

## The genomes of several plant species contain endogenous geminiviral sequences

Denis Filloux<sup>1</sup>, Sasha Murrell<sup>2,3</sup>, Maneerat Koohapitagtam<sup>1,4</sup>, Michael Golden<sup>2</sup>, Charlotte Julian<sup>1</sup>, Serge Galzi<sup>1</sup>, Marilyne Uzest<sup>1</sup>, Marguerite Rodier-Goud<sup>5</sup>, Angélique D'Hont<sup>5</sup>, Marie-Stéphanie Vernerey<sup>1</sup>, Paul Wilkin<sup>6</sup>, Michel Peterschmitt<sup>1</sup>, Stephan Winter<sup>7</sup>, Ben Murrell<sup>2,8</sup>, Darren P. Martin<sup>2</sup>, and Philippe Roumagnac<sup>1</sup>

<sup>1</sup> CIRAD-INRA-SupAgro, UMR BGPI, Campus International de Montpellier-Baillarguet, 34398 Montpellier Cedex-5, France

<sup>2</sup> Computational Biology Group, Institute of Infectious Diseases and Molecular Medicine, University of Cape Town, Cape Town 4579, South Africa

<sup>3</sup> Department of Integrative Structural and Computational Biology, The Scripps Research Institute, La Jolla, CA 92037, USA

<sup>4</sup> Department of Pest Management, Faculty of Natural Resources, Prince of Songkla University, Hat Yai campus, Thailand 90120

<sup>5</sup> CIRAD, UMR AGAP, TA A-108/03, Avenue Agropolis, F-34398 Montpellier Cedex 5, France

<sup>6</sup> Royal Botanic Gardens, Kew, Richmond, Surrey, TW9 3AB, UK

<sup>7</sup> DSMZ Plant Virus Department, Messeweg 11/12, 38102, Braunschweig, Germany

<sup>8</sup> Department of Medicine, University of California, San Diego, La Jolla, CA

Endogenous viral sequences are essentially ‘fossil records’ that can sometimes reveal the genomic features of long extinct virus species. Although numerous known instances exist of single-stranded DNA (ssDNA) genomes becoming stably integrated within the genomes of bacteria and animals, there remain very few examples of such integration events in plants. The best studied of these events are those which yielded the geminivirus-related DNA elements (GRD) and the geminivirus-like elements (EGV) found respectively within the nuclear genomes of several *Nicotiana* species and various *Dioscorea* spp. of the *Enantiophyllum* clade.

Those two new classes of endogenous plant virus sequence are apparently derived from ancient geminiviruses in the genus *Begomovirus*. GRD and EGV sequences likely became integrated millions years ago. Interestingly, we found evidence of natural selection actively favouring the maintenance of EGV-expressed replication-associated protein (Rep) amino acid sequences, which clearly indicates that functional EGV Rep proteins were probably expressed for prolonged periods following endogenization.

We recently found using *in silico* searches that other ssDNA virus-like sequences are included within complete or draft genomes of various plant species, including apple tree (*Malus domestica*), black cottonwood (*Populus trichocarpa*), various *Coffea* spp, eggplant (*Solanum melongena*), lettuce (*Lactuca sativa*), and Tepary bean (*Phaseolus acutifolius*), which suggests that endogenous geminiviruses may be more common in plant genomes than has previously been appreciated.