

Tutorial of *STRUCTURE* software

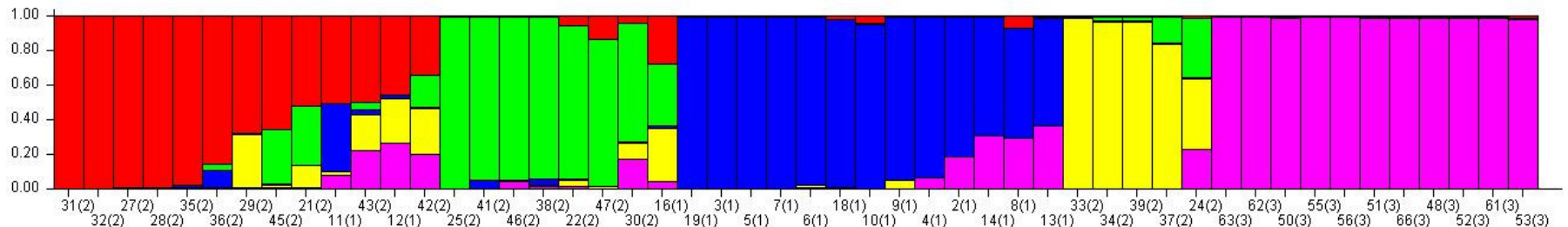
Sung-Chur Sim

Tomato Genetics and Breeding program

The Ohio State Univ., OARDC

STRUCTURE software

- A model-based clustering method (Pritchard et al. 2000)
 - Free software (http://pritch.bsd.uchicago.edu/software/structure2_1.html)
 - Bayesian approach (MCMC: Markov Chain Monte Carlo)
 - Detects the underlying genetic population among a set of individuals genotyped at multiple markers
 - Computes the proportion of the genome of an individual originating from each inferred population (quantitative clustering method)



Input data

🍅 A matrix where the data for individuals are in rows, the loci are in column

- **n consecutive rows** have the data for each individual of n -ploid species
- **Integer** should be used for coding genotype
- Missing data should be indicated by **a number** which doesn't occur elsewhere in the data (e.g. -1)
- The data file should be a **text file (.txt)** not an excel file (.xls) for running STRUCTURE

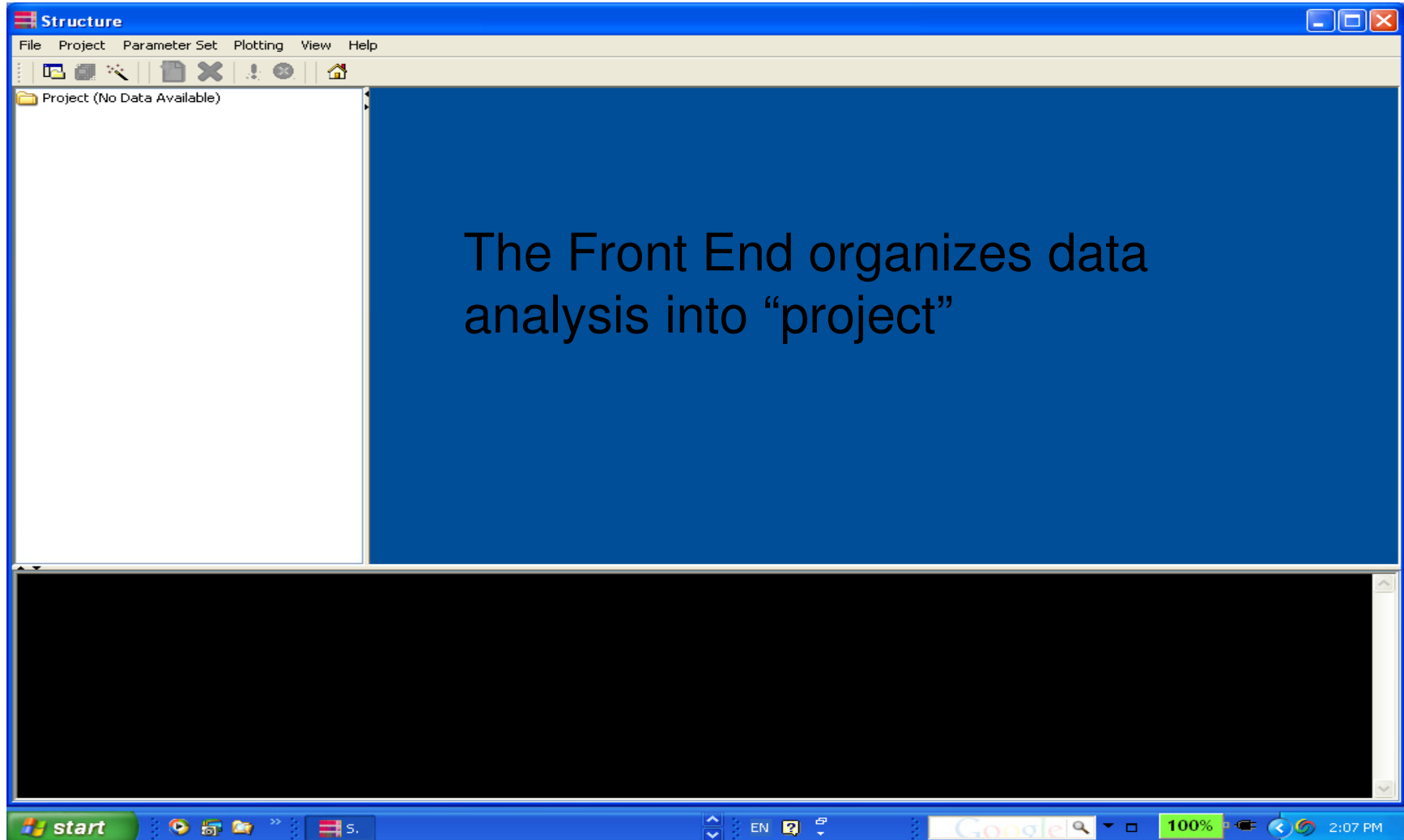
Information of user-defined populations (market class)

2 consecutive rows
for alleles

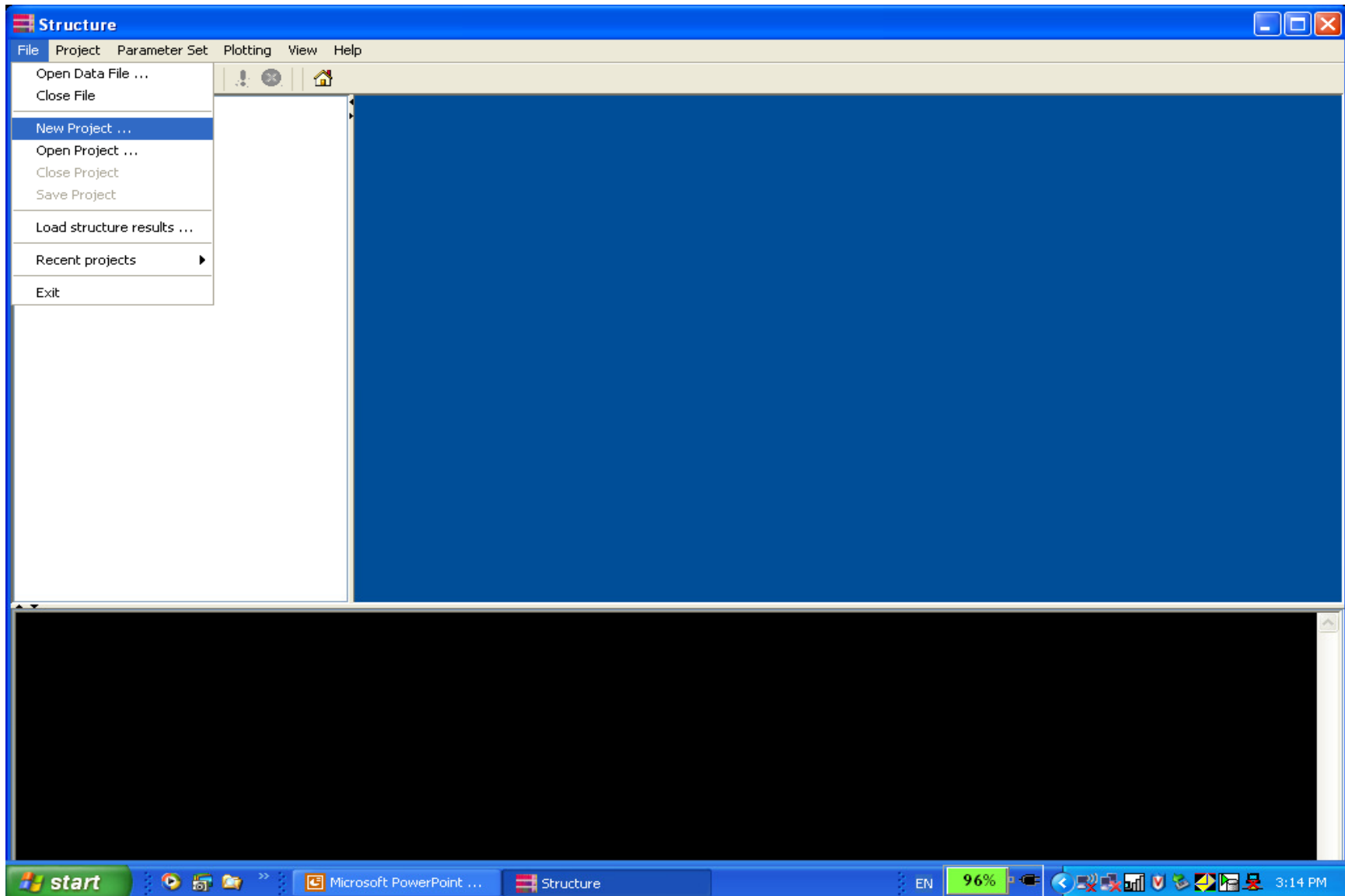
	A		C	D	E	F	G	H	I	J	K	L	M	N	O
1			CT10153	CT10162	CT10184	CT10187	CT10238	CT10242	CT10386	CT10396	CT10526	CT10554	CT10556	CT10649	CT1073
2	Campbell28	1	14	12	-1	13	13	12	14	13	13	11	13	11	
3	Campbell28	1	14	12	-1	13	13	12	14	13	13	11	13	11	
4	Fla7010	1	12	13	12	13	13	12	14	13	13	11	13	11	
5	Fla7060	1	12	13	12	13	13	12	14	13	13	11	13	11	
6	Fla7547	1	12	12	12	13	13	12	14	12	13	11	13	11	
7	Fla7547	1	12	12	12	13	13	12	14	12	13	11	13	11	
8	Fla7771	1	14	12	12	13	13	12	14	12	13	11	13	11	
9	Fla7771	1	14	12	12	13	13	12	14	12	13	11	13	11	
10	Fla7775	1	14	13	12	13	13	12	14	12	13	11	-1	11	
11	Fla7775	1	14	13	12	13	13	12	14	12	13	11	-1	11	
12	Fla7600	1	14	12	13	13	13	12	14	13	13	11	13	11	
13	Fla7600	1	14	12	13	13	13	12	14	13	13	11	13	11	
14	Floradade	1	14	12	12	13	13	12	14	12	13	11	13	11	
15	Floradade	1	14	12	12	13	13	12	14	12	13	11	13	11	
16	HC23E-2(93)	1	14	12	13	13	13	12	14	13	13	14	13	11	
17	HC23E-2(93)	1	14	12	13	13	13	12	14	13	13	14	13	11	
18	HC353-1	1	12	13	13	13	13	12	14	13	-1	11	13	-1	
19	HC353-1	1	12	13	13	13	13	12	14	13	-1	11	13	-1	
20	HC84173	1	12	13	12	13	13	12	14	13	13	11	13	11	
21	HC84173	1	12	13	12	13	13	12	14	13	13	14	13	11	
22	HC98248	1	14	12	13	13	13	12	14	13	13	11	13	11	
23	HC98248	1	12	12	13	13	13	12	14	13	13	11	13	11	
24	HC99471-3	1	12	12	13	13	13	12	14	13	13	11	13	14	
25	HC99471-3	1	12	12	13	13	13	12	14	13	13	11	13	14	
26	HCCEBR2	1	14	12	13	13	13	12	14	13	13	-1	13	11	
27	HCCEBR2	1	14	12	13	13	13	12	14	13	13	-1	13	11	
28	Ohio-MR13	1	14	12	12	13	13	12	14	13	13	11	13	11	
29	Ohio-MR13	1	14	12	12	13	13	12	14	13	13	11	13	11	
30	Ohio11	1	14	12	13	13	13	12	14	13	13	11	13	11	
31	Ohio11	1	14	12	13	13	13	12	14	13	13	11	13	11	

Missing data

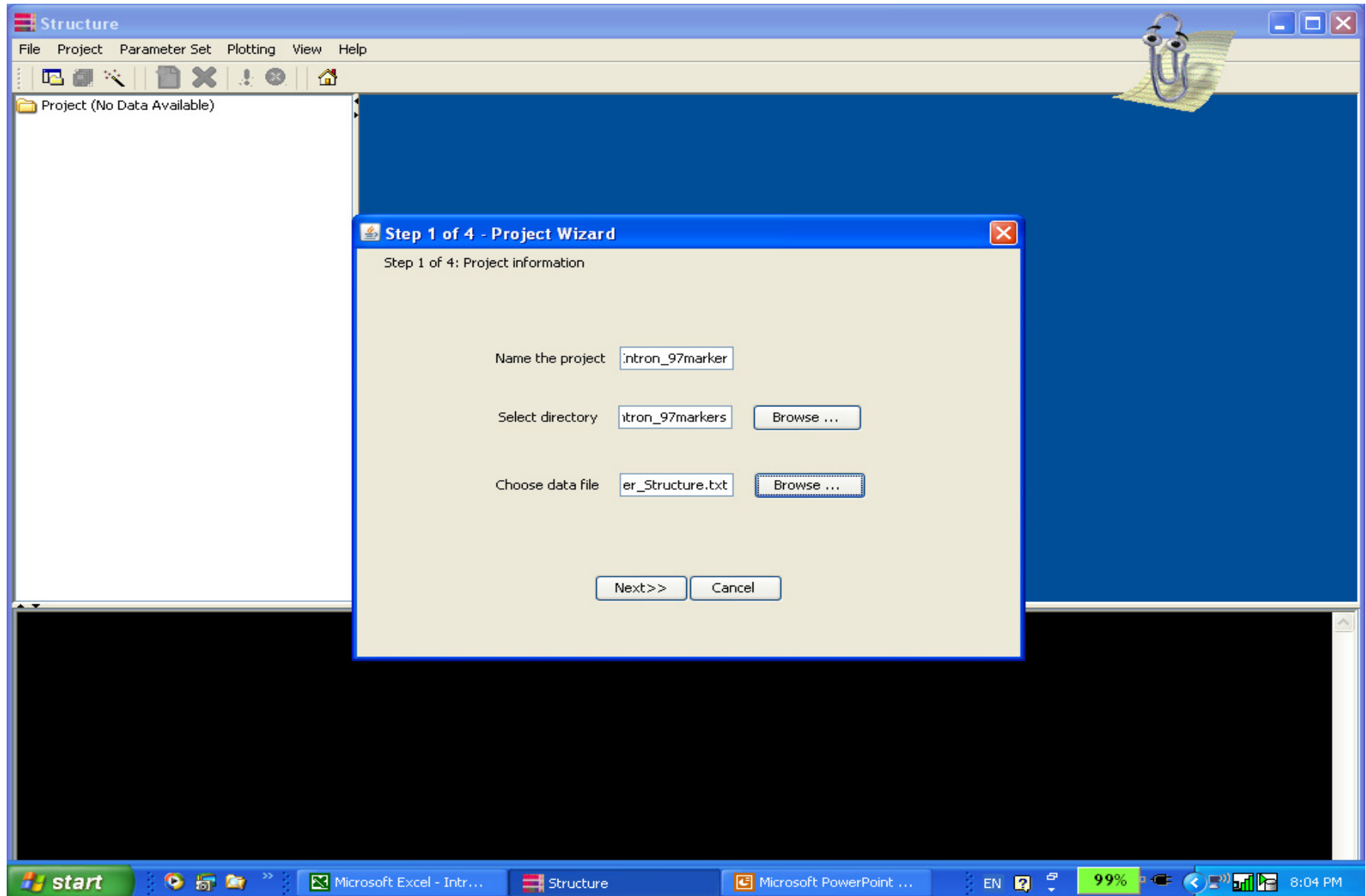
Running STRUCTURE from a graphical interface, Front End



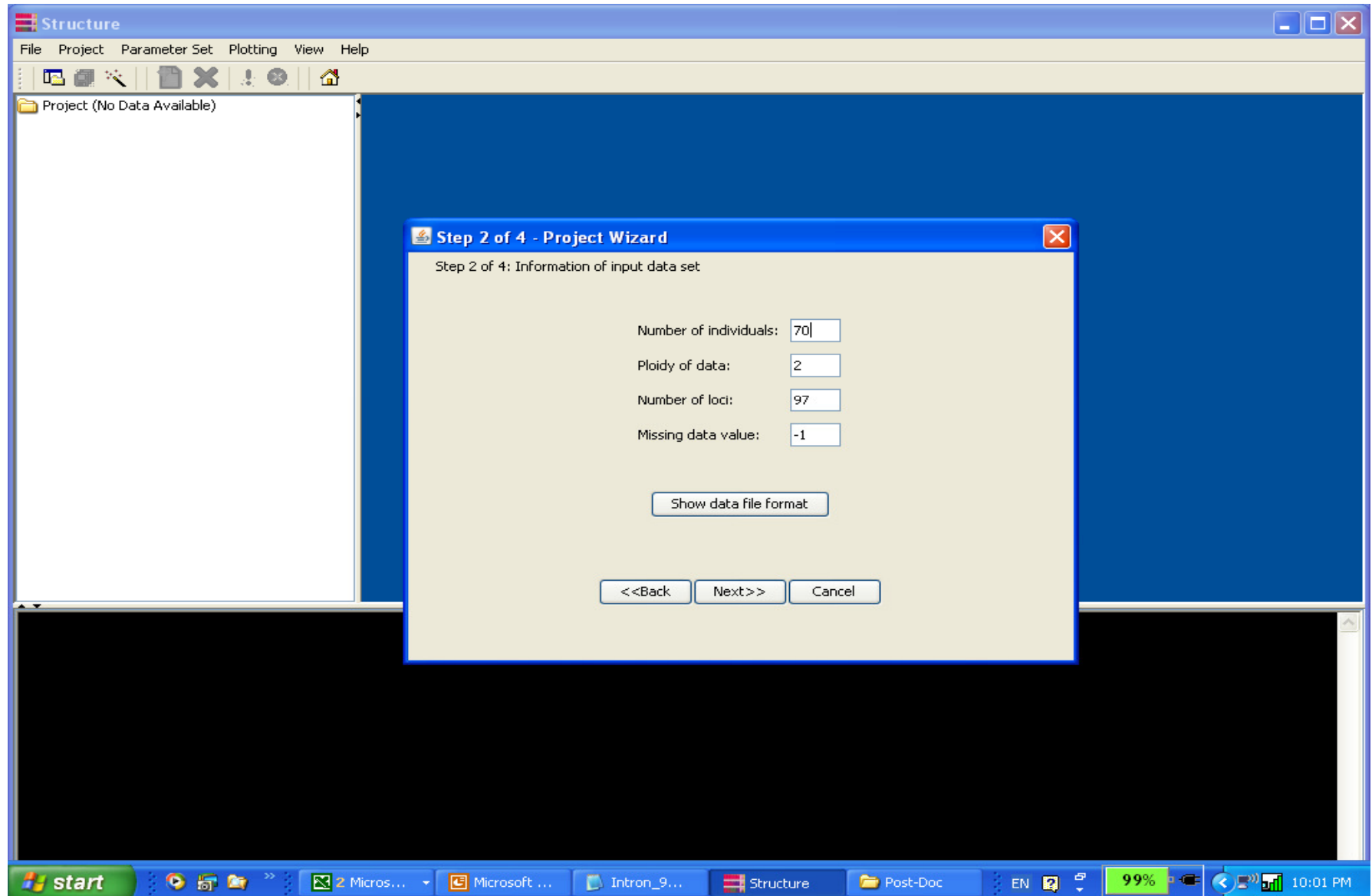
Importing a input data into a project



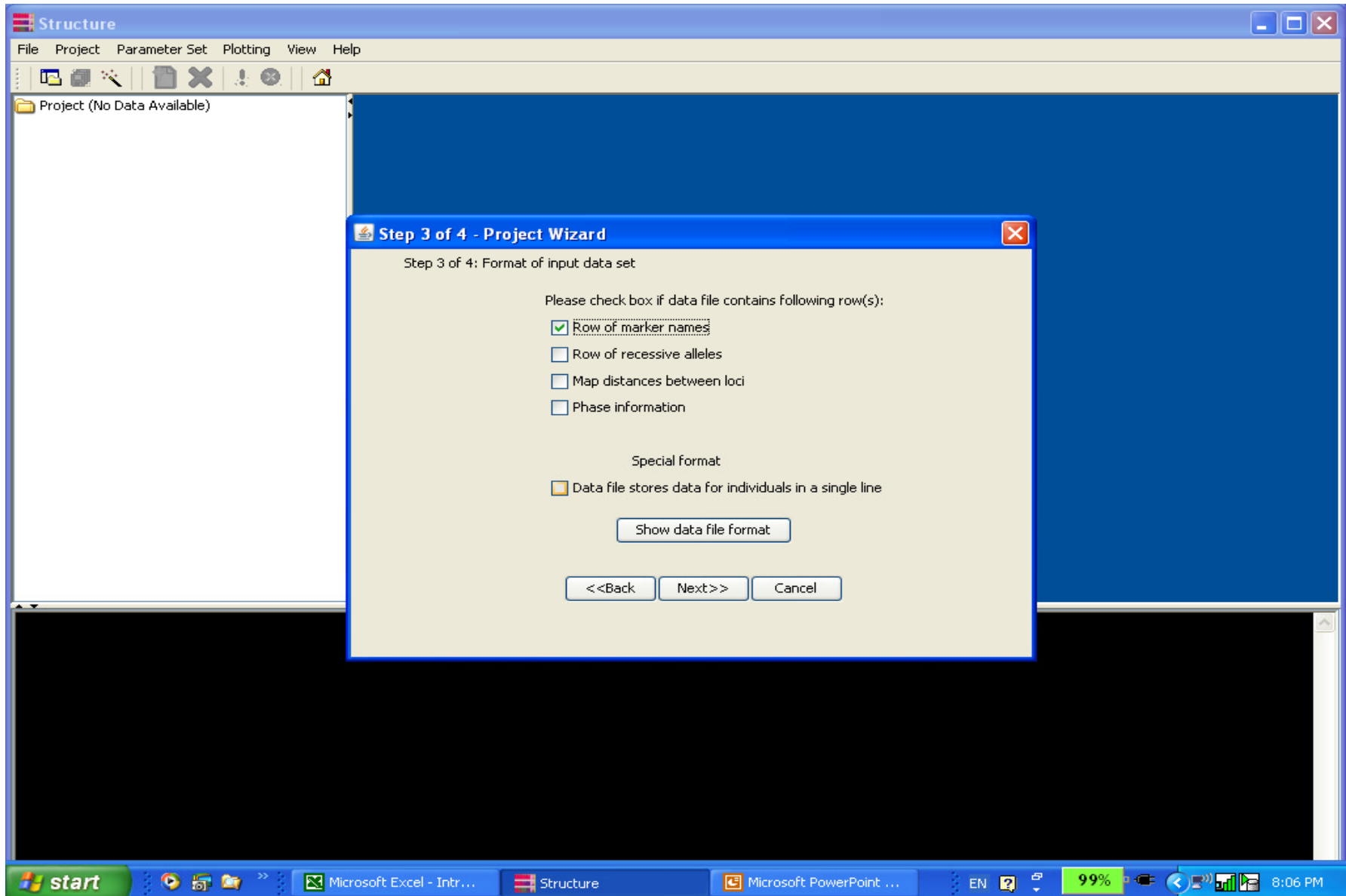
Importing a input data into a project (cont.)



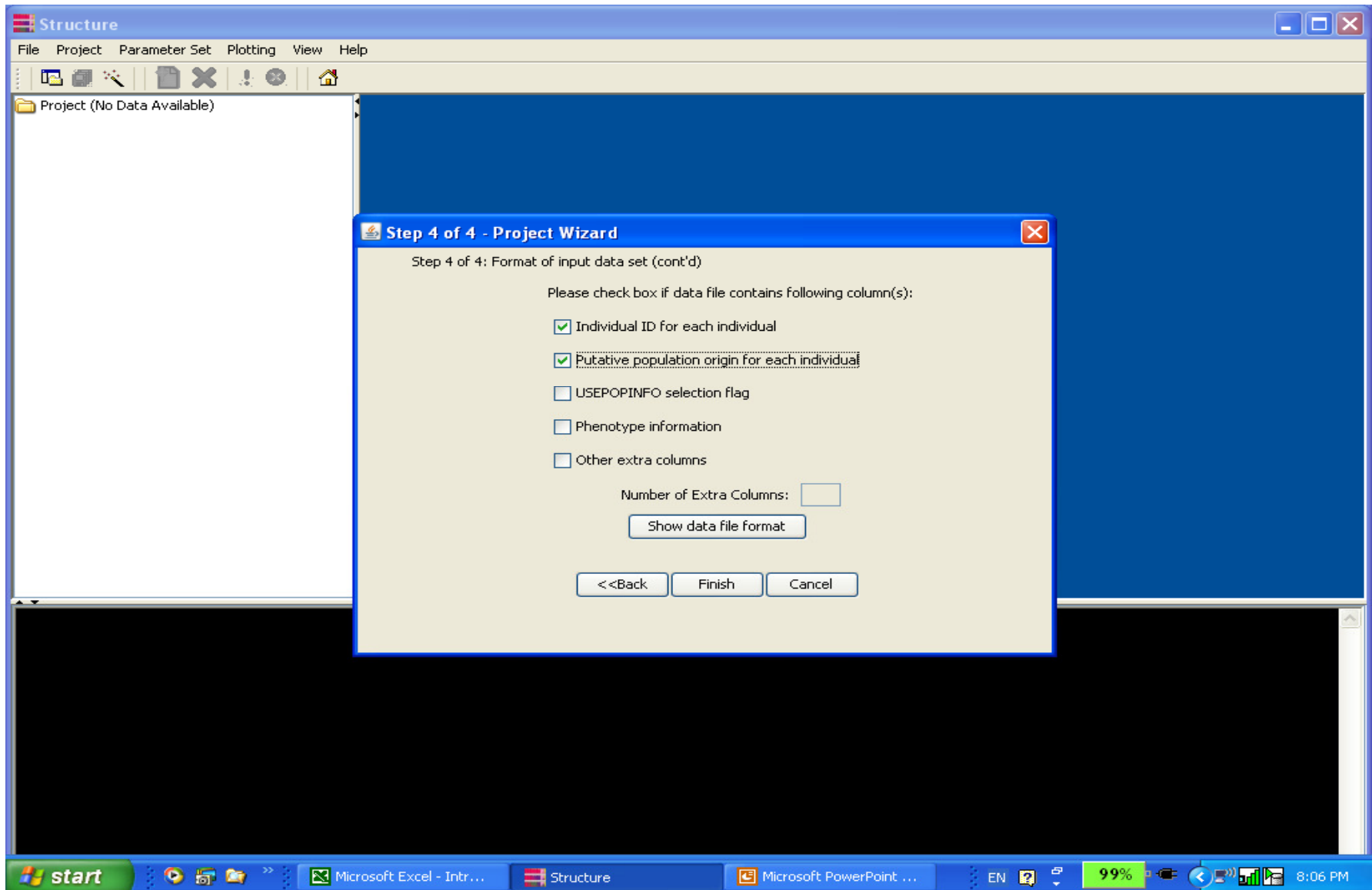
Importing a input data into a project (cont.)



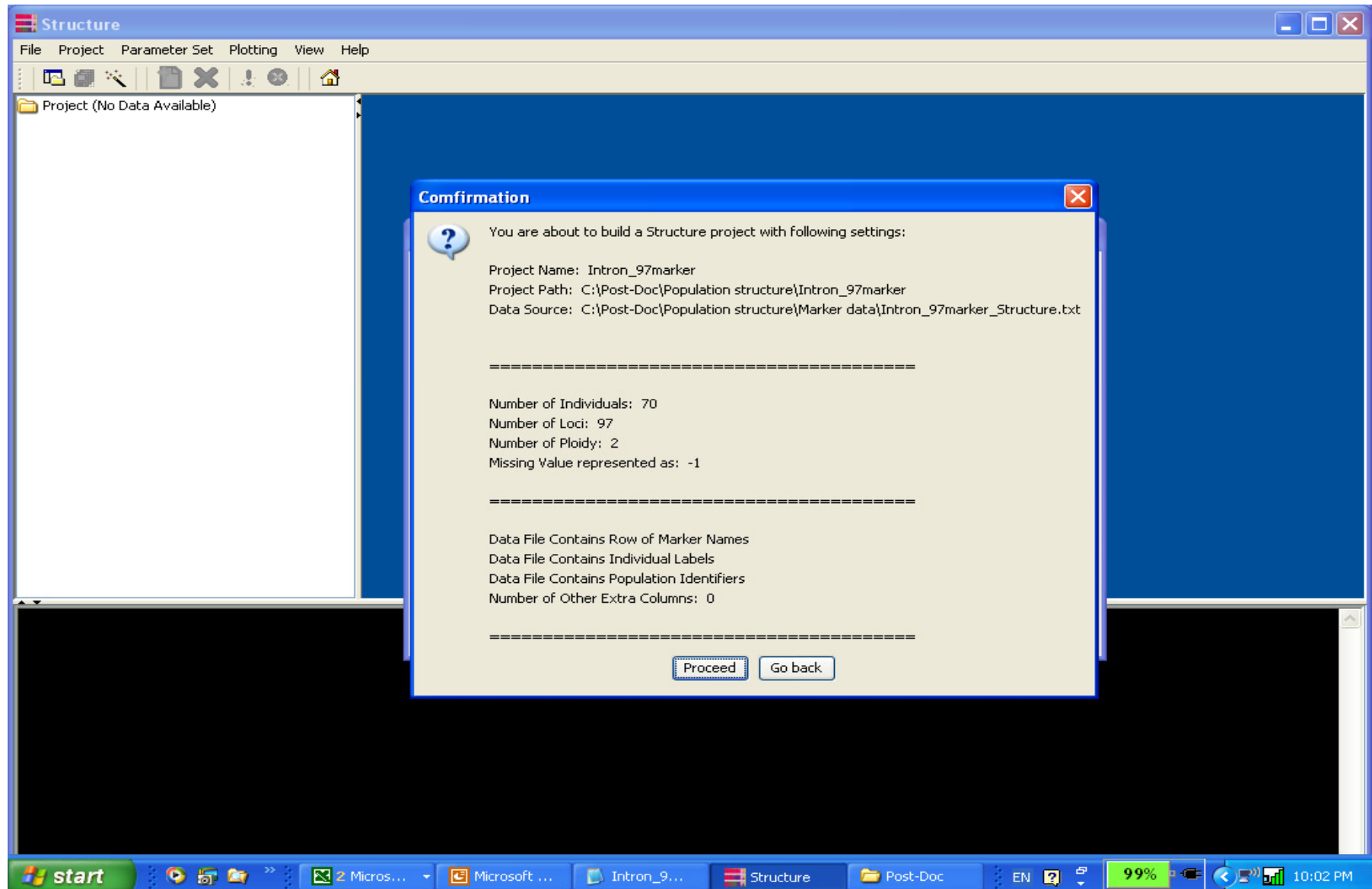
Importing a data file into a project (cont.)



Importing a input data into a project (cont.)



Importing a input data into a project (cont.)

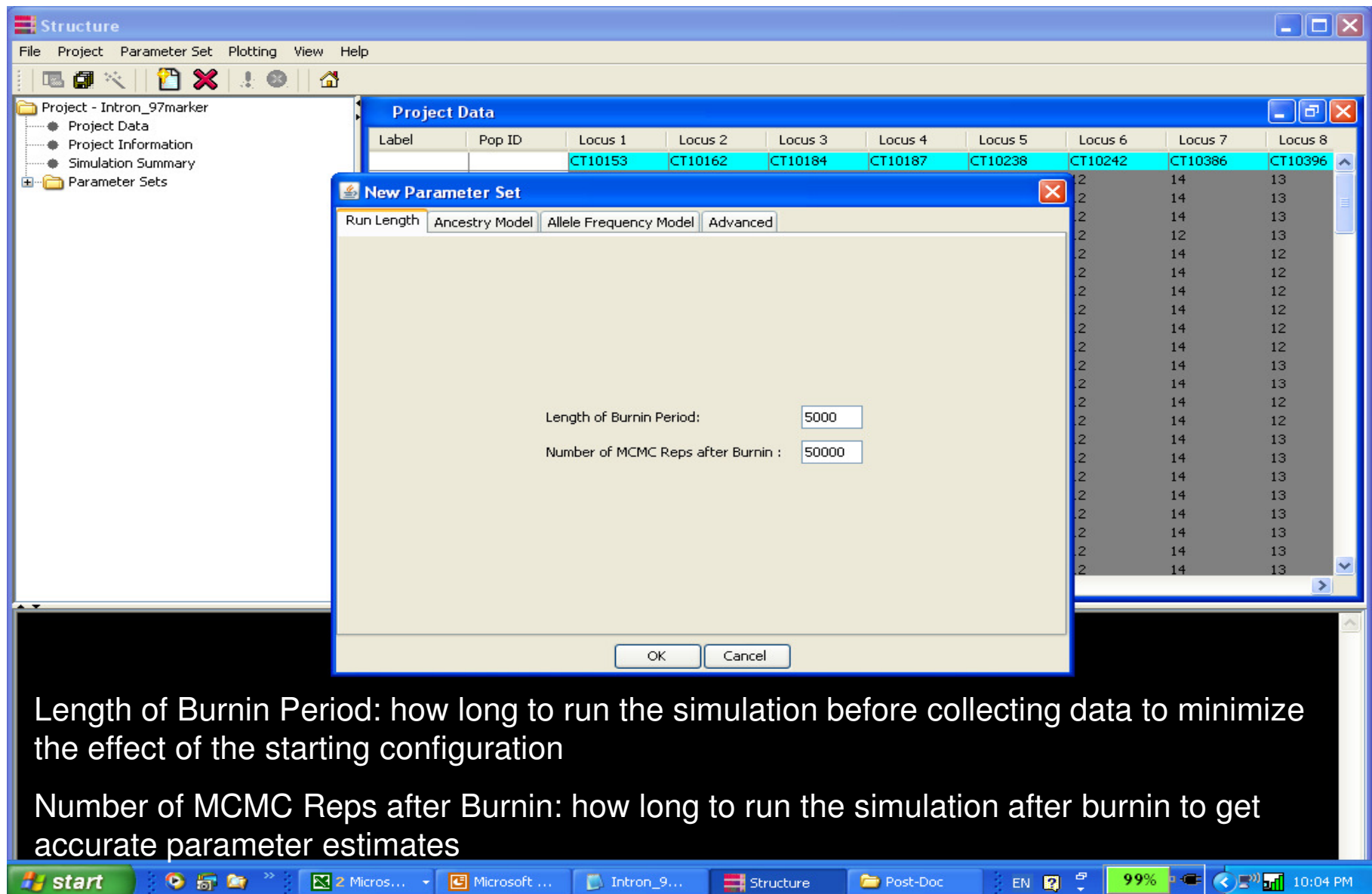


Configuring a parameter set

The screenshot displays the Structure software interface. The 'Parameter Set List' menu is open, showing options: 'Modify current set ...', 'New ...', 'Remove Parameter Set ...', 'Run', and 'Stop'. The 'Project Data' table is visible, listing various loci and their values across different populations.

Label	Pop ID	Locus 1	Locus 2	Locus 3	Locus 4	Locus 5	Locus 6	Locus 7	Locus 8
		CT10153	CT10162	CT10184	CT10187	CT10238	CT10242	CT10386	CT10396
C28	1	14	12	-1	13	13	12	14	13
C28	1	14	12	-1	13	13	12	14	13
F7060	1	12	13	12	13	13	12	14	13
F7060	1	12	13	12	13	13	12	12	13
F7547	1	12	12	12	13	13	12	14	12
F7547	1	12	12	12	13	13	12	14	12
F7771	1	14	12	12	13	13	12	14	12
F7771	1	14	12	12	13	13	12	14	12
F7775	1	14	13	12	13	13	12	14	12
F7775	1	14	13	12	13	13	12	14	12
FL7600	1	14	12	13	13	13	12	14	13
FL7600	1	14	12	13	13	13	12	14	13
Floradade	1	14	12	12	13	13	12	14	12
Floradade	1	14	12	12	13	13	12	14	12
NC23E-2	1	14	12	13	13	13	12	14	13
NC23E-2	1	14	12	13	13	13	12	14	13
NC353	1	12	13	13	13	13	12	14	13
NC353	1	12	13	13	13	13	12	14	13
NC84173	1	12	13	12	13	13	12	14	13
NC84173	1	12	13	12	13	13	12	14	13
NC98248	1	14	12	13	13	13	12	14	13
NC98248	1	12	12	13	13	13	12	14	13

Configuring a parameter set (cont.)



Length of Burnin Period: how long to run the simulation before collecting data to minimize the effect of the starting configuration

Number of MCMC Reps after Burnin: how long to run the simulation after burnin to get accurate parameter estimates

Configuring a parameter set (cont.)

The screenshot shows the Structure software interface. The main window displays a 'Project Data' table with columns for Label, Pop ID, and eight Loci (Locus 1 to Locus 8). The table contains two rows of data for 'C28' with Pop ID 1. A 'New Parameter Set' dialog box is open in the foreground, allowing configuration of the parameter set. The dialog has tabs for 'Run Length', 'Ancestry Model', 'Allele Frequency Model', and 'Advanced'. The 'Ancestry Model' tab is selected, and the 'Use Admixture Model' option is chosen. There are 'Advanced ...' buttons next to each model option and a 'Default Setting' button at the bottom. The 'OK' and 'Cancel' buttons are at the bottom of the dialog.

Project Data

Label	Pop ID	Locus 1	Locus 2	Locus 3	Locus 4	Locus 5	Locus 6	Locus 7	Locus 8
		CT10153	CT10162	CT10184	CT10187	CT10238	CT10242	CT10386	CT10396
C28	1	14	12	-1	13	13	12	14	13
C28	1	14	12	-1	13	13	12	14	13

New Parameter Set

Run Length | **Ancestry Model** | Allele Frequency Model | Advanced

Select ONE from the following:

- ☐ Use No Admixture Model
- ☒ Use Admixture Model
- ☐ Use Linkage Model
- ☐ Use Population Information

Advanced ...

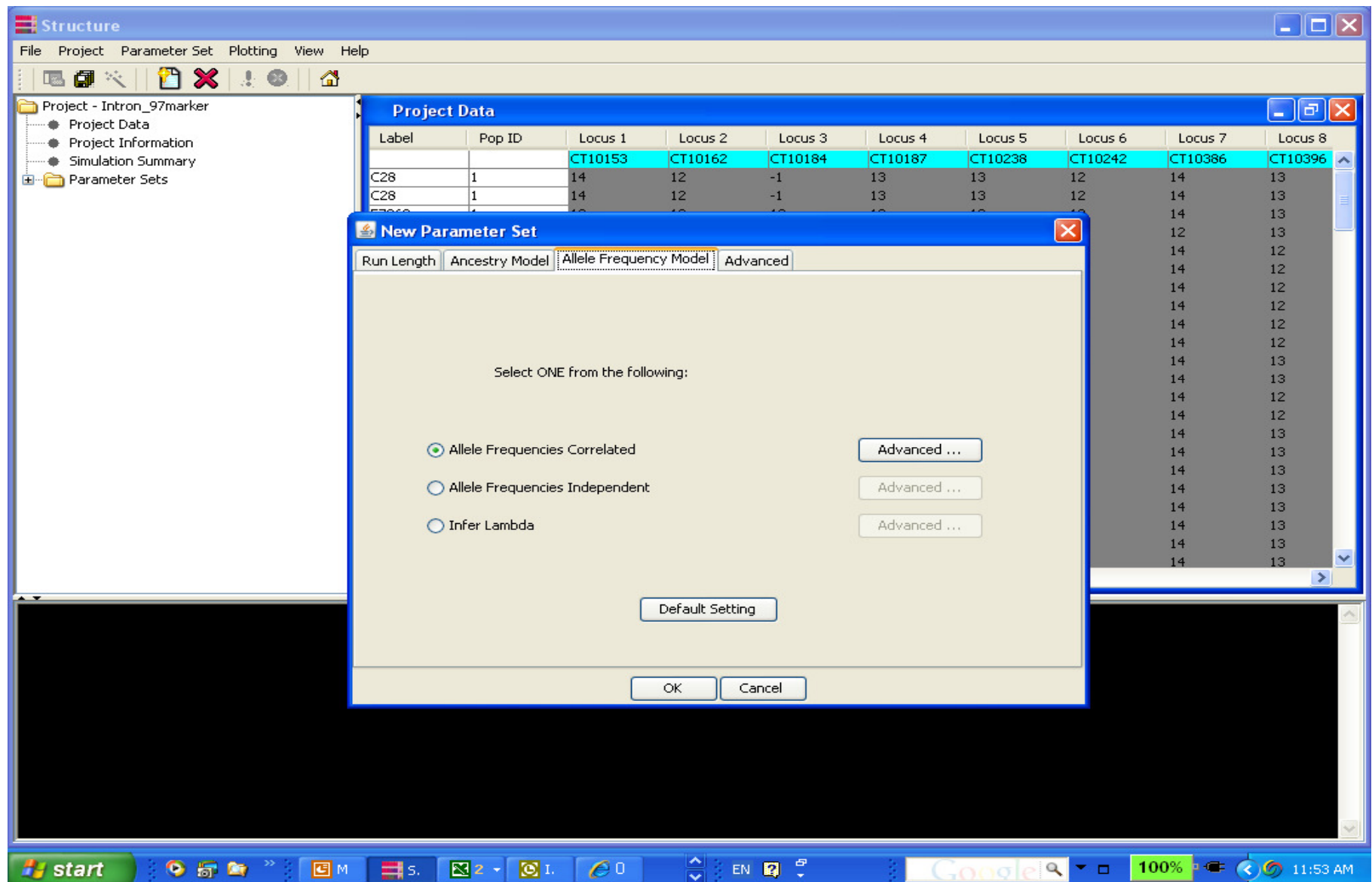
Advanced ...

Advanced ...

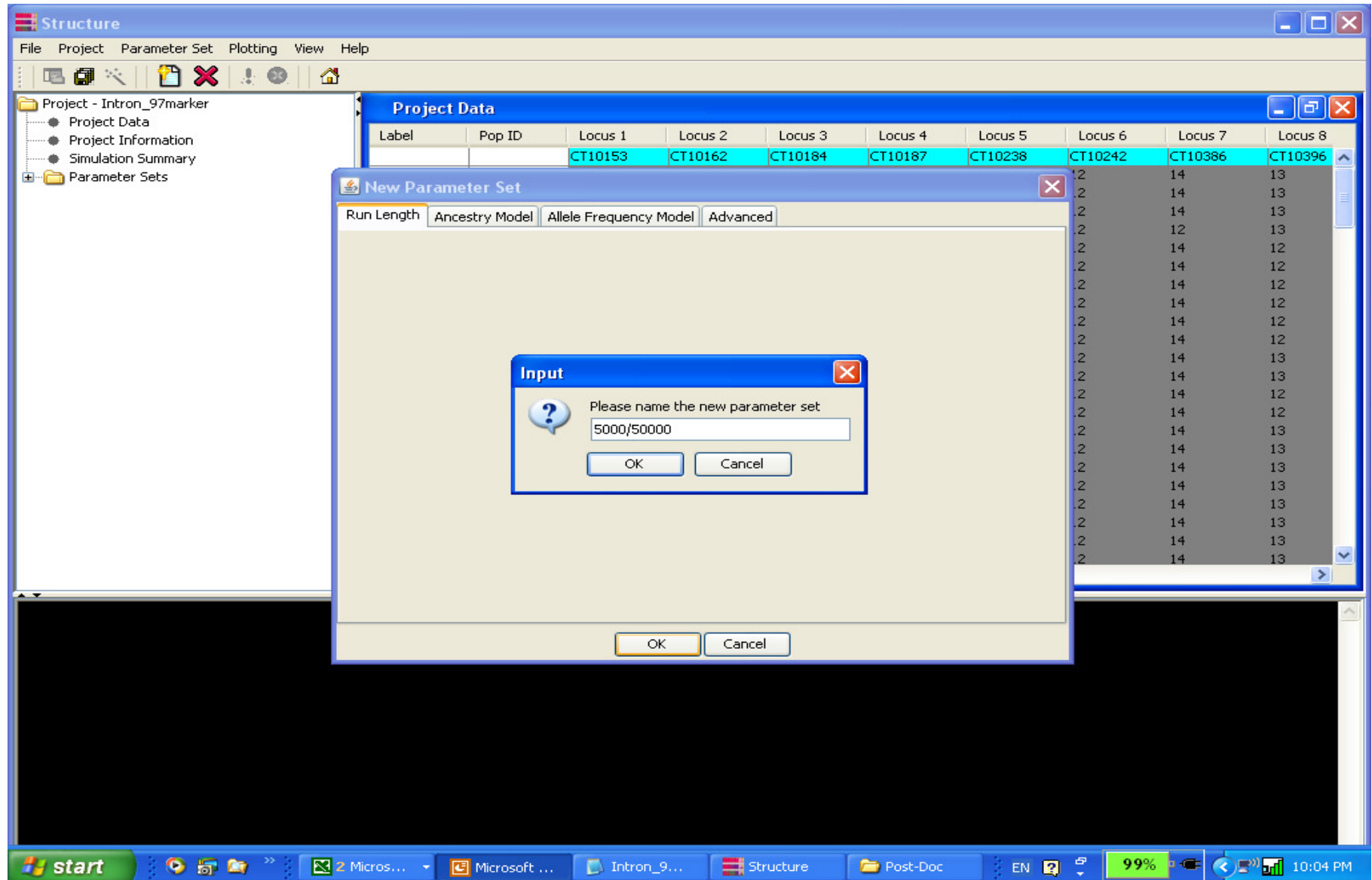
Default Setting

OK Cancel

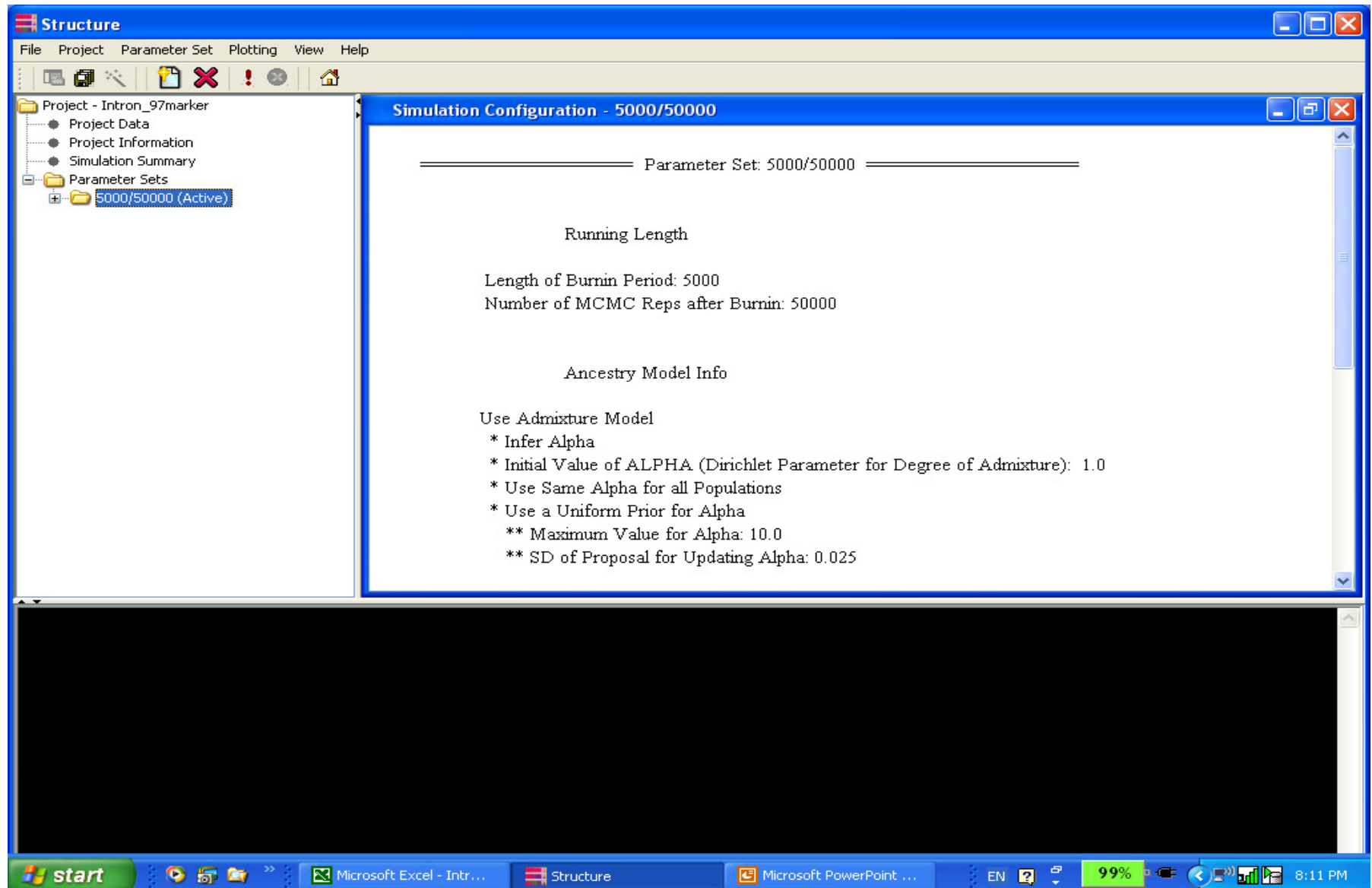
Configuring a parameter set (cont.)



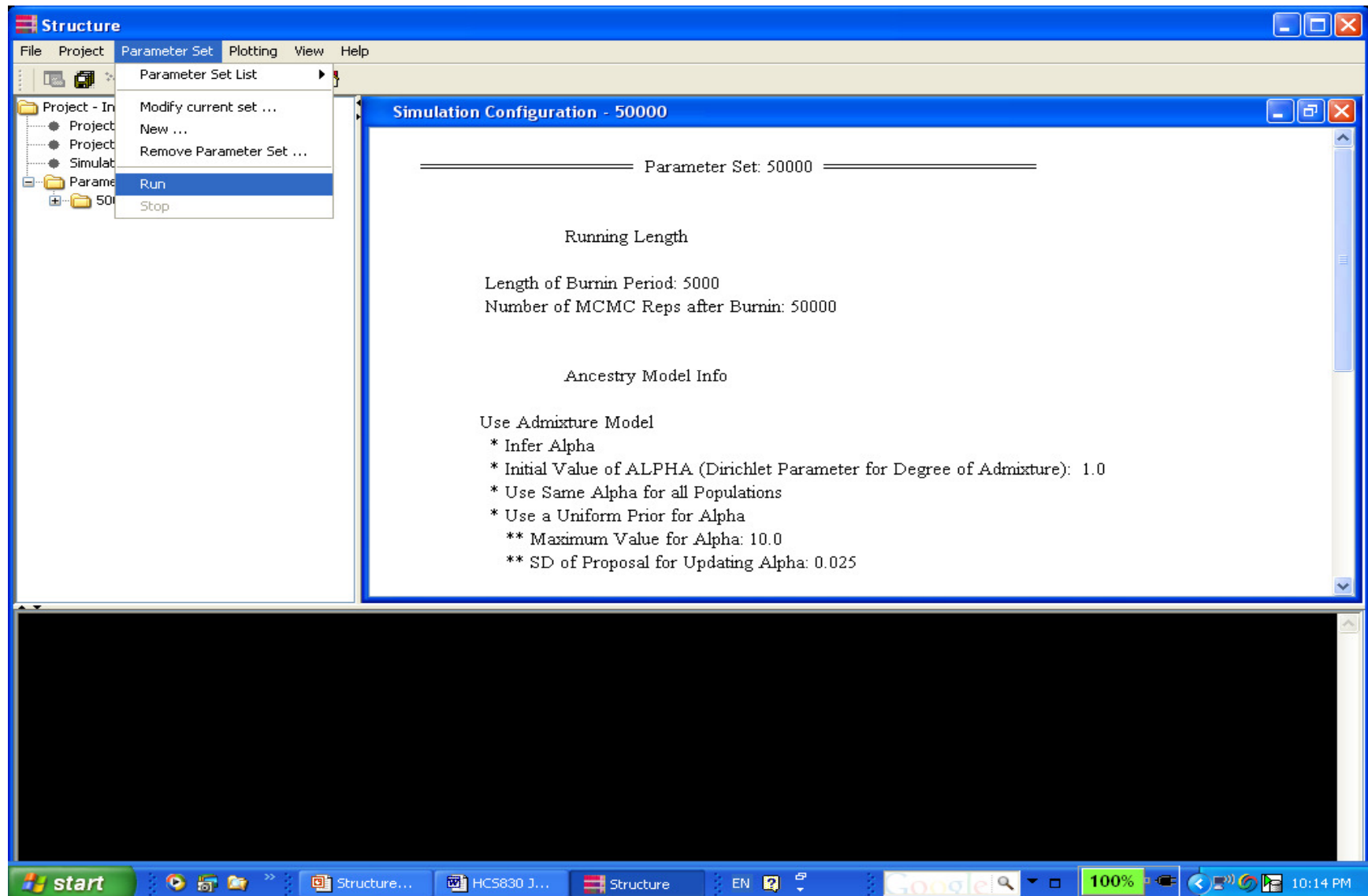
Configuring a parameter set (cont.)



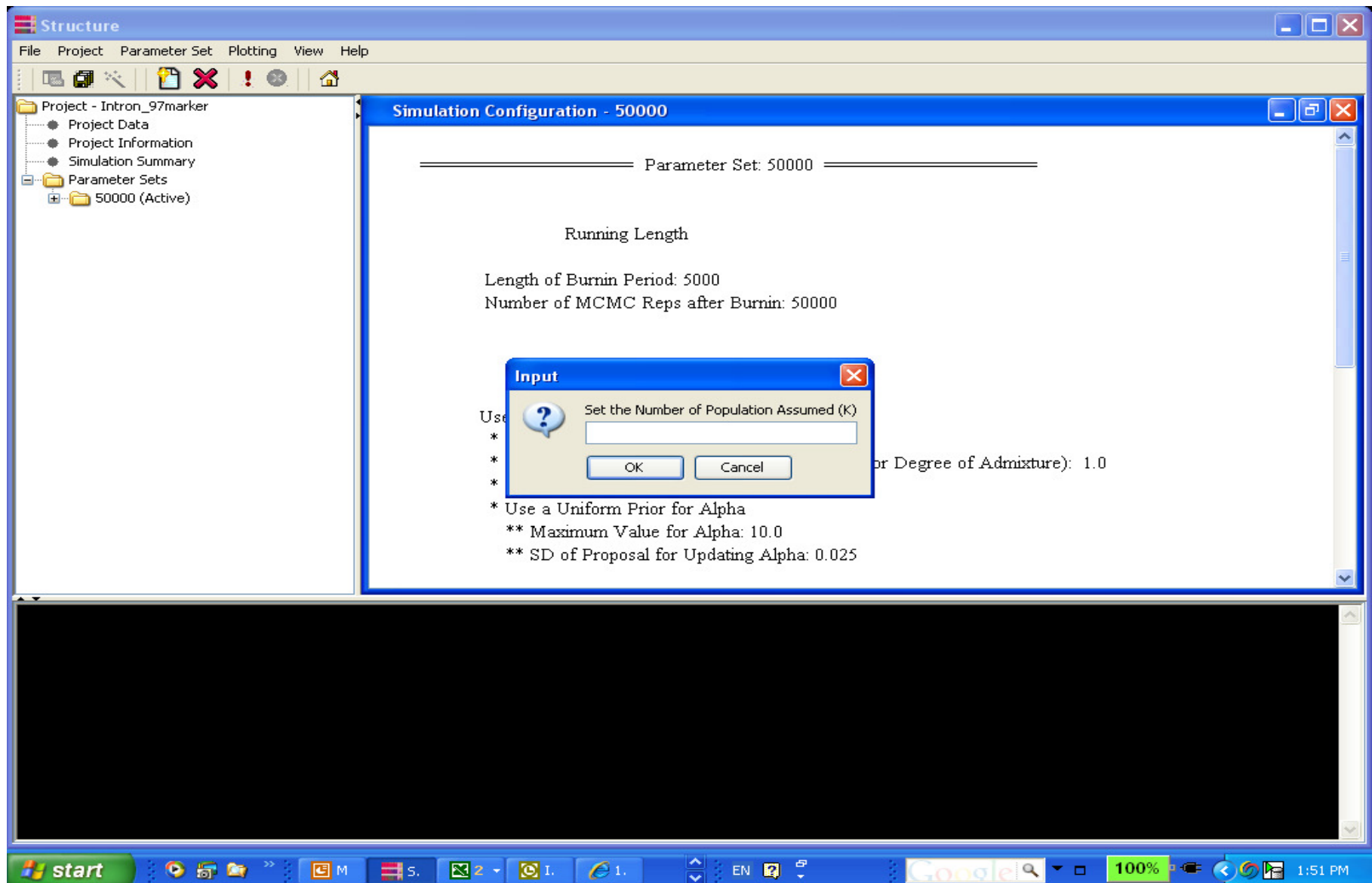
Configuring a parameter set (cont.)



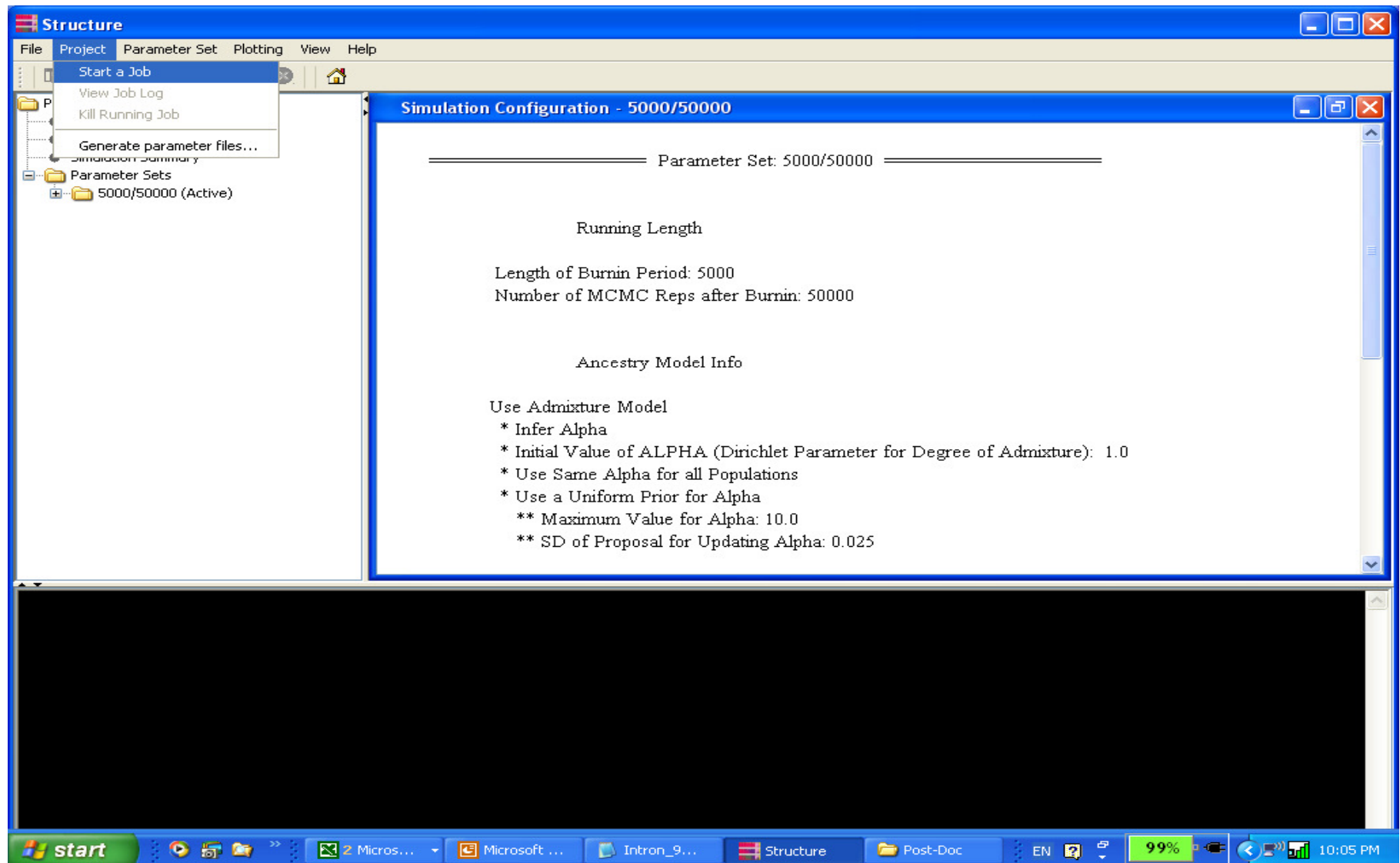
Running STRUCTURE: a single run



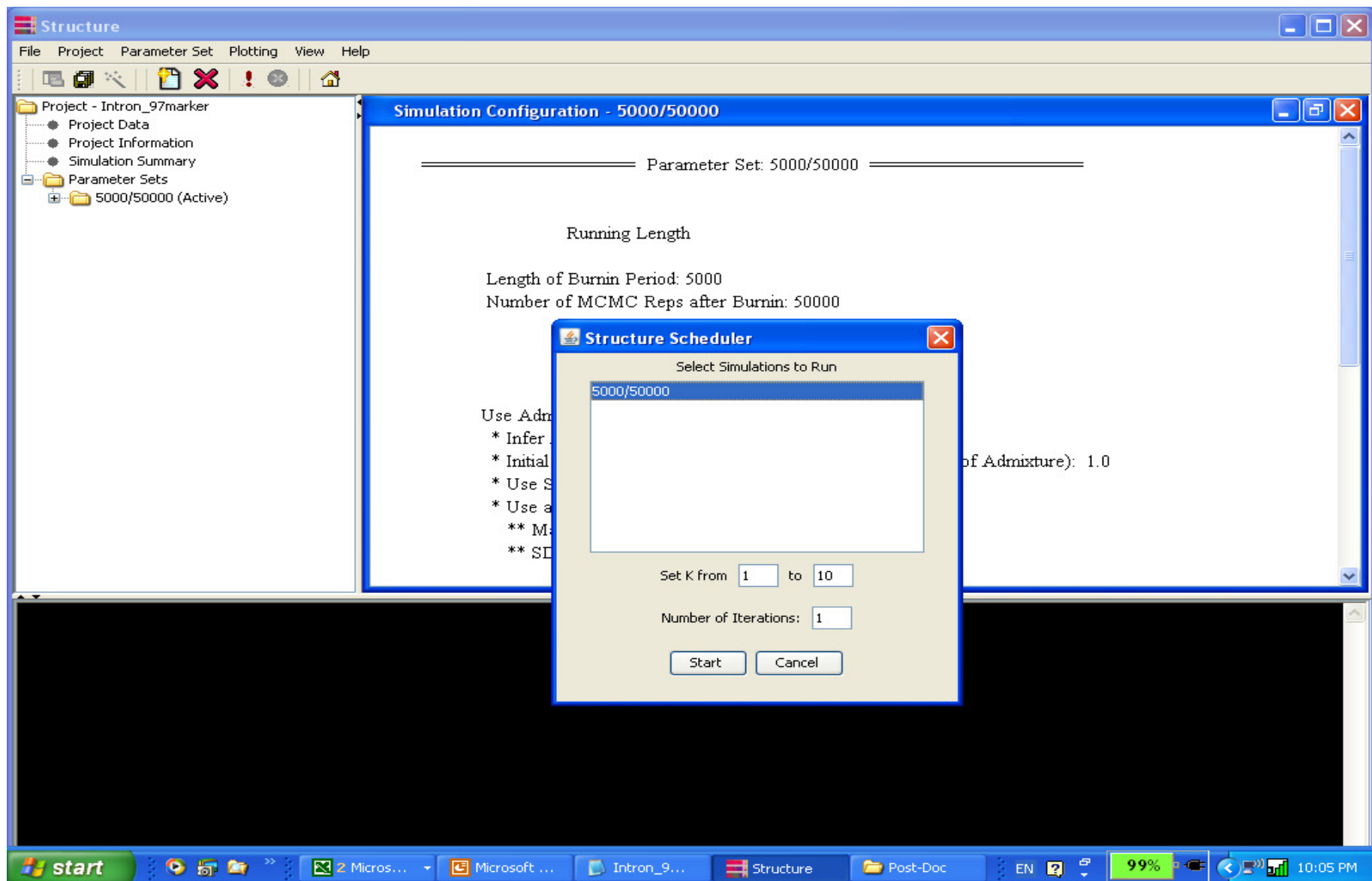
Running STRUCTURE: a single run (cont.)



Running STRUCTURE: a batch run



Running STRUCTURE: a batch run (cont.)



Structure

File Project Parameter Set Plotting View Help

Project - Intron_97marker

- Project Data
- Project Information
- Simulation Summary
- Parameter Sets
 - 50000 (Active)
 - Settings
 - Results
 - 50000_run_1 (K=1)
 - 50000_run_2 (K=2)
 - 50000_run_3 (K=3)
 - 50000_run_4 (K=4)
 - 50000_run_5 (K=5)
 - 50000_run_6 (K=6)
 - 50000_run_7 (K=7)
 - 50000_run_8 (K=8)
 - 50000_run_9 (K=9)
 - 50000_run_10 (K=10)

Summary of Project Intron_97marker

File

Summary of Simulations

Paramete...	Run Name	K	Ln P(D)	Var[LnP(D)]	α1	Fst_1	Fst_2	Fst_3	Fst_4
50000	50000_run_1	1	-3488.9	37.2	-	0.0108	-	-	-
50000	50000_run_2	2	-2856.2	126.2	0.0625	0.4607	0.4006	-	-
50000	50000_run_3	3	-2543.1	222.9	0.0566	0.6034	0.5273	0.4799	-
50000	50000_run_4	4	-2259.5	268.8	0.0556	0.6209	0.7202	0.5863	0.5201
50000	50000_run_5	5	-2170.8	398.6	0.0557	0.7496	0.7331	0.6807	0.6257
50000	50000_run_6	6	-2109.1	432.9	0.0439	0.8151	0.6628	0.7395	0.7836
50000	50000_run_7	7	-2137.4	606.9	0.0394	0.7687	0.8512	0.7599	0.6162
50000	50000_run_8	8	-2627.7	1561.0	0.0358	0.8267	0.7914	0.2132	0.7008
50000	50000_run_9	9	-2236.0	790.3	0.0330	0.8686	0.1899	0.7721	0.7099
50000	50000_run_10	10	-2173.0	808.3	0.0350	0.8679	0.7649	0.7080	0.2204

Ln P(D): Estimated probability of Ks

Proportion of membership of each pre-defined population in each of the 10 clusters

Given Pop	Inferred Clusters										Number of Individuals
	1	2	3	4	5	6	7	8	9	10	
1:	0.183	0.012	0.005	0.023	0.015	0.013	0.049	0.569	0.073	0.057	19
2:	0.183	0.142	0.277	0.017	0.010	0.004	0.274	0.007	0.069	0.016	28
3:	0.584	0.009	0.003	0.003	0.007	0.004	0.006	0.010	0.003	0.372	19
4:	0.149	0.005	0.004	0.079	0.069	0.078	0.016	0.016	0.017	0.566	4

Final results printed to file C:\Post-Doc\Population structure\Intron_97marker\Intron_97marker\50000\Results\50000_run_10_f

start

Microsoft PowerPoint ... Structure

EN 99% 11:33 PM

Inference of true K (number of populations)

🍌 The log likelihood for each K, $\ln P(D) = L(K)$

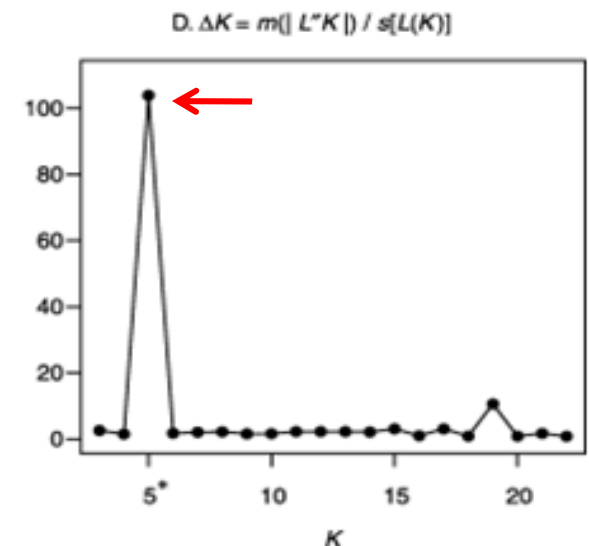
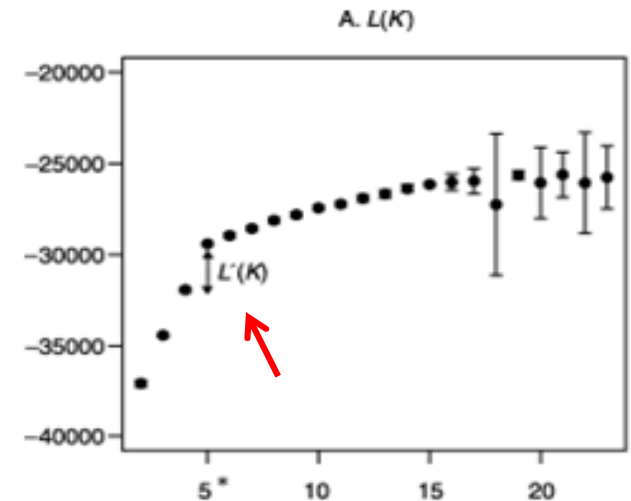
🍌 Two approaches to determine the best K

1. Use of $L(K)$: When K is approaching a true value, $L(K)$ plateaus (or continues increasing slightly) and has high variance between runs (Rosenberg et al. 2001, Evanno et al. 2005).

⇒ **Nonparametric Kruskal-Wallis test**

2. Use of an ad hoc quantity (ΔK): Calculated based on the second order rate of change of the likelihood (ΔK) (Evanno et al. 2005). The ΔK shows a clear peak at the true value of K.

⇒ $\Delta K = m([L''K])/s[L(K)]$



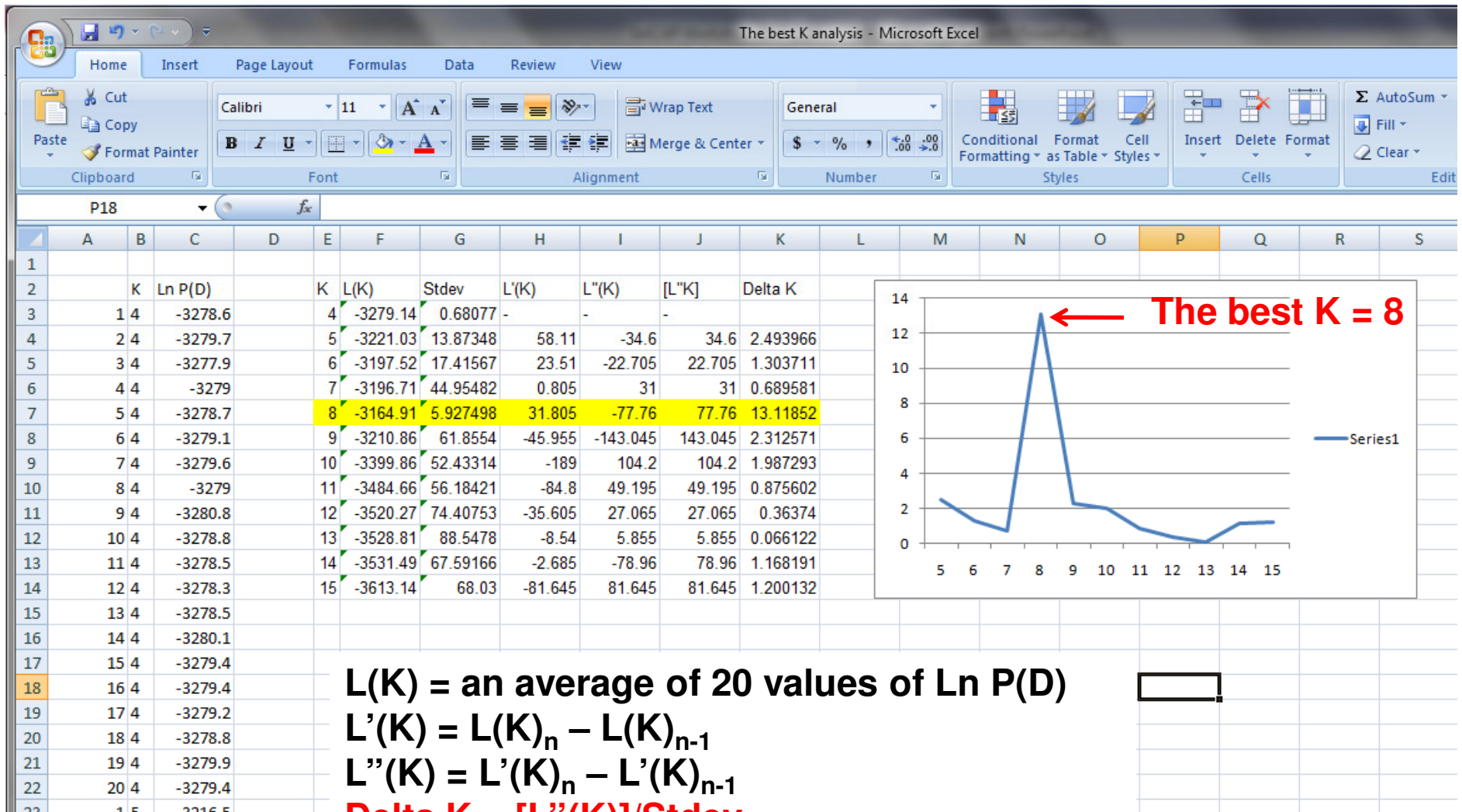
SAS code for nonparametric Kruskal-Wallis test

The screenshot shows a Microsoft Excel spreadsheet titled "Wilcoxontest_Intron&ESTmarker (8-7-07).xls". The spreadsheet contains SAS code and data for a Wilcoxon test. The code is as follows:

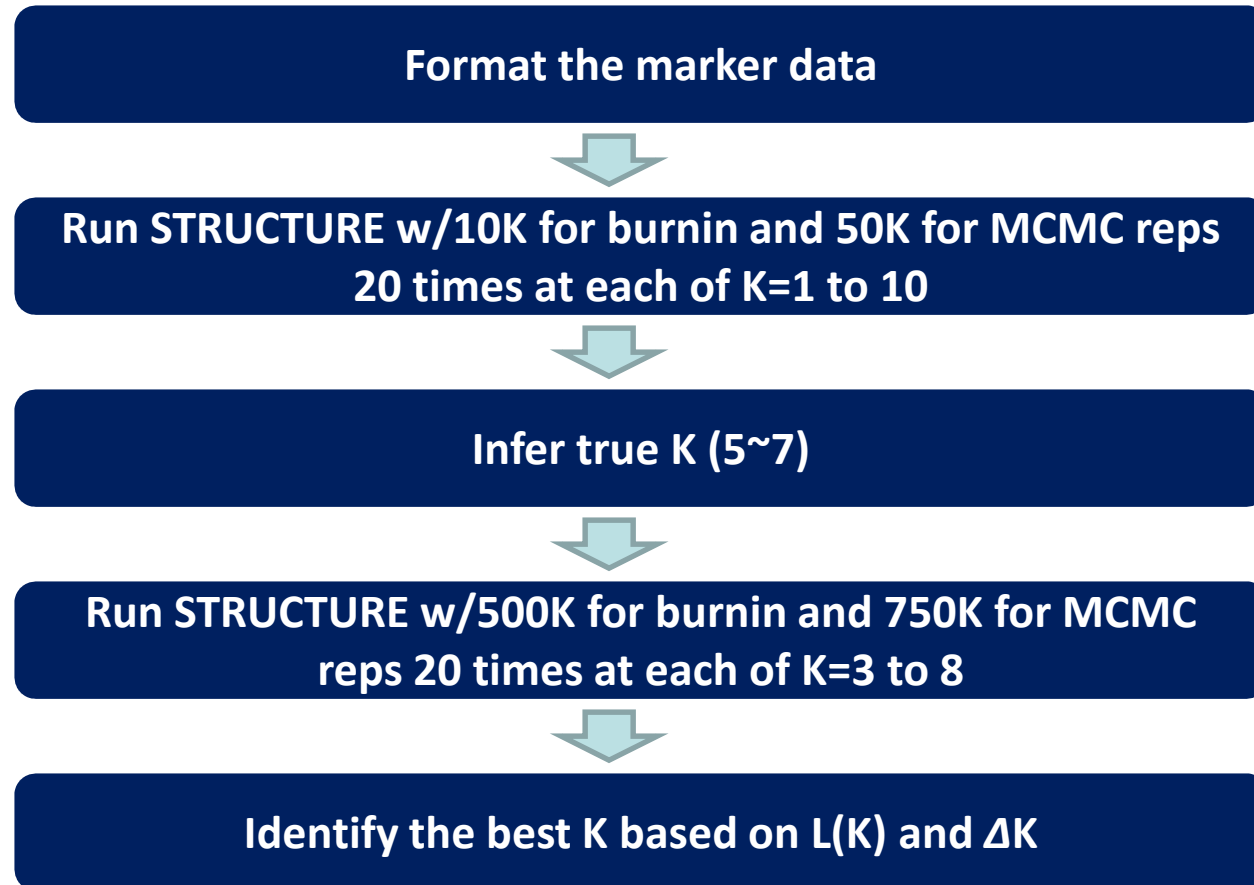
```
1 data tomato1;  
2 input K Ln @@;  
3 datalines;  
4 4 -1879.5 4 -1827.3 4 -1889.5 4 -1957.8 4 -1933.5 4 -1927.4 4 -1881.2 4 -1873.5 4 -1941.2 4 -1940.7  
5 4 -1908.7 4 -1954.0 4 -1931.5 4 -1903.8 4 -1927.0 4 -1923.1 4 -1866.5 4 -1826.1 4 -1920.0 4 -1869.8  
6 5 -1766.9 5 -1826.4 5 -1758.7 5 -1758.0 5 -1813.2 5 -1761.3 5 -1755.4 5 -1756.2 5 -2015.5 5 -1783.4  
7 5 -1760.2 5 -1758.5 5 -1893.2 5 -1751.7 5 -1738.2 5 -1834.6 5 -1732.2 5 -1852.8 5 -1818.9 5 -1802.4  
8 ;  
9 proc npar1way wilcoxon data=tomato1;  
10 class K;  
11 var Ln;  
12 exact;  
13 run;  
14  
15  
16 data tomato2;  
17 input K Ln @@;  
18 datalines;  
19 5 -1766.9 5 -1826.4 5 -1758.7 5 -1758.0 5 -1813.2 5 -1761.3 5 -1755.4 5 -1756.2 5 -2015.5 5 -1783.4  
20 5 -1760.2 5 -1758.5 5 -1893.2 5 -1751.7 5 -1738.2 5 -1834.6 5 -1732.2 5 -1852.8 5 -1818.9 5 -1802.4  
21 6 -1665.3 6 -1661.7 6 -1658.5 6 -1675.0 6 -2650.8 6 -1705.5 6 -1640.3 6 -1812.6 6 -1670.6 6 -2047.9  
22 6 -1668.8 6 -1659.3 6 -1664.9 6 -1667.9 6 -1671.5 6 -1638.1 6 -1656.1 6 -1662.7 6 -1660.1 6 -1634.9  
23 ;  
24 proc npar1way wilcoxon data=tomato2;  
25 class K;  
26 var Ln;  
27 exact;  
28 run;  
29  
30  
31 data tomato3;
```

The data is organized into three groups (tomato1, tomato2, tomato3) with columns for K (class) and Ln (var). The bottom of the spreadsheet shows the Windows taskbar with the Start button and several open applications: Microsoft PowerPoint, Structure, OSU Inbox - Mi..., Microsoft Excel..., and a battery indicator showing 96%.

Inference of best K using the delta K method

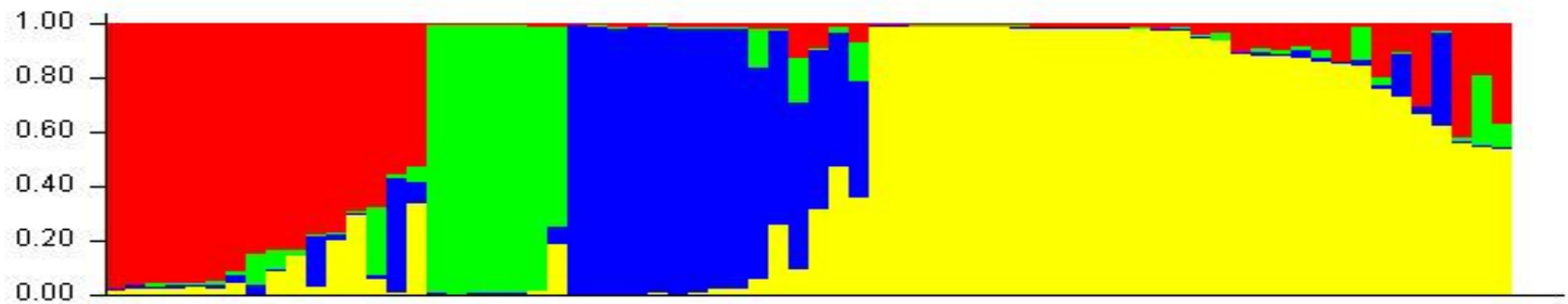


An example of steps to identify the best K



We may not always be able to know the TRUE value of K , but we should aim for the smallest value of K that captures the major structure in the data

Pritchard et al. (2000)



Enjoy running STRUCTURE

