

Supplementary Figure 1. (A) Annotated 2.6 Kb endogenous geminivirus nucleotidic sequence of *D. alata*. (B). Annotated 2.1 Kb endogenous geminivirus nucleotidic sequence of *D. alata*.

Supplementary Figure 2. Alignment of the putative Rep sequences of EGV1 and EGV2 with the Rep of six representative geminiviruses (Tomato golden mosaic virus (TGMV), Maize streak virus (MSV), Eragrostis curvula streak virus (ECSV), Bean golden mosaic virus (BGMV), Beet curly top virus (BCTV) and Tomato pseudo-curly top virus (TPCTV)). RCR motifs (I, II, III and GRS) and SF3 helicase motifs (Walker A, B and C), which are conserved in geminiviruses, are highlighted in red.

Supplementary Figure 3. Rolling circle amplification. Lanes: 1) Ladder; 2) pUC19; 3) African cassava mosaic virus (ACMV); 4) *Dioscorea badnavirus*; 5) Water; 6) *Dioscorea alata* (acc. #313, seedlings); 7) *Dioscorea transversa* (acc. #336).

Supplementary Figure 4. Maximum likelihood tree of eighteen 402 bp partial rep sequences obtained using the “ren-ren” and “rep-rep” primer pairs.

Supplementary Figure 5. (A) Western blot on the total protein extracts of several *Dioscorea* species, using an antibody directed to a Rep peptide of EGV1 and (B) nitrocellulose membrane stained with Ponceau S dye for protein detection during western blotting. Lanes: 1) Proteins extracted from a turnip plant infected by Cauliflower mosaic virus (CaMV); 2) a Tomato plant infected by Tomato yellow leaf curl virus (TYLCV); 3) Ladder; 4) *Dioscorea trifida* (accession#64); 5) *Dioscorea sansibarensis* (accession#269); 6) *Dioscorea bulbifera* (accession#272); 7) *Dioscorea dumetorum* (accession#47); 8) *Dioscorea togoensis* (accession#114, seedling); 9) *Dioscorea praeheasilis* (accession#255); 10) *Dioscorea rotundata* (accession#118, seedling); 11) *Dioscorea nummularia* (accession#335); 12) *Dioscorea alata* (accession#297, seedling); 13) *Dioscorea alata* (accession#313 seedling); 14) *Dioscorea alata* (accession#402 seedling). No cross-reactivity was detected in the turnip/CaMV, tomato/TYLCV, *D. trifida* and *D. sansibarensis* samples.

Supplementary Figure 6. Maximum likelihood trees of replication enhancer protein (Ren) amino acid sequences (JTT substitution model) of EGV1 and representative geminiviruses, from the genera begomovirus, topocovirus and curtovirus. Numbers associated with branches indicate degrees of bootstrap support (100 replicates) for those branches.

Supplementary Figure 7. Duplication and diversification of the yam EGVs following their integration. Two scenarios resulting in contrasting patterns of endogenous sequence diversity are expected and presented. (A) The endogenous sequences were only duplicated (indicated by red arrows) either at or very close to the time when the integration event occurred (indicated by green arrow), or (B) The endogenous sequences are duplicated for prolonged periods post-integration such that duplication and speciation events (indicated by blue arrows) are interspersed. In (A) endogenous sequences sampled from one species will usually be more closely related to endogenous sequences from different species than they are to other genetically distinct endogenous sequences sampled from the

same species. In (B) genetically distinct endogenous sequences sampled from the same species should frequently be more closely related to one another than they are to endogenous sequences sampled from different species.

A- EGV1 annotated sequence from *D. alata* acc. 313

Virion strand origin of replication Inverted repeat sequence

TAATATTACACGGATGGCCGCCCCGACTGTACTGCCACGTGGCGCATTTTGATAGGTCCACGCTAAAGATGCTGTACCC 80
GATTGTCTCCACGGTAAAGATGCTGTACCCGATTGTCTCCACCGTACATCCACGCGCCATCATGGCATTGGTCCATT 160
AAATACAAACCGCGCCCGAAATTTAATATGCCGAGAAAAATAAATAGCGCCCGAATTGATGTGTTCCTTTCTCTGGCG 240
TGGGTCCCACTATATTCAAAAATAATAATAAACAAACCATTTAATAGTGTGTGGTACAAGAGTGGCCAAATTGTTGA 320
AATATATTCAAACACTCTACAGCGCCGTATGATGACTACATCAATATCGAGAACTTGGACCGGAAGGCTAATTATTATA 400
TAAGGCGATGAACCAGAGTCTCCTCTCAGGCAAGATGGATAAGGAAGAAATAATGGCCGTATACCTGCTCTTCTCAGAT 480
ACGGACCATCAGAAATAATTAATGTTTGCACGGCATTTCGCATCGCCGAAGCAATTGAGGAAAATCCACTACCTCTTCTG 560
GAATATGTCAACTCACTTCAGAAGCATCAAGAAGCCGTGTACAGAATGCTCCGCGAAGCCGGGAAAAC TAGTAGTGAAGT 640
AGAAGAAGTCACTGAGGCTTATACCGTCGTTGACAGATCCGGCAAGAAGAAAAAATGAAAGTGAAGACCTCGGCGAAG 720
TTGGTACCTCGAAGGGAGTTGTTTCTACCGGATGGTACTGTCATACATCCGGCGTCGTATATTTATATTAGTATTAGTA 800
TTGTAATCAGTAGTCTGTACTGCTTATGTTAATTAATGTTTTCGAAAGAAATGAACTATTGTTCAGAAATATATAATCAA 880

Ren gene stop codon

ACGAAGATTAATACATAACA TTAATACAAGTTGAACTGAATATGGGCCGACTGTTGTACAGATAACATATATCTATCAA 960
TCGGCCCACGAAATGTTTTACAGCGCTACTACAGCATTAAATAGAAATTACACCGAAATCATACAAATACTGTAAAATAC 1040
GATTGCGAAATACATGTAGAAAACCACTAGATGATGGCCGTAGAGTTGTCATATGGTGAAAGTAAGCCAGCATTATGT 1120
ATCTGTAGAGCACGTCGAAGATTGTAGTTGAACCGTATTTCAACGGTTATACGGT CCCAGTTTTCCCGCACGGTCTGTC 1200

Potential PolyA signal

CAAA TGATCCTTTATACGAAAATACAATGGGTTGGATAATTCCCAGTAAAATACACCCCGGTTTGCCTGTTCAACTGTGA 1280

Ren gene start codon Ren gene stop codon

TGTCATCACCGGTTCTGTAAAT CATTTCAGATGCCGGAATTTGTTGAAGCTGAATTTGTACAAGGC TTA ACTTGAGAAC 1360
TAGGACCTATGCTGTCTCGGTCGACTCAGGCGTTTGGTACAGAGGGCAATATATGAAATCAAATACCGCATTCTTCAA 1440
GTCCATTACGTAAACCAAACCTGGTCGGGTCTATCAAGATAATCTTTGTAGCTGCAGTCATTACCAGGATTACACAGAAC 1520
TATTGACGGTACTTCTCCCTTGATCCTGACAGGCTTGCCATATTTGCAATTAGTTAGCCAATGCCGTTGAGACCCGATCA 1600
AATGTTTCCAATGCTTCATACGGAGATAATTTGGCTCGATGTCGTCATTACATTGTAATAACATCATTGTTGAATACC 1680
TTCGGATTGAAGTCCAGATGACCACTTATGTAGTTATGGGGGCCAAAGAACGAGCCCATGCTGTTTTTCCAATCCGACT 1760
GGGCCCTTCAAGAATCAAACCTCGTTGGCCGGTCCAAACCAACCTTTGGGACTATTATATTATCTCCATTAGGCCGCGCAG 1840
CGTCATTGATCTGAAAATATGATTAATCCAATCATTATCACAGGAGTTATCACAAACGTAGAATAGTTCCAATTAGAT 1920
ATATATTTTTCAGGAGCTTTGTAAATATCCTGTCATAGTTAGCCTGTAGATTATGATATGTAAATGTAAGCTCTGGG 2000
ATCCCTCTCTCTAATTAATTGCAAAGCATCTTCAGCAGAACCAAAATTTAAAGCATCTGCATATACGGTTGATAGACTAC 2080
GTCTTCCATCCCTTGTAGATCTACCATCAACCTGGAACCTCTCCCAGTCCATGTAATCCCTCCTTTCTCGATGTATTCC 2160

A-T Rich Region

TTAACGTCGGAGCTGGATTAGCTGACTGAATATTTGGATGG AATTTTATAGAA GATGAA GGTGAGTGAAGATCGAATAA 2240
CCTGTTAATTTGTTACCTGTGCTCTGCCCTCCAACCTGTATTAATA CATGTAGATGTGGAGAACCATCCTCATGTAACTCTC 2320
TAGCAACCCGATGTACTTCTTGTGCTTGCAGG GCACACTTCTAGCAATTGCTGTAATGCCTCCTCCTTTGTTCAGAGAA 2400

Ren gene start codon

CACCTGGCGTAGGTCAGGAAGTAATTTCTAGCATTTAATCTGAACTGCCTCGGTCTGACAGG CATATTTGCTAAATAATA 2480

Potential Itron sequences (type B) [5]

TA **BAACACT** AAGCTAATGGACCGGCAGTAATGTAAATGAGTCTTACACACTTCTGCAACAGGGAAATTTCAAAAACCCCT 2560
 Potential complementary- and virion-sense gene TATA box
 Inverted repeat sequence [2]
 TAGCAATCGGTGTATGTGTGTATGT **TATATAGGCGACCC** CAGAAGGCTCAGAAGGCTCAGAAGGCCACA **AGCCATCC** 2640
GT 2642

A- EGV1 annotated sequence from *D. nummularia* acc. 313

Virion strand origin of replication
 Inverted repeat sequence
TAATATTACACGGATGGCC GCCCCGTCCGTACTGCAACGTGTCGAGTTTGGATTGGTCTTCGAAAAAGAAGAAGTACTTG 80
 CTGTCTCCCTTCGTACGTCCACGTGGCTTGTCTCATGAGCCATTCAAATTTAAACCGGCGCCGAATCCTTTTGCCT 160
 CTTTCTCAGTTTGTGGGACCCACTAATTTCAAAAATAAATATTAATAAAGCAGAAAAATAATGGAGTGGGTGCGAGACAATA 240
 TTGTTTGTACATAAAGTGGCCCATAGAGGGGAAATAAATTTCAAAAACCCCTACAGGGCCGTATGATGACCGGTTACAGTTAT 320
 TGAGAACTTGGATTGGAAGAATTAATAATTATATAAAGTGCCGAACCGGAGTTACCTCTCAGGCAAGATGGAGAAGGAAG 400
 AAATCATGGCCGTATACCTGCTCTTCCCTCAGATACGGTCTATCTGCAATAATTAAGAGTGTACATGGTATTTGTATTGCT 480
 GAAGCCATCAAGGAAATCCCTTACCTCTCCTTGAGTATGTCAACTCACTGTAGAAGCACCAAGATGCTATGTACCGAAT 560
 Ren aene stop codon
 GCTCCGCGAAGCCAGCAAGTCCAGCAAGGATGTACAGGAGGTTA **CCGATGCATATACAGTGGTCGATAGAGCCAGTAATA** 640
 AAAAAAATGAACTAGAGGACCTCGGTGAAGTTGGGACCTCTAAGGGTGTGGTGTCCACTGTTTGGTGAAGTACATGTA 720
 Ren aene start codon
 TTTGTATAAC **CAT** ATATGGAATGAAGTAATGTATCACGTATGTATATTTGTGCTTATTATCAGTTATTCCAAAAATAAA 800
 GAATGTCAATTTGGGAAAAATTATATGTATGTTAAGACATAACCGGTGCCTTATATAAAAAATAAAACACATATTTGCATA 880
 ATAAATAATTTTAAATACAGTTGAACATGACATGGGCCGACTGTTGTACATATAGAATAAACCCATCAAACCGGCCACG 960
 Rep aene stop codon
 AAATGAGTAACAACACATATTACAGTATTAATAAAAATTACACCAAATTTATTCAAATACTGCAAGATGCGA **TTA** CAAAA 1040
 TACATGAAGGAATCCCGCAGATGAAGGCCGTAGATCGAGTTATCCACAGGGTGAAGATAGCCAGCATTGTGTATTTCGT 1120
 AGTGCACGACGTAAATTTAGTTGAACCATATTTGACGGTTATACGGTCCAGTTGTTATTGAACGGTCTGTCTATATG 1200
 ATCCTTTTATACGAAATAACAATGGCTTGACAATTTCCAGTAAATACGCCTCTGTTTGCCTGTACCACAGTAATGTCTT 1280
 CACCTGTTCTGAAATCCATCCGAGATACGGTTATGGAATCAAACCTCGGCATTCCTTCAGAGTCCATTCACGCAAACCTAAG 1360
 TTGTGGGTGCGATCAAGATAATCCTTGTAGCTTGAATCATTTCCAGGATTACACAGAACTATAGACGGGACTCCACCCTT 1440
 GATCTTGACAGGCTTCCATACTTGCAATTCGTTTGCCAATCCCGTTGAGACCAATCAAATGTTTCCAGTGCTTGAGCC 1520
 GTAGATAACTTGGTTCTACGTCAATGATGTTGTAGAAGATATCATTTGGAATGTCTTCGGATTGAAGTCCAGATGA 1600
 CCACTGATGTAATTTATGGGGCCCAATGACCGAGCCATGCGGTTTCCCAACCGACTGGGCCCTTCAAGGATCAAGCT 1680
 GGTGGCCCTGTCCAGCCCAACCTTAGGGATTATATATATCTCCATTAGGCCGTGCAGCGTCATTAATATGAAAATTAT 1760
 Rep gene start codon
 TGTTAATCCAAT **CAT** TCATCACAGGATTTACATGAGGATGGTTCTCCACATCTAGATG **GGGAGAAG** CATCCTCACATAAC 1840
 TCCCTTGCAACCCGTATGTACTTCTTGTGCTTGAAGAGACACCTCCAGTAATTGCCGTAATGCCTCCTACTTGGTTAA 1920
 Potential iteron sequences (type B) [5]
AGAAC ACTTGGCATAAGTCAGGAAGTAATCTTATCATTCATCTGAACTGTTTAGGTCTGACAGGCATCTCTGCGAAAA 2000
 AA **GGGAGAA** GAATGGCACTTCTGTAAATAATGGCTTACACACTTCTCCAACAGTATTTTGCAAAAACCCATATCTATGT 2080
 Potential complementary- and virion-sense gene TATA box

GTGTAAGTGTGTATCTTATATAAGCGAACCCCAAAAAGCCCCAGAAAGCCATCCGT

Figure S1

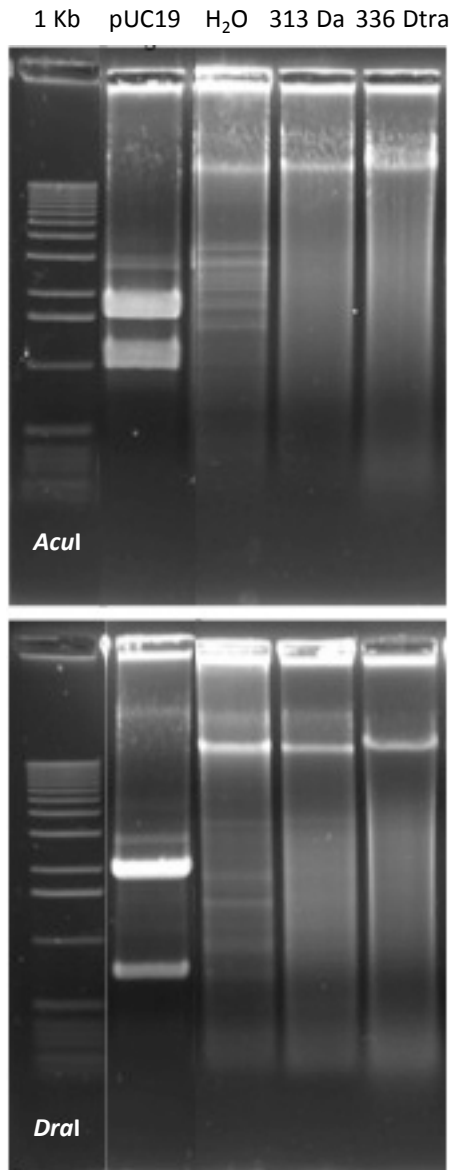


Figure S3

Rolling circle amplification using a cocktail of random and specific primers used for detecting EGV1 or EGV2.

1. Ladder
2. pUC19
3. H₂O
4. *Dioscorea alata* (acc. #313, seedlings)
5. *Dioscorea transversa* (acc. #336)

Maximum likelihood tree of eighteen 402 bp partial *rep* sequences obtained using the “*ren-ren*” and “*rep-rep*” primer pairs.

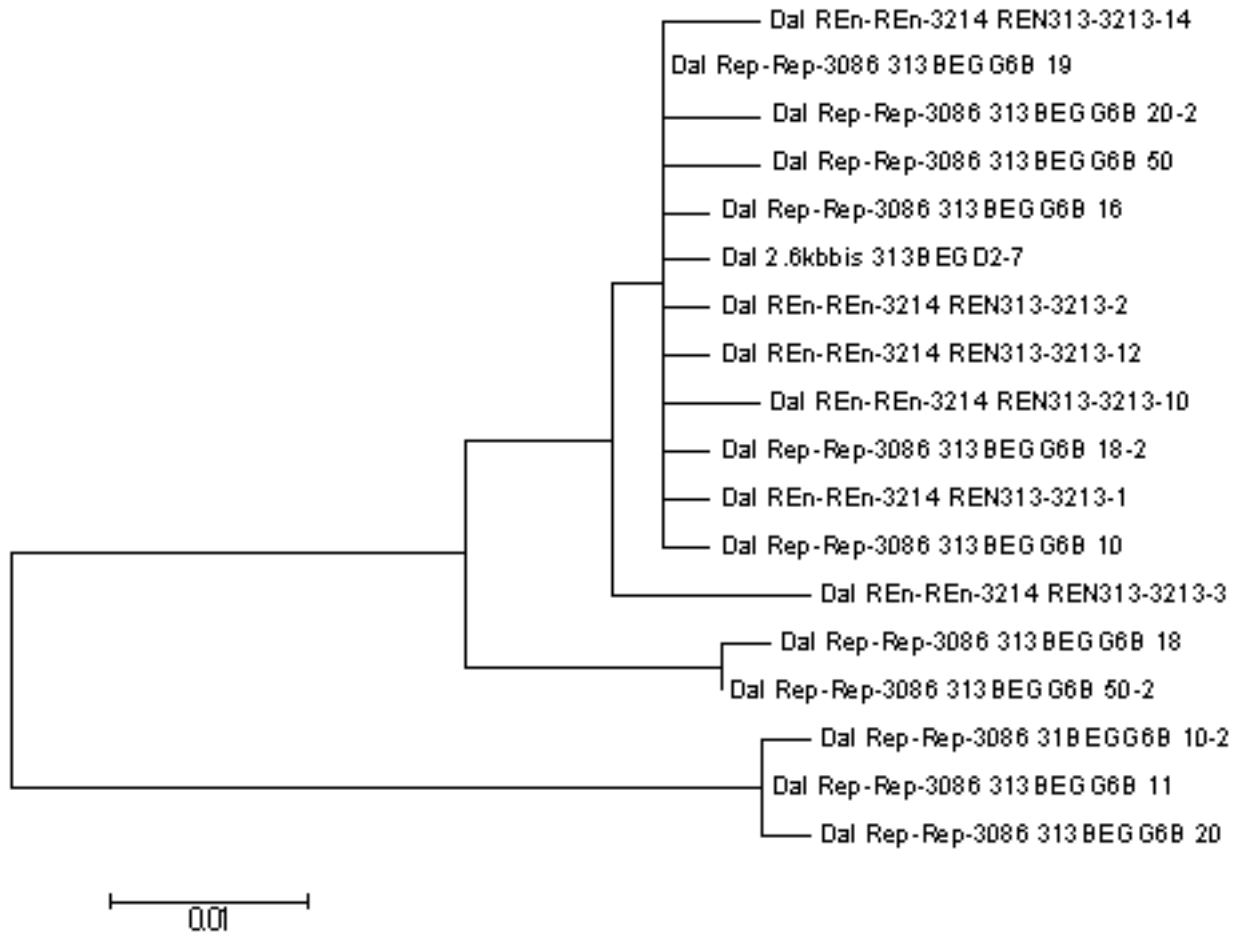


Figure S4

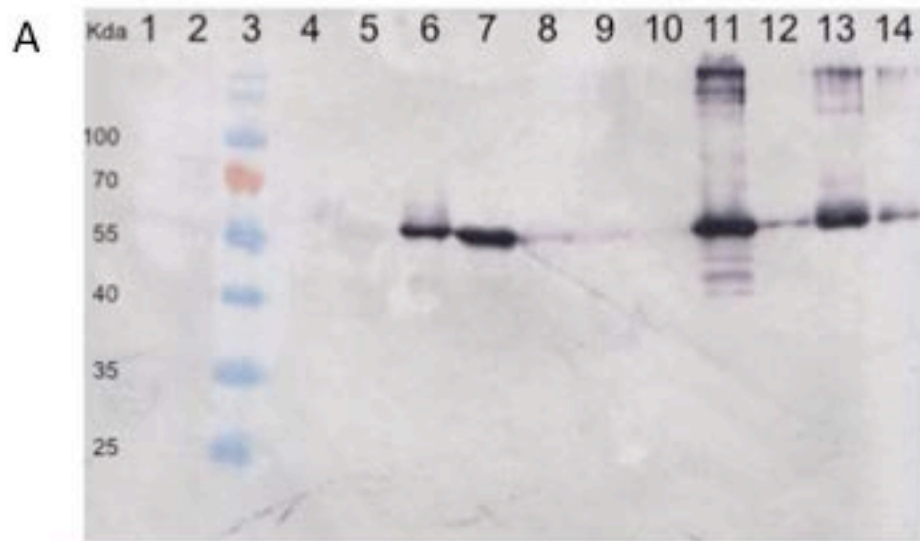


Figure S5

A: Western blot on the total protein extracts of several *Dioscorea* species, using an antibody directed to a Rep peptide of EGV1

B: The nitrocellulose membrane stained with Ponceau S dye for protein detection during western blotting.

Lanes:

- 1) Proteins extracted from a turnip plant infected by CaMV
- 2) Tomato plant infected by TYLCV
- 3) Ladder
- 4) *Dioscorea trifida* (accession#64)
- 5) *Dioscorea sansibarensis* (accession#269)
- 6) *Dioscorea bulbifera* (accession#272)
- 7) *Dioscorea dumetorum* (accession#47)
- 8) *Dioscorea togoensis* (accession#114, seedling)
- 9) *Dioscorea praehensilis* (accession#255)
- 10) *Dioscorea rotundata* (accession#118, seedling)
- 11) *Dioscorea nummularia* (accession#335)
- 12) *Dioscorea alata* (accession#297, seedling)
- 13) *Dioscorea alata* (accession#313 seedling)
- 14) *Dioscorea alata* (accession#402 seedling)..

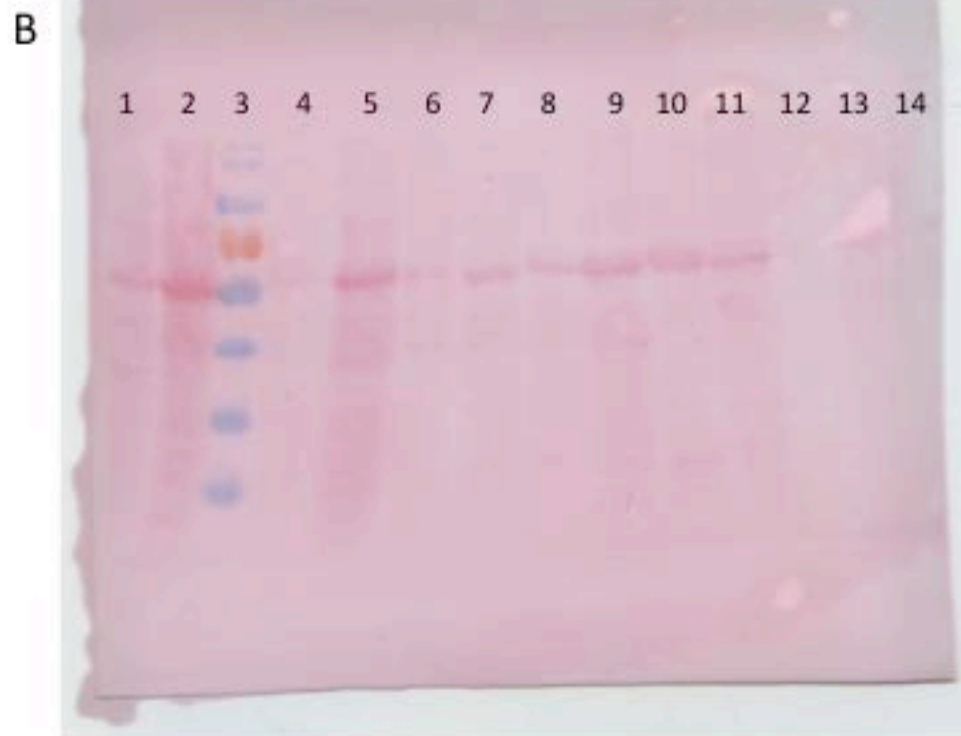
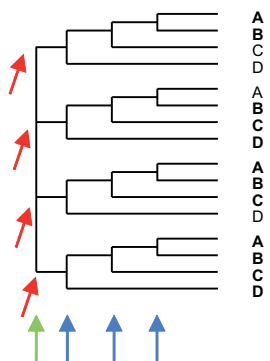


Figure S6

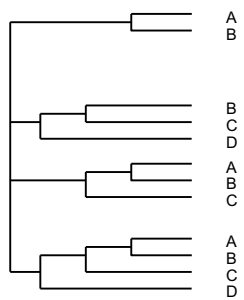


0.2

A Sequence duplication only at integration



Sampling

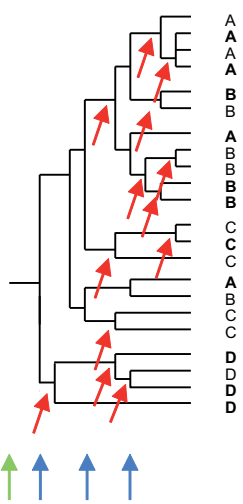


- ← Speciation events
- ← Integration event
- ← Sequence duplication events

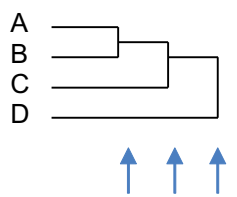
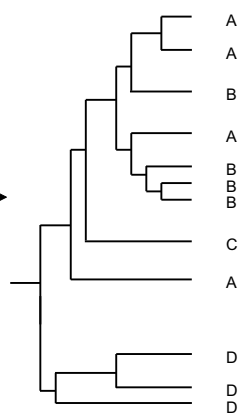
Real endogenous sequence tree

Tree inferred after sampling

B Repeated sequence duplications post-integration



Sampling



Chloroplast genome tree

Real endogenous sequence tree

Tree inferred after sampling

Figure S7

Supplementary Table 1. List of primers used in this study

Material and method part	Primer name	Sequence	
Full genome amplification of EGV1	YLCV2F	AGATGTGGAGAACCATCCTC	
	YLCV1R	GGTTATTCGATCTTCACTCAC	
Whole EGV1 and EGV2 <i>rep</i> genes amplification	<i>rep</i> _EGV1_D_alataF	GCCGGAATTTGTTGAAGCTG	
	Primer Set1	<i>rep</i> _EGV1_D_alataR	GGTCCATTAGCTTAGTGTCT
	Primer Set2	<i>rep</i> _EGV1_D_nummulariaF	TCTCACCTGTTCTGAAATCC
		<i>rep</i> _EGV1_D_nummulariaR	GTTGGAGAAGTGTGTAAGCC
	Primer Set3	<i>rep</i> _EGV2_F	GGAGACGTATACCGAACAAAC
		<i>rep</i> _EGV2_R	TGGGTTATGAGTTGTTACTTGT
	Characterization of multiple EGV1 repeats from <i>D.alata</i>	<i>ren-ren</i> _D_alataF	GTCATACATCCGGCGTCGT
<i>ren-ren</i> _D_alataR		CAGCTTCAACAAATTCGGC	
<i>rep-rep</i> _D_alataF		AATCCAATCATTATCACAGGA	
<i>rep-rep</i> _D_alataR		GAGGATGGTTCTCCACATCT	
Distribution of EGVs within the genomes of <i>Dioscoreacea</i> sp.		EGV1_Detection_1F	GTGAGTGAAGATCGAATAACC
		EGV1_Detection_1R	GAGGATGGTTCTCCACATCT
	EGV1_Detection_2F	CGCGCAGCRTCATTRATCTG	
	EGV1_Detection_2R	TGGGGWGAGTTYCAGGTTGA	
	EGV2_Detection_1F	AGGAATGGAAGTCAAGTCGTA	
	EGV2_Detection_1R	AGGTTTCAGATTCCAGCTATTC	
	EGV2_Detection_2F	TTCCAGGTGTTCTTCTATCTC	
	EGV2_Detection_2R	GTCACTACCAASRAAAYGCTTC	

Supplementary Table 2. Prevalence of EGV1 and EGV2 sequences among a collection of 9 *D. alata* plants collected worldwide, including two *D. alata* plants grown from true seeds under virus-free conditions.

Species	Accession number	Sampling location	EGV1 292 bp	EGV2 1274 bp
<i>D. alata</i>	301	Brazil	+	+
	440	Costa-Rica	+	+
	71	Cuba	+	+
	408	Ghana	+	+
	314	Haiti	+	+
	313, 402	India (grown from true seeds)	+	+
	MDG009	Madagascar	+	+
	167	Vanuatu	+	+

Supplementary Table 3. Viral sequences recovered from partially purified virions and EST screening.

Procedure/plant sample	Virus-like sequence		Best hit obtained from NCBI Blast (BlastN or BlastX)								
	Plant vs EST	Name	Length (bp)	Family	Genus	Accession Number	Species	Identity (%)	E-value	Gene	Blast
Virion - <i>D. alata</i> - VU564a		EGV1	209	<i>Geminiviridae</i>	<i>Begomovirus</i>	DQ665866	SiGMBV	71	3e-08	Rep	N
ESTs - <i>D. alata</i> (dbEST Id : 71472229)		HO836974	274	<i>Caulimoviridae</i>	<i>Badnavirus</i>	L14546	CSSV	72	2e-11	Polyprotein	N
		HO838291 (EGV2)	316	<i>Geminiviridae</i>	<i>Begomovirus</i>	HM585443	SiMBoV2	72	1e-20	Rep	N
		Contig3987	2306	<i>Potyviridae</i>	<i>Macluravirus</i>	AB044386	ChYNMV	72	1e-75	CP	N
		HO833648	325	<i>Secoviridae</i>	Unassigned	NP733954	SMoV	45	1e-17	RdRP	X
		Contig1995	2308	<i>Secoviridae</i>	<i>Sadwavirus</i>	NP620567	SDV	21	3e-6	Polyprotein	X
		Contig2044	2535	<i>Secoviridae</i>	Unassigned	NP733984	SMoV	22	1e-16	CP	X
		Contig2047	2027	<i>Secoviridae</i>	<i>Sadwavirus</i>	NP620567	SDV	23	5e-16	Polyprotein	X
		Contig2049	2577	<i>Secoviridae</i>	<i>Sadwavirus</i>	BAD12076	SDV	21	5e-19	CP	X
	Contig2216	2950	<i>Secoviridae</i>	Unassigned	EU419645	BRNV	84	3e-08	RdRP	N	

Supplementary Table 4. Selection analyses results (*N. benthamiana*). For each test, each clade was assigned to a foreground partition (FG), or one of two background partitions (BG1 and BG2). We report p-values for significance tests for purifying selection along foreground lineages, and the parameter point-estimates for the foreground partition under the alternative model (which allows purifying selection).

GRD3	GRD5	Connecting branch	p-value for neg. selection	ω_1	p_1
BG 1	BG 2	FG	7.9×10^{-12}	0.043	1
BG 1	FG	BG 2	0.44	0.75	1
FG	BG 1	BG 2	0.995	0.94	0.65
FG	FG	BG	0.55	0.82	1