



OPEN

DATA DESCRIPTOR

The phased *Solanum okadae* genome and *Petota* pangenome analysis of 23 other potato wild relatives and hybrids

S. R. Achakkagari¹, I. Bozan¹, J. C. Camargo-Tavares¹, H. J. McCoy², L. Portal³, J. Soto³, B. Bizimungu⁴, N. L. Anglin^{3,5}, N. Manrique-Carpintero^{3,6}, H. Lindqvist-Kreuze³, H. H. Tai⁴ & M. V. Strömvik¹✉

Potato is an important crop in the genus *Solanum* section *Petota*. Potatoes are susceptible to multiple abiotic and biotic stresses and have undergone constant improvement through breeding programs worldwide. Introgression of wild relatives from section *Petota* with potato is used as a strategy to enhance the diversity of potato germplasm. The current dataset contributes a phased genome assembly for diploid *S. okadae*, and short read sequences and *de novo* assemblies for the genomes of 16 additional wild diploid species in section *Petota* that were noted for stress resistance and were of interest to potato breeders. Genome sequence data for three additional genomes representing polyploid hybrids with cultivated potato, and an additional genome from non-tuberizing *S. etuberosum*, which is outside of section *Petota*, were also included. High quality short reads assemblies were achieved with genome sizes ranging from 575 to 795 Mbp and annotations were performed utilizing transcriptome sequence data. Genomes were compared for presence/absence of genes and phylogenetic analyses were carried out using plastome and nuclear sequences.

Background & Summary

Potato is currently the third most important crop for human consumption with increasing production in the developing world and it is an important sustainable global food security crop for climate-smart agriculture^{1,2}. Breeding for improved nutritious potato varieties suitable for sustainable production and climate change resilience will be aided by introgression of diverse genes from native Andean cultivars and landraces along with wild *Solanum* relatives of potato³. The genebank at the International Potato Center (CIP) and other genebanks globally, carry diverse germplasm used by breeding programs for genetic improvement of the crop^{4,5}. The CIP genebank has potato cultivars, landraces, and wild *Solanum* germplasm collected from a wide range of environments in North and South America with variation in abiotic and biotic stress. Native Andean cultivars and landraces also carry beneficial nutrients and bioactive compounds for human health^{6,7}.

A reference genome sequence for a doubled monoploid potato clone DM1-3 516 R44 (DM) was released in 2011 and later improved⁸⁻¹⁰. Low depth sequencing of potato varieties and wild species aligned against the DM reference was used to assess genome diversity, copy number variation (CNVs), introgressions, and selective sweeps¹¹⁻¹³. Since then, the genomes of additional potato clones, diploid as well as tetraploid, have been sequenced and *de novo* assembled¹⁴⁻²⁰.

Potato and its wild relatives are in the genus *Solanum* section *Petota*, which consists of over 100 species that form tubers²¹. Increasing knowledge on the genetic potential of species in section *Petota* will increase capacity for their use in enhancing potato germplasm. Genomes of several wild species are also available²²⁻²⁵ and a *Petota* super pangenome representing 60 species was recently released²⁶. The current data set includes a phased

¹Department of Plant Science, McGill University, Sainte-Anne-de-Bellevue, QC, Canada. ²Department of Chemistry, University of New Brunswick, Fredericton, NB, Canada. ³International Potato Center (CIP), Lima, Peru. ⁴Agriculture and Agri-Food Canada Fredericton Research and Development Centre, Fredericton, NB, Canada. ⁵USDA ARS Small Grains and Potato Germplasm Research, Aberdeen, ID, USA. ⁶Present address: Alliance of Bioversity International and International Center for Tropical Agriculture (CIAT), Cali, Colombia. ✉e-mail: martina.stromvik@mcgill.ca

Short ID	Species (Taxonomic group, Spooner)	CIP/ID number	other ID	ploidy	Source location	Important traits
oka15	<i>S. okadae</i> (Clade 4)	—	—	2n = 24	AAFC Fredericton- CPGR	resistance: CPB, moderate late blight; drought, cold tolerance
B1595106	<i>S. brevicaulis</i> (= <i>S. oplocense</i>) X <i>S. tuberosum</i> (Clade 4)	PGR-15951-06	—	2n = 48	AAFC Fredericton- CPGR	<i>S. brevicaulis</i> (= <i>S. oplocense</i>) X <i>S. tuberosum</i> backcross clone; resistance: CPB
blv1353	<i>S. boliviense</i> (Clade 4)	CIP761353.009	—	2n = 24	CIP	drought tolerance
buk0368	<i>S. bukasovii</i> (Clade 4)	CIP760368.015	—	2n = 24	CIP	drought tolerance
chc0917	<i>S. chacoense</i> (Clade 4)	CIP760917.1	PI 197760	2n = 24	CIP	resistance: bacterial wilt
chq2573	<i>S. chiquidenum</i> (Clade 3)	CIP762573.219	OCHS 12566	2n = 24	CIP	resistance: late blight
cmm1080	<i>S. commersonii</i> (Clade 4)	CIP761080.201	—	2n = 24	CIP	resistance: bacterial wilt
cxax70P5	<i>S. commersonii</i> X <i>S. andigena</i> (Clade 4)	H6S170P5	—	2n = 48	CIP	cultivar- Winay; frost tolerant variety
etb0161	<i>S. etuberosum</i> (outgroup)	—	PI 245939	2n = 24	USDA Aberdeen, ID	resistance: insect, PLRV; cold tolerance, self compatible
grc1498	<i>S. gracilifrons</i> (Clade 4)	CIP761498.010	—	2n = 24	CIP	drought tolerance
ifd1359	<i>S. megistacrolobum</i> (Clade 4)	CIP761359.032	—	2n = 24	CIP	drought tolerance
lgl2830	<i>S. lignicaule</i> (Clade 4)	CIP762830.057	—	2n = 24	CIP	drought tolerance
mga1403	<i>S. megistacrolobum</i> (Clade 4)	CIP761403.208	OCH 12032	2n = 24	CIP	resistance: bacterial wilt
pcs2126	<i>S. paucisectum</i> (Clade 3)	CIP762126.217	OCHS 14818	2n = 24	CIP	resistance: late blight, CPB
pur1868	<i>S. piurae</i> (Clade 3)	CIP761868.202	OCH 13959.6	2n = 24	CIP	resistance: late blight, CPB
S12	<i>S. tarijense</i> (Clade 4)	DPM-S12	PI 473243	2n = 24	AAFC Fredericton- CPGR	secretory trichome A hairs; resistance: peach aphid, CPB
S15	<i>S. gandarillasii</i> (Clade 4)	DPM-S15	PI 545864	2n = 24	AAFC Fredericton- CPGR	drought tolerance
S3	<i>S. commersonii</i> (Clade 4)	DPM-S3	PI 472837	2n = 24	AAFC Fredericton- CPGR	leaf roll susceptible, cold tolerance
S7	<i>S. pinmatisectum</i> (Clade 1 + 2)	DPM-S7	PI 275236	2n = 24	AAFC Fredericton- CPGR	resistance: insect, PVY; PVX susceptible, Verticillium wilt susceptible, frost susceptible, drought tolerance
SH7_18_3	<i>S. tarnii</i> X <i>S. tuberosum</i> (Clade 1 + 2/4)	PGR-7/18/3	—	2n = 48	AAFC Fredericton- CPGR	F1 <i>tarnii</i> somatic hybrid (<i>S. tarnii</i> x <i>tbr</i> (4x))
spl0147	<i>S. sparsipilum</i> (Clade 4)	CIP760147.7	PI 1230502	2n = 24	CIP	resistance: bacterial wilt
tcn8662	<i>S. tacnaense</i> (Clade 4)	CIP762866.026	—	2n = 24	CIP	drought tolerance
tcn8663	<i>S. tacnaense</i> (Clade 4)	CIP762866.038	—	2n = 24	CIP	drought tolerance
trp2833	<i>S. tarapatatum</i> (Clade 4)	CIP762833.025	—	2n = 24	CIP	drought tolerance

Table 1. A list of the 24 potato wild relatives (*Solanum* sp) accessions selected for this study with their detailed descriptions.

genome assembly of the wild species *S. okadae*, for which short reads and the mitogenome were previously published^{26,27}. In addition, the study contributes Illumina short read sequences (50–100x sequencing depth) and *de novo* assemblies (3,328–59,223 N50) for genomes of 16 wild species in section *Petota* that were noted for stress resistance and of interest to potato breeders. An *S. etuberosum* genome from outside section *Petota* was also included in this data set, as were three polyploid hybrids. For most of the species in this project, these are the first genome sequences released of their species. Taxonomic groupings of the species of section *Petota* into Clade 1 + 2, Clade 3 and Clade 4 were represented in the genomes sequenced.

Transcriptome sequence data was generated and used for annotation. Both nuclear and organellar genome sequences are provided. Comparative analysis of the presence/absence variation against the *Petota* super pangenome provided insight into the diversity of the species in the phylogenetic analysis.

Methods

***Solanum okadae* (OKA15) genome sequencing.** Botanical *S. okadae* seed sourced from US Genebank NRSP-6 (PI 458367) were sown on agar media. Individual clones were propagated *in vitro*. Seventeen lines were planted in the greenhouse and self-pollinated. While abundant flowering was noted, a single line, OKA15, was selected for further study^{26,27}. The 10X Genomics short reads were previously published^{26,27}. High molecular weight DNA for *S. okadae* OKA15 was prepared using a modified CTAB procedure²⁸ and used for Hi-C (chromatin interactions), Pacific Biosciences (PacBio) and Nanopore sequencing. The Hi-C library preparation was done using Phase Genomics Proximo Plant Hi-C version 4.0, and sequenced at Génome Québec's Sequencing centre, Montreal, Canada, on Illumina NovaSeq. 6000 platform in PE 100 bp mode with ~3B reads, while both the PacBio libraries (sequenced on Sequel II system using the SMRT technology (HiFi) at a depth of 100X,) and the Nanopore ONT libraries (sequenced on the PromethION at a depth of 100X) were prepared and sequenced by Novogene (Novogene Co., Ltd, Beijing, China).

Genome and transcriptome sequencing. A total of 23 additional accessions representing potato wild relatives or hybrids with agronomically important traits such as disease resistance, cold and drought tolerance,

Short ID	Species	Sequencing mode	WGS			RNA-seq	
			Tissue	Sequencing depth (X)	Raw data (Gbp)	Tissues	Raw data (Gbp)
oka15	<i>S. okadae</i>	PE150/10X	Leaf	100	66.2	—	—
oka15	<i>S. okadae</i>	PacBio	Leaf	100	64.6	—	—
oka15	<i>S. okadae</i>	Nanopore	Leaf	100	60.3	—	—
oka15	<i>S. okadae</i>	Hi-C	Leaf	—	300	—	—
oka15	<i>S. okadae</i>	PE150	Leaf	100	—	Leaf, Tuber, Sprout	35.9
B1595106	<i>S. brevicaulis</i> (= <i>S. oplocense</i>) X <i>S. tuberosum</i>	PE150	Leaf	50	41.5	—	—
blv1353	<i>S. boliviense</i>	PE150	Leaf	100	81.7	Leaf	9
buk0368	<i>S. bukasovii</i>	PE150	Leaf	100	89.5	Leaf	7.1
chc0917	<i>S. chacoense</i>	PE150	<i>In vitro</i> plantlets	100	81.5	Leaf, Shoot, Flower, Tuber	68.9
chq2573	<i>S. chiquidenum</i>	PE150	<i>In vitro</i> plantlets	100	92	Leaf, Shoot	45.3
cmm1080	<i>S. commersonii</i>	PE150	<i>In vitro</i> plantlets	50	43.5	Leaf, Shoot, Tuber	53.1
cxa70P5	<i>S. commersonii</i> X <i>S. andigena</i>	PE150	Leaf	50	40.3	Shoot	6.4
etb0161	<i>S. etuberosum</i>	PE150	Leaf	100	80.3	Leaf	9.2
grc1498	<i>S. gracilifrons</i>	PE150	Leaf	50	41.2	Leaf	7.5
ifd1359	<i>S. megistacrolobum</i>	PE150	Leaf	100	88.1	Leaf	7.8
lgl2830	<i>S. lignicaule</i>	PE150	Leaf	100	84.1	Leaf	8.7
mga1403	<i>S. megistacrolobum</i>	PE150	<i>In vitro</i> plantlets	100	82.5	Leaf, Shoot, Flower	54.5
pcs2126	<i>S. paucisectum</i>	PE150	<i>In vitro</i> plantlets	100	80.7	Leaf, Shoot	48.7
pur1868	<i>S. piurae</i>	PE150	<i>In vitro</i> plantlets	100	80.3	Leaf, Shoot	50.8
S12	<i>S. tarijense</i>	PE150	Leaf	100	82.9	Leaf, Tuber	13.2
S15	<i>S. gandarillasii</i>	PE150	Leaf	100	84.9	Leaf, Tuber	20.8
S3	<i>S. commersonii</i>	PE150	Leaf	50	44.8	Leaf, Tuber	14.4
S7	<i>S. pinnatisectum</i>	PE150	Leaf	100	80.8	Leaf	21.9
SH7_18_3	<i>S. tarnii</i> X <i>S. tuberosum</i>	PE150	Leaf	50	41.6	—	—
spl0147	<i>S. sparsipilum</i>	PE150	<i>In vitro</i> plantlets	100	88.7	Leaf, Shoot	46
tcn8662	<i>S. tacnaense</i>	PE150	Leaf	100	80.4	Leaf	8
tcn8663	<i>S. tacnaense</i>	PE150	Leaf	100	80.3	Leaf	8.7
trp2833	<i>S. tarapatanum</i>	PE150	Leaf	50	40.3	Leaf, Shoot, Flower	53.2

Table 2. Details of the tissues used and amount of data generated for WGS and RNA-Seq. The raw data is calculated as: *number of reads***sequence length* (150 bp).

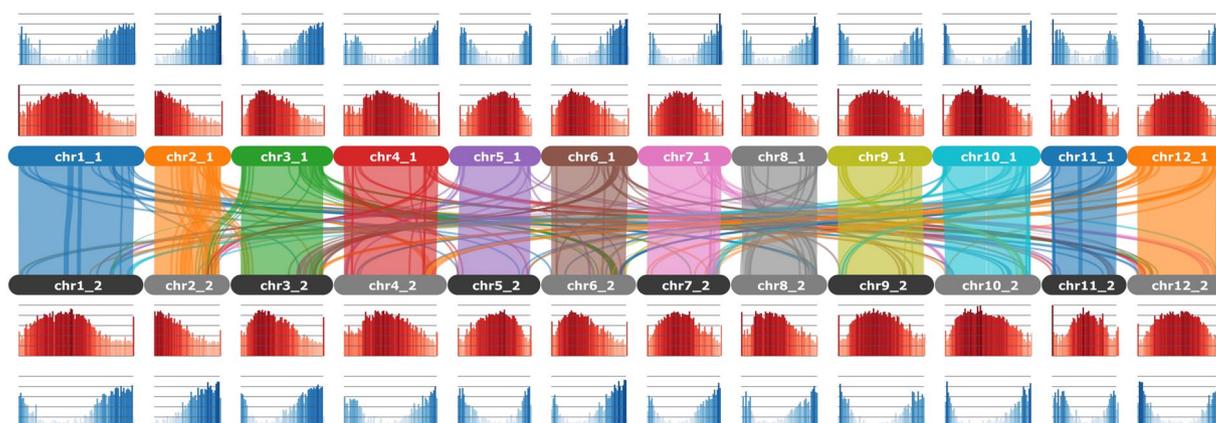


Fig. 1 The haplotypes of the phased *Solanum okadae* (OKA15) genome.

and other biotic and abiotic stresses were selected for this study (Table 1). The genomic DNA of these accessions was isolated from *in vitro* or greenhouse grown plantlets or leaves. The library construction, quality control, and sequencing were done by Novogene (Novogene Co., Ltd, Beijing, China). The genomic DNA was fragmented, ligated with adapters, and PCR amplified to construct single libraries with an insert size of 350 bps. The sequencing was performed on the Illumina NovaSeq. 6000 platform, in a paired-end mode (2 X 150 bp) at ~50x or ~100x depth (Table 2). Similarly, total RNA was extracted from these accessions, including from OKA15, from leaf, shoot, flower, and/or tuber tissues. The RNA-Seq library construction used standard protocol (Novogene) and sequencing was done on the Illumina NovaSeq. 6000 platform (PE150) (Table 2).

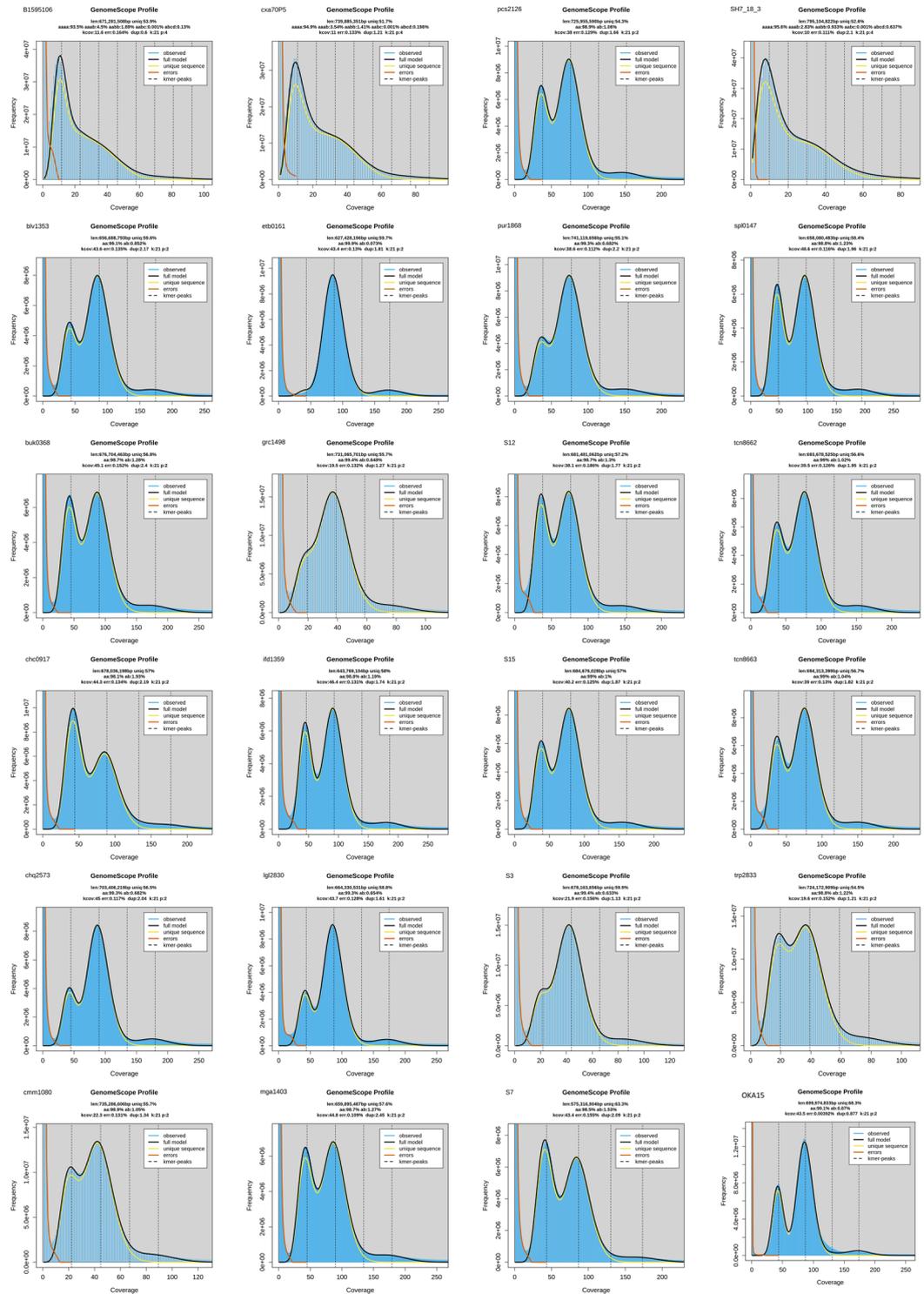


Fig. 2 The k-mer coverage plots of each of 23 genome sequences from potato wild relatives (*Solanum* sp).

***S. okadae* OKA15 phased genome assembly and quality control.** The Nanopore reads were adapter trimmed using *porechop v0.2.4* (<https://github.com/rrwick/Porechop>) and default parameters. *hifiasm v0.16.1-r375*²⁹ was used to generate a haplotype-resolved assembly using the PacBio HiFi and the Hi-C reads (Fig. 1). Organellar genomes were removed from both resulting haplotype assemblies. The HiFi and Nanopore reads were used to scaffold the haplotype assemblies. Hi-C reads were trimmed using *fastp v0.23.1*³⁰, and mapped to haplotype assemblies individually using the Arima mapping pipeline (https://github.com/ArimaGenomics/mapping_pipeline) and default parameters. The alignments were filtered to keep only the unique alignments (*samtools view -q 40*)³¹. The scaffolding was performed using *YaHS*³². Hi-C maps were generated using *Juicer*³³ and manually reviewed using *JBAT*³³. Any duplicated contigs from the assemblies were removed using *mummer*³⁴.

Short ID	Species	%Heterozygosity (Max)	Estimated genome size (bp)	Final assembly size (bp)	N50	#Contigs	Largest contig (bp)
OKA15	<i>S. okadae</i>	0.88	699,974,833	728,982,852 (hap1) 725,850,496 (hap2)	58,550,880 (hap1) 58,363,667 (hap2)	214 (hap1) 131 (hap2)	84,709,123 (hap1) 83,182,665 (hap2)
B1595106	<i>S. brevicaule</i> (= <i>S. oplocense</i>) <i>X S. tuberosum</i>	11.0	671,281,508	1,007,750,739	3,547	480,556	233,393
blv1353	<i>S. boliviense</i>	0.9	656,688,793	666,810,905	23,940	66,445	286,025
buk0368	<i>S. bukasovii</i>	1.3	676,704,463	716,989,807	22,551	76,513	238,352
chc0917	<i>S. chacoense</i>	2.0	678,036,198	765,361,945	6,068	344,542	100,658
chq2573	<i>S. chiquidenum</i>	0.7	703,406,219	693,400,501	23,961	86,120	269,173
cmm1080	<i>S. commersonii</i>	1.1	735,286,606	720,550,625	17,735	105,239	183,894
cxa70P5	<i>S. commersonii X S. andigena</i>	10.0	739,885,351	957,824,389	3,340	483,218	224,152
etb0161	<i>S. etuberosum</i>	0.1	627,428,196	610,576,211	59,223	34,061	559,139
grc1498	<i>S. gracilifrons</i>	0.7	731,065,701	669,865,430	20,516	75,069	267,424
ifd1359	<i>S. megistacrobolobum</i>	1.2	643,769,104	684,052,840	18,429	84,701	211,877
lgl2830	<i>S. lignicaule</i>	0.7	664,330,531	662,942,901	23,327	68,006	254,715
mga1403	<i>S. megistacrobolobum</i>	1.3	659,895,487	695,075,751	18,244	82,893	188,119
pcs2126	<i>S. paucisectum</i>	1.1	725,955,590	734,795,577	20,376	96,165	234,965
pur1868	<i>S. piurae</i>	0.7	741,119,656	720,192,443	23,712	84,700	248,207
S12	<i>S. tarijense</i>	1.3	681,481,062	711,715,600	19,861	87,536	236,814
S15	<i>S. gardarillasii</i>	1.0	684,676,028	686,228,213	20,834	84,293	308,125
S3	<i>S. commersonii</i>	0.7	678,163,656	660,741,360	23,691	62,463	238,718
S7	<i>S. pinnatisectum</i>	1.5	575,316,904	619,463,363	28,631	49,365	339,432
SH7_18_3	<i>S. tarnii X S. tuberosum</i>	10.3	795,104,822	1,069,548,688	3,328	548,546	189,557
spl0147	<i>S. sparsipilum</i>	1.3	658,080,483	697,910,295	21,584	76,602	190,640
tcn8662	<i>S. tacnaense</i>	1.0	693,678,525	714,416,403	20,975	77,730	177,444
tcn8663	<i>S. tacnaense</i>	1.1	694,313,399	728,989,011	22,321	75,214	237,438
trp2833	<i>S. tarapatatum</i>	1.3	724,172,909	717,717,986	15,266	98,790	164,540

Table 3. *De novo* genome assembly statistics and the heterozygosity information.

Gaps were closed using *tsgapcloser*³⁵ with PacBio and Nanopore reads. A custom repeat library was constructed using *RepeatModeler*³⁶ and repeats were masked using *RepeatMasker*³⁷ and default parameters. RNA-Seq reads were filtered using *Kraken2*³⁸ and trimmed using *fastp*³⁰, then aligned to the assembly using *hisat2*³⁹. Structural annotations were identified using *Braker3*⁴⁰, and curated using *gfacs*⁴¹, and functional annotations were generated using *ahrd* (<https://github.com/groupschoof/AHRD>) and *interproscan*⁴².

Quality control and analysis of WGS and RNA-Seq for additional potato wild relatives. Quality control was performed on the raw sequencing data (WGS and RNA-Seq) using *FastQC v0.11.9*⁴³. An adaptor removal and read quality trimming was performed on the raw reads using *Trimmomatic v0.39*⁴⁴. Genome heterozygosity and estimated size were determined using the trimmed WGS reads. First, k-mer frequencies were calculated for each genome ($k = 21$) and a k-mer histogram was generated using *Jellyfish v2.3.0*⁴⁵. The k-mer histogram was used to generate various genome characteristics with *GenomeScope 2.0*⁴⁶. The polyploid genomes of *S. brevicaule* (= *oplocense*) *X S. tuberosum* (B1595106), *S. tarnii X S. tuberosum* (SH7_18_3), and *S. commersonii X S. andigena* (cxa70P5) hybrids showed higher rates of heterozygosity (Fig. 2), as seen in previous potato polyploid genomes²⁰, while the *S. etuberosum* (etb0161) genome showed very low levels of heterozygosity. The estimated haploid genome sizes ranged from 575 Mbp in *S. pinnatisectum* (S7) to 795 Mbp in *S. tarnii X S. tuberosum* (SH7_18_3).

***De novo* assembly and annotation.** Each genome was assembled into contigs from the raw Illumina reads using the *de novo* assembler *MaSuRca v4.0.5*⁴⁷. Organelle sequences were removed from the resulting assemblies using publicly available potato plastome^{48,49} and mitochondrial sequences^{27,50}. The *de novo* assemblies were aligned against the organellar reference using *nucmer* from *Mummer v4.0.0beta2* utilities³⁴ and the alignments were filtered for 95% identity and 200 bp alignment length. Any contigs with >95% coverage against the organelles were removed, as well as organellar sequences at the start and ends of contigs. The resulting filtered contigs were queried against the non-redundant nucleotide database using *blast+* v2.11.0⁵¹ to remove contaminant sequences. Any contig with a reliable match (90% query converge with 90% sequence identity) to organisms outside of green plants were removed from the assembly. Finally, contigs that are less than 200 bp in length were removed and the remaining contig ids were modified to create a final clean assembly file using *BBmap v38.86*⁵² (Table 3). The plastomes of each genome were assembled and annotated from raw Illumina reads using the *Plastaumatic*⁵³ pipeline. The quality of each *de novo* assembly was calculated using *QUAST v5.0.2*⁵⁴ to determine the assembly lengths, N50 values etc, and *BUSCO v5.2.2*⁵⁵ to check for the completeness of the assembly by looking for the presence of single copy orthologs from *Viridiplantae* (Fig. 3). An assembly was aligned against its reference genome²² when available using *nucmer* from *mummer v4.0.0beta2*³⁴ and the alignment statistics were generated using *dnadiff*.

Short ID	Species	Size (bp)	Bases masked (bp)	LTR elements (length occupied)	SINEs (length occupied)	LINEs (length occupied)	DNA TE (length occupied)	Small RNA
oka15	<i>S. okadae</i>	1,454,833,348	916,145,551	382,327,733	311,969	31,610,488	27,977,851	8,910,891
B1595106	<i>S. brevicaula</i> (= <i>S. oplocense</i>) <i>X S. tuberosum</i>	1,007,750,739	701,616,563	350,301,794	1,149,341	21,372,177	27,305,055	1,460,888
blv1353	<i>S. boliviense</i>	666,810,905	435,060,802	208,807,085	869,239	16,084,967	21,192,414	990,869
buk0368	<i>S. bukasovii</i>	716,989,807	480,158,888	235,333,428	848,893	17,276,972	20,948,217	1,046,511
chc0917	<i>S. chacoense</i>	765,361,945	513,233,181	243,991,092	954,158	17,069,642	23,354,232	1,076,340
chq2573	<i>S. chiquidenum</i>	693,400,501	466,348,914	223,602,249	768,785	16,231,519	20,580,712	940,668
cmm1080	<i>S. commersonii</i>	720,550,625	485,952,044	243,982,820	814,492	16,889,446	22,470,598	1,067,960
cxa70P5	<i>S. commersonii X S. andigena</i>	957,824,389	653,608,360	315,996,868	1,157,382	20,509,112	26,519,625	1,559,458
etb0161	<i>S. etuberosum</i>	610,576,211	407,532,810	215,984,646	695,999	13,262,589	11,258,462	1,037,743
grc1498	<i>S. gracilifrons</i>	669,865,430	444,824,908	216,390,612	826,044	15,725,904	20,096,452	1,035,345
ifd1359	<i>S. megistacrolobum</i>	684,052,840	444,642,488	211,138,196	890,857	16,205,334	20,757,278	1,033,001
lgl2830	<i>S. lignicaule</i>	662,942,901	435,332,101	210,672,488	851,994	16,142,336	20,366,503	935,951
mga1403	<i>S. megistacrolobum</i>	695,075,751	456,966,566	220,900,148	862,821	16,633,307	20,828,617	1,129,864
pcs2126	<i>S. paucissectum</i>	734,795,577	499,861,459	244,892,952	808,303	17,179,519	21,268,375	933,269
pur1868	<i>S. piurae</i>	720,192,443	491,426,301	242,004,208	781,453	16,664,457	20,933,611	921,302
S12	<i>S. tarijense</i>	711,715,600	470,513,861	226,540,751	884,250	16,607,481	22,555,724	1,044,403
S15	<i>S. gandarillasii</i>	686,228,213	451,836,344	215,877,305	865,194	16,770,897	21,598,771	981,756
S3	<i>S. commersonii</i>	660,741,360	434,370,170	214,113,765	783,717	16,270,941	21,546,608	923,547
S7	<i>S. pinnatisectum</i>	619,463,363	401,435,659	195,095,305	743,115	15,709,814	18,201,493	910,108
SH7_18_3	<i>S. tarnii X S. tuberosum</i>	1,069,548,688	738,444,227	363,846,598	1,246,886	23,381,848	29,029,059	1,560,804
sp10147	<i>S. sparsipilum</i>	697,909,412	457,896,220	218,926,860	883,367	16,834,545	22,095,633	992,710
tcn8662	<i>S. tacnaense</i>	714,416,403	474,138,047	230,699,916	882,481	17,044,624	21,513,973	1,203,905
tcn8663	<i>S. tacnaense</i>	728,989,011	489,162,880	242,451,878	871,253	17,339,593	21,361,744	1,024,620
trp2833	<i>S. tarapatanum</i>	717,717,986	478,288,796	235,624,996	866,639	16,996,387	21,120,093	969,990

Table 4. Results of the repeat analysis with different classes of TEs and their size.

The genome assemblies were annotated using the RNA-Seq data as the evidence for gene prediction along with homology-based prediction. The two accessions (B1595106 and SH7_18_3) with no RNA-Seq data were annotated only with the protein sequences. A repeat masking was performed on each genome assembly using *RepeatMasker v4.1.2-pl³⁷* with a custom repeat library of potato reference sequences constructed by concatenating repeat libraries from the *Petota* super pangenome²⁶ and DMv6.1⁹. Table 4 details the number of bases masked in each genome and size of different types of transposable elements found in them. Then the structural annotation was performed using *BRAKER v2.1.6^{40,56,57}* in two runs, one with the RNA-Seq data as evidence and another run with protein sequences as evidence. The trimmed RNA-Seq reads of each accession were aligned to its genome assembly using *HiSAT2 v2.2.1³⁹* for *BRAKER1⁵⁷*. Most of the accessions have an optimal RNA-Seq alignment rate (>65%) against their short read *de novo* assemblies. The alignment files were sorted and indexed using *SAMtools v1.16.1³¹*. Then the structural annotation was performed using the *BRAKER* pipeline^{58–63} on the masked genome assembly with the alignment file. The *BRAKER* pipeline performed gene prediction by *AUGUSTUS v3.4.0⁶⁴* and *GeneMark-ET^{65,66}* to generate gene structural annotations. For the homology-based prediction (*BRAKER2*), *Solanales* protein sequences from *OrthoDB* along with the protein sequences from potato reference genomes^{9,18,67} were combined to create a protein database for *BRAKER*. The *BRAKER1* and *BRAKER2* results were combined to select transcripts based on support by extrinsic evidence using *TSEBRA v1.0.3⁶⁸* to get a high-confidence gene set. The annotations for the accessions with no RNA-Seq data (B1595106 and SH7_18_3) were obtained from the *BRAKER2* run with protein sequences as evidence. Since no RNA-Seq evidence was available for these two, all the predicted structural genes were added to the final gene set. Figure 4 shows the number of genes and transcripts predicted in these accessions by the *BRAKER* pipeline. A total of 31,247–47,705 high confidence genes (i.e. genes with RNA-Seq evidence) were found in these genomes, with *S. etuberosum* having the lowest number of genes and the *S. commersonii X S. andigena* hybrid having the greatest number of genes. A high confidence gene set was not generated for B1595106 and SH7_18_3 genomes due to lack of the RNA-Seq data. A functional annotation was performed using *Interproscan v5.52-86⁴²* against the PFAM database. Approximately 65–70% of the annotated genes in all the genomes were found to have PFAM domains, except for the SH7_18_3, and B1595106 genomes (for which RNA-Seq data was not available).

Presence/absence variation analysis. The trimmed reads of each genome were aligned to the *Petota* super pangenome²⁶ using *bwa v0.7.17⁶⁹*. The resulting alignments were filtered to keep only properly paired (*-f2*) and remove secondary alignments (*-F 2048*), then sorted and indexed using *SAMtools v1.13³¹*. The gene presence/absence variations (PAVs) were identified in each genome using *SGSGeneLoss v0.1⁷⁰* with *minCov = 2* and *lost-Cutoff = 0.2*. These PAVs were used to generate a maximum-likelihood phylogenetic tree using *IQ-TREE v2.1.3⁷¹* with GTR2 + FO + R4 as a substitution model and 1000 bootstrap replicates and *S. etuberosum* set as an outgroup

BUSCO Assessment Results

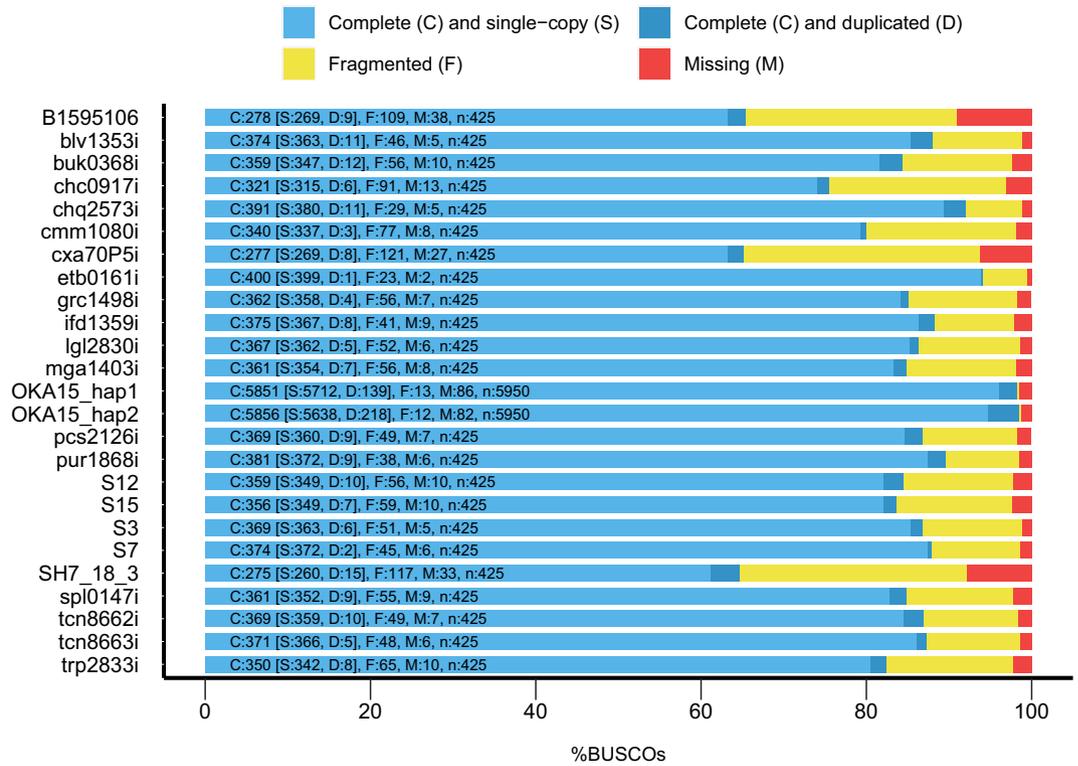


Fig. 3 BUSCO quality assessment results of short reads *de novo* genome assemblies for 24 potato wild relatives or hybrid clones (*Solanum* sp). The bar plot shows the percent of BUSCO genes (*Viridiplantae* dataset) present in each genome assembly. Species information and accession numbers are detailed in Table 1.

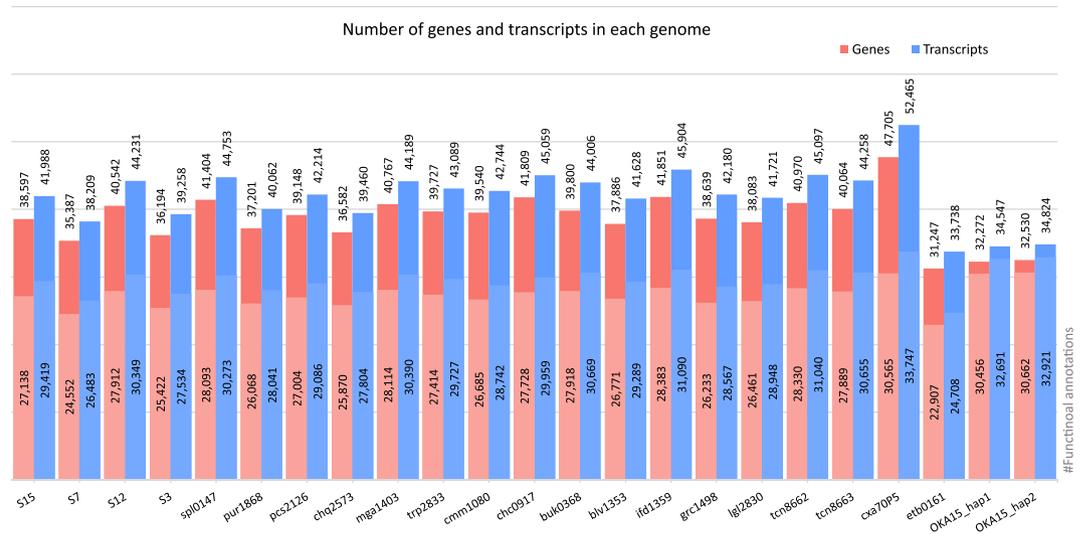


Fig. 4 The number of high-confidence genes and transcripts found in 22 potato wild relatives accessions (*Solanum* sp). The two columns for each genome represent the total number of predicted genes and transcripts (larger number on top of bar) along with their proportions of functional annotations (smaller number inside bar). Species information and accession numbers are detailed in Table 1.

(Fig. 5). A multiple sequence alignment of the 23 plastomes was made using *MAFFT* v7⁷², followed by a plastome based parsimonious phylogenetic tree constructed using *paup* v4.0a⁷³ with 1000 bootstrap replicates and *S. etuberosum* as an outgroup (Fig. 6). The phylogenetic trees were visualized in *FigTree* v1.4.4 (<http://tree.bio.ed.ac.uk/software/figtree/>).

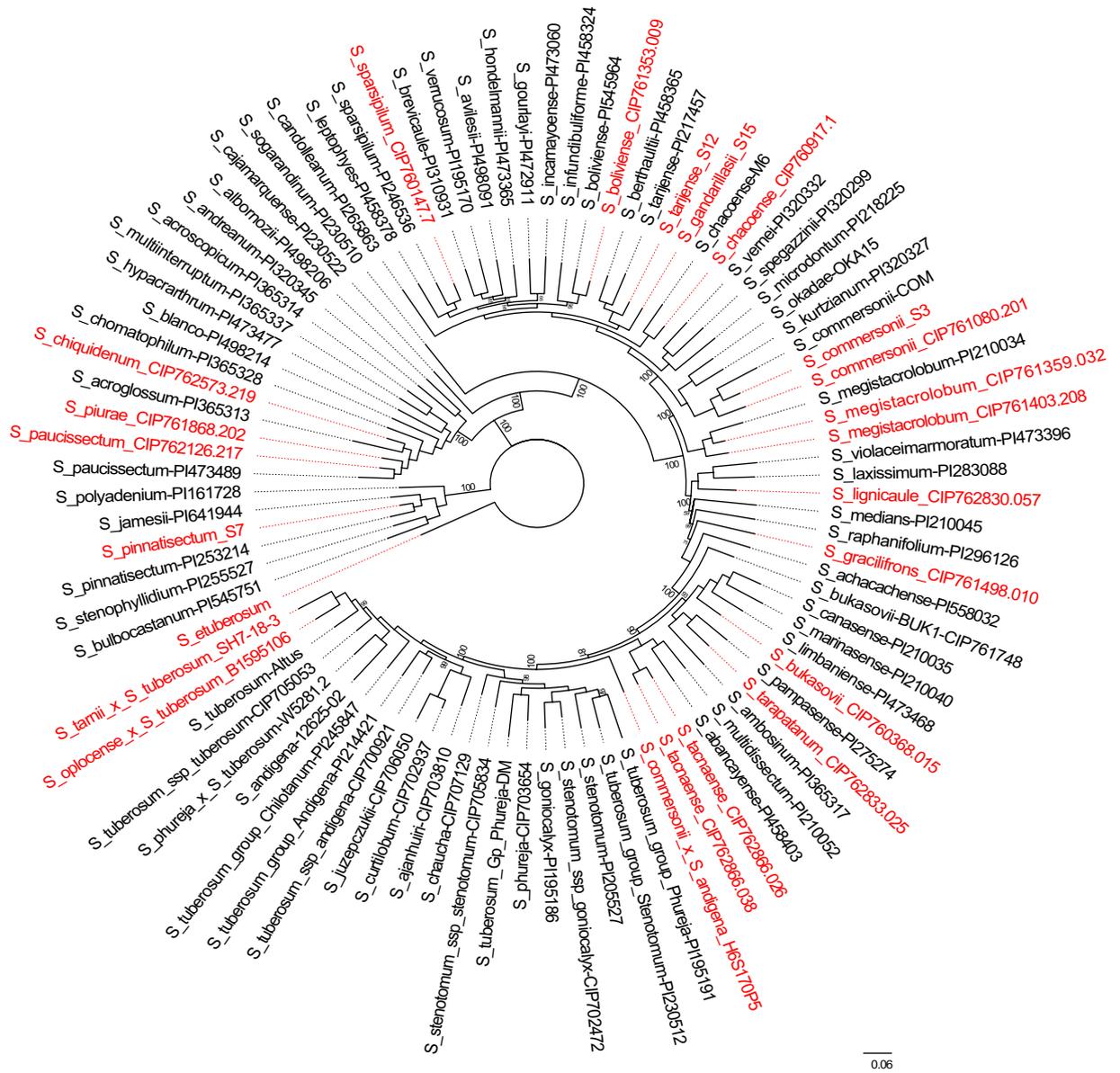


Fig. 5 A whole genome presence-absence variation (PAV) based phylogenetic tree showing the relatedness of potato wild relatives (*Solanum* sp.). The *S. etuberosum* genome sequence was used as an outgroup. The numbers at the nodes represent bootstrap support values. For clarity of the figure, only key values are given, and the ones omitted are all 100.

Data Records

S. okadae (OKA 15) genome sequence data is deposited at NCBI under BioProject PRJNA684565⁷⁴ and BioSample SAMN17860560⁷⁵; the Hi-C data under SRR26081972⁷⁶; Nanopore SRR20870051⁷⁷; PacBio SRR20870052⁷⁸; Illumina reads SRR14482384⁷⁹; and the phased haplotype genome assemblies are available under BioProject PRJNA1018115⁸⁰ (haplotype 1) JAWDCX000000000⁸¹ and PRJNA1018115⁸² (haplotype 2) JAWDCY000000000⁸³. The OKA 15 RNA-Seq data BioSamples are available under SAMN37429684- SAMN37429686⁸⁴⁻⁸⁶ and the RNA-Seq Illumina reads are available under SRR26082554-SRR26082556⁸⁷⁻⁸⁹.

All other data used in this study is deposited in NCBI and available under the BioProject PRJNA779368⁹⁰. The sample descriptions of the 23 genomes used in the genome sequencing are available under the BioSample accession numbers (SAMN23440977⁹¹, SAMN23440980 – SAMN23440983⁹²⁻⁹⁵, SAMN23440986 – SAMN23440990⁹⁶⁻¹⁰⁰, SAMN23440993 – SAMN23440998¹⁰¹⁻¹⁰⁶, SAMN23441000 – SAMN23441002¹⁰⁷⁻¹⁰⁹, SAMN23441004 – SAMN23441007¹¹⁰⁻¹¹³). The Illumina WGS reads of each genome were deposited with their respective Biosample IDs under the Sequence Read Archive (SRA) submission (SRR17078416 – SRR17078418¹¹⁴⁻¹¹⁶, SRR17078420¹¹⁷, SRR17078422 – SRR17078424¹¹⁸⁻¹²⁰, SRR17078426 – SRR17078432¹²¹⁻¹²⁷, SRR17078435 – SRR17078439¹²⁸⁻¹³², SRR17078441 – SRR17078444¹³³⁻¹³⁶). The final genome assemblies were deposited under the WGS assembly projects with the accession numbers JAJONR000000000

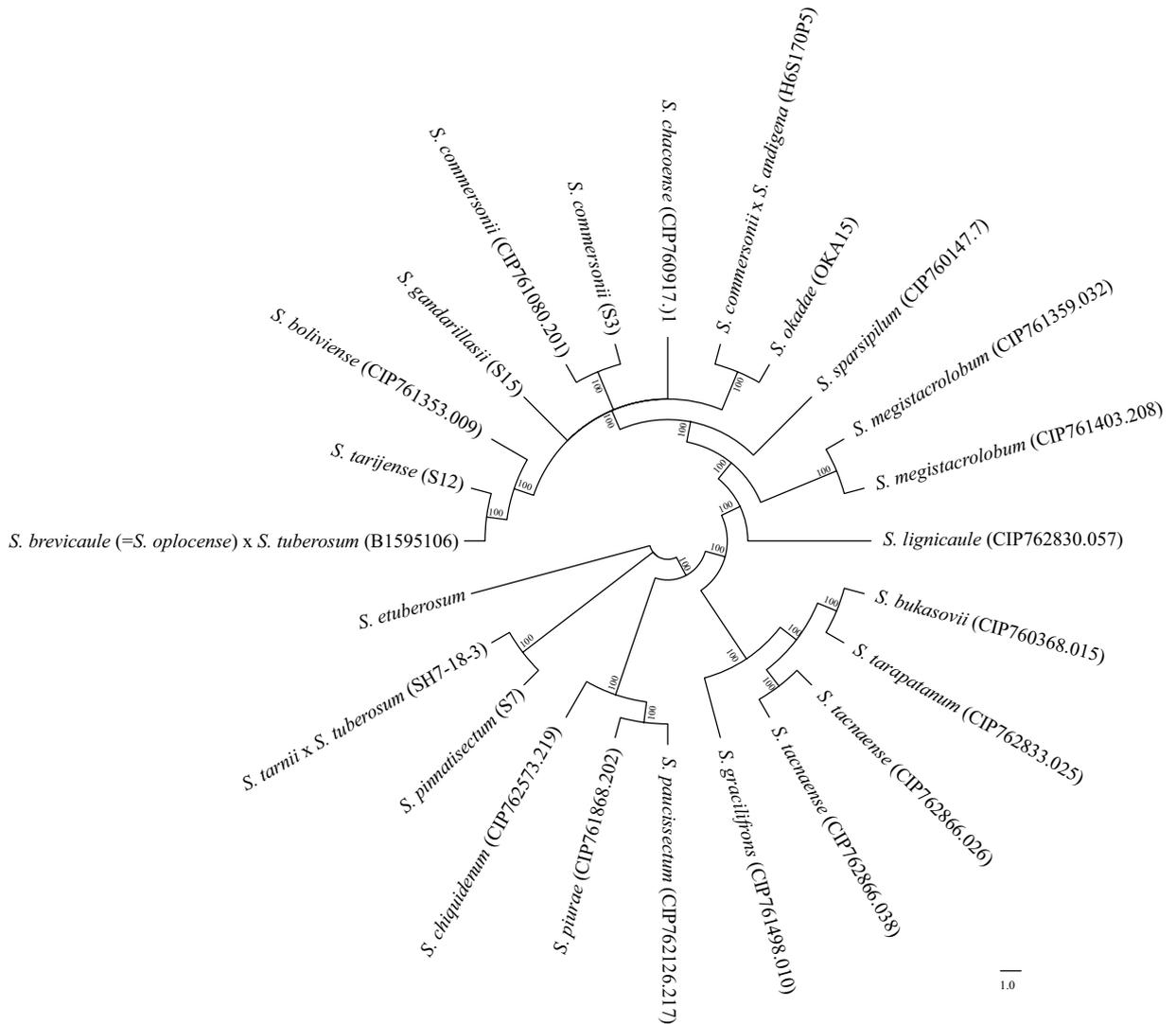


Fig. 6 A plastome sequence based phylogenetic tree of 24 potato wild relatives (*Solanum* sp.). The *S. etuberosum* was used as an outgroup. The numbers at the nodes represent bootstrap support values.

– JAJONU000000000^{137–139}, JAJONW000000000 – JAJOOD000000000^{140–148}, JAJOOF000000000
 – JAJOOJ000000000^{149–153}, JAJOOM000000000 – JAJOOP000000000^{154–157}, JAJOOS000000000¹⁵⁸,
 JAKRZT000000000¹⁵⁹ and the 23 plastome sequences were deposited in the GenBank database with the accession numbers OM638053-OM638057^{160–164}, OM638059-OM638064^{165–170}, OM638066-OM638068^{171–173},
 OM638071-OM638072^{174,175}, OM638074-OM638077^{176–179}, OM638081-OM638083^{180–182}. A complete list of WGS raw data and assemblies are available in Table 5. The sample descriptions of the RNA-Seq samples are available under the BioSample accession numbers SAMN32886741 – SAMN32886746^{183–188}, SAMN32886748 – SAMN32886751^{189–192}, SAMN32886756 – SAMN32886763^{193–200}, SAMN32886766 – SAMN32886771^{201–206}, SAMN32886782 – SAMN32886784^{207–209}, SAMN32886786²¹⁰, SAMN32886787²¹¹, SAMN32886790 – SAMN32886796^{212–218}, SAMN32886798²¹⁹, SAMN32886799²²⁰ and their raw sequencing reads are available under the SRA accession numbers SRR23225904 – SRR23225909^{221–226}, SRR23225913²²⁷, SRR23225914²²⁸, SRR23225916 – SRR23225918^{229–231}, SRR23225929²³², SRR23225930²³³, SRR23225933²³⁴, SRR23225938 – SRR23225940^{235–237}, SRR23225942 – SRR23225944^{238–240}, SRR23225947 – SRR23225951^{241–245}, SRR23225953 – SRR23225955^{246–248}, SRR23225960 – SRR23225962^{249–251}, SRR23225964²⁵², SRR23225966 – SRR23225971^{253–258} (Table 6). The gene annotations are available at the following link: <https://potatogenomeportal.org/download>

Technical Validation

Sequencing data quality control.

The quality of the raw WGS and RNA-Seq reads were checked using *FastQC*⁴³. Various metrics such as per base sequence quality, overall read quality score, GC content, N content, and overrepresented sequences were analyzed to detect low quality sequences or biases in the data. We have found no issues or bias in the WGS reads, whereas the RNA-Seq data is found to have overrepresented sequences in some of the datasets, which is not unusual for RNA-Seq data. We have also checked for viral genome sequences in the RNA-Seq data since potato plants are frequently found to have various types of potato viruses. The reads were

Short ID	Species	BioSample	Nuclear genome assembly	WGS reads	Plastome
oka14	<i>S. okadae</i>	SAMN17860560	JAWDCX000000000/ JAWDCY000000000	SRR26081972 (Hi-C) SRR20870051 (Nanopore) SRR20870052 (PacBio) SRR14482384 (10X)	MW592006
chq2573	<i>S. chiquidenum</i>	SAMN23440977	JAJOOS000000000	SRR17078430	OM638057
chc0917	<i>S. chacoense</i>	SAMN23440980	JAJOOP000000000	SRR17078444	OM638056
pur1868	<i>S. piurae</i>	SAMN23440981	JAJOOO000000000	SRR17078443	OM638068
mga1403	<i>S. megistacrolobum</i>	SAMN23440982	JAJOON000000000	SRR17078442	OM638066
cmm1080	<i>S. commersonii</i>	SAMN23440983	JAJOOM000000000	SRR17078441	OM638059
blv1353	<i>S. boliviense</i>	SAMN23440986	JAJOJ000000000	SRR17078439	OM638054
buk0368	<i>S. bukasovii</i>	SAMN23440987	JAJOOI000000000	SRR17078438	OM638055
grc1498	<i>S. gracilifrons</i>	SAMN23440988	JAJOOH000000000	SRR17078437	OM638062
ifd1359	<i>S. megistacrolobum</i>	SAMN23440989	JAJOOG000000000	SRR17078436	OM638063
lgl2830	<i>S. lignicaule</i>	SAMN23440990	JAJOOF000000000	SRR17078435	OM638064
tcn8663	<i>S. tacnaense</i>	SAMN23440993	JAJOOD000000000	SRR17078432	OM638082
tcn8662	<i>S. tacnaense</i>	SAMN23440994	JAJOOC000000000	SRR17078431	OM638081
trp2833	<i>S. tarapatanum</i>	SAMN23440995	JAJOOB000000000	SRR17078429	OM638083
cxax70P5	<i>S. commersonii</i> X <i>S. andigena</i>	SAMN23440996	JAJOOA000000000	SRR17078428	OM638060
pcs2126	<i>S. paucissectum</i>	SAMN23440997	JAJOZ000000000	SRR17078427	OM638067
spl0147	<i>S. sparsipilum</i>	SAMN23440998	JAKRZT000000000	SRR17078426	OM638077
S15	<i>S. gandarillasii</i>	SAMN23441000	JAJOY000000000	SRR17078424	OM638072
S7	<i>S. pinnatisectum</i>	SAMN23441001	JAJOX000000000	SRR17078423	OM638075
S12	<i>S. tarijense</i>	SAMN23441002	JAJOV000000000	SRR17078422	OM638071
S3	<i>S. commersonii</i>	SAMN23441004	JAJOV000000000	SRR17078420	OM638074
B1595106	<i>S. brevicaulis</i> (= <i>S. oplocense</i>) X <i>S. tuberosum</i>	SAMN23441005	JAJOV000000000	SRR17078418	OM638053
SH7-18-3	<i>S. tarnii</i> X <i>S. tuberosum</i>	SAMN23441006	JAJOV000000000	SRR17078417	OM638076
etb0161	<i>S. etuberosum</i>	SAMN23441007	JAJOV000000000	SRR17078416	OM638061

Table 5. A detailed list of genomic data descriptor records. It includes BioSample, SRA, genome, and plastome assembly accession numbers of each genome.

classified using *kraken2* v2.1.2³⁸ by searching against their own genome assemblies and potato virus sequences downloaded from NCBI. Most genomes have very few hits to the viral sequences, and we removed reads that have hits to viruses from each sample and used the clean reads in the following analyses.

Quality assessment of the genome assemblies. The *de novo* genome assembly quality was assessed by *QUAST*⁵⁴ and *BUSCO*⁵⁵ where contiguity and completeness of the assemblies were analyzed. The OKA15 genome assembly haplotype 1 and 2 have a complete BUSCO score of 98.3% and 98.5%, respectively, when compared to the (solanales_odb10). *Mercury*²⁵⁹ was used to estimate the completeness (QV). The complete OKA15 assembly was 99.4%, while haplotype 1 and haplotype 2 were 65 and 65% complete individually (meaning 15–30% of the reads map equally well to both haplotypes). Haplotype 1 (hap1) has 214 contigs with a total length of 729 Mb. The largest contig is 84.7 Mb and the N50 is 58.6 Mb. Haplotype 2 (hap2) has 131 contigs with a total length of 726 Mb, the largest contig is 83.2 Mb and the N50 58.4 Mb. The heterozygosity of the oka15 genome was calculated using *jellyfish*⁴⁵ and found to be 0.87%.

The *S. etuberosum* genome, which has been self-pollinated for several generations, has the best assembly of all with an N50 value of 59,223, 34,061 contigs, and 610 Mbp genome size. The three tetraploids (B1595106, SH7_18_3, and cxa70P5), and a diploid *S. chacoense* (chc0917) have fragmented assemblies (Table 2). Nine genome assemblies from this study were also compared with long-read assemblies of the same species from a recent study²². Overall, the amount of sequences that mapped to their individual reference ranged from ~85–99% with ~79–98% of the assembly coverage (Table 7). The BUSCO assessment of the assemblies revealed the majority of them are complete with presence of >80% complete genes (Fig. 2). Overall, the *S. etuberosum* genome has the highest percent of core plant orthologous genes with 94.1% of genes present as complete copies, followed by the chq762573 and pur1868 genomes with 92% and 89.6% of complete BUSCO genes. The SH7_18_3, cxa70P5, and B1595106 genomes have a higher percentage of fragmented genes among the 23 genomes. High rates of heterozygosity and higher ploidy levels are the major contributing factors to highly fragmented assemblies³⁸.

Phylogenetic inference. Previous studies classified section *Petota* (tuber-bearing) into major Clades and subgroups^{13,26,260,261}. Here we have reconstructed phylogenetic trees from the PAV data (Fig. 5) as well as plastome sequences (Fig. 6) to understand the relationship between these genomes. The results have shown similar groupings as previous studies where *S. pinnatisectum* (Clade 1 + 2), *S. paucissectum*, *S. piurae*, and *S. chiquidenum* (Clade 3), and the remaining accessions (Clade 4) separated into different clades. Within Clade 4, the *S. megistacrolobum* accessions formed a sister clade to Clade 4 South, something which was also seen in the *Petota*

Short ID	Species	SRA accession	BioSample accession	BioSample name
oka15	<i>S. okadae</i>	SRR26082556	SAMN37429684	OKA15_Leaf
oka15	<i>S. okadae</i>	SRR26082555	SAMN37429685	OKA15_sprout
oka15	<i>S. okadae</i>	SRR26082554	SAMN37429686	OKA15_tuber
chq2573	<i>S. chiquidenum</i>	SRR23225967	SAMN32886745	chq2573_Leaf
chq2573	<i>S. chiquidenum</i>	SRR23225966	SAMN32886746	chq2573_Shoot
chc0917	<i>S. chacoense</i>	SRR23225970	SAMN32886741	chc0917_Flower
chc0917	<i>S. chacoense</i>	SRR23225969	SAMN32886742	chc0917_Leaf
chc0917	<i>S. chacoense</i>	SRR23225968	SAMN32886743	chc0917_Shoot
chc0917	<i>S. chacoense</i>	SRR23225971	SAMN32886744	chc0917_Tuber
pur1868	<i>S. piurae</i>	SRR23225948	SAMN32886762	pur1868_Leaf
pur1868	<i>S. piurae</i>	SRR23225947	SAMN32886763	pur1868_Shoot
mga1403	<i>S. megistacrolobum</i>	SRR23225955	SAMN32886756	mga1403_Flower
mga1403	<i>S. megistacrolobum</i>	SRR23225954	SAMN32886757	mga1403_Leaf
mga1403	<i>S. megistacrolobum</i>	SRR23225953	SAMN32886758	mga1403_Shoot
cmm1080	<i>S. commersonii</i>	SRR23225964	SAMN32886748	cmm1080_Leaf
cmm1080	<i>S. commersonii</i>	SRR23225962	SAMN32886749	cmm1080_Shoot
cmm1080	<i>S. commersonii</i>	SRR23225961	SAMN32886750	cmm1080_Tuber
blv1353	<i>S. boliviense</i>	SRR23225907	SAMN32886792	blv1353_Leaf
buk0368	<i>S. bukasovii</i>	SRR23225906	SAMN32886793	buk0368_Leaf
grc1498	<i>S. gracilifrons</i>	SRR23225905	SAMN32886794	grc1498_Leaf
ifd1359	<i>S. megistacrolobum</i>	SRR23225904	SAMN32886795	ifd1359_Leaf
lgl2830	<i>S. lignicaule</i>	SRR23225933	SAMN32886796	lgl2830_Leaf
tcn8663	<i>S. tacnaense</i>	SRR23225930	SAMN32886798	tcn8663_Leaf
tcn8662	<i>S. tacnaense</i>	SRR23225929	SAMN32886799	tcn8663_Leaf
trp2833	<i>S. tarapatanum</i>	SRR23225918	SAMN32886782	trp2833_Flower
trp2833	<i>S. tarapatanum</i>	SRR23225917	SAMN32886783	trp2833_Leaf
trp2833	<i>S. tarapatanum</i>	SRR23225916	SAMN32886784	trp2833_Shoot
cx70P5	<i>S. commersonii</i> X <i>S. andigena</i>	SRR23225960	SAMN32886751	cx70P5_Shoot
pcs2126	<i>S. paucisectum</i>	SRR23225951	SAMN32886759	pcs2126_Leaf
pcs2126	<i>S. paucisectum</i>	SRR23225950	SAMN32886760	pcs2126_Shoot
spl0147	<i>S. sparsipilum</i>	SRR23225939	SAMN32886770	spl0147_Leaf
spl0147	<i>S. sparsipilum</i>	SRR23225938	SAMN32886771	spl0147_Shoot
S15	<i>S. gandarillasii</i>	SRR23225944	SAMN32886766	S15_Leaf
S15	<i>S. gandarillasii</i>	SRR23225943	SAMN32886767	S15_Tuber
S7	<i>S. pimatisectum</i>	SRR23225942	SAMN32886768	S7_Leaf
S7	<i>S. pimatisectum</i>	SRR23225940	SAMN32886769	S7_Tuber
S12	<i>S. tarijense</i>	SRR23225914	SAMN32886786	S12_Leaf
S12	<i>S. tarijense</i>	SRR23225913	SAMN32886787	S12_Tuber
S3	<i>S. commersonii</i>	SRR23225909	SAMN32886790	S3_Leaf
S3	<i>S. commersonii</i>	SRR23225908	SAMN32886791	S3_Tuber
etb0161	<i>S. etuberosum</i>	SRR23225949	SAMN32886761	PI245939_Leaf

Table 6. A detailed list of transcriptomic data descriptor records. It includes BioSample and their respective SRA accession numbers of individual samples.

pangenome²⁶. *S. okadae* placed in the Clade 4 South as also seen in the *Petota* pangenome. The *S. tarnii* X *S. tuberosum* (SH7_18_3) genome, which is a somatic hybrid between *S. tarnii* (Clade 1 + 2 species) and *S. tuberosum* (Clade 4) obtained by *in vitro* fusions of protoplasts, placed differently in the PAV and plastome phylogenetic trees. In the plastome phylogenetic tree, it is grouped with the Clade 1 + 2 species, whereas, in the PAV tree it is grouped with another *S. tuberosum* hybrid (Clade 4) showing the phylogenetic difficulty with hybrids and the importance of including both nuclear and organellar data. Moreover, genomes that are known to be closely related are grouped together, and genomes representing the same species grouped together, which validates the use of PAVs in determining phylogenetic relationships²⁶.

Usage Notes

The *Solanum* section *Petota* has over 100 cultivated and wild species. The potato crop wild relatives are an excellent source of genetic variation; however, not many sequencing efforts have been carried out to study these wild species. Here, we present the 24 phased chromosomes for *S. okadae*, and genome and transcriptome data for an additional 23 potato wild relative accessions. Though most of these are short-read assemblies, they are a

Short ID	Species	Reference	% Aligned Seqs	% Aligned Bases	Avg Identity	SNPs
etb0161	<i>S. tuberosum</i>	<i>S. tuberosum</i> (PG0019)	96.74	97.79	98.47	5,125,086
S7	<i>S. pinnatisectum</i>	<i>S. pinnatisectum</i> (PG1013)	95.06	96.37	98.71	4,034,460
pcs2126	<i>S. paucisectum</i>	<i>S. paucisectum</i> (PG3022)	84.83	78.76	91.98	16,608,649
pur1868	<i>S. piurae</i>	<i>S. piurae</i> (PG3023)	86.62	84.15	94.06	15,723,546
lg12830	<i>S. lignicaule</i>	<i>S. lignicaule</i> (PG4017)	96.82	97.17	99.19	2,848,239
chc0917	<i>S. chacoense</i>	<i>S. chacoense</i> (PG4042)	97.79	94.8	97.51	9,402,305
cmm1080	<i>S. commersonii</i>	<i>S. commersonii</i> (PG4049)	98.96	97.73	98.6	5,454,247
S3	<i>S. commersonii</i>	<i>S. commersonii</i> (PG4049)	98.39	95.18	96.97	10,870,964
blv1353	<i>S. boliviense</i>	<i>S. boliviense</i> (PG5076)	96.88	88.4	94.95	14,093,976

Table 7. Alignment statistics of nine genomes from this study compared against reference genomes.

useful tool for exploring genetic diversity that can provide insight in potato pre-breeding. The WGS reads can be used in various analyses such as determining SNPs, structural variations, and PAVs to understand the genetic diversity that exists within the section *Petota*. The RNA-Seq data of these wild species with important agronomic traits is especially crucial to understand their transcriptome profile. *S. okadae* foliage has been shown to contain the glycoalkaloid tomatine, and undetectable levels of solanine and chaconine, and carries Colorado potato beetle resistance, drought and cold tolerance (but not frost) and moderate late blight resistance^{262–265}.

Code availability

All software with their specific version used for data processing are clearly described in the methods section. If no specific variable or parameters are mentioned for a software, the default parameters were used.

Received: 24 February 2023; Accepted: 23 April 2024;

Published online: 04 May 2024

References

- Devaux, A., Kromann, P. & Ortiz, O. Potatoes for sustainable global food security. *Potato Res.* **57**, 185–199 (2014).
- Hancock, R. D. *et al.* Physiological, biochemical and molecular responses of the potato (*Solanum tuberosum* L.) plant to moderately elevated temperature. *Plant Cell Environ.* **37**, 439–450 (2014).
- Fumia, N. *et al.* Wild relatives of potato may bolster its adaptation to new niches under future climate scenarios. *Food Energy Secur.* <https://doi.org/10.1002/fes3.360> (2022).
- Jansky, S. H. *et al.* A case for crop wild relative preservation and use in potato. *Crop Sci.* **53**, 746–754 (2013).
- Bradshaw, J. E., Bryan, G. J. & Ramsay, G. Genetic resources (including wild and cultivated *Solanum* species) and progress in their utilisation in potato breeding. *Potato Research* **49**, 49–65 (2006).
- Bellumori, M. *et al.* A study on the biodiversity of pigmented Andean potatoes: Nutritional profile and phenolic composition. *Molecules* **25**, (2020).
- Giusti, M. M., Polit, M. F., Ayvaz, H., Tay, D. & Manrique, I. Characterization and quantitation of anthocyanins and other phenolics in native Andean potatoes. *J. Agric. Food Chem.* **62**, 4408–4416 (2014).
- Potato Genome Sequencing Consortium *et al.* Genome sequence and analysis of the tuber crop potato. *Nature* **475**, 189–195 (2011).
- Pham, G. M. *et al.* Construction of a chromosome-scale long-read reference genome assembly for potato. *Gigascience* **9**, (2020).
- Yang, X. *et al.* The gap-free potato genome assembly reveals large tandem gene clusters of agronomical importance in highly repeated genomic regions. *Mol. Plant* **16**, 314–317 (2023).
- Hardigan, M. A. *et al.* Genome reduction uncovers a large dispensable genome and adaptive role for copy number variation in asexually propagated *Solanum tuberosum*. *Plant Cell* **28**, 388–405 (2016).
- Hardigan, M. A. *et al.* Genome diversity of tuber-bearing *Solanum* uncovers complex evolutionary history and targets of domestication in the cultivated potato. *Proc. Natl. Acad. Sci. USA* **114**, E9999–E10008 (2017).
- Li, Y. *et al.* Genomic analyses yield markers for identifying agronomically important genes in Potato. *Mol. Plant* **11**, 473–484 (2018).
- Achakkagari, S. R. *et al.* Genome sequencing of adapted diploid potato clones. *Front. Plant Sci.* **13**, 954933 (2022).
- Bao, Z. *et al.* Genome architecture and tetrasomic inheritance of autotetraploid potato. *Mol. Plant* **15**, 1211–1226 (2022).
- Sun, H. *et al.* Chromosome-scale and haplotype-resolved genome assembly of a tetraploid potato cultivar. *Nat. Genet.* **54**, 342–348 (2022).
- Hoopes, G. *et al.* Phased, chromosome-scale genome assemblies of tetraploid potato reveal a complex genome, transcriptome, and predicted proteome landscape underpinning genetic diversity. *Mol. Plant* **15**, 520–536 (2022).
- Zhou, Q. *et al.* Haplotype-resolved genome analyses of a heterozygous diploid potato. *Nat. Genet.* **52**, 1018–1023 (2020).
- van Lieshout, N. *et al.* *Solyntus*, the new highly contiguous reference genome for potato (*Solanum tuberosum*). *G3* **10**, 3489–3495 (2020).
- Kyriakidou, M. *et al.* Genome assembly of six polyploid potato genomes. *Sci Data* **7**, 88 (2020).
- Ovchinnikova, A. *et al.* Taxonomy of cultivated potatoes (*Solanum* section *Petota*: Solanaceae). *Bot. J. Linn. Soc.* **165**, 107–155 (2011).
- Tang, D. *et al.* Genome evolution and diversity of wild and cultivated potatoes. *Nature* **606**, 535–541 (2022).
- Kyriakidou, M. *et al.* Structural genome analysis in cultivated potato taxa. *Theor. Appl. Genet.* **133**, 951–966 (2020).
- Leisner, C. P. *et al.* Genome sequence of M6, a diploid inbred clone of the high-glycoalkaloid-producing tuber-bearing potato species *Solanum chacoense*, reveals residual heterozygosity. *Plant J.* **94**, 562–570 (2018).
- Aversano, R. *et al.* The *Solanum commersonii* genome sequence provides insights into adaptation to stress conditions and genome evolution of wild potato relatives. *Plant Cell* **27**, 954–968 (2015).
- Bozan, I. *et al.* Pangenome analyses reveal impact of transposable elements and ploidy on the evolution of potato species. *Proc. Natl. Acad. Sci. USA* **120**, e2211117120 (2023).

27. Achakkagari, S. R. *et al.* The complete mitogenome assemblies of 10 diploid potato clones reveal recombination and overlapping variants. *DNA Res.* **28**, (2021).
28. Hamilton, J. P. *et al.* Chromosome-scale genome assembly of the ‘Munstead’ cultivar of *Lavandula angustifolia*. *BMC Genomic Data* **24**, 75 (2023).
29. Cheng, H. *et al.* Haplotype-resolved de novo assembly using phased assembly graphs with hifiasm. *Nat Methods* **18**, 170–175 (2021).
30. Chen, S. *et al.* *fastp*: an ultra-fast all-in-one FASTQ preprocessor. *Bioinformatics* **34**, i884–i890 (2018).
31. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).
32. Zhou, C. *et al.* YaHS: yet another Hi-C scaffolding tool. *Bioinformatics* **39**, btac808 (2023).
33. Durand, N. *et al.* *Juicer* provides a one-click system for analyzing loop-resolution Hi-C experiments. *Cell Systems* **3**, (2016).
34. Marçais, G. *et al.* MUMmer4: A fast and versatile genome alignment system. *PLoS Comput. Biol.* **14**, e1005944 (2018).
35. Xu, M. *et al.* TGS-GapCloser: A fast and accurate gap closer for large genomes with low coverage of error-prone long reads. *GigaScience* **9**, giaa094 (2020).
36. Flynn, J. M. *et al.* “RepeatModeler2 for Automated Genomic Discovery of Transposable Element Families.”. *Proc. Natl. Acad. Sci. USA* **117**, 9451–57 (2020).
37. Smit, A. *et al.* RepeatMasker Open-4.0. <http://www.repeatmasker.org> (2015).
38. Wood, D. E. *et al.* Improved metagenomic analysis with Kraken 2. *Genome Biol.* **20**, 257 (2019).
39. Kim, D. *et al.* Graph-based genome alignment and genotyping with HISAT2 and HISAT-genotype. *Nat. Biotechnol.* **37**, 907–915 (2019).
40. Brůna, T. *et al.* BRAKER2: Automatic Eukaryotic Genome Annotation with GeneMark-EP+ and AUGUSTUS Supported by a Protein Database. *NAR Genomics & Bioinformatics* **3** (2021).
41. Caballero, M. & Wegrzyn, J. gFACs: Gene Filtering, Analysis, and Conversion to Unify Genome Annotations Across Alignment and Gene Prediction Frameworks. *Genomics, Proteomics & Bioinformatics* **17**, 305–10 (2019).
42. Quevillon, E. *et al.* InterProScan: Protein Domains Identifier. *Nucleic Acids Res.* **33**, W116–20 (2005).
43. Andrews, S. FastQC: a quality control tool for high throughput sequence data. Version 0.11.9. Babraham Bioinformatics, Babraham Institute, UK. <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/> (2019).
44. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics* **30**, 2114–2120 (2014).
45. Marçais, G. & Kingsford, C. A fast, lock-free approach for efficient parallel counting of occurrences of k-mers. *Bioinformatics* **27**, 764–770 (2011).
46. Ranallo-Benavidez, T. R., Jaron, K. S. & Schatz, M. C. GenomeScope 2.0 and Smudgeplot for reference-free profiling of polyploid genomes. *Nat. Commun.* **11**, 1432 (2020).
47. Zimin, A. V. *et al.* The MaSuRCA genome assembler. *Bioinformatics* **29**, 2669–2677 (2013).
48. Achakkagari, S. R. *et al.* Complete plastome assemblies from a panel of 13 diverse potato taxa. *PLoS One* **15**, e0240124 (2020).
49. Achakkagari, S. R., Tai, H. H., Davidson, C., Jong, H. D. & Strömvik, M. V. The complete plastome sequences of nine diploid potato clones. *Mitochondrial DNA B Resour* **6**, 811–813 (2021).
50. Achakkagari, S. R. *et al.* Complete mitogenome assemblies from a panel of 13 diverse potato taxa. *Mitochondrial DNA B Resour* **6**, 894–897 (2021).
51. Camacho, C. *et al.* BLAST+: architecture and applications. *BMC Bioinformatics* **10**, 421 (2009).
52. Bushnell, B. BMap: A fast, accurate, splice-aware aligner. <https://www.osti.gov/biblio/1241166> (2014).
53. Chen, W., Achakkagari, S. R. & Strömvik, M. Plastaumatic: Automating plastome assembly and annotation. *Front. Plant Sci.* **13**, 1011948 (2022).
54. Gurevich, A., Saveliev, V., Vyahhi, N. & Tesler, G. QUASt: quality assessment tool for genome assemblies. *Bioinformatics* **29**, 1072–1075 (2013).
55. Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V. & Zdobnov, E. M. BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics* **31**, 3210–3212 (2015).
56. Hoff, K. J., Lomsadze, A., Borodovsky, M. & Stanke, M. Whole-Genome Annotation with BRAKER. *Methods Mol. Biol.* **1962**, 65–95 (2019).
57. Hoff, K. J., Lange, S., Lomsadze, A., Borodovsky, M. & Stanke, M. BRAKER1: Unsupervised RNA-Seq-Based Genome Annotation with GeneMark-ET and AUGUSTUS. *Bioinformatics* **32**, 767–769 (2016).
58. Barnett, D. W., Garrison, E. K., Quinlan, A. R., Strömberg, M. P. & Marth, G. T. BamTools: a C++ API and toolkit for analyzing and managing BAM files. *Bioinformatics* **27**, 1691–1692 (2011).
59. Lomsadze, A., Ter-Hovhannissyan, V., Chernoff, Y. O. & Borodovsky, M. Gene identification in novel eukaryotic genomes by self-training algorithm. *Nucleic Acids Res.* **33**, 6494–6506 (2005).
60. Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using DIAMOND. *Nat. Methods* **12**, 59–60 (2014).
61. Stanke, M., Diekhans, M., Baertsch, R. & Haussler, D. Using native and syntenically mapped cDNA alignments to improve *de novo* gene finding. *Bioinformatics* **24**, 637–644 (2008).
62. Iwata, H. & Gotoh, O. Benchmarking spliced alignment programs including Spaln2, an extended version of Spaln that incorporates additional species-specific features. *Nucleic Acids Res.* **40**, e161 (2012).
63. Gotoh, O. A space-efficient and accurate method for mapping and aligning cDNA sequences onto genomic sequence. *Nucleic Acids Res.* **36**, 2630–2638 (2008).
64. Stanke, M., Schöffmann, O., Morgenstern, B. & Waack, S. Gene prediction in eukaryotes with a generalized hidden Markov model that uses hints from external sources. *BMC Bioinformatics* **7**, 62 (2006).
65. Lomsadze, A., Burns, P. D. & Borodovsky, M. Integration of mapped RNA-Seq reads into automatic training of eukaryotic gene finding algorithm. *Nucleic Acids Res.* **42**, e119 (2014).
66. Brůna, T., Lomsadze, A. & Borodovsky, M. GeneMark-EP+: eukaryotic gene prediction with self-training in the space of genes and proteins. *NAR Genom Bioinform* **2**, (2020).
67. Petek, M. *et al.* Cultivar-specific transcriptome and pan-transcriptome reconstruction of tetraploid potato. *Scientific Data* **7**, 1–15 (2020).
68. Gabriel, L., Hoff, K. J., Brůna, T., Borodovsky, M. & Stanke, M. TSEBRA: transcript selector for BRAKER. *BMC Bioinformatics* **22**, 566 (2021).
69. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv [q-bio.GN]* (2013).
70. Golicz, A. A. *et al.* Gene loss in the fungal canola pathogen *Leptosphaeria maculans*. *Funct. Integr. Genomics* **15**, 189–196 (2015).
71. Minh, B. Q. *et al.* IQ-TREE 2: New models and efficient methods for phylogenetic inference in the genomic era. *Mol. Biol. Evol.* **37**, 1530–1534 (2020).
72. Katoh, K., Rozewicki, J. & Yamada, K. D. MAFFT online service: multiple sequence alignment, interactive sequence choice and visualization. *Brief. Bioinform.* **20**, 1160–1166 (2019).
73. Swofford, D. L. PAUP: phylogenetic analysis using parsimony (and other methods), 4.0 beta. <http://paup.csit.fsu.edu/>.
74. NCBI BioProject <https://identifiers.org/bioproject:PRJNA684565> (2020).
75. NCBI BioSample <https://identifiers.org/biosample:SAMN17860560> (2021).
76. NCBI Sequence Read Archive <https://identifiers.org/insdc.sra:SRR26081972> (2023).

239. NCBI Sequence Read Archive <https://identifiers.org/insdc.sra:SRR23225943> (2023).
240. NCBI Sequence Read Archive <https://identifiers.org/insdc.sra:SRR23225944> (2023).
241. NCBI Sequence Read Archive <https://identifiers.org/insdc.sra:SRR23225947> (2023).
242. NCBI Sequence Read Archive <https://identifiers.org/insdc.sra:SRR23225948> (2023).
243. NCBI Sequence Read Archive <https://identifiers.org/insdc.sra:SRR23225949> (2023).
244. NCBI Sequence Read Archive <https://identifiers.org/insdc.sra:SRR23225950> (2023).
245. NCBI Sequence Read Archive <https://identifiers.org/insdc.sra:SRR23225951> (2023).
246. NCBI Sequence Read Archive <https://identifiers.org/insdc.sra:SRR23225953> (2023).
247. NCBI Sequence Read Archive <https://identifiers.org/insdc.sra:SRR23225954> (2023).
248. NCBI Sequence Read Archive <https://identifiers.org/insdc.sra:SRR23225955> (2023).
249. NCBI Sequence Read Archive <https://identifiers.org/insdc.sra:SRR23225960> (2023).
250. NCBI Sequence Read Archive <https://identifiers.org/insdc.sra:SRR23225961> (2023).
251. NCBI Sequence Read Archive <https://identifiers.org/insdc.sra:SRR23225962> (2023).
252. NCBI Sequence Read Archive <https://identifiers.org/insdc.sra:SRR23225964> (2023).
253. NCBI Sequence Read Archive <https://identifiers.org/insdc.sra:SRR23225966> (2023).
254. NCBI Sequence Read Archive <https://identifiers.org/insdc.sra:SRR23225967> (2023).
255. NCBI Sequence Read Archive <https://identifiers.org/insdc.sra:SRR23225968> (2023).
256. NCBI Sequence Read Archive <https://identifiers.org/insdc.sra:SRR23225969> (2023).
257. NCBI Sequence Read Archive <https://identifiers.org/insdc.sra:SRR23225970> (2023).
258. NCBI Sequence Read Archive <https://identifiers.org/insdc.sra:SRR23225971> (2023).
259. Rhie, A. *et al.* Merqury: reference-free quality, completeness, and phasing assessment for genome assemblies. *Genome Biol.* **21**, 245 (2020).
260. Spooner, D. M., Ghislain, M., Simon, R., Jansky, S. H. & Gavrilenko, T. Systematics, diversity, genetics, and evolution of wild and cultivated potatoes. *Bot. Rev.* **80**, 283–383 (2014).
261. Huang, B., Ruess, H., Liang, Q., Colleoni, C. & Spooner, D. M. Analyses of 202 plastid genomes elucidate the phylogeny of *Solanum* section *Petota*. *Sci. Rep.* **9**, 4454 (2019).
262. Karki, H. S., Jansky, S. H. & Halterman, D. A. Screening of Wild Potatoes Identifies New Sources of Late Blight Resistance. *Plant disease* **105**, 368–376 (2020).
263. Pelletier, Y., Clark, C. & Tai, G. C. Resistance of three wild tuber-bearing potatoes to the Colorado potato beetle. *Entomologia Experimentalis et Applicata* **100**, 31–41 (2001).
264. Watanabe, K. N., Kikuchi, A., Shimazaki, T. & Asahina, M. Salt and drought stress tolerances in transgenic potatoes and wild species. *Potato Res.* **54**, 319–324 (2011). (2011).
265. Vega, S. E. & Bamberg, J. B. Screening the US potato collection for frost hardiness. *American Pot. J.* **72**, 13–21 (1995).

Acknowledgements

The authors would like to acknowledge the potato breeding and genebank teams at Agriculture and Agri-Food Canada (Sylvia Steeves, Emily McCoy, Martin Lagüe); USDA, Aberdeen, ID, USA; and at The International Potato Center (CIP), Lima, Peru for maintenance and propagation of clones used in the study; Thiago Mendes, Mariela Aponte and Edith Poma (CIP) for providing plant material; Daniel Philippe Matton (Université de Montreal, Canada) and Ramona Thieme (Julius Kühn-Institut, Germany) for germplasm transferred to the Canadian Potato Gene Resources Genebank. Funding for the study was provided by a grant from Agriculture and Agri-Food Canada (International Collaboration Program) “Genome sequencing of wild *Solanum* diploids” to CIP; a grant from the Agriculture and Agri-Food Canada Genomics Research and Development Initiative to H.H.T.; an AAFC-Génome Québec joint call Management Driven Genomics award “Revolutionizing potato variety development for climate smart potato” to HHT and MVS (J-002367/GQ-AAC-2019-2); a Compute Canada RPP award (Research Portals and Platforms: “The Potato Genome Diversity Portal”) and a Compute Canada Resources for Research Groups award (“Structural variation analyses of complex plant genomes in search of climate smart adaptations”) to MVS; and a 10-year project to strengthen food and nutrition security worldwide by supporting the conservation and use of crop diversity supported by the Government of Norway (NORAD) (“Biodiversity for Opportunities, Livelihoods and Development (BOLD)”) (Project CONT-0791 “CWR-derived potatoes integrated in breeding pipelines for climate change resilience of farming communities of Ecuador, Kenya and Peru” to CIP), managed by the Crop Trust and implemented in partnership with national and international genebanks and seed system actors and researchers around the world.

Author contributions

H.H.T., N.L.A., N.M.C., H.L.K., B.B., M.V.S. - designed the overall project and supervised experiments. S.R.A., I.B., J.C.C.T., L.P., J.S., H.J.M., H.H.T., N.L.A., N.M.C., H.L.K., B.B., - performed experiments and/or provided germplasm and data. H.J.M. prepared the *S. okadae* materials and genome sequencing libraries. S.R.A., I.B., J.C.C.T., H.H.T. analyzed data. S.R.A., J.C.C.T., H.H.T., N.L.A., H.L.K., N.M.C., M.V.S. prepared the manuscript. All authors have read and approved the final version of the manuscript.

Competing interests

The authors declare no competing interests.

Additional information

Correspondence and requests for materials should be addressed to M.V.S.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher’s note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2024