

The Poplar Genome Project

Lieven Sterck, Stéphane Rombauts, Pierre Rouzé, and Yves Van de Peer

Department of Plant Systems Biology, Flanders Interuniversity Institute for Biotechnology (VIB), Ghent University, Technologiepark 927, B-9052 Ghent, Belgium



Introduction

Poplar trees have been used all over the world to produce a large variety of wood-based products such as timber, pulp, and paper. Besides their great economical value, poplars are also rapidly becoming the model organism for forest biology and tree biotechnology. It is therefore not surprising that, in 2001, poplar was selected as the first woody plant to have its genome sequenced.

Now that the poplar genome is publicly available, genome-wide bioinformatics analyses of the genome sequence can take a start.

Genome annotation

The annotation of the genome was performed by 3 different labs: 2 in the USA (using FgenesH and GraiL_EXP) and 1 in Belgium (using Eugene). From these 3 different annotations one consensus annotation was constructed using a rule-based system that for a given locus selected the most plausible gene-model. Before the post-processing steps (like removing putative transposons) we ended up with +- 58000 genes.

method	gene count	
Eugene	25018	43.11%
FgenesH	22546	38.85%
GraiL_exp	7352	12.66%
genewise	3120	5.38%

Table 1. Contributions of the different gene finders to the final annotation

Table 1 gives an overview of the selection of genes from the different gene finders. It is clear that our tool performs better than all others used in this annotation effort.

Although the genome is not yet published, the annotations can be viewed and queried on the JGI Genome Browser.



Figure 1. Screenshot of the Poplar Genome Browser

The Genome Browser is available at : <http://genome.jgi-psf.org/Poptr1/Poptr1.home.html>

Exploring the duplication past of poplar

For all paralogous genes in the poplar genome a Ks-estimation was performed. After corrections and filtering a histogram was constructed by plotting the number of paralogs against their age (Ks). As can be seen on the histogram there are many paralogs with the similar Ks-values (0.20-0.30) which indicate the fact that a large number of genes were created in a short period of time. The most plausible explanation for this observation is that a genome duplication has occurred in the poplar genome.

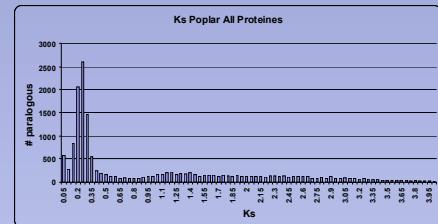


Figure 1. Ks age distribution of the paralogous genes in the poplar genome

We performed the same procedure on ESTs of 7 different Poplar species and came to the same conclusion for all of them. This suggests that the genome duplication has occurred before the divergence of all poplar species.

The most striking result of this analysis is not the fact that there was a genome duplication, but rather when it has happened. When we infer the time in My from the Ks distribution we dated the duplication event some 8 My ago, which clearly contradicts with the fossil records on poplar.

This indicates that the average Ks-substitution rate used for plants is probably not applicable to poplar. Due to its much longer generation time it is very likely that poplar has a much smaller Ks-substitution rate than all other plants studied so far.

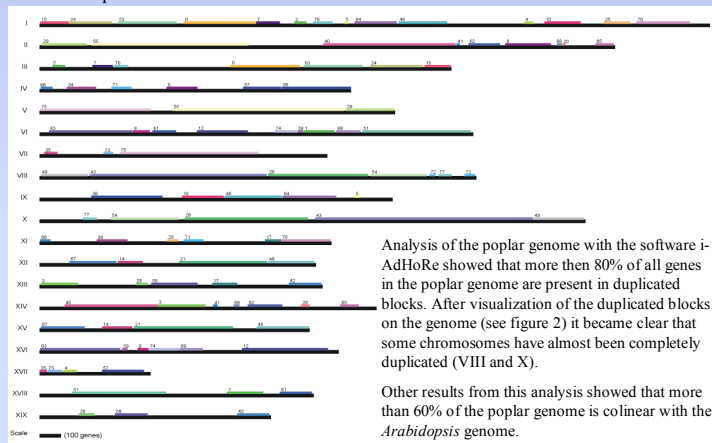


Figure 2. Map with duplicated blocks in the poplar genome

Analysis of the poplar genome with the software i-AdHoRe showed that more than 80% of all genes in the poplar genome are present in duplicated blocks. After visualization of the duplicated blocks on the genome (see figure 2) it became clear that some chromosomes have almost been completely duplicated (VIII and X).

Other results from this analysis showed that more than 60% of the poplar genome is colinear with the *Arabidopsis* genome.

Gene family comparison

TF class	Poplar	Arabidopsis	Oryza
01	114	114	114
02	416	283	283
03	207	207	207
04	207	207	207
05	143	143	143
06	134	134	134
07	30	30	30
08	104	104	104
09	73	73	73
10	82	82	82
11	69	69	69
12	24	24	24
13	49	49	49
14	29	29	29
15	27	27	27
16	24	24	24
17	24	24	24
18	24	24	24
19	24	24	24
20	24	24	24
21	24	24	24
22	24	24	24
23	24	24	24
24	24	24	24
25	24	24	24
26	24	24	24
27	24	24	24
28	24	24	24
29	24	24	24
30	24	24	24
31	24	24	24
32	24	24	24
33	24	24	24
34	24	24	24
35	24	24	24
36	24	24	24
37	24	24	24
38	24	24	24
39	24	24	24
40	24	24	24
41	24	24	24
42	24	24	24
43	24	24	24
44	24	24	24
45	24	24	24
46	24	24	24
47	24	24	24
48	24	24	24
49	24	24	24
50	24	24	24
51	24	24	24
52	24	24	24
53	24	24	24
54	24	24	24
55	24	24	24
56	24	24	24
57	24	24	24
58	24	24	24
59	24	24	24
60	24	24	24
61	24	24	24
62	24	24	24
63	24	24	24
64	24	24	24
65	24	24	24
66	24	24	24
67	24	24	24
68	24	24	24
69	24	24	24
70	24	24	24
71	24	24	24
72	24	24	24
73	24	24	24
74	24	24	24
75	24	24	24
76	24	24	24
77	24	24	24
78	24	24	24
79	24	24	24
80	24	24	24
81	24	24	24
82	24	24	24
83	24	24	24
84	24	24	24
85	24	24	24
86	24	24	24
87	24	24	24
88	24	24	24
89	24	24	24
90	24	24	24
91	24	24	24
92	24	24	24
93	24	24	24
94	24	24	24
95	24	24	24
96	24	24	24
97	24	24	24
98	24	24	24
99	24	24	24
100	24	24	24

Table 1. No. of proteins belonging to a specific TF-class.

Table 1 shows a comparison between Poplar and *Arabidopsis* of the number of transcription factors. For most of the TFs the numbers are comparable although there are some differences, eg. the Myb TFs, which seem to have undergone lineage specific expansion in poplar. Also the pathogen-related TFs are more abundant in Poplar, which is perhaps correlated with the perennial nature of Poplar.

Repetitive sequences and transposons

- Repeats occupy 28% of the assembled genome and more than 50% of the unassembled scaffolds.
- Small scaffolds and unassembled reads contained substantially more repetitive sequence than large scaffolds.

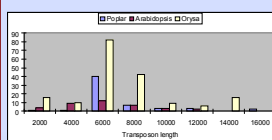
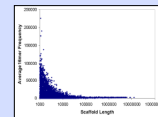


Figure 1. Comparison of length transposons

When we compare the length and frequency of LTR-transposons in *Arabidopsis*, poplar and rice it is clear that in rice there are more and longer transposons, although the longest ones are identified in poplar.

This comparison included 37 Ath, 58 poplar and 181 rice transposons.

Conclusions

- Outstanding international collaboration has placed poplar on an equivalent footing with other model species
- Assembly was heavily complicated by repeats, but facilitated by genetic maps
- Unprecedented opportunities to explore molecular basis of commercial and environmental traits
- Most of the hard work and many exciting discoveries are still to come!

References:

Salamov A., Solovyev V. (2000) **EgenesH**: Ab initio gene finding in *Drosophila* genomic DNA. *Genome Res.*, 10,516-522

Schiex T., Moisan, A. and Rouzé, P. (2001) **EnGene** : An Eucaryotic Gene Finder that combines several sources of evidence. *Computational Biology*, Eds. O. Gascuel and M-F. Sagot, LNCS 2066, pp. 111-125

Simillion, C., Vandepoel, K., Saey, Y. & Van de Peer, Y. (2004) Building genomic profiles for uncovering segmental homology in the twilight zone. *Genome Res.* 14, 1095-1106

Sterck, L., Rombauts, S., Jansson, S., Sterck, F., Rouzé, P. & Van de Peer, Y. (2005) EST data suggest that poplar is an ancient polyploid. *New Phytol.* (In Press)

Yang Xu, Manesh Shah, Doug Hyatt, Richard Mural, and Edward C. Uberbacher. **GRAIL-EXP**. Multiple Gene Modeling Using Pattern Recognition and Homology. *The Seventh Department of Energy Contractor and Grantee Workshop, January 1999*