# Red clover (*Trifolium pratense* L.) draft genome provides a platform for trait improvement

Jose J. De Vega[1,2*], Sarah Ayling[1], Matthew Hegarty[2], Dave Kudrna[3], Jose L. Goicoechea[3], Åshild Ergon[4], Odd A. Rognli[4], Charlotte Jones[2], Martin Swain[2], Rene Geurts[5], Chunting Lang[5], Klaus F. X. Mayer[6], Stephan Rössner[6], Steven Yates[2,7], Kathleen J. Webb[2], Iain S. Donnison[2], Giles E. D. Oldroyd[8], Rod A. Wing[3], Mario Caccamo[1], Wayne Powell[2,9], Michael T. Abberton[2,10], Leif Skøt[2*]

[1]The Genome Analysis Centre, Norwich Research Park, Norwich, NR4 7UH, UK

[2]Institute of Biological, Environmental and Rural Sciences, Aberystwyth University, Gogerddan, Aberystwyth, Ceredigion SY23 3EB, UK

[3]Arizona Genomics Institute, School of Plant Sciences, University of Arizona, Tucson AZ 85721, USA

[4]Norwegian University of Life Sciences, Department of Plant Sciences, N-1432, Ås, Norway

[5]Laboratory of Molecular Biology, Department of Plant Science, Wageningen University, Droevendaalsesteeg 1, 6708PB Wageningen, The Netherlands

[6]MIPS/Institute for Bioinformatics and Systems Biology, Helmholtz Center Munich, Ingolstädter Landstrasse 1, Neuherberg, Germany

[7]Forage Crop Genetics, Institute of Agricultural Sciences, ETH Zurich, CH-8092, Zurich, Switzerland

[8]Department of Disease and Stress Biology, John Innes Centre, Norwich NR4 7UH, UK

[9]CGIAR Consortium Office 1000, Avenue Agropolis, F-34394, Montpellier, Cedex 5, France

[10]International Institute of Tropical Agriculture (IITA), PMB 5320, Oyo Road, Ibadan, Nigeria

*e-mail: Jose.DeVega@tgac.ac.uk; lfs@aber.ac.uk

**Supplementary Material**

- Figures 1-23

- Tables 1-6

- Supplementary References

## Supplementary Figures

**Figure 1.** Cumulative length (y-axis) plot of the red clover scaffolds sorted by length (x-axis) in the final assembly produced with Platanus 1.2.1 (Black line) and an alternative assembly assembled with ABySS and scaffolds with SOAP2 (Red line). (Top) All the scaffolds, and (Bottom) the 20000 longest scaffolds.
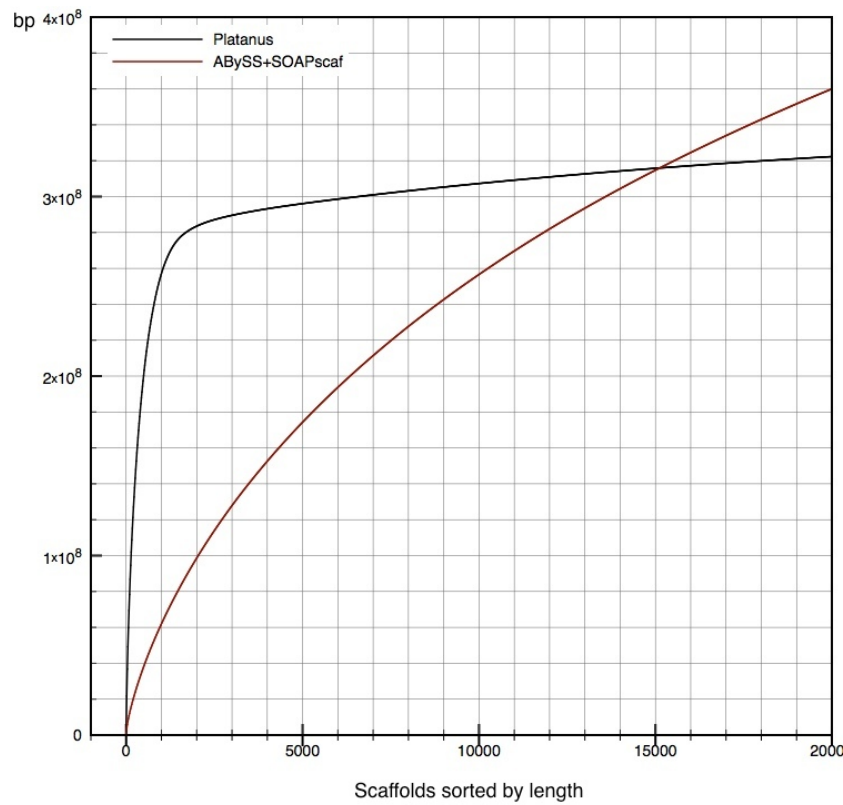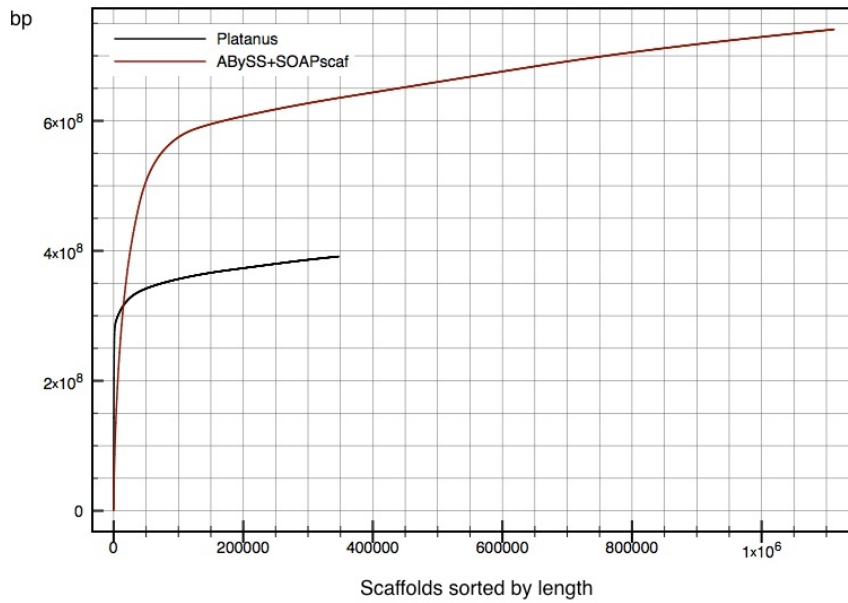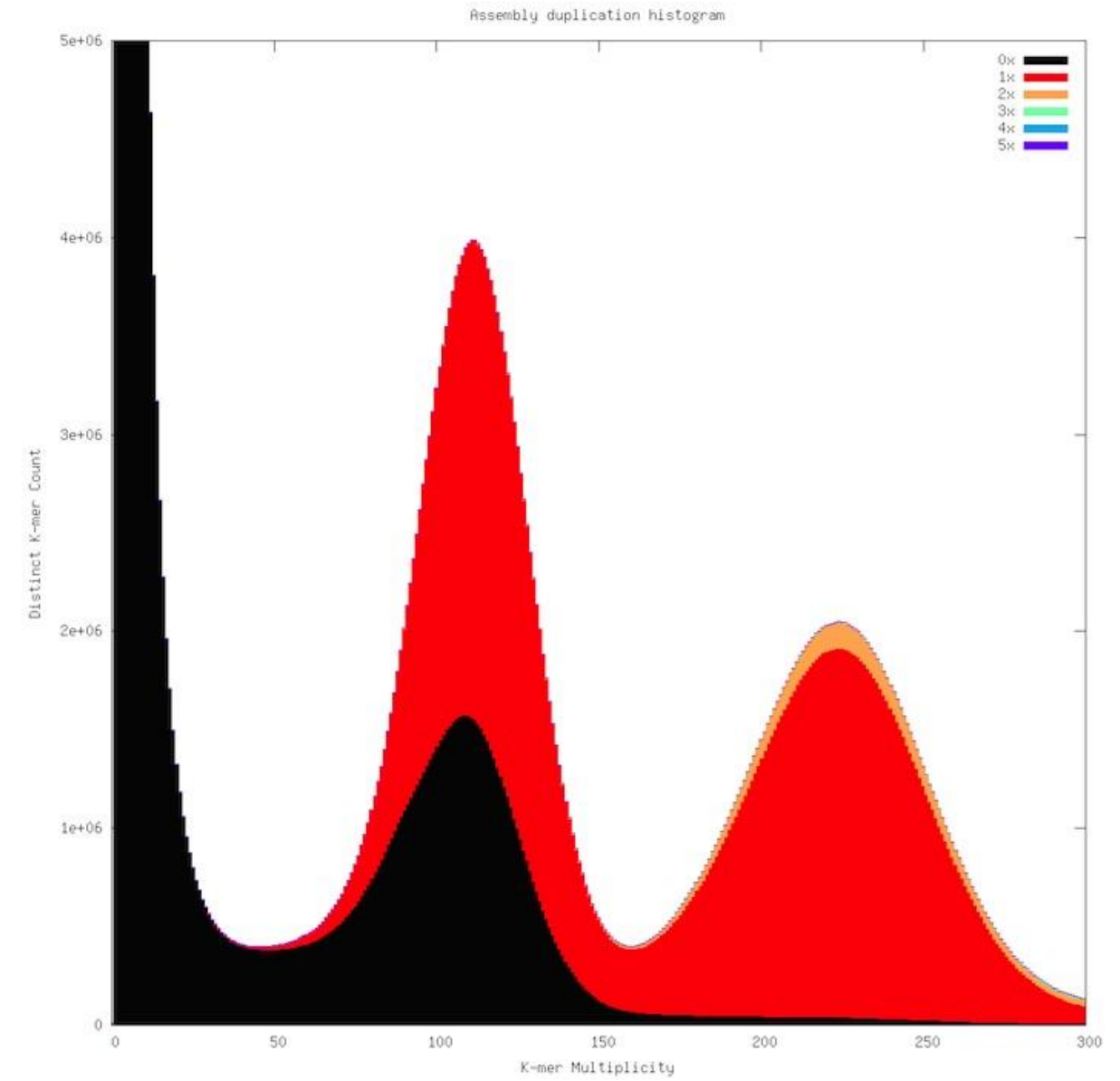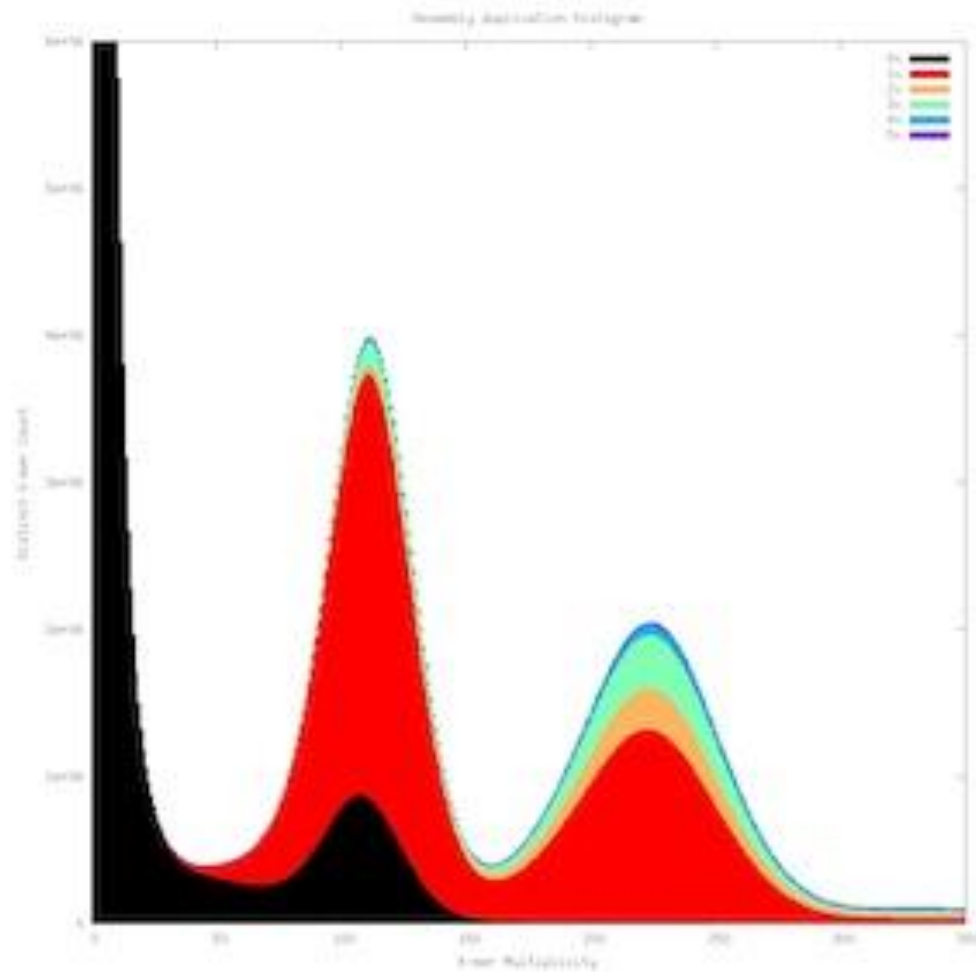
**Figure 2**. **Comparison of the Kmer spectra of four assembly strategies.** The K-mer spectra of the paired-end and single-end reads versus the different output of the assemblers. The area under the curve of the Kmer spectra has been coloured according to the number of times that such K-mers appear in the assembly: none in back, once in red, twice in orange, etc.

# Platanus (Final reference)

# ABySS

De Vega *et al.* (2015). Red clover (*Trifolium pratense* L.) draft genome provides a platform for trait improvement.

# SOAP2



Assembly duplication histogram

# SOAP2 + GAPCLOSER



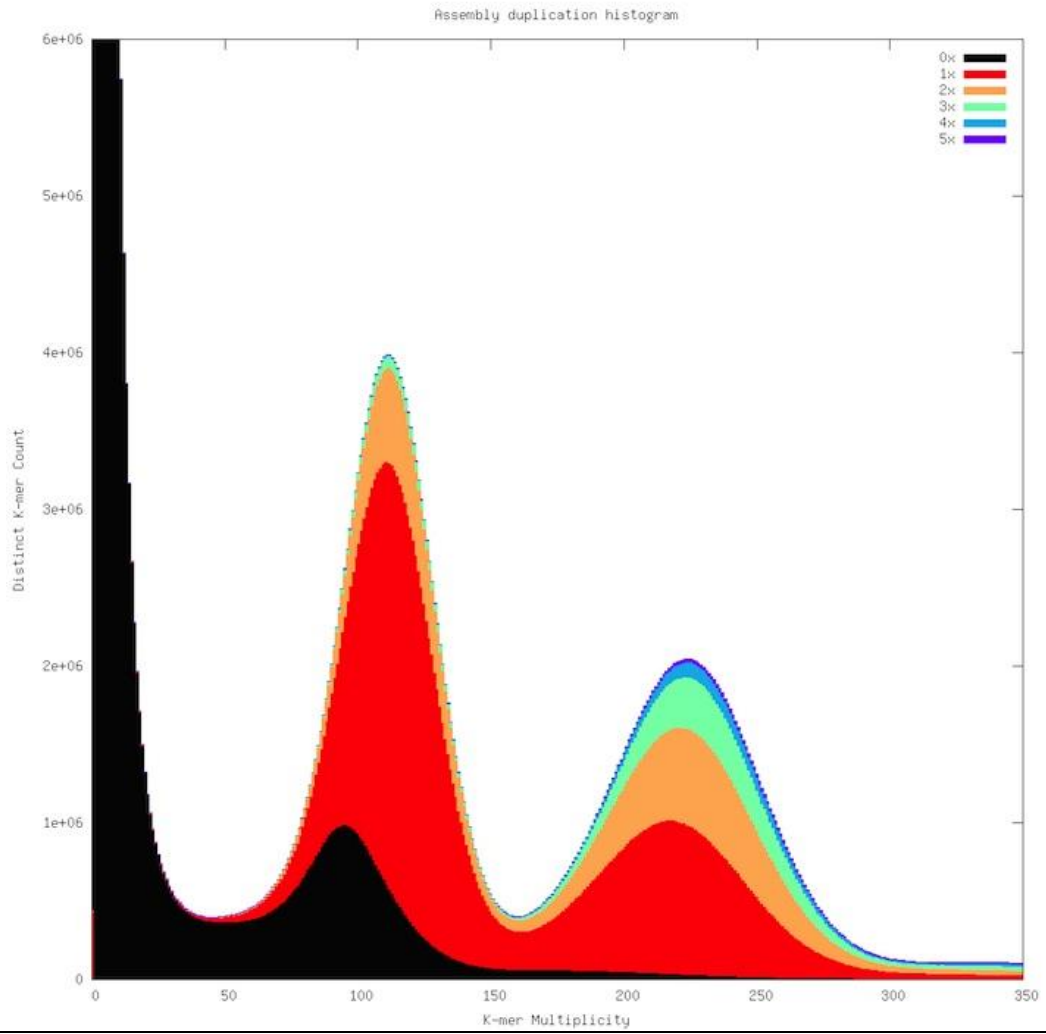Assembly duplication histogram

**Figure 3**. K-mer spectra of the red clover assembly described previously[1] (see also Results section).
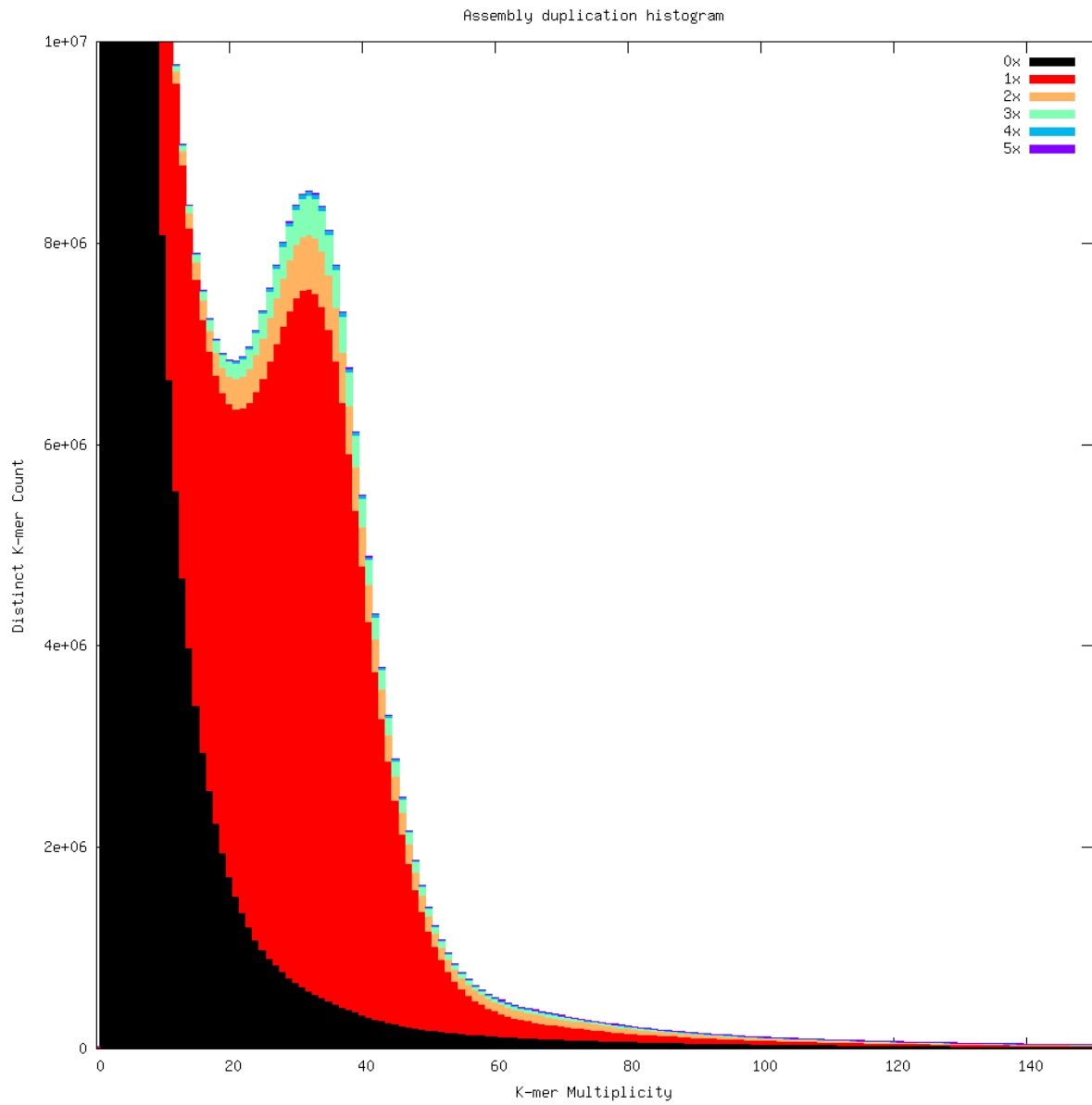
**Figure 4. Genetic map of red clover Milvus x Britta F1 population.** Genetic distances on the left of each linkage group (LG) are in cM.
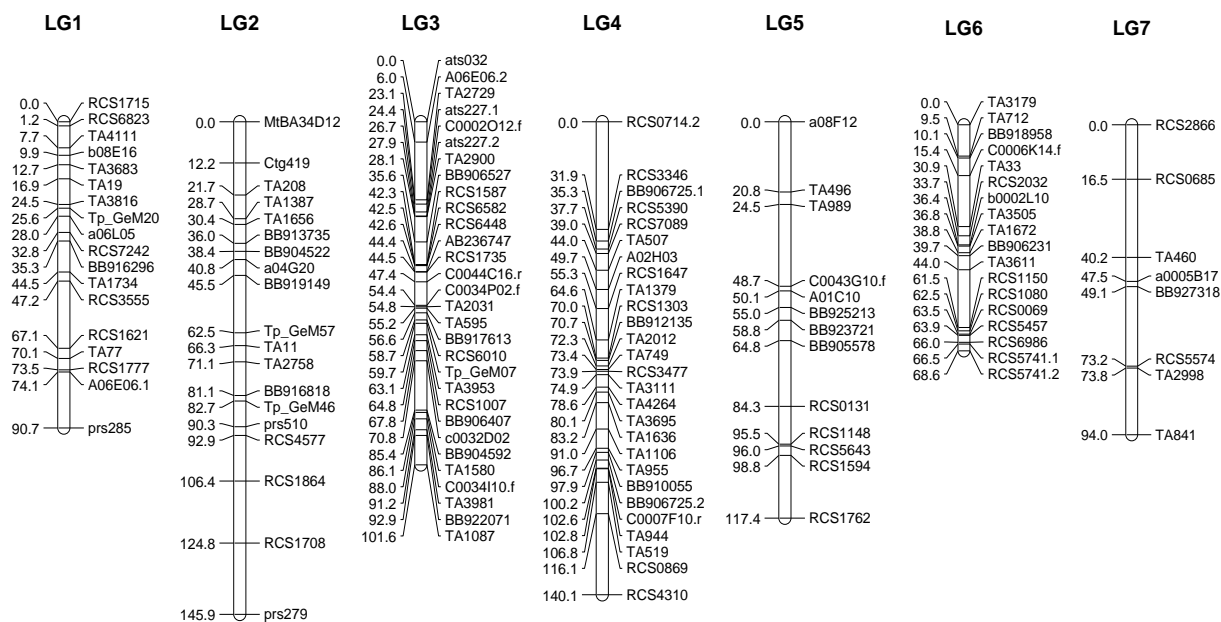
**LG1**

| 0.0 | RCS1715 |
| 1.2 | RCS6823 |
| 7.7 | TA4111 |
| 9.9 | b08E16 |
| 12.7 | TA3683 |
| 16.9 | TA19 |
| 24.5 | TA3816 |
| 25.6 | Tp_GeM20 |
| 28.0 | a06L05 |
| 32.8 | RCS7242 |
| 35.3 | BB916296 |
| 44.5 | TA1734 |
| 47.2 | RCS3555 |
| 67.1 | RCS1621 |
| 70.1 | TA77 |
| 73.5 | RCS1777 |
| 74.1 | A06E06.1 |
| 90.7 | prs285 |

**LG2**

| 0.0 | MtBA34D12 |
| 12.2 | Ctg419 |
| 21.7 | TA208 |
| 28.7 | TA1387 |
| 30.4 | TA1656 |
| 36.0 | BB913735 |
| 38.4 | BB904522 |
| 40.8 | a04G20 |
| 45.5 | BB919149 |
| 62.5 | Tp_GeM57 |
| 66.3 | TA11 |
| 71.1 | TA2758 |
| 81.1 | BB916818 |
| 82.7 | Tp_GeM46 |
| 90.3 | prs510 |
| 92.9 | RCS4577 |
| 106.4 | RCS1864 |
| 124.8 | RCS1708 |
| 145.9 | prs279 |

**LG3**

| 0.0 | ats032 |
| 6.0 | A06E06.2 |
| 23.1 | TA2729 |
| 24.4 | ats227.1 |
| 26.7 | C0002O12.f |
| 27.9 | ats227.2 |
| 28.1 | TA2900 |
| 35.6 | BB906527 |
| 42.3 | RCS1587 |
| 42.5 | RCS6582 |
| 42.6 | RCS6448 |
| 44.4 | AB236747 |
| 44.5 | RCS1735 |
| 47.4 | C0044C16.r |
| 54.4 | C0034P02.f |
| 54.8 | TA2031 |
| 55.2 | TA595 |
| 56.6 | BB917613 |
| 58.7 | RCS6010 |
| 59.7 | Tp_GeM07 |
| 63.1 | TA3953 |
| 64.8 | RCS1007 |
| 67.8 | BB906407 |
| 70.8 | c0032D02 |
| 85.4 | BB904592 |
| 86.1 | TA1580 |
| 88.0 | C0034I10.f |
| 91.2 | TA3981 |
| 92.9 | BB922071 |
| 101.6 | TA1087 |

**LG4**

| 0.0 | RCS0714.2 |
| 31.9 | RCS3346 |
| 35.3 | BB906725.1 |
| 37.7 | RCS5390 |
| 39.0 | RCS7089 |
| 44.0 | TA507 |
| 49.7 | A02H03 |
| 55.3 | RCS1647 |
| 64.6 | TA1379 |
| 70.0 | RCS1303 |
| 70.7 | BB912135 |
| 72.3 | TA2012 |
| 73.4 | TA749 |
| 73.9 | RCS3477 |
| 74.9 | TA3111 |
| 78.6 | TA4264 |
| 80.1 | TA3695 |
| 83.2 | TA1636 |
| 91.0 | TA1106 |
| 96.7 | TA955 |
| 97.9 | BB910055 |
| 100.2 | BB906725.2 |
| 102.6 | C0007F10.r |
| 102.8 | TA944 |
| 106.8 | TA519 |
| 116.1 | RCS0869 |
| 140.1 | RCS4310 |

**LG5**

| 0.0 | a08F12 |
| 20.8 | TA496 |
| 24.5 | TA989 |
| 48.7 | C0043G10.f |
| 50.1 | A01C10 |
| 55.0 | BB925213 |
| 58.8 | BB923721 |
| 64.8 | BB905578 |
| 84.3 | RCS0131 |
| 95.5 | RCS1148 |
| 96.0 | RCS5643 |
| 98.8 | RCS1594 |
| 117.4 | RCS1762 |

**LG6**

| 0.0 | TA3179 |
| 9.5 | TA712 |
| 10.1 | BB918958 |
| 15.4 | C0006K14.f |
| 30.9 | TA33 |
| 33.7 | RCS2032 |
| 36.4 | b0002L10 |
| 36.8 | TA3505 |
| 38.8 | TA1672 |
| 39.7 | BB906231 |
| 44.0 | TA3611 |
| 61.5 | RCS1150 |
| 62.5 | RCS1080 |
| 63.5 | RCS0069 |
| 63.9 | RCS5457 |
| 66.0 | RCS6986 |
| 66.5 | RCS5741.1 |
| 68.6 | RCS5741.2 |

**LG7**

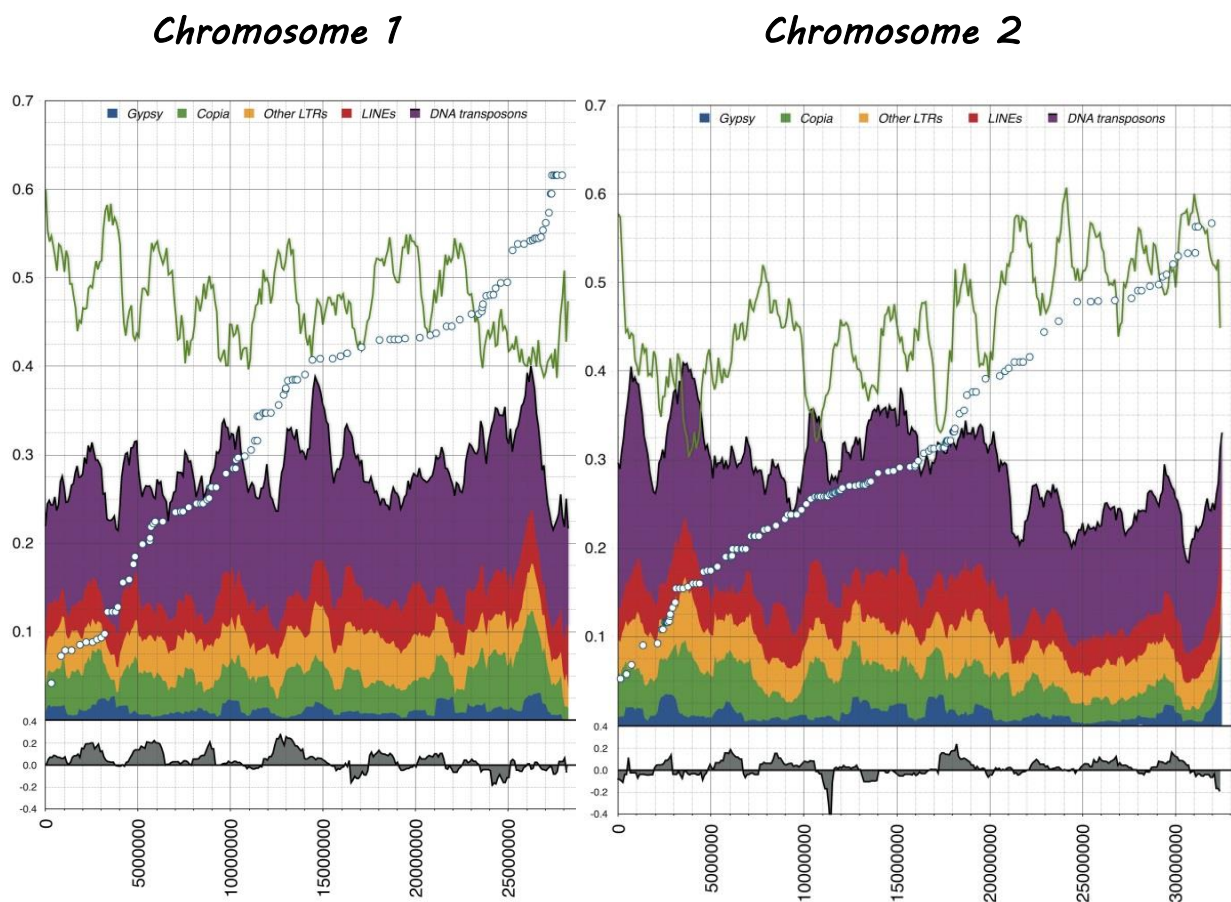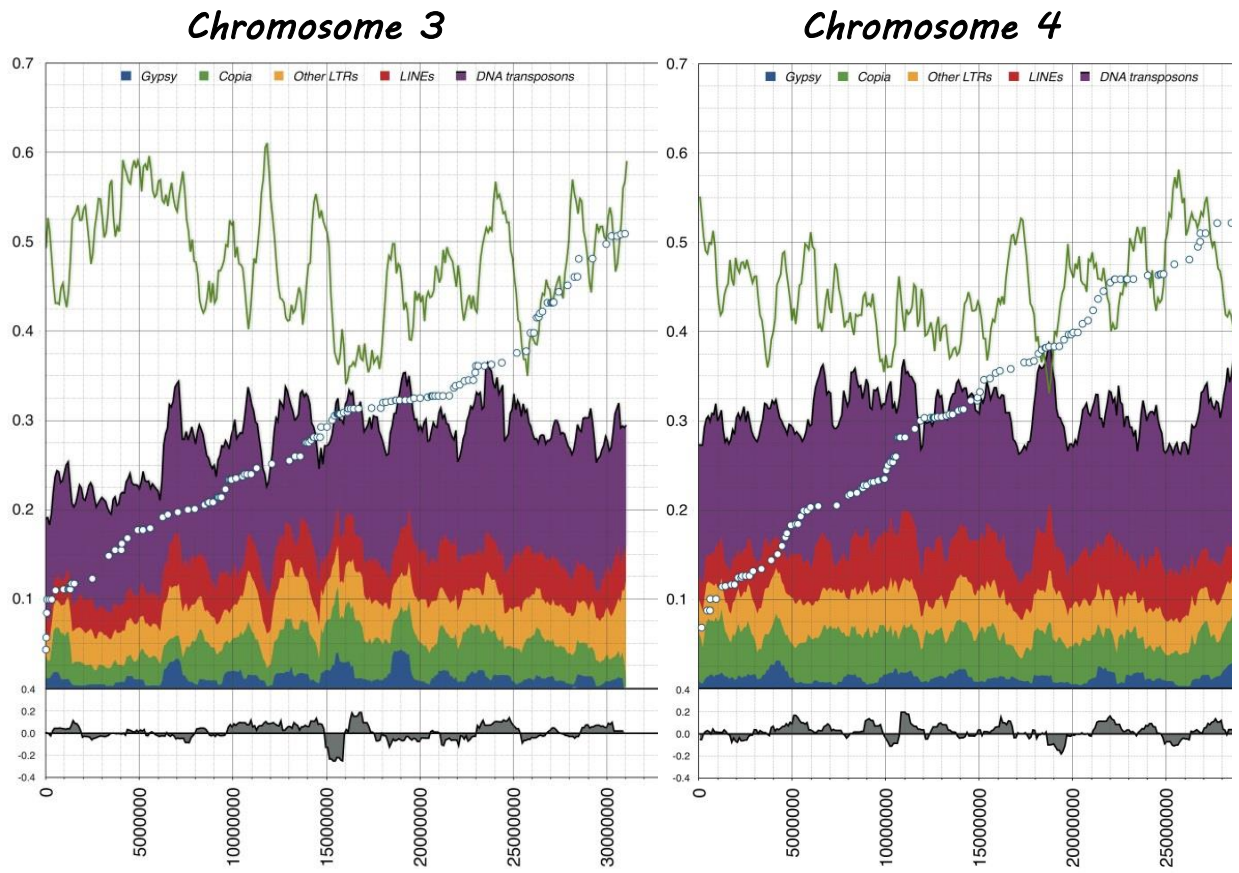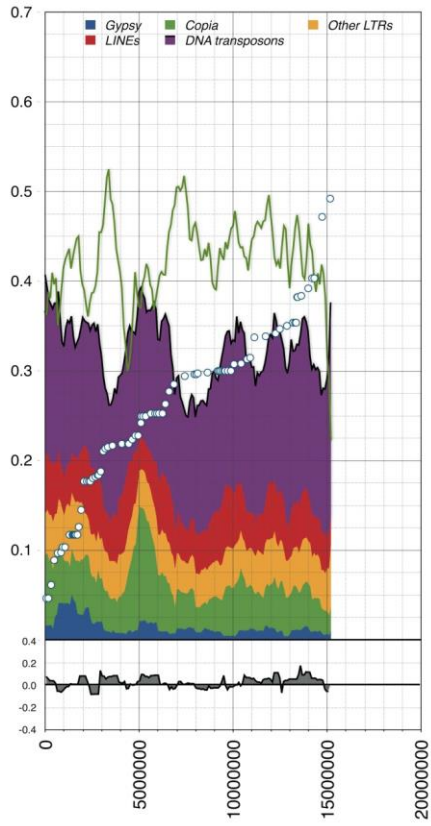| 0.0 | RCS2866 |
| 16.5 | RCS0685 |
| 40.2 | TA460 |
| 47.5 | a0005B17 |
| 49.1 | BB927318 |
| 73.2 | RCS5574 |
| 73.8 | TA2998 |
| 94.0 | TA841 |

**Figure 5. Landscapes of the red clover chromosomes.** The landscapes represent the proportion (0-1) of content along the chromosome in 10Kb intervals of each of the elements in different categories. Top panel: Density of various types of repetitive elements. Circular symbols represent genetic markers. Green line is gene density per Mb. So, a density of 0.5 means that the gene content (exons and introns) occupied 500 Kb of each 1 Mb in each window. Bottom panel: Expected minus Observed Heterozygosity.
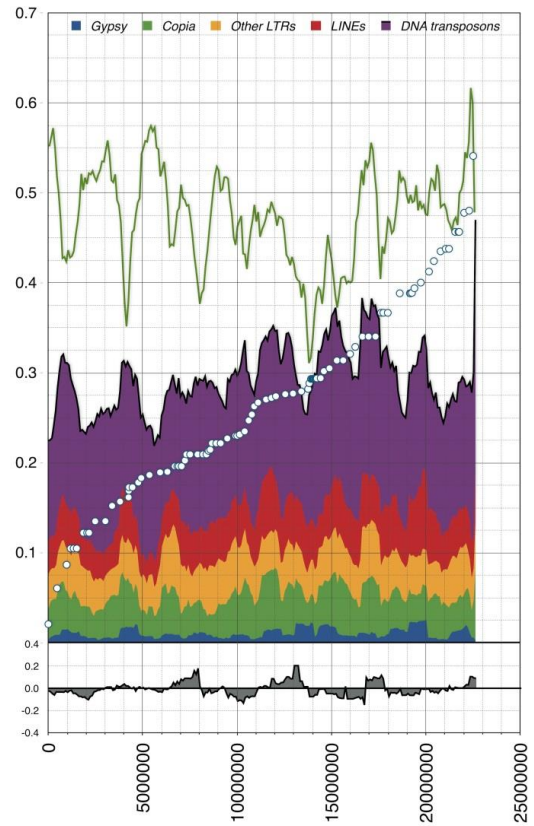
De Vega *et al.* (2015). Red clover (*Trifolium pratense* L.) draft genome provides a platform for trait improvement.

De Vega *et al.* (2015). Red clover (*Trifolium pratense* L.) draft genome provides a platform for trait improvement.



Chromosome 5

Chromosome 6

## Chromosome 7

**Figure 6. Comparison of gene features.** Length and frequency of the features (genes, transcripts, CDS and exons/introns) in the red clover (Tp), common bean (Pv), *M. truncatula* (Mt), *L. japonicus* (Lj), soybean (Gm) and *A. thaliana* (At) genomes.

**Figure 7. GO terms over-represented in red clover gene clusters that are expanded in comparison with *M. truncatula*.** (A) Scatter-plot that represents the number of genes in each cluster in *M. truncatula* (X-axis) and red clover (Y-axis). A cluster is expanded when the number of members in red clover is at least twice the number of members in *M. truncatula*. (B) GO-terms overrepresented and proportion of sequences observed in the expanded clusters (blue columns) versus expected/total (red columns). (C&D) Treemaps with the biological process GO terms (C) and Molecular activity GO terms (D). The area is proportional to the enrichment, so bigger areas represent highly overrepresented functions in red clover.
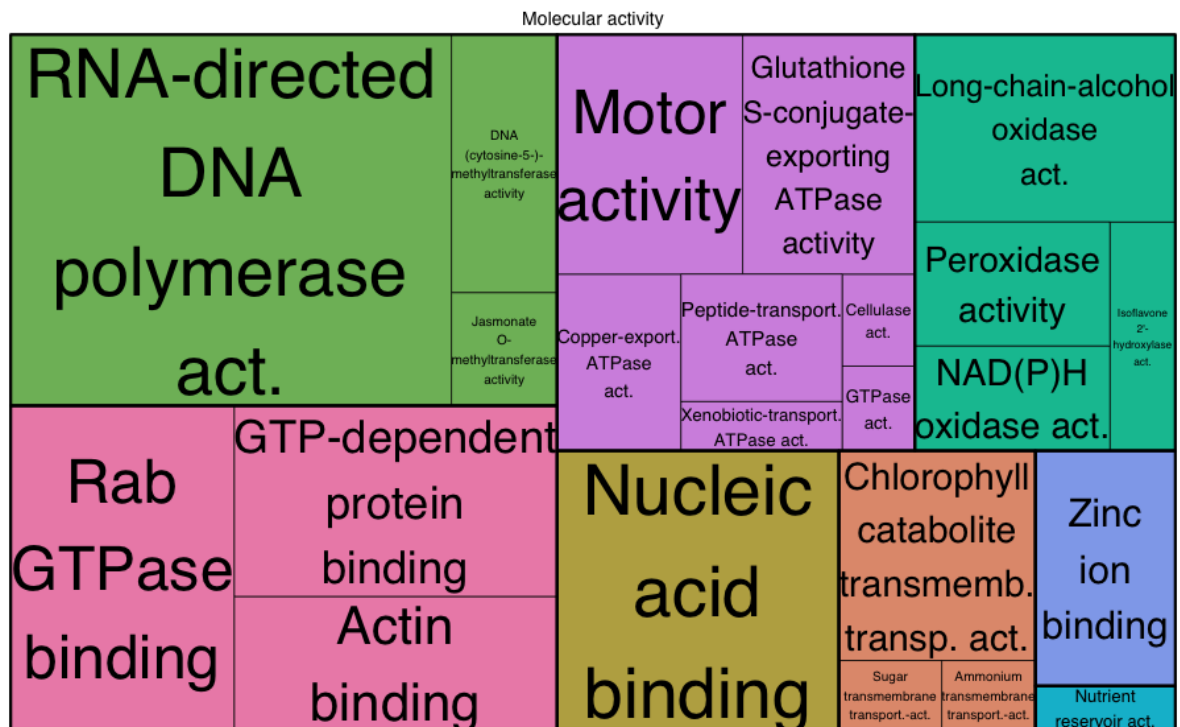
7A

De Vega *et al.* (2015). Red clover (*Trifolium pratense* L.) draft genome provides a platform for trait improvement.

7B



Differential GO-term Distribution

De Vega *et al.* (2015). Red clover (*Trifolium pratense* L.) draft genome provides a platform for trait improvement.

7C



7D

**Figure 8. Red clover gene duplication events.** Close events are represented together as thicker lines.

**Figure 9. Gene duplication events in red clover.** (A) Frequency of Gene duplication events with time (Kimura rates) between chromosomes in the red clover genome. (B) Frequency of Gypsy duplication events with time (Kimura rates) in the red clover and *M. truncatula* genomes.
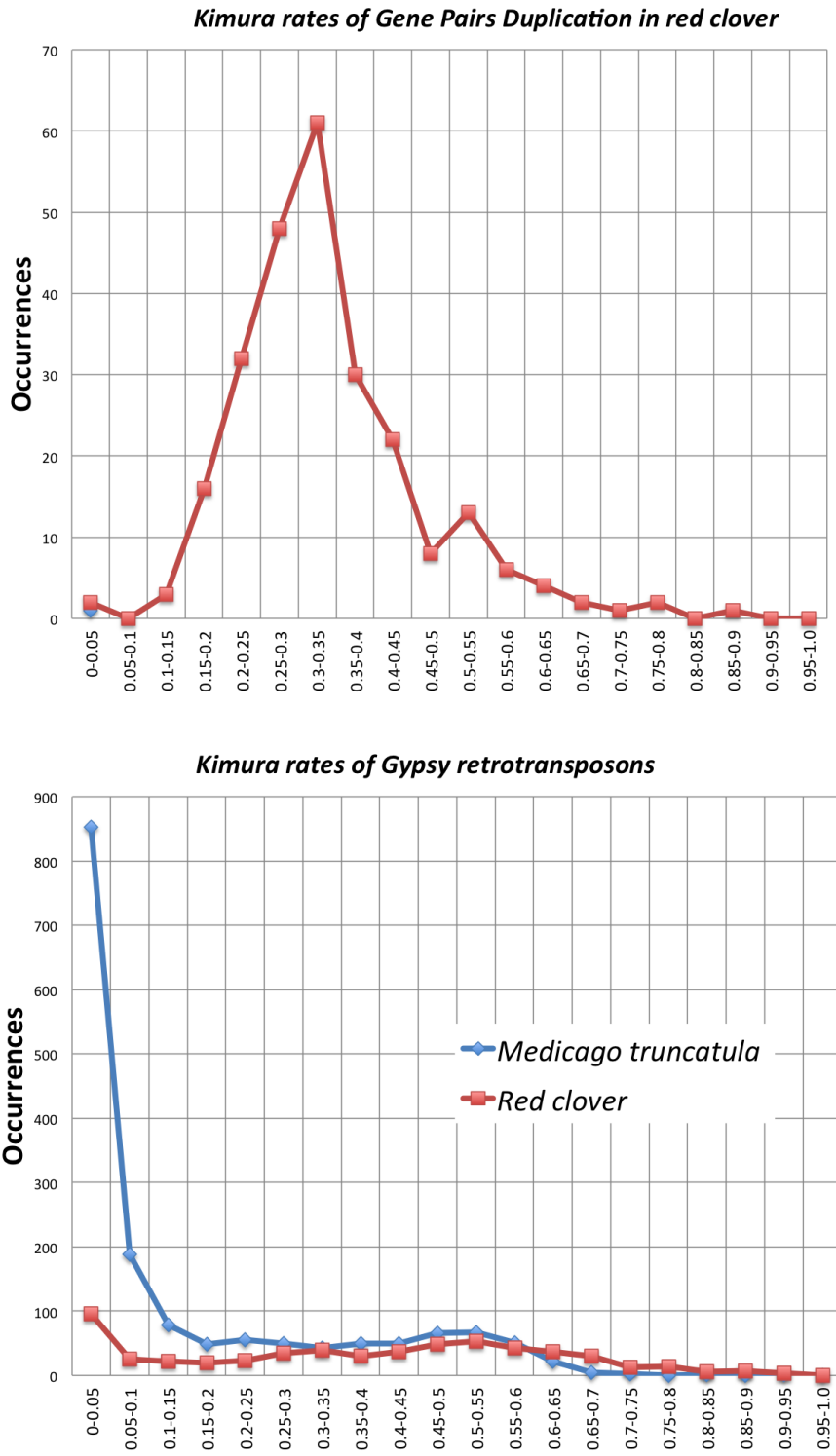
**Figure 10. Distribution of the distance from red clover genes to the closest *Gypsy* or *Copia* element, for the whole set of genes (Top panel) or duplicated gene pairs (Bottom panel).**
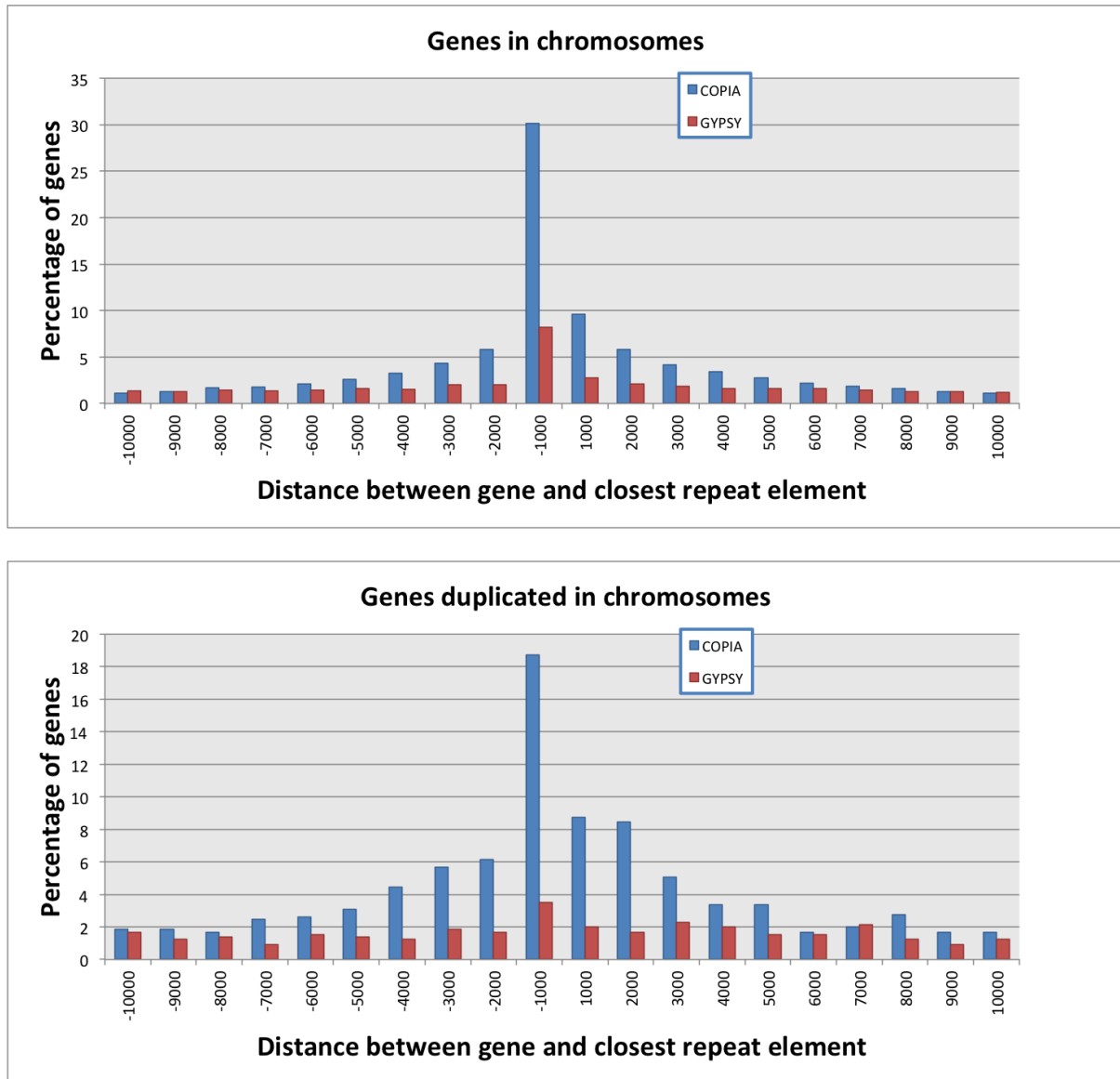
**Figure 11. Biosynthesis of formononetin in plants.** Schematic representation of the (A) formononetin biosynthesis pathway and (B) formononetin interconversion pathway.
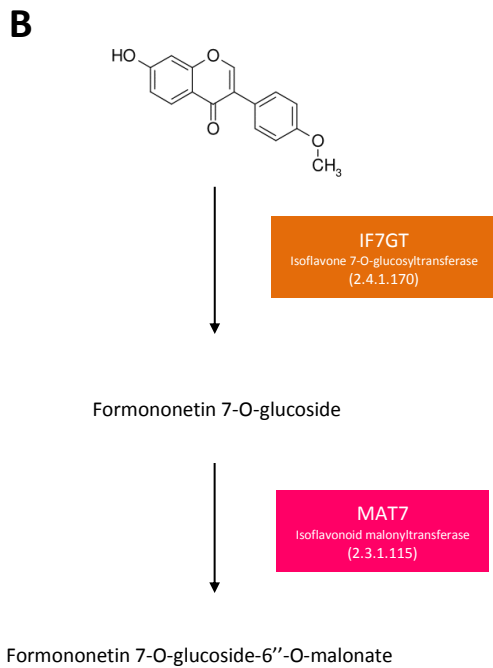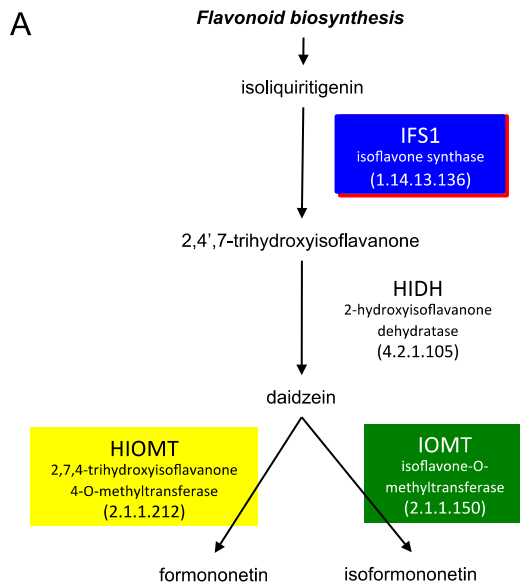
A

Flavonoid biosynthesis

isoliquiritigenin

IFS1
isoflavone synthase
(1.14.13.136)

2,4',7-trihydroxyisoflavanone

HIDH
2-hydroxyisoflavanone
dehydratase
(4.2.1.105)

daidzein

HIOMT
2,7,4-trihydroxyisoflavanone
4-O-methyltransferase
(2.1.1.212)

IOMT
isoflavone-O-
methyltransferase
(2.1.1.150)

formononetin          isoformononetin

B

IF7GT
Isoflavone 7-O-glucosyltransferase
(2.4.1.170)

Formononetin 7-O-glucoside

MAT7
Isoflavonoid malonyltransferase
(2.3.1.115)

Formononetin 7-O-glucoside-6''-O-malonate

**Figure 12. Phylogenetic tree of 2-hydroxyisoflavanone-dehydratase (HIDH)** (4.2.1.105) in red clover, soybean (Glyma), common bean (Phvu) and *M. truncatula* (Medtr). The outgroup is represented by *A. thaliana* carboxylesterase.
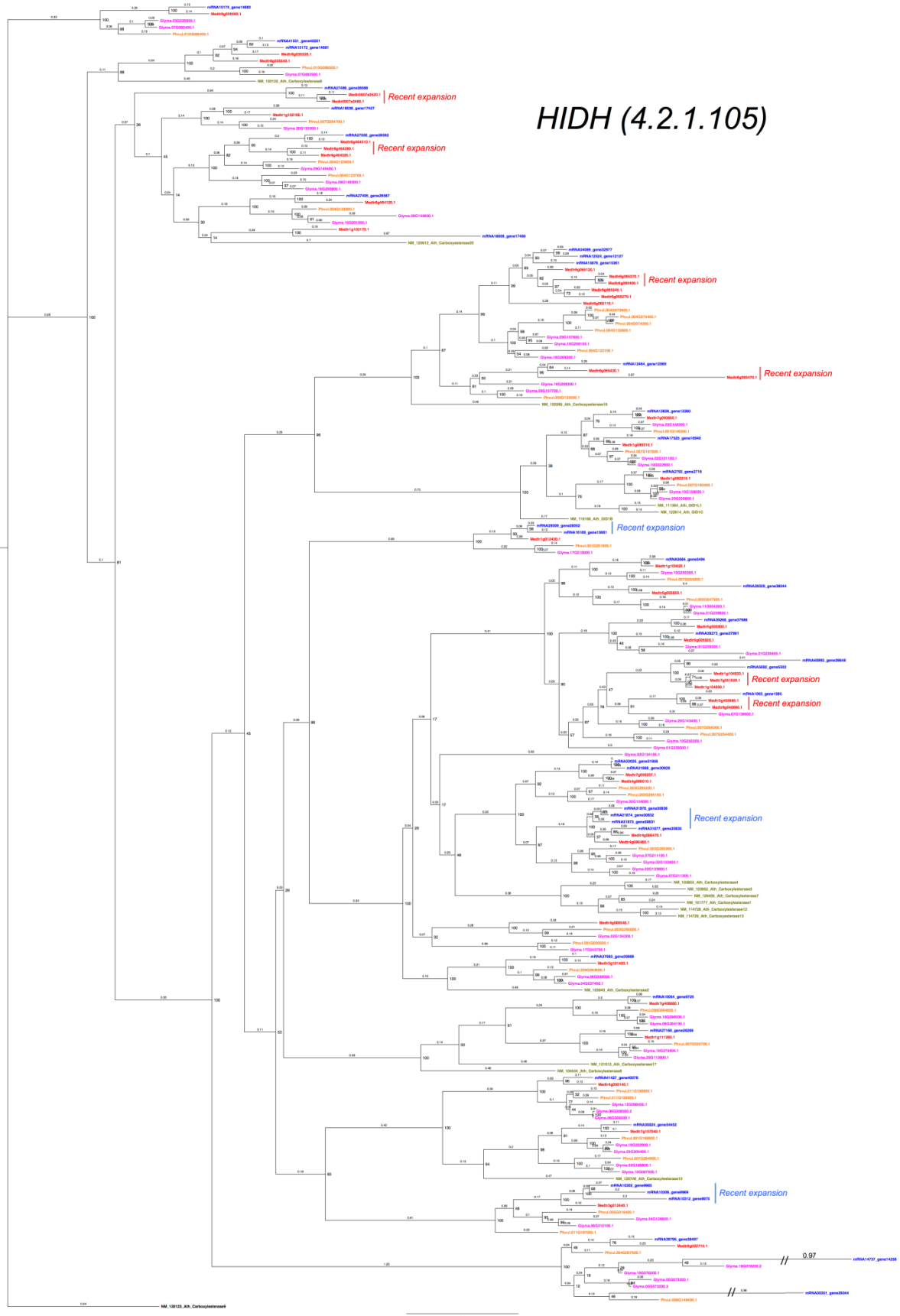


*HIDH (4.2.1.105)*

**Figure 13. Spatial distribution of four clusters of genes** encoding enzymes involved in formononetin biosynthesis in red clover and *M. truncatula.*
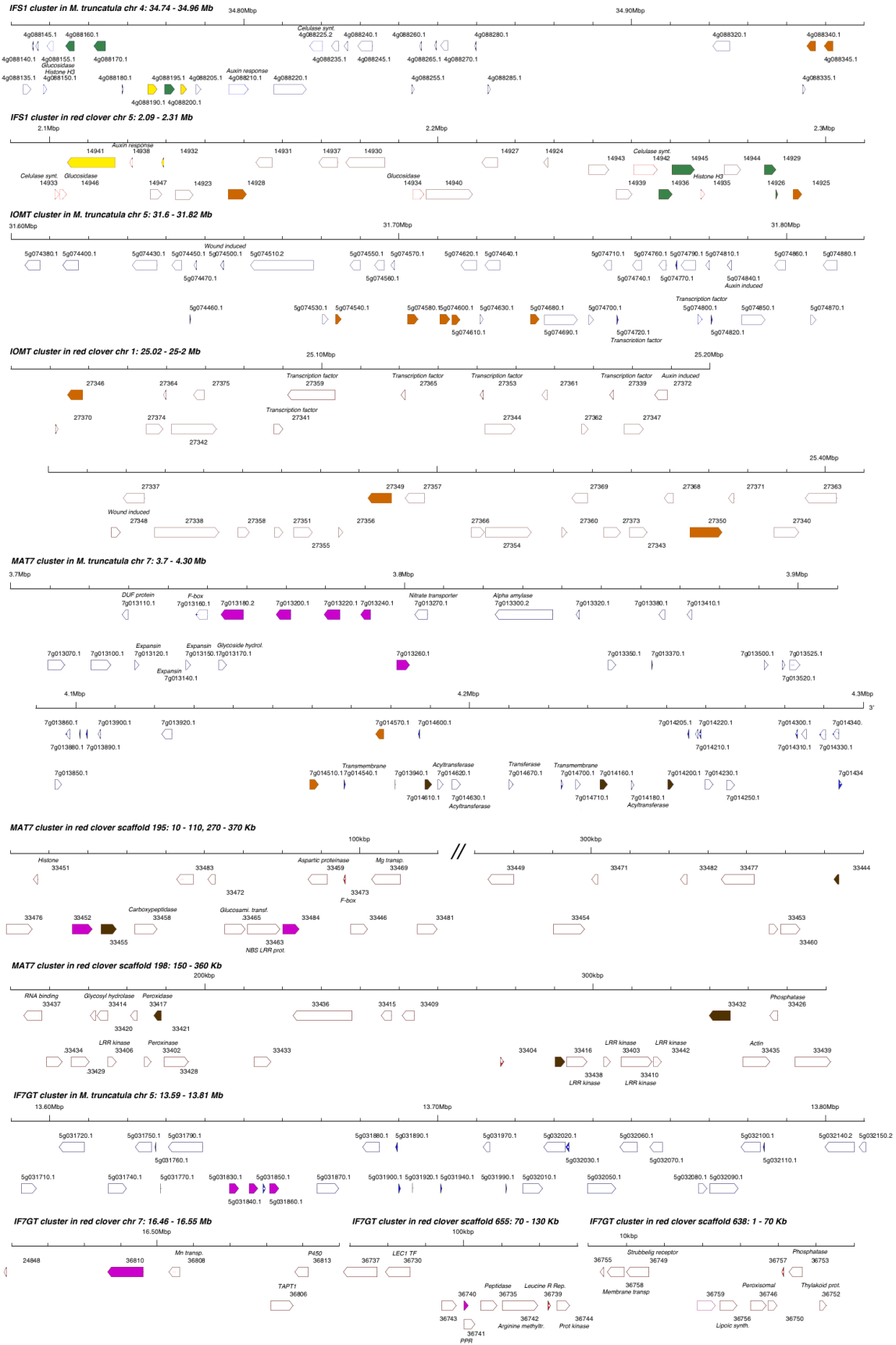
**Figure 14. Phylogenetic tree of isoflavone synthase (IFS1)** (1.14.13.136) in red clover, soybean (Glyma), common bean (Phvu) and *M. truncatula* (Medtr). The AF532999 gene from *Pisum sativum* provides the outgroup.
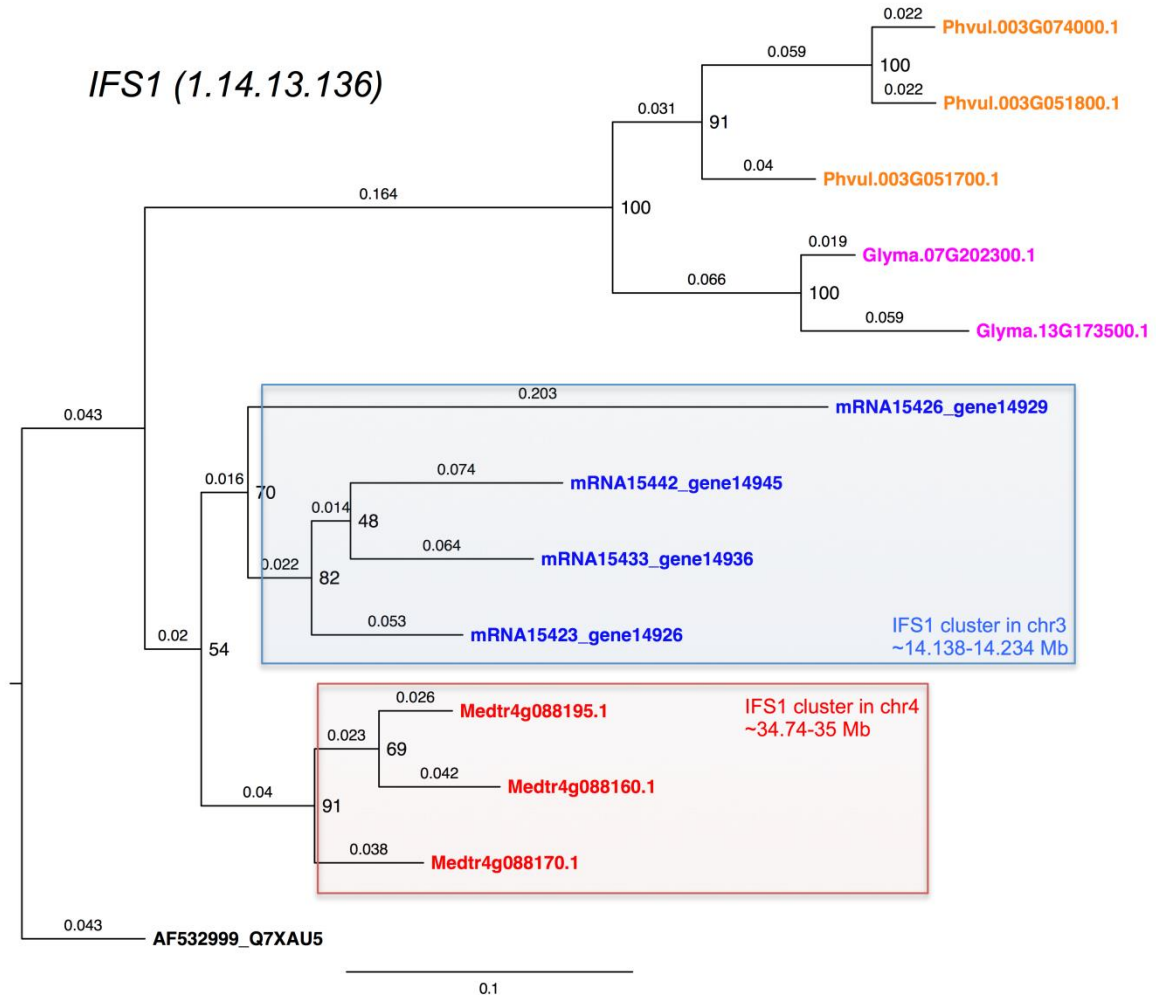
**Figure 15. Phylogenetic tree of isoflavone O-methyltransferase (IOMT)** (2.1.1.150) in red clover, soybean (Glyma), common bean (Phvu) and *M. truncatula* (Medtr). AT4G35150.1 and AT4G35160.1 are outgroups from *A. thaliana.*
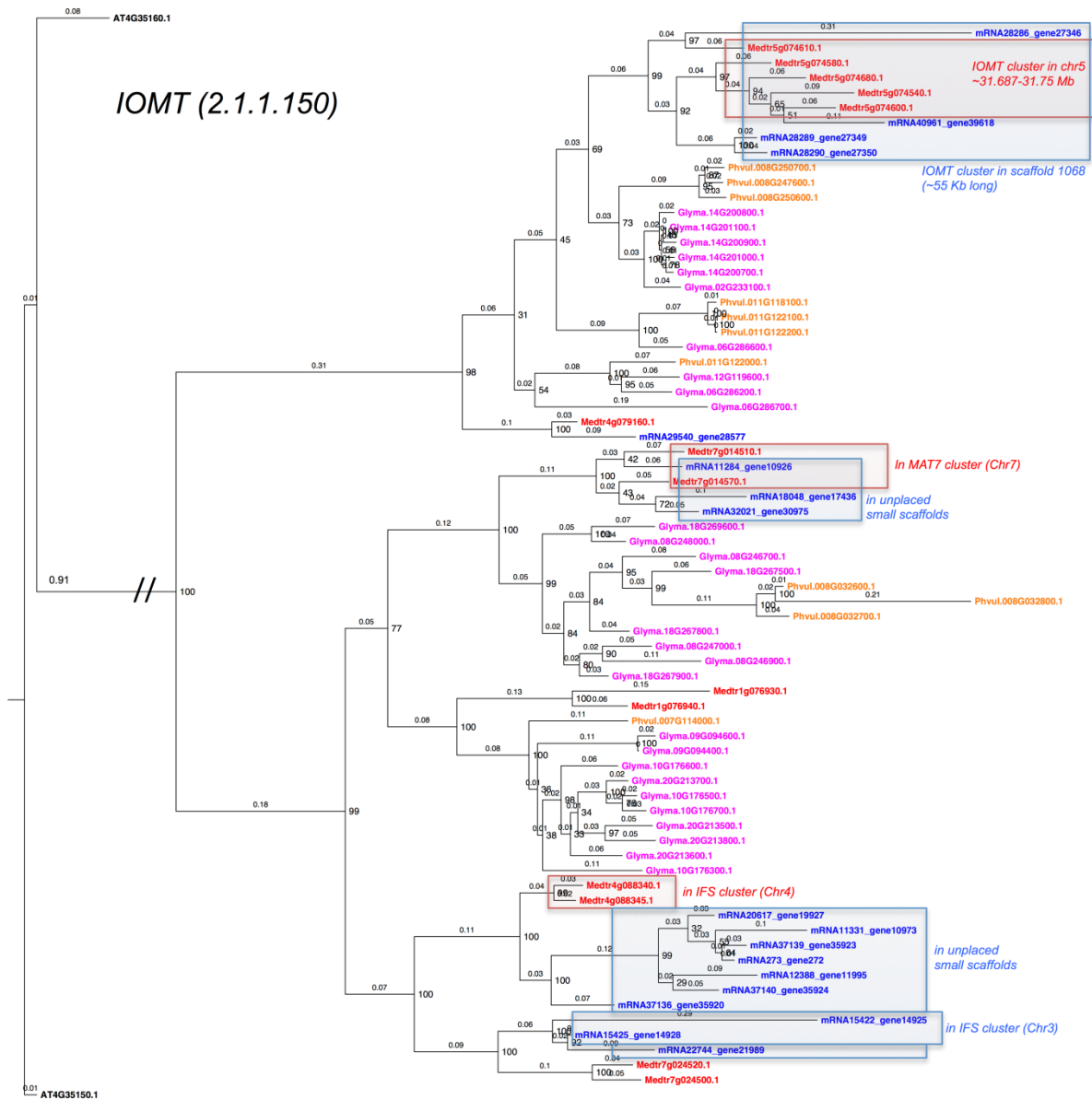
**Figure 16. Phylogenetic tree of 2,7,4'-trihydroxyisoflavanone 4'-O-methyltransferase (HIOMT)** (2.1.1.212) in red clover, soybean (Glyma), common bean (Phvu) and *M. truncatula* (Medtr). AB091686_mRNA_Q84KK4_extraction from *L. japonicus* is the outgroup.
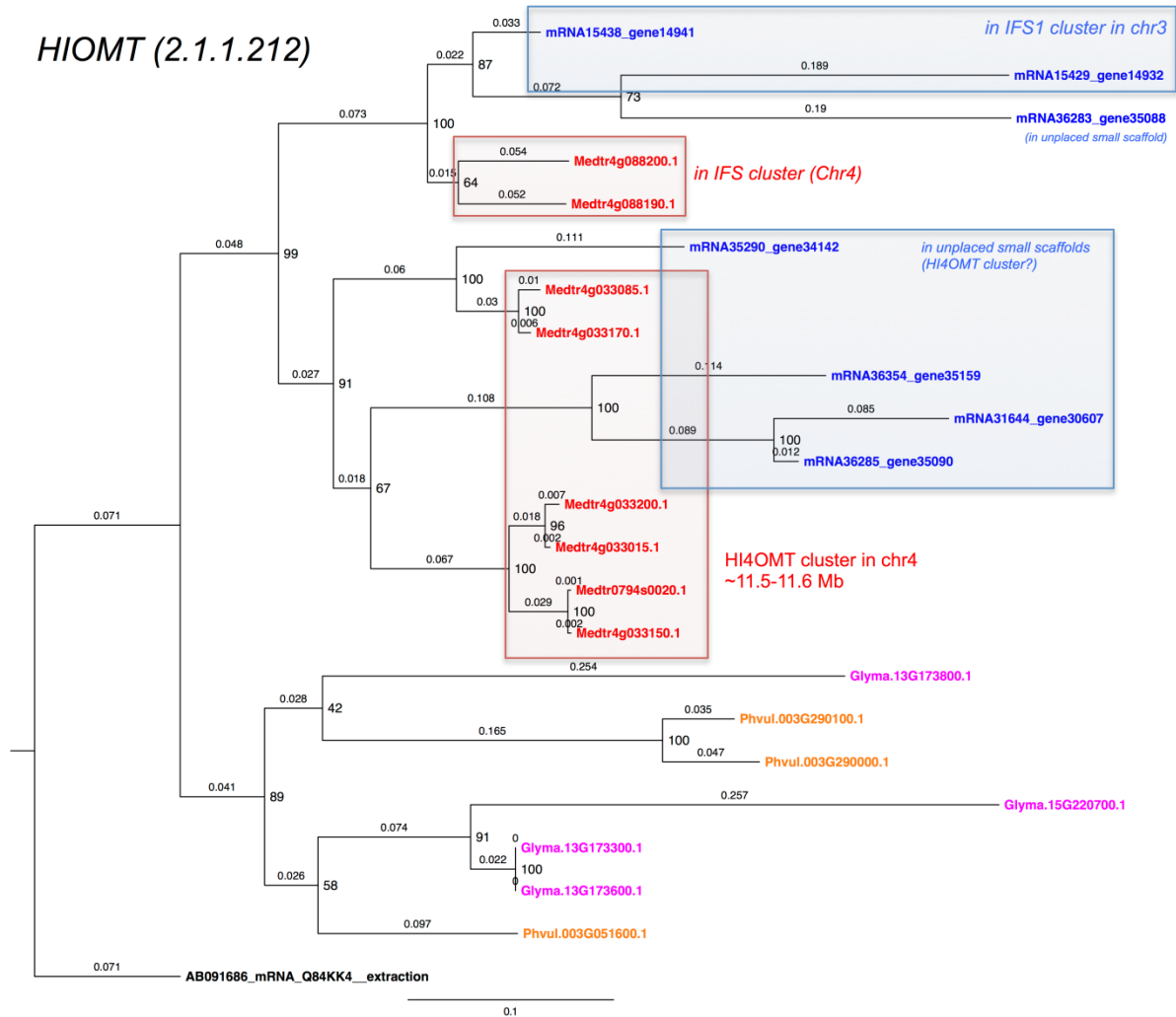
**Figure 17. Phylogenetic tree of isoflavone 7-O-glucosyltransferase (IF7GT)** (2.4.1.170) in in red clover, soybean (Glyma), common bean (Phvu) and *M. truncatula* (Medtr). The outgroup is AT1G1040 from *A. thaliana*.
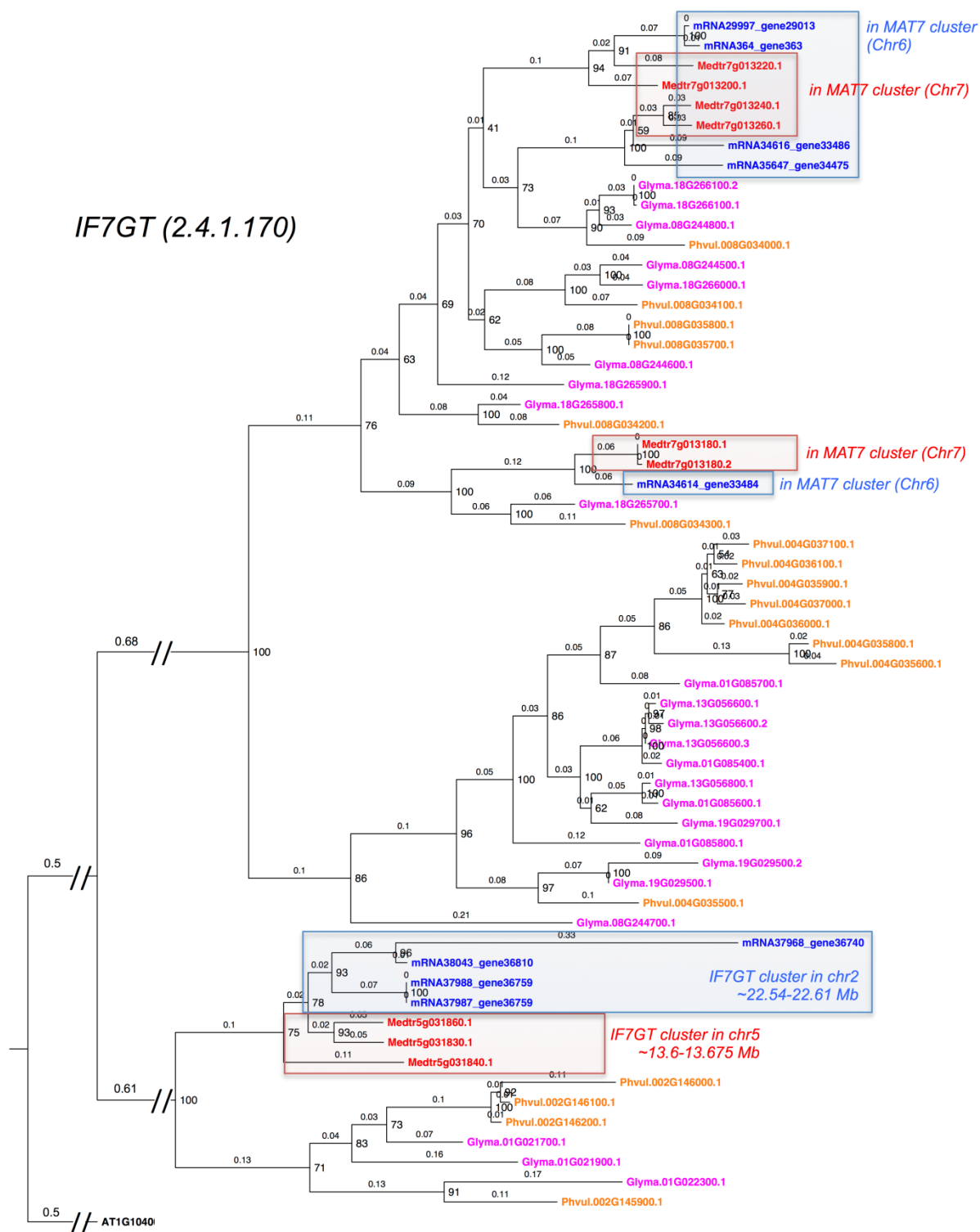
**Figure 18. Phylogenetic tree of Isoflavonoid malonyltransferase (MAT7)** (2.3.1.115) in red clover, soybean (Glyma), common bean (Phvu) and *M. truncatula* (Medtr). The *Ricinus communis* gene is the outgroup.
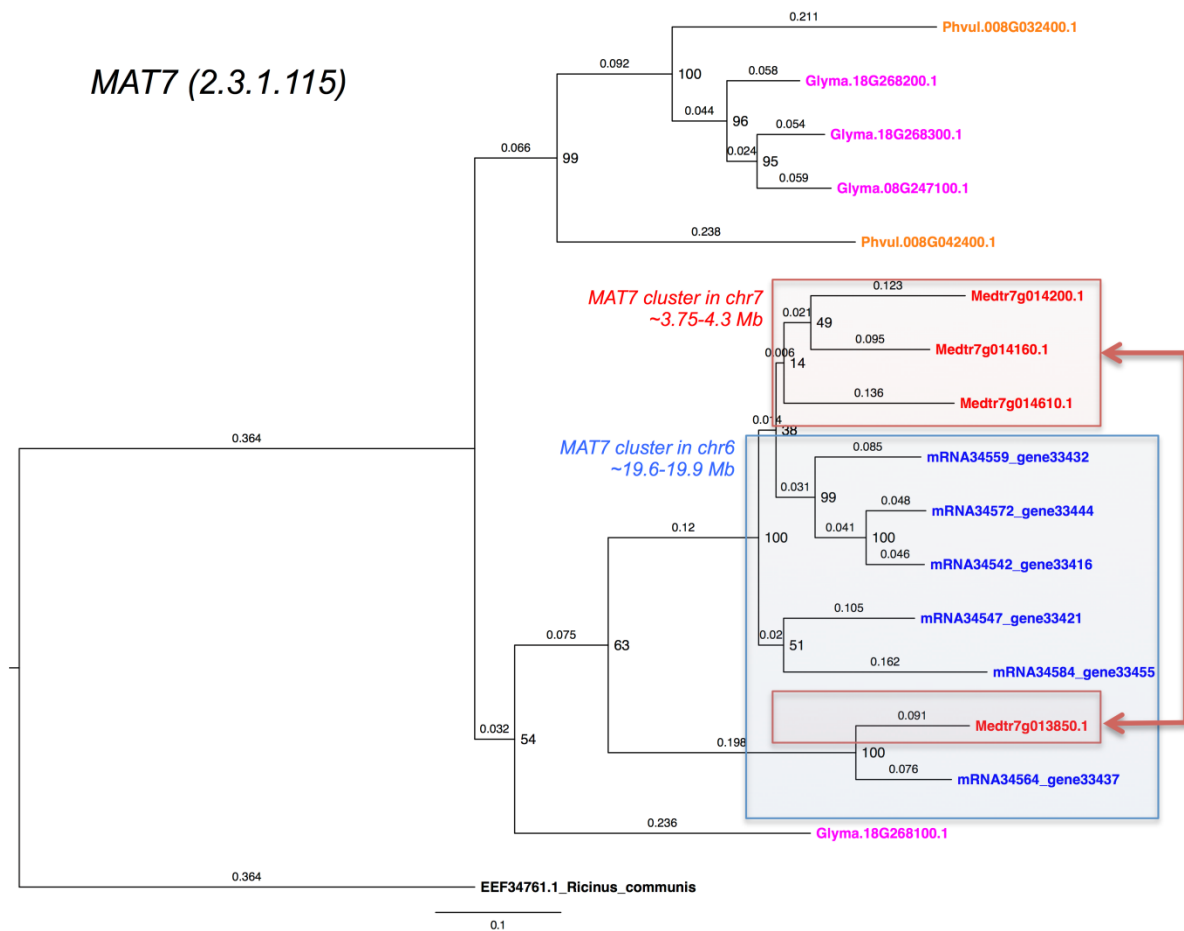
**Figure 19. Phylogenetic tree of PPO protein family in legume species.** The three transcripts isolated from red clover by Sullivan et al (2004)[2] are in purple, the three extra transcripts proposed by Winters et al (2009)[3] are in orange, and the five genes in the red clover genome are in red. PPO genes identified in *M. truncatula* are indicated in blue. PPO2 genes from *Populus balsamifera* provide outgroups.
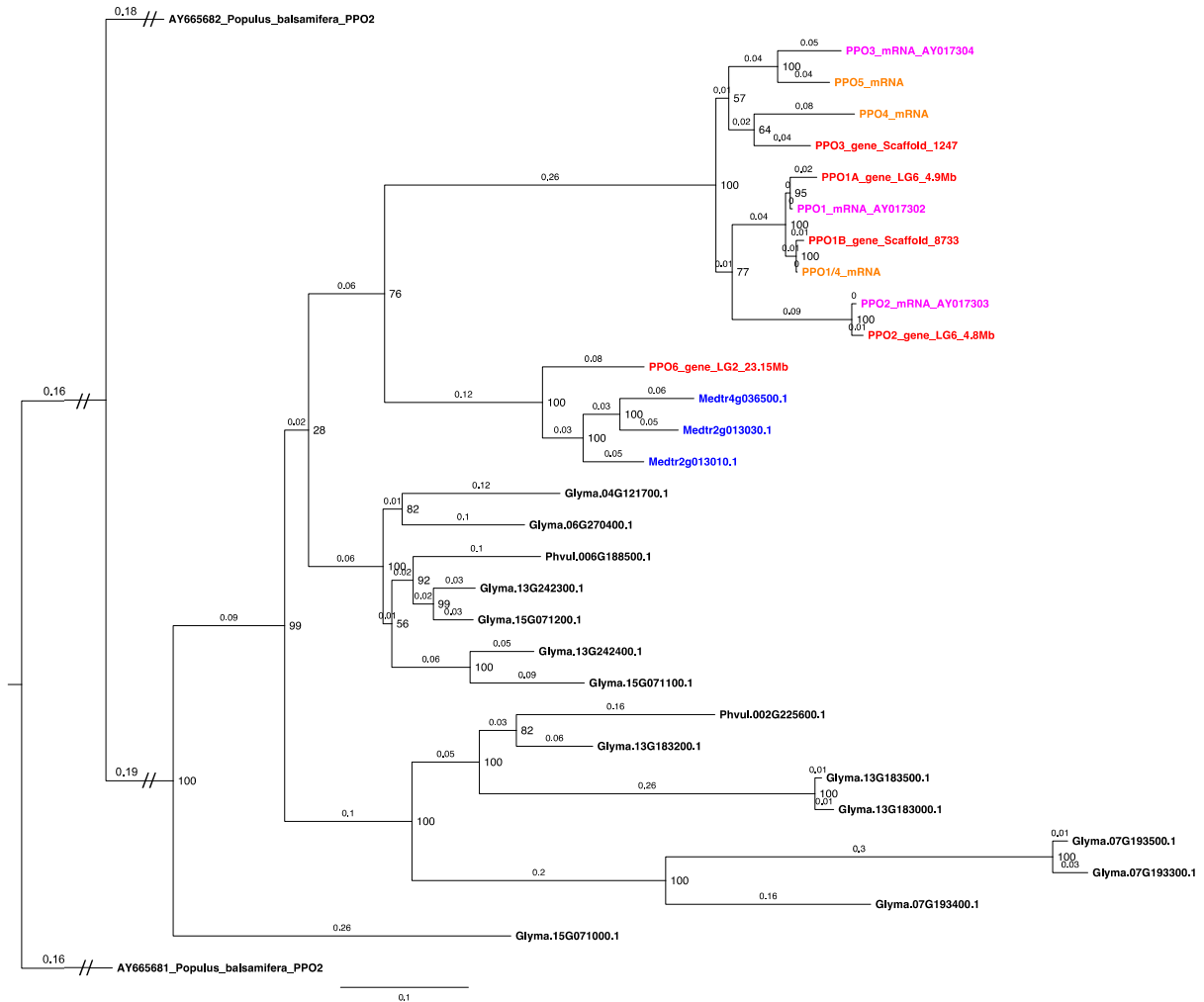
**Figure 20. Microsynteny between *M. truncatula* and red clover in the PPO region.**
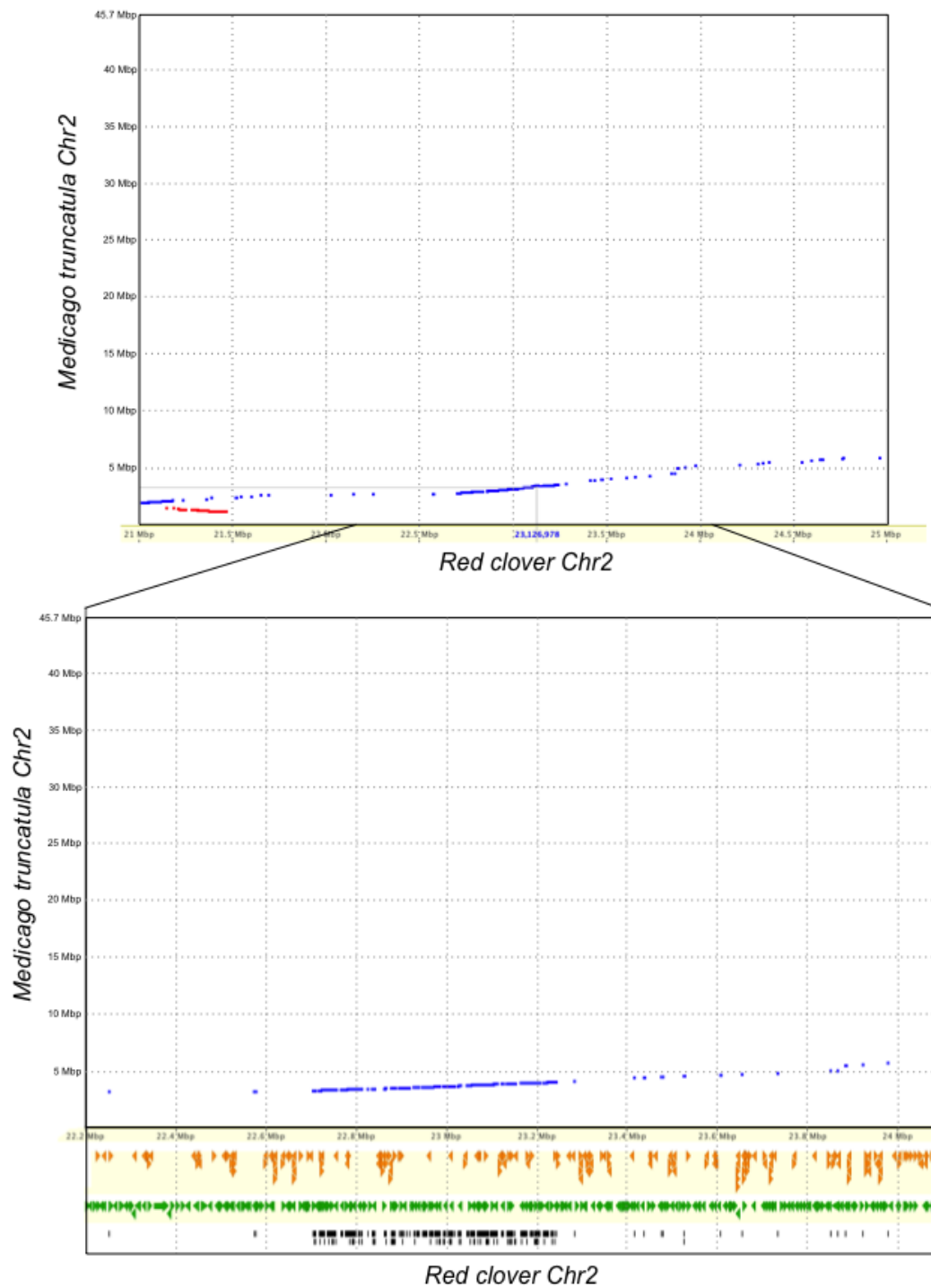
**Figure 21. Linkage disequilibrium in a synthetic population of red clover.** A: Predicted decay of LD in the seven chromosomes of red clover. B: LD heatmap and landscape plots of the red clover chromosomes. The graphs were generated as described[50].

De Vega *et al.* (2015). Red clover (*Trifolium pratense* L.) draft genome provides a platform for trait improvement.

**Figure 22. PCA analysis of red clover synthetic population.** The analysis is based on GBS data for the 86 genotypes of the red clover variety "Lea". The two populations were deduced from analysis using the STRUCTURE programme[4].

**Figure 23.** Outline of the annotation pipeline.

## Supplementary Tables

**Table 1**. Statistics of the WGS, final assembly and pseudo-molecules.

| Version | Sequences (Total) | Seqs >N50 | N80 | N50 | N20 | Max | Total (Mb) | Placed (Mb) |
|---|---|---|---|---|---|---|---|---|
| **WGS (unfiltered)** | 347,062 | 493 | 1,320 | 167,200 | 485,424 | 1,636,054 | 362.7 | 0 |
| **WGS (>500bp)** | 39,904 | 353 | 25,874 | 223,063 | 533,285 | 1,636,054 | 309 | 0 |
| **Final** | 39,051 | 7 | 29,842 | 13.02 Mb | 26.5 Mb | 28.17 Mb | 309 | 164.2 |
| **Pseudo-molecules** | 7 | 4 | 24.71Mb | 25.1 Mb | 27 Mb | 28.17Mb | 164.2 | 164.2 |

**Table 2.** Analysis of the Kmers shared or unique between and within different assembly strategies.

| Assembly | Reference: Platanus | Alternative: ABySS |
|---|---|---|
| Kmers | 349,123,999 | 522,570,659 |
| *... of the previous, are unique/distinct* | *...279,627,096 (80.1%)* | *...319,174,661 (61.1%)* |
| Kmers found in the other assembly | 330,745,871 (94,7%) | 455,166,458 (87.1%) |
| Kmers NOT found in the other assembly | 18,378,128 (5.3%) | 67,404,201 (12.9%) |
| *... of the previous, are unique/distinct* | *...18,113,106 (98.6%)* | *...57,660,671 (85.5%)* |

**Table 3.** Position and functional annotation of the genome (MS Excel file).

- Sheet1: Gene/transcripts annotation using Blast2GO, Uniprot homologous, Interpro, eggNOG, and ORF description
- Sheet2: Transcripts correspondence in eggNOG gene clusters, cluster function, cluster family.
- Sheet3: 1,253 transcripts in clusters expanded in red clover in comparison to *M. truncatula.*

**Table 4.** Classification and proportion of repetitive content in the red clover genome and comparison with other legume genomes.

| | Superfamily | Red clover | | | Medicago truncatula | | Phaseolus vulgaris | | Glycine max | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | Coverage (bp) | Fraction genome (% of 309 Mbp) | Fraction in pseudo-molecules (% of 164 Mbp) | Coverage (bp) | Fraction genome (% of 389 Mbp) | Coverage (bp) | Fraction genome (% of 472 Mbp) | Coverage (bp) | Fraction genome (% of 955 Mbp) |
| Class 1 TEs | Gypsy | 7566938 | 2.45 | 1.65 | 38954999 | 10.01 | 131786780 | 27.92 | 269298878 | 28.19 |
| | Copia | 24391421 | 7.89 | 5.99 | 30381006 | 7.81 | 79334021 | 16.81 | 154704057 | 16.19 |
| | SINEs | 1480120 | 0.48 | 0.47 | 3085624 | 0.79 | 164280 | 0.03 | 1340657 | 0.14 |
| | LINEs | 19637581 | 6.35 | 5.57 | 25763983 | 6.62 | 50090542 | 10.61 | 25568379 | 2.68 |
| | Total Class 1 | 63552730 | 20.57 | 16.29 | 112990132 | 29.05 | 286983344 | 60.8 | 479305937 | 50.19 |
| Class 2 (DNA) TEs | hAT | 9500001 | 3.07 | 2.99 | 5920232 | 1.52 | 6014080 | 1.27 | 17306917 | 1.81 |
| | Harbinger/PIF | 4926009 | 1.59 | 1.70 | 7055779 | 1.81 | 3268021 | 0.69 | 2389726 | 0.25 |
| | MULE | 11052932 | 3.58 | 3.56 | 29629191 | 7.62 | 4795780 | 1.02 | 30491071 | 3.19 |
| | Stowaway | 3803945 | 1.23 | 1.26 | 3139504 | 0.81 | 0 | 0.00 | 494603 | 0.05 |
| | Pogo | 1359180 | 0.44 | 0.45 | 0 | 0.00 | 0 | 0.00 | 0 | 0.00 |
| | CMC_EnSpm | 1220979 | 0.40 | 0.33 | 1942964 | 0.50 | 11950619 | 2.53 | 28217241 | 2.95 |
| | Helitron | 4505094 | 1.46 | 1.35 | 5911950 | 1.52 | 1590726 | 0.34 | 3162793 | 0.33 |
| | Other | 22539391 | 7.29 | 7.02 | 19516885 | 5.02 | 22539087 | 4.78 | 57559831 | 6.03 |
| | Total Class 2 | 58907531 | 19.06 | 18.65 | 73116505 | 18.80 | 50158313 | 10.63 | 139622182 | 14.62 |
| UnclassifTE | Unclassif TE | 5116317 | 1.66 | 1.56 | 2222114 | 0.57 | 2465270 | 0.52 | 3676862 | 0.38 |
| Non TEs | Simple reps | 1516187 | 0.49 | 0.46 | 1056782 | 0.27 | 657868 | 0.14 | 1652138 | 0.17 |
| | Satellites | 110986 | 0.036 | 0.03 | 36149 | 0.01 | 1126 | 0.00 | 764888 | 0.08 |
| | TOTAL | 129226537 | 41.82 | 37.01 | 189498274 | 48.71 | 340265921 | 72.09 | 626039060 | 65.55 |

**Table 5.** Average pairwise linkage disequilibrium ($r^2$) in the seven red clover chromosomes at three distances.

| Chromosome | 0.076 Mb | 0.1 Mb | 0.5 Mb |
|:---:|:---:|:---:|:---:|
| 1 | 0.31 | 0.25 | 0.06 |
| 2 | 0.21 | 0.17 | 0.04 |
| 3 | 0.19 | 0.15 | 0.03 |
| 4 | 0.24 | 0.20 | 0.05 |
| 5 | 0.21 | 0.17 | 0.04 |
| 6 | 0.25 | 0.20 | 0.05 |
| 7 | 0.22 | 0.18 | 0.04 |

**Table 6.** Shotgun short-read libraries used in the WGS.

| Library | Reads | EBI-ENA Accession |
|---|---|---|
| **PE (100 bp)** | 101.4e6 | ERX946106 |
| **PE (150 bp)** | 306.3e6 | ERX946107 |
| **SE (150 bp)** | 254.6e6 | ERX946108, ERX946109 |
| **Overlap PE (150 bp)** | 25.6e6 | ERX946110 |
| **3 Kb MP (100 bp)** | 224.8e6 | ERX946085, ERX946084 |
| **5 Kb MP (65 bp)** | 4.2e6 | ERX946083 |
| **7 Kb MP (100 bp)** | 20.4e6 | ERX946086 |
| **7 Kb MP (150 bp)** | 88.0e6 | ERX946087 |

All shotgun reads and the assembly are deposited in the European Nucleotide Archive (accession PRJEB9186). The genome assembly and annotation can also be downloaded as files (http://dx.doi.org/10.5281/zenodo.17232) or browsed online (http://tgac-browser.tgac.ac.uk/trifolium_pratense).

**Supplementary References**

1      Ištvánek, J., Jaroš, M., Křenek, A. & Řepková, J. Genome assembly and annotation for red clover (Trifolium pratense; Fabaceae). *American Journal of Botany* **101**, 327-337 (2014).

2      Sullivan, M. L., Hatfield, R. D., Thoma, S. L. & Samac, D. A. Cloning and Characterization of Red Clover Polyphenol Oxidase cDNAs and Expression of Active Protein in Escherichia coli and Transgenic Alfalfa. *Plant Physiology* **136**, 3234-3244 (2004).

3      Winters, A. *et al.* Identification of an extensive gene cluster among a family of PPOs in Trifolium pratense L. (red clover) using a large insert BAC library. *BMC Plant Biology* **9**, 94 (2009).

4      Pritchard, J. K., Stephens, M. & Donnelly, P. Inference of Population Structure Using Multilocus Genotype Data. *Genetics* **155**, 945-959 (2000).