
4 advance	5 backup	6 cut	9 del. char
change case			substitute
1 move 10 char	2 end of line	3 mv 1 char	ENTER
open New Line, Resol			
Beginning of line select			ends

MOVING

- ARROWS ... 1 character
- WORD ... 10 characters
- PAGE ... 50-100
- BEG. LINE
- END. LINE
- GOLD # ARROW = repeat # arrow
- GOLD FIND GAATC RETURN = move to GAATC
- ADVANCE ... set direction R
- BACKUP ... " " L
- ctrl F ... move to next residue
- GOLD ... 2nd FUNCTION

GOLD ~~command~~  
 puts it in  
 other things

GOLD command  
 pack = compact

FOUNTS

BKG Cuck ... color



data: weighting mask

GOLD COMMAND WRITE -- saves

GOLD COMMAND QUIT ... QUIT ~~NO~~ SAVE

GOLD COMMAND EXIT -- SAVES

GOLD COMMAND SIMILARITY/WEIGHT = SEQ# [#..#]

/OUTPUT = NAME /tree

CTRLZ TO COMMAND

WRITE/PHYLIPI/SEQUENCES = [#..#]

- remove F

WRITE/PAUP/SEQ = [#..#]

WRITE/PAUP/WEI=2/SEQ = [2..7]

# Phylip Notes

## PROGRAMS

PROTPARS

DNAPARS

DNAMOVE ... reconst. ancestral

DNAPENNY ... branch & bound

DNACOMP ... compatibility

DNAINVAR ...

DNAML ... max. lik

DNAMLK ... assumes mol. clock

DNA DIST

PROTDIST ... makes dist. matrix

RESTML ... restriction sites

SEQBOOT ... bootstrapping

COALIKE

FITCH ... takes dist. matrix & makes tree; no clock

KITSCH ... " " " ; clock

NEIGHBOR ... UPGMA

CONTML ... gene freqs → tree

GENDIST ... gene freqs → dist.

CONTRAST ...

MIX ...

MOVE ...

PENNY ...

DOLLUP ...

DOLMOVE ...

DOLPENNY ...

CLIQUE ...

FACTOR ...

DRAWGRAM ... plots rooted trees

DRAWTREE ... plots unrooted trees

CONSENSE ... consensus tree

RETRACE ... reads in tree



Highly likely

likely

likely

likely

likely

likely

likely

likely

likely

likely

likely

likely

likely

likely

likely

likely

likely

likely

likely

likely

likely

likely

likely

likely

likely

likely

likely

likely

likely

likely

likely

likely

likely

W. Fitch

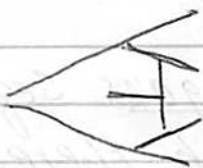
Definitions

- ① Tree ... in math. is a type of graph  
 ... graphs have points & edges



... a tree is a connected graph ... minimally

- ② Network ... e.g. recombination trees



← not a minimal graph

- ③ Rooting - most recent ancestor of all in group

... parsimony ... uniformly weighted trees  
 can be rooted anywhere

... many ways to root tree

① distance from tip

② outgroup

- must be useful outgroup

- too distant  $\cong$  random

- random roots the longest branch

- so want one that is very close

- but not w/in group

## Bifurcating trees

### Trifurcating trees

- usually used to indicate ambiguity about deep branches

### Data

① homologous alignment

② similarity  $\neq$  homology

③ gene tree  $\neq$  species tree

④ orthologous sequences are those for which there is an <sup>exact</sup> correspondence betw. gene & organism tree.

⑤ paralogous -- not same as organism tree  
- e.g. duplications (e.g. lysozyme)



## Methods Distances

- ① hybridization
- ② immunology
- ③ sequence

- but problems:

① back/d. ary mutations

Jukes-Cantor = ~~for DNA~~

DNA

$$d = -\frac{3}{4} \ln(1 - \frac{4}{3}\lambda) \quad \lambda = \% \text{ sites that differ}$$

- only works if all sites variable

- all sites can accept all changes

AA

$$d = -\frac{19}{20} \ln(1 - \frac{20}{19}\lambda)$$

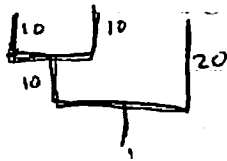
- but assumes all 20 aa are possible

~~UPGMA~~

a) UPGMA

unweighted pair group method using averages

	A	B	C
A		20	38
B			42
C			



Makes all branches from a node equal (averages distances)

b) Fitch-Margoliash

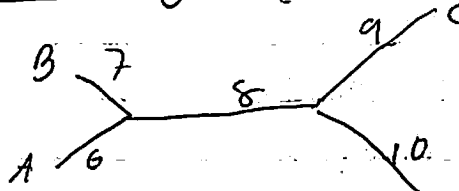
- similar but allows variable length from node

c) Farris-Wagner distance

- length =  $l_{min} = \sum \text{branch lengths}$

- so minimizes branch lengths as long as none are smaller than distance

d) Neighbor-joining (Fitch  $\frac{1}{2}$  ??) (NPS)



true #s = distance

	A	B	C	D
A	///	13	23	24
B		///	24	25
C			///	19
D				///

$$d_{AB} + d_{CD} < d_{AC} + d_{BD} = d_{AD} + d_{BC}$$

$$(13 + 19) < 23 + 25 = 24 + 24$$

SHOULD ALWAYS HOLD

METRICS

$$d(AA) = 0$$

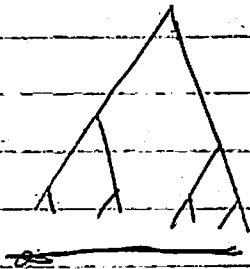
$$d(AB) = d(BA)$$

$$d(AB) + d(BC) \geq d(AC) \quad \text{triangle inequality}$$

why doesn't it always hold  
- homoplasy = multiple events

## 2) Sequence methods

a) Maximum likelihood Felsenstein et al

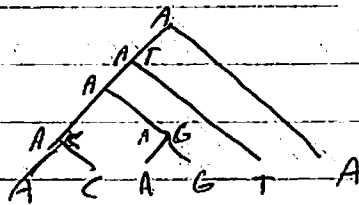


transformation matrix = prob. of nucleotide changes

- pick tree that gives maximum likelihood of getting that data

b) Maximum parsimony Hennig et al

- designed for morphology



- what are minimal # changes

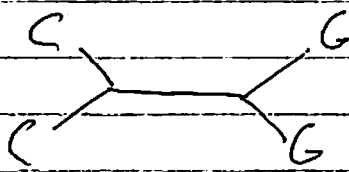
- pick intersections

- if no consensus put both

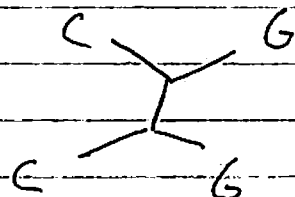
- do again

- can get an estimate of ancestral sequence

c) Evolutionary parsimony U Lake



parsimony would link this way



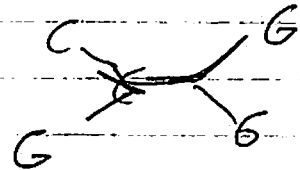
what if this is true?



• @ only uses transversions

- ASSUMES EACH TRANSVERSION  
IS EQUALLY LIKELY

- SUBTRACT OUT LIKELIHOOD  
OF GETTING



WHEN IT IS WRONG

- EXPECTED VALUE FOR ~~ADDER~~ ~~WR~~  
WRONG TREES IS ZERO

-  $\sigma^2$  CALLED INVARIANCE

- But need lots of data

## CONFIDENCE

### 1) Bootstrapping or jack-knifing

jack-knifing - sample w/o replacement

bootstrapping - sample w/ replacement (one side can be done twice)  
- 37% ( $e^{-1}$ ) are NOT sampled  
- allows scoring of ~~branches/nodes~~ clades  
reliability

CLADE = monophyletic group

- appears to be conservative

- prob of being correct  $>$  bootstrap value  
if bootstrap  $>$  75%

### 2) Likelihood

- can find value for diff. trees  
but NOT groups w/in tree

## FAILURES IN TREES

① alignment error - most are progressive ... pw comps  
- then grouping

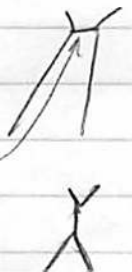
- progressivity is BAD for trees

- LAKE showed order biases tree

② homoplasy - multiple events in one position

③ unequal rates

- kills maximum parsimony if tree is  
because two long branches  
become alike  
especially if deep branch short



## Tree Failures

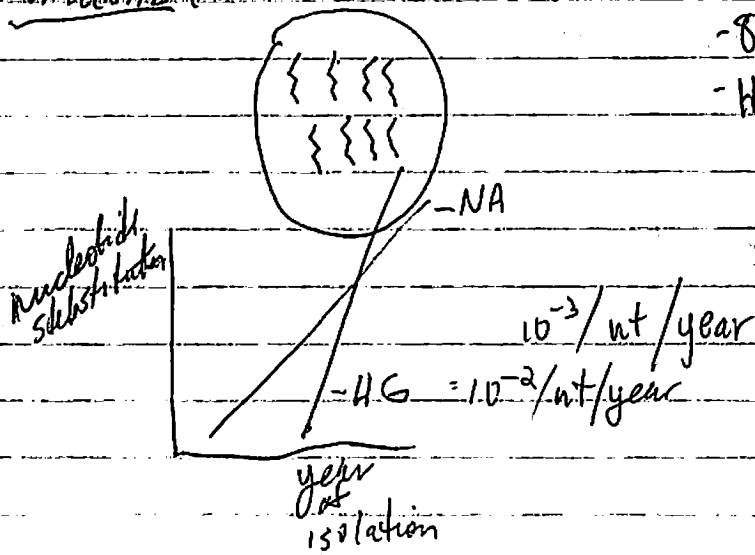
### 4) Composition

- e.g. GC content - changes v. rapidly



# Fitch - Applications

## Influenza



- 8 RNAs
- HA, NA, NS

NUCLEOTIDE

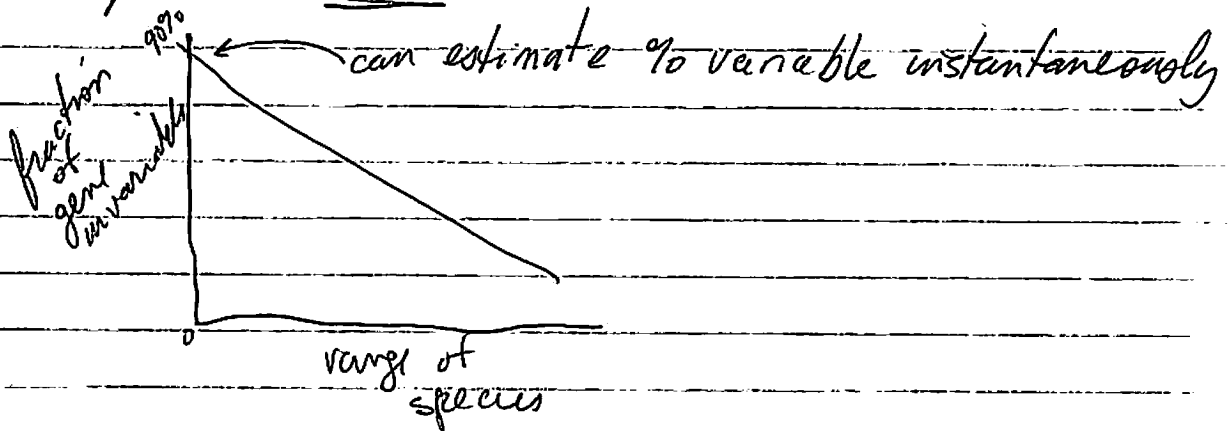
ANTIGENIC SITE  
vs  
NON-ANTIGENIC SITE

nucleotide rate high in all  
aa rate high only in mammals

## Cytochrome C Rates

- can estimate # positions that should have NO change using Poisson

- compare to real



- Fixation Rate Per Covariation (sites free to vary)  
is relatively constant

- SOD

- ISs in E. coli

Mitch Sogin -

- ① Plants vs. Animals  
but - Euglena
- ② Protists vs Plants vs Animals
- ③ 5 Kingdoms

PROTISTS - V Diverse

AMOEBOIDS - FLAGELLATES - CILIATES -

Amoebas  
- not monophyletic

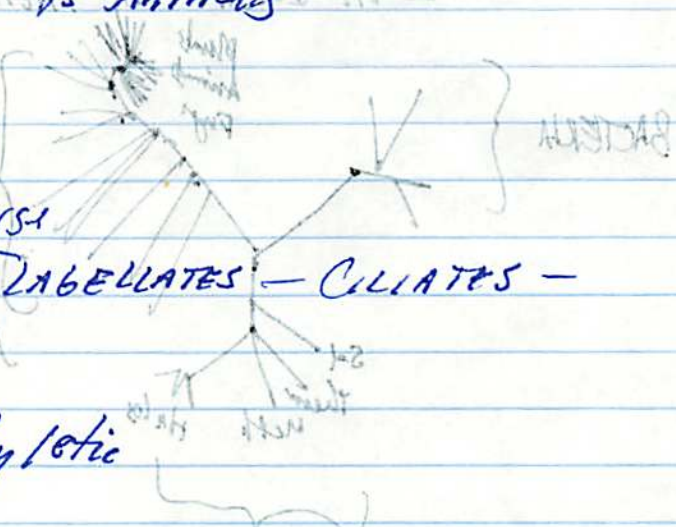
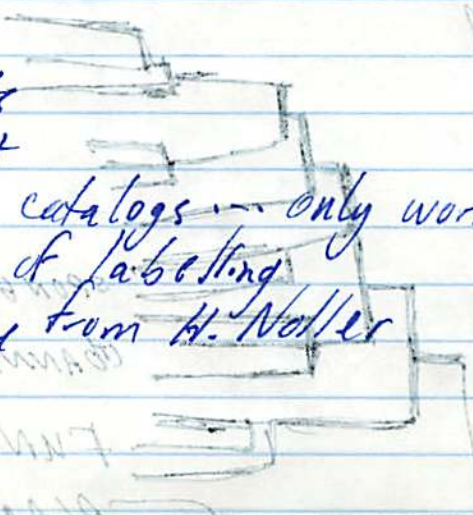
Chlorophytes ...  
Dinoflagellates ... diff. chlorophyll than chlorophytes

Flagellates  
- Euglena  
- Trypanosomes

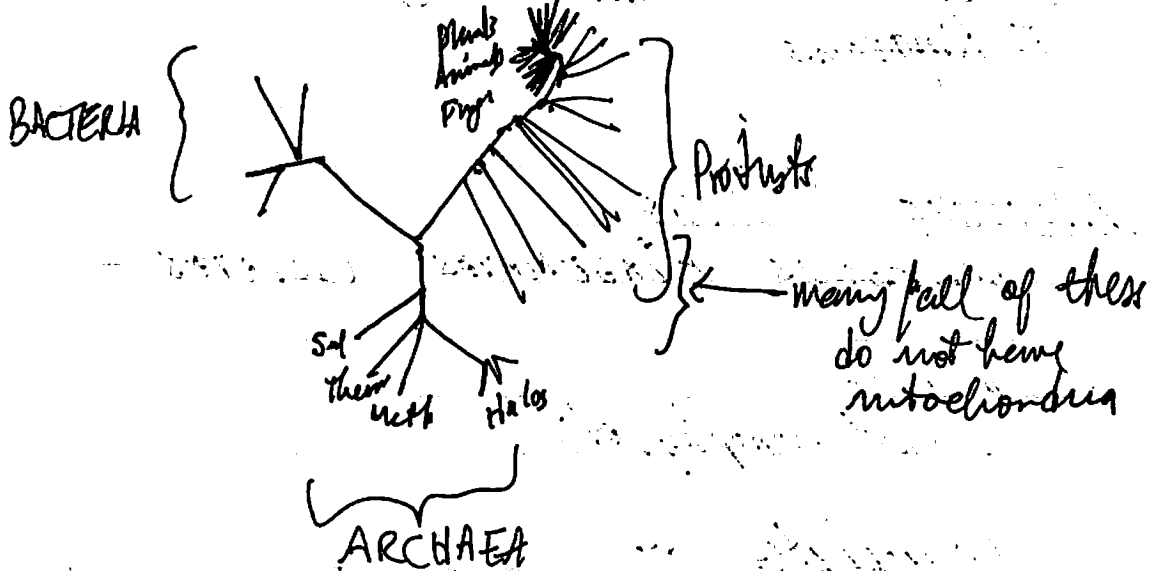
Ciliated

TREES

- rDNA
- Carl Woese ...
- 20 Fingerprints
- 55 ... too short
- Oligonucleotide catalogs ... only worked for bacteria because of labelling
- 1st sequence from H. Noller



- © Clone
- © Alignment ... does by eye
  - ... cannot align hypervariable regions
  - ... so he doesn't use them

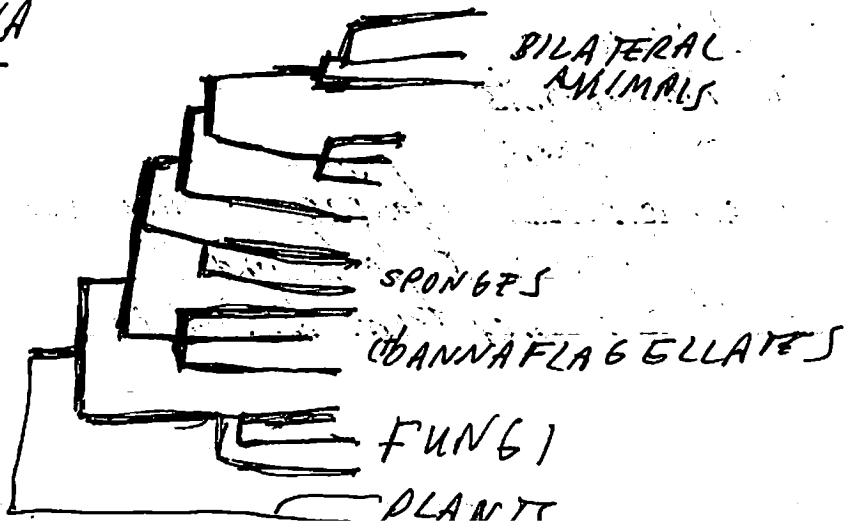


CROWN OF TREE  
Animals/Plants/Fungi

SOGIN 2 Field et al - Molecular Phylogeny of the Animal Kingdom

- © Animals in two clades
- © But Mitch says rt sequencing was wrong  
 & interior nodes are ov. vv. limited # of positions

NEW rRNA



- Hasegawa also? did this w/ EF.
- Baldauf & Palmer similar

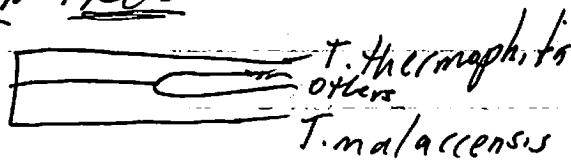
Phylogenetic w/in closely related groups

- use variable regions
- Fernandez et al '1

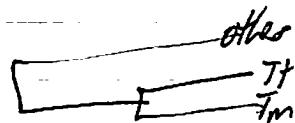
### Tetrahymena

- 13 species rDNA sequences
- rDNA RFLP's
- 20 cytoskeletal proteins (even more resolution)
- micronucleus has  $\approx 1$  rRNA gene
- similar branching w/ mt DNA

### INTRON TREE



### PDNA



### Conclusion

- lateral transfer of introns

### Introns

- type I introns commonly in most conserved domains if anywhere

- variable regions pick up long new stretches

### Other

- length not correlated w/ phylogeny
- organization not correlated either



Plasmodium berghei.

- two diff. tx. units
- A gene is colinear w/ rRNA
- C gene is not colinear w/ rRNA
- suggests C gene is pseudogene
- but not random variation
  - variation maps to variable regions
  - so is it a pseudogene?
  - discovered expression in diff hosts!

### ATTINE ANTS

- leaf harvester's fungal gardens
- 200 diff. ant species

QUESTION

J. WETTERER

- ANT & FUNGI PHYLOGENY ARE PARALLEL

### DEPTH OF BRANCHING

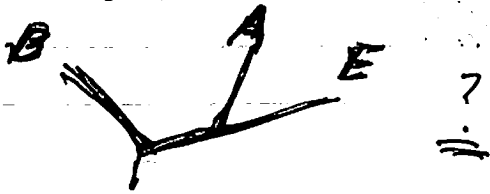
One explanation

- eukaryotes v.v. old
- but why not much biochem. diversity?
  - maybe because it hasn't been looked for
- but why not Euk fossils for v.v. long?
  - suggests definition of Euk. fossil may be bad (> 1u)
  - suggests may not form fossils



# Sogin contd

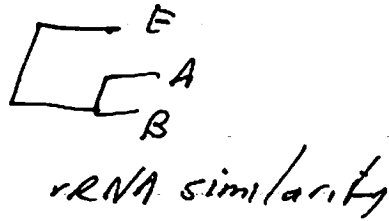
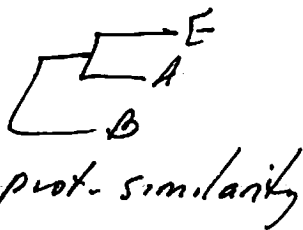
Where is the root?



Why this tree?

- ⊙ proteins (RNA pol's) more similar
- ⊙ duplicated genes as outgroup - Gogarten, Iwabe

Mitch -- doesn't like this because rRNA sequences in eukaryotes would have had to incr. at beginning of euk. evolution & then slow



How explain?

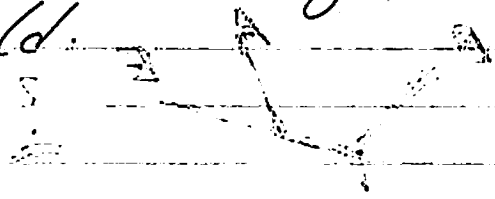
- ⊙ Woese says progenote
  - micells
  - small genes
  - RNA based system

The progenote isolated itself from micelle world

~~sphere~~ ~~cytoskeleton~~ eukaryotes → suggests cytoskeleton allowed for endosymbiosis and this then took in Archaeal genomes for its nucleus.

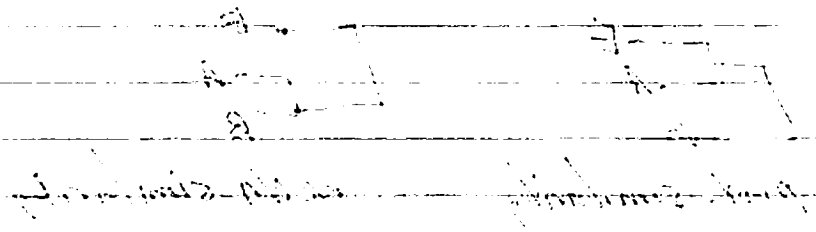
cell wall → Archaea  
Bacteria

Thus shouldn't find cytoskeleton homologs  
in the bacterial world.



... things ... from ...

... the ...



... ..

... ..

... ..

... ..

... = everything below

## VMS

DCL = language to talk to VMS

Getting out of commands

LEAST SEVERE - ctrl-Z sends exit signal  
↓  
MOST - ctrl-C sends "cancel" "interrupt"  
- ctrl-Y sends interrupt

TYPE ... view files

DIR ... list

MAKE

CASE ... insensitive

CREATE

OPTIONS ... use /

TYPE

/page 1 line at a time

RECALL ... shows last command

SHOW

-DEF ... shows current location [DIR.DIR.FILE]

- ... '...'

SET

DEF ... moves to new location [ ]

CREATE/DIR [name] ... makes dir. below where you are

COPY ...

RENAME ... [ ]



edit  
- rchange -- shows whole screen  
- \* exit -- saves

#c.e.

.com

- .command files
- all lines begin w/ \$/
- "\*" marks comment line

ove - extended vax editor

- ctrl Z -- immediately out

- use DO key -- commands, either says DO or TR keypad

- HELP

!out = creates file

# Dave's Swafford & Maddison

## D. Maddison

① Known phylogeny uses

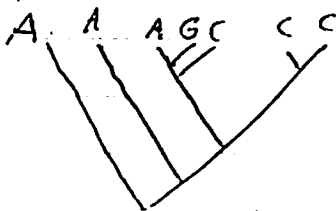
② infer ancestral state

③ types of changes/process

- Markovian model of state change
- ⇒ prob. of change betw. nodes or betw. generations
- depend only on state

A C T G  
A  
C  
T  
G

② example - Most PARSIMONIOUS RECONSTRUCTION



MPR = most parsimonious reconstruction of ancestral state given a tree & data

- PROBLEM - parsimony underestimates #'s of changes & homogenizes regions

- TESTS - HOLMQUIST & TUTENO

HILLIS ET AL W/ VIRUS

97% ACCURACY W/ PARSIMONY

③ W. Maddison method

- imagines a fully symmetrical tree w/ 64 taxa, binary character
- markovian matrix

	0	1
0	0.9	0.1
1	0.1	0.9

- prob. of accurate reconstruction is worse in deeper branches



How does one find MPR?

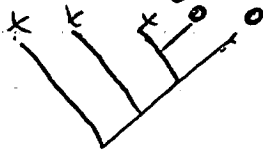
- Farris & Fitch
- Sankoff et al

⊛ Objective function - some value that is optimized  
- eg # of steps

Variation in MPR's

- # of changes along branch
- evolution of a character

How choose among MPR's



~~Deltran~~ ① Deltran vs. Acetran

Extremes  
much of the  
time but  
not always

Deltran = delayed transformation  
= changes farthest from root

Acetran = opposite

② Cost Matrices

- change weight of diff. changes
- influence MPR's

TRELENGTH - count # of steps for MPR's

INFERRING MATRICES

⊛ use MLE to compare matrices but hard to <sup>choose</sup> matrix

BUT PARSIMONY ② reconstruct ancestral states & make  
MAY HOMOGENEOUS matrix of freq of Δ's

# SWAFFORD

## CLASSIFICATION OF METHODS

### ALGORITHM

Average-linkage  
4-PT. Metric  
NJ

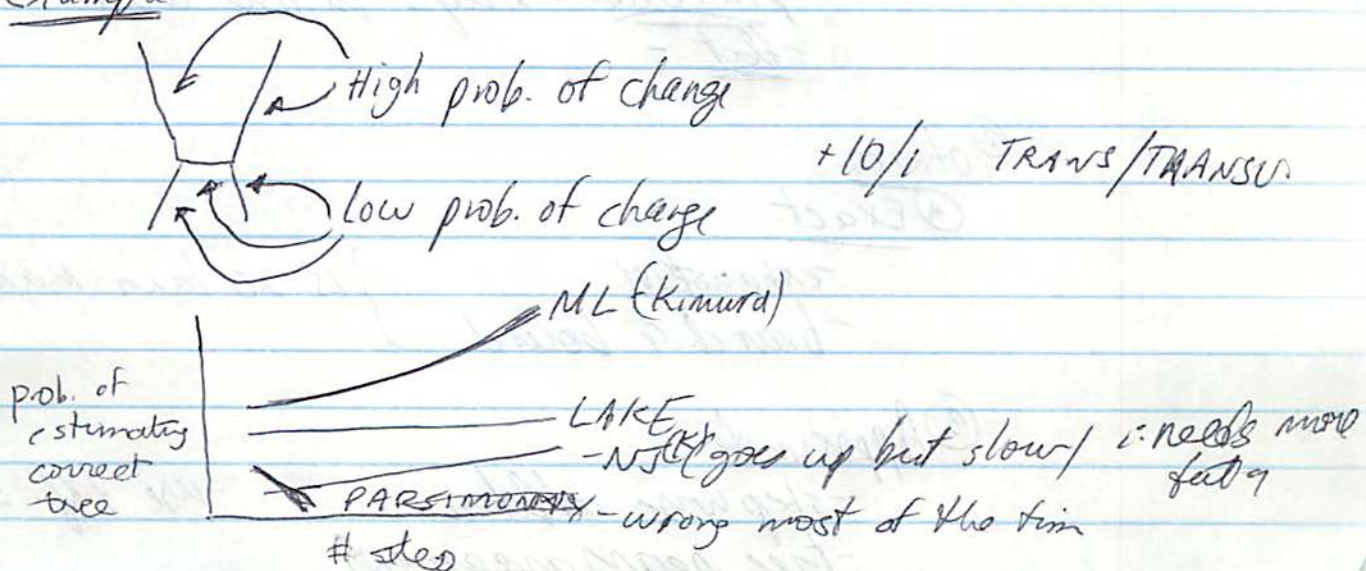
### CRITERION

Max-parsimony  
Additive tree distance  
MLE  
Spectral analysis

## QUALIFICATION OF METHODS

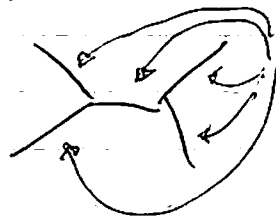
- Consistency - dependent on model
- Robustness - how is it affected by violation of assumptions
- Efficiency - how much data needed to get X% accuracy
- Predictability -
- Practicality -

### Example

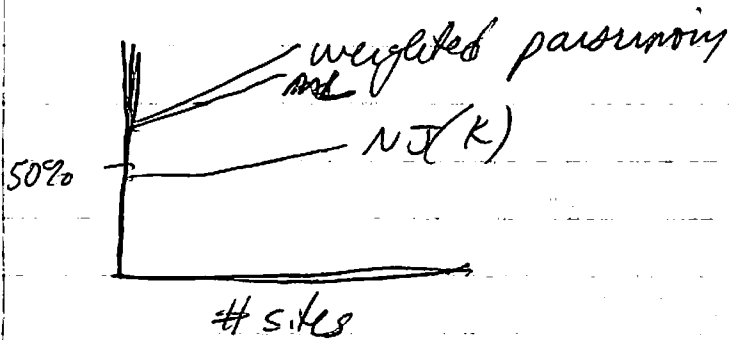


So methods are consistent but needs lots of data

### Example 2



All very high rates of change  
+ transition bias



### CONCLUSION

- Diff conditions --- diff methods best

### SEARCHING FOR TREES

- ① Hillclimbing - go up when can  
but - bad w/ multiple optima  
- plateaus - stays in one area  
- flat -

### ② Others

#### ① Exact

- exhaustive
  - branch & bound
- } 15-20 taxa max

#### ② Approximate

- step wise addition
  - tree rearrangements
- use diff. starts

## Exhaustive



} - get nice distribution patterns  
- some information

## Branch & Bound

- cuts off those searches when get too high a tree length

## Stepwise Addition

- only take one path at each step
- but can get stuck in local optima

## Nearest-Neighbor exchange

- too limited

## Tree Bisection Reconnection

- still get stuck on islands
- solution

- use diff. starting points

Example - MtDNA

### Minimum evolution method

- ① construct NJ tree
- ② get other trees
- ③ evaluate trees w/ ME score

Swafford/Maddison

MacClade → sets up files for PAUP  
Tree

Characters

- Quanta  
- PowerPC

- good alignments
- block selection

NEXUS
TEXT
PHYLIP
NBRF

Character status

- coding position
- translate

Type Edit --- matrix

PAUP

3.1.1.000

- SINNAUER IS NEXT DISTRIBUTOR
- ~3 MONTHS
- C COMPILER VERSION

Editor

- #NEXUS
- #INFO
- BEGIN DATA

AMBIGUITY CODES = EQUATE MACROS  
 $R \Rightarrow (A, B)$

Execute

can do from edit file or directly



## Search

- Branch & bound --
- Then ~~slow~~ examine

## ROOTING

- doesn't affect tree length
- default chooses 1st taxa
- select defines outgroup

- can do midpt. rooting

SHOW TREE

EXAMINE TREE

DESCRIBE TREE

cladogram  
phylogram

tree length

consistency index

PRINTING - multiple trees/page

## WHY? CHANGING CHARACTER TYPES (MATICES)

- using step matrices tree slows down

## GAPS

- using?

o program assumes each char. is independent

o can code gap as a separate character

## UNINFORMATIVE CHARS

o use for bootstrapping

INCLUDE/EXCLUDE CHARS

SET CHAR WEIGHTS

PAUPER

PAUPET

DL

PAUPER

## Paup & MacClade II



- Save tree in PAUP

- MacClade 3.04 \$75  
SINNAUER

publish.sinauer.com

- Open tree in MacClade

ROOTING ...  tool

CUTTING ...  (option + ) = cut all below


 = remove branch

 = rotate

 = search for trees



## FOOTNOTES - ANY CELL IN MATRIX

paintbrush  = force ancestors to have certain chars.

TRACE ALL CHANGES ... COLORIZES BY # OF  
CHANGES ON BRANCH

LIST TOOL -- SHOWS CHANGES ON a BRANCH

## CHART

- can show just transversions

Distances

PAUP PAUPOL  
PAUP PAUPOL

PAUP

Changes

D = distance

Compare trees

Save tree for PAUP

L = likelihood  
M = max likelihood

- Create new tree file
- Save as many trees as you want

Loading constraint trees

Load constraints

- Ⓞ backbone - only some of the taxa
- Ⓞ monophyly constraints

Compute consensus

- strict - shows only those that are shared

search

- include constraints (are compat)
- or converse constraints (are NOT compat)

Decay Index - # of steps it takes to decay a group

~~Exhaustive searches~~

Exhaustive searches

- shape of distribution (61 statistic)

must use  
random  
addition

## Bootstrap

- # replicates
- weights

1 2 3 4 5 6 7 8

→ ~~\*\*\*~~ 10 means taxa 1, 2, 3 are in one set  
& 4-8 on another in 10 bootstraps

## Random Addition Sequence

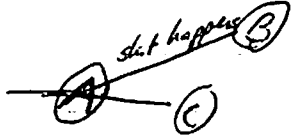
- choose stepwise addition panel
- select random
- close search --- 15/level info



GARY OLSEN

GARY'S WORLD VIEW

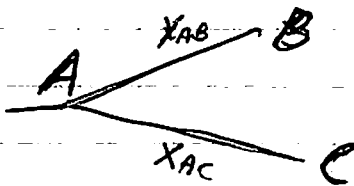
① there was a tree / history



- A is most recent common ancestor
- B & C diverging; possibly one more than other

EVOLUTIONARY DISTANCE

- number of fixed genetic events separating two genes



$$X_{BC} = X_{AB} + X_{AC}$$

MOLECULAR PHYLOGENY

- history of a portion of the genome
- not always the organism tree
- what region to choose to get organism tree?
- should do with multiple genes

WHY?

- ① organism phylogeny / evolution
- ② gene
- ③ process of sequence change
- ④ similarity searches
- ⑤ function

## Homology I

- common ancestry of gene - similarity in excess to that expected by chance

## Homology II

- common ancestry of residues  
(what about when there is a deletion then reappearance)

## Alignment

- expressing hypothesis about homology

### ① Terminal length variation

④ AUCGG  
AUCGG

↓  
- don't know what to do with this residue  
- so leave it out

### ② Internal length variation

there no sacrifices  
to the hypen

- indels - insertion in one / deletion in other  
(hypothesis of a genetic event)

- multiple possibilities for some ~~residues~~  
indels

- where does a gap go?

- if you don't know ... don't use that region

### ③ Aligning to 2ary structure

- e.g. tRNA

- may have no identical 1<sup>o</sup> structure

## why use 2ary structure

- ① most likely correspond to phylog. positions
- ② assumes unlikely to have non-pt. mutations
- ③ helps alignment

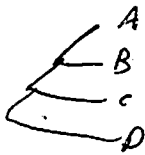
## METHODS

- ① BY HAND
- ② PAIRWISE
- ③ MULTIPLE

STAR -- ALIGN ALL TO REFERENCE  
BLAST -- CONSERVED REGIONS  
GLOBAL -- BEGINNING TO END

All to alignment  
as a tree

PROGRESSIVE (PILEUP)  
CLUSTAL  
TREE ALIGN



- ① generate tree from PWSimilarity
- ② align most similar
- ③ parsimoniously create consensus
- ④ align of next

- but the alignment will reproduce the history generated in step 1

- treealign can search alternative trees

## CAVEATS

omit regions  
of uncertain  
alignment

- ① optimal doesn't mean correct
- ② multiple alignments often assume a specific tree
- ③ if a region of an alignment is wrong it contributes nonsense to the tree

All alignments  
are the same  
then OK

Gary Olsen II

ALIGNMENT REVIEW

- M McClure - MBE??

METHOD

- ① align
- ② mask
- ③ trees

TREES METHODS

① algorithms - just make tree

② optimality - make tree most consistent w/ data

<del>ALGORITHMS</del>	CLOCK	NO CLOCK
ALGORITHMS	UPGMA WPGMA	NJ
OPTIMALITY CRITERIA	PARSIMONY	
		MAX. PARS.
	LEAST SQUARES	
	KITSCH	FITCH
	MAX. LIKELIHOOD	
	DNA MLK	DNA ML FAST DNA ML PROT ML
	INVARIANT	EVOLUT. PARS



## Distance

~~Don't know them~~  
Only observe differences

∴ must make models about  
relationship betw. observed  
diff. & # events

→ for distance to ML

MODELS	JC	K2	Fels.	Parsimo
Homologous positions	+	+	+	+
Changes along tree branches indep	+	+	+	+
Positions indep	?	?	?	+/
Indels?	NO	NO	NO	+/
All sites same rate	+	+	-	-
Some base freq	+	+	-	-
Trans = Transver	+	-	-	-
Transitions all =	+	+	+	-
Transvers all =	+	+	+	-
Changes all rare	-	-	-	+

v. fast

- doesn't allow vary change

## Non-equal distance

- still adding

## Distance

- only observe differences  
∴ must make models about relationships betw.  
observed difference & # of events

• use to segs

which the expected #  
be observed

probability of giving  
differences observed today

- ① groups most similar
- ② averages similarity to all others
- ③ takes next group

Gary Olsen II

ALIGNMENT REVIEW

- M McLure - MBE??

METHOD

- ① align
- ② mask
- ③ trees

TREES METHODS

① algorithms - just make tree

② optimality - make tree most consistent w/ data

<del>ALGORITHMS</del>	CLOCK	NO CLOCK
ALGORITHMS	UPGMA WPGMA	NJ
OPTIMALITY CRITERION	PARSIMONY	
		MAX. PARS.
	LEAST SQUARES	
	KITSCH	FITCH
	MAX. LIKELIHOOD	
	DNA MLK	DNA ML FAST DNA ML PROT ML
	INVARIANT	
		EVOLUT. PARS

Parsimony - fewest events to give rise to seqs

Least sq. dist - want the history for which the expected # of differences best fits the observed

ML - the history w/ the highest probability of giving rise to exactly the sequences observed today

### TERMS

- ① similarity = DNA vs. prot. difference
- ② dissimilarity = 1-S
- ③ distance = 0 ~ Identical  
higher = increase

### UPGMA

- cluster analysis
- groups by similarity
- tree vs. clock-like
- v. fast
- doesn't allow vary change

- ① groups most similar
- ② averages similarity to all others
- ③ takes next group

### Non-equal distance

- still additive

### Distance

- only observe differences
- ∴ must make models about relationships between observed difference & # of events

# Distance

- ~~① don't know them~~
- ② only observe differences

~~∴ must make models about relationship betw. observed diff. of events~~

→ for distance & ML

<u>MODELS</u>	<u>JC</u>	<u>K2</u>	<u>Fels.</u>	<u>Parsimony</u>
Homologous positions	+	+	+	+
<sup>changes along</sup> Tree branches indep.	+	+	+	+
Positions indep.	?	?	?	?
Indels?	NO	NO	NO	+/ -
All sites same rate	+	+	-	-
Same base freq	+	+	- <sup>x</sup>	-
Trans = Transver	+	-	-	-
Transitions all =	+	+	+	-
Transvers all =	+	+	+	-
Changes all rare	-	-	-	+



## TERMS

- ① similarity : DNA vs. prot. difference
- ② dissimilarity : 1-S
- ③ distance : 0 = ID  
higher = increase

## UPGMA

- cluster analysis
- groups by similarity
- tree is clock-like
- v. fast
- doesn't allow unequal change

① groups most similar

② averages similarity to all other

③ takes next group

## Non-equal distances

- still additive



# MAXIMUM LIKELIHOOD

- ① same models as distances
- ② traces evolution of trees

## WHAT'S IN AN ARCHAEA GENOME

### ① RANDOM DNA SEQS

② Tata binding protein

- i-like EUKS

- branches deeply

- internal duplication

- using T. celis method

T. celis branches

most deeply

- PROT-ML

- Hasegawa program

- v. slow

- see if it do any branch sweeps

# LEAST SQUARES

$$\text{ERROR} = \sum_{\alpha = \text{seq. pairs}} (X_{E\alpha} - X_{O\alpha})^2$$

expected      observed

# WEIGHTED LS

$$\text{ERROR} = \sum_{\alpha} W_{\alpha} (X_{E\alpha} - X_{O\alpha})^2$$

- ① Fitch-Margoliash =  $\frac{(X_E - X_O)^2}{X_O X}$
- ② Cavalli & Edwards = none
- ③ Inverse uncertainty of the differences =  $\frac{(X_E - X_O)^2}{\sigma_0^2}$

DISTANCE IS FROM THE BIOLOGY

- compartmented from tree
- esp. useful for non-char distances

## Tc. celer

- prot sequence has changed  
less than EUKS

- distance betw. 1st & 2nd  
half w/in species

- Tc. celer is lowest suggest.  
it is the least recombined

Bary Olsen - Segeedit

dir data:

sequence

#

#

- load/locus = \* "file" all loci in genbank format

- name of sequence set = default file name

-o seq = sequence file

- write = save



Joe Felsenstein

LIKELIHOOD & PHYLOGENY

- ① Ideally one should have weighted parsimony that
  - ② takes into account non-parsimonious as well as parsimonious reconstruction
  - ③ weighted differently in branches of diff. length
  - ④ weight diff kinds of events

He Says -- maximum likelihood does this

Maximum Likelihood - Invented by Fisher one of the most unpleasant people that ever lived

① Disadvantages

- need model for change
- computationally slower than parsimony
  - ① probabilities
  - ② Joe wrote it

① Justification

- H<sub>1</sub> hypothesis 1 -- that the branch length is 0.1
- H<sub>2</sub> hypothesis 2 -- that the branch length is 0
- D data
- I symbol for given

$$\frac{P(H_1 | D)}{P(H_2 | D)} = \frac{P(H_1 \& D) / P(I)}{P(H_2 \& D) / P(I)} = \frac{\text{Prob } D | H_1}{\text{Prob } D | H_2} \frac{P(H_1)}{P(H_2)}$$

↑
↑  
 posterior odds ratio                      likelihood ratio                      prior odds ratio



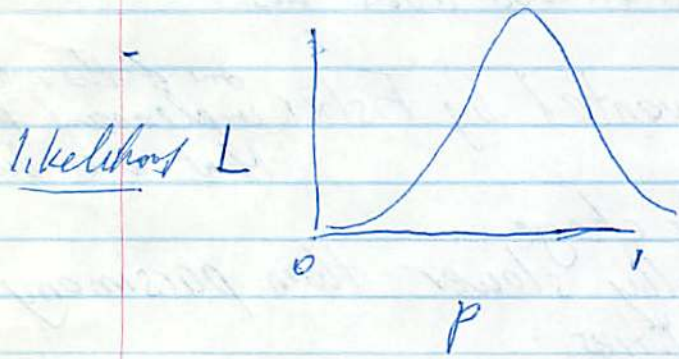
- prior prob. rarely known
- but if enough data likelihood ratios outweigh it

Coin tossing

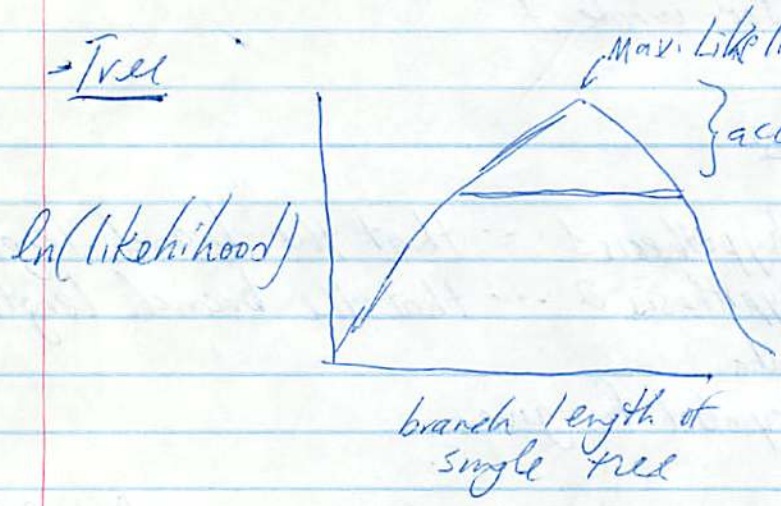
= H H T H T T T H T T H

$p$  = prob. heads  
 $1-p$  = prob. of tail

- prob of this data is  $p^5(1-p)^6$



Tree

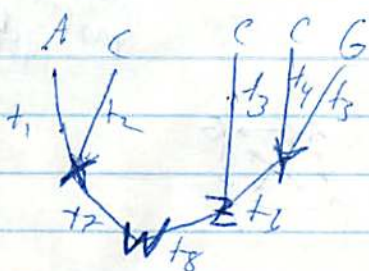


- 1/2 the value of  $\chi^2$  chi-square w/ 1 d.f. significant at 95%

( $\sim$  multiplying by 0.12)

Likelihood for Trees





### ① Need

- alignment
- probability model for how bases change
- trees w/ branch lengths to evaluate

### ② Assume

- independent ~~align~~ evolution of sites

$$\text{Likelihood} = P(D)_{\text{site 1}} \times P(D)_{\text{site 2}}$$

### Probability models

- Markov models (past state not important)

① compute prob. of data for each possible path of state changes

$W \rightarrow X \rightarrow A$   
 $W \rightarrow X \rightarrow C$

$$L_i = P(W) \cdot P(X|W, t_7) \cdot P(A|X, t_1) \dots$$

x given W

$$L = \sum \text{for all possible } x, y, w, z$$

because want likelihood for all that tree for all unknowns



Origin of  
numerical  
taxonomy  
1964.

## MODELS OF DNA CHANGE

① Jukes-Cantor 1969

- all equal  $\frac{1}{4}(1 - e^{-4\alpha t})$   
- P-change =  $\frac{1}{4}(1 - e^{-4\alpha t})$



② Kimura 2 Parameter 1980

- transitions diff. from transversions

③ Felsenstein

- allow unequal base freqs  
- transitions vs. transversion

④ Hasegawa et al 1984

~~All models of~~

Branch length

- are product of rates and time

But not all sites change at same rate

So he has a method

- suppose there are three rates  
- can do MLE for tree

but what  
about  
changes  
over  
tree

$$D = \left( \sum_{\text{paths}} \text{prob. data given tree \& path} \right) (\text{Prob path})$$

transition/transversion  
ratio



to use the likelihood machine and crank the handle and something will drop out

## Evolution w/in species

① Likelihood based method

② random mating method (w/ selfing allowed)

- coalescence
- trees of genes
- depth depends on effective population size

③ what do we want to know - effective pop. size

④ trees of genes

⑤ Kingman's coalescence

- randomly branching tree
- from present - past lineages come together one after another



$t_k$  = time to coalescence is drawn from an exponential distribution

$$= 4N_e^{\text{effective pop. size}} / k(k-1) \quad k = \# \text{ of branches}$$

time to coalescence increases as you go back because there are fewer taxa



How estimate  $N_e$ ?

① Two levels of variability

② randomness of coalescence of lineages  
means you ~~cannot~~ don't get similar trees

- can improve this by looking at more loci & more gene copies

③ randomness of substitutions  
" " " " mutation

- can use more sites if all part of same tree

$$\text{Likelihood} = L = \text{Prob}(\text{Data} | \Theta)$$

$$= \sum_{\text{genealogies}} \text{Prob}(\text{tree} | \Theta) \cdot \text{Prob}(\text{Data} | \text{tree}) \quad \Theta = \frac{4N_e\mu}{\text{time}}$$

∴ all they can estimate is product of size and rate

$$= \sum \text{Prob}(\text{Tree} | \Theta) \cdot \text{Prob}(\text{Data} | \text{Tree})$$

known prior

likelihood of tree

many to do



How do this?

① sum of all trees can be reduced by subsampling trees

- subsample trees by Hastings-Metropolis method  
→ method samples preferentially those of interest by using likelihood

SIMULATED  
ANNEALING

→ method gradually lowers the maximum, so that it ends up on highest peak

Recombination

PROGRAM

LAMARC -- Likelihood Analysis Metropolis

↳ COALESCE

C program -- available where phylog is

How do this?

① sum of all trees can be reduced by subsampling trees

- subsample trees by Hastings-Metropolis method  
→ method samples preferentially those of interest by using likelihood

SIMULATED ANNEALING → method gradually lowers the maximum, so that it ends up on highest peak

Recombination

PROGRAM

LAMARC -- Likelihood Analysis Metropolis

↳ COALESCE

c program -- available where phylog is

# Computers

## New Data Set

- 21 segs
- two tricks
- trim down to do MLE

0835

David Swofford

## PAUP ON EVOZY

- evol4

paup - need to be in directory w/ data files

help (lists commands) command name

execute (runs file) file name

set pause - will pause every screen full

set nopause

showtrees -

showtrees all

- TAXA block

- some commands only need taxa info  
(e.g. consensus)

- delete 7- (delete 7-end)

- criterion likelihood



## Parse

Mania - Multiple alignment never is automatic

## PHYUP SUN

infile

outfile

treefile

23 dnapars type the letter to toggle options

INTERLEAVED - like an alignment

- DNA DIST - calculates distance matrix

- FITCH - Fitch tree

- DRAW TREE set to techtronics emulator

## Problem

① choose masks ...

DNA } SAME MASK ... ALL GAPPED POSITIONS  
PROT } PLUS 80

② copy files to all accounts

③ PHYLIP

- RUN PHYLIP

- SAVE EDIT

← fastest

MOLE(7) 8 - α's

EVOL3 -- v. fast

.  
AGCT  
0111

WEDPM - PCR



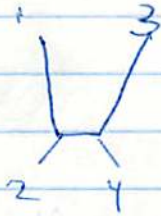
Peter Waddell -

Calculation of probabilities of sequence patterns  
for a specified weighted tree.

edge weight = branch length

① Take Tree

eg



② Break up into path sets

- path set - set of non  
intersecting distances  
among even # taxa

$d_{14}$  = #  
 $d_{24}$  = #  
 $d_{1234}$  = path set  
⋮

matrix eg purnies/pyramindines

	1	2	3	4
state 1	1	0	0	0
2	0	0	0	1
3	0	1	0	1
4	0	0	1	1



$\rho$  true Paths

step ① ↓ break down to more accurate dists

$r$  Smaller paths

step ② ↓

$s(t)$  pattern probabilities

### STEP 1

after expressing tree as vector calc. all paths equally

### STEP 2

shrinking tree generalized distances into observed distances

### STEP 3

using observed mismatches counts to solve for pattern probabilities

## LOG DET

- Lockhart et al 1994 } sites are evolving identically w/ reasonable independence } assumes all sites evolve at same rate

- M is any 4x4 rate matrix

- Good at dealing w/ GC content

-  $\ln(\det(F_{ij}))$  = corrected dist

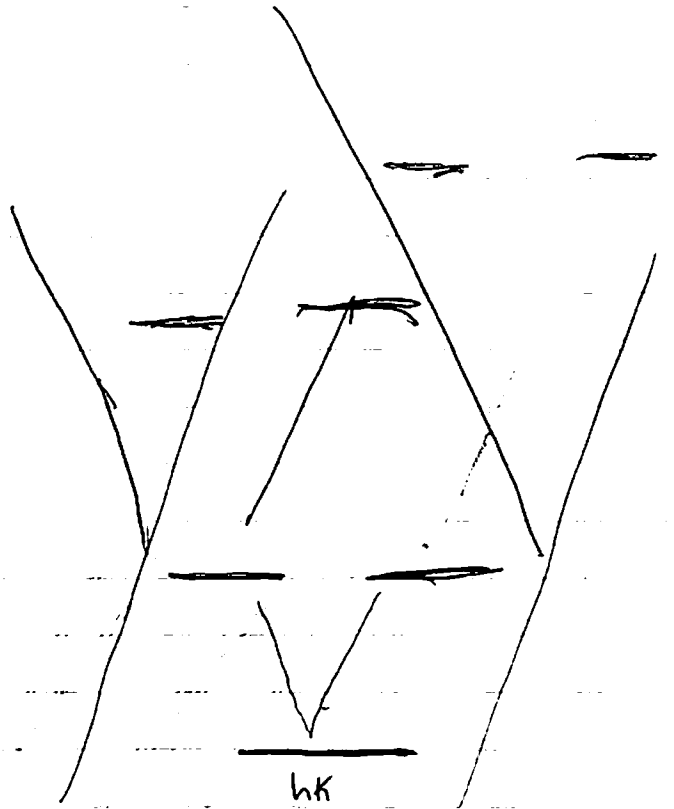
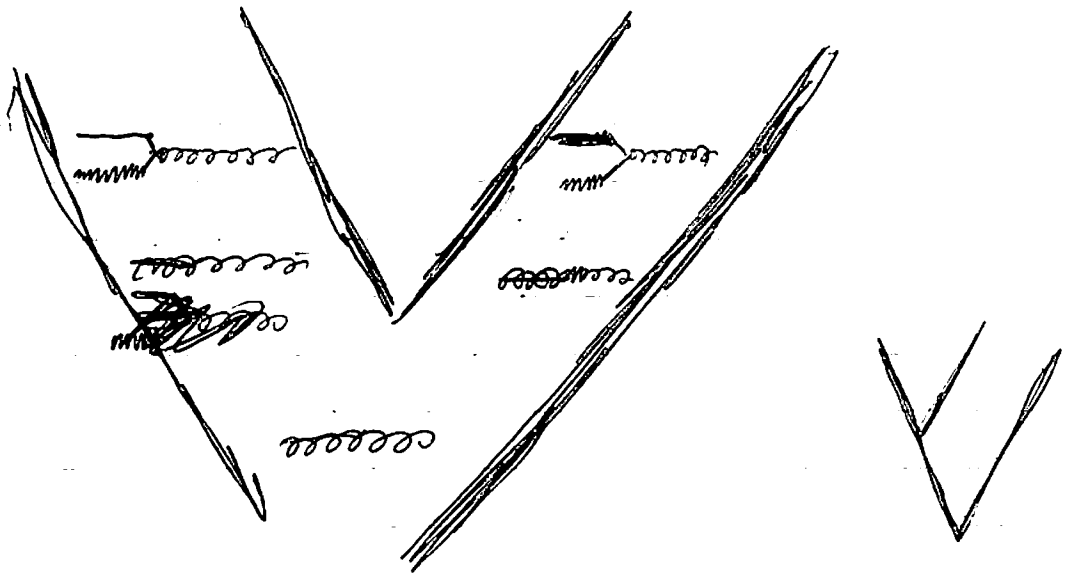
= matrix of observed divergences

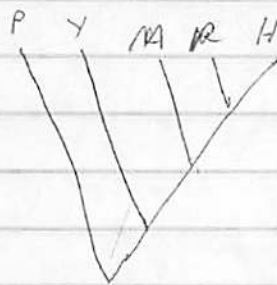
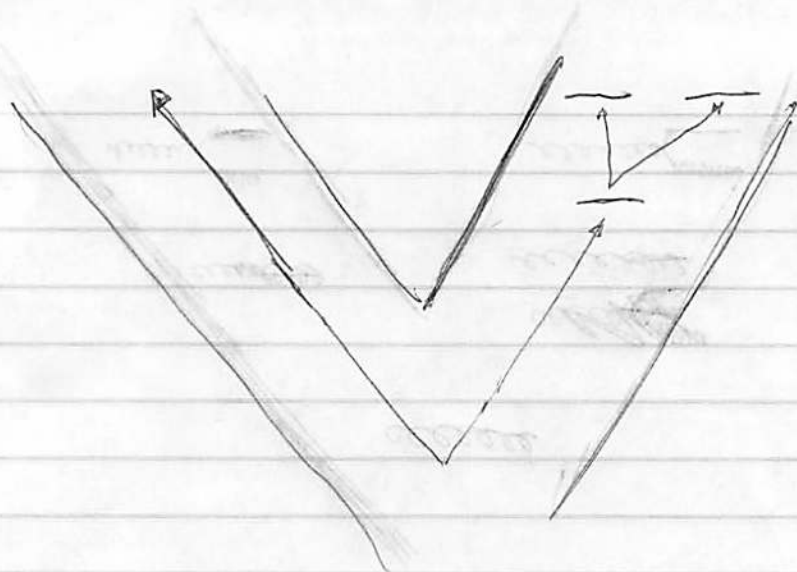
→ might have high variance

$F_{ij}$  = 

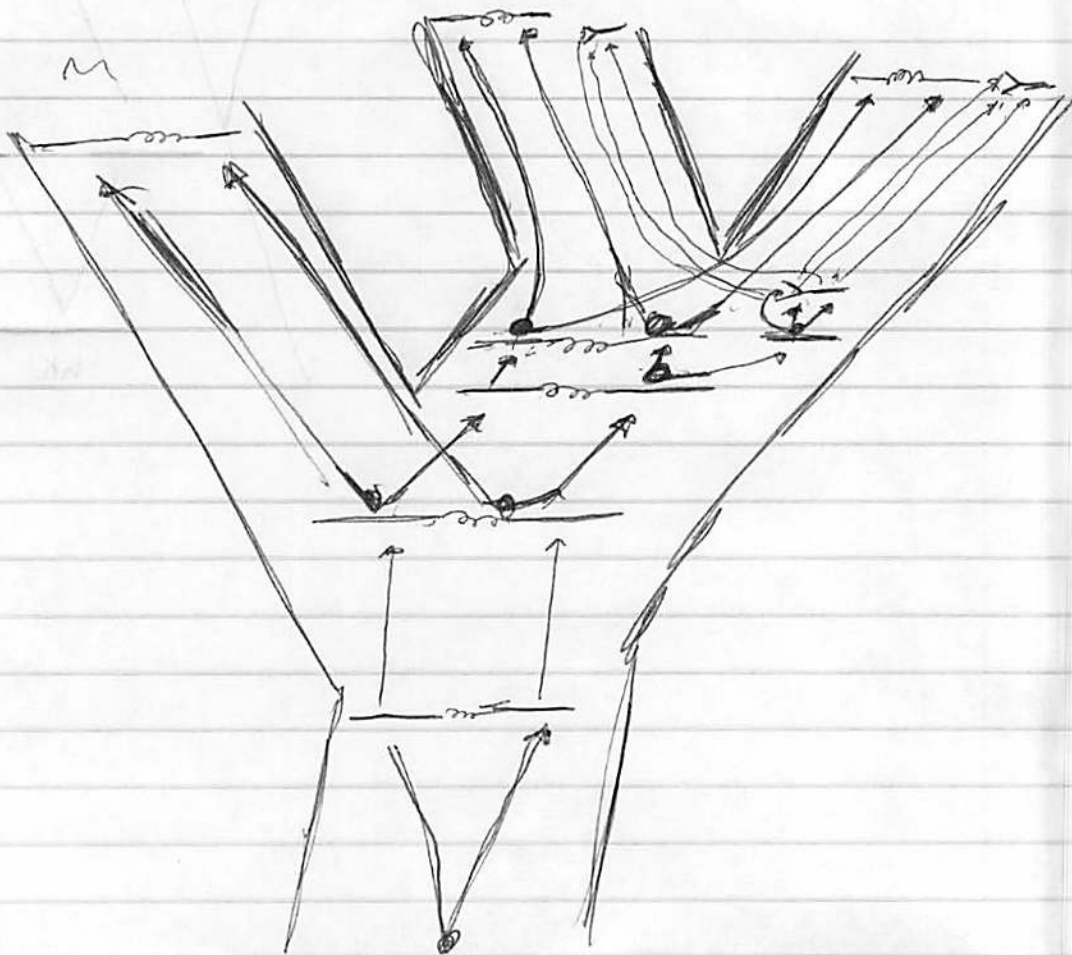
	A	C	G	T
A				
C				
G				
T				

 } sums to 1

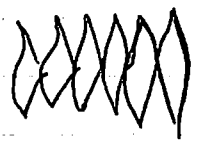




overlay L term  
 overlay R term



Bruce Walsh



- Coalescence: vs. Drift in genealogies

- coalescence is dependent on drift and population size

- it can be decoupled from mutation

- Mutation of alleles

- ① infinite alleles } work for DNA but not great
- ② infinite sites
- ③ finite sites
- ④ single event stepwise

Drift & Mutation & Selection

① Selection

- too many possibilities to model in the real world

Selection vs. Pop. Size

e.g. Codon bias ... how happens in bacteria

① Eukaryotes - strong selection

② Bacteria - large pop. size

Substitution Rate

$$K = 2NuU(2N)^{-1} = \text{subst rate}$$

neutral mutation rate is dependent on population size but pop. size determines if a mutation is neutral



"Large organisms have slow generations & small pop. size"



∴ clocks that are correlated w/ time ~~and~~ not generations ~~but~~ generations may be so because those w/ longer generations have smaller pop. size.

### Clock problems

- ① can say that due to pop. size diff.
- ② can say that due to physiology
- ③

~~Steve O'Brien~~ Steve O'Brien

Conservation Genetics

Cats

- (like butterflies) lots of field observation
- started w/ cheetah breeding problem
- cheetahs

- 60-70 mph ~ 100 yards

- v. popular historically

- but never really able to breed them

- ~~obs~~ - approached O'Brien

- infant mortality v. high

- cheetahs in SA

- abnormal spermatozoa

- tons of defects

- 71% ~~of~~ morphological abnormal

- isozymes

- cheetahs v. low heterozygosity (0 variants)

- other genetic characters same

- mtDNA

- microsatellites - none in some

- skin grafting

from self

from cat

from other cheetah

} - ~~ok~~ ok  
 } - ok  
 } - reject

- MHC

- o'



No more prizes for predicting rain  
Prizes only for building ark

### Separate problem

- homogenization of immune system  
leads to great susceptibility to virulent pathogens

~~copy~~

- virus that normally kills v. few cats  
killed ~~to~~ 50-70% of the cheetahs

### Hypothesizes bottleneck

When?

- ① hunting
- ② ~10,000 years ago many mammals died  
- cheetahs in NA died

### Rapidly evolving genes/loci

- multilocus AFLPs num. satellites
- using human time estimates

### Variation Lions

- High - ① Serengeti ~ 3000 lions
- Medium - ② Ngoro goro crater - 100 lions  
- known bottleneck
- None - ③ lions in India - known bottleneck

Same pattern for - sperm morphology  
- testosterone levels

### Florida Panther/Puma

Hybrids are fine

- same story -- v. little variation
- Everglades -- different characters  
came from Everglades park

Demographic  
collapse  
probability  
eg. log growth  
chaos problem

Observing  
increases in  
congenital  
defects

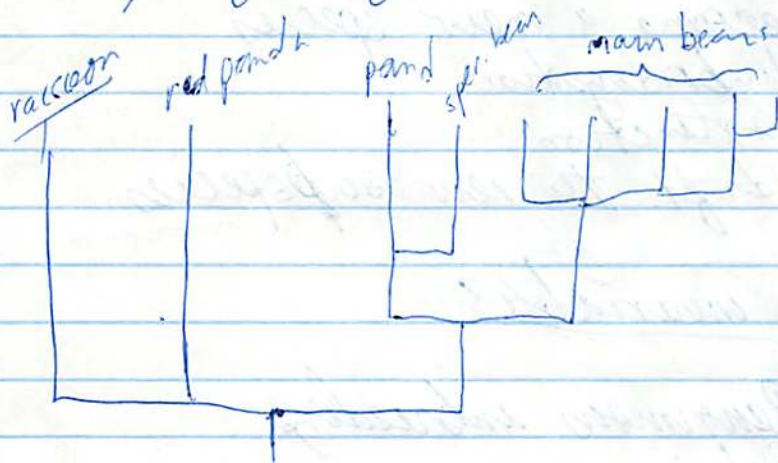


Steve O'Brien

## Phylogenetics for Ucers

②.9 Panda

① morphologically ~~was~~ somewhat like red pandas



④ 125 mtDNA

⑤ jackknifing - removing the outgroups one by one

② Cats

- mtDNA - RFLPs were too big
- a piece of mt like DNA is in nucleus
- chromo  $\alpha$
- integrated mtDNA sequence

- Variation

- can compare rates in mt & nucleus

USFWS

- hybrid policy
- does an endangered species that has crossed w/ another species deserve protection

## Hybrid Policy

- based on Biological Species Concept
- but should not be applied to subspecies

## Potential Fates of Sub-species

- ① become a new species
- ② hybridization
- ③ extinction
- ④ drift to new subspecies

## Why so invariable?

- ① enormous inbreeding
- ② drive

inbreeding



8/13/94

Steve O'Brien

Retrovirus

- HTLV I
- HIV
- RSU
- Hep B

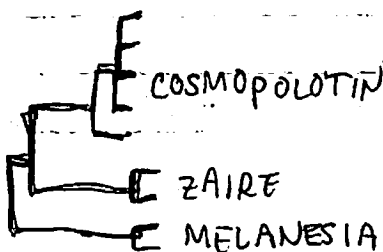
HTLV I

- sequence variation in env
- occurs in many parts of world

- found in many simian species

ENV TREE

- 600 bp - humans



simians

- ⇒ viruses broken up - not specific for host
- asian viruses deep branching
- when include humans get interleaved branches

LENTIVIRUS

HIV, SIV

- SIV doesn't kill African Macaques
- SIV does "Asian"
- FII - not in Asia
- evolves v. rapidly
- looked at pol

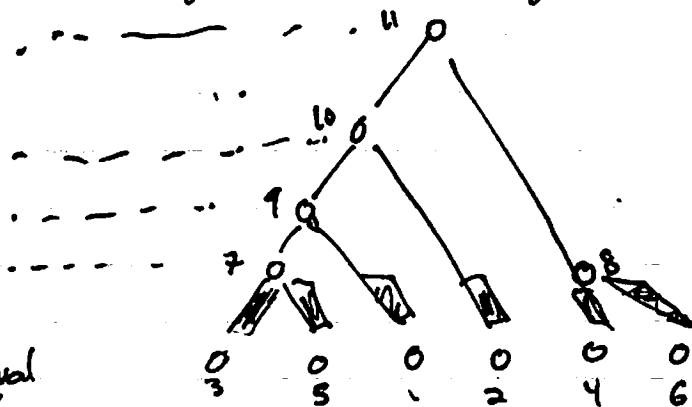
00001

ADL  
HLA

- in MHC region
- synonymous vs. non-synonymous subst.

## SUMMARY - uses of phylogenetics

- ① taxonomic recognition
- ② rank endangerment
- ③ evaluating flux adaptation
- ④ hybridization
- ⑤ subspecies vs species
- ⑥ virus divergence
- ⑦ genomic vestiges

Dick Hudson - CoalescenceInferring Population Genetic Parameterspop. genet models  $\rightarrow$  gene trees  $\rightarrow$  genetic variation

- what is topology?
- what are branch lengths?
- what is sum of branch lengths?
- what is total time?

T6

T6

= time interval  
in which there  
are (#) lineages

$T(\#)$  - have statistical distribution  
# = in generations

- $E(\# \text{ diff. betw. pairs of seqs})$
- $E(\# \text{ polymorphic sites in a sample of seqs})$
- frequency spectrum

$n_1$  - # of ~~mutational~~ <sup>polym. sites</sup> in which 1 seq. has derived state  
 $n_2$  - # " " " " 2 " " " "

Assumptions

① Neutral mutation process ... assumes all variation is due to neutral mutations -- no selection on this variation

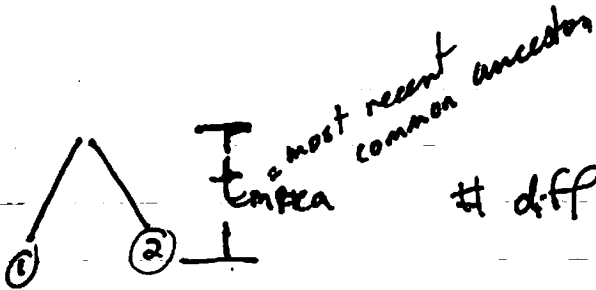
② Offspring copies of the gene differ from parents by a Poisson distributed # of neutral mutations with mean ( $\mu$ ) - mutation for whole region

③ Mutations in diff. indiv. & at diff. times are independent

- ④ 'u' is constant
- ⑤ ∞-site (no back mutations, multiple hits)

① Example

- GIVEN TREE



# diff.  $\rightarrow$  = poisson dist. w/ mean 2 ext

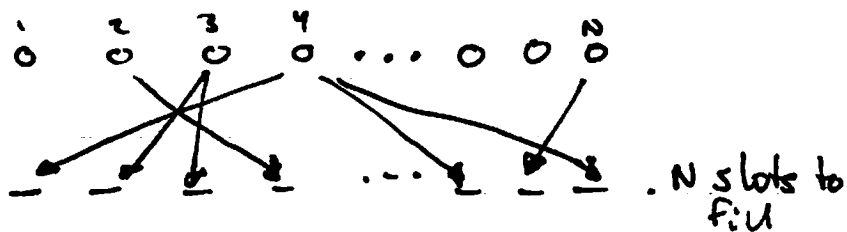
- CAN DO WITH ANY GIVEN TREE

② POP. GENETICS

② Wright-Fisher sampling (to produce succeeding generations)

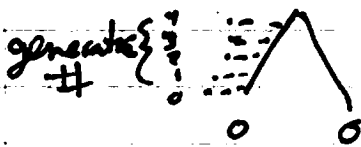
- Haploid  
- For diploids change everything to  $2N$ .

parent generation  
offspring generation



- fill N slots
- sampling randomly w/ replacement (distrib. of offspring is poisson w/ mean 1) if N is v. large)
- many generations

③ Sampling



= time of most recent common ancestor

$$P(t_{mrca} = 1) = \frac{1}{2N}$$

$$P(t_{mrca} > 1) = 1 - \frac{1}{2N}$$

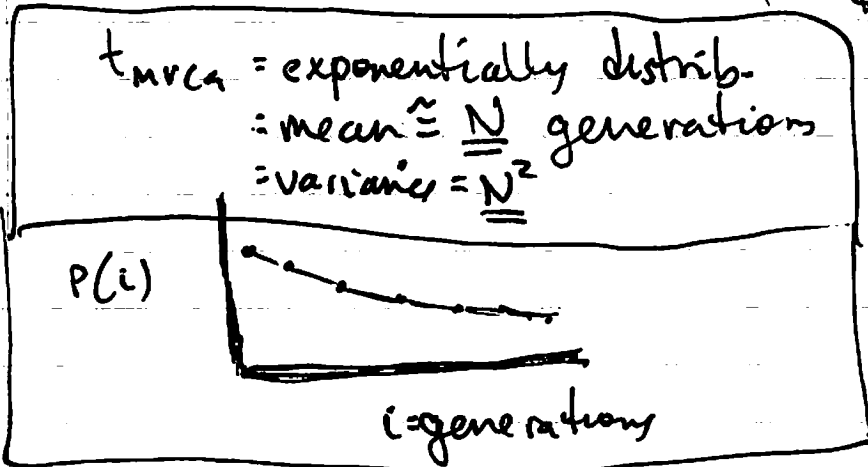
SQ

- if N is large  $P(t_{mrca} = 1)$  is low



$$P(t_{mrca} > t) = \left(1 - \frac{1}{N}\right)^t$$

$$P(t_{mrca} = i) = \left(1 - \frac{1}{N}\right)^{i-1} \frac{1}{N} = \left(\text{prob. going back to } i-1\right) \left(\text{prob. of an } i-1 \text{ of coalescence at } i\right)$$



$$E(\# \text{ diff betw. 2 seqs}) = 2\mu E(t_{mrca})$$

$$= 2\mu N$$

$$= \Theta$$

### SEQUENCE DIFF

(A)  $P(2 \text{ seqs have no diff}) = e^{-2\mu t}$  (given  $t$ )  
 $E(e^{-2\mu t_{mrca}}) = \frac{1}{1+\Theta}$

(B)  $P(2 \text{ seqs have no diff})$



- follow lineages back and count likelihood of mutation  $\neq$  common ancestor

- estimate of likelihood of common ancestor or mutation  $\approx \frac{1}{N} + 2\mu$

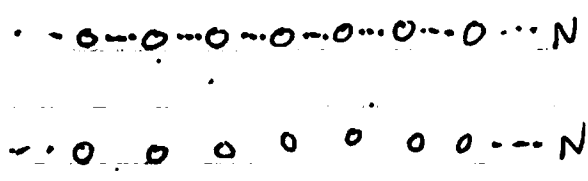
- estimation of likelihood that that was a mutation is  $\frac{2\mu}{\frac{1}{N} + 2\mu}$

- estimation of likelihood that 1st event is ~~mutation~~ coalescence is  $\frac{\frac{1}{N}}{\frac{1}{N} + 2\mu} = \frac{1}{1+\Theta}$

$$\text{prob. } j \text{ diff} = \left(\frac{2u}{N+2u}\right)^j \left(\frac{1}{N+2u}\right)$$

$$= \left(\frac{\theta}{1+\theta}\right)^j \left(\frac{1}{1+\theta}\right)$$

SAMPLE SIZES BIGGER THAN 2



what is prob. ~~that~~ if we have  $N$  indiv. that they all have distinct parents one generation back?

$N=4$

$$= (\text{Prob. that } 1 \text{ distinct from } 2) (\text{Prob. that } 3 \text{ dist. from } 1,2)$$

$$(\text{Prob. that } 4 \text{ dist. from } 1,2,3)$$

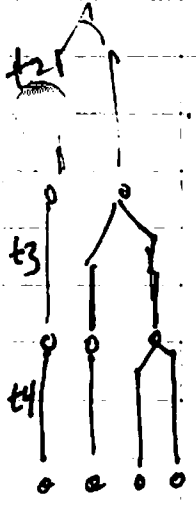
$$= \left(1 - \frac{1}{N}\right) \left(1 - \frac{2}{N}\right) \left(1 - \frac{3}{N}\right) = 1 - \left(\frac{1}{N} + \frac{2}{N} + \frac{3}{N}\right) + \left(\frac{a}{N_2} + \frac{b}{N_2} + \frac{c}{N_2}\right)$$

$$= 1 - \frac{\binom{N}{2} N_{\text{choose } 2}}{N}$$

$$\binom{n}{i} = \frac{n!}{i!(n-i)!}$$

suggests you can ignore this most of the time

$$= 1 - \frac{N!}{2!(N-2)!}$$



$t(4) \cong$  exponentially variable w/ mean =

$$\frac{N}{\binom{4}{2}} = \frac{N}{6}$$

$t(3) \cong$  " " " =

$$\frac{N}{\binom{3}{2}} = \frac{N}{3}$$

$t(2) \cong$  " " " =

$$t_{\text{mean}} = t_4 + t_3 + t_2$$

$$E(t_{\text{mean}}) = 2N \left(1 - \frac{1}{N}\right)$$

$n = \text{sample size}$

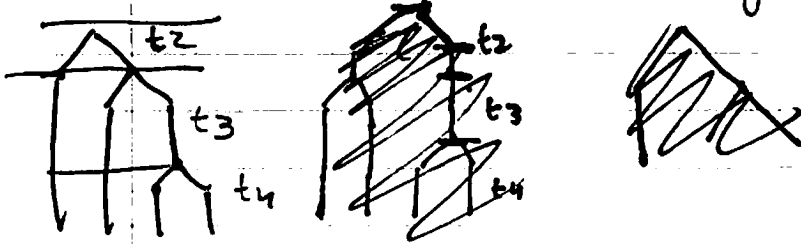
## Offspring #

- Wright Fisher assumes 1 thing
- can repeat w/  $N_e = \frac{N}{\sigma^2} = \frac{\#}{\text{variance in offspring \#}}$

## Mutations

- E # of mutations on tree given tree =  $\mu \cdot \sum_{\text{branches}} \text{distances}$

$$= \mu \cdot E(\sum \text{distances})$$

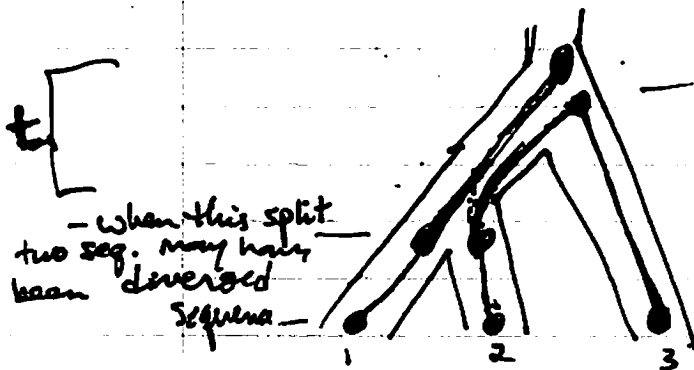


$$= \mu [4E(t_4) + 3E(t_3) + 2E(t_2)]$$

total distance at  $t_2 = 2t_2$   
 " " "  $t_3 = 3t_3$

$$E(\# \text{ polymorphic sites}) = \theta \sum_{j=1}^{n-1} \frac{1}{j}$$

## BETWEEN SPECIES



- possible that 1 & 2 may not coalesce until 2 & 3 coalesce

- prob. that (2,3) coalesce = prob. that 1 & 2 remain distinct & that then 2 & 3 coalesce

- prob. non-concordance of gene & species tree =  $\frac{2}{3} e^{-t/N}$  - this is affected by bottlenecks, selection



can you use  
this for lat. transfer?

## Make tree

- uniform random variable
- =  $-\ln(\text{uniform}) = \text{exponential}$

## Variable Population Size

- ~~short~~ time intervals by fractions  
that  $N$  changes



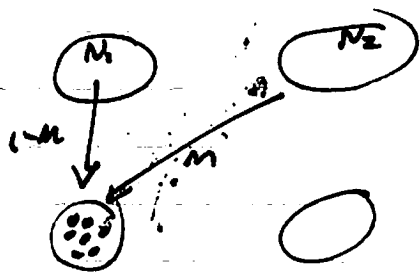
vs



= leads to lots of rare  
variants because more  
like a star tree

## Migration

- essentially use WF model
- but prob. that indiv. comes from other population



$$\text{prob of migration} = \frac{2m}{N+2m}$$

$$\text{mean } t = \frac{1}{2m + \frac{1}{N}}$$

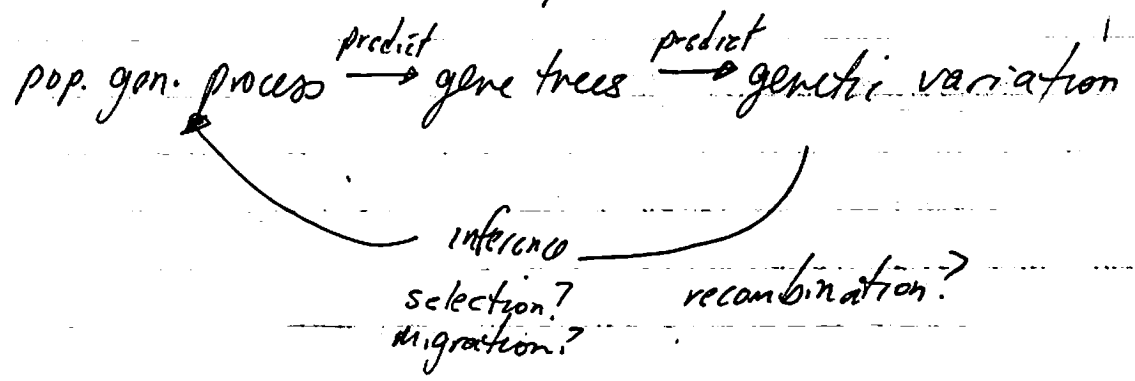
of migration event

$$E(t_s) = 2N \quad \text{- variance not constant w/ } m$$

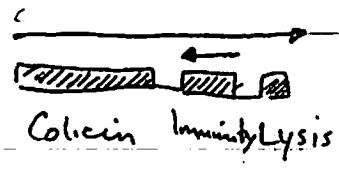
# Peg Riley: My Hero

## Citations

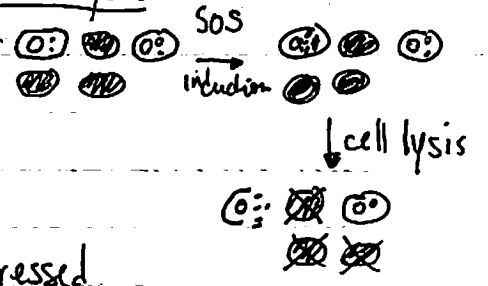
Riley 1993 MBE 10:1048-1059 } Colicin  
 Riley 1993 MBE 10:1380-1395 }  
 Sharp et al PNAS 89:9836-9840 - Restr. Modif.  
 Ohta & Biston, Murphy Evol 1:87-90 } Serine Proteases  
 Ohta MBE 1994 in press }



## Colicin gene cluster



## Life cycle



CELLS ARE  
RESISTANT  
TO THE COLICINS  
THEY CARRY

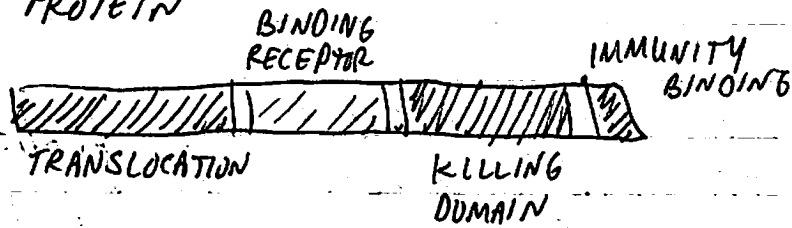
- ① Always on plasmids
- ② Immunity protein usually constit. expressed
  - Col & Lys. SOS inducible
- ③ Colicin recognizes cell surface receptor
  - IF NO IMMUNITY PROTEIN -

- either
- ④ forms ion channel
  - ⑤ enters cyto. → degrades rRNA
  - ⑥ enters cyto. → degrades DNA
  - ⑦ considers them analogous to viruses

## Phenotypic screen

- lawn of E. coli sensitive to all colicins
- toothpick onto lawn
- grow o/n → plaques form

# COLICIN PROTEIN

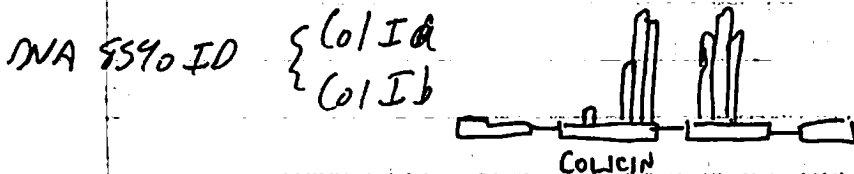
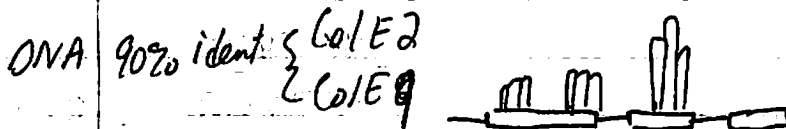
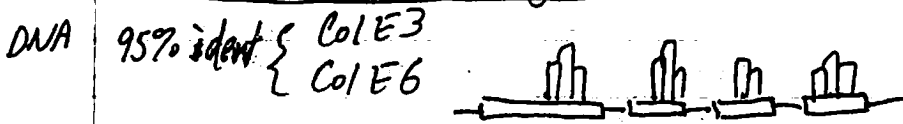


## 5 Classes

- betw. classes no a.a seq. similarity

## Pattern of substitution

### 3 Pairs of Closely Related Colicins



substitutions

- suggest clustering of substitutions
- G test for homogeneity suggests clustering is real

- mutation  
2ary structure

Why variation?

① functional constraint

- but should not have variation in synonym. subst.

② recombinational diversification

- suggests NOT likely bec.
  - too many events
  - need 2nd coli in popn.
  - need 2 in individ
  - need recomb. too

- what about

- is there any info on mutation rate?

- what about the induction?

### ③ mutator

- genome wide (alter trans. transversion)
- locus specific
- seq. specific

### ④ diversification - selection

- like Fitch's virus stuff

- what is avg gpt mut
- what is info on mutation
- what variation in rate would be needed to get this pattern?

## RESTRICTION-MODIFICATION

### TYPE I)

- EcoB
  - EcoK
  - EcoD
  - StySB
  - StySg
- } K Family

- chromosomally encoded
- same location as A

- EcoA
  - EcoE
  - CfrA
- } A Family

- chromosomal
- same location as K

- EcoRII
  - EcoXXI
- } R Family

- plasmid

3 Polypeptides

- R = endonuclease
- M = methylation
- S = sequence recognition

3 possibilities

- if dimethyl - fall off
- if hemimethyl - methyl
- if no methyl - cut



Sharp

- M locus comparison

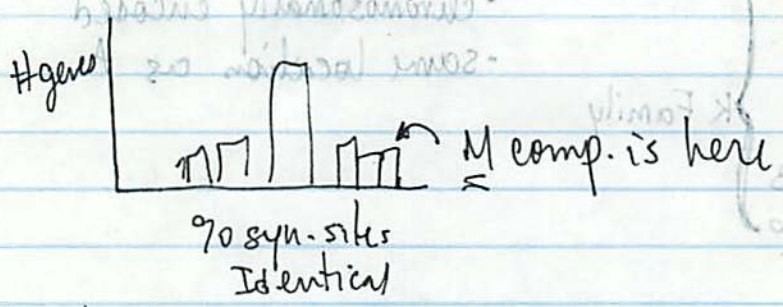
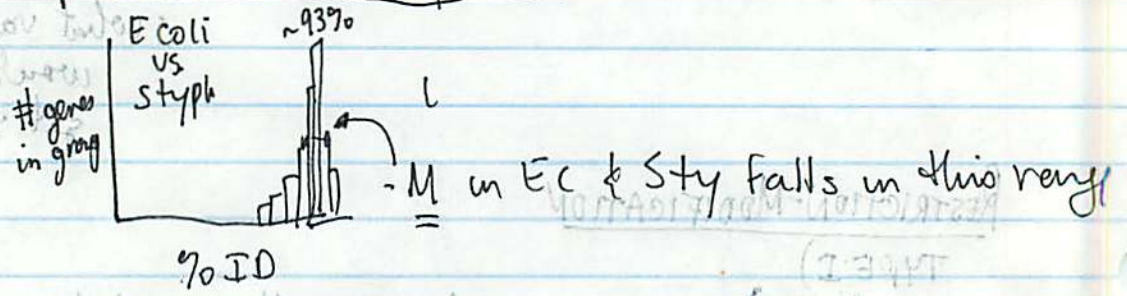
- Intraspecific w/in family - v. highly similar
- Interspecific betw. " - not v. highly similar
- Interspecific betw. " - v. highly similar

Assumes

- ~25-30% ID represents ~~low~~ homology

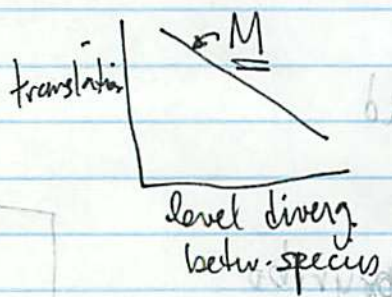
Interspecific vs. Intraspecific

Interspecific



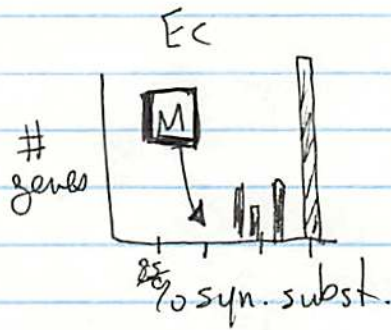
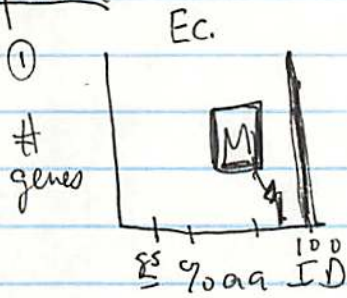
Codon bias

- highly tx genes  $\Rightarrow$  high codon bias  $\Rightarrow$  low divergence betw. species

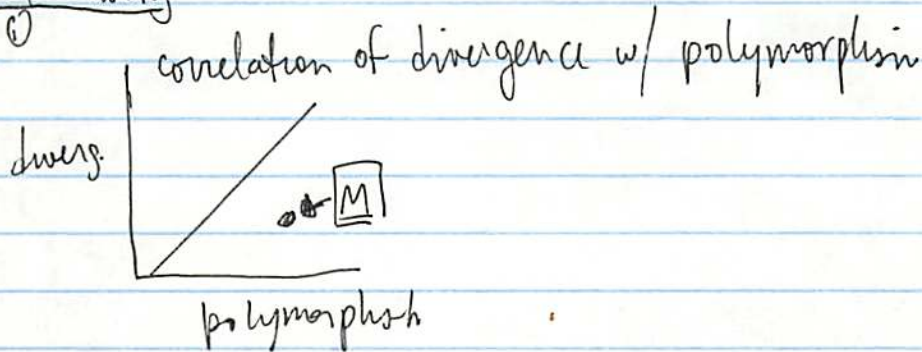


Handwritten notes in the bottom left corner, including "2002-2003" and other illegible text.

Intraspecific



Neutral theory



② Explanation

- ① recombination generates variation
- ② selection at linked site
- ③ selection at M-locus

} either of these could be due to selection for new RM systems

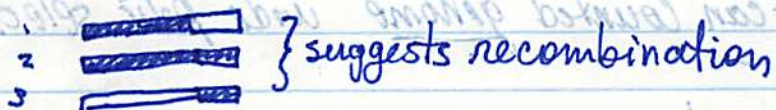
Are there sites at recognition sites?  
 what about variation w/in RM system?  
 would you expect most bias near selection?



8/16/90

# Roger Milkman Evolution of the E. coli chromosome

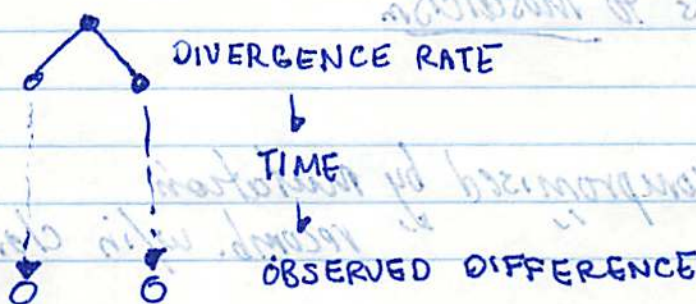
## Mosaicism



## POPULATION GENETICS/EVOLUTION

"gene" in gene tree is essentially a recombination unit

E. coli



## Periodic Selection

- 1951 Atwood & Broyer?

- ① Every once in a while a rare favorable mutation arises
- ② Begins to spread
- ③ Can spread all over world - (w/o recomb. this is 1 genome)
- ④ Could replace all existing forms (essentially a bottleneck)

actually doesn't eliminate all competitors

Confirmed in batch transfers & chemostats

## Mutation

- $\sim 3 \times 10^{-10}$  / nucl. - generation
  - in population of  $10^9$  get most mutations in short time
  - $\therefore$  rare events would be
    - double substitutions
    - triple substitutions
- } esp. if each alone are harmful



9/10/18

*Thermococcus*  
*Sulfolobus acidophilus*

Selection vs. Drift will affect likelihood of spread

Recombination - can count genome wide pdic selection

① in bacteria -

- conjugal
- phage mediated ~50-100 kbp

② rare but lots of time

③ leads to Mosaicism

Clones

① not compromised by mutation  
" " " recomb. w/in clone

can get replaced

② is compromised by recomb. from out of group

③ can be nested



④ can be applied to parts of the genome

Trp Operon

- multiple clones in wild pop. detected by sequence comparison

SAYS E.COW  
15 100 MYO

- calculates divergence times from



- reversible first order reaction



## Recombination

- suggest transformation is likely to not occur  
in *E. coli* in wild

- conjugation - does happen

- restriction enzymes - recombinogenic

8/16/94

## Robin Gutell - rRNA structure

Started w/ H. Noller

- 1st 16s rRNA sequence

### Methodology

A) dot plot against self w/ highlighting potential 2ary structure

⊙ do same thing but do comparison w/ other species

B) Folding algorithms

- searching for stable energies
- look for compensating changes in other species

Woese et al 1980  
NAR

C) Red-Dot Green-Dot

- transitions get red dots
- transversions get green dots
- look for inverse patterns
- continue w/ new sequences
- see if compensatory changes have occurred multiple times

D) Covariation in columns in alignment

- transform nucl. pattern into  $\Delta$  pattern
- search for residues w/ sim. pattern
- ignores WC base-pairs

Gutell  
et al 1985

- COVARIATION MAY BE DUE TO PAIRING  
BUT MAY BE DUE TO SELECTION  
FOR ANTPAIRING

## Triplet Variation

- looking for tricorrelations

- take 2 nucleotides

and search for singles that vary

with this

G C C C C C T T T T A A A A  
G C A T G C A T G C A T G C A T

~~G  
C  
A  
T  
G  
C  
A  
T  
G  
C  
A  
T  
G  
C  
A  
T~~

## Variation in 16S rRNA structure

# seqs	<u>16S</u>	<u>23S</u>
Arch	120	17
Bact	2130	72
Euc	700	42
Organ <del>mt</del>	110	89
CPST	40	32

- w/in groups only a little variation  
in length except in mt & nuclear

- some w/ huge inserts in mature RNA

- Cryptosporidium mt DNA

- fragmented

- Physarum - edited



8/17/94

# NORM FACE: USE OF PHYLOGENY - WHO'S OUT THERE?

Why know so little about natural world?

① most is non-cultivable

- normally can culture  $1/10^4$  of organisms

## Environments

① Open ocean

- seem to be sterile because unculturable

- but lots of stuff  $< 2 \mu\text{M}$

## Molecular Phylogeny

① can provide "definition" of organism by placement on the tree

② can infer physiology/biochemistry...

③ identify cultivatable model

④ physical separation  
⑤

## Methods

① shotgun cloning

- water  $\xrightarrow{\text{filter}}$  biomass  $\rightarrow$  DNA  $\rightarrow$  library  $\rightarrow$  rDNA

- 50% of marine picoplankton are in cyanobacteria and chloroplast group

② PCR

② Epulopisium

- giant bacteria in fish surgeon fish stomach

- no nuclear membrane

- selected organism

- PCR ... groups w/ Gram + 's (close to Clostridium)

- ~~two~~ daughter cells w/in cell - no binary fission and then are released

③ Fluorescent probes

- phylogenetic stains

e.g. - *Nasonia* wasp. SON-KILLER

- closest to *Proteus vulgaris*

④ Survey biodiversity

High Temperature

Chydrothermal vents

② Yellowstone

Octopus Springs

@ 92°C

② lots of biomass

③ three dominant species

H<sub>2</sub> oxidizing

- related to *sq. pyrophilus*

- related to *Th. maritima*

most deeply diverging bacterial group

Sue Barns

Black Pool

@ 80°C - 93°C

② lots of organisms

③ single dNucleotide tracts

④ sequence of small regions

⑤ sequence completes of uniques (>85% diff)

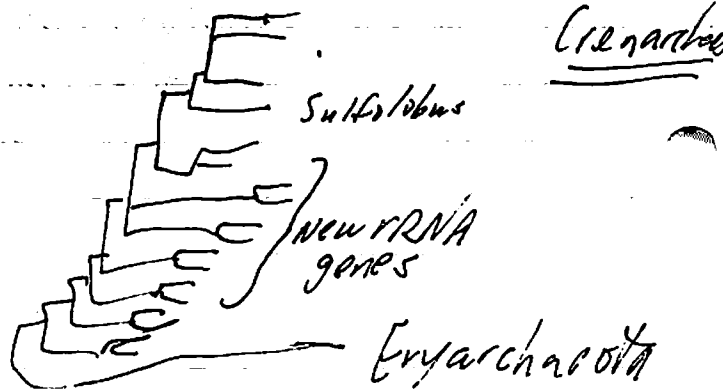
⑥ many new ones are very deeply branching

⑦ new sequences bases

Lake's tree

- Deep branches are unreliable

- (can support w/ signature seqs)



## Hydrothermal vents

① Believes there is a "girdle of slime" circling the earth below the crust

### ② Vents

- some v.v. hot (black smokers)

- ~~most~~ 20-40°C

## High Temperatures

### Contact slides

- if they grow there might get microcolonies

### Juan de Fuca vent caps

- looks like some colonies of *Pyrodictum*  
relative

## RNAse P

- 120 aa protein

- 400 nt RNA

P-Protein

P-RNA

- catalytic

① believe the role of the protein is electrostatic because RNA can fold on own in high ionic solution

## Phylogeny to Predict Structure

RNA secondary structure ... two or more continuous canonical base pairs to give rise to A-form RNA

RNA 3° structure ... single bp

RNAse P cuts here



## RNAse P

① secondary structure

② too many possible 2ary structures

(e.g. DOT PLOT vs COMPLEMENT)

③ so... which are correct

① predict w/ computer (free energy)

② chemical/enzyme assays

- good for 3ary structure

- but can't get inside

③ phylogenetic comparisons

SAYS YOU WOULDNT PUBLISH A CONJECTURAL SEQUENCE



"I'm after the essence of P"

## Phylogenetic Comparison

Look for covariation that preserves complementarity

## RNase P

① B. subtilis vs. E. coli very diff. to align

- ∴ 1<sup>o</sup> seq. evolves more rapidly than rRNA
- solution use more close relatives

② Solved structure somewhat -

- ③ ... a lot of variation in structure
- ④ ... core of about 270 molecules

## Consensus

① problem ... what do you do in variable regions?

- use sequence from those w/o the region

② Commonalities should point to  $f(x)$  regions

## CONING

① heterologous hybridization

② consensus probes

- good for mutagenesis

- similar to hammerhead

- can use as RNA probes

## New Sequences ...

- want them for 3<sup>ary</sup> structure

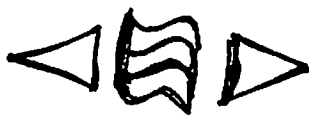
- so using PCR to pull out RNase P from natural populations

- probably related to stacking

- photoaffinity crosslinking

- use RT elongation

- length variation useful to include



George Fox

## Sequence Space

suggests the evolutionary process maintains structural constraints that are (le) necessary

## Structure Space

- a subset of seq. space consisting of only those seqs. consistent w/ biological (le)
- can determine # of seqs that could have a particular structure

## Evolutionary Rates

- rate of events changing sequence
- prob. that event accepted

## Structural evolution

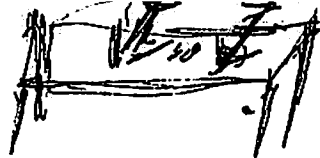
- a seq. change that leads to a change in structure space

### two types

Equilibrium - a sequence can belong (almost) to two structures.  
- single change may favor either

### (c) direct change

SAYS THAT "RARE" EVENTS AREN'T REALLY "RARE" JUST THAT THERE ARE LOTS OF POSSIBLE RARE EVENTS



# EXPLORING SEQUENCE SPACE

10/10/10

Underwater construction using materials that are light  
and strong. The use of concrete and steel.

Use of concrete and steel in underwater construction.  
The use of concrete and steel in underwater construction.

Use of concrete and steel in underwater construction.  
The use of concrete and steel in underwater construction.

Use of concrete and steel in underwater construction.  
The use of concrete and steel in underwater construction.

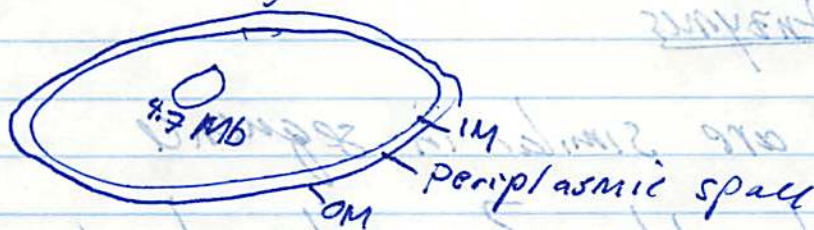
Use of concrete and steel in underwater construction.  
The use of concrete and steel in underwater construction.

Use of concrete and steel in underwater construction.  
The use of concrete and steel in underwater construction.

Use of concrete and steel in underwater construction.

Use of concrete and steel in underwater construction.  
The use of concrete and steel in underwater construction.

Monica Riley



Dan  
Andy  
Jim  
Marc

chemoautotroph  
CHOPKINS - negd

ECOBYE

- knowledge base

Number of ancestral sequences

- How many

- How many also in Archaea/Eukes

OPW comp for all E.coli prots (~2000)

- 899 have at least 1 partner  
= using PAM250

- w/20%  $\gg$  aa IO

- 97% of prot. pairs have similar  $f(x)$

	# prots	# sets
singles	967	967
doubles	252	126
triples	93	31
4-6	384	77
>6	170	14



# Doubled Enzymes

are similar in sequence  
 as it are by shbang? one big duplication?

No relationship betw. position of duplicate.



Number of genes in sequences  
 this many also in Arabidopsis / E. coli  
 829 pair of genes at least 1 partner  
 1200 pairs  
 - w/ 90% >> or 10  
 diff of prot pairs (not similar) [x]

# genes	# pairs	# genes
107	107	107
136	136	136
177	177	177
190	190	190



Laura Landweber

## Genetic Diversity

- Early eukaryotes -

- have interesting & unique nuclear metabolism

Genes

- Static data (DNA sequence) vs. Processed Data

## RNA editing

Definition - any form of sequence modification after transcription

## History

Fagan et al 1988

- Cryptic gene modified in huge consent

- Cyt-c<sub>1b</sub> - edited

## Organisms

- T. Brucei mt have editing

- maybe other euk's do



## Models

① editing only w/in group

② editing ancient - may reflect fossil of early mt



Why?

## Distribution of RNA editing

### ① Herpetomonas

- kinetoplast DNA contains 1000's of minicircles & dozens of maxicircles
- cox III - same location but diff. amt's of editing

### ② amplified w/ flanking genes

#### ② PCR cDNA

- lots of size variation,

#### ② amplify cDNA

- little size variation

←→ U

#### ③ v.v. low levels of G in DNA

v. low

C in RNA

- encoded U's are deleted
- other d's are added
- lots in ORF and in upstream

## Mechanism

### - guide RNAs

- small minicircles

- complementary to edited genes

- synonymy, subs. limited bec. still have to pair w/ RNA not just make prot.

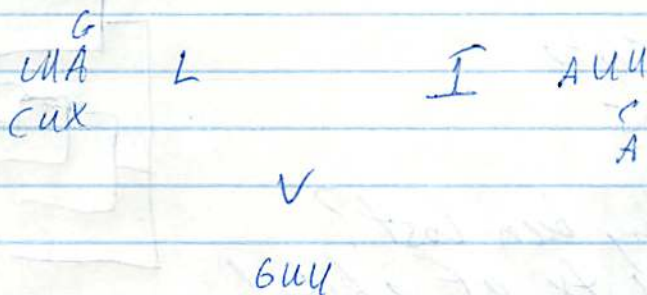


8/18/94

what is mutation rate of guide RNAs?

- what about tRNA selection?

AA subst



- in kinetoplasts L  $\rightarrow$  V changes occur more freq. by UUA  $\rightarrow$  GUU (not CUU  $\rightarrow$  GUU) by a frameshift

DNA

- highly conserved
- almost all A's & G's
- Cs could pair w/ G's & confuse editing

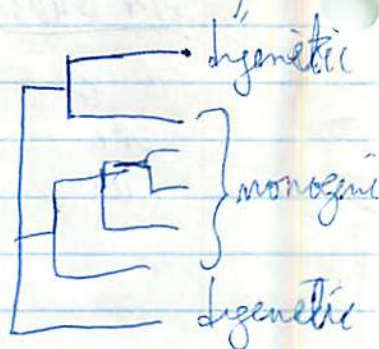
AA

- relatively highly conserved
- variable regions correlate w/ frameshift regions
- similarity among species is low in species w/ editing



## Phylogeny of RNA editing

Nuclear 16S-like  
Mt 94125



How would RNA editing be lost?

- perhaps reverse tx of edited gene can recombine using conserved template

## Mechanism of RNA editing

- e.g. MURF4

① Edited at 5' end

- 1<sup>st</sup> guide RNA corresponds w/ never edited 3' end

② next guide is better paired bec. less wobble

## Why edit?

① allows translational control

② can control AUG - like in humans

- usu. v. last site

## PCR assay

## Other editing

① Plant editing

prob. by deamination → C → U in land mt  
U → G at some sites in fern mt

prob. by deamination A → I - brain glut. receptor

② Physarum mt 16S-like

C insertions + A insertions

Are these bygone systems  
present simply bec. enzymes  
that could do the forward  
reaction were present and of  
course could do reverse

## DNA Editing

Actin I gene in Ciliated protozoa

- two types of nuclei
- micronuclei = germ line
- macronuclei = contained cut DNA
- extensive rearrangements needed
- including reordering

## Mechanism

seems to be guided by flanking regions



Paup

email publish@sinauer.com

FTP

-copy ftpservers.txt (evolv.mbl.edu)  
/home/mbs4/dbd/Internet

## DNA Patterns

### Comparisons w/o alignment

- 4 base window
- count occurrence of each 4 combination
- normal distribution
- compare distribution w/ other organisms
- w/ correlations
- compare catalog strand (6b coding strand)
- correlation of strands v. high
- not due to palindromes
- selected non-coding sequences  $\frac{1}{2}$  similar

### - Correlation

- 17 v. strange

- ~~book~~

18 organisms

mic

hymen

halobits

Electi

Mt -- two strands don't look alike  
yeast its ok