

# A proposed naming convention for genes in Rosaceae species



Rosaceae Gene Name Standardization Subcommittee of RosEXEC and RosIGI: Sook Jung<sup>1</sup>, Carole Bassett<sup>2</sup>, Doug Bielenberg<sup>3</sup>, Chun-Huai Cheng<sup>1</sup>, Chris Dardick<sup>2</sup>, Dorrie Main<sup>1</sup>, Lee Meisel<sup>4</sup>, Janet Slovin<sup>5</sup>, Michela Troggio<sup>6</sup> and Robert J. Schaffer<sup>7</sup>

<sup>1</sup>Washington State University,Pullman, WA, USA , <sup>2</sup>USDA ARS Appalachian Fruit Research Station, Kearneysville, WV, USA, ARS, <sup>3</sup>Clemson University, Clemson SC, USA, <sup>4</sup>University of Chile, Santiago, Chile, <sup>5</sup>USDA, Beltsville, MD, USA, <sup>6</sup>Istituto Agrario San Michele all'Adige, San Michele all'Adige, Italy, <sup>7</sup>Plant & Food Research Ltd., Auckland, New Zealand

## Introduction

The importance of a consistent hierarchical naming convention in science was pioneered by Linnaus, who was the father of bringing order into biological species. Through linking similarity in species he generated a naming convention that is still used today. Naming gives clarity as well as and a small amount of biological information to researchers, with genes in model species often being named after their mutant phenotype. The convention has been to generate a 2-3 letter acronym for each gene and when comparing across different species usually a 2 letter species abbreviation is used as a prefix. For example, the prefix for *Arabidopsis thaliana* genes would be At. The species prefix is not required when presenting results from a single species but it is necessary when genes from multiple species are being compared.

Naming rights are typically given to the “discoverer”, and like all things in history, sometimes things are discovered more than once. In model species such as Arabidopsis and tomato, multiple names have been used in literature for a single gene. For example ETHYLENE INSENSITIVE 5 has also been published as AIN1 and XRN4. This has only been resolved by the genome naming convention of giving genes a chromosomal name and gene number (in this case AT1G54490) for the predicted genes from whole genome sequencing. The Rosaceae community is significantly smaller than the Arabidopsis community, and therefore we have the chance to set out naming conventions that do not fall into such traps.

## Proposed Naming Convention

### 1. Species abbreviation

The comparison of genes across Rosaceae species gives valuable insights into the way that different species have evolved. Whole genome sequence is available for strawberry (Shulaev et al. 2010), apple (Velasco et al. 2010), peach (International Peach Genome Initiative 2013), and pear (Wu et al. 2013), with more to come. This allows cross species comparisons at the genome level and within a gene family. Considering the 3000 species within Rosaceae, however, it is clear that two-letter species abbreviations will not give clear distinction between species. It is likely that the majority of researchers will be focusing their gene based research on commercially cultivated varieties, however, even within these there are issues as both *Prunus persica* and *Pyrus pyrofolia* would become Pp. We therefore propose standard naming of major Rosaceae species using the following abbreviations (Table 1). Using a 3-letter prefix resulted in two conflicts: both *Prunus cerasus* and *Prunus cerasifera* would be Pce following our convention, so we recommend Pci for *Prunus cerasifera*. Likewise, both *Prunus mume* and *Prunus munsoniana* would be Pmu. We recommend Pmn for *Prunus munsoniana* (highlighted in the table). For taxonomy studies across non commercial species, abbreviations will be needed to be longer to distinguish the species. For these papers we suggest that researchers take a UNIPROT approach (<http://www.uniprot.org/docs/speclist> ) using 5 letter abbreviations: 3 for the Genus name and 2 for the species with *Malus x domestica* becoming Maldo and *Prunus persica* becoming Prupe. This convention is still not sufficient for the complete list, and in these papers a clear nomenclature needs to be stated in the paper by the researcher. **However, we recommend that authors do not include a species prefix in the gene symbol when they submit the gene data to NCBI, GDR or any other databases, to minimize the creation of duplicated names due to the differences in the species prefix.** In GDR, the gene symbol with the recommended three-letter prefix will be stored along with the gene symbol without the prefix for genes from the major varieties.

**Table 1. Proposed Species Abbreviation for major Rosaceae species.**

Genus	Species	Prefix
Aronia	melanocarpa	Ame
Chaenomeles	japonica	Cja
Cydonia	oblonga	Cob
Eriobotrya	japonica	Eja
Fragaria	x ananassa	Fan
Fragaria	chiloensis	Fch
Fragaria	vesca	Fve
Fragaria	virginiana	Fvi
Malus	x domestica	Mdo
Mespilus	germanica	Mge
Malus	pumila	Mpu
Malus	sylvestris	Msy
Prunus	americana	Pam
Prunus	angustifolia	Pan
Prunus	armeniaca	Par
Prunus	avium	Pav
Pyrus	calleryana	Pca
Prunus	cerasifera	Pci
Prunus	cerasus	Pce
Pyrus	communis	Pco
Prunus	domestica	Pdo
Prunus	dulcis	Pdu
Prunus	hortulana	Pho
Prunus	mume	Pmu
Prunus	munsoniana	Pmn
Prunus	nigra	Pni
Prunus	persica	Ppe
Pyrus	pyrifolia	Ppy
Prunus	salicina	Psa
Prunus	serotina	Pse
Prunus	simonii	Psi
Prunus	spinosa	Psp
Pyrus	ussuriensis	Pus
Rubus	idaeus	Rid
Rubus	occidentalis	Roc

## 2. Gene symbol

For the actual name we encourage using a ‘root’ symbol for members of a gene family together with a hierarchical numbering system. It is ideal to design the symbol so that it can be associated with biological function. The Arabidopsis community has set out naming guidelines that can be found at (<http://www.arabidopsis.org/portals/nomenclature/guidelines.jsp>). When giving a gene a name we encourage people to check that the gene they want to name (and publish) does not already have a name assigned by searching the current literature and gene databases (GenBank and GDR). Co-naming with Arabidopsis may be difficult due to clade expansions between the two species. In spite of this we recommend using similar naming and numbering when possible.

**No specific naming convention is proposed for homeologues, but we recommend that they be named sequentially within the gene families. We propose the following convention for naming splice variants and alleles.**

- When the gene name does not contain a number eg. *PG*
  - PG1*, *PG2*... (homeologues and homologues)
- When the gene has a number in the name eg *DHN3*
  - DHN3.1*, *DHN3.2* (homeologues and homologues)
  - DHN3.1\_1*, *DHN3.1\_2* (splice variant)
  - DHN3.1a*, *DHN3.1b* (allele)

## Facilitation of Gene Naming Standardization

### Submission of gene data to GDR prior to publication

Researchers are encouraged to submit their gene data to GDR in addition to NCBI. While NCBI does not accept named genes that do not come from single molecule sequencing, the GDR database will. We also recommend researchers provide the corresponding gene models from whole genome sequencing, the gene family name, and the level of confidence based on the criteria we developed (see below). **When new gene models, gene sequences and/or splice variants are identified that are different from the whole genome data, we will recommend that users submit FASTA file in addition to the information in the table below.**

Species	Species prefix	Gene Symbol	Gene Name	Gene Family	Synonyms	Gene Model	Genbank ID	Description	Submitter	Level of confidence	Comments

**Defining the 'level of confidence' for various sequencing methods and subsequent analysis.**

- A. Single molecule bi-directional cDNA sequence
- B. Compilation DNA sequencing
- C. Compilation of RNA-seq
- D. Computational evidence

**A GDR page for browsing expert-curated Rosaceae genes/gene families**

GDR   Genome Database for Rosaceae											
List of Gene Names submitted by researchers											
Rosaceae Gene Name Standardization Subcommittee of RosEXEC and RosIGI recommends researchers to submit gene names to GDR prior to publication (gene name submission form) to minimize the duplication of gene names. GDR curators will check for any existing gene names and notify researchers. Three letter prefix for species and three letter gene symbol are recommended. Refer to the documents for the details of recommended gene naming convention. Below is the gene names that follows the convention submitted by researchers. Current gene names from NCBI can be found in gene search page.											
Species	Species prefix	Gene Symbol	Gene Name	Gene Family	Synonyms	Gene Model	Genbank ID (version)	Description	Submitter	Level of Confidence	Submitted
Prunus persica	Ppe	DH12	dehydri	dehydri	Xero1_LEA	ppa011537m	AY45376.1	similar to Xero1 and Raa18 of Arabidopsis thaliana	CL Bassett	A	Note: NCBI Uniprot has Ppe_19087 designation for this gene
Malus DomCica	Mdo	DH13.1	dehydri	dehydri	LT29: ERD10	MDP000077049:spnd	X59814.1 NAL_180616		CL Bassett	A	MDCH13.1 and 3.2 are identical at amino acid level
Malus DomCica	Mdo	DH13.3	dehydri	dehydri	LT29: ERD10	MDP000052900:spnd	X59814.1 NAL_180616		CL Bassett	D	
Malus DomCica	Mdo	PG1	dehydri	GLYCOSYLHYDROLASE 28 FAMILY		MDP000032673:spnd	P48978	Ripening associated polygalacturonase	Robert Schaffer	A	Abrinson et al. 1994, Abrinson et al 2012
Prunus persica	Ppe	CKX1	Cytokinin oxidase 1	Cytokinin oxidase	PpCKX1	ppa024442		Cytokinin oxidase/dehydrogenase, cytokinin catabolism, oxidoreductas FAD-binding domain	Lee Meisel	C	Vizzoso et al 2009; Immananen et al 2013
Prunus persica	Ppe	CKX2	Cytokinin oxidase 2	Cytokinin oxidase	PpCKX2	ppa021859		Cytokinin oxidase/dehydrogenase, cytokinin catabolism, oxidoreductas FAD-binding domain	Lee Meisel	C	Vizzoso et al 2009; Immananen et al 2013
Prunus persica	Ppe	CKX3	Cytokinin oxidase3	Cytokinin oxidase	PpCKX3	ppa021417		Cytokinin oxidase/dehydrogenase, cytokinin catabolism, oxidoreductas FAD-binding domain	Lee Meisel	C	Vizzoso et al 2009; Immananen et al 2013
Prunus persica	Ppe	CKX5	Cytokinin oxidase 5	Cytokinin oxidase	PpCKX5	ppa003895		Cytokinin oxidase/dehydrogenase, cytokinin catabolism, oxidoreductas FAD-binding domain	Lee Meisel	C	Vizzoso et al 2009; Immananen et al 2013

## Conclusions

Here we propose a standardised naming convention for Rosaceae species. We respectfully urge the Rosaceae community to follow these suggestions and by doing so, we will benefit from simplified literature reading, and less confusion when looking at genes. Please send feedback on these suggestions to any of the co-authors.