

Numerous software programs are available for assessing genetic diversity. Most are freely available through Internet. Many perform similar tasks, with the main differences being in the user interface, type of data input and output, and platform. Thus, choosing which to use depends heavily on individual preferences. In this section, we describe some of the programs available, noting specific options that users may find preferable.

Peature Diversity Heterozygosity (observed) Expected heterozygosity No alleles per locus Frective alleles (no.) Percentage polymorphic loci Shannon-Weaver Population structure = statistics	TFPGA <sup>1</sup> x x x	Arlequin <sup>1</sup> * x x x x	GDA <sup>1</sup> x x x	Program GENEPOP <sup>1</sup> *	GeneStrut x x	POPGENE*+ x
Heterozygosity (observed)       Expected heterozygosity       No alleles per locus       Effective alleles (no.)       Percentage polymorphic loci       Shannon-Weaver       Population structure	x x	x x x x x	x x		111x000	H 01x 000
Heterozygosity (observed)       Expected heterozygosity       No alleles per locus       Effective alleles (no.)       Percentage polymorphic loci       Shannon-Weaver       Population structure	×	x x x	x	00110110		
Expected heterozygosity Vo alleles per locus Effective alleles (no.) Percentage polymorphic loci Shannon-Weaver Population structure	00000	x x		00110110	101 Y 000	
Effective alleles (no.) Percentage polymorphic loci Shannon-Weaver Population structure	00000	×	x			x
Percentage polymorphic loci Shannon-Weaver Population structure	×				x	x
Shannon-Weaver Population structure	01010		11001		×	×
Population structure	01010	×	x	00440440	0444404	×
		100000	1100	00110110	UT110	1010100
	01101	010100	10000	10110100	011111	1000101
- statistics			1.737.1737			
	x	×	x	×	×	×
G-statistics	01001	×	×	00011001	×	×
Rho-statistics	10110	×	*	×	111000	14 04 0004
Homogeneity	x	-	1100	×	TH TUUU	x
Vigration	01001	×	11001	×	011111	x
solation-by-distance	i nana	oppose a	1 ana	×	a 4 a mm m	100000000
	TOOM	TOOD LITE	1000		111000	0101010101
Equilibrium	01001	000011	1000/	00110110	121000	0440004
Hardy-Weinberg	x	x	x	x	x	x
Two-locus	10100	x	x	x	x	x
Aultilocus	001000	400446	x	001100440	111000	1010100
J-test	00100	100110	1.1.1.92	×	1110000	1010100
	10110	1100110	1000	00110110	11510000	21010001
Genetic distance Nei's	01110	1011	0100	00000000	000000	1010011
Rogers'	x	×	×		x	x
Pairwise Fsr	×	×	×	20111110	*	H110001
	00001	- î	00111	00011000	011010	1101001
Clustering						
Neighbour-joining	11011	11001	×	00011000	011010	/101110
JPGMA	x	000100	x	100011001	×	x
	00000	×				
Neutrality test					WELL PROPERTY AND	x

Joanne Labate (2000) wrote an excellent review of six programs: TFPGA (Miller, 1997), Arlequin (Schneider et al., 1997), GDA (Lewis and Zaykin, 1999), GENEPOP (Raymond and Rousset, 1995), GeneStrut (Constantine et al., 1994), and POPGENE (Yeh et al., 1997). Her review includes the particular options of each program, a table of functions available in each, and Web sites where they can be downloaded. To avoid redundancy, we have included only the Arlequin, which is possibly the most widely used program of the six.

For full references to these six programs and selected others, see Appendix 9.

# Reference

Labate, J.A. 2000. Software for population genetic analyses of molecular marker data. Crop Sci. 40:1521-1528.

	sites,	programs, their authors,
10000110	0001101001101101	101010100010100100001100001101001101101
Name	Author	Available from:
Arlequina	Laurent Excoffier	http://lgb.unige.ch/arlequin
DnaSP	Julio and Ricardo Rozas	http://www.ub.es/dnasp
PowerMarker	Kejun (Jack) Liu	http://www.powermarker.net/
MEGA2	S. Kumar and others	http://www.megasoftware.net
PAUP*	David Swofford	http://paup.csit.fsu.edu/
TFPGAª	Mark Miller	http://bioweb.usu.edu/mpmbio/index.htm
GDAª	Paul Lewis, Dmitri Zaykin	http://lewis.eeb.uconn.edu/lewishome/software.html
GENEPOP <sup>a</sup>	Michel Raymond, Francois Rousset	ftp://ttp.cefe.cnrs-mop.fr/pub/PC/MSDOS/GENEPOP/Genepop.zip also at <http: genepop="" wbiomed.curtin.edu.au=""></http:>
NTSYSpc	F.J. Rohlf	http://www.exetersoftware.com/cat/ntsyspc/ntsyspc.html
structure	Jonathan K. Pritchard	http://pritch.bsd.uchicago.edu/
GeneStruta	Constantine, Hobbs & Lymbery	http://wwwvet.murdoch.edu.au/vetschl/imgad/GenStrut.htm
POPGENE <sup>a</sup>	F.C. Yeh, RC. Yang, T. Boyle	http://www.ualberta.ca/-fyeh/index.htm
MacClade	David R. & Wayne P. Maddison	http://phylogeny.arizona.edu/macclade
PHYLIP	Joe Felsentein	http://evolution.genetics.washington.edu/phylip.html
SITES	Jody Hey	http://lifesci.rutgers.edu/~heylab/ProgramsandData/Programs/SITES/SITES_Documentation.htm#Contents
CLUSTAL W	Thompson, Higgins & Gibson	http://www.ebi.ac.uk/clustalw
MALIGN	D. Janies and W.C. Wheeler	http://research.amnh.org/users/djanies/
Discussed in Laba	ate (2000)	

We review other programs, selecting them for their wide use and giving priority to those that are available for no charge (except for PAUP\*). Web sites listing and linking many other available programs are also given in this and following slides. We include information on authors, costs, platform specificities and Web sites. Note that, while these programs are sometimes made specifically for only one platform (usually Windows or Macintosh), with the recent advent of 'emulators' (such as SoftWindows, VirtualPC), most programs can be run on any computer, regardless of platform. Although we note the cases where a program has been successfully used with one of these emulators, we do not say that all listed programs can be used across platforms; simply that we know for sure that these have been successfully used. Sometimes, using emulators can cause the program to run more slowly or create other problems. Where possible, prefer using the platform for which the program was designed.

Programs listed in **bold face** in the slide are discussed in this paper. Web sites were available as of 28 February, 2003.

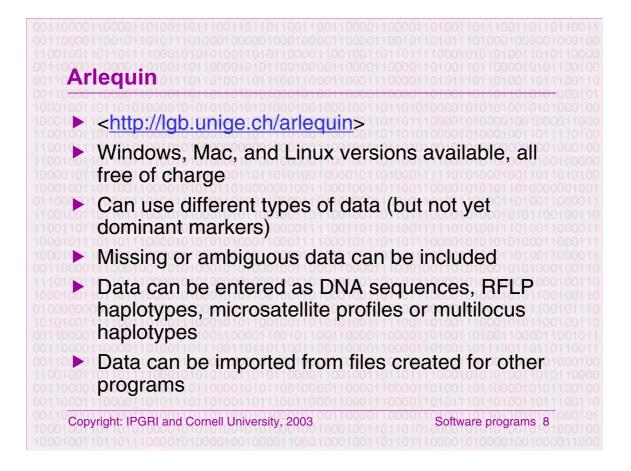
	0101010101010101	31000101001001100	0110110100	0101010101001010
0110 Name 1000	001000010010			
	Windows	Macintosh	Other	Cost (US\$)
Arlequin	X	Х	Х	0
DnaSP	X	C	11010001111	0
PowerMarker	X	0001001100010	111011001100	010101010100
MEGA2	X	1000010011000	CONTRACTOR	1110100101001
Arlequin	010100 X 01011	010011X110010	X	100
TFPGA	010100 X 10100	1000011110011	011011001111	1010100-0 101001
GDA	010101 X 01000	1000111100010	1101 X 1100	001010100010001
GENEPOP	0011000110100	1000011100010	X1100	00110000101001
NTSYSpc	10100 X 01010	01001100011000	001110001001	230-300
structure	X 00010	01001101110100	X	01010100000100
GeneStrut	0101000101011	X	011011011100	01010000010110
POPGENE	X	010111 <b>c</b> 010000	000001000000	01010010101110
MacClade	01110001010111	010011X101010	011111001101	125
PHYLIP	0 1 1 0 1 0 X 1 0 0 0 0	100101 <b>X</b> 001100	001100 X 011010	001101000100001
SITES	X	011001X001100	0011000011010	0 101 0
CLUSTAL W	11010 X 10000	10001 X 00100	Х	11010010000
MALIGN	0101010100110	10110000110010	Х	0

c = program runs well with an emulator such as SoftWindows or VirtualPC.

To save space, references to the programs discussed are given in Appendix 9.

Copyright: IPGRI and Cornell Ur	niversity 2003	Software programs 7
IVILUA		
▶ MEGA		
000000100000000101001110		
PAUP*		
DnaSP		
001101101110000101000101		
PowerMarker		
Arlequin		
011101101010010101000000		
000110000110100110110110	1010100010100110000	11000011010011011011010101010
Five software pr	ograms in de	etail 0110101110110100110111
00110110111100010101010100	1010110010000110000	1010011100010101010100101010101010101010

Five software programs were selected to show detail. The choice was made after informally surveying users on the programs they use most and their opinions as to the most useful or representative. Users included graduate students, postdoctorates and research associates, as well as faculty. A list was compiled of the most-mentioned programs. In cases of doubt, those that were freely available or seemed more widely used were chosen. To be more representative, an additional criterion was to choose those programs that did different things.



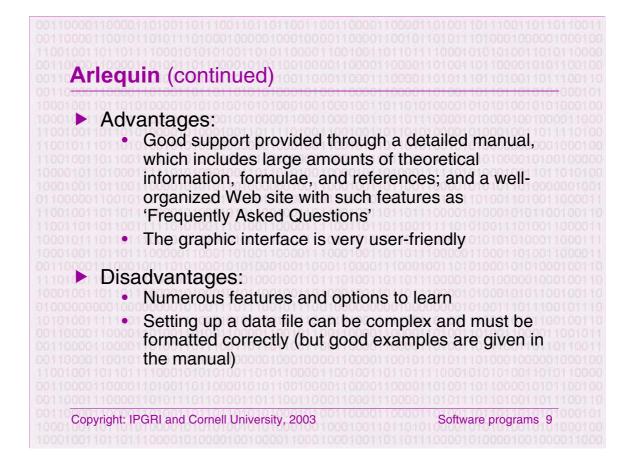
Released in 1997, Arlequin (current version 2.001) is still very popular. It is 'an exploratory population genetics software environment able to handle large samples of molecular data (RFLPs, DNA sequences, microsatellites), while retaining the capacity of analysing conventional genetic data (standard multi-locus data or mere allele frequency data).'

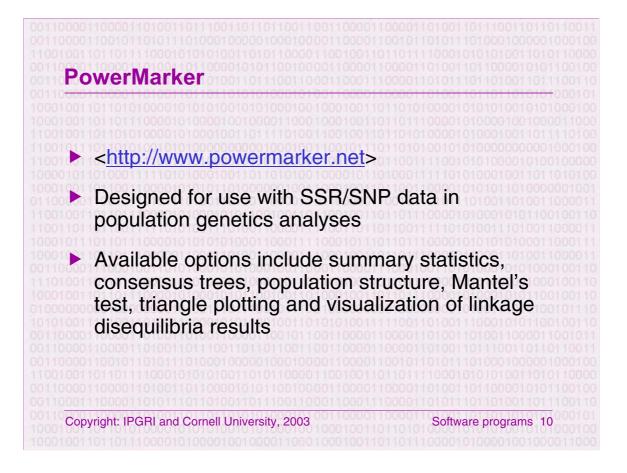
Arlequin can use many different types of data, such as molecular data and genotype or haplotype frequencies, including codominant or recessive data but not yet dominant data. Molecular data can be entered as DNA sequences, RFLP haplotypes, microsatellite profiles, or multilocus haplotypes. The data format is specified in an input file. The user can create a data file from scratch, using a text editor and appropriate keywords, or use the 'Project Outline Wizard'. Data can be imported from files created for other programs, including MEGA, BIOSYS, GENEPOP, and PHYLIP. Missing or ambiguous data can be included. A very detailed user manual is available, which includes a large amount of theoretical information, formulae, and references. A large number of data can be analysed, and a Batch Files option is available.

Authors: Laurent Excoffier, Stefan Schneider and David Roessli, University of Geneva, Switzerland.

Reference

Schneider, S., D. Roessli and L. Excoffier. 2000. Arlequin: A Software for Population Genetics Data Analysis, Version 2.000. Genetics and Biometry Laboratory, Dept. of Anthropology, University of Geneva, Switzerland.



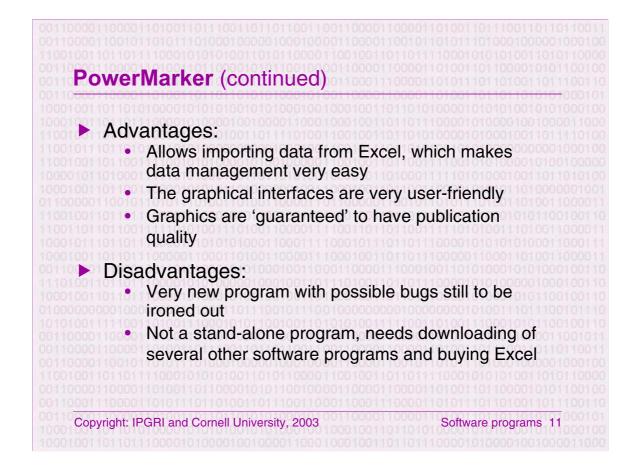


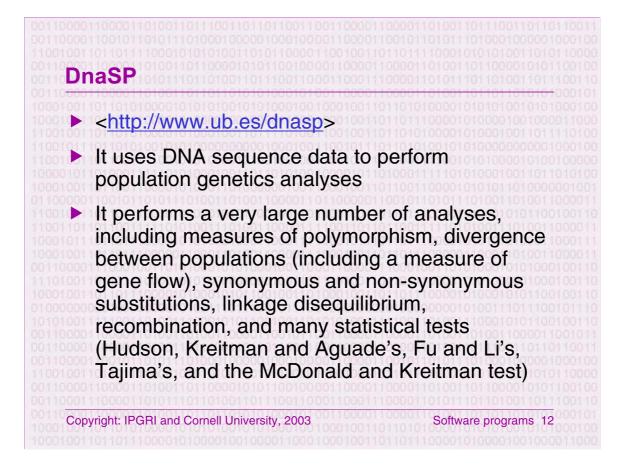
PowerMarker is a new program, with the first official version released in January 2004. It was designed specifically for the use of SSR/SNP data in population genetics analyses. Data can be imported from Excel or other formats, making data set-up very easy. Data can also be exported to NEXUS and Arlequin formats. It includes a '2D viewer' for linkage disequilibrium visualization. The user can edit graphics within PowerMarker or export them for publication. The program has been tested extensively for accuracy and efficieny. Full documentation is included. Several new modules for association study are included in the package. Several demonstration datasets available to get started. The program is free, but requires having PHYLIP, TreeView and the Microsoft.net framework system (all freely available) and Excel 2000 (not free). Another disadvantage is that it is available only for Windows 98 and above (not for Macintosh or other systems). Email support (registered user only): powermarker@hotmail.com

Author: Kejun (Jack) Liu, North Carolina State University.

## Reference

Liu, K. 2003. PowerMarker: New Genetic Data Analysis Software, Version 3.0. Free program distributed by the author over Internet at <http://www.powermarker.net>



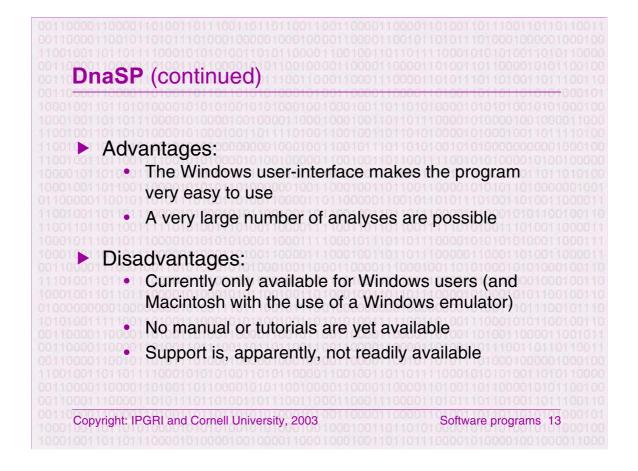


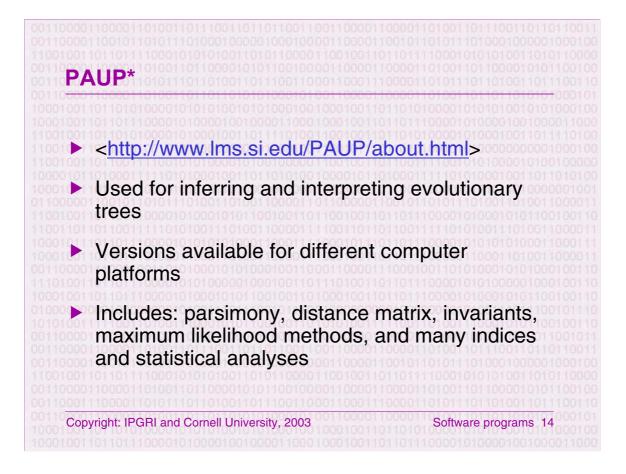
DnaSP, for DNA Sequence Polymorphism, uses DNA sequence data. This program is widely used for sequence analysis because it does all the necessary analyses and at the same time is easy to use. It was written exclusively for the Windows operating system, but can be run on a Macintosh using SoftWindows or VirtualPC software emulators. DnaSP can import and export several types of data formats, including FASTA and NEXUS, which is very convenient, and can handle large numbers of long sequences, depending on your computer's memory. The authors are currently working on version 4. It is freely available, downloadable from the Web site. Although no manual is available, a Help file is incorporated into the program. In addition, the Web site includes much explanatory material, as well as many references. The authors have several publications about the program (e.g. see citations below).

Authors: Julio and Ricardo Rozas

### References

- Rozas, J. and R. Rozas. 1995. DnaSP, DNA sequence polymorphism: an interactive program for estimating population genetics parameters from DNA sequence data. Comput. Appl. Biosci. 11:621-625.
- Rozas, J. and R. Rozas. 1997. DnaSP version 2.0: a novel software package for extensive molecular population genetics analysis. Comput. Appl. Biosci. 13:307-311.
- Rozas, J. and R. Rozas. 1999. DnaSP version 3: an integrated program for molecular population genetics and molecular evolution analysis. Bioinformatics 15:174-175.



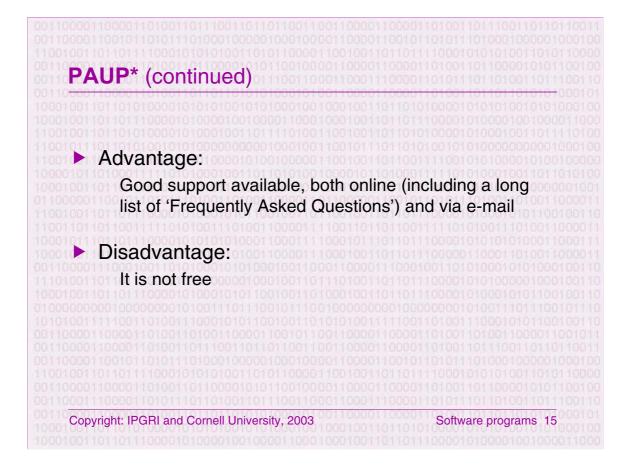


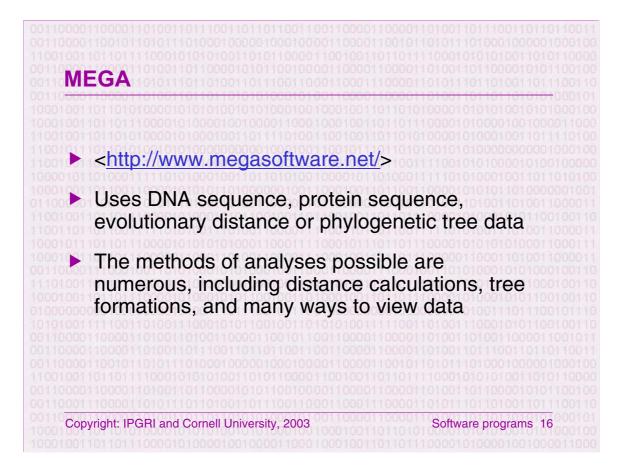
PAUP\* is widely used for inferring and interpreting evolutionary trees. It originally meant Phylogenetic Analysis Using Parsimony, but now has many other options. PAUP\* is available from Sinauer Associates, Sunderland, MA, at <a href="http://www.sinauer.com/detail.php?id=8060">http://www.sinauer.com/detail.php?id=8060</a>>. Although not free, it is relatively inexpensive (US\$100 at writing). A new version, 4.0 beta, has been released as a provisional version. Macintosh, PowerMac, Windows and Unix/OpenVMS versions are available; the Mac version has some extra features. PAUP\* is closely compatible with MacClade (another program available from Sinauer), since they use a common data format (NEXUS, Maddison et al. 1997).

Author: David Swofford, Laboratory of Molecular Systematics, National Museum of Natural History, Smithsonian Institution, Washington, DC.

Reference

Swofford, D.L. 2002. PAUP\*: Phylogenetic Analysis Using Parsimony (\*and Other Methods), Version 4. Sinauer Associates, Sunderland, MA.





MEGA (Molecular Evolutionary Genetics Analysis) software has been widely used since its creation in 1993; MEGA2 has since come out. It uses DNA sequence, protein sequence, evolutionary distance or phylogenetic tree data. The authors' goal was to take advantage of advances in computer power and graphic user interfaces to make available a 'flexible and easy-to-use genetic data analysis workbench'. Although it was designed for the Windows platform, it runs well on Macintosh with a Windows emulator, Sun workstation (with SoftWindows95) or Linux (with Windows by VMWare). The newest version, 2.1, has many important additions, such as the ability to import data from NEXUS or CLUSTAL W, unlimited dataset sizes, and many others.

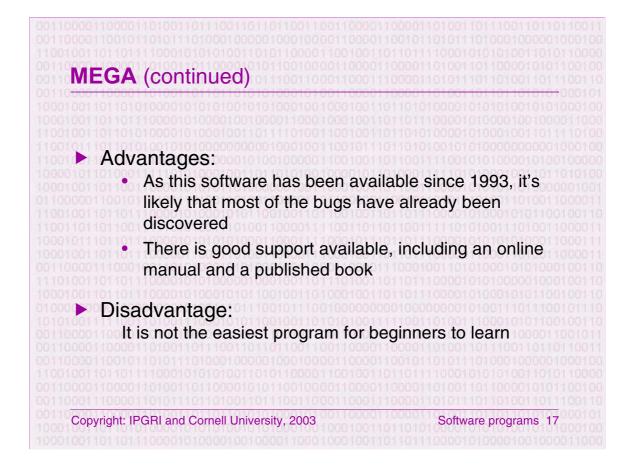
A book by the software authors Nei and Kumar (2000) includes theoretical information about statistical analyses and how to interpret results from both their software and other software programs. Online, a thorough manual is available (although the format is not easy to page through), together with a bulletin board for users to interact with each other.

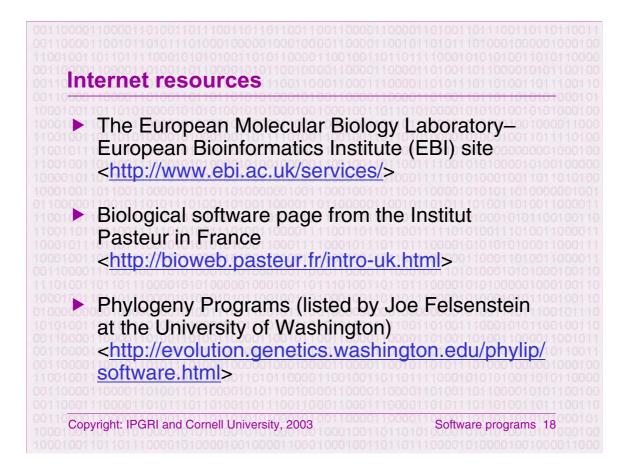
Authors: Sudhir Kumar, Koichiro Tamura, Ingrid Jakobsen and Masatoshi Nei.

## References

Nei, M. and S. Kumar. 2000. Molecular Evolution and Phylogenetics. Oxford University Press, NY.

Sudhir, K., T. Koichiro, I.B. Jakobsen and M. Nei. 2001. MEGA2: Molecular Evolutionary Genetics Analysis software. Bioinformatics 12(17):1244-1245.



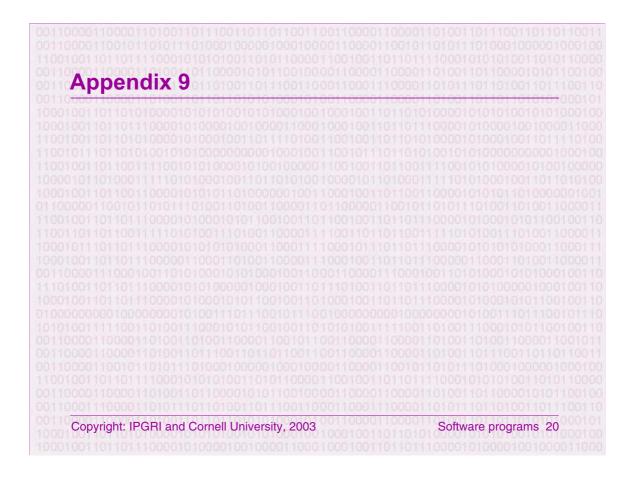


In this and the next slide, we give a sample list of Internet resources that you may find useful for locating information related to, for example, genetic diversity analysis, population genetics, other software available, and links to useful extra information. For each resource presented, we briefly describe their contents in the notes below.

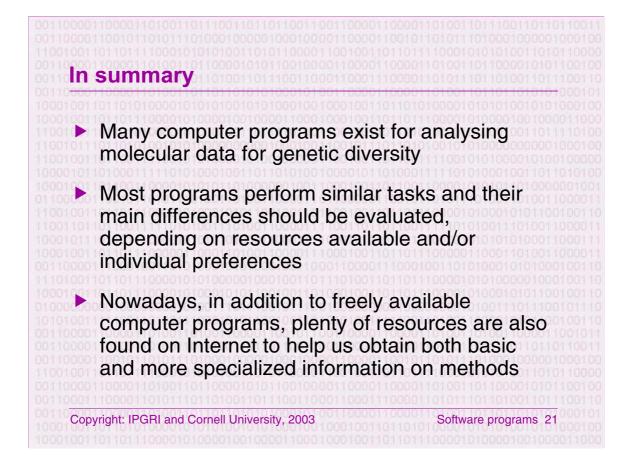
- The European Molecular Biology Laboratory–European Bioinformatics Institute (EBI) site: not only does it contain links to many useful programs and other sites but it also describes what they do and is therefore a good source of general information as well.
- Biological software page from the Institut Pasteur in France: although some pages are in French only, it provides a very comprehensive list of software programs available online, including links, and is current (last updated December 2002). It also contains links to many programs developed at the Institut Pasteur.
- Phylogeny Programs: this is the longest list of phylogeny programs we have seen, counting 194. The author adds the caveat that he has not tried to assess their quality or cost. Nor has the list been updated since 2001, but it contains so many links to programs, and sorts them in various ways (e.g. by methods, system used) that it is still very useful.

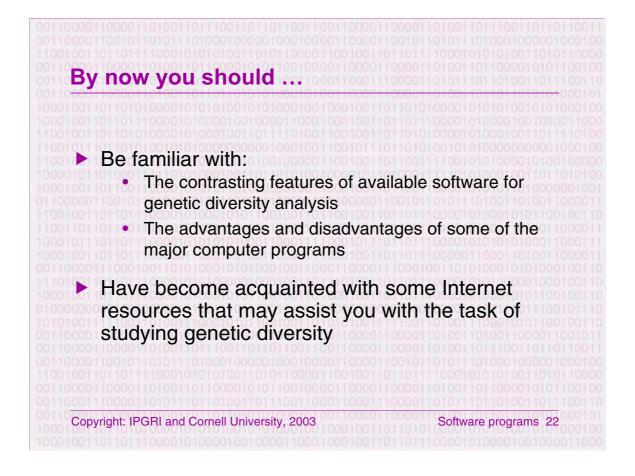


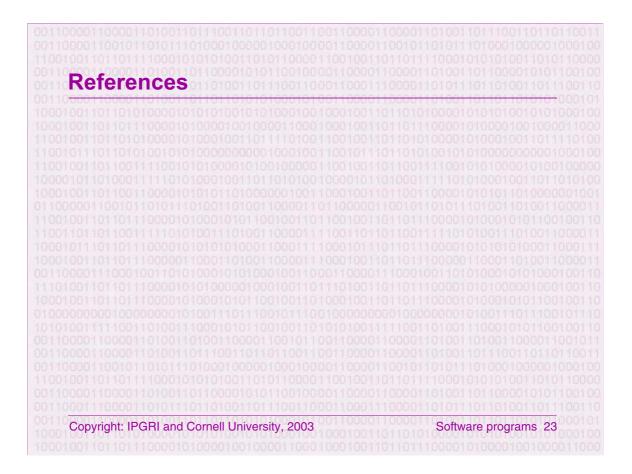
- Dr. Ed Buckler's Maize Genetics site: contains freely available software programs developed by his laboratory group. Although some information is specific to maize, the site also contains useful information about genomics, links to many journals, and PDF versions of Dr. Buckler's publications.
- Kent Holsinger's site at the University of Connecticut: has links to many software programs, including ones for biology, programming and statistics. The author does not post things he does not use regularly, so some guarantee of good quality exists. At writing, this was also the most recently updated.
- Claire Constantine's site at Murdoch University: although not updated recently, this site contains links to the most used programs for population genetics analyses. Even more useful, it includes a comparison table of the kinds of statistics available in each of 7 commonly used programs (Arlequin, GENEPOP, POPGENE, GDA, GeneStrut, DnaSP and SITES).
- Software page of the Institute of Forest Genetics and Forest Tree Breeding, University of Göttingen, Germany: this site contains just 4 software programs, developed in-house, but they are freely available, include good descriptions, and the page is regularly updated.



Appendix 9. References to software programs







#### References

- Labate, J.A. 2000. Software for population genetic analyses of molecular marker data. Crop Sci. 40:1521-1528.
- Liu, K. 2003. PowerMarker: New Genetic Data Analysis Software, Version 1.0. Free program distributed by the author over Internet at <a href="http://www.powermarker.net">http://www.powermarker.net</a>>
- Maddison, D. R., D.L. Swofford and W.P. Maddison. 1997. NEXUS: an extensible file format for systematic information. Syst. Biol. 46:590–621.
- Nei, M. and S. Kumar. 2000. Molecular Evolution and Phylogenetics. Oxford University Press, NY.
- Rozas, J. and R. Rozas. 1995. DnaSP, DNA sequence polymorphism: an interactive program for estimating population genetics parameters from DNA sequence data. Comput. Appl. Biosci. 11:621-625.
- Rozas, J. and R. Rozas. 1997. DnaSP, version 2.0: a novel software package for extensive molecular population genetics analysis. Comput. Appl. Biosci. 13:307-311.
- Rozas, J. and R. Rozas. 1999. DnaSP, version 3: an integrated program for molecular population genetics and molecular evolution analysis. Bioinformatics 15:174-175.
- Schneider, S., D. Roessli and L. Excoffier. 2000. Arlequin: A Software for Population Genetics Data Analysis, Version 2.000. Genetics and Biometry Laboratory, Dept. of Anthropology, University of Geneva, Switzerland.
- Sudhir, K., T. Koichiro, I.B. Jakobsen and M. Nei. 2001. MEGA2: Molecular Evolutionary Genetics Analysis software. Bioinformatics 12(17):1244-1245.
- Swofford, D.L. 2002. PAUP\*, Phylogenetic Analysis Using Parsimony (\*and Other Methods), Version 4. Sinauer Associates, Sunderland, MA.

