

1 **An integrated multi-level comparison highlights common**  
2 **aspects and specific features between distantly-related**  
3 **species: Tomato and Grapevine**

4 Luca Ambrosino<sup>1</sup>, Hamed Bostan<sup>1</sup>, Valentino Ruggieri<sup>1</sup>, Maria Luisa Chiusano<sup>1</sup>

5 <sup>1</sup> Department of Agriculture, University of Naples Federico II, Portici (NA), Italy

6

7 Corresponding Author:

8 Maria Luisa Chiusano<sup>1</sup>

9 Via università 100, 80155 Portici (NA), Italy

10 Email address: [chiusano@unina.it](mailto:chiusano@unina.it)

11 **Abstract**

12

13 **Motivation.** Even after years from the first completion of genomes by sequencing,  
14 comparative genomics still remains a challenge, also enhanced by the availability of numerous  
15 draft genomes with still poor annotation quality. The detection of ortholog genes between  
16 different species is a key approach for comparative genomics. For example, ortholog gene  
17 detection may support investigations on mechanisms that shaped the organization of the  
18 genomes, highlighting on gain or loss of function and on gene annotation. On the other hand,  
19 the detection of paralog genes is fundamental for understanding the evolutionary  
20 mechanisms that drove gene function innovation and support gene families analyses. Here we  
21 report on the gene comparison between two distantly related plants, *Solanum lycopersicum*  
22 (Tomato) (The Tomato Genome Consortium 2012) and *Vitis vinifera* (Grapevine) (Jaillon et al.  
23 2007), considered as economically important species from asterids and rosids clades,  
24 respectively. The strategy was accompanied by integration of multilevel analyses, from  
25 domain investigations to expression profiling, to get to the most reliable results and to offer  
26 powerful resources, in order to understand different useful aspects of plant evolution and  
27 physiology and to dissect traits and molecular aspects that could provide novel tools for  
28 agriculture applications and biotechnologies.

29 **Methods.** In order to predict best putative orthologs and paralogs between Tomato and  
30 Grapevine, and to overcome possible annotation issues, all-against-all sequence similarity  
31 searches between genes, mRNAs and proteins collections of both species were performed. A  
32 Bidirectional Best Hit approach was implemented to detect the best orthologs between the  
33 two species. Moreover we developed a dedicated algorithm in Python programming language  
34 able to define more extended alignments between mRNA sequences. NetworkX package  
35 (Hagberg et al. 2008) was used to define networks of paralogs and orthologs. Proteins domain  
36 prediction was carried out on the entire Tomato and Grapevine protein collection by using  
37 InterProScan program (Jones et al. 2014). The enzyme classification was obtained by sequence  
38 similarity searches between Tomato and Grapevine mRNA collections and the entire UniProt  
39 reviewed protein collection (UniProt consortium 2015). The metabolic pathways associated  
40 to the detected enzymes were identified exploiting the KEGG Database (Kanehisa and Goto

41 2000). Expression level of three developmental stages of Tomato (2 cm fruit, breaker and  
42 mature red) and the corresponding stages of Grapevine (post-setting, veraison, mature berry)  
43 was defined on the basis of the iTAG loci (Shearer et al. 2014) and v1 vitis loci, respectively.  
44 The expression was normalized by Reads Per Kilobases per Million (RPKM) for each  
45 tissue/stage. The identification of similar expression profiles was performed by the K-means  
46 clustering method (Soukas et al. 2000), using the Pearson correlation coefficient as distance  
47 metric. For each cluster a subsequent clusterization by the Hierarchical Clustering (HCL) (Eisen  
48 et al. 1998) using the Euclidean distance grouped genes also on the basis of expression levels.  
49 Both the clustering methods used are those from implemented in the MultiExperiment Viewer  
50 (MeV) software.

51 **Results.** Although Tomato and Grapevine are phylogenetically distant species, they are both  
52 model species for understanding fleshy fruit formation. Comparative analyses, though the  
53 available annotations are still preliminary, are essential to understand fruit development. We  
54 predicted the presence of a strong core of orthologs genes, exploiting an appropriate  
55 approach and overcoming the annotation limits. Networks of ortholog/paralog genes were  
56 built between the compared species, offering resources to support studies about the  
57 organization and the evolution of gene families in different organisms. By this approach, we  
58 detected gene families of one species that underwent an expansion/reduction in the number  
59 of their elements when compared to the other species. Species-specific genes of Tomato and  
60 Grapevine were also detected. The protein domains common to both species, as the ones  
61 exclusively detected in Tomato and Grapevine, and the common and the distinctive enzymatic  
62 classes associated to related metabolic pathways, were also predicted for the two compared  
63 species supporting structure and functional annotations. Furthermore, the association of  
64 RNA-seq data offered an additional information level for comparing gene functionalities from  
65 the two species. Thanks to this core collection, we report on similarities and peculiarities  
66 between the two genomes.