

APL Pipeline

These scripts are written in APL for the APL+Win Version 2.0.00 software. Because APL commands can be executed individually, as well as by script, some of the steps in our analyses were done with individually executed steps rather than with scripts. When this was done, it is explained below. Scripts are listed in APLPLUS font, whereas descriptive material is in Times New Roman (this font).

I. Reading SNPs from file JG1 (in script READCALLS) and SNP indexing

A. The script READCALLS produces 3 variables:

1. ΔSCORES, an N X 62 X 2 variable. Each 62 X 2 submatrix holds the SNP alleles for 62 individuals.
2. ΔTIGS, a vector with contig number corresponding to each SNP
3. ΔPOS, a vector with contig position corresponding to each SNP.

It reads calls from a Table file that was produced using the GATK VariantsToTable tool. In the script, this file is referred to as 'JG1'.

READCALLS

```
A THIS PROGRAM READS IN DATASET JG1 AND CALCULATES ALLELE COUNTS SEPARATELY FOR
A THE DIFFERENT POPULATIONS
A ALSO, MAKES THREE VARIABLES:
ΔSCORES←0 62 2ρ' ' A VARIABLE TO HOLD GENOTYPE CALLS
ΔTIGS←0ρ0 A VARIABLE TO HOLD CORRESPONDING CONTIG NUMBERS
ΔPOS←0ρ0 A VARIABLE TO HOLD CORRESPONDING GENOME POSITION

'C:\JG1' UNTIE -1
SIZE←ONSIZE -1
FILE←ONREAD (-1,82,(SIZE+1000),0) A READ IN FILE
ONUNTIE -1
HEADER←843↑FILE A SEPARATE HEADER
DATA←843↓FILE A REMOVE HEADER FROM DATA
SAMPLES←79↓HEADER
DATA2←DATA,'tig' A APPEND 'TIG' TO END OF DATA
CHET←LHET←AUSTHET←WHET←0ρ0 A INITIALIZE VARIABLES FOR INDIVIDUAL
SAMPLE HETEROZYGOSITIES
TOTC←TOTL←0 5ρ0 A SET UP VARIABLES FOR TOTAL COUNTS FOR
R C AND L FOR 5 NUCLEOTIDES (INCLUDING ★)

CINDEX←(1+131),62 A SET UP INDICES FOR ROWS OF DATA2 TO DISTINGUISH
SAMPLES FROM DIFFERENT POPULATIONS
AUSTINDEX←1
LININDEX←32+129

I←0
RETI:I←I+1
```

```

TEMP6←((I÷1000)-(⌈(I÷1000)))  A REPORT SNP NUMBER EVERY 100 SNPS
⌊(TEMP6=0)/'␣TCLF ⋄ I'

A PICK DATA FOR NEXT SNP--STRIP OFF LEADING INFO AND LEAVE JUST NUC
LEOTIDES FOR DIFFERENT
A SAMPLES. TEMP5 HAS FORMAT G/G G/G G/C . . .
TEMP←5000↑ DATA2
TEMP2←TEMP ␣SS 'tig'
TEMP3←TEMP2/⌊TEMP2
TEMP4←(TEMP3[2]-1)↑TEMP
TEMP5←(TEMP3[1]-1)↓TEMP4

A IDENTIFY NUCLEOTIDE IN REFERENCE GENOME (NOT CURRENTLY USED)
IND0←(TEMP5=␣AV[10])/⌊TEMP5
TIG←IND0[1]↑TEMP5
TIG2←(TIG∈'0123456789')/TIG
TIG3←␣FI TIG2
ΔTIGS←ΔTIGS,TIG3

POS←-1↓(IND0[1]↓IND0[2]↑TEMP5)
POS2←␣FI POS
ΔPOS←ΔPOS,POS2

A REFORMAT TEMP5 AS 62 X 3 MATRIX OF NUCLEOTIDES IN FORMAT G/G
IND1←TEMP5='/'
IND2←IND1/⌊IND1
SCORES←(IND2[1]-2)↓TEMP5
SCORES2←62 4ρSCORES
SCORES3←SCORES2[;13]
CHECK←(+/SCORES3[;2]='/')=62  A CHECK TO ENSURE PROGRAM IS NOT OUT
OF ALIGNMENT
⌊(CHECK=0)/''BAD DATA FOR SNP ',I ⋄ ␣TCLF ⋄ SCORES3 ⋄→0'
ΔSCORES←ΔSCORES,[1]SCORES3[;1 3]

DATA2←(TEMP3[2]-1)↓DATA2  A DROP CURRENT SNP FROM DATA2

TEST←+/(DATA2 ␣SS 'tig')

→(TEST>1)/RETI
␣TCLF
'PROGRAM COMPLETE'
' HETEROZYGOSITIES IN VARIABLES CHET, LHET, AUSTHET, AND WHET'
' NUCLEOTIDE COUNTS IN VARIABLES TOTC AND TOTL'
' CALLS IN VARIABLE ''ΔSCORES''
' CORRESPONDING CONTIGS IN VARIABLE ''ΔTIGS''
' CORRESPONDING CONTIG POSITION IN VARIABLE ''ΔPOS''
'NOTE: THESE VARIABLES SHOULD BE SAVED TO AN APL COMPONENT FILE'
*****

```

A series of index variables are created manually. These are vectors corresponding to rows in the 62 X 2 submatrices of ΔSCORES that correspond to particular subsets of samples. For example, ΔCINDEX and ΔLINDEX correspond to *I. cordatotriloba* and *I. lacunosa* samples. ΔCALLO3 and ΔLALLO3 correspond to known allopatric samples for the two species. \

B. A series of scripts is run to identify codon and codon position of SNPs using the LAC genome and produce a number of indexing variables. The following is an overview of these scripts and how they are used:

1. PROGRAM 'READLACGFF' READS IN THE I. LACUNOSA GFF3 FILE AND CONVERTS IT TO A GFF3-LIKE MATRIX ('ΔLACCCDS') THAT KEEPS ONLY LINES WITH COLUMN 3 = 'CDS'. 'ΔLACCCDS' IS SAVED IN THE COMPONENT FILE 'C:\LACCCDS'
2. 'ΔLACCCDS' IS REORDERED IN ORDER OF SCAFFOLD NUMBER AND PUT IN 'ΔLACCCDS2'. WITHIN A GIVEN SCAFFOLD, THE START AND END POSITIONS OF THE CDS ELEMENT ARE NOT ORDERED
3. 'CONVLAC' IS USED TO REORDER BEGINNING POSITIONS WITHIN SCAFFOLDS AND ELIMINATE DUPLICATE ENTRIES. IT PRODUCES 'ΔLACCCDS4'

THE COLUMNS OF 'ΔLACCCDS4' ARE:

COLUMN1 CONTIG NO.
COLUMN2 STARTPOS IN CONTIG
COLUMN3 ENDPOS IN CONTIG
COLUMN4 STRAND
COLUMN5 PHASE

4. 'FIXLACCCDS4' IS USED TO ELIMINATE ADDITIONAL DUPLICATES FROM 'ΔLACCCDS4' AND PRODUCE 'ΔLACCCDS5'

THE COLUMNS OF 'LACCCDS5' ARE THE SAME AS FOR 'ΔLACCCDS4'

5. RUN 'INDEXLAC' TO MAKE INDEX OF SCAFFOLDS THIS PROGRAM INDEXES THE SCAFFOLDS OF THE I. LACUNOSA GENOME. IT CREATES THE VARIABLE 'ΔLACINDEX', WHICH HAS THE FOLLOWING COLUMNS:

COLUMN1 SCAFFOLD NUMBER
COLUMN2 START POSITION OF SCAFFOLD IN GENOME FASTA FILE
COLUMN3 END POSITION OF SCAFFOLD IN GENOME FASTA FILE

6. 'GETLACSNPS2' READS IN SNPS FROM A TABLES FILE AND MERGES THE INFORMATION WITH 'ΔLACCCDS5' TO PRODUCE THE VARIABLE 'ΔLACSNPS', WHICH HAS THE FOLLOWING INFORMATION:

COLUMN1 CONTIG NO.
COLUMN2 START POSITION OF CDS FEATURE
COLUMN3 POSITION OF SNP
COLUMN4 END POSITION OF CDS FEATURE
COLUMN5 STRAND (0 = -, 1 = +)
COLUMN6 PHASE (CODON POSITION (0,1,2) OF START POSITION OF CDS FEATURE
COLUMN7 ALLELE 1
COLUMN8 ALLELE 2

'GETLACSNPS2' ALSO MAKES 'ΔLACCODONS', WHICH HAS THE FOLLOWING COLUMNS:

COLUMN1 ALTERNATIVE CODON 1
COLUMN2 ALTERNATIVE CODON 2
COLUMN3 WHETHER DIFFERENCE BETWEEN CODONS IS SYNONYMOUS (S) OR NON-

SYNONYMOUS (N)
COLUMN4 SNP NUMBER

'ΔLACSNPS' AND 'ΔLACCODONS' ARE STORED IN COMPONENT FILE 'LACNPS'

7. AFTER RUNNING GETLACSNPS2, RUN 'SEPSNPS' TO MAKE FOLLOWING VARIABLES:

'ΔSYNSCORES' (FROM 'ΔSCORES') SNP DATA ON JUST SYNONYMOUS SNPS

'ΔNONSCORES' (FROM 'ΔSCORES') SNP DATA ON JUST NON-SYNONYMOUS SNPS

'ΔOSCORES' (FROM 'ΔSCORES') DATA ON JUST 'OTHER' SNPS

8. RUN 'SPLICE4' TO CALCULATE NUMBER OF SYN AND NON-SYN SITES. PRODUCES TWO VALUES OF COUNTS:

'ΔTOTSYN' AND 'ΔTOTNONSYN'.

SCRIPTS:

READLACGFF

A THIS PROGRAM READS IN PARTS OF THE I. LAC GFF3 FILE IN 10000000 BY
TE CHUNKS

A FOR EACH CHUNK, IT ASCERTAINS WHERE THE CDS FEATURES AND SAVES THE
LINE GIVING

A INFO OF THAT FEATURE IN VARIABLE NEWGTF2.

A PROCESSING IS AS FOLLOWS:

A 1. FIRST CHUNK IS READ IN.

A 2. ALL LINES WITH FEATURE = 'CDS' ARE KEPT AND APPENDED TO 'N
EWGTF2'

A 3. 'NEWGTF2' IS SAVED TO A COMPONENT OF FILE 'GFF3.SF'

A 4. CHANGE J←0 TO J←1 IN PROGRAM AND COMMENT OUT 'GTF←0p0' AND
RESTART

A 5. NEXT CHUNK IS READ IN AND STEPS 2 AND 3 ARE REPEATED,

A 6. STEPS 4 AND 5 REPEATED FOR ALL CHUNKS

A 7. THE DIFFERENT MATRICES (DIFFERENT NEWGTF2) IN 'GFF3.SF' AR
E MANUALLY CATENATED TO

A FORM MATRIX 'ΔLACCCDS' AND SAVED TO 'GFF3.SF' AS ANOTHER
COMPONENT

A 8. 'ΔLACCCDS' IS MANUALLY SORTED BY SCAFFOLD NUMBER TO FORM 'Δ
LACCCDS2', WHICH IS

A SAVED TO ANOTHER COMPONENT OF 'GFF3.SF'

FLAG←0

J←4

RETJ:J←J+1

DIM←pGTF

'C:\LACGFF1' ONTIE -1

GTF←GTF,ONREAD -1 82 10000000 ((J-1)×10000000)

ONUNTIE -1

ITCLF

'SEGMENT ',J,' READ'

```

  1((10000000+DIM)>ρGTF)/'FLAG←1  ♦  ΠTCLF  ♦  ''FLAG SET TO 1''

  A NEWGTF2←NEWGTF←0 82ρ' '
  IND←(GTF=ΠTCLF)/1ρGTF
  GTF←IND[1]↓GTF

  I←0
  RETI:I←I+1

  TEST←(I÷100)=(I(I÷100))
  1(TEST=1)/''J = '' ,J, ''  I= '' ,I, ''  MBYTES LEFT: '' , (10 5π(ρGTF)÷1
000000)'
  IND←(GTF=ΠTCLF)/1ρGTF
  →(0=ρIND)/0
  LINE←IND[1]↑GTF
  LINE←-1↓LINE
  IND2←(LINE=ΠAV[10])/1ρLINE
  COL3←-1↓(IND2[2]↓(IND2[3]↑LINE))
  TEST←^/COL3[13]='CDS'
  →(TEST=0)/DOWN

  COL1←-1↓(IND2[1]↑LINE)
  COL2←-1↓(IND2[1]↓(IND2[2]↑LINE))

  COL4←-1↓(IND2[3]↓(IND2[4]↑LINE))
  COL5←-1↓(IND2[4]↓(IND2[5]↑LINE))
  COL6←-1↓(IND2[5]↓(IND2[6]↑LINE))
  COL7←-1↓(IND2[6]↓(IND2[7]↑LINE))
  COL8←-1↓(IND2[7]↓(IND2[8]↑LINE))
  COL9←IND2[8]↓LINE
  COL3A←15↑(COL3,15ρ' ')
  COL9A←40↑(COL9,40ρ' ')
  COL4A←10↑(COL4,10ρ' ')
  COL5A←10↑(COL5,10ρ' ')
  NEWLINE←COL1, ' ', COL2, ' ', COL3A, ' ', COL4A, ' ', COL5A, ' ', COL6, '
', COL7, ' ', COL8, ' '
  NEWGTF2←NEWGTF2,[1]NEWLINE

  DOWN: IND←(GTF=ΠTCLF)/1ρGTF
  GTF←IND[1]↓GTF
  →((FLAG=0)^(0=ρIND))/RETJ
  →RETI

  A ΔLACCDs←NEWGTF2

  ΠTCLF
  'PROGRAM READLACGFF COMPLETE.  CDS DATA IN VARIABLE ΔLACCDs'
  'THIS VARIABLE NEEDS TO BE MANUALLY REORDERED

```

CONVLAC

Ⓐ THIS PROGRAM CONVERTS THE CHARACTER VARIABLE ΔLACCD_S2 TO THE NUMERIC VARIABLE ΔLACCD_S3

Ⓐ SORTS IT AND ELIMINATES DUPLICATE ENTRIES, PRODUCING THE VARIABLE ΔLACCD_S4

ΔLACCD_S3←0 5ρ0

DIM←1↑ρΔLACCD_S2

□TCLF

DIM

START←0

I←0

RETI:I←I+1

TEST←(I÷1000)=(⌈(I÷1000))

⊕(TEST=1)/'I'

LINE←ΔLACCD_S2[I;]

CONTIG←□FI LINE[3+18]

STARTPOS←□FI LINE[49+111]

ENDPOS←□FI LINE[61+112]

STRAND←0

⊕(LINE[77]='+')/'STRAND←1'

PHASE←□FI LINE[80]

NEWLINE←CONTIG,STARTPOS,ENDPOS,STRAND,PHASE

ΔLACCD_S3←ΔLACCD_S3,[1]NEWLINE

→(I<DIM)/RETI

□TCLF

'STARTING CONVLAC2'

CONVLAC2

□TCLF

'STARTING CONVLAC3'

CONVTRIF3

ΔLACCD_S4←ΔLACCD_S3[INDEX9;]

□TCLF

'PROGRAM FINISHED. DATA ON I. LACUNOSA CODING SEQUENCES IN VARIABLE
'ΔLACCD_S4''.'

FIXLACCD_S4

Ⓐ THIS PROGRAM ELIMINATES FURTHER DUPLICATES IN ΔLACCD_S4 AND OVERLAPPING CDS FEATURES

ΔLACCD_S5←0 5ρ0

DIM←ρΔUNIQSCAFFS

DIM1←1↑ρΔLACCD_S4

Ⓐ I LOOP LOOPS THROUGH SCAFFOLDS

I←0

RETI:I←I+1

TEST←(I÷10)=(⌈(I÷10))

⊕(TEST)/'I'

SCAFF←ΔUNIQSCAFFS[I]

IND←(ΔLACCD_S4[;1]=SCAFF)/1DIM1

PART←ΔLACCD_S4[IND;]

```

PART←PART[ΔPART[;2];]
DIMPART←1↑ρPART
TEST←1=DIMPART
Δ(TEST)/'ΔLACCD55←ΔLACCD55,[1]PART ◇ →DOWN2'

A J LOOP REMOVES CDS FEATURES WITH DUPLICATE BEGINNING POSITIONS
DUPS←0ρ0
DIM2←1↑ρPART
KEEP←0 5ρ0
J←0
RETJ:J←J+1

LINE1←,PART[J;]
LINE2←,PART[J+1;]

Δ(LINE1[2]≠LINE2[2])/'KEEP←KEEP,[1]LINE1 ◇ →DOWN1'
DUPS←DUPS,J
IND←(PART[;2]=LINE1[2])/ιDIM2

NEQUAL←ρIND
LINES←PART[IND;]

IND2←1↑((ι/LINES[;3])=LINES[;3])/ιρIND
ADDLINE←LINES[IND2;]
KEEP←KEEP,[1]ADDLINE
J←J+(NEQUAL-1)

DOWN1: →(J<(DIM2-1))/RETJ

A K LOOP REMOVES DUPLICATES WITH SAME END POSITION

KEEP←KEEP[ΔKEEP[;3];] A REORDER KEEP IN ASCENDING ORDER OF END POS
ITIONS

DUPS2←0ρ0
DIM3←1↑ρKEEP

Δ(DIM3=1)/'ΔLACCD55←ΔLACCD55,[1]KEEP ◇ →DOWN2'
KEEP2←0 5ρ0
K←0
RETK:K←K+1

LINE1←,KEEP[K;]
LINE2←,KEEP[K+1;]

Δ(LINE1[3]≠LINE2[3])/'KEEP2←KEEP2,[1]LINE1 ◇ →DOWN3'
DUPS2←DUPS,K
IND←(KEEP[;3]=LINE1[3])/ιDIM3

NEQUAL←ρIND
LINES←KEEP[IND;]

IND2←1↑((ι/LINES[;3])=LINES[;3])/ιρIND
ADDLINE←LINES[IND2;]
KEEP2←KEEP2,[1]ADDLINE

```



```

K←K+(NEQUAL-1)
DOWN3:→(K<DIM3-1)/RETK

ΔLACCD5←ΔLACCD5,[1]KEEP2

DOWN2:→(I<DIM)/RETI
*****
INDEXLAC ;STARTBYTE
A THIS PROGRAM INDEXES THE SCAFFOLDS OF THE I. LACUNOSA GENOME. IT
  CREATES THE VARIABLE
A   ΔLACINDEX, WHICH HAS THE FOLLOWING COLUMNS:
A       COLUMN1      SCAFFOLD NUMBER
A       COLUMN2      START POSITION OF SCAFFOLD IN GENOME FASTA FILE
A       COLUMN3      END POSITION OF SCAFFOLD IN GENOME FASTA FILE

STARTBYTE←0
I←0
RETI:I←I+1

TEST←(I÷100)=(↑(I÷100))
↓(TEST=1)/'I'

'C:\LACGENOME' ONTIE -1
PARTSCAFF←ONREAD -1 82 2000000 STARTBYTE
ONUNTIE -1

DIM←ρPARTSCAFF
→(DIM=0)/0

IND←(PARTSCAFF = '>')/1ρPARTSCAFF
ENDSCAFF←IND[2]-1

TEMPSCAFF←ENDSCAFF↑PARTSCAFF
IND←(TEMPSCAFF=↑TCLF)/1ρTEMPSCAFF
HEAD←IND[1]↑TEMPSCAFF
SEQ←IND[1]↓TEMPSCAFF
SEQ←(SEQ≠↑TCLF)/SEQ
ΔSCAFF←↑FI HEAD[4+18]

SEQ[110]
READLACSCAFF ΔSCAFF
ΔSEQUENCE[110]
TRSH←□

A   ΔLACINDEX←ΔLACINDEX,[1](ΔSCAFF,(STARTBYTE),(STARTBYTE+ENDSCAFF))

STARTBYTE←STARTBYTE+ENDSCAFF

→RETI
*****

```

GETLACSNPS2

A THIS PROGRAM READS IN SNPS FROM DATASET OF SNP CALLS (I. LAC REFERENCE; E.G. JG1-LIKE) AND

A COMBINES WITH DATA FROM ΔLACCD5 TO FORM ΔLACSNPS, WHICH HAS FOLLOWING COLUMNS:

A
A COLUMN1 CONTIG
A COLUMN2 START POSITION OF CDS FEATURE
A COLUMN3 POSITION OF SNP
A COLUMN4 END POSITION OF CDS FEATURE
A COLUMN5 STRAND (0 = -, 1 = +)
A COLUMN6 PHASE (CODON POSITION (0,1,2) OF START POSITION OF CDS FEATURE
A COLUMNS7-9 CODONINDEX (POSITIONS OF CODON CONTAINING SNP ; ON + STRAND)
A COLUMN10 ΔSNPNUM (SNP NUMBER, AS COUNTED BY THIS PROGRAM)

ΔLACSNPS←0 10p0

ΔLACCODONS←0 19p' ' A THIS VARIABLE CONTAINS CODON INFORMATION FOR THE CORRESPONDING SHP

A COLUMN1 ALTERNATIVE ALLELES
A COLUMN2: ALTERNATIVE CODON 1
A COLUMN3: ALTERNATIVE CODON 2
A COLUMN4: SYNONYMOUS(S) OR NON-SYNONYMOUS

(N) DIFFERENCE

ΔTScores←0 62 2p' ' A THIS VARIABLE CONTAINS ALL OF THE SNP CALLS IN SAME ORDER AS ΔLACSNPS AND ΔLACCODONS

PROBLEMS←0 3p0

NOTFOUND←0p0

ΔSNPINFO←0 3p0

ΔALTALLELES←0 2p' '

ΔUNEQUAL2←0

'C:\JG1' ONTIE -1

SIZE←ONSIZE -1

FILE←ONREAD (-1,82,(SIZE+1000),0) A READ IN FILE

ONUNTIE -1

HEADER←846↑FILE A SEPARATE HEADER

DATA←846↓FILE A REMOVE HEADER FROM DATA

ASAMPLES←79↓HEADER

DATA2←DATA,'tig' A APPEND 'TIG' TO END OF DATA

A THE FOLLOWING STATEMENT IS CURRENTLY ACTIVE IN 'ANALYSNPS'

A WHEN THIS PROGRAM IS SET UP TO CALL 'ANALYSNPS', THE STATEMENT SHOULD BE

A DE-ACTIVATED IN 'ANALYSNPS' AND ACTIVATED HERE

A ΔSNPS←0 4p0 A THIS VARIABLE WILL CONTAIN INFORMATION ON SNP COD ON POSITION

A COLUMNS ARE: (1) SCAFFOLD NUMBER (2) POSITION ON SC AFFOLD (3) CODON POSITION (1,2,OR 3) (4) SYNONYMOUS (0) OR NON-SYN ONYMOUS (1)

```

I←0
RETI:I←I+1

TEST←(I÷100)=(↑(I÷100))      A REPORT SNP NUMBER EVERY 100 SNPS
↓(TEST=1)/'I'

ΔSNPNUM←I      A SNIP NUMBER

→(ΔSNPNUM>MAXSNP)/0

A PICK DATA FOR NEXT SNP--STRIP OFF LEADING INFO AND LEAVE JUST NUC
LEOTIDES FOR DIFFERENT
A SAMPLES. TEMP5 HAS FORMAT G/G G/G G/C . . .
TEMP←5000↑ DATA2
TEMP2←TEMP □SS 'tig'
TEMP3←TEMP2/↑ρTEMP2
TEMP4←(TEMP3[2]-1)↑TEMP
TEMP5←(TEMP3[1]-1)↑TEMP4
INDA←(TEMP5=□AV[10])/↑ρTEMP5
TEMP7←TEMP5[3+↑8]
ΔSCAFF←□FI TEMP7
SCORES←INDA[14]↑TEMP5      A GET CALLED BASES FOR SNP
CALLEDBASE←1↑(INDA[4]↑INDA[5]↑TEMP5)

ALLELE1←-1↑INDA[4]↑INDA[5]↑TEMP5      A CHOOSE ALLELES FROM DATA
ALLELE2←-1↑INDA[5]↑INDA[6]↑TEMP5

ALTALLELES←ALLELE1,ALLELE2

A PROCESS SCORES
IND←(SCORES='/' )/↑ρSCORES
SCORES[IND]←' '
IND←(~SCORES∈(' ',□AV[10]))/↑ρSCORES
SCORES2←SCORES[IND]
SCORES3←62 2ρSCORES2

NONAUSTSCORES←,SCORES3[ΔNONAUSTINDEX;]      A TEST FOR VARIATION BEST
DES AUSTINII
NUCS←'ACGT★'

TEST1←+/NUCS<NONAUSTSCORES
↓(TEST1=1)/' →DOWN4'      A IF NO VARIATION AFTER REMOVE AUSTINII, SK
IP SNP

TEMP8←-1↑INDA[1]↑INDA[2]↑TEMP5
ΔSCAFFPOS←□FI TEMP8      A SCAFFOLD POSITION OF SNP

A DETERMINE WHICH ELEMENT OF ΔLACCD5 IS RELEVANT
IND1←ΔLACCD5[;1]=ΔSCAFF      A BOOLEAN, 1 IF SCAFFOLD
IND2←ΔLACCD5[;2]≤ΔSCAFFPOS      A BOOLEAN, 1 IF SNP POSITION ≥ STAR
POSITION OF CDS FEATURE
IND3←ΔLACCD5[;3]≥ΔSCAFFPOS      A BOOLEAN, 1 IF SNP POSITION ≤ END
POSITION OF CDS FEATURE

```

```

IND←(IND1^IND2^IND3)/1ρIND1      A INDEX OF ENTRY THAT SATISFIES EA
CH OF ABOVE CONDITIONS
FLAG←ρIND
→(FLAG≠0)/DOWN20      A IF CDS FOUND, GO TO DOWN20
A IF NO CDS FOUND, PROCESS SNP

NEWLINE←ΔSCAFF,-1,ΔSCAFFPOS,-1,-1,-1,-1,-1,ΔSNPNUM
ΔLACSNPS←ΔLACSNPS,[1](NEWLINE)  A ADD NEW LINE TO ΔLACSNPS

CODON1←'XXX'
CODON2←'XXX'
STATUS←'O'      A THIS STATUS INDICATES NO CDS FOUND
CLINE←' ',CODON1,' ',CODON2,' ',STATUS,7 0ΔSNPNUM  A MAKE LINE
LISTING CODONS AND STATUS (SYN VS. NON-SYN)
ΔLACCODONS←ΔLACCODONS,[1]CLINE
→DOWN4      A PROCESS NEXT SNP

DOWN20:      A START PROCESSING SNPS FOR WHICH CDS IS F
FOUND

IND4←1↑IND      A IN CASE MORE THAN ONE ENTRY, PIC
K FIRST
LINE←,ΔLACCDSS[IND4;]      A PICK APPROPRIATE LINE FROM ΔLACCD
S5
STARTPOS←LINE[2]
END←LINE[3]
STRAND←LINE[4]
PHASE←LINE[5]

A NEXT, PICK OUT CODON CORRESPONDING TO SNP

READLACSCAFF ΔSCAFF      A READ IN SCAFFOLD SEQUENCE, IN VARIABLE Δ
SEQUENCE
A DIFF←ΔSCAFFPOS-STARTPOS  A DIFFERENCE BETWEEN SNP POSITION AND S
TART POSITION OF CDS FEATURE

Δ(CALLEDBASE≠ΔSEQUENCE[ΔSCAFFPOS])/'ΔUNEQUAL2←ΔUNEQUAL2+1'

→(STRAND=0)/DOWN1      A GO TO DOWN1 IF STRAND IS NEGATIVE

A CALCULATE CODON POSITIONS FOR + STRAND
STARTFRAME←STARTPOS+PHASE  A POSITION OF FIRST CODON AFTER STARTIN
G POSITION

DIFF←ΔSCAFFPOS-STARTFRAME  A NUCLEOTIDES BETWEEN STARTFRAME AND SC
AFFOLD POSITION
POSFRAME←3|DIFF      A CODON POSITION OF SNP
STCOD←ΔSCAFFPOS-POSFRAME  A START POSITION OF CODON CONTAINING TH
E SNP
CODONINDEX←STCOD,(STCOD+1),(STCOD+2)  A POSITIONS OF ALL THREE NUCS
OF CODON
CODON←ΔSEQUENCE[CODONINDEX]  A CODON EXTRACTED FROM I. LACIDA SEQU
ENCE
CODPOS←(ΔSCAFFPOS=CODONINDEX)/13  A VARIABLE POSITION IN CODON

```

```

A ADD INFORMATION TO ΔLACSNPS
NEWLINE←ΔSCAFF,STARTPOS,ΔSCAFFPOS,END,STRAND,PHASE,CODONINDEX,ΔSNPNUM
M
ΔLACSNPS←ΔLACSNPS,[1]NEWLINE

TEST←','<ALLELE2      A TEST FOR WHETHER MULTIPLE ALLELES

→(TEST=0)/DOWN2      A IF ONLY 2 ALLELES, GO TO DOWN2

A IF > 2 ALLELES, PROCESS

STATUS←'M'           A INDICATES SNP HAS > 2 ALLELES
CODON1←'XXX'
CODON2←'XXX'
STATUS←'O'           A THIS STATUS INDICATES NO CDS FOUND
CLINE←' ',CODON1,' ',CODON2,' ',STATUS,7 0ΔSNPNUM      A MAKE LINE
LISTING CODONS AND STATUS (SYN VS. NON-SYN)
ΔLACCONDONS←ΔLACCONDONS,[1]CLINE

→DOWN4              A PROCESS NEXT SNP

DOWN2:  A DETERMINE WHETHER VARIATION IS NON-SYNONYMOUS OR SYNONYMOUS
S

A MAKE ALLELES REVERSE COMPLEMENT IF STRAND IS NEGATIVE
Δ(STRAND=0)/'ALLELE1←REVCOMP ALLELE1 Δ ALLELE2←REVCOMP ALLELE2'

CODON1←CODON2←CODON      A INITIALIZE VARIABLES
CODON1[CODPOS]←ALTALLELES[1]  A INSERT ONE ALTERNATIVE ALLELE
CODON2[CODPOS]←ALTALLELES[2]  A INSERT OTHER ALTERNATIVE ALLELE

A DETERMINE WHETHER CODONS PRODUCE SAME AA
TRANSLATE CODON1
TEMP1←AA1
TRANSLATE CODON2
TEMP2←AA1
TEST←TEMP1=TEMP2
STATUS←'N'
Δ(TEST=1)/'STATUS←'S''
A 'SYNONYMOUS VS NON-SYNONYMOUS: ',STATUS
CLINE←ALTALLELES,' ',CODON1,' ',CODON2,' ',STATUS,7 0ΔSNPNUM      A
MAKE LINE LISTING CODONS AND STATUS (SYN VS. NON-SYN)
ΔLACCONDONS←ΔLACCONDONS,[1]CLINE      A ADD LINE TO ΔLACCONDONS

DOWN4:DATA2←(TEMP3[2]-1)↓DATA2      A DROP CURRENT SNP FROM DATA2

TEST←+/(DATA2 DSS 'tig')

→(TEST>1)/RETI
□TCLF
'PROGRAM COMPLETE.  SNP DATA IN VARIABLES ''ΔLACSNPS'' AND ''ΔLACCONDONS''.'

```

```
*****
```

SEPSNPS

A THIS FUNCTION SEPARATES SNPS INTO SYNONYMOUS, NON-SYNONYMOUS AND OTHER

```
STATUS←ΔLACCODONS[;12]      A COLUMN OF N'S, S'S, O'S
SNPNUMS←ΔLACCODONS[;12+17]  A SNP NUMBERS FROM ΔLACCODONS IN TEXT
FORMAT
BLANKS←((ρSTATUS)ρ' ')      A ADD ONE COLUMN OF SPACES
SNPNUMS←,(SNPNUMS,BLANKS)
SNPNUMS2←□FI SNPNUMS        A CHANGE SNP NUMBERS TO NUMERIC
```

A GET SYN SNPS

```
IND1←(STATUS='S')/1ρSTATUS  A INDEX OF WHICH ROWS OF ΔLACODONS COR
RESPOND TO SYNONYMOUS SNPS
```

```
SYNSNPNUMS←SNPNUMS2[IND1]   A PICK SNP NUMBERS CORRESPONDING TO SY
NONYMOUS SITES
```

```
DIMS←ρΔSNPS1
ΔSYNSCORES←0 62 2ρ0
I←0
RETI:I←I+1
```

```
TEST←(I÷1000)=(⌈(I÷1000))
⌊(TEST)/'I'
TEST←ΔSNPS1[I]∈SYNSNPNUMS
⌊(TEST)/'ΔSYNSCORES←ΔSYNSCORES,[1]ΔSCORES[I;;]'
→(I<DIMS)/RETI
```

A GET NON-SYN SNPS
'PROCESSING NON-SYN SNPS'

```
IND2←(STATUS='N')/1ρSTATUS
NONSNPNUMS←SNPNUMS2[IND2]
```

```
DIMN←ρΔSNPS1
ΔNONSCORES←0 62 2ρ0
J←0
RETJ:J←J+1
```

```
TEST←(J÷1000)=(⌈(J÷1000))
⌊(TEST)/'J'
TEST←ΔSNPS1[J]∈NONSNPNUMS
⌊(TEST)/'ΔNONSCORES←ΔNONSCORES,[1]ΔSCORES[J;;]'
→(J<DIMN)/RETJ
```

DOWN:'PROCESSING O SNPS'
IND3←(STATUS='O')/1ρSTATUS

```
OSNPNUMS←SNPNUMS2[IND3]
```

```
DIMO←ρΔSNPS1
ΔOSCORES←0 62 2ρ0
K←0
```

RETK:K←K+1

TEST←(K÷1000)=(r(K÷1000))
⊥(TEST)/'K'
TEST←ΔSNPS1[K]∈OSNPNUMS
⊥(TEST)/'ΔOSCORES←ΔOSCORES,[1]ΔSCORES[K;;]'
→(K<DIMO)/RETK

SPLICE4

⌘ THIS FUNCTION CALCULATES THE NUMBERS OF SYNONYMOUS AND NON-SYNONYM
OUS SITES INT THE TRANSCRIPTS
⌘ BASED ON MATCHES TO I. LACUNOSA CODING REGIONS

ΔTOTSYN←ΔTOTNONSYN←0
DIMTRAN←1↑ρΔLACTRANDATA

LACSCAFFS←ΔLACINDEX[;1]

COUNTER←1 ⌘ THIS IS A COUNTER FOR TRANSCRIPT NUMBER

ΔLACTRANIND←0ρ' ' ⌘ THIS IS THE VARIABLE THAT WILL HOLD SUCCES
SIVE SPLICED TRANSCRIPT INDICIES

⌘ FORMAT IS > TRANSCRIPTNUMBER SCAFFNUM POSI
TIONINDEX ⌘TCLF
⌘ TRANSCRIPTNUMBER IS NUMBER OF TRANSCRIPT (I
.E. 'COUNTER' IN THIS PROGRAM)
⌘ SCAFFNUM IS 8 0# SCAFFOLD NUMBER
⌘ POSITIONINDEX IS VECTOR OF POSITIONS OF
SPLICED TRANSCRIPT IN 10 0# FORMAT

ΔLACTRANINDDATA←0 8ρ0 ⌘ THIS VARIABLE HOLDS START AND END POSITION
S OF THE TRANSCRIPT
⌘ FORMAT FOR EACH ROW IS TRANSCRIPTNUMBER SC
AFFNUM, STRAND, STARTPOS, ENDPOS,
⌘ START POSITION OF TRANSCRIPT DATA IN ΔL
ACTRANIND,END POSITION OF TRANSCRIPT DATA IN ΔLACTRANIND

LACTRANLENGTH←0

I←0
RETI:I←I+1
TEST←(I÷100)=(r(I÷100))
⊥(TEST=1)/'I'

LINE1←ΔLACTRANDATA[I;] ⌘ READ IN DATA FOR TRANSCRIPT I
SCAFF←LINE1[2] ⌘ LACUNOSA SCAFFOLD CORRESPONDING TO TRANSCR
IPT

⊥(~SCAFF∈LACSCAFFS)/'⌘TCLF ⌘ 'NO CONTIG FOUND' ⌘ TRSH←⌘ ⌘ →DOWN'

READLACSCAFF SCAFF

```

A PICK OUT I. LACUNOSA CDS FEATURES CONTAINED IN THE TRANSCRIPT
IND1←ΔLACCD5[;1]=SCAFF
IND2←ΔLACCD5[;2]≥LINE1[3]      A TEST WHETHER CDS FEATURE CONTAINED I
N TRANSCRIPT
IND3←ΔLACCD5[;3]≤LINE1[4]      A DITTO
IND←IND1∧IND2∧IND3

↓(0=+/IND)/' →DOWN'          A SKIP IF TRANSCRIPT CONTAINS NO
I. LACCD5 FEATURES

IND←IND/↓ρIND

PARTCDS4←ΔLACCD5[IND;]        A PICK I. LAC CDS FEATURES CONTAINED I
N TRANSCRIPT
DIMPART←1↑ρPARTCDS4
IND1←(PARTCDS4[;4]=1)/↓DIMPART  A INDEX FOR CDS ON POSITIVE STRAND
IND2←(PARTCDS4[;4]=0)/↓DIMPART  A INDEX FOR CDS ON NEGATIVE STRAND

PARTPOS←PARTCDS4[IND1;]        A PICK CDS ON POSITIVE STRAND

→(0=1↑ρPARTPOS)/DOWN

IND←ΔPARTPOS[;2]
PARTPOS←PARTPOS[IND;]          A SORT CDS ON POSITIVE STRAND IN O
RDER OF INCREASING BEGINNING POSITION
PARTNEG←PARTCDS4[IND2;]        A PICK CDS ON NEGATIVE STRAND
IND←ΔPARTNEG[;2]
PARTNEG←PARTNEG[IND;]          A SORT CDS ON NEGATIVE STRAND IN O
RDER OF DECREASING BEGINNING POSITION

BEGPHASE←PARTPOS[1;5]

A SPLICE TRANSCRIPT FOR CDS ON POSITIVE STRAND

FLAG←0                          A FLAG= 0 INDICATES THE NEXT CDS IS THE F
IRST CDS IN A GENE
SPLICEDTRAN←0ρ' '
DIM1←1↑ρPARTPOS
→(DIM1=0)/DOWN10

A START POSITIVE STRAND LOOP

CDSIND2←CDSIND←0ρ0

K←0
RETK:K←K+1
LINE2←PARTPOS[K;]              A LAC CDS DATA
PHASE←LINE2[5]

STRAND←LINE2[4]

CDSSTARTPOS←LINE2[2]
CDSENDPOS←LINE2[3]

```



```

CDSIND←CDSIND,CDSSTARTPOS,CSENDPOS
CDSIND2←CDSIND2,((CDSSTARTPOS-1)+1(1+(CSENDPOS-CDSSTARTPOS)))

```

```

DOWN1:→(K<DIM1)/RETK

```

```

BEGPHASE←PARTPOS[1;5]
A FIRSTPOS←CDSIND[1]
  ±(BEGPHASE=0)/'CDSIND2+2↓CDSIND2'
  ±(BEGPHASE=2)/'CDSIND2+1↓CDSIND2'

```

```

STRAND←1
LINE5←'> ',(8 0↯COUNTER),(8 0↯SCAFF),(10 0↯CDSIND),□TCLF

```

```

ΔLACTRANIND←ΔLACTRANIND,LINE5

```

```

LACTRANSTART←LACTRANLENGTH+1
LACTRANLENGTH←LACTRANEND←LACTRANLENGTH+ρLINE5

```

```

LINE6←COUNTER,SCAFF,STRAND,BEGPHASE,(1↑CDSIND),(-1↑CDSIND),LACTRA
NSTART,LACTRANEND
ΔLACTRANINDDATA←ΔLACTRANINDDATA,[1]LINE6
COUNTER←COUNTER+1

```

```

COUNTSYN SEQ      A CALL SCRIPT COUNTSYN
ΔTOTSYN←ΔTOTSYN+SYN
ΔTOTNONSYN←ΔTOTNONSYN+NONSYN

```

```

→(0=1↑ρPARTNEG)/DOWN

```

```

DOWN10:

```

```

A SPLICED TRANSCRIPT FOR CDS ON NEGATIVE STRAND
  PH

```

```

FLAG←0      A FLAG= 0 INDICATES THE NEXT CDS IS THE F
IRST CDS IN A GENE
SPlicedTRAN←0ρ' '
DIM1←1↑ρPARTNEG
→(DIM1=0)/DOWN
CDSIND2←CDSIND←0ρ0
J←0
RETJ:J←J+1
LINE2←PARTNEG[J;]
PHASE←LINE2[5]

```

```

CDSSTARTPOS←LINE2[2]
CSENDPOS←LINE2[3]
INDB←(CDSSTARTPOS-1)+1(1+(CSENDPOS-CDSSTARTPOS))
CDSIND←CDSIND,CDSSTARTPOS,CSENDPOS
CDSIND2←CDSIND2,((CDSSTARTPOS-1)+1(1+(CSENDPOS-CDSSTARTPOS)))

```

```

DOWN2:→(J<DIM1)/RETJ

BEGPHASE←PARTNEG[1;5]

SPLICETRAN2←ΔSEQUENCE[CDSIND2]
SPLICETRAN←REVCOMP SPLICETRAN2
TRANSLATE SPLICETRAN

COUNT1←+/AA1='★'
COUNT2←+/AA2='★'
COUNT3←+/AA3='★'

COUNTS←COUNT1,COUNT2,COUNT3
MIN←L/COUNTS
INDC←(COUNTS=MIN)/13
Δ((ρINDC)>1)/'→DOWN'

CDSIND2←(-1×(INDC-1))+CDSIND2

STRAND←0
LINE5←'> ',(8 0#COUNTER),(8 0#SCAFF),(10 0#CDSIND),DTCLF
ΔLACTRANIND←ΔLACTRANIND,LINE5

LACTRANSTART←LACTRANLENGTH+1
LACTRANLENGTH←LACTRANEND←LACTRANLENGTH+ρLINE5

LINE6←COUNTER,SCAFF,STRAND,BEGPHASE,(1+CDSIND),(-1+CDSIND),LACTRANST
ART,LACTRANEND

ΔLACTRANINDDATA←ΔLACTRANINDDATA,[1]LINE6
COUNTER←COUNTER+1

SEQ←ΔSEQUENCE[CDSIND2]
SEQ2←REVCOMP SEQ
COUNTSYN SEQ2
ΔTOTSYN←ΔTOTSYN+SYN
ΔTOTNONSYN←ΔTOTNONSYN+NONSYN

DOWN:→(I<DIMTRAN)/RETI

*****
COUNTSYN X
# THIS PROGRAM CALCULATES THE NUMBER OF SYNONYMOUS AND NON-SYNONYMOU
S SITES IN A SEQUENCE X

SEQ←X
LENGTH←ρX
MOD←3|LENGTH
SEQ←(-1×MOD)+SEQ

# DROP LAST CODON IF IT IS STOP CODON
ENDCOD←-3+SEQ

```

```

TRANSLATE ENDCOD
 $\downarrow$  (AA1='*')/'SEQ $\leftarrow$ -3 $\downarrow$ SEQ'

```

```

LENGTH2 $\leftarrow$ ( $\rho$ SEQ) $\div$ 3
MAT $\leftarrow$ (LENGTH2,3) $\rho$ SEQ
MAT $\leftarrow$  $\mathbb{Q}$ MAT

```

```

TEMP $\leftarrow$ +/( $\Delta$ CODONS $\wedge$ .=MAT)
TEMP $\leftarrow$ (( $\rho$ TEMP),1) $\rho$ TEMP

```

```

TEMP2 $\leftarrow$ TEMP,TEMP,TEMP

```

```

TEMP3 $\leftarrow$  $\Delta$ DEGENMATRIX $\times$ TEMP2

```

```

SYN $\leftarrow$ +//TEMP3
TOT $\leftarrow$ LENGTH2 $\times$ 3
NONSYN $\leftarrow$ TOT-SYN
*****

```

Several of the above scripts call the script TRANSLATE, which translate nucleotide sequences into amino-acid sequences:

```

*****

TRANSLATE X;MAX;CODONS;AA
MAX $\leftarrow$ L( $\rho$ X) $\div$ 3
CODONS $\leftarrow$  $\mathbb{Q}$ (MAX,3) $\rho$ X
FIXFN
AA $\leftarrow$ , (1 65 $\rho$ 165)+. $\times$ ( $\Delta$ CODE $\wedge$ .=CODONS)
AA1 $\leftarrow$ , $\Delta$ SYMB[AA;]
X $\leftarrow$ 1 $\downarrow$ X
MAX $\leftarrow$ L( $\rho$ X) $\div$ 3
CODONS $\leftarrow$  $\mathbb{Q}$ ((MAX,3) $\rho$ X)
FIXFN
AA $\leftarrow$ , (1 65 $\rho$ 165)+. $\times$ ( $\Delta$ CODE $\wedge$ .=CODONS)
AA2 $\leftarrow$ , $\Delta$ SYMB[AA;]
X $\leftarrow$ 1 $\downarrow$ X
MAX $\leftarrow$ L( $\rho$ X) $\div$ 3
CODONS $\leftarrow$  $\mathbb{Q}$ ((MAX,3) $\rho$ X)
FIXFN
AA $\leftarrow$ , (1 65 $\rho$ 165)+. $\times$ ( $\Delta$ CODE $\wedge$ .=CODONS)
AA3 $\leftarrow$ , $\Delta$ SYMB[AA;]
*****

```

The script TRANSLATE requires the variables ΔCODE (a 65 x 3 character matrix) and ΔSYMB (a 65 x 1 character vector):

Transpose of ΔCODE =

TTTTTTTTTTTTTTTCCCCCCCCCCCCCAAAAAAAAAAAAAAGGGGGGGGGGGGGG-
TTTTCCCAAAAGGGGTTTTCCCCAAAAGGGGTTTTCCCCAAAAGGGGTTCCTCCCAAAGGGG-
TCAGTCAGTCAGTCAGTCAGTCAGTCAGTCAGTCAGTCAGTCAGTCAGTCAGTCAGTCAG-

Transpose of $\Delta \text{SYMB} =$

FLLSSSSYY**CC*WLLLPPPPHHQRRRRIIIMTTTNNKKSSRRVVVVAAAADDEEGGGG-

II. Indexing the *I. lacunosa* genome. The script INDEXLAC indexes contigs in the LAC genome FASTA file. It produces the variable Δ LACINDEX, which is an N X 3 matrix. Each row of the matrix consists of the following elements:

1. scaffold number
2. start position of scaffold sequence in FASTA file
3. end position of scaffold sequence in FASTA file.

Δ LACINDEX is used primarily for retrieving contig sequences using the script READLACSCAFF (also below)

```

*****
INDEXLAC ;STARTBYTE
@ THIS PROGRAM INDEXES THE SCAFFOLDS OF THE I. LACUNOSA GENOME. IT
CREATES THE VARIABLE
@    $\Delta$ LACINDEX, WHICH HAS THE FOLLOWING COLUMNS:
@       COLUMN1      SCAFFOLD NUMBER
@       COLUMN2      START POSITION OF SCAFFOLD IN GENOME FASTA FILE
@       COLUMN3      END POSITION OF SCAFFOLD IN GENOME FASTA FILE

 $\Delta$ LACINDEX $\leftarrow$ 0 300

STARTBYTE $\leftarrow$ 0
I $\leftarrow$ 0
RETI:I $\leftarrow$ I+1

TEST $\leftarrow$ (I $\div$ 100)=( $\uparrow$ (I $\div$ 100))
 $\uparrow$ (TEST=1)/'I'

'C:\LACGENOME' ONTIE -1 @ OPEN LAC GENOME FASTA FILE

PARTSCAFF $\leftarrow$ ONREAD -1 82 5000000 STARTBYTE
ONUNTIE -1

DIM $\leftarrow$ 0PARTSCAFF
 $\rightarrow$ (DIM=0)/0

IND $\leftarrow$ (PARTSCAFF = '>')/ $\uparrow$ 0PARTSCAFF
ENDSCAFF $\leftarrow$ IND[2]-1

TEMPSCAFF $\leftarrow$ ENDSCAFF $\uparrow$ PARTSCAFF
IND $\leftarrow$ (TEMPSCAFF=0TCLF)/ $\uparrow$ 0TEMPSCAFF
HEAD $\leftarrow$ IND[1] $\uparrow$ TEMPSCAFF
SEQ $\leftarrow$ IND[1] $\uparrow$ TEMPSCAFF
SEQ $\leftarrow$ (SEQ $\neq$ 0TCLF)/SEQ
 $\Delta$ SCAFF $\leftarrow$ 0FI HEAD[17+15]

 $\Delta$ LACINDEX $\leftarrow$  $\Delta$ LACINDEX,[1]( $\Delta$ SCAFF,(STARTBYTE),(STARTBYTE+ENDSCAFF))
STARTBYTE $\leftarrow$ STARTBYTE+ENDSCAFF

 $\rightarrow$ RETI
*****

```

READLACSCAFF SCAFFN

A BEFORE RUNNING THIS, MAKE SURE TO RUN 'CONVLACINDEX' TO CONVERT ΔL
ACINDEX FROM A CHARACTER
A TO A NUMERIC MATRIX

A THIS PROGRAM READS IN I. LAC SCAFFOLD FROM FILE C:\LACGENOME
A SCAFFN IS THE SCAFFOLD NUMBER TO READ

```
SCAFFNUMS←ΔLACINDEX[;1]
IND←(SCAFFN=SCAFFNUMS)/1ρSCAFFNUMS
TEMP1←,ΔLACINDEX[IND;]
STARTBYTE←TEMP1[2]
ENDBYTE←TEMP1[3]

'C:\LACGENOME' ONTIE -1
SCAFFOLD←ONREAD -1 82 (ENDBYTE-STARTBYTE) STARTBYTE
ONUNTIE -1
SCAFFNAME←24↑SCAFFOLD
ΔSEQUENCE←24↓SCAFFOLD
ΔSEQUENCE←(ΔSEQUENCE≠0TCLF)/ΔSEQUENCE
*****
```

READLACSCAFF requires running 'CONVLACINDEX' before running:

CONVLACINDEX

A THIS PROGRAM CONVERTS ΔLACINDEX, A CHARACTER MATRIX, TO ΔLACINDEX,
A A NUMERIC MATRIX

```
DIM←1↑ρΔLACINDEX
TEMP1←ΔLACINDEX,(DIM,3)ρ' '
TEMP2←,TEMP1
IND←(TEMP2=0AV[10])/1ρTEMP2
TEMP2[IND]←' '
TEMP3←0FI TEMP2
ΔLACINDEX←(DIM,3)ρTEMP3
*****
```

III. Calculating and bootstrapping π values

A. Calculate average pairwise π values for all samples; as listed below this is done for all SNPs, but the script can be modified to do just synonymous, just non-synonymous, or just non-coding SNPs.

```
PICALC2 X;I
A THIS PROGRAM CALCULATES PAIRWISE PI VALUES FOR EITHER THE FULL DAT
ASET (ALL SNPS)
A OR FOR SUBSETS OF DATA (E.G. SYNONYMOUS, NON-SYNONYMOUS SITES)
A IT PRODUCES A MATRIX  $\Delta\pi_2$  THAT HAS THE AVERAGE PAIRWISE PI VALUES
FOR EACH PAIR OF SAMPLES
A PROGRAM READS IN VARIABLE  $\Delta$ SCORES CREATED BY READCALLS AS X
A (CAN ALSO READ IN  $\Delta$ SYNSCORES OR  $\Delta$ NONSCORES OR  $\Delta$ RSCORES)

A DIVIDER $\leftarrow$  $\Delta$ TOTSYN A CHANGE THIS IF CALCULATING PI FOR NON-SYN SIT
ES
A DIVIDER $\leftarrow$  $\Delta$ TOTNONSYN A FOR USE WITH  $\Delta$ RSCORES
DIVIDER $\leftarrow$ 30036768 A TOTAL NUMBER OF SITES IN TRANSCRIPTOME

DIM $\leftarrow$ 1 $\uparrow$  $\rho$ X

PICUM $\leftarrow$ 62 62 $\rho$ 0 A SET UP PI MATRIX- WILL EVENTUALLY HAVE NUMBER OF P
AIRWISE DIFFERENCES ACROSS ALL VARIABLE SNPS
PICOUNT $\leftarrow$ 62 62 $\rho$ 0 A SET UP MATRIX FOR CUMULATIVE COUNTS (EXCLUDES MI
SSING VALUES)

I $\leftarrow$ 0
RETI:I $\leftarrow$ I+1 A SNP LOOP

TEMP6 $\leftarrow$ ((I $\div$ 1000)-(I $\div$ 1000))) A REPORT SNP NUMBER EVERY 100 SNPS
 $\&$ (TEMP6=0)/'DTCLF  $\diamond$  I'

A PICK DATA FOR NEXT SNP--STRIP OFF LEADING INFO AND LEAVE JUST NUC
LEOTIDES FOR DIFFERENT
A SAMPLES. TEMP5 HAS FORMAT G/G G/G G/C . . .

SCORES3 $\leftarrow$ X[I;;] A NOTE: MODIFY WHICH  $\Delta$ SCORES TO USE (E.G.
ALL, SYNONYMOUS, ETC.)

A CALCULATE PAIRWISE PI VALUES FUR CURRENT SNP
TEMP7 $\leftarrow$ SCORES3 $\circ$ .. $\circ$ SCORES3
TEMP8 $\leftarrow$ +/[2]TEMP7
TEMP9 $\leftarrow$ (+/[3]TEMP8) $\div$ 4
TEMP9 $\leftarrow$ 1-TEMP9
TIND $\leftarrow$ (SCORES3[;1]='. ')/162
COUNTMAT $\leftarrow$ 62 62 $\rho$ 1
COUNTMAT[TIND;] $\leftarrow$ 0
COUNTMAT[;TIND] $\leftarrow$ 0
PI $\leftarrow$ TEMP9 $\times$ COUNTMAT
A ADD CURRENT PAIRWISE PI VALUES TO CUMULATIVE VALUES
PICUM $\leftarrow$ PICUM+PI
```

```
PICOUNT←PICOUNT+COUNTMAT
```

```
→(I<DIM)/RETI
```

```
ΔPI2←PICUM÷DIVIDER
```

```
□TCLF
```

```
'PROGRAM PICALC2 FINISHED.'
```

```
'PAIRWISE PI VALUES IN VARIABLE ''ΔPI2''.'
```

```
*****
```

B. Calculate pairwise π values between and within species. This script uses ΔPI2 from PICALC2

```
*****
```

```
PIANAL2
```

```
  A THIS PROGRAM CALCULATES AVERAGE PAIRWISE DIFFERENCES (PI) FOR WITH  
  IN EACH SPECIES AND BETWEEN SPECIES
```

```
  A IT USES THE VARIABLE 'ΔPI2', WHICH IS PICUM÷ΔTOTNUC, WHERE PICUM I  
  S INTERMEDIATE
```

```
  A MATRIX FROM 'PICALC2', AND ΔTOTNUC IS TOTAL NUMBER OF SITES IN T  
  RANSCRPTOMES, AS CALCULATED BY 'SPLICE'
```

```
CINDEX←ΔCINDEX      A INDEX FOR WHICH CORDATATRILOBA SAMPLES BEING US  
ED
```

```
LINDEX←ΔLINDEX      A INDEX FOR WHICH LACUNOSA SAMPLES BEING USED  
                    A THESE INDICES CAN BE CHANGED TO LOOK AT ONLY AL  
LOPATRIC OR ONLY SYMPATRIC SAMPLES
```

```
  A MAKE MASKS FOR PUTTING 0 ON DIAGONALS FOR WITHIN SPECIES SAMPLES
```

```
TMP←ρCINDEX
```

```
MASKC←(TMP,TMP)ρ1
```

```
I←0
```

```
RETI:I←I+1
```

```
MASKC[I;I]←0
```

```
→(I<TMP)/RETI
```

```
TMP←ρLINDEX
```

```
MASKL←(TMP,TMP)ρ1
```

```
J←0
```

```
RETJ:J←J+1
```

```
MASKL[J;J]←0
```

```
→(J<TMP)/ρRETJ
```

```
  A CALCULATE WITHIN-CORDAT AVERAGE PI
```

```
CW1←ΔPI2[CINDEX;CINDEX]      A CHOOSE SUBSET OF ΔPI2 CORRESPONDING TO  
CORDAT
```

```
CW2←CW1×MASKC      A MAKE DIAGONAL ELEMENTS 0
```

```
DIM1←ρCINDEX      A NUMBER OF CORDAT SAMPLES
```

```
CW3←(+/+/CW2)÷(DIM1×(DIM1-1))      A AVERAGE PI FOR WITHIN CORDAT
```



```

A CALCULATE WITHIN-LAC AVERAGE PI
LW1←ΔPI2[LINDEX;LINDEX]
LW2←LW1×MASKL
DIM2←ρLINDEX
LW3←(+ / + / LW2) ÷ (DIM2 × (DIM2 - 1))

A CALCULATE BETWEEN-SPECIES AVERAGE PI
B←ΔPI2[CINDEX;LINDEX]
B2←(+ / + / B) ÷ (DIM1 × DIM2)

A COMMENT NEXT STATEMENTS IF RUNNING PIBOOT4

```

```

□TCLF
'AVERAGE PI WITHIN CORDAT: ',CW3
'AVERAGE PI WITHIN LAC: ',LW3
'AVERAGE PI BETWEEN SPEC: ',B2
*****

```

C. Bootstrap π values from PIANAL2 and test for species differences. By modifying indices used in third and fourth line, can test for differences using different sample sets, e.g. all CORDAT vs. all LAC, allopatric CORDAT vs. allopatric LAC, etc.

```

*****
PIBOOT4
A THIS PROGRAM BOOTSTRAPS DIFFERENCES IN PI THE TWO SPECIES

CINDEX←ΔCINDEX      A INDEX FOR WHICH CORDAT SAMPLES TO USE
LINDEX←ΔLINDEX      A INDEX FOR WHICH LACUNOSA SAMPLES TO USE

BOOTDATA2←0 2ρ0      A VARIABLE TO HOLD BOOTSTRAP DATA

MASK4←62 62ρ1

DIM←1↑ρΔSCORES
II←0
RETII:II←II+1      A BOOTSTRAP LOOP
TEST←(II÷1000)=(⌈(II÷1000))
Δ(TEST) / 'II'      A PRINT II EVERY 1000 SNPS
II
A CREATE BOOTSTRAP DATASET
A FIRST BOOTSTRAP SAMPLES
SCORES←ΔSCORES      A COPY SNP GENOTYPES INTO VARIABLE 'SCORES'
CSCORES←SCORES[;CINDEX;]      A SNP GENOTYPES FOR CORDAT
LScores←SCORES[;LINDEX;]      A SNP GENOTYPES FOR LAC

DIMC←ρCINDEX
DIML←ρLINDEX
A BOOTSTRAP CORDAT AND LAC SAMPLES
RANDC←?DIMCρDIMC
SCORES[;CINDEX;]←SCORES[;CINDEX[RANDC];]

```

```

RANDL←?DIMLρDIML
SCORES[;LINDEX;]←SCORES[;LINDEX[RANDL];]

# NOW BOOTSTRAP SNPS
DIMSnpS←1↑ρSCORES
RANDSnpS←?DIMSnpSρDIMSnpS
SCORES←SCORES[RANDSnpS;;]

PICALC2 SCORES          # CALL PICALC2 AND PIANAL2 TO CALCULATE
PI VALUES
PIANAL2

BOOTDATA2←BOOTDATA2,[1](CW3,LW3)  # APPEND BOOTSTRAP DATA FOR PI FOR
CORDAT AND LAC SAMPLES

→(II<1000)/RETI

# CALCULATE STATISTICS
DIFF←BOOTDATA2[;1]-BOOTDATA2[;2]
PROP←+/DIFF≤0
ΠTCLF
'PROPORTION OF 1000 BOOTSTRAP SAMPLE DIFFS ≤0: ',PROP

CVALS←,BOOTDATA2[;1]
CVALS2←CVALS[ΔCVALS]
LVALS←,BOOTDATA2[;2]
LVALS2←LVALS[ΔLVALS]
CONF C←CVALS2[25,975]
CONFL←LVALS2[25,975]
ΠTCLF
'CONF INTERVAL FOR I. CORD: ',CONF C
'CONF INTERVAL FOR I. LAC: ',CONFL
ΠTCLF
'PROGRAM PIBOOT4 COMPLETED. DATA IN VARIABLE BOOTDATA2'
*****

```

D. Calculate average between-species π values for sympatric vs. allopatric comparison.

```

*****
PICONTRAST2
# THIS PROGRAM CALCULATES PI WITHIN AND BETWEEN SYMP AND ALLOPATRIC
SAMPLES
# IT USES THE ΔPI2 MATRIX CALCULATED BY 'PICALC2'

MASK4←62 62ρ1

I←0
RETI:I←I+1
MASK4[I;I]←0
→(I<62)/RETI

PI←MASK4×ΔPI2

```

```

PIBETWALLO←PI[ΔCALLO3;ΔLALLO3]
PIBETWALLOCLOSE←PI[ΔCALLOCLOSE;ΔLALLOCLOSE]
PIBETWSYMP←PI[ΔCSYMPINDEX;ΔLSYMPINDEX]
PICBETWALLOSYMP←PI[ΔCALLO3;ΔCSYMPINDEX]
PILBETWALLOSYMP←PI[ΔLALLO3;ΔLSYMPINDEX]
PICBETWALLOCLOSESYMP←PI[ΔCALLOCLOSE;ΔCSYMPINDEX]
PILBETWALLOCLOSESYMP←PI[ΔLALLOCLOSE;ΔLSYMPINDEX]

AVEBETWSYMP←(+/+PIBETWSYMP)÷((ρΔCSYMPINDEX)×((ρΔLSYMPINDEX)))
AVEBETWALLO←(+/+PIBETWALLO)÷((ρΔCALLO3)×((ρΔLALLO3)))
AVEBETWALLOCLOSE←(+/+PIBETWALLOCLOSE)÷((ρΔCALLOCLOSE)×((ρΔLALLOCLOSE)))
AVECBETWALLOSYMP←(+/+PICBETWALLOSYMP)÷((ρΔCALLO3)×(ρΔCSYMPINDEX))
AVELBETWALLOSYMP←(+/+PILBETWALLOSYMP)÷((ρΔLALLO3)×(ρΔLSYMPINDEX))

AVECBETWALLOCLOSESYMP←(+/+PICBETWALLOCLOSESYMP)÷((ρΔCALLOCLOSE)×(ρΔCSYMPINDEX))
AVELBETWALLOCLOSESYMP←(+/+PILBETWALLOCLOSESYMP)÷((ρΔLALLOCLOSE)×(ρΔLSYMPINDEX))

ΔBETWDIFF1←AVEBETWALLO-AVEBETWSYMP
ΔBETWDIFF2←AVEBETWALLOCLOSE-AVEBETWSYMP
ΔBETWDIFF3←AVECBETWALLOSYMP-AVELBETWALLOSYMP
ΔBETWDIFF4←AVECBETWALLOCLOSESYMP-AVELBETWALLOCLOSESYMP

ΔOBSPIVECTOR←AVEBETWALLO,AVEBETWALLOCLOSE,AVEBETWSYMP,ΔBETWDIFF1,ΔBETWDIFF2,AVECBETWALLOSYMP,AVELBETWALLOSYMP,AVECBETWALLOCLOSESYMP,AVELBETWALLOCLOSESYMP,ΔBETWDIFF3,ΔBETWDIFF4  a  OBSERVED VALUES

□TCLF
'BETWEEN SPECIES COMPARISONS'
□TCLF
'a. AVERAGE PI BETWEEN KNOWN ALLO: ',AVEBETWALLO
'b. AVERAGE PI BETWEEN CLOSE ALLO: ',AVEBETWALLOCLOSE
'c. AVERAGE PI BETWEEN SYMP: ',AVEBETWSYMP
'd. DIFFERENCE a. - c.: ',(ΔBETWDIFF1)
'e. DIFFERENCE b. - c.: ',(ΔBETWDIFF2)
□TCLF
'WITHIN SPECIES COMPARISONS'
□TCLF
'f. I. CORD BETW ALLO(KNOWN) AND SYMP ',(AVECBETWALLOSYMP)
'g. I. LAC BETW ALLO(KNOWN) AND SYMP ',(AVELBETWALLOSYMP)
'h. I. CORD BETW ALLO(CLOSE) AND SYMP ',(AVECBETWALLOCLOSESYMP)
'i. I. LAC BETW ALLO(CLOSE) AND SYMP ',(AVELBETWALLOCLOSESYMP)

'j. DIFFERENCE f. - g.: ',(ΔBETWDIFF3)
'k. DIFFERENCE h. - i.: ',(ΔBETWDIFF4)

□TCLF
'PROGRAM PICONTRAST2 FINISHED'
*****

```

E. Bootstrap within and between species differences. There are two scripts. PIBOOT is the main script, which makes bootstrap samples then calls script PICONTRAST3.

PIBOOT;K

A THIS PROGRAM BOOTSTRAPS DIFFERENCES IN P BETWEEN ALLOPATRIC AND SYMPATRIC POPULATIONS OF THE TWO SPECIES

COMBCINDEX←ΔCALLO3,ΔCSYMPINDEX A COMBINED INDEX OF CORDAT KNOWN ALLOPATRIC AND SYMPATRIC SAMPLES

COMBLINDEX←ΔLALLO3,ΔLSYMPINDEX A COMBINED INDEX OF LAC KNOWN ALLOPATRIC AND SYMPATRIC SAMPLES

COMBCINDEX2←ΔCALLOCLOSE,ΔCSYMPINDEX A COMBINED INDEX OF CORDAT CLOSE ALLOPATRIC AND SYMPATRIC SAMPLES

COMBLINDEX2←ΔLALLOCLOSE,ΔLSYMPINDEX A COMBINED INDEX OF LAC CLOSE ALLOPATRIC AND SYMPATRIC SAMPLES

CONTRDATA←0 11ρ0 A INITIALIZE MARIX TO HOLD OUTPUT DATA

MASK4←62 62ρ1

I←0

RETI:I←I+1

MASK4[I;I]←0

→(I<62)/RETI

J←0

RETJ:J←J+1 A BOOTSTRAP SAMPLE LOOP

TEMP←(J÷100)=(r(J÷100))

↓(TEMP)/'J'

PITEMP←PITEMP2←ΔPI2

A MAKE BOOTSTRAP SAMPLE FOR KNOWN ALLOW AND SYMP SAMPLES

CIND←?(ρCOMBCINDEX)ρ(ρCOMBCINDEX)

LIND←?(ρCOMBLINDEX)ρ(ρCOMBLINDEX)

CIND2←COMBCINDEX[CIND]

LIND2←COMBLINDEX[LIND]

PITEMP[COMBCINDEX;COMBCINDEX]←PITEMP[CIND2;CIND2]

PITEMP[COMBLINDEX;COMBLINDEX]←PITEMP[LIND2;LIND2]

PITEMP[COMBCINDEX;COMBLINDEX]←PITEMP[CIND2;LIND2]

PITEMP[COMBLINDEX;COMBCINDEX]←PITEMP[LIND2;CIND2]

A MAKE BOOTSTRAP SAMPLE FOR CLOSE ALLOW AND SYMP SAMPLES

CIND←?(ρCOMBCINDEX2)ρ(ρCOMBCINDEX2)

LIND←?(ρCOMBLINDEX2)ρ(ρCOMBLINDEX2)

CIND2←COMBCINDEX2[CIND]

LIND2←COMBLINDEX2[LIND]

```

PITEMP2[COMBCINDEX2;COMBCINDEX2]←PITEMP2[CIND2;CIND2]
PITEMP2[COMBLINDEX2;COMBLINDEX2]←PITEMP2[LIND2;LIND2]
PITEMP2[COMBCINDEX2;COMBLINDEX2]←PITEMP2[CIND2;LIND2]
PITEMP2[COMBLINDEX2;COMBCINDEX2]←PITEMP2[LIND2;CIND2]

```

```

PICONTRAST3      @ CALL SCRIPT PICONTRAST3

```

```

→(J<1000)/RETJ

```

```

@ CALCULATE CONFIDENCE SETS

```

```

□TCLF

```

```

'PROPORTION OF 1000 BOOTSTRAP SAMPLES ≥ OBSERVED PI CONTRASTS'
'  FIRST NUMBER IS OBSERVED, SECOND IS  PROPORTION'

```

```

□TCLF

```

```

'a. AVERAGE PI BETWEEN LAC AND CORD ALLOKNOWN POPS: ',ΔOBSPIVECTOR[1
]
'b. AVERAGE PI BETWEEN LAC AND CORD ALLOCLOSE POPS: ',ΔOBSPIVECTOR[2
]
'c. AVERAGE PI BETWEEN LAC AND CORD SYMP      POPS: ',ΔOBSPIVECTOR[3
]
'd. DIFFERENCE a. - c.                                ',ΔOBSPIVECTOR[4
], ' ',+/(CONTRDATA[;4]≥ΔOBSPIVECTOR[4])÷1000
'e. DIFFERENCE b. - c.                                ',ΔOBSPIVECTOR[5
], ' ',+/(CONTRDATA[;5]≥ΔOBSPIVECTOR[5])÷1000

```

```

□TCLF

```

```

'f. AVERAGE PI BETWEEN I. CORD ALLOKNOWN AND SYMP: ',ΔOBSPIVECTOR[6
]
'g. AVERAGE PI BETWEEN I. LAC  ALLOKNOWN AND SYMP: ',ΔOBSPIVECTOR[7
]
'h. AVERAGE PI BETWEEN I. CORD ALLOCLOSE AND SYMP: ',ΔOBSPIVECTOR[8
]
'i. AVERAGE PI BETWEEN I. LAC  ALLOCLOSE AND SYMP: ',ΔOBSPIVECTOR[9
]
'j. DIFFERENCE f. - g.                                ',ΔOBSPIVECTOR[1
0], ' ',+/(CONTRDATA[;10]≥ΔOBSPIVECTOR[10])÷1000
'k. DIFFERENCE h. - i.                                ',ΔOBSPIVECTOR[1
1], ' ',+/(CONTRDATA[;11]≥ΔOBSPIVECTOR[11])÷1000

```

```

□TCLF

```

```

'PROGRAM PIBOOT FINISHED.'

```

```

*****

```

```

PICONTRAST3

```

```

@ THIS PROGRAM PERFORMS BOOTSTRAP ON DIFFERENCES IN PI WITHIN AND BE
TWEEN SYMP AND ALLOPATRIC SAMPLES

```

```

@   CALLED BY 'PIBOOT'

```

```

@ BEFORE RUNNING THIS PROGRAM, RUN PICONTRAST2

```

```

PI←PITEMP×MASK4

```

```

PI2←PITEMP2×MASK4

```

```

N1←ρΔCSYMPINDEX
N2←ρΔLSYMPINDEX
N3←ρΔCALLO3
N4←ρΔLALLO3
N5←ρΔCALLOCLOSE
N6←ρΔLALLOCLOSE

PIBETWALLO←PI[ΔCALLO3;ΔLALLO3]
PIBETWALLOCLOSE←PI2[ΔCALLOCLOSE;ΔLALLOCLOSE]
PIBETWSYMP←PI[ΔCSYMPINDEX;ΔLSYMPINDEX]
PIBETWSYMP2←PI2[ΔCSYMPINDEX;ΔLSYMPINDEX]    a FOR USE IN ALLOCLOSE C
OMPARISONS

PIBETWCKNOWNC SYMP←PI[ΔCALLO3;ΔCSYMPINDEX]
PIBETWCCLOSECSYMP←PI2[ΔCALLOCLOSE;ΔCSYMPINDEX]

PIBETWLKNOWNLSYMP←PI[ΔLALLO3;ΔLSYMPINDEX]
PIBETWLCLOSELSYMP←PI2[ΔLALLOCLOSE;ΔLSYMPINDEX]

AVEBETWSYMP←(+ / + / PIBETWSYMP) ÷ (N1×N2)    a c.
AVEBETWALLO←(+ / + / PIBETWALLO) ÷ (N3×N4)    a a.
AVEBETWALLOCLOSE←(+ / + / PIBETWALLOCLOSE) ÷ (N5×N6)    a b.

AVEBETWCKNOWNC SYMP←(+ / + / PIBETWCKNOWNC SYMP) ÷ (N1×N3)    a f.
AVEBETWCCLOSECSYMP←(+ / + / PIBETWCCLOSECSYMP) ÷ (N1×N5)    a h.
AVEBETWLKNOWNLSYMP←(+ / + / PIBETWLKNOWNLSYMP) ÷ (N2×N4)    a g.
AVEBETWLCLOSELSYMP←(+ / + / PIBETWLCLOSELSYMP) ÷ (N2×N6)    a i.

BETW1←AVEBETWALLO-AVEBETWSYMP    a = ΔBETWDIFF1 IN PICONTRAST2; d.
BETW2←AVEBETWALLOCLOSE-AVEBETWSYMP    a = ΔBETWDIFF2    e.

BETWDIFF1←AVEBETWCKNOWNC SYMP-AVEBETWLKNOWNLSYMP    a = ΔBETWDIFF3 IN
PICONTRAST2; j. = f. - g.
BETWDIFF2←AVEBETWCCLOSECSYMP-AVEBETWLCLOSELSYMP    a = k. = h. - i.

TRSH←(AVEBETWALLO, AVEBETWALLOCLOSE, AVEBETWSYMP, BETW1,BETW2, AVEBE
TWCKNOWNC SYMP,AVEBETWLKNOWNLSYMP,AVEBETWCCLOSECSYMP,AVEBETWLCLOSELSY
MP,BETWDIFF1,BETWDIFF2 )
CONTRDATA←CONTRDATA,[1]TRSH
*****

```

IV. Calculating and bootstrapping D values (allele frequency differences between species)

A. Calculate average D values for allopatric and sympatric samples

FREQCONTRAST5

A THIS PROGRAM COMPARES AVERAGE PAIRWISE BETWEEN-SPECIES FREQUENCY DIFFERENCES FOR SYMPATRIC AND KNOWN ALLOPATRIC SAMPLES

A FOR EACH SNP, FREQUENCY OF ALLELE WITH HIGHEST FREQUENCY IN ALL SAMPLES IS CALCULATED FOR 4 GROUPS:

A (1) ALLOPATRIC I. CORDAT SAMPLES (2) ALLOPATRIC I. LAC SAMPLES (3) SYMPATRIC I. CORDAT SAMPLES

A (4) SYMPATRIC I. LAC SAMPLES

A THEN THREE FREQUENCY DIFFERENCES ARE CALCULATED: (1) (I CORDAT ALLO - I. SYMP ALLO) (2) (I. CORDAT SYMP - I. LAC SYMP)

A (3) ((2) - (1))

A THESE THREE FREQUENCY DIFFERENCES ARE THEN AVERAGED OVER ALL SNPS

ADIFFS←0 5p0

ADIFFS2←0 5p0

DIM←1↑pΔSCORES

I←0

RETI:I←I+1

TEST←(I÷1000)=(I(I÷1000))

±(TEST)/'I'

SCORES←ΔSCORES[I;;]

IND←161

IND2←(IND≠27)/IND

SCORES2←SCORES[IND;]

IND←(SCORES2[;1]≠'.')/11↑pSCORES2

SCORES3←SCORES2[IND;]

COUNTS←,(+/[1])(,SCORES3)◦.=ΔALLELES))

MAX←r/COUNTS

IND←(COUNTS=MAX)/15

MAXALLELE←'CTGA★'[IND]

MAXALLELE←1↑MAXALLELE

CASCOES←,SCORES[ΔCALLO3;]

CSSCOES←,SCORES[ΔCSYMPINDEX;]

LASCOES←,SCORES[ΔLALLO3;]

LSSCOES←,SCORES[ΔLSYMPINDEX;]

FCA←(+/CASCOES=MAXALLELE)÷(pCASCOES)

FCS←(+/CSSCOES=MAXALLELE)÷(pCSSCOES)

FLA←(+/LASCOES=MAXALLELE)÷(pLASCOES)

FLS←(+/LSSCOES=MAXALLELE)÷(pLSSCOES)

DIVIDER←-1 A SHOULD HAVE CALLED THIS 'MULTIPLIER'

±(FCA>FLA)/'DIVIDER←1'

SYMPDIFF←(FCS-FLS)×DIVIDER

ALLODIFF←(FCA-FLA)×DIVIDER

CDIFF1←(FCA-FCS)×DIVIDER

```

LDIFF1←(FLS-FLA)×DIVIDER
DIFF←ALLODIFF-SYMPDIFF
ΔDIFFS←ΔDIFFS,[1](ALLODIFF,SYMPDIFF,DIFF,CDIFF1,LDIFF1)

ALLODIFF←FCA-FLA
→(ALLODIFF=0)/DOWN

CDIFF←(FCA-FCS)÷ALLODIFF
LDIFF←(FLS-FLA)÷ALLODIFF

ALLOMIDPOINT←(FCA+FLA)÷2
ADJALLOMID←(FCA-ALLOMIDPOINT)÷ALLODIFF

ADJSYMPMID←(LDIFF+(1-CDIFF))÷2
MIDDIFF←ADJSYMPMID-ADJALLOMID
ΔDIFFS2←ΔDIFFS2,[1](CDIFF,LDIFF,ADJALLOMID,ADJSYMPMID,MIDDIFF)

DOWN:→(I<DIM)/RETI

SUM←+/[1]ΔDIFFS
AVES←SUM÷(1↑ρΔDIFFS)
□TCLF
'AVERAGE DIFFERENCES ACROSS SNPS'
'  ALLOPATRIC DIFFERENCE: ',AVES[1]
'  SYMPATRIC DIFFERENCE: ',AVES[2]
'  ALLO DIFF - SYMP DIFF: ',AVES[3]
'  CALLO - CSYMP          ',AVES[4]
'  LALLO - LSYP          ',AVES[5]
SUM2←+/[1]ΔDIFFS2
AVES2←SUM2÷(1↑ρΔDIFFS2)
□TCLF
'AVERAGE RELATIVE ALLO-SYMP DIFF FOR C: ',AVES2[1]
'AVERAGE RELATIVE ALLO-SYMP DIFF FOR L: ',AVES2[2]
'AVERAGE ALLO MIDPOINT: ',AVES2[3]
'AVERAGE SYM MIDPOINT: ',AVES2[4]
'ALLO - SYM MIDPOINTS: ',AVES2[5]

FREQCONTDIST  @  GENERATE DATA FOR SAS PLOT OF VARIABLES IN ΔDIFFS

□TCLF
'END PROGRAM FREQCONTRAST5'
*****

```

B. Bootstrapping differences. The following script calculates bootstrap 95% confidence intervals for means calculated in FREQCONTRAST5

```

*****
FREQCONTRAST6
@ THIS PROGRAM DOES BOOTSTRAPPING FOR KNOWN ALLO AND SYMP DIFFERENCE
S BETWEEN SPECIES
@ IT IS DERIVED FROM PROGRAM 'FREQCONTRAST'

```



```

A IT COMPARES AVERAGE PAIRWISE BETWEEN-SPECIES FREQUENCY DIFFERENCE
S FOR SYMPATRIC AND ALLOPATRIC SAMPLES
A FOR EACH SNP, FREQUENCY OF ALLELE WITH HIGHEST FREQUENCY IN ALL
SAMPLES IS CALCULATED FOR 4 GROUPS:
A (1) ALLOPATRIC I. CORDAT SAMPLES (2) ALLOPATRIC I. LAC SA
MPLES (3) SYMPATRIC I. CORDAT SAMPLES
A (4) SYMPATRIC I. LAC SAMPLES
A THEN THREE FREQUENCY DIFFERENCES ARE CALCULATED: (1) |(I CORDAT
ALLO - I. SYMP ALLO) | (2) |(I. CORDAT SYMP - I. LAC SYMP)
A (3) ((2) - (1))
A THESE THREE FREQUENCY DIFFERENCES ARE THEN AVERAGED OVER ALL SNP
S

```

```

A SHUFFLE ALLO AND SYMP SAMPLES WITHIN SPECIES

```

```

MAXJ←1000
ΔAVES←0 5ρ0

```

```

J←0
RETJ:J←J+1
'J= ',J

```

```

ΔSCORES2←ΔSCORES
IND1←?28ρ28
ΔSCORES2[;ΔCINDEX3;]←ΔSCORES2[;ΔCINDEX3[IND1];]
IND2←?28ρ28
ΔSCORES2[;ΔLINDEX3;]←ΔSCORES2[;ΔLINDEX3[IND2];]

```

```

ΔDIFFS←0 5ρ0
ΔDIFFS2←0 5ρ0
DIM←1↑ρΔSCORES2
I←0
RETI:I←I+1
TEST←(I÷10000)=(L(I÷10000))
⊡(TEST)/'I'

```

```

SCORES←ΔSCORES2[I;;]
IND←161
IND2←(IND≠27)/IND
SCORES2←SCORES[IND;]
IND←(SCORES2[;1]≠'.')/11↑ρSCORES2
SCORES3←SCORES2[IND;]
COUNTS←,(+/[1])(,SCORES3)◦.=ΔALLELES))
MAX←↑/COUNTS
IND←(COUNTS=MAX)/15
MAXALLELE←'CTGA★'[IND]
MAXALLELE←1↑MAXALLELE

```

```

CASCORES←,SCORES[ΔCALLO3;]
CSSCORES←,SCORES[ΔCSYMPINDEX;]
LASCORES←,SCORES[ΔLALLO3;]
LSSCORES←,SCORES[ΔLSYMPINDEX;]

```

```

FCA←(+/CASCORES=MAXALLELE)÷(ρCASCORES)

```

```

FCS←(+ /CSSCORES=MAXALLELE)÷(ρCSSCORES)
FLA←(+ /LASCORES=MAXALLELE)÷(ρLASCORES)
FLS←(+ /LSSCORES=MAXALLELE)÷(ρLSSCORES)
DIVIDER←-1
↓(FCA>FLA) / 'DIVIDER←1'

SYMPDIFF←(FCS-FLS)×DIVIDER
ALLODIFF←(FCA-FLA)×DIVIDER
CDIFF1←(FCA-FCS)×DIVIDER
LDIFF1←(FLS-FLA)×DIVIDER
DIFF←ALLODIFF-SYMPDIFF
ΔDIFFS←ΔDIFFS,[1](ALLODIFF,SYMPDIFF,DIFF,CDIFF1,LDIFF1)

DOWN:→(I<DIM)/RETI

SUM←+ / [1] ΔDIFFS
AVES←SUM÷(1↑ρΔDIFFS)
ΔAVES←ΔAVES,[1]AVES
□TCLF
'AVERAGE DIFFERENCES ACROSS SNPS'
'  ALLOPATRIC DIFFERENCE: ',AVES[1]
'  SYMPATRIC DIFFERENCE: ',AVES[2]
'  ALLO DIFF - SYMP DIFF: ',AVES[3]
'  CALLO - CSYMP          ',AVES[4]
'  LALLO - LSYP           ',AVES[5]

→(J<MAXJ)/RETJ

A  REPORT CONFIDENCE INTERVALS

LOW←L.025×1000
UP←U.975×1000

Q←ΔAVES[;1]
QQ←Q[ΔQ]
□TCLF
'95 PERCENT CONF INTERVAL FOR ALLOPATRIC FREQ DIFF BETWEEN SPECIES:'
'      (',(6 3*QQ[LOW]),', ' ,(6 3*QQ[UP]),')'

Q←ΔAVES[;2]
QQ←Q[ΔQ]
□TCLF
'95 PERCENT CONF INTERVAL FOR SYMPATRIC FREQ DIFF BETWEEN SPECIES:'
'      (',(6 3*QQ[LOW]),', ' ,(6 3*QQ[UP]),')'

Q←ΔAVES[;3]
QQ←Q[ΔQ]
□TCLF
'95 PERCENT CONF INTERVAL FOR ALLO - SYM FREQ DIFF BETWEEN SPECIES:'
'      (',(6 3*QQ[LOW]),', ' ,(6 3*QQ[UP]),')'

Q←ΔAVES[;4]
QQ←Q[ΔQ]
□TCLF

```

```
'95 PERCENT CONF INTERVAL FOR CORD ALLO - SYMP FREQ DIFF BETWEEN SPE
CIES:'
'          (',(6 3*QQ[LOW]),', ' ',(6 3*QQ[UP]),',')'
```

```
Q←ΔAVES[;5]
```

```
QQ←Q[ΔQ]
```

```
□TCLF
```

```
'95 PERCENT CONF INTERVAL FOR LAC ALLO-SYM FREQ DIFF BETWEEN SPECIES
:'
'          (',(6 3*QQ[LOW]),', ' ',(6 3*QQ[UP]),',')'
```

```
□TCLF
```

```
'PROGRAM FREQCONTRAST6 COMPLETED. DATA IN VARIABLE ΔAVES'
```

```
*****
```

V. M-K Tests

A. M-K test for all cordat samples vs all lac samples. This analysis uses two scripts: RUNLARGEDIFFS calls LARGEDIFFS. The former establishes a "cutoff". If the cutoff is, say, 0.9, then it treats all SNPs with frequency differences between species ≥ 0.9 but less than 1 as "fixed" differences, and all SNPs with freq difference < 0.9 as polymorphisms, then performs standard M-K test. RUNLARGEDIFFS establishes different cutoffs, then calls LARGEDIFFS to perform the corresponding M-K analysis.

RUNLARGEDIFFS

␣ THIS PROGRAM RUNS 'LARGEDIFFS' FOR DIFFERENT CUTOFF VALUES

CUTOFFS←11 2␣1 1.1 .9 1 .8 .9 .7 .8 .6 .7 .5 .6 .4 .5 .3 .4 .2 .3 .1
.2 0 .1

ΔOUTDATA←0 5␣0

II←0

RETII:II←II+1

CUTOFF←CUTOFFS[II;]

LARGEDIFFS

LINE←SYNFOCUS, NONFOCUS, SYNLOWER, NONLOWER, G

ΔOUTDATA←ΔOUTDATA, [1]LINE

→(II<10)/RETII

PI←ΔOUTDATA[;2]÷ΔOUTDATA[;1]

G←10 1␣ΔOUTDATA[;5]

␣TCLF

'PIN/PIS VALUES FOR DIFFERENT FREQUENCY-DIFFERENCE BINS'

TEMP←CUTOFFS[110;], G

TEMP

LARGEDIFFS

␣ THIS PROGRAM IDENTIFIES SNPS WITH LARGE DIFFERENCES BETWEEN SPECIES FOR SCORES VARIABLE X

␣ AND PERFORMS MK TEST FOR SELECTION

MAXI←1␣ΔSYNSCORES

NUCS←'ACGT★'

ΔFIXEDSNPS←0␣0 ␣ THIS VARIABLE CONTAINS SNP NUMBERS OF SNPS SHOWING FIXED DIFFS BETWEEN SPECIES

ΔFIXEDSNPPROPS←0 14␣' ' ␣ THIS IS A CHARACTER MATRIX WITH INFORMATION ON FIXED SNPS

CINDEX←ΔCINDEX ␣ USE ΔCINDEX FOR ALL SAMPLES, ΔCALLO3 FOR ONLY KNOW ALLOPATRIC SAMPLES, ETC.

LINDEX←ΔLINDEX

FIXEDSYN←FIXEDNON←0

```

A CUTOFF←1 1.05          A COMMENTED OUT WHEN RUN WITH PROGRAM RUNLARG
EDIFFS
SYNCOUNT←0 0
I←0
RETI:I←I+1
A TEST←(I÷100)=(I(I÷100))
A  $\Delta$ (TEST)/'I'
SCORES← $\Delta$ SYNSCORES[I;;]
SNPNUM← $\Delta$ SNPS1[I]

CSCORES1←,SCORES[CINDEX;]
LSCORES1←,SCORES[LINDEX;]

CSCORES←(CSCORES1  $\in$  'ACGT')/CSCORES1
LSCORES←(LSCORES1 $\in$  'ACGT')/LSCORES1

ALLELEIND←NUCS $\in$ (CSCORES,LSCORES)

ALLELES←ALLELEIND/NUCS
 $\Delta$ (1= $\rho$ ALLELES)/'FIXEDSYN←FIXEDSYN+1  $\diamond$  →DOWN1'

ALLELE1←ALLELES[1]
ALLELE2←ALLELES[2]
FREQ1←(+/(CSCORES=ALLELE1))/( $\rho$ CSCORES)
FREQ2←(+/(LSCORES=ALLELE1))/( $\rho$ LSCORES)
DIFF←|(FREQ1-FREQ2)
 $\Delta$ ((DIFF $\geq$ CUTOFF[1]) $\wedge$ (DIFF<CUTOFF[2]))/'SYNCOUNT[1]←SYNCOUNT[1]+1'
 $\Delta$ (DIFF $\geq$ CUTOFF[2])/'SYNCOUNT[2]←SYNCOUNT[2]+1'

DOWN1:→(I<MAXI)/RETI

MAXJ←1+ $\rho$  $\Delta$ NONSCORES
NONCOUNT←0 0
J←0
RETJ:J←J+1
SCORES← $\Delta$ NONSCORES[J;;]
CSCORES←,SCORES[CINDEX;]
LSCORES←,SCORES[LINDEX;]

CSCORES←(CSCORES  $\in$  'ACGT')/CSCORES
LSCORES←(LSCORES $\in$  'ACGT')/LSCORES

ALLELEIND←NUCS $\in$ (CSCORES,LSCORES)
ALLELES←ALLELEIND/NUCS

ALLELES←ALLELEIND/NUCS
 $\Delta$ (1= $\rho$ ALLELES)/'FIXEDNON←FIXEDNON+1  $\diamond$  →DOWN2'

ALLELE1←ALLELES[1]
ALLELE2←ALLELES[2]
FREQ1←(+/(CSCORES=ALLELE1))/( $\rho$ CSCORES)
FREQ2←(+/(LSCORES=ALLELE1))/( $\rho$ LSCORES)
DIFF←|(FREQ1-FREQ2)
 $\Delta$ ((DIFF $\geq$ CUTOFF[1]) $\wedge$ (DIFF<CUTOFF[2]))/'NONCOUNT[1]←NONCOUNT[1]+1'
 $\Delta$ (DIFF $\geq$ CUTOFF[2])/'NONCOUNT[2]←NONCOUNT[2]+1'

```

```
DOWN2:→(J<MAXJ)/RETJ
```

```
SYNFOCUS←SYNCOUNT[1]
NONFOCUS←NONCOUNT[1]
SYNLOWER←(1+ρΔSYNSCORES)-(+/SYNCOUNT)+FIXEDSYN)
NONLOWER←(1+ρΔNONSCORES)-(+/NONCOUNT)+FIXEDNON)
□TCLF
'M-K TABLE FOR CUTOFF ',CUTOFF
□TCLF
'
          SYN      NONSYN'

'FIXED    ',SYNFOCUS, NONFOCUS
'POLY     ',SYNLOWER, NONLOWER
□TCLF
'RATIOS    ',(SYNFOCUS÷SYNLOWER),(NONFOCUS÷NONLOWER)

ALPHA←1-(SYNFOCUS×NONLOWER)÷(NONFOCUS×SYNLOWER)
NUMBER←ALPHA×NONFOCUS
□TCLF
'ALPHA, NUMBER ',ALPHA,NUMBER
□TCLF
MATRIX←2 2ρSYNFOCUS, NONFOCUS, SYNLOWER, NONLOWER
GTEST1 MATRIX
*****
```

B. M-K tests for allopatric-sympatric analysis. This analysis treats as “fixed” differences SNPs with allopatric frequency differences between CUTOFF1 and CUTOFF1 + 0.1 and sympatric frequency differences > CUTOFF2. All other SNPs are polymorphic SNPs.

```
*****
```

```
MKTEST1
```

```
⌘ THIS PROGRAM PERFORMS M-K TEST.  'FIXED' SAMPLES ARE THOSE FOR WHI
CH ALLOPATRIC ALLELE FREQ DIFFS
⌘ > CUTOFF1, AND SYMPATRIC ALLELE FREQ DIFFS > CUTOFF2
```

```
NUCS←'ACGT★'
```

```
CINDEX←ΔCALLO3      ⌘ USE EITHER ΔCALLO3 OR ΔCALLOCLOSE
LINDE←ΔLALLO3      ⌘ DITTO
```

```
CCOMBINDEX←CINDEX,ΔCSYMPINDEX    ⌘ INDEX OF ALL C SAMPLES
LCOMBINDEX←LINDE←ΔLSYMPINDEX    ⌘ INDEX OF ALL L SAMPLES
```

```
FIXEDSYN←FIXEDNON←0  ⌘ COUNTERS FOR NONVARIABLE SNPS AFTER '★' IS RE
MOVED
```

```
CUTOFF1←.9          ⌘ ADJUST THESE AS NEEDED
CUTOFF2←.9          ⌘ ADJUST THESE AS NEEDED
```

```
SYNCOUNT←0          ⌘ VARIABLE FOR COUNT OF SYN SNPS MEETING CUTO
FF CRITERIA
```

```

MAXI←1↑ρΔSYNSCORES
I←0
RETI:I←I+1      A LOOP FOR SYNONYMOUS SNPS
TEST←(I÷1000)=(↑(I÷1000))
↓(TEST)/'''I = ''',I'

SYNSCORES←ΔSYNSCORES[I;;]      A PICK SCORES FOR ITH SYN SNP

CALLOSCORES←,SYNSCORES[CINDEX;]      A PICK OUT SCORES FOR CALLO
SAMPLES
LALLOSCORES←,SYNSCORES[LINDEX;]      A PICK OUT SCORES FOR LALLO
SAMPLES
CSYMPSCORES←,SYNSCORES[ΔCSYMPINDEX;]      A PICK OUT SCORES FOR CSYMP
SAMPLES
LSYMPSCORES←,SYNSCORES[ΔLSYMPINDEX;]      A PICK OUT SCORES FOR LSYMP
SAMPLES

CALLOSCORES←(CALLOSCORES∈'ACGT')/CALLOSCORES      A REMOVE SCORES THAT A
RE '★' OR '.'
LALLOSCORES←(LALLOSCORES∈'ACGT')/LALLOSCORES      A DITTO
CSYMPSCORES←(CSYMPSCORES∈'ACGT')/CSYMPSCORES
LSYMPSCORES←(LSYMPSCORES∈'ACGT')/LSYMPSCORES

ALLSCORES←CALLOSCORES,LALLOSCORES,CSYMPSCORES,LSYMPSCORES
ALLELEIND←'ACGT'∈(ALLSCORES)

↓(1=+/ALLELEIND)/' FIXEDSYN←FIXEDSYN+1 ♦ →DOWN1'      A IF ONLY 1 ALLE
LE, SKIP

NUMA←+/ALLSCORES='A'      A COUNT NUMBERS OF EACH NUCLEOTIDE IN
TOTAL SAMPLE
NUMC←+/ALLSCORES='C'
NUMG←+/ALLSCORES='G'
NUMT←+/ALLSCORES='T'
NUMS←NUMA,NUMC,NUMG,NUMT
MAX←↑/NUMS      A PICK OUT ALLELE WITH LARGEST COUNT
IND←1↑(NUMS=MAX)/14
ALLELE1←'ACGT'[IND]

FREQ1←(+/(CALLOSCORES=ALLELE1))/ρCALLOSCORES      A CALC ALLELE FREQ
UENCIES IN DIFFERENT SAMPLES
FREQ2←(+/(LALLOSCORES=ALLELE1))/ρLALLOSCORES
FREQ3←(+/(CSYMPSCORES=ALLELE1))/ρCSYMPSCORES
FREQ4←(+/(LSYMPSCORES=ALLELE1))/ρLSYMPSCORES

DIFF1←|(FREQ1-FREQ2)      A CALCULATE ABSOLUTE VALUE OF DIFFERENCE
IN ALLOPATRIC FREQUENCIES
DIFF2←FREQ3-FREQ4      A CALCULATE DIFFERENCE IN SYMPATRIC FREQU
ENCIES
↓((FREQ1-FREQ2)<0)/'DIFF2←-1×DIFF2'      A ADJUST DIFF IN SYMP FREQS
TO CORRECT FOR WHICH SPECIES HAS LARGER FREQ

↓((DIFF1≥CUTOFF1)^(DIFF2≥CUTOFF2))/'SYNCOUNT←SYNCOUNT+1'      A ADD 1 T
O SYNCOUNT IF CRITERION MET

```

40


```

1((DIFF1>=CUTOFF1)^(DIFF2>=CUTOFF2))/'NONCOUNT<-NONCOUNT+1'      A ADD 1 T
O NONNCOUNT IF CRITERION MET

DOWN2:-(J<MAXJ)/RETJ

SYNFOCUS<-SYNCOUNT
NONFOCUS<-NONCOUNT
SYNLOWER<-(1+ρΔSYNSCORES)-(SYNCOUNT+FIXEDSYN)
NONLOWER<-(1+ρΔNONSCORES)-(NONCOUNT+FIXEDNON)
□TCLF
'M-K TABLE FOR CUTOFFS 1 AND 2 ',CUTOFF1,CUTOFF2
□TCLF
'
      SYN      NONSYN '

'FIXED      ',SYNFOCUS, NONFOCUS
'POLY       ',SYNLOWER, NONLOWER
□TCLF
'RATIOS      ',(SYNFOCUS÷SYNLOWER),(NONFOCUS÷NONLOWER)

ALPHA<-1-(SYNFOCUS×NONLOWER)÷(NONFOCUS×SYNLOWER)
NUMBER<-ALPHA×NONFOCUS
□TCLF
'ALPHA, NUMBER ',ALPHA,NUMBER
□TCLF
MATRIX<-2 2ρSYNFOCUS, NONFOCUS, SYNLOWER, NONLOWER
GTEST1 MATRIX
□TCLF
'END PROGRAM MKTEST1'
*****

```

C. Messer-Petrov analysis of α .

1. This analysis requires identification of ancestral allele at each SNP. To do this, the scripts 'PICKTRIF' and 'PICKTRILO' were used to identify alleles in *I. trifida* and *I. triloba* corresponding to the previously identified SNPs in the transcriptome. This was done using maf files from the alignment of the *I. lacunosa* genome to the *I. trifida* and *I. triloba* genomes. Below is the listing of 'PICKTRIF'. The script 'PICKTRILO' is identical except for it referral to the *I. triloba* genome. The programs produce the vectors 'TRIFNUCS' and 'TRILONUCS', which contain ancestral nucleotides in position corresponding to appropriate SNP (e.g. position corresponding to positions in ΔTIGS and ΔPOS).

```

*****
PICKTRIF
A THIS PROGRAM READS IN LAC ALIGNMENT TO TRIFIDA FROM C:\TRIFLAC.TX
T AND DETERMINS THE TRIF NUCLEOTIDE
A CORRESPONDING TO LAC SNPS.

STARBYTE<-0
LOWER<-'actg'

```

UPPER←'ACTG'

BADCONTIG← 0 3ρ0

DIM←1↑ρΔSCORES

TRIFNUCS←DIMρ' ' A THIS WILL HOLD TRIFIDA NUCLEOTIDES; ORDER COR
RESPONDS TO ORDER IN ΔSCORES

TYPESNP←DIMρ' ' A THIS WILL HOLD TYPE OF SNP (E.G. SYN, NONSYN,
OTHER)

CURRCONTIG←0

I←0

A READ IN FIRST PART OF FILE

'C:\TRIFLAC.TXT' ONTIE -1

BUFFER←ONREAD (-1,82,3000000,STARBYTE) A READ IN 3M CHARACTERS
ONUNTIE -1

IND11←(BUFFER=OAV[11])/1ρBUFFER

A MAKE AN INDEX OF WHERE

OAV[11]'S ARE

HEADER←IND11[4]↑BUFFER

A HEADER

BUFFER←IND11[4]↓BUFFER

A STRIP HEADER FROM BUFFE

R

UP2:

READBUFF A CALL 'READBUFF' TO PICK FIRST SEQUENCE SEGMENT FROM BUFF
ER

DOWN6: A SET INFO TO 'OLD' INFO

OLDSEQUENCE←SEQUENCE

OLDSOURCESIZE←SOURCE SIZE

OLDSTRAND←STRAND

OLD SIZE←SIZE

OLDSTART←START

OLDSPECIES←SPECIES

OLDCONTIG←CONTIG

OLDKEEPINFO←KEEPINFO

UP1: READBUFF A CALL 'READBUFF' TO PICK NEXT SEQUENCE SEGMENT FROM
BUFFER

I←I+1

→(SPECIES='T')/DOWN1 A SKIP TO DOWN 1 IF READ SEQUENCE IS FROM TRI
FIDA

A IF ANOTHER LAC SEQUENCE SET INFO TO 'OLD' INFO

OLDSEQUENCE←SEQUENCE

OLDSOURCESIZE←SOURCE SIZE

OLDSTRAND←STRAND

OLD SIZE←SIZE

OLDSTART←START

OLDSPECIES←SPECIES

OLDKEEPINFO←KEEPINFO

→UP1 A GO UP AND READ IN NEXT SEQUENCE AND INFO

DOWN1:

A PICK OUT TRIF NUCS

IND1←ΔTIGS=OLDCONTIG A INDEX FOR ALL SNPS WITH CURRE
NT CONTIG

OLDEND←(−1+OLDSTART+OLDSIZE) A POSITION OF END OF SEQUENCE O
N CONTIG

IND2←(ΔPOS≥OLDSTART)^(ΔPOS≤OLDEND) A INDEX FOR ALL SNPS WITH POSIT
ION BETWEEN BEGINNING AND END OF SEQUENCE

IND3←IND1^IND2 A LOGICAL AND FOR IND1 AND IND1
: ALL SNPS WITHIN SEQUENCE

IND3←IND2^IND3

NUMS←(IND3)/1ρIND3 A SNP NUMBER FOR THOSE SNPS (I.
E. POSITION IN ΔSCORES)

→(0=ρNUMS)/DOWN2 A SKIP SEQUENCE IF NO SNPS

A PROCESS SNPS

READLACSCAFF OLDCONTIG A CALL READLACSCAFF TO READ IN
LAC CONTIG SEQUENCE

IND4←(OLDSEQUENCE≠'−') A INDEX OF '−'S
LACSEQ←IND4/OLDSEQUENCE A COMPRESS −S OUT OF LAC SEQUEN
CE

TRIFSEQ←IND4/SEQUENCE A COMPRESS CORRESPONDING POSITI
ONS OUT OF TRIF SEQUENCE

DIM2←ρNUMS

K←0

RETK:K←K+1

□TCLF

'*****'
*****'

□TCLF

'PROCESSING SNP'

CURNUM←NUMS[K]

CURPOS←ΔPOS[CURNUM]

LACSEQ2←(OLDSTARTρ'−'),LACSEQ
TRIFSEQ2←(OLDSTARTρ'−'),TRIFSEQ

UP3: □TCLF

'LAC, TRIF, AND GENOME SEQS'

SEQ3←OLDSTART+ΔSEQUENCE

SIZE2←1/(OLDSIZE,200)

LACSEQ[1SIZE2]

TRIFSEQ[1SIZE2]

```

SEQ3[1:SIZE2]

SC←,ΔSCORES[CURNUM;;]
LACSNP←LACSEQ2[CURPOS]
TRIFSNP←TRIFSEQ2[CURPOS]
GENSNP←ΔSEQUENCE[CURPOS]

DTCLF
'CONTIG ',OLDCONTIG,' PROCESSING ',(1+(STARBYTE÷3000000)), 'TH 3MB S
EGMENT'
BPPROC←(3000000-ρBUFFER)
' ',BPPROC,' BASES OF SEGMENT PROCESSED'
DTCLF
'SNP SCORES'
SC
DTCLF
'LAC, TRIF, AND GENOME NUCLEOTIDE'
LACSNP
TRIFSNP
GENSNP

IND21←(TRIFSEQ≠'-')
LACSEQ3←IND21/LACSEQ
TRIFSEQ3←IND21/TRIFSEQ

PCTEQ←(+/LACSEQ3=TRIFSEQ3)÷(ρLACSEQ3)
'PCTEQ ',PCTEQ
⊥(PCTEQ<.5)/'TRIFNUCS[CURNUM]←'M' ⋄ 'M INSERTED' ⋄ →DOWN4'
      A INSERTS M INTO TRIFNUCS TO SIGNIFY UNRELIABLE ALIGNMENT

TEST3←LACSNP≠SC
'TEST3 ',TEST3
⊥(TEST3=0)/'TRIFNUCS[CURNUM]←'N' ⋄ 'N INSERTED' ⋄ →DOWN4'
      A INSERTS N INTO TRIFNUCS TO SIGNIFY MISMATCH BETW ALLELES AND
      LACSNP

SEQ4←SEQ3[1:OLDSIZE]
PCTEQ2←(+/LACSEQ=SEQ4)÷(ρLACSEQ)
'PCTEQ2 ',PCTEQ2
LINEINFO←(CURNUM,CURPOS,OLDCONTIG)

⊥(PCTEQ2<0.9)/'BADCONTIG←BADCONTIG,[1]LINEINFO ⋄ 'BADCONTIG' '
      A INDICATES MISAL
IGNEMENT BETW LAC SEQ AND GENOME SCAFFOLD
TRIFNUCS[CURNUM]←TRIFSNP

DTCLF
'INSERTED NUCLEOTIDE: ',TRIFSNP

DOWN4: TYPESNP[CURNUM]←TYPE

TRSH←DDL 2
→(K<DIM2)/RETK

DOWN2: A 'BUFFER BEFORE RETURN'

```

```

DOWN4:→(100000<ρBUFFER)/DOWN5
STARBYTE←STARBYTE+3000000
'C:\TRIFLAC.TXT' UNTIE -1
BUFFER←BUFFER,UNREAD (-1,82,3000000,STARBYTE)      A READ IN 3M MORE
CHARACTERS
UNUNTIE -1

```

```

DTCLF
'3 M MORE BYTES ADDED TO 'BUFFER'. ',(STARBYTE÷3000000),'TH SEGMENT
READ.'

```

```

DOWN5:→UP2

```

```

DTCLF
'PROGRAM TRIFNUCS FINISHED. DATA IN VARIABLE 'TRIFNUCS'.'
*****

```

The following scripts are called by 'PICTRIF' and 'PICKTRILO':

a. READBUFF

```

*****

```

```

READBUFF

```

```

A THIS PROGRAM CALLED BY PICKTRIF
A IT READS IN NEXT ALIGNMENT FROM VARIABLE 'BUFFER'

```

```

UP1:

```

```

IND11←2↑((BUFFER=DAV[11])/ρBUFFER)      A POSITION OF FIRST
FIRST DAV[11] IN BUFFER
INFO←-1↓(IND11[1]↑BUFFER)      A GET INFO FOR NEXT ALIGNMENT
KEEPINFO←INFO

```

```

BUFFER←IND11[1]↓BUFFER      A DROP INFO FROM BUFFER

```

```

IND10←(INFO=DAV[10])/ρINFO      A MAKE INDEX OF WHERE DAV[10]'S ARE IN INFO
REST←IND10[6]↑INFO      A REST OF INFO BESIDES SEQUENCE
SEQUENCE←(IND10[6]↓INFO)      A PUT SEQUENCE IN 'SEQUENCE'
SEQUENCE←(UPPER,DAV)[(LOWER,DAV)↓SEQUENCE]
END←-1↑SEQUENCE
↓(END<(DAV[10 11]))/'SEQUENCE←-1↓SEQUENCE'

```

```

SOURCE SIZE←DFI -1↓(IND10[5]↓REST)      A SIZE OF ENTIRE SOURCE SEQUENCE, NOT JUST PARTS INVOLVED IN ALIGNMENT
REST←IND10[5]↑REST      A REST OF INFO BESIDES SOURCE SIZE
STRAND←1↑(IND10[4]↓REST)      A STRAND (+ OR -)

```

REST<-IND10[4]↑REST RAND	A REST OF INFO BESIDES ST
SIZE<-OFI -1↓(IND10[3]↓REST) IN LAC	A SIZE OF ALIGNING REGION
REST<-IND10[3]↑REST ZE	A REST OF INFO BESIDES SI
DASH CHARACTERS	A EQUAL TO NUMBER OF NON-
START<-OFI -1↓(IND10[2]↓REST) ING REGION ON LAC CONTIG	A START POSITION OF ALIGN
REST<-IND10[2]↑REST RT	A REST OF INFO BESIDES STA
TEMP<-IND10[1]↓REST	A INFO ON CONTIG NUMBER
CONTIG<-OFI (TEMP<'0123456789')/TEMP	A CONTIG NUMBER
TEST<-+/(TEMP OFS 'lac') ↓(TEST=1)/'SPECIES<'L''	
TEST<-+/(TEMP OFS 'trif') ↓(TEST=1)/'SPECIES<'T''	
↓(BUFFER[1]='a')/'BUFFER<2↓BUFFER'	
↓(BUFFER[2]='a')/'BUFFER<3↓BUFFER'	

b. READLACSCAFF. This program reads in the sequence of a particular lacunosa genome contig (SCAFFN) into variable Δ SEQUENCE

READLACSCAFF SCAFFN

␣ BEFORE RUNNING THIS, MAKE SURE TO RUN 'CONVLACINDEX' TO CONVERT Δ L
ACINDEX FROM A CHARACTER
␣ TO A NUMERIC MATRIX

␣ THIS PROGRAM READS IN I. LAC SCAFFOLD FROM FILE C:\LACGENOME
␣ SCAFFN IS THE SCAFFOLD NUMBER TO READ

SCAFFNUMS \leftarrow Δ LACINDEX[;1]
IND \leftarrow (SCAFFN=SCAFFNUMS)/1␣SCAFFNUMS
TEMP1 \leftarrow , Δ LACINDEX[IND;]
STARTBYTE \leftarrow TEMP1[2]
ENDBYTE \leftarrow TEMP1[3]

'C:\LACGENOME' ␣NTIE -1
SCAFFOLD \leftarrow ␣NREAD -1 82 (ENDBYTE-STARTBYTE) STARTBYTE
␣NUNTIE -1
SCAFFNAME \leftarrow 24↑SCAFFOLD
 Δ SEQUENCE \leftarrow 24↓SCAFFOLD
 Δ SEQUENCE \leftarrow (Δ SEQUENCE \neq ␣TCLF)/ Δ SEQUENCE

2. Ancestral alleles in TRIFNUCS and TRILONUCS are merged with script 'MERGE'.
Ancestral alleles from TRILONUCS take precedence over those from TRIFNUCS when both are
identified for a given SNP.

X MERGE Y

␣ THIS PROGRAM MERGES RESULTS FROM TRIFIDA AND TRILOBA
␣ X IS TRIFNUCS, Y IS TRILONUCS

MAXI \leftarrow ␣Y

MERGED \leftarrow 0␣' '

I \leftarrow 0
RETI:I \leftarrow I+1

X1 \leftarrow X[I]
Y1 \leftarrow Y[I]

⧿(X1 \in 'ACGT')/'MERGED \leftarrow MERGED,X1 ⧿ →DOWN'
⧿(Y1 \in 'ACGT')/'MERGED \leftarrow MERGED,Y1 ⧿ →DOWN'

```

MERGED←MERGED, ' '
DOWN:→(I<MAXI)/RETI

□TCLF
'TOTAL SNPS IDENTIFIED = ',+/MERGED∈'ACGT'
□TCLF

'PROGRAM MERGE FINISHED. DATA IN VARIABLE 'MERVED'.'.
*****

```

3. Ancestral nucleotides are separated for synonymous, non-synonymous, and non-coding SNPs, and then for each SNP pair the Messer-Petrov α is calculated, as well as distance separating the two SNPs. The script produces a matrix 'MATRIX' that has 5 columns:

Col1: midpoint of distance bin (distance between SNPS)
Col2: $\alpha(d)$ for each distance bin for fixed differences (non-syn vs. syn SNPs)
Col3: $\alpha(d)$ for each distance bin for nearly fixed differences (non-syn vs. syn SNPs)
Col4: $\alpha(d)$ for each distance bin for fixed differences (non-coding vs. syn SNPs)
Col5: $\alpha(d)$ for each distance bin for nearly fixed differences (non-coding vs. syn SNPs)

It also produces vectors A1A, A1B, A2A and A2B corresponding to cols 2 – 5.

The matrix is used for analysis in SAS. The vectors are used for analysis in Mathematica

```

*****
SEPMERGEDNUCS

A THIS PROGRAM SEPARATES MERGED NUCS INTO SYN, NON-SYN, AND REG ANC
ESTRAL NUCS
A THEN ANALYZES THEM FOR MESSER AND PETROV MK TESTS
A USES ONLY ALL SAMPLES

MAXI←ρΔMERGED A ΔMERGED IS SAME AS 'MERGED' PRODUCED BY PROGRAM '
MERGE'

NUMS←□FI ,ΔLACCODONS[;13+16] A SNP NUMBERS FROM ΔLACCODONS

CSYNFREQS←LSYNFREQS←CNONFREQS←LNONFREQS←CRFREQS←LRFREQS←0ρ0 A HOL
D COUNTS

I←0
RETI:I←I+1

SC←ΔSCORES[I;;] A READ IN SCORES FOR SNP I
CSC←,SC[ΔCINDEX;] A PICK OUT CORD SCORES
CSC←(CSC∈'ACGT')/CSC A GET RID OF '.'
LSC←,SC[ΔLINDEX;] A PICK OUT LAC SCORES
LSC←(LSC∈'ACGT')/LSC A GET RID OF '.'
BOTH←CSC,LSC A COMBINE CORD AND LAC SCORES

TNUC←ΔMERGED[I] A PICK CORRESPONDING MERGED NUC
→(TNUC∈' NM')/DOWN A IF NUC∈NM, SKIP

```



```

SNP←ΔLACSNPS[I;10]      A READ SNP NUMBER FROM ΔLACSNPS
TEST←SNP∈NUMS            A TEST IF SNP NUMBER IN ΔLACCODONS
→(TEST=0)/DOWN           A IF NOT, SKIP

IND←(NUMS=SNP)/ιρNUMS    A POSITION OF SNP IN ΔLACCODONS
TYPE←ΔLACCODONS[IND;12] A GET TYPE (NON-SYN, SYN, ETC.)

→(TYPE='O')/DOWN        A SKIP IF TYPE='O'

TEST←(TNUC∈BOTH)         A TEST IF MERGED ALLELE IN LAC OR CORD (AN
CESTRAL ALLELE)
→(TEST=0)/DOWN           A SKIP IF NOT

CFREQ←(+/CSC=TNUC)÷ρCSC  A ANCESTRAL ALLELE FREQ IN CORD
LFREQ←(+/LSC=TNUC)÷ρLSC  A ANCESTRAL ALLELE FREQ IN LAC

CFREQ←1-CFREQ            A FREQ OF NEW ALLELE IN CORD
LFREQ←1-LFREQ            A FREQ OF NEW ALLELE IN LAC

A APPEND NEW ALLELE FREQ TO APPROPRIATE VECTOR

⊕(TYPE='N')/'CNONFREQS←CNONFREQS,CFREQ ⊙ LNONFREQS←LNONFREQS,LFREQ'
⊕(TYPE='S')/'CSYNFREQS←CSYNFREQS,CFREQ ⊙ LSYNFREQS←LSYNFREQS,LFREQ'
⊕(TYPE='R')/'CRFREQS←CRFREQS,CFREQ ⊙ LRFREQS←LRFREQS,LFREQ'

DOWN:→(I<MAXI)/RETI

A BIN THE FREQUENCIES

NUMBINS←50      A CHANGE THIS TO WHATEVER
BINS←(ιNUMBINS)÷NUMBINS
BINS←BINS-BINS[1]

NBINS←+/CNONFREQS°.≥BINS
CNONTOT←+/[1]NBINS°.=(ιNUMBINS)

NBINS←+/LNONFREQS°.≥BINS
LNONTOT←+/[1]NBINS°.=(ιNUMBINS)

NBINS←+/CSYNFREQS°.≥BINS
CSYNTOT←+/[1]NBINS°.=(ιNUMBINS)

NBINS←+/LSYNFREQS°.≥BINS
LSYNTOT←+/[1]NBINS°.=(ιNUMBINS)

NBINS←+/CRFREQS°.≥BINS
CRTOT←+/[1]NBINS°.=(ιNUMBINS)

NBINS←+/LRFREQS°.≥BINS
LRTOT←+/[1]NBINS°.=(ιNUMBINS)

NONTOT←CNONTOT+LNONTOT      A SUM NON-SYN SNPS OVER SPECIES

```

```

SYNTOT←CSYNTOT+LSYNTOT      A SUM SYN SNPS OVER SPECIES

RTOT←CRTOT+LRTOT            A SUM NON-CODING SNPS OVER SPECIES

A CALCULATE ALPHA VECTORS FOR MESSER AND PETROV MK
A NON-SYN FREQ DIFF = 1
DS←169
DN←190
ALPHA1A←1 - (DS÷DN)×NONTOT÷SYNTOT

A NON-SYN FREQ DIFF [0.9, 1)
DS←699
DN←644
ALPHA1B←1 - (DS÷DN)×(NONTOT÷SYNTOT)

A NON-CODING FREQ DIFF=1
DR←107
DS←169
ALPHA2A←1 - (DS÷DR)×(RTOT÷SYNTOT)

A NON-CODING FREQ DIFF [0.9, 1)
DS←699
DR←333
ALPHA2B←1 - (DS÷DR)×(RTOT÷SYNTOT)

A CONVERT DATA TO MATRIX FOR SAS
MIDPOINTS←((1NUMBINS)÷NUMBINS)-(1÷(2×NUMBINS))  A MIDPOINTS OF BINS

COUNT←+/BINS<.9      A NUMBER OF BINS WITH FREQ ≤ .9

MATRIX←(COUNT 1ρMIDPOINTS),(COUNT 1ρALPHA1A),(COUNT 1ρALPHA1B),(COUNT 1ρALPHA2A),(COUNT 1ρALPHA2B)

A CONVERT DATA TO MATRICES FOR MATHEMATICA

TMP←(COUNT 1ρMIDPOINTS),(COUNT 1ρALPHA1A)
CONVMATH2 TMP
A1A←CONVDATA

TMP←(COUNT 1ρMIDPOINTS),(COUNT 1ρALPHA1B)
CONVMATH2 TMP
A1B←CONVDATA

TMP←(COUNT 1ρMIDPOINTS),(COUNT 1ρALPHA2A)
CONVMATH2 TMP
A2A←CONVDATA

TMP←(COUNT 1ρMIDPOINTS),(COUNT 1ρALPHA2B)
CONVMATH2 TMP
A2B←CONVDATA

SKIP:

□TCLF

```

```
'PROGRAM SEPMERGEDNUCS2 FINISHED. DATA FOR SAS IN 'MATRIX'. DATA
FOR MATHEMATICA IN A1A, A1B,A2A,A2B.'
*****
```

4. Estimating $\alpha(1)$ (asymptotic value of α) using SAS. The data in 'MATRIX' produced by SEPMERGEDNUCS (immediately above) is cut and pasted into the following SAS program to estimate non-linear regression coefficients. This example is for an analysis of non-synonymous vs synonymous SNPs from all samples, and fixed differences between species.

```
*****
data alpha1;
  input x ala alb a2a a2b;
  cards;
[data MATRIX]
  run;

proc nlin data=alpha1;
  parms a=.042 b=.937 c=.8353;
  model ala=1-(a+b*exp(-c*x));
  title nonlin fit for all samples, freq diff = 1, non-syn;
run;

*****
```

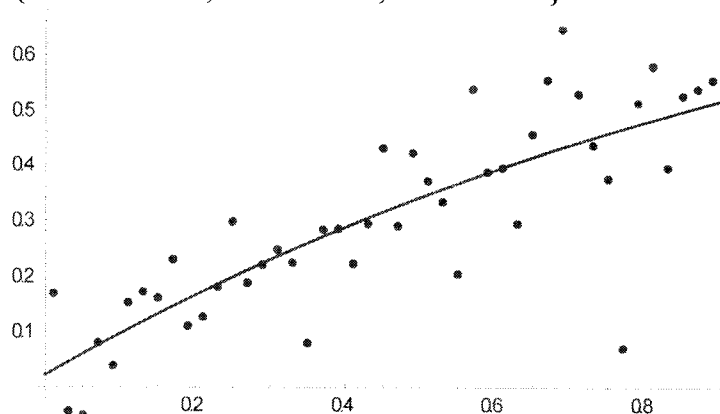
5. Estimating $\alpha(1)$ (asymptotic value of α) using MATHEMATICA. The data in vectors A1A, A1B, A2A and A2B are cut and pasted into the MATHEMATICA program (highlighted in yellow below), which produces estimates of the three regression parameters, and also plots the data and the fitted regression.

```
*****
(* all samples , non_syn, freq diff = 1 *)
x={{0.01,0.1682589312},{0.03,-0.04384515007},{0.05,-
0.05024709661},{0.07,0.08014811619},{0.09,0.03959845735},{0.11,0.1521155831},{0.13,0.1
730674342},{0.15,0.1619181287},{0.17,0.2295005028},{0.19,0.1105263158},{0.21,0.1266399
695},{0.23,0.180632616},{0.25,0.29825443},{0.27,0.1890092879},{0.29,0.2209437387},{0.31
,0.247871517},{0.33,0.2245614035},{0.35,0.07985480944},{0.37,0.283110762},{0.39,0.28533
61728},{0.41,0.222109036},{0.43,0.2945553539},{0.45,0.4317251462},{0.47,0.2900165211},
{0.49,0.4218421053},{0.51,0.3719776715},{0.53,0.3337152845},{0.55,0.2039964318},{0.57,0
.5381578947},{0.59,0.3870853605},{0.61,0.3942847917},{0.63,0.2945553539},{0.65,0.45542
4275},{0.67,0.5552631579},{0.69,0.6475670308},{0.71,0.5305555556},{0.73,0.4359435173},
{0.75,0.3749644381},{0.77,0.07099415205},{0.79,0.5129072682},{0.81,0.5797762122},{0.83,
0.3960363873},{0.85,0.5256140351},{0.87,0.5385437277},{0.89,0.5552631579}};
x1=x[[All,{1,2}]];
nlm=NonlinearModelFit[x1,(1-(a+b Exp[-c y])),{a,b,c},y]
nlm[[1,2]]
Show[ListPlot[x1],Plot[nlm[x],{x,0,.9}]]
nlm[1] (* value of alpha(1) *)
```

(* output *)

FittedModel[$0.958508 - 0.936118e^{-0.835267y}$]

{a->0.0414925,b->0.936118,c->0.835267}



0.552458

6. Estimating 95% Confidence Intervals for $\alpha(1)$. This is done by bootstrapping over SNPs within a SNP type (e.g. non-synonymous, synonymous, non-coding) using scripts NONLIN and BOOT1 (the former calls the latter). The script produces four matrices (MMAT1A, MMAT1B, MMAT2A, and MMAT2B) corresponding, respectively, to analyses of 1. Fixed differences, non-syn vs. syn; 2. Nearly fixed differences, non-syn vs. syn; 3. Fixed differences, non-coding vs. syn; and 4. Nearly fixed differences, non-coding vs. syn. These matrices are converted to character vectors for reading into MATHEMATICA. These vectors are saved as text files (APL Native Files), which are read by MATHEMATICA.

Each matrix consists of n columns and 1001 rows, where n is the number of allele frequency bins. The first column consists of the midpoints of each bin. Each of the remaining columns corresponds to one bootstrap sample. For each sample, the values are the $\alpha(d)$ values corresponding to the appropriate bin. The MATHEMATICA program calculates $\alpha(1)$ for each bootstrap sample and then calculates the confidence interval by ordering the $\alpha(1)$'s and taking the 25th and 975th values.

NONLIN

@ THIS PROGRAM RUNS NON-LINEAR REGRESSION ON BOOTSTRAP ALPHA DATA

@ NOTE: HAVE MIDPOINTS FROM PREVIOUSLY RUN SEPMERGEDNUCS

COUNT<+ /MIDPOINTS<.9

MAT1A<MAT1B<MAT2A<MAT2B<(COUNT 1ρMIDPOINTS)

MAXI<1000

I<0

RETI:I<I+1

```

TEST←(I÷100)=(↑(I÷100))
⚡(TEST=1)/' 'I = ',I'

BOOT1      ⚡ CALL BOOT1 TO PROCESS EACH BOOTSTRAP SAMPLE

MAT1A←MAT1A,MATRIX[;2]
MAT1B←MAT1B,MATRIX[;3]
MAT2A←MAT2A,MATRIX[;4]
MAT2B←MAT2B,MATRIX[;5]

→(I<MAXI)/RETI

⚡ CONVERT MATRICES TO TEXT AND SAVE TO NATIVE FILES
MMAT1A←,⚡MAT1A
IND←(MMAT1A=' ')/⌈ρMMAT1A
MMAT1A[IND]←'-'

'C:\MMAT1A' ⓀNCREATE -1
MMAT1A ⓀNAPPEND -1
ⓀNUNTIE -1

MMAT1B←,⚡MAT1B
IND←(MMAT1B=' ')/⌈ρMMAT1B
MMAT1B[IND]←'-'

'C:\MMAT1B' ⓀNCREATE -1
MMAT1B ⓀNAPPEND -1
ⓀNUNTIE -1

MMAT2A←,⚡MAT2A
IND←(MMAT2A=' ')/⌈ρMMAT2A
MMAT2A[IND]←'-'

'C:\MMAT2A' ⓀNCREATE -1
MMAT2A ⓀNAPPEND -1
ⓀNUNTIE -1

MMAT2B←,⚡MAT2B
IND←(MMAT2B=' ')/⌈ρMMAT2B
MMAT2B[IND]←'-'

'C:\MMAT2B' ⓀNCREATE -1
MMAT2B ⓀNAPPEND -1
ⓀNUNTIE -1

ⓀTCLF
'PROGRAM NONLIN FINISHED.  DATA STORED IN NATIVE FILES MMAT1A, MMAT1
B, MMAT2A, MMAT2B.'

*****
BOOT1

⚡ THIS PROGRAM BOOTSTRAPS THE ALPHA VALUES

DIM←ρCNONFREQS

```

```

IND←?DIMρDIM
BCNONFREQS←CNONFREQS[IND]  A NOTE: CNONFREQS, ETC. FROM PREVIOUS RUN
  OF SEPMERGEDNUCS
BLNONFREQS←LNONFREQS[IND]

```

```

DIM←ρCSYNFREQS
IND←?DIMρDIM
BCSYNFREQS←CSYNFREQS[IND]
BLSYNFREQS←LSYNFREQS[IND]

```

```

DIM←ρCRFREQS
IND←?DIMρDIM
BCRFREQS←CRFREQS[IND]
BLRFREQS←LRFREQS[IND]

```

```

A BIN THE FREQUENCIES  NOTE: BINS AND NUMBINS FROM PROGRAM SEPMERGE
DNucs

```

```

NBINS←+/BCNONFREQS°.>BINS
CNONTOT←+/[1]NBINS°.= ( 1NUMBINS)

```

```

NBINS←+/BLNONFREQS°.>BINS
LNONTOT←+/[1]NBINS°.= ( 1NUMBINS)

```

```

NBINS←+/BCSYNFREQS°.>BINS
CSYNTOT←+/[1]NBINS°.= ( 1NUMBINS)

```

```

NBINS←+/BLSYNFREQS°.>BINS
LSYNTOT←+/[1]NBINS°.= ( 1NUMBINS)

```

```

NBINS←+/BCRFREQS°.>BINS
CRTOT←+/[1]NBINS°.= ( 1NUMBINS)

```

```

NBINS←+/BLRFREQS°.>BINS
LRTOT←+/[1]NBINS°.= ( 1NUMBINS)

```

```

NONTOT←CNONTOT+LNONTOT      A SUM NON-SYN SNPS OVER SPECIES
SYNTOT←CSYNTOT+LSYNTOT     A SUM SYN SNPS OVER SPECIES
RTOT←CRTOT+LRTOT           A SUM NON-CODING SNPS OVER SPECIES

```

```

A CALCULATE ALPHA VECTORS FOR MESSER AND PETROV MK
A NON-SYN FREQ DIFF = 1
DS←169
DN←190
ALPHA1A←1 - (DS÷DN)×NONTOT÷SYNTOT

```

```

A NON-SYN FREQ DIFF [0.9, 1)
DS←699
DN←644
ALPHA1B←1 - (DS÷DN)×(NONTOT÷SYNTOT)

```

```

A →DOWN3  A ONLY FOR CLOSE ALLOPATRIC

```

```

A NON-CODING FREQ DIFF=1
DR←107

```

```

DS←169
ALPHA2A←1 - (DS÷DR)×(RTOT÷SYNTOT)

# NON-CODING FREQ DIFF [0.9, 1)
DS←699
DR←333
ALPHA2B←1 - (DS÷DR)×(RTOT÷SYNTOT)

DOWN3:

# CONVERT DATA TO MATRIX FOR SAS

MATRIX←(COUNT 1ρMIDPOINTS),(COUNT 1ρALPHA1A),(COUNT 1ρALPHA1B),(COUNT 1ρALPHA2A),(COUNT 1ρALPHA2B)
# NOTE: COUNT IS FROM PROGRAM SEPMERGEDNUCS

# CONVERT DATA TO MATRICES FOR MATHEMATICA

TMP←(COUNT 1ρMIDPOINTS),(COUNT 1ρALPHA1A)
CONVMATH TMP
A1A←CONVDATA

TMP←(COUNT 1ρMIDPOINTS),(COUNT 1ρALPHA1B)
CONVMATH TMP
A1B←CONVDATA

TMP←(COUNT 1ρMIDPOINTS),(COUNT 1ρALPHA2A)
CONVMATH TMP
A2A←CONVDATA

TMP←(COUNT 1ρMIDPOINTS),(COUNT 1ρALPHA2B)
CONVMATH TMP
A2B←CONVDATA
*****

```

The MATHEMATICA program for the estimation of bootstrapped CI's:

```

*****
SetDirectory["C:\Users\mrausher\Desktop"]
bootn=1000
x=ReadList["MMAT1A",Table[Number,{bootn+1}]];
list={};
For[i=2,i<(bootn+2),i++,x1=x[[All,{1,i}]];
  nlm=NonlinearModelFit[x1,(1-(a+b Exp[-c y])),{a,b,c},y];
  alpha1=nlm[1];AppendTo[list,alpha1>(*Print[Show[ListPlot[x1],Plot[nlm[x],{x,0,.9}]]]*);
list2=Sort[list];
confint={list2[[bootn .05]],list2[[bootn .95]]}
*****

```

VI. Analysis of admixture linkage disequilibrium (ALD). This is done with two scripts, ALD and COUNTALDB. For each *I. lacunosa* genome contig, ALD identifies synonymous SNPs on that contig and calculates pairwise distances and pairwise ald using the formula in the text. COUNTALDB bins these pairwise values according to pairwise distances and averages ald for all pairs within a bin. COUNTALDB calls script MAKEBINS, which defines the distance bins, and produces matrix 'ALDBINS', which is an N x 5 matrix with the following columns:

Col1 distance bin midpoint
 Col2 ald(d)
 Col3 count of SNP pairs in distance bin
 Col4 mean absolute weightings of SNP pairs in distance bin
 Col5 mean covariance between SNP pairs in distance bin

The versions of the scripts below are for calculating ald in the *I. cordatotriloba* sympatric samples.

ALD

A THIS PROGRAM CALCULATES ADMIXTURE LD FOR CORDAT SYMP POPULATION

SNPNUM←1↑ρΔSYNSCORES

ALD1←0 2ρ0

DIMSYMP←ρΔCSYMPINDEX

DATA←0 4ρ0

MAXI←840

I←0

RETI:I←I+1 A LOOP FOR CONTIGS

DTCLF

'I = ',I

TIG←ΔUNIQUE TIGS[I] A PICK OUT CONTIG NUMBER

SCINDEX←(ΔSYNSNPS[;1]=TIG)/1SNPNUM A INDEX OF SNPS ON CONTIG

→(1≥ρSCINDEX)/DOWN

SC←ΔSYNSCORES[SCINDEX;;] A PICK SCORES FOR EACH SNP ON CONTIG

SPOS←ΔSYNSNPS[SCINDEX;3] A PICK POSITIONS OF EACH SNP ON CONTIG

MAXSNP←ρSCINDEX A NUMBER OF SNPS ON CONTIG

J←0

RETJ:J←J+1 A SNP1 LOOP

SC1←SC[J;;]

SC1CA←SC1[ΔCALLO3;]

SC1LA←SC1[ΔLALLO3;]

UNSC1CS←SC1CS←SC1[ΔCSYMPINDEX;]

TEST←+/('ACGT' ∈ UNSC1CS)


```

COL1←SC1CA[;1]      A CONVERT GENOTYPES TO X'S FOR CORD ALLOPATRIC
IND1←COL1∈'ACGT'
COL2←SC1CA[;2]
IND2←COL2∈'ACGT'
IND3←IND1^IND2
IND4←(IND3)/1ρIND3
SC1CA←SC1CA[IND4;]

ALLELE1←SC1CA[1;1]

Q←SC1CA=ALLELE1
XCA←+/Q      A EACH GENOTYPE CONVERTED TO SUM OF NUMBER OF ALL
ELES GIVEN BY 'ALLELE'
FREQ1CA←(+/XCA)÷(2×ρXCA)  A ALLELE FREQUENCY IN CORDAT ALLOPATRIC

COL1←SC1LA[;1]      A CONVERT GENOTYPES TO X'S FOR LAC ALLOPATRIC
IND1←COL1∈'ACGT'
COL2←SC1LA[;2]
IND2←COL2∈'ACGT'
IND3←IND1^IND2
IND4←(IND3)/1ρIND3
SC1LA←SC1LA[IND4;]
Q←SC1LA=ALLELE1
XLA←+/Q      A EACH GENOTYPE CONVERTED TO SUM OF NUMBER OF ALL
ELES GIVEN BY 'ALLELE'
FREQ1LA←(+/XLA)÷(2×ρXLA)  A ALLELE FREQUENCY IN LAC ALLOPATRIC

COL1←SC1CS[;1]      A CONVERT TENOTYPES TO X'S FOR CORD SYMPATRIC
IND1←COL1∈'ACGT'
COL2←SC1CS[;2]
IND2←COL2∈'ACGT'
IND3←IND1^IND2
IND4←(IND3)/1ρIND3
SC1CS←SC1CS[IND4;]
Q←SC1CS=ALLELE1
XCS←+/Q      A EACH GENOTYPE CONVERTED TO SUM OF NUMBER OF ALL
ELES GIVEN BY 'ALLELE'
FREQ1CS←(+/XCS)÷(2×ρXCS)  A ALLELE FREQUENCY IN CORD SYMPATRIC

ALL1←,(SC1CA,[1]SC1LA),[1]SC1CS
TEST1←+/'ACGT'∈ALL1
→(TEST1≤1)/DOWN      A SKIP SNP IF NOT POLYMORPHIC

MEANX1←2×FREQ1CS

K←J
RETK:K←K+1      A SNP 2 LOOP

SC2←SC[K;;]
SC2CA←SC2[ΔCALLO3;]
SC2LA←SC2[ΔLALLO3;]
UNSC2CS←SC2CS←SC2[ΔCSYMPINDEX;]

```

```

COL1←SC2CA[;1]      A CONVERT GENOTYPES TO X'S FOR CORD ALLOPATRIC
IND1←COL1∈'ACGT'
COL2←SC2CA[;2]
IND2←COL2∈'ACGT'
IND3←IND1^IND2
IND4←(IND3)/1pIND3
SC2CA←SC2CA[IND4;]

ALLELE2←SC2CA[1;1]

Q←SC2CA=ALLELE2
XCA2←+/Q      A EACH GENOTYPE CONVERTED TO SUM OF NUMBER OF AL
LELES GIVEN BY 'ALLELE'
FREQ2CA←(+/XCA2)÷(2×pXCA2)  A ALLELE FREQUENCY IN CORDAT ALLOPATRI
C

COL1←SC2LA[;1]      A CONVERT TENOTYPES TO X'S FOR LAC ALLOPATRIC
IND1←COL1∈'ACGT'
COL2←SC2LA[;2]
IND2←COL2∈'ACGT'
IND3←IND1^IND2
IND4←(IND3)/1pIND3
SC2LA←SC2LA[IND4;]
Q←SC2LA=ALLELE2
XLA2←+/Q      A EACH GENOTYPE CONVERTED TO SUM OF NUMBER OF AL
LELES GIVEN BY 'ALLELE'
FREQ2LA←(+/XLA2)÷(2×pXLA2)  A ALLELE FREQUENCY IN LAC ALLOPATRIC

COL1←SC2CS[;1]      A CONVERT TENOTYPES TO X'S FOR CORD SYMPATRIC
IND1←COL1∈'ACGT'
COL2←SC2CS[;2]
IND2←COL2∈'ACGT'
IND3←IND1^IND2
IND4←(IND3)/1pIND3
SC2CS←SC2CS[IND4;]
Q←SC2CS=ALLELE2
XCS2←+/Q      A EACH GENOTYPE CONVERTED TO SUM OF NUMBER OF AL
LELES GIVEN BY 'ALLELE'
FREQ2CS←(+/XCS2)÷(2×pXCS2)  A ALLELE FREQUENCY IN CORD SYMPARIC

ALL1←,(SC2CA,[1]SC2LA),[1]SC2CS
TEST2←+/'ACGT'∈ALL1
→(TEST2≤1)/DOWN3      A SKIP SNP IF NOT POLYMORPHIC
TEST←+/'ACGT'∈UNSC2CS)

MEANX2←2×FREQ2CS

A CALCULATE COVAR(SNP1,SNP2)

SUM←0
COUNT←0
L←0

```

```

RETL:L←L+1
G1←UNSC1CS[L;]
G2←UNSC2CS[L;]

TEST←('★'∈(G1,G2))∨('.'∈(G1,G2))
→(TEST=1)/DOWN2

X1←+/G1=ALLELE1
X2←+/G2=ALLELE2
PROD←(X1-MEANX1)×(X2-MEANX2)
SUM←SUM+PROD
COUNT←COUNT+1

DOWN2:→(L<DIMSYP)/RETL
COV←SUM÷(COUNT-1)

‡(COV≥1.2)/'DTCLF ♦ ''COV GREATER THAN 1.2'' ♦ →0'

‡ CALC ALPHA (ald)

W←(FREQ1CA-FREQ1LA)×(FREQ2CA-FREQ2LA)
ALPHA←COV×W
‡ CALC DISTANCE

DIST←SPOS[K]-SPOS[J]

DATA←DATA,[1](ALPHA,DIST,COV,W)
DOWN3:→(K<MAXSNP)/RETK

DOWN:→(J<(MAXSNP-1))/RETJ

→(I<MAXI)/RETI

'PROGRAM ALD FINISHED. DATA IN VARIABLE ''DATA''.'
*****
COUNTALDB
‡ THIS PROGRAM TAKES DATA CREATED BY PROGRAM ALD, BINS THEM, AND CA
LCULATES ALD(D) USING ALTERNATE FORMULA
‡ BINS INCREASE IN SIZE EXPONENTIALLY
‡ CORDAT SYMPATRIC SAMPLES

MAKEBINS 1.5

BINS←ΔBINS

MAXI←1+ρBINS
ALTMEANW←0ρ0

ALDBINS←0 5ρ0
X←ΔCALDDATAB2 ‡ CHANGE FILE ACCORDINGLY
COUNTS←0ρ0

I←0
RETI:I←I+1

```

```

TEST←(I÷100)=(↑(I÷100))
⊡(TEST=1)/'I'

COUNT←IND←+/(X[;2]≤(BINS[I;2]))
COUNTS←COUNTS,IND
PART←X[⊡IND;]

X←(IND,0)⊕X
BINSIZE←BINS[I;2]-BINS[I;1]
A←(+/(PART[;1])÷(COUNT)

MEANW←(+/(PART[;4]) ÷COUNT

REDPART←(⊡,PART[;4])

IND←(REDPART≠0)/⊡REDPART
ALTM←(+/(REDPART[IND])÷IND
ALTMEANW←ALTMEANW,ALTM

MEANCOV←(+/,PART[;3])÷COUNT

MID←↑(+/(BINS[I;])÷2

ALDBINS←ALDBINS,[1](MID,A,COUNT,MEANW,MEANCOV)

DOWN:→(I<MAXI)/RETI

'PROGRAM COUNTALDA FINISHED. DATA IN VARIABLE ''ALDBINS''.'
*****
MAKEBINS X

ΔBINS←1 2ρ(1,1000)
I←0
RETI:I←I+1
START←ΔBINS[I;2]+1
END←↑(START-1)+1000×X★I
ΔBINS←ΔBINS,[1](START,END)
→(END<4500000)/RETI
*****

```

VII. Admixture and selfing rate estimates were performed using the script LIKS. It produces two matrices, ΔLNCLIK and ΔLNLLIK , that contain log-likelihood values for combinations of selfing rate and admixture proportion. Script MAKESAS2 converts these matrices into SAS form for graphics.

LIKS X

A THIS PROGRAM CALCULATES LIKELIHOOD OF ADMIXTURE PARAMETER HANIS ET AL. 1986. AM. J. PHYS. ANTHROPOLOGY

A 70: 433-441

A X IS ΔSCORES , OR $\Delta\text{SYNSCORES}$, OR $\Delta\text{NONSCORES}$, ETC.

A THIS PROGRAM IS MODIFIED TO TAKE INTO ACCOUNT SELFING RATE, S

MS← .01 .1 .2 .3 .4 .5 .6 .7 .8 .9 .99 A VALUES M (GENE FLOW) CAN TAKE ON

A MS←.21 .22 .23 .24 .25 .26 .27 .28 .29

A MS←.31 .32 .33 .34 .35 .36 .37 .38 .39

A MS←.71 .72 .73 .74 .75 .76 .77 .78 .79 .81 .82 .83 .84 .85 .86 .87 .88 .89

SS←.1 .2 .3 .4 .5 .6 .7 .8 .9 .99

A SS←.91 .92 .93 .94 .95 .96 .97 .98

A SS←.81 .82 .83 .84 .85 .86 .87 .88 .89 .91 .92 .93 .94 .95 .96 .97 .98

MAXII←ρSS

MAXJ←1↑ρMS

$\Delta\text{LNCLIK} \leftarrow \Delta\text{LNLLIK} \leftarrow (\text{MAXII}, \text{MAXJ}) \rho 0$

II←0

RETII:II←II+1 A SELFING RATE LOOP

S←SS[II]

J←0

RETJ:J←J+1 A M LOOP (GENE FLOW)

M←MS[J]

DTCLF

'STARTING MVALUE = ',M,' S VALUE = ',S

MAXI←1↑ρX

CLIK←LLIK←1 A INITIAL LIKELIHOOD

LNCLIK←LNLLIK←0 A INITIAL LOG-LIKELIHOOD

I←0

RETI:I←I+1 A SNP LOOP

GENS←X[I;;] A PICK GENOTYPES FOR SNP I

CALLO←GENS[ΔCALLO3 ;] A PICK GENOTYPES FOR ALLOPATRIC I. CORDAT

LALLO←GENS[ΔLALLO3 ;] A PICK GENOTYPES FOR ALLOPATRIC I. LAC

CSYMP←GENS[$\Delta\text{CSYMPINDEX}$;] A PICK GENOTYPES FOR SYMPATRIC I. CORDAT

LSYMP←GENS[$\Delta\text{LSYMPINDEX}$;] A PICK GENOTYPES FOR SYMPATRIC I. LAC

```

→(('★'∈,CALLO)∨('★'∈,LALLO)∨('★'∈,CSYMP)∨('★'∈,LSYMP))/DOWN1

TEST←(,CALLO),(,LALLO)
TEST2←+/'ACGT'∈TEST
→(TEST2=1)/DOWN1      a SKIP SNP IF CALLO AND LALLO FIXED FOR SAME ALL
ELE

IND←(,CALLO[;1]∈'ACGT')/11↑ρCALLO
CALLO←CALLO[IND;]      a GET RID OF '.'S AND '★'S
IND←(,LALLO[;1]∈'ACGT')/11↑ρLALLO
LALLO←LALLO[IND;]
IND←(,CSYMP[;1]∈'ACGT')/11↑ρCSYMP
CSYMP←CSYMP[IND;]
IND←(,LSYMP[;1]∈'ACGT')/11↑ρLSYMP
LSYMP←LSYMP[IND;]

ALLELE1←CALLO[1;1]

PCA←+/(ALLELE1=,CALLO)÷(ρ,CALLO) a ALLELE 1 FREQ IN I. CORD ALLOPAT
RIC
QCA←1-PCA      a ALLELE 2 FREQ IN I. CORD ALLOPATRIC
PLA←+/(ALLELE1=,LALLO)÷(ρ,LALLO)
QLA←1-PLA
PCS←+/(ALLELE1=,CSYMP)÷(ρ,CSYMP)
QCS←1-PCS
PLS←+/(ALLELE1=,LSYMP)÷(ρ,LSYMP)
QLS←1-PLS

TEST←1-(PCA-PLA)
→(TEST≥0.9)/DOWN1      a SKIP SNP IF ABSOLUTE FREQ DIFFERENCE IN ALL
OPATRY ≥ 0.9 ; CAN ALTER

MAXK←1↑ρCSYMP

K←0
RETK:K←K+1      a LOOP FOR I. CORDAT SYMP INDIVIDUALS

IND←CSYMP[K;]      a PICK KTH INDIVIDUAL FROM I. CORDAT SYMP POPULATION
TEST←+/(IND=ALLELE1

P1←(M×PCA) + ((1-M)×PLA)
Q1←1-P1
‡(Q1<.0000001)/'Q1←0'
CORRECTION←S×P1×Q1÷(2×(1-(S÷2)))

‡(TEST=2)/'P←(P1★2)+CORRECTION'
‡(TEST=1)/'P←(2×P1×Q1)-(2×CORRECTION)'
‡(TEST=0)/'P←(Q1★2)+CORRECTION'

CLIK←CLIK×P
LNCLIK←LNCLIK+⊗P

→(K<MAXK)/RETK

```

```

MAXL←1↑ρLSYMP

L←0
RETL:L←L+1    A LOOP FOR I. LAC SYMP INDIVIDUALS

IND←LSYMP[L;]  A PICK KTH INDIVIDUAL FROM I. CORDAT SYMP POPULATION
TEST←+/IND=ALLELE1

P1←(M×PCA) + ((1-M)×PLA)
Q1←1-P1
CORRECTION←S×P1×Q1÷(2×(1-(S÷2)))

⊥(TEST=2)/'P←(P1★2)+CORRECTION'
⊥(TEST=1)/'P←(2×P1×Q1)-(2×CORRECTION)'
⊥(TEST=0)/'P←(Q1★2)+CORRECTION'

LLIK←LLIK×P
LNLLIK←LNLLIK+⊗P

→(L<MAXL)/RETL

DOWN1:→(I<MAXI)/RETI

ΔLNCLIK[II;J]←LNCLIK
ΔLNLLIK[II;J]←LNLLIK

□TCLF
'LNCLIK, LNLLIK ',LNCLIK,LNLLIK

→(J<MAXJ)/RETJ

→(II<MAXII)/RETII

□TCLF
'PROGRAM LIK FINISHED. DATA IN VARIABLES ΔLNCLIK, ΔLNLLIK'
*****

MAKESAS X
A THIS PROGRAM MAKES LIST OF M, S, AND LN LIKLIHOOD FOR SAS ANALYSIS
A X IS ΔLNCLIK OR ΔLNLLIK FROM PROGRAM LIKS

A MS← .01 .1 .2 .3 .4 .5 .6 .7 .8 .9 .99    A VALUES M (GENE FLOW) C
AN TAKE ON
A MS←.21 .22 .23 .24 .25 .26 .27 .28 .29
A MS←.31 .32 .33 .34 .35 .36 .37 .38 .39

MS←.71 .72 .73 .74 .75 .76 .77 .78 .79 .81 .82 .83 .84 .85 .86 .87 .
88 .89
A SS←.1 .2 .3 .4 .5 .6 .7 .8 .9 .99
A SS←.91 .92 .93 .94 .95 .96 .97 .98
SS←.81 .82 .83 .84 .85 .86 .87 .88 .89 .91 .92 .93 .94 .95 .96 .97 .
98

```

```

MAXI←pSS
MAXJ←pMS
DATA←0 3p0

```

```

I←0
RETI:I←I+1

```

```

J←0
RETJ:J←J+1

```

```

LINE←MS[J],SS[I],X[I;J]
DATA←DATA,[I]LINE

```

```

→(J<MAXJ)/RETJ
→(I<MAXI)/RETI

```

```

DTCLF
'PROGRAM MAKESAS FINISHED.  OUTPUT IN VARIABLE ''DATA''.'

```

Log-likelihoods of selfing rate estimates for allopatric samples were calculated using script SELFEST, followed by MAKESAS2 to convert to SAS format for graphing.

SELFEST X

```

A THIS PROGRAM CALCULATES LOG-LIKELIHOODS OF ESTIMATES OF SELFING RATES IN ALLOPATRIC
A I. CORDAT AND I. LAC SAMPLES
A X IS ΔSCORES, OR ΔSYNSCORES, OR ΔNONSCORES, ETC.

```

```

pSS←.01 .1 .2 .3 .4 .5 .6 .7 .8 .9 .99
SS←.1 .2 .3 .4 .5 .6 .7 .71 .72 .73 .74 .75 .76 .77 .78 .79 .80 .81
.82 .83 .84 .85 .86 .87 .88 .89 .9 .91 .92 .93 .94 .95 .96 .97 .98 .
99
A SS←.91 .92 .93 .94 .95 .96 .97 .98
MAXII←1+pSS
ΔSELFCL←SELFCL+0p0
ΔLNCLIK←ΔLNLLIK+0p0  A INITIAL LOG-LIKELIHOOD

```

```

II←0
RETI:II←II+1  A SELFING RATE LOOP
S←SS[II]

```

```

LNCLIK←LNLLIK+0

```

```

DTCLF
'STARTING S VALUE = ',S

```

```

MAXI←1+pX

```

```

I←0
RETI:I←I+1  A SNP LOOP

```

```

GENS←X[I;;]  A PICK GENOTYPES FOR SNP I

```



```

CALLO←GENS[ΔCALLOCLOSE;]  A PICK GENOTYPES FOR ALLOPATRIC I. CORDAT
LALLOSAVE←LALLO←GENS[ΔLALLOCLOSE;]  A PICK GENOTYPES FOR ALLOPATRIC
I. LAC

IND←(,CALLO[;1]∈'ACGT')/11↑ρCALLO
CALLO←CALLO[IND;]  A GET RID OF '.'S AND '★'S
IND←(,LALLO[;1]∈'ACGT')/11↑ρLALLO
LALLO←LALLO[IND;]

TEST1←1↑ρCALLO

⊡(TEST1=0)/'SKIPC←'Y''
⊡(TEST1>0)/'SKIPC←'N''  ◇ ALLELE1C←CALLO[1;1]  ◇ PCA←+/(ALLELE1C=,CAL
LO)÷(ρ,CALLO)  ◇ QCA←1-PCA  '
A ALLELE 1 FREQ IN I. CORD ALLOPATRIC
A ALLELE 2 FREQ IN I. CORD ALLOPATRIC

TEST2←1↑ρLALLO
⊡(TEST2=0)/'SKIPL←'Y''
⊡(TEST2>0)/'SKIPL←'N''  ◇ ALLELE1L←LALLO[1;1]  ◇ PLA←+/(ALLELE1L=,LA
LLO)÷(ρ,LALLO)  ◇ QLA←1-PLA  '

MAXK←1↑ρCALLO

→(SKIPC='Y')/DOWN2  A SKIP SNP IF NO INDIVIDUALS

TEST←+/'ACGT'∈,CALLO
→(TEST=1)/DOWN2  A SKIP SNP IF NO VARIATION

K←0
RETK:K←K+1  A LOOP FOR I. CORDAT ALLO INDIVIDUALS

IND←CALLO[K;]  A PICK KTH INDIVIDUAL FROM I. CORDAT SYMP POPULATION

TEST←+ /IND=ALLELE1C
CORRECTION←S×PCA×QCA÷(2×(1-(S÷2)))

⊡(TEST=2)/'P←(PCA★2)+CORRECTION'
⊡(TEST=1)/'P←(2×PCA×QCA)-(2×CORRECTION)'
⊡(TEST=0)/'P←(QCA★2)+CORRECTION'

LNCLIK←LNCLIK+⊗P

→(K<MAXK)/RETK

DOWN2:MAXL←1↑ρLALLO

→(SKIPL='Y')/DOWN1

TEST←+/'ACGT'∈,LALLO
→(TEST=1)/DOWN1

L←0
RETL:L←L+1  A LOOP FOR I. LAC ALLO INDIVIDUALS

```

IND←LALLO[L;] A PICK KTH INDIVIDUAL FROM I. CORDAT SYMP POPULATION

TEST←+/IND=ALLELE1L

CORRECTION←S×PLA×QLA÷(2×(1-(S÷2)))

Δ(TEST=2)/'P←(PLA★2)+CORRECTION'

Δ(TEST=1)/'P←(2×PLA×QLA)-(2×CORRECTION)'

Δ(TEST=0)/'P←(QLA★2)+CORRECTION'

LNLLIK←LNLLIK+⊗P

→(L<MAXL)/RETL

DOWN1:→(I<MAXI)/RETI

ΔLNCLIK←ΔLNCLIK, LNCLIK

ΔLNLLIK←ΔLNLLIK, LNLLIK

□TCLF

'LNCLIK, LNLLIK ', LNCLIK, LNLLIK

→(II<MAXII)/RETII

□TCLF

'PROGRAM SELFEST FINISHED. DATA IN VARIABLES ΔLNCLIK, ΔLNLLIK'

MAKESAS2

A THIS TAKES ΔLNCLIK AND ΔLNLLIK FROM PROGRAM 'SELFEST' AND MAKES DATA FOR SAS

A VARIABLE SS WAS JUST FORMED BY 'SELFEST'

MAXI←ρSS

DATA←0 3ρ0

I←0

RETI:I←I+1

S←SS[I]

LINE←S, ΔLNCLIK[I], ΔLNLLIK[I]

DATA←DATA, [1]LINE

→(I<MAXI)/RETI

'PROGRAM MAKESAS2 FINISHED. OUTPUT IN VARIABLE ''DATA''.'

VIII. Simulating gene flow. The script BREAKPOS divides all contigs into blocks of size 10 kb. This is run first. It produces the variable Δ BLOCKS, each row of which contains the first and last position of block's SNPs in Δ SCORES.

The script SWAPFIT3 swaps different proportions of alleles in 100 kb contig segments from allopatric *I. lacunosa* into allopatric *I. cordatotriloba* samples to create simulated *I. cordatotriloba* sympatric samples. For each proportion, it evaluates the sum of squared differences between the simulated allele frequencies and observed allele frequencies in the actual *I. cordatotriloba* sympatric samples. The program also calculates mean of sympatric between-species π across 20 replicate simulations for each proportion. The proportion with the value of π closest to the observed value is taken as the best estimate of the true proportion.

BREAKPOS

```

A THIS DIVIDES ALL CONTIGS INTO BLOCKS OF SIZE BLOCKSIZE KB OR LESS
BLOCKSIZE= 100000  A SIZE OF BLOCK IN KILOBASES
 $\Delta$ BLOCKS=0 2p0      A WILL HOLD START AND END INDEX OF BLOCKS
DIM= $\rho$  $\Delta$ UNIQUETIGS    A NUMBER OF UNIQUE CONTIGS
I=0
RETI:I=I+1
CONTIGNUMS=( $\Delta$ UNIQUETIGS[I]= $\Delta$ TIGS)/1 $\rho$  $\Delta$ TIGS  A PICK LOCUS NUMBERS F
OR CONTIG I
POS= $\Delta$ POS[CONTIGNUMS]  A PICK POSITIONS ON CO
NTING

RET:STARTNUM=CONTIGNUMS[1]  A THIS LOOOP DIVIDES CO
NTIG INTO 100 KB BLOCKS
STARTPOS=POS[1]  A FIRST CONTIG POSITION
OF BLOCK

IND3=POS<(STARTPOS+BLOCKSIZE)  A PICK OUT ALL POSIT
IONS WITHIN BLOCKSIZE
NUM=+/IND3  A NUM IS LIST OF PO
SITIONS IN BLOCK

BLOCKNUMS=NUM+CONTIGNUMS  A TAKE LOCUS NUMBER
S FROM CONTIGNUMS
BLOCK=(1+BLOCKNUMS), (-1+BLOCKNUMS)  A DEFINE FIRST AND
LAST LOCI NUMS FOR BLOCK
 $\Delta$ BLOCKS= $\Delta$ BLOCKS,[1]BLOCK  A ADD BLOCK TO BL
OCKS
CONTIGNUMS=NUM+CONTIGNUMS  A DROP LOCI NUMBER
S FROM CONTIGNUMS
POS=NUM+POS  A DROP POSITIONS F
ROM  $\Delta$ POS
 $\rightarrow$ (0< $\rho$ CONTIGNUMS)/RET
 $\rightarrow$ (I<DIM)/RETI
*****

```

Program SEPScores creates an index, Δ FIXEDSCORES that indicates whether a SNP is highly diverged or less diverged in allopatry for either known or close allopatric sites

SEPScores X

```

A THIS FUNCTION IDENTIFIES SNPS THAT ARE HIGHLY DIVERGENT (FREQ DIFF
  IN ALLOPATRY  $\geq$  0.9) OR LESS DIVERGENT (FREQ DIFF IN ALLOPATRY  $<$  0.9
)
A IT CREATES A VARIABLE  $\Delta$ FIXEDSCORES THAT HAS 2 ROWS AND 66,720 COLU
MNS, ONE COLUMN FOR EACH SNP
A A 1 IN FIRST ROW INDICATES THAT SNP IS HIGHLY DIVERGENT FOR KNOWN
  ALLOPATRIC SAMPLES
A A 1 IN SECOND ROW INDICATES THAT SNP IS HIGHLY DIVERGENT FOR CLOS
E ALLOPATRIC SAMPLES
A A 0 IN EITHER ROW INDICATES SNP IS LESS DIVERGENT
A TYPICALLY X WILL BE  $\Delta$ SCORES, BUT IT COULD BE A SUBSET OF  $\Delta$ SCORE
S (E.G. SYNONYMOUS SNPS, NON-SYNONYMOUS SNPS, ETC.)

```

Δ FIXEDSCORES \leftarrow 2 0p0

MAXI \leftarrow 1+pX

I \leftarrow 0 A SNP LOOP

RETI:I \leftarrow I+1

SCORES \leftarrow X[I;;]

CALLOK \leftarrow ,SCORES[Δ CALLO3;] A PICK OUT GENOTYPES CORRESPONDING
TO I. CORDAT KNOWN ALLOPATRIC SITES

LALLOK \leftarrow ,SCORES[Δ LALLO3;] A PICK OUT GENOTYPES CORRESPONDING
TO I. LAC KNOWN ALLOPATRIC SITES

CALLOK \leftarrow ,SCORES[Δ CALLOCLOSE;] A DITTO FOR CLOSE ALLOPATRIC SITES
LALLOK \leftarrow ,SCORES[Δ LALLOCLOSE;] A DITTO

CALLOK1 \leftarrow (CALLOK \in 'ACGT')/CALLOK A ELILMINATE MISSING DATA
LALLOK1 \leftarrow (LALLOK \in 'ACGT')/LALLOK
CALLOK1 \leftarrow (CALLOK \in 'ACGT')/CALLOK
LALLOK1 \leftarrow (LALLOK \in 'ACGT')/LALLOK

ALLELE \leftarrow 1+CALLOK1 A DESIGNATE FOCAL ALLELE FOR FREQUE
NCY CALCULATIONS

DIFF1 \leftarrow 1((+/CALLOK1=ALLELE) \div (pCALLOK1))-((+/LALLOK1=ALLELE) \div (pLALLOK1
)) A KNOWN ALLOPATRIC FREQUENCY DIFFERENCE

DIFF2 \leftarrow 1((+/CALLOK1=ALLELE) \div (pCALLOK1))-((+/LALLOK1=ALLELE) \div (pLALLOK1
)) A CLOSE ALLOPATRIC FREQUENCY DIFFERENCE

ADD \leftarrow 2 1p0 A INITIALIZE VARIABLE TO HOLD 1 OR 0 VALUES

$\&$ (DIFF1 \geq 0.9)/'ADD[1;1] \leftarrow 1' A IF KNOWN ALLOPATRIC DIFFERENCE \geq 0.
9, PLACE 1 IN FIRST ROW OF 'ADD'

```

1(DIFF2≥0.9)/'ADD[2;1]←1'      A IF CLOSE ALLOPATRIC DIFFERENCE ≥0.
9, PLACE 1 IN SECOND ROW OF 'ADD'

```

```

ΔFIXEDSCORES←ΔFIXEDSCORES,ADD  A APPEND RESULTS TO ΔFIXEDSCORES

```

```

→(I<MAXI)/RETI

```

```

□TCLF

```

```

'PROGRAM SEPCORES FINISHED. DATA IN VARIABLE ''ΔFIXEDSCORES''.'

```

```

*****

```

Program MIXPI. This program simulates gene flow by replacing a proportion of the genotypes in the "newly sympatric" *I. cordatotriloba* population with randomly chosen genotypes from the newly sympatric *I. lacunosa* population. Data saved in variable OUTPUT.

```

MIXPI

```

```

A THIS FUNCTION SIMILATES GENE FLOW AND CALCULATES RESULTING PI FOR
DIFFERENT VALUES OF PROPORTIONS

```

```

A ADMIXTURE (PROPORTION OF ALLELES REPLACED IN ''NEWLY SYMPATRIC''
POPULATION).

```

```

PROPNUMS←.45 .46 .47 .48 .49 .51 .52 .53 .54 .55      A ADMIXTURE PR
OPORTIONS TO USE; CAN CHANGE ACCORDINGLY
OUTPUT←0 1200

```

```

MAXII←ρPROPNUMS

```

```

II←0      A PROPORTION LOOP

```

```

RETI:II←II+1

```

```

PROP←PROPNUMS[III]  A PICK ADMIXTURE PROPORTION

```

```

MAXJ←20

```

```

J←0

```

```

RETJ:J←J+1  A REPLICATE LOOP FOR A GIVEN ADMIXTURE PROPORTION

```

```

II,J,PROP

```

```

SWAPTIG3MA  A PRODUCES PSEUDOSCORES BY SWAPPING GENOTYPES FROM ALLO
PATRIC TO SYMPATRIC POPULATION

```

```

PICALC2B PSEUDOSCORES  A CALCULATES ΔPI

```

```

PICONTRAST2      A CALCULATES AVERAGE PI BETWEEN SYMPATRIC SA
MPLES (CONTAINED IN ΔOBSPIVECTOR)

```

```

OUTPUT←OUTPUT,[1](PROP,ΔOBSPIVECTOR)  A CATENATES ΔOBSPIVECTOR TO O
UTPUT

```

```

→(J<MAXJ)/RETJ

```

```

→(II<MAXII)/RETI

```

```

□TCLF

```

'PROGRAM MIXPI FINISHED. DATA IN VARIABLE ''OUTPUT''.'

Program SWAPTIG3MA. This program is called by program MIXPI. It performs the genome swapping so simulate gene flow from *I. lacunosa* to *I. cordatotriloba* in sympatry. It produces a variable PSEUDOSCORES

SWAPTIG3MA

A THIS FUNCTION SIMULATES GENE FLOW FROM I. LACUNOSA INTO I. CORDATO TRILOBA IN SYMPATRY.

A IT FIRST CREATES 'NEW' SYMPATRIC POPULATION BY SUBSTITUTING RANDOM DRAWS FROM CORD ALLO AND LAC ALLO

A INTO CORD SYMP AND LAC SYMP. SUBSTITUTIONS ARE BY BLOCK

A IT THE SUBSTITUTES A RANDOMLY CHOSEN NUMBER OF LAC ALLO INDIVIDUALS INTO THE NEW CORD SYMPATRIC

A POPULATION, REPLACING A RANDOMLY CHOSEN SET OF INDIVIDUALS.

A THE NUMBER OF INDIVIDUALS SUBSTITUTED IS DETERMINED BY 'PROP'

PSEUDOSCORES←0 62 2ρ' ' A NEW VARIABLE TO HOLD DATA

DIM1←1+ρΔSCORES

BLOCKDIM←1+ρΔBLOCKS A NUMBER OF BLOCKS

I←0 A BLOCK LOOP

RETI:I←I+1

BLOCK←ΔBLOCKS[I;] A CHOOSE ENDPOINTS OF BLOCK I

LOCNUMS←(BLOCK[1]-1)+1(1+BLOCK[2]-BLOCK[1]) A CALCULATE LOCI NUMBERS

PARTSCORES←ΔSCORES[LOCNUMS;;]

A NEXT PART SUBSTITUTES L ALLO INTO L SYMP

LALLOSCORES←PARTSCORES[;ΔLALLO3;] A NOTE CAN CHANGE TO ΔLALLOCLOSE

RAND←?13ρ16 A INDEX FOR RANDOMLY CHOSEN LAC ALLO INDIVIDUALS

PARTSCORES[;ΔLSYMPINDEX;]←LALLOSCORES[;RAND;] A MAKE SUBSTITUTION

A NEXT PART SUBSTITUTES C ALLO INTO C SYMP

CALLOSCORES←PARTSCORES[;ΔCALLO3;]

RAND2←?11ρ14

PARTSCORES[;ΔCSYMPINDEX;]←CALLOSCORES[;RAND2;]

A NEXT PART SUBSTITUTES RANDOM INDIVIDUAL FROM LALLO INTO CSYMP

NUMSUB←1/((1((PROP×(ρΔCSYMPINDEX))+((?10000)÷10000))), (ρΔCSYMPINDEX)) A THIS IS NUMBER OF INDIVIDUALS IN 'NEW' LSYMP TO CHOOSE

RAND3←?NUMSUBρ13 A RANDOM INDEX FOR PICKING SUBS

SUBS←PARTSCORES[;ΔLSYMPINDEX[RAND3];] A PICK SUBS FROM NEW LSYMP

RAND4←NUMSUB?11 A RANDOM INDEX FOR CSYMP INDIVIDUALS TO BE REPLACED

PARTSCORES[;ΔCSYMPINDEX[RAND4];]←SUBS

```
PSEUDOSCORES←PSEUDOSCORES,[1]PARTSCORES  A ADD TO PSEUDOSCORES
```

```
→(I<BLOCKDIM)/RETI
```

```
IND←(~ΔFIXEDSCORES[2;])/11↑PSEUDOSCORES  A PICKS OUT FIXED AND  
NEARLY FIXED SCORES
```

```
PSEUDOSCORES←PSEUDOSCORES[IND;;]
```

```
*****
```

Program PICALC2B. This program is called by program MIXPI. It calculates pairwise between-species PI* for sympatric samples and saves them in variable ΔPI2

```
PICALC2B X;I
```

```
A THIS PROGRAM CALCULATES PAIRWISE PI VALUES FOR EITHER THE FULL DAT  
ASET (ALL SNPS)
```

```
A OR FOR SUBSETS OF DATA (E.G. SYNONYMOUS, NON-SYNONYMOUS SITES)
```

```
A IT PRODUCES A MATRIX ΔPI2 THAT HAS THE AVERAGE PAIRWISE PI VALUES  
FOR EACH PAIR OF SAMPLES
```

```
A PROGRAM READS IN VARIABLE X (PSEUDOSCORES) CREATED BY SWAPTIG3MA
```

```
PICUM←62 62P0  A SET UP PI MATRIX- WILL EVENTUALLY HAVE NUMBER OF P  
AIRWISE DIFFERENCES ACROSS ALL VARIABLE SNPS
```

```
PICOUNT←62 62P0  A SET UP MATRIX FOR CUMULATIVE COUNTS (EXCLUDES MI  
SSING VALUES)
```

```
DIM←1↑P X
```

```
I←0
```

```
RETI:I←I+1  A SNP LOOP
```

```
A PICK DATA FOR NEXT SNP--STRIP OFF LEADING INFO AND LEAVE JUST NUC  
LEOTIDES FOR DIFFERENT
```

```
A SAMPLES. TEMP5 HAS FORMAT G/G G/G G/C . . .
```

```
SCORES3←X[I;;]  A NOTE: MODIFY WHICH ΔSCORES TO USE (E.G.  
ALL, SYNONYMOUS, ETC.)
```

```
A CALCULATE PAIRWISE PI VALUES FUR CURRENT SNP
```

```
TEMP7←SCORES3*.=SCORES3
```

```
TEMP8←+/[2]TEMP7
```

```
TEMP9←(+/[3]TEMP8)÷4
```

```
TEMP9←1-TEMP9
```

```
TIND←(SCORES3[;1]='. ')/162
```

```
COUNTMAT←62 62P1
```

```
COUNTMAT[TIND;]←0
```

```
COUNTMAT[;TIND]←0
```

```
PI←TEMP9×COUNTMAT
```

```
A ADD CURRENT PAIRWISE PI VALUES TO CUMULATIVE VALUES
```

```
PICUM←PICUM+PI
```

```
PICOUNT←PICOUNT+COUNTMAT
```

```
→(I<DIM)/RETI
```

```
MAXPICOOUNT←( / , , PICOOUNT  
CORRCOUNT←MAXPICOOUNT÷PICOOUNT  
# ΔPI2←PICUM×CORRCOUNT÷DIVIDER  
# ΔPI2←PICUM÷DIVIDER  
ΔPI2←PICUM÷PICOOUNT
```

```
□TCLF
```

```
'PROGRAM PICALC2 FINISHED.'  
'PAIRWISE PI VALUES IN VARIABLE ''ΔPI2''.'
```

```
*****
```

Program PICONTRAST2. This program is called by MIXPI and calculates *average* between-species sympatric PI (PI*) values and saves them in ΔOBSPIVECTOR.

```
PICONTRAST2
```

```
# THIS PROGRAM CALCULATES PI WITHIN AND BETWEEN SYMP AND ALLOPATRIC  
SAMPLES
```

```
# IT USES THE ΔPI2 MATRIX CALCULATED BY 'PICALC2'
```

```
MASK4←62 62#1
```

```
I←0
```

```
RETI:I←I+1
```

```
MASK4[I;I]←0
```

```
→(I<62)/RETI
```

```
#PI←MASK4×ΔPI×(1÷#ΔSCORES2)÷ΔTOTNUC # THIS STEP CONVERTS THE PI VAL  
UES CALCULATED BY 'PICALC2'
```

```
# TO TRUE PI VALUES
```

```
PI←MASK4×ΔPI2
```

```
PIBETWALLO←PI[ΔCALLO3;ΔLALLO3]
```

```
PIBETWALLOCLOSE←PI[ΔCALLOCLOSE;ΔLALLOCLOSE]
```

```
PIBETWSYMP←PI[ΔCSYMPINDEX;ΔLSYMPINDEX]
```

```
PICBETWALLOSYP←PI[ΔCALLO3;ΔCSYMPINDEX]
```

```
PILBETWALLOSYP←PI[ΔLALLO3;ΔLSYMPINDEX]
```

```
PICBETWALLOCLOSESYMP←PI[ΔCALLOCLOSE;ΔCSYMPINDEX]
```

```
PILBETWALLOCLOSESYMP←PI[ΔLALLOCLOSE;ΔLSYMPINDEX]
```

```
AVEBETWSYMP←(+//PIBETWSYMP)÷((#ΔCSYMPINDEX)×(#ΔLSYMPINDEX))
```

```
AVEBETWALLO←(+//PIBETWALLO)÷((#ΔCALLO3)×(#ΔLALLO3))
```

```
AVEBETWALLOCLOSE←(+//PIBETWALLOCLOSE)÷((#ΔCALLOCLOSE)×(#ΔLALLOCLOSE))
```

```
AVECBETWALLOSYP←(+//PICBETWALLOSYP)÷((#ΔCALLO3)×(#ΔCSYMPINDEX))
```

```
AVELBETWALLOSYP←(+//PILBETWALLOSYP)÷((#ΔLALLO3)×(#ΔLSYMPINDEX))
```

```
AVECBETWALLOCLOSESYMP←(+//PICBETWALLOCLOSESYMP)÷((#ΔCALLOCLOSE)×(#ΔCSYMPINDEX))
```



```

AELBETWALLOCLOSESYMP←(+/+/PILBETWALLOCLOSESYMP)÷((ρΔLALLOCLOSE)×(ρΔ
LSYMPINDEX))

```

```

ΔBETWDIFF1←AVEBETWALLO-AVEBETWSYMP
ΔBETWDIFF2←AVEBETWALLOCLOSE-AVEBETWSYMP
ΔBETWDIFF3←AVECBETWALLOSYMP-AVELBETWALLOSYMP
ΔBETWDIFF4←AVECBETWALLOCLOSESYMP-AVELBETWALLOCLOSESYMP

```

```

ΔOBSPIVECTOR←AVEBETWALLO,AVEBETWALLOCLOSE,AVEBETWSYMP,ΔBETWDIFF1,ΔBE
TWDIFF2,AVECBETWALLOSYMP,AELBETWALLOSYMP,AVECBETWALLOCLOSESYMP,AEL
BETWALLOCLOSESYMP,ΔBETWDIFF3,ΔBETWDIFF4  A OBSERVED VALUES

```

```

□TCLF

```

```

'BETWEEN SPECIES COMPARISONS'

```

```

□TCLF

```

```

'a. AVERAGE PI BETWEEN KNOWN ALLO: ',AVEBETWALLO
'b. AVERAGE PI BETWEEN CLOSE ALLO: ',AVEBETWALLOCLOSE
'c. AVERAGE PI BETWEEN SYMP: ',AVEBETWSYMP
'd. DIFFERENCE a. - c.: ',(ΔBETWDIFF1)
'e. DIFFERENCE b. - c.: ',(ΔBETWDIFF2)

```

```

□TCLF

```

```

'WITHIN SPECIES COMPARISONS'

```

```

□TCLF

```

```

'f. I. CORD BETW ALLO(KNOWN) AND SYMP ',(AVECBETWALLOSYMP)
'g. I. LAC BETW ALLO(KNOWN) AND SYMP ',(AELBETWALLOSYMP)
'h. I. CORD BETW ALLO(CLOSE) AND SYMP ',(AVECBETWALLOCLOSESYMP)
'i. I. LAC BETW ALLO(CLOSE) AND SYMP ',(AELBETWALLOCLOSESYMP)

'j. DIFFERENCE f. - g.: ',(ΔBETWDIFF3)
'k. DIFFERENCE h. - i.: ',(ΔBETWDIFF4)

```

```

□TCLF

```

```

'PROGRAM PICONTRAST2 FINISHED'

```

```

*****

```

Program SUMMARIZEPI. This program takes variable OUTPUT produced by program MIXPI and calculates mean PI* and its standard deviation for each value of admixture proportion.

```

SUMMARIZEPI X

```

```

A THIS FUNCTION READS IN X, WHICH IS VARIABLE "OUTPUT" FROM MIXPI.
A IT PRODUCES VARIABLE DATA, WHICH HAS 3 COLUMNS:

```

```

A      ADMIXTURE PROP      MEAN PI*      SD PI*

```

```

A EACH ROW CORESPONDS TO A DIFFERENT ADMIXTURE PROPORTION

```

```

DATA←0 3ρ0
ROWS←1↑ρX
CATS←ROWS÷20

```

```

I←0

```

```

RETI:I←I+1

```

```

IND←((I-1)×20)+120

```

```

PART←X[IND;]
PROP←PART[1;1]
B←,PART[;4]
MEAN←+/B÷20
VAR←(+/(B-MEAN)*2)÷19
DATA←DATA,[1](PROP,MEAN,(2×VAR*.5))
→(I<CATS)/RETI

```

□TCLF

'PROGRAM SUMMARIZEPI FINISHED. DATA IN VARIABLE 'DATA'.'

Program PISIMBOOT. This program bootstraps between-species sympatric PI values of observed data to estimate 95% credible intervals for PI for either highly-divergent or less-divergent SNPs. Bootstrap PI values accumulated in variable ΔPIBETW. Mean and credible interval calculated manually from this variable.

PISIMBOOT

⌘ THIS PROGRAM BOOTSTRAPS THE OBSERVED BETWEEN-SPECIES SYMPATRIC PI SCORES FOR
 ⌘ EITHER HIGHLY DIVERGED OR LESS-DIVERGED SNPS

SCORES←ΔSCORES ASSIGN SNP GENOTYPES TO VARIABLE 'SCORES'

ΔPIBETW←0␣0 ⌘ INITIALIZE VARIABLE TO CONTAIN PI VALUES FOR DIFFERENT SNPS

IND←(~ΔFIXEDSCORES[2;])/11␣PScores ⌘ INDEX FOR FIXED AND NEARLY FIXED SCORES
 ⌘ OR LESS DIVERGENT SCORES (DEPENDS IF '~' PRESENT BEFORE 'ΔFIXEDSCORES')
 ⌘ ROW 1 CORRESPONDS TO KNOWN ALLOPATRIC SAMPLES, ROW 2 TO CLOSE ALLOPATRIC SAMPLES

SCORES←SCORES[IND;;] ⌘ PICKS OUT HIGHLY DIVERGENT OR LESS DIVERGENT SCORES (GENOTYPES)

CSYMP←ΔCSYMPINDEX ⌘ INDEX OF CORDAT SYMPATRIC SAMPLES
 LSYMP←ΔLSYMPINDEX ⌘ INDEX OF LAC SYMPATRIC SAMPLES

II←0
 RETII:II←II+1 ⌘ SNP LOOP
 TEST←(II÷100)=(LII÷100)
 ⌘(TEST=1)/'II'

TSCORES←SCORES ⌘ BOOTSTRAP LOOP-- BOOTSTRAPPING OVER SAMPLES
 IND1←CSYMP[?(␣CSYMP)␣␣CSYMP] ⌘ INDEX FOR CORDAT SAMPLES PICKED RANDOMLY WITH REPLACEMENT
 TSCORES[;CSYMP;]←SCORES[;IND1;] ⌘ REPLACE OBS SAMPLE GENOTYPES

IND2←LSYMP[?(␣LSYMP)␣␣LSYMP] ⌘ DITTO FOR LAC SAMPLES

```

TSCORES[;LSYMP;]←SCORES[;IND2;]

PICALC2A TSCORES      # PRODUCES ΔPI2

PICONTRAST2B          # CALCULATES AVERAGE BETWEEN-SPECIES PI VALUE
S

PIBETW←ΔOBSPIVECTOR[3]

ΔPIBETW←ΔPIBETW,PIBETW      # APEND PI VALUE TO VECTOR OF PI VALUES.

→(II<1000)/RETII

# MANUALLY COMPUTE MEAN AND 95 PERCENT CREDIBLE INTERVAL FOR PI

```