## Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)

1    **Authors:**

2    Qihua Liang, María Muñoz-Amatriaín, Shengqiang Shu, Sassoum Lo, Xinyi Wu, Joseph W.
3    Carlson, Patrick Davidson, David M. Goodstein, Jeremy Phillips, Nadia M. Janis, Elaine J. Lee,
4    Chenxi Liang, Peter L. Morrell, Andrew D. Farmer, Pei Xu, Timothy J. Close, Stefano Lonardi

5    **Author contributions:** Biological materials (SaL, XW, PX, MMA, TJC), Assembly (StL), SNP
6    and structural variant calling (QL, StL), Gene Annotation (SS), CowpeaPan (JWC, PD, DMG, JP),
7    GO enrichment analysis (MMA, ADF, TJC), VeP annotation (NMJ, EJL, CL, PLM), gene size
8    comparisons (NMJ, EJL, CL, PLM, ADF, TJC), synteny view (ADF), manuscript preparation
9    (TJC, StL, MMA, QL, PLM, AF)

10   **Keywords:** Cowpea, *Vigna unguiculata*, pan-genome, genetic variation, SNPs, presence/absence
11   variants (PAV), inversions, IT97K-499-35, Suvita-2, Sanzi, CB5-2, UCR779, TZ30, ZN016

12   **Abstract:**

13   Cowpea, *Vigna unguiculata L.* Walp., is a diploid warm-season legume of critical importance as
14   both food and fodder in sub-Saharan Africa. This species is also grown in Northern Africa, Europe,
15   Latin America, North America, and East to Southeast Asia. To capture the genomic diversity of
16   domesticates of this important legume, *de novo* genome assemblies were produced for
17   representatives of six sub-populations of cultivated cowpea identified previously from genotyping
18   of several hundred diverse accessions. In the most complete assembly (IT97K-499-35), 26,026
19   core and 4,963 noncore genes were identified, with 35,436 pan genes when considering all seven
20   accessions. GO-terms associated with response to stress and defense response were highly
21   enriched among the noncore genes, while core genes were enriched in terms related to transcription
22   factor activity, and transport and metabolic processes. Over 5 million SNPs relative to each
23   assembly and over 40 structural variants >1 Mb in size were identified by comparing genomes.
24   Vu10 was the chromosome with the highest frequency of SNPs, and Vu04 had the most structural
25   variants. Noncore genes harbor a larger proportion of potentially disruptive variants than core
26   genes, including missense, stop gain, and frameshift mutations; this suggests that noncore genes
27   substantially contribute to diversity within domesticated cowpea.

28

**Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)**

1    **Article Summary (80 words maximum)**

2    This study reports annotated genome assemblies of six cowpea accessions. Together with the
3    previously reported annotated genome of IT97K-499-35, these constitute a pan-genome resource
4    representing six subpopulations of domesticated cowpea. Annotations include genes, variant calls
5    for SNPs and short indels, larger presence or absence variants, and inversions. Noncore genes are
6    enriched for loci involved in stress response and harbor many genic variants with potential effects
7    on coding sequence.

8    **Introduction:**

9    Individuals within a species vary in their genomic composition. The genome of any individual
10   does not include the full complement of genes contained within the species. A pan-genome
11   includes genes core to the species (shared among all individuals) and those absent from one or
12   more individuals (noncore, dispensable, or variable genes). This pan-genome concept started to be
13   applied to plants by Morgante *et al*. (2007) but began in bacterial species (reviewed by Golicz *et*
14   *al.,* 2020). Due to the complexity of plant genomes, the first studies exploring gene presence-
15   absence variation (PAV) in plants used reduced-representation approaches, including array
16   comparative genomic hybridization (CGH) and sequencing of transcriptomes (e.g., Springer *et al.*
17   2009, Muñoz-Amatriaín *et al.* 2013; Hirsch *et al.* 2014). Once sequencing of multiple plant
18   genomes became feasible, several pan-genomes of variable degrees of completeness were
19   generated, and it was soon understood that PAV is prevalent in plants and that the pan-genome of
20   any plant species is larger than the genome of any individual accession (reviewed by Lei *et al.*
21   2021). Moreover, many of the genes absent in reference accessions have functions of potential
22   adaptive or agronomic importance, such as time to flowering, and response to abiotic and biotic
23   stresses (Gordon *et al.* 2017; Montenegro *et al.* 2017; Bayer *et al.* 2020), making the construction
24   of a pan-genome a crucial task for crops of global importance.

25   Cowpea is a diploid ($2n = 22$) member of the family Fabaceae tribe Phaseoleae, closely related to
26   mung bean, common bean, soybean, and several other warm-season legumes. Cowpea was
27   domesticated in Africa, but its cultivation has spread throughout most of the globe (Herniter *et al*.,
28   2020). The inherent resilience of the species to drought and high temperatures (Hall 2004),
29   together with its nutritional value as a reliable source of plant-based protein and folic acid, position
30   cowpea favorably as a component of sustainable agriculture in the context of global climate
31   change. Most cowpea production and consumption presently occur in sub-Saharan Africa,
32   especially in the Sudano-Sahelian Zone, with production mainly by smallholder farmers, often as
33   an intercrop with maize, sorghum, or millet (Boukar *et al*., 2019). Tender green seeds are often
34   consumed during the growing season, and immature pods are eaten as a vegetable, especially in

2

**Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)**

1    East and Southeast Asia. In addition, fresh leaves are sometimes consumed, and dry haulms are
2    harvested and sold as fodder for livestock. Spreading varieties are also utilized as cover crops to
3    prevent soil erosion and weed control.

4    A reference genome sequence of cowpea cv. IT97K-499-35 was previously generated (Lonardi *et*
5    *al*., 2019). Preliminary sequence comparisons using whole genome shotgun (WGS) data of 36
6    accessions suggested that extensive SNP and structural variation exists within domesticated
7    cowpea (Lonardi *et al*., 2019). Cowpea also displays a wide range of phenotypic variation, and
8    genetic assignment approaches have identified six subpopulations within cultivated cowpea
9    germplasm (Muñoz-Amatriaín *et al*., 2021). These observations support the need to develop
10    cowpea pan-genome resources based on diverse cowpea accessions.

11    This study reports *de novo* assemblies of six cultivated cowpea accessions. Each accession was
12    annotated using transcriptome sequences from the accession along with *ab initio* methods. These
13    genome sequences, together with the previously reported sequence of IT97K-499-45 (Lonardi *et*
14    *al*., 2019), constitute a pan-genome resource for domesticated cowpea. Using annotations for the
15    seven genomes, including genes, along with variant calls for SNPs and short indels, and larger
16    structural variants, the following questions were addressed: (i) What proportion of genes are core
17    and noncore, and do core and noncore genes differ in size or functional class? (ii) What proportion
18    of large-effect variants are created by single nucleotide variants versus structural variants
19    (including indels), and do the proportions of large-effect variants differ among core and noncore
20    genes? (iii) To what extent are gene content and gene order consistent across accessions within the
21    species *V. unguiculata* and across species within the genus *Vigna* and the tribe Phaseoleae? The
22    results suggest that both extensive structure differences among individual accessions and the
23    nature of variation in noncore genes are important considerations in efforts to identify genetic
24    variation with adaptive potential.

25    **Materials & Methods:**

26    **Cowpea accessions selected for sequencing (Supplemental Table S01)**

27    Accessions chosen for sequencing and *de novo* assembly represented the six subpopulations of
28    domesticated cowpea described in Muñoz-Amatriaín et al. (2021), as indicated in Figure 1. The
29    intention of choosing accessions that cover each subpopulation was to maximize the discovery of
30    genetic variations relevant to cultivated cowpea using a small number of samples. As shown by
31    Gordon *et al.* (2017) in *Brachypodium distachyon*, the addition of individuals from subpopulations
32    not previously sampled contributes much more to increasing the pan-genome size than adding
33    closely related individuals.

**Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)**

1   IT97K-499-35 is a blackeye variety with resistance to the parasitic plants Striga and Alectra,
2   developed at the International Institute of Tropical Agriculture in Ibadan, Nigeria (Singh *et al*.,
3   2006) and provided by Michael Timko (U Virginia, Charlottesville, Virginia, USA) to the
4   University of California Riverside (UCR) in 2006. The sequence assembly and annotation of
5   IT97K-499-35 were described in Lonardi *et al*. (2019). CB5-2 is a fully inbred isolate closely
6   related to CB5, the predominant Blackeye of the US Southwest for several decades. CB5
7   (Blackeye 8415) was bred by WW Mackie at the University of California (Mackie, 1946) to add
8   resistances to Fusarium wilt and nematodes to a California Blackeye landrace, and provided to
9   UCR by K Foster, University of California, Davis, in 1981. Suvita-2, also known as Gorom Local
10  (IITA accession TVu-15553, US NPGR PI 583259), is somewhat resistant to bruchids and certain
11  races of Striga and is relatively drought tolerant. This landrace was collected from a local market
12  by VD Aggarwal at the Institut de l'Environnement et de Recherches Agricoles (INERA) in
13  Burkina Faso (Aggarwal *et al*.,1984) and provided to UCR by VD Aggarwal in 1983. Sanzi is an
14  early flowering, small-seeded landrace from Ghana with resistance to flower bud thrips (Boukar
15  *et al*., 2013), provided by KO Marfo, Nyankpala Agricultural Experiment Station, Tamale, Ghana
16  to UCR in 1988. UCR779 (PI 583014) is a landrace from Botswana (de Mooy, 1985; Ehlers *et al*.,
17  2002) that was provided to UCR as B019-A in 1987 by CJ de Mooy of Colorado State University.
18  Yardlong bean or asparagus bean (cv.-gr. Sesquipedalis), the vegetable type of cowpea, is widely
19  grown in Asian countries for the consumption of tender long pods. TZ30 is an elite Chinese variety
20  with a pod length of around 60 cm. ZN016 is a landrace originating from southeastern China with
21  a pod length of about 35 cm and showing resistance to multiple major diseases of cowpea. TZ30
22  and ZN016 were used previously as parents of a mapping population to study the inheritance of
23  pod length (Xu *et al*., 2017).

24  **DNA sequencing and *de novo* assembly of seven cowpea accessions**

25  The annotated genome (v1.0) of African variety IT97K-499-35 was assembled from Pacific
26  Biosciences (Menlo Park, California, USA) long reads, two Bionano Genomics (San Diego,
27  California, USA) optical maps and ten genetic linkage maps as described previously (Lonardi *et*
28  *al.,* 2019). The six additional *de novo* assemblies were produced by Dovetail Genomics (Scotts
29  Valley, California, USA) using Illumina (San Diego, California, USA) short reads (150x2). DNA
30  was extracted by Dovetail Genomics from seedling tissue of CB5-2, TZ30, and ZN016, and seeds
31  of CB5-2, Suvita-2, Sanzi, and UCR779. Meraculous (Chapman *et al*., 2011) was used to assemble
32  the reads, then sequences from Dovetail Chicago® and Dovetail Hi-C® libraries were added
33  (using their proprietary pipeline) to resolve misassemblies and increase contiguity. These
34  assemblies were further refined using ALLMAPS (Tang *et al*., 2015). This analysis used ten
35  previously reported genetic linkage maps to relate assemblies to the standard orientations and

**Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)**

1    numbering of the eleven cowpea chromosomes, as described in Lonardi *et al*. (2019) for IT97K-
2    499-35. See "Data Availability Statement" for access to raw data and assemblies.

3    **Calling of SNPs, indels, and structural variants**

4    SNPs and indels were called using each reference genome versus the reads from the six other
5    accessions. Reads of each accession described above for genome assemblies, plus short-read
6    sequences produced by 10X Genomics from IT97K-49-35, were mapped to all assemblies using
7    BWA (Li *et al.,* 2009). SNPs and indels were called using the GATK 4.2.0 pipeline in GVCF mode
8    for each accession. All the per-sample GVCFs were gathered in joint genotyping to produce a set
9    of joint-called SNPs and indels. Both per-sample SNPs and joint-called SNPs were filtered with
10   the same parameters of 'QD < 2.0 || FS > 60.0 || MQ < 40.0 || MQRankSum < -12.5 ||
11   ReadPosRankSum < -8.0 || SOR > 4.0'. Indels were filtered with 'QD < 2.0 || FS > 200.0 ||
12   ReadPosRankSum < -20.0 || SOR > 10.0'.

13   Each pair of individual genomes was aligned using minimap2 (Li, 2018), producing $\binom{7}{2} = 21$
14   alignment files. Structural variants, including inversions and translocations, were identified from
15   the alignment files using SyRI (Goel *et al*., 2019). Figures were produced using PlotSR (Goel *et*
16   *al*., 2022). Depth analyses were carried out using Mosdepth (Pedersen & Quinlan 2018). The
17   average nucleotide diversity within and between populations was calculated from a VCF file using
18   Pixy (Korunes *et al*., 2021).

19   **Annotation of genes and repeats**

20   All genomes were annotated using the JGI plant genome annotation pipelines (Shu *et al.,* 2014),
21   integrated gene call (IGC), and gene model improvement (GMI). Both IGC and GMI are evidence-
22   based gene call pipelines. In IGC, a gene locus was defined by peptide alignments of related
23   organism homologous peptides and with alignments of within-organism transcriptome assemblies.
24   Genes were predicted by homology-based gene prediction programs FGENESH+ (Salamov and
25   Solovyev, 2000), FGENESH_EST, and GenomeScan (Yeh *et al.,* 2001), and a JGI in-house
26   homology-constrained transcriptome assembly ORF finder. Homologous proteomes included
27   *Arabidopsis thaliana* and those from common bean (*Phaseolus vulgaris*), soybean (*Glycine max*),
28   barrel medic (*Medicago truncatula*), poplar (*Populus trichocarpa*), rice (*Oryza sativa*), grape
29   (*Vitis vinifera*) and Swiss-Prot. For transcript-based annotations of the six new assemblies, RNA
30   for RNA-seq was extracted using Qiagen RNeasy Plant (Hilden, Germany) from each accession
31   from well-hydrated and drought-stressed young seedling root and leaves, immature flower buds,
32   and pods five days after pollination, and from developing seeds of Suvita-2, TZ30 and ZN016 (not
33   CB5-2, Sanzi or UCR779) 13 days after pollination. RNA quality was assessed, and concentrations

**Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)**

1    were determined using an Agilent 2100 BioAnalyzer (Santa Clara, California, USA) and the

2    Agilent RNA 6000 Nano Kit. The RNA-seq short reads from each accession were assembled using

3    a JGI in-house genome-guided assembler, PERTRAN (Shu *et al*., 2013), using each genome

4    assembly. Each short-read-based assembly and UNIGENE sequences (P12_UNIGENES.fa from

5    harvest.ucr.edu) were fed into PASA (Haas *et al.*, 2003) to produce transcriptome assemblies. The

6    best gene per locus (based on evidence) was defined using PASA from alignment of transcriptome

7    assemblies for splicing correctness, alternative transcripts, and UTR addition. The PASA genes

8    were filtered to obtain the final gene set, including an automated repeat coding sequence (CDS)

9    overlap filter, a manual low-quality gene filter, and an automatic filter from transposable element

10    (TE) protein domain assignments. This process was repeated once with one additional homology

11    seeding of non-self, high-confidence gene models.

12    **Determination of core and noncore genes among seven accessions**

13    Core and noncore genes were determined by running the GET_HOMOLOGUES-EST tool

14    (https://github.com/eead-csic-compbio/get_homologues) on the primary transcripts of the seven

15    cowpea accessions provided in nucleotide and protein formats. GET_HOMOLOGUES-EST was

16    run in orthoMCL-mode, as suggested by the authors for pan-genome analyses (Contreras-Moreira

17    *et al.*, 2017). The other GET_HOMOLOGUES-EST options "**-M -c -z -t 0 –A -L**" were used to

18    obtain orthoMCL gene clusters, which had genes in 1-7 accessions. The term "core" means that a

19    matching gene was identified in all seven accessions and "noncore" means that a matching copy

20    gene was identified in less than all seven accessions.

21    GO-term enrichment analyses were performed in agriGO v2.0 (Tian *et al.*, 2017) for core and

22    noncore genes using GO terms available from the Legume Information System

23    (https://www.legumeinfo.org/). Given the large number of GO terms in both the core and noncore

24    gene sets, GO slims (Onsongo *et al.*, 2008) were extracted and used for Figure 3. The full list of

25    core and noncore genes, with GO and other annotations, is available from the Google Drive noted

26    in the Data Availability Statement.

27    **Annotation of variants in core and noncore genes**

28    To test if variants in noncore genes have been subject to reduced selective constraint, Variant

29    Effect Predictor (VeP) (McLaren *et al.*, 2016) was used to annotate variants identified in the

30    primary transcripts of core and noncore genes. Gene annotations for IT97K-499-35 were used to

31    identify intervals that overlap core and noncore genes, and filtering of the VCF file used BEDtools

32    intersect (Quinlan & Hall, 2010) with variants called relative to the IT97K-499-35 assembly using

33    the six other assemblies. Scripts used for these analyses are at

**Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)**

1  https://github.com/MorrellLAB/Cowpea_Pangenome. VeP was run separately for SNPs and
2  indels, reporting classes of variants with potentially large effects, including missense, stop gains,
3  start or stop changes, and frameshifts. The numbers of synonymous changes and in-frame indels
4  are also reported.

5  **Relative size of core and noncore genes**

6  The physical sizes of core and noncore genes were compared in the total annotated length and the
7  length of the coding portion of the primary transcript of each gene. The length of each gene was
8  extracted from the general feature format (GFF) annotations. The CDS length was calculated based
9  on the primary transcript identified in Phytozome annotations (https://phytozome-
10  next.jgi.doe.gov/cowpeapan/info/Vunguiculata_v1_2). The full list of core and noncore genes,
11  with gene and CDS sizes indicted, is available from the Google Drive noted in the Data
12  Availability Statement.

13  **Nucleotide sequence diversity in cowpea**

14  Tajima's (1983) estimate of $\theta = 4Ne\mu$ was used to determine the level of sequence diversity in the
15  pangenome accessions. "Callable" regions were identified based on coverage estimates in
16  mosdepth (Pederson & Quinlan, 2018), with "callable" regions defined as those with coverage
17  between 5x and 400x. This estimate was derived from a sample with ~200X average coverage.
18  The callable regions were used to create a BED file used for filtering genomic regions. This
19  approach is intended to avoid variant calls in regions with inadequate sequence depth or regions
20  where very high coverage may indicate non-unique mapping of sequence reads. The callable
21  regions and the VCF file of filtered variants mapped to the IT97K-499-35 reference were used
22  with pixy (Korunes & Samuk, 2020), a tool designed to deal with missing data in genome-level
23  resequencing datasets.

24  **Physical locations of SNPs from genotyping platforms**

25  The physical positions of SNPs in the Illumina iSelect Cowpea Consortium Array (Muñoz-
26  Amatriaín *et al*., 2017), whose positions in the IT97K-499-35 genome were provided in Lonardi
27  *et al.* (2019), were mapped using BWA MEM (Li *et al*., 2009) within each of the seven assemblies
28  using the contextual sequence that flanked each variant. The resulting alignment file was processed
29  with SAMtools (Li *et al*., 2009) and SNP_Utils (https://github.com/MorrellLAB/SNP_Utils) to
30  report positions in a VCF file. The positions of iSelect SNPs relative to all seven genome
31  assemblies are provided in Supplemental Table S02, and an updated summary map for the 51,128
32  iSelect SNPs is in Supplemental Table S03. The positions identified for iSelect SNPs relative to

7

**Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)**

1    the IT97K-499-35 assembly were used to annotate the variants. The annotation used variant effect

2    predictor (VeP) (McLaren *et al*., 2016) with the GFF file provided by Phytozome

3    (https://genome.jgi.doe.gov/portal/) and SNP positions in VCF files

4    (https://github.com/MorrellLAB/cowpea_annotation/blob/main/Results/IT97K-499-

5    35_v1.0/iSelect_cowpea.vcf; see Data Availability Statement).

6    **Synteny analysis among genome assemblies**

7    To assess the conservation of gene content and ordering between genome assemblies from diverse

8    species, MCScanX (Wang *et al*., 2012) was run for every genome pair, using default settings and

9    homologous gene pairings derived from gene family assignments defined as the best match of the

10    longest protein product with an E-value of 1e-10 or better from hmmsearch (Eddy 2011) applied

11    to the legfed_v1_0 families (Stai *et al.*, 2019).

12    **Results and Discussion:**

13    **Development of six *de novo* assemblies and pan-genome construction**

14    Summary statistics for the seven assemblies (assembly characteristics, repetitive content, genes,

15    BUSCO completeness) are reported in Table 1. More detailed statistics of the intermediate

16    assembly steps are reported in Supplemental Table S04. The contiguity of the new six assemblies,

17    as indicated by their N50s, is comparable to the PacBio assembly for IT97K-499-35 despite being

18    based on short-read sequences. In all six new assemblies, each of the eleven chromosomes of

19    cowpea is represented by a single scaffold. These six assembled genomes are similar to each other

20    in size, ranging from 447.58 Mb to 453.97 Mb, with a mean of 449.91 Mb. IT97K-499-35 had a

21    ~15% larger (more complete) assembled size (519.44 Mb) than these six accessions, with the

22    difference attributable to long-read sequencing and optical mapping providing a more complete

23    assembly. Assemblies of the six additional accessions share the same percentage of repetitive

24    content of about 45-46% (Table 1 and Supplemental Figure S1). The IT97K-499-35 assembly has

25    a somewhat higher repetitive content than the assemblies of these six accessions. This may be

26    attributable to more complete resolution of unique positions of repetitive sequences within long

27    sequence reads than is possible from only short reads. A difference between the sequencing

28    methods in the resolution of repetitive sequences is evident in centromeric regions, which are

29    typically abundant in repetitive sequences, where some chromosomes of the six newly sequenced

30    accessions appear to be missing from the assemblies. Centromeric regions were defined based on

31    a 455-bp tandem repeat previously identified by fluorescence in situ hybridization (Iwata-Otsubo

32    *et al*., 2016). Supplemental Table S05 shows the coordinates of the putative centromeric regions

33    in IT97K-499-35 for all eleven chromosomes for a total span of 20.18 Mb, in CB5-2 on five

**Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)**

1   chromosomes for a total span of 5.6 Mb, in Sanzi on one chromosome for a total span of 0.59 Mb,

2   in ZN016 on four chromosomes for a total of 7.13 Mb and TZ30 on one chromosome for 1.32 Mb.

3   The tandem repeat was not found in any assembled chromosome of Suvita-2 or UCR779, nor in

4   the other chromosome assemblies where coordinates are not listed.

5   RNA was prepared from each accession to support gene annotation, and the same annotation

6   protocol was applied to each accession (see Materials & Methods). This is important when

7   comparing genomes at the gene level, as it reduces the technical variability that can otherwise

8   obfuscate the interpretation of results (Lei *et al.* 2021). The number of genes annotated in the six

9   new assemblies ranged from 27,723 to 28,562, with a mean of 28,222 (Table 1). IT97K-499-35

10  had ~13% more annotated genes, with a total of 31,948, reflecting deeper transcriptome

11  sequencing and, to some extent, the more complete assembly of its genome. Supplemental Table

12  S06 summarizes the number of alternative transcripts, exon statistics, gene model support, and

13  ontology annotations (Panther, PFam, KOG, KEGG, and E.C.). The number of alternative

14  transcripts in the six new assemblies ranged from 15,088 to 17,115. Again, IT97K-499-35 had a

15  higher number of alternative transcripts, a total of 22,536, than the other six accessions. The

16  average number of exons was 5.4 in each of the six new assemblies and 5.2 in IT97K-4899-35,

17  with a median length ranging from 162 to 169 bp. Gene and repeat density were computed in 1Mb

18  non-overlapping sliding windows along each chromosome and each accession (Supplemental

19  Figure S1). All chromosomes have a higher gene density in their more recombinationally active

20  regions, while repeat density peaks in the low-recombination centromeric and pericentromeric

21  regions (see also Supplemental Figure S8 in Lonardi *et al*., 2019). All seven accessions have

22  similar gene and repeat density, and high BUSCO v4 completeness at the genome, transcript, and

23  protein levels (Supplemental Table S07), with somewhat higher numbers for IT97K-499-35 than

24  the six new assemblies.

25  As stated above (Materials and Methods), genes annotated in the seven genomes were classified

26  as core if a matching gene was present in all accessions and noncore if absent in one or more of

27  the seven accessions. In IT97K-499-35, a total of 26,026 core genes (in 24,476 core clusters) and

28  4,963 noncore genes (in 4,285 noncore clusters) were identified (Supplemental Table S08). When

29  considering all seven accessions itemized in Supplemental Table S08, a total of 26,494 core genes

30  and  9,042 noncore genes (in 8,157 noncore clusters) were identified, resulting in a total of 35,536

31  pan genes in 32,633 pan gene clusters.

32  To determine if adding accessions significantly changed the numbers and proportions of core and

33  noncore genes, we took advantage of the analysis results produced by GET_HOMOLOGUES-

34  EST. GET_HOMOLOGUES-EST produces pan or core genome growth simulations by adding

35  accessions in random order, using twenty permutations. Figure 2 shows the growth of core and

**Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)**

1    pan genomes for an increasing number of accessions. A fitted Tettelin function (Tettelin *et al.*,

2    2005) is plotted in green. As expected, the number of pan genes increases as additional accessions

3    are "added" to the pan-genome, while the number of core genes decreases. However, the fact that

4    the core gene plot is flattening considerably (approaching an asymptotic limit) for six and seven

5    accessions indicates that most core genes have been identified with these seven diverse accessions.

6    In contrast, the pan-genome plot has not flattened, indicating that there may be many more noncore

7    genes not included among these seven accessions. Figure 2 provides an estimated 29,659 pan gene

8    clusters and an estimated 24,439 core gene clusters as the output of GET_HOMOLOGUES-EST

9    from 20 random samplings. Roughly, it appears that the pan-genome defined by the seven

10    cultivated cowpea accessions is comprised of about 80% core genes, constituting nearly the entire

11    set of core genes in cultivated cowpea, and 20% noncore genes. Clearly, more noncore genes

12    would be revealed with a larger number of accessions.

13    A GO term enrichment analysis was performed for genes within the two components of the pan-

14    genome (core and noncore) using agriGO v2 (Tian *et al.*, 2017). Many GO terms for all three

15    ontology aspects (biological process, cellular component, and molecular function) were

16    significantly enriched in both core and noncore genes (Supplemental Table S09). Given the high

17    number of significant GO terms, GO Slim terms (Onsongo *et al.*, 2008) were extracted and used

18    for Figure 3. Terms enriched in the core genes were related to transport and some metabolic

19    processes and molecular functions involving DNA-binding transcription factor activity (Figure 3;

20    Supplemental Table S09). This supports the idea that the core genome contains genes that perform

21    essential cellular functions that are highly conserved at the species level. The output was quite

22    different for the noncore genes, with very high enrichment of the GO term "response to stress"

23    (Figure 3), in particular "defense response" ($-\log_{10}q = 123.7$; Supplemental Table S09). This is

24    consistent with previous research showing that the "dispensable" genome encodes genes involved

25    in defense response and other beneficial functions for some individuals (Golicz *et al.*, 2016;

26    Gordon *et al.*, 2017; Montenegro *et al.*, 2017).

27    **Genetic variation analysis**

28    In addition to identifying gene PAVs (presence-absence variants), the seven assemblies were used

29    to identify other types of variation. Variants were detected using two different software pipelines,

30    depending on their size. SNPs and indels of length up to 300 nucleotides, both considered small

31    variants, were detected using GATK (see Materials & Methods). Larger structural variations,

32    including deletions, duplications, inversions, and translocations, were detected using SyRI (Goel

33    *et al.*, 2019).

10

**Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)**

1   Across all "callable" regions of the genome, average $\theta_\pi$ = 0.0111 (± 0.0549). At the
2   pseudomolecule level, average diversity was highest on Vu05, with $\theta_\pi$ = 0.0155 (± 0.0723), and
3   lowest on Vu10, with $\theta_\pi$ = 0.0095 (± 0.0447) (Supplemental Table S10). A mean diversity of ~1%
4   is higher than many grain crops, such as barley (Morrell et al., 2014; Schmid et al., 2018) and
5   roughly comparable to maize (Tittes et al., 2021). The observed diversity in the cowpea
6   pangenome sample is above average for herbaceous plants (Miller & Gross, 2011; Leffler et al.,
7   2012; Corbett-Detig et al., 2015).

8   For SNPs and indels, the genome of each accession was used in turn as the "reference," mapping
9   the reads for each of the six other accessions against that genome. For each, the six SNP sets
10  produced by GATK were merged by taking the union of the SNPs based on their location (i.e., a
11  SNP in two accessions was counted only once if it appeared in the same genomic position).
12  Supplemental Table S11 summarizes the number of SNPs detected, where the reference genome
13  is listed on each row. For instance, using Suvita-2 as the reference, 1,489,850 SNPs were detected
14  using mapped reads from CB5-2, compared to 2,625,678 SNPs using the reads from UCR779.
15  Combining the SNPs by counting all distinct SNPs in the union of the six sets of SNPs, the number
16  of SNPs for Suvita-2 was 5,292,933.

17  When UCR779 was used as the reference, a much higher number of SNPs was detected in every
18  pairwise comparison, indicating that UCR779 is the most divergent among these seven accessions.
19  Conversely, CB5-2 (a California cultivar) has fewer SNPs in pairwise comparisons to TZ30 or
20  ZN016 (both from China) than in pairwise comparisons to other accessions. This suggests that
21  CB5-2 is more similar to these two accessions than to the other four accessions. This is consistent
22  with genetic assignment analyses reported by Muñoz-Amatriaín *et al.* (2021) and historical
23  considerations discussed in Herniter *et al.* (2020). Supplemental Table S12 provides a similar
24  analysis for indels, where again, UCR779 stands out as the most different among the seven
25  accessions. Summary statistics for SNPs and indels for each chromosome and each accession can
26  be found in the file "SNPs_indels_stats.xlsx," available from the Google Drive indicated in the
27  Data Availability Statement below.

28  GATK requires a minimum coverage of 5X to call SNPs. Coverage analysis with Mosdepth
29  indicated that the average read coverage of IT97K-499-35 is very high (e.g., about ~190X when
30  mapping CB5-2 reads to IT97K-499-35), thus a very high fraction of IT97K-499-35 chromosomes
31  was covered by at least five reads. The lowest was Vu10 with 85.1%, the highest was Vu07 with
32  98.6%, and the overall percentage of SNPs in IT97K-499-35 that were in a "callable" region (i.e.,
33  with coverage 5x-400x) was 88.96%. The frequency of SNPs, as the number of unique SNPs
34  identified (Supplemental Table S11) divided by the size of the assembled genome (Table 1), ranges
35  from one in 139 to one in 309 bp, and the indel frequency (Supplemental Table S12) ranges from

**Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)**

1    one in 486 to one in 529 bp. Circos plots for SNP density (SNPs per Mb) on each chromosome
2    using each accession as the reference are in Supplemental Figure S2 (A-G), where it is evident, for
3    example, that Vu04 and Vu10 have the highest SNP frequency. In contrast, Vu05 and Vu09 have
4    the lowest. This was observed previously when mapping nearly one million SNPs on the IT97K-
5    499-35 reference genome (Lonardi *et al*., 2019). Also, when using UCR779 as the reference
6    (Supplemental Figure S2-E), the number of SNPs on Vu04 and Vu10 is significantly higher than
7    when any other accession is used as the reference, again consistent with UCR779 being the most
8    different among the seven accessions.

9    Structural variations were identified using SyRI (Goel *et al*., 2019) from the alignment of each pair
10   of individual genomes and visualized using PlotSR (Goel *et al*., 2022) (Figure 4). The visualization
11   shows a relatively large number of apparent structural rearrangements between the seven cowpea
12   genomes, which are more abundant in the centromeric and pericentromeric regions of all
13   chromosomes. Vu04 is the chromosome with the highest abundance of structural variants (Figure
14   4). A summary of all the structural variants identified in all pairs of accessions is reported in
15   Supplementary Table S13. The table shows that Suvita-2 versus UCR779 had the largest number
16   of inversions (2,008) and translocations (1,822). This intuitively makes sense since these two
17   accessions belong to two different genetic subpopulations separated by the first principal
18   component (Figure 1).

19   Inversions are a common type of rearrangement with important consequences for cross-over
20   frequency and distribution, as they suppress recombination in heterozygotes (Kirkpatrick, 2010).
21   While inversion can be important to maintaining locally adaptive variants (Kirkpatrick & Barton,
22   2006), crossover inhibition can impede plant breeding efforts. Table 2 summarizes the genomic
23   coordinates of all inversions larger than 1 Mbp. For example, the first column of Table 2,
24   corresponding to IT97K-499-35, shows 27 inversions that were identified by comparing the
25   reference genome against the other six accessions. The same inversion can appear in multiple sub-
26   tables. For instance, the ~4.2 Mb inversion on chromosome 3 previously described in (Lonardi *et*
27   *al*., 2019) occurs in the same orientation in six accessions and the opposite orientation only in
28   IT97K-499-35, so it is listed six times in the column for IT97K-499-35.

29   Similarly, the inversions on Vu04 and Vu05 are detected against five accessions. The ~9.0 Mb
30   inversion on Vu06 is the largest inversion found by SyRI, and its orientation is unique to Suvita-
31   2. However, this inversion appears to be due to an assembly imperfection. It is reported as
32   unoriented in the ALLMAPS output (Supplemental Table S14), and comparisons between optical
33   maps derived from Suvita-2 and another cowpea accession not included here indicate a non-
34   inverted orientation in Suvita-2 (unpublished). Also, as shown in Lonardi *et al.* (2019) and
35   Supplemental Figure S3, this entire region has a very low recombination rate and comprises nearly

**Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)**

1   the entire short arm of acrocentric chromosome 6 (Iwata-Otsubo *et al*., 2016). These factors can

2   account for a spurious orientation assignment for this region in the Suvita-2 Vu06 assembly.

3   The positions of the largest inversions shown in Figure 4 are provided in Table 2, e.g., the

4   inversions on Vu03 in IT97K-499-35 reported by Lonardi *et al*. (2019), and the inversion on Vu06

5   in Suvita-2 likely due to a mis-assembly, as discussed above. It should be noted that regions with

6   apparently low synteny within several chromosomes are low-recombination centromeric and

7   pericentromeric regions (Lonardi *et al*., 2019), which are notoriously hard to assemble due to their

8   high repetitive content and hard to orient due to a paucity of mapped and recombinationally

9   ordered SNPs. In these regions, it is expected to find compressed contigs, gaps, and misassemblies,

10  any of which might be flagged as apparent structural variations. The number of false-positive

11  structural variations can likely be reduced by increasing the completeness of the assemblies within

12  these regions using long-read sequencing and optical mapping. Supplemental Figure S4 (A-U)

13  shows all 21 SyRi+PlotSR alignments between all pairs of cowpea accessions.

14  **Further characterization of core and noncore genes**

15  Partitioning SNPs into those found in core versus noncore genes in IT97K-499-35 resulted in

16  702,073-SNPs in core genes and 239,100 SNPs in noncore genes. The indel comparison involves

17  161,900 indels in core genes and 39,845 in noncore genes. The numbers of variants with potential

18  consequences are summarized in Figure 5 and Supplementary Table S15. Counting both SNPs and

19  indels, there are 80,693 potentially benign variants among core genes (3.10 per gene) and 36,519

20  in noncore genes (7.36 per gene), which is a 2.37-fold higher frequency in noncore versus core

21  genes. Likewise, potentially harmful variants, including missense, stop gained, start or stop

22  change, and frameshift total 95,465 among core genes (3.67 per gene) and 75,048 in noncore genes

23  (15.12 per gene),which is a 4.12-fold higher incidence in noncore versus core genes. Among these,

24  noncore genes have a much higher incidence of frameshift variants (1.48 per gene) than do core

25  genes (0.23 per gene), this being a 6.43-fold difference. In each of these comparisons, noncore

26  genes contribute proportionally a larger number of variants than do core genes, whether benign or

27  potentially harmful.

28  Based on the gene annotations, core gene primary transcripts are longer than noncore gene primary

29  transcripts, with a mean length of 4,226.08 (± 4,047.234) for IT97K-499-35 core genes versus

30  2,341.32 bp (± 3,190.67) for IT97K-499-35 noncore genes (with median lengths of 3,292 and

31  1,347 bp, respectively). This difference is significant based on a non-parametric, two-sample

32  Wilcoxon rank sum test, with p-value $< 2.2$ e$^{-16}$. For IT97K-499-35, primary transcripts from core

33  genes cover 110.9 Mb of the genome, while primary transcripts from noncore genes cover 11.6

34  Mb. These differences in lengths could result from either longer coding regions or longer or more

13

**Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)**

1  abundant introns within the primary transcripts. When considering only the coding sequence
2  (CDS) for each IT97K-499-35 gene, the mean length of the CDS in core genes is greater than in
3  noncore genes, with a mean of 1,319.14 (± 960.61) for core versus 792.97 bp (± 915.98) for
4  noncore (with median lengths of 1,113 and 426 bp). Based on the Wilcoxon test, the difference in
5  length of the coding sequence is significant, with a p-value $< 2.2e^{-16}$. The explanation for this CDS
6  length difference is unknown.

7  **Presence-absence variation of genes controlling black seed coat color**

8  To facilitate the community's use of the cowpea pan-genome, all the genomes and their annotations
9  have been included as resources in the Legume Information System (LIS; www.legumeinfo.org;
10  Dash *et al.*, 2016). As an example of a use case for pan-genomics, the Genome Context Viewer
11  (GCV) is an application that enables dynamic comparison of genomes based on their gene content,
12  using assignments of genes to families as the basis for computation and visualization of conserved
13  gene order and structural variation with potential impact on function, e.g., copy number variation
14  (CNV) and presence-absence variation (PAV) (Cleary and Farmer, 2018). Figure 6A shows the
15  results of a query centered on a region from the reference cowpea genome that features a cluster
16  of tandemly duplicated MYB transcription factor genes in which presence-absence variation was
17  previously determined to be associated with seed coat pigmentation (Herniter et al., 2018). The
18  colors of the genes in this "beads on a string" representation reflect the gene family assignments;
19  here, the brown triangles in the center of the region represent the MYB genes with varying copy
20  numbers in the different cowpea accessions, with a maximum of five copies in the reference
21  accession to as few as a single copy in UCR779. Outside the CNV region, there is strong
22  conservation of gene content, with one other region showing some evidence of reordering among
23  the cowpea accessions. The viewer facilitates comparison not only within but across species, and
24  one can see evidence of similar CNV in the corresponding region of several *Phaseolus* spp.
25  genomes (Schmutz *et al.* 2014, Moghaddam *et al.* 2021), as well as an inversion of the segment
26  containing the genes relative to cowpea, soybean (Valliyodan *et al.* 2019) and other *Vigna* species
27  (Sakai *et al.* 2015, Kang *et al.* 2014). Two corresponding homoeologous regions evidence the most
28  recent whole genome duplication in soybean. The region serves as a breakpoint for the syntenic
29  block in Gm09, which, taken together with the other structural variation, suggests that the
30  expansion of gene copy number here has had consequences for the stability of the chromosome in
31  these regions over evolutionary time (Hastings et al., 2009).

32  Although the GCV view shows good evidence for CNV, there are some limitations to what may
33  be inferred from that alone. First, since the viewer only has access to gene family assignment
34  information, it cannot determine which elements among those in tandem arrays have the highest
35  sequence similarity and provide insight into which copies have been deleted. Second, because it

**Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)**

1   relies on the surrounding genomic context of each gene to place it into correspondence, it will have
2   limited capability for finding genes that are present in the assembly but are largely isolated on
3   small scaffolds that were not incorporated into the main pseudomolecules. Another tool at LIS that
4   provides a complementary view based on the underlying sequence identity of the different copies
5   of the expanded gene family is shown in Figure 6B. Here, the InterMine (Kalderamis *et al*., 2014)
6   instance for cowpea (https://mines.legumeinfo.org/cowpeamine/begin.do) was used to collect all
7   protein sequences for cowpea genes assigned to the given family. A dynamic tree construction
8   procedure invoked based on hmmalign-derived
9   (http://www.csb.yale.edu/userguides/seq/hmmer/docs/node18.html; Eddy, 2011) additions of
10  these genes to the multiple sequence alignment for the founding members of the family. The
11  resulting tree (a subtree of which is shown) allows the user to determine the best correspondences
12  of the copies in each genome and pulls in two additional genes on unanchored contigs that likely
13  belong to the region.

14  **Pangenome core genes and cross-species synteny**

15  To explore the question of how within-species gene content conservation compares with gene
16  content shared between species in other species and genera, we used the LIS gene family
17  assignments to define homology pairings between all members of each gene family, then used the
18  resulting data to determine collinearity blocks among all pairwise comparisons of the cowpea
19  genomes, as well as to soybean and representative genomes from *Vigna* and *Phaseolus* spp. The
20  counts of genes participating in at least one collinear block were tallied for each genome in each
21  pairwise comparison. As expected, intra-specific comparisons between cowpea accessions yield
22  higher numbers of conserved collinear genes than inter-specific comparisons. On the other hand,
23  there is no appreciable difference in the extent of conserved collinearity when comparing cowpea
24  genomes to other species within the *Vigna* genus versus species from *Phaseolus* or *Glycine* genera
25  (Supplemental Figure S5). Because soybean has an additional whole genome duplication relative
26  to all other species in the comparison, the total number of soybean genes found in collinear blocks
27  is higher than in other comparisons. Comparisons between all species and the *Vigna radiata*
28  version 6 genome (Kang *et al.* 2014) show fewer conserved collinear genes, but this is presumably
29  due to missing data in that assembly, given that all other interspecific comparisons are similar.

30

15

**Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)**

**Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)**

**Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)**

## Literature Cited

Aggarwal VD, Muleba N, Drabo I, Souma J, Mbewe M (1984) Inheritance of *Striga gesnerioides* resistance in cowpea. In: Proceedings of 3rd International Symposium on Parasitic Weeds (Parker C, Musselman LJ, Polhill RM, Wilson AK eds.) ICARDA, Aleppo, Syria, pp. 143-147

Bayer PE, Golicz AA, Scheben A, Batley J, Edwards D (2020) Plant pan-genomes are the new reference. Nature Plants 6:914-920.

Boukar O, Belko N, Chamarthi S, Togola A, Batieno J, Owusu E, Haruna M, Diallo S, Umar ML, Olufajo O, Fatokun C (2019) Cowpea (*Vigna unguiculata*): Genetics, genomics and breeding. Plant Breeding 138:415-424, doi.org/10.1111/pbr.12589

Boukar O, Bhattacharjee R, Fatokun C, Kumar PL, Gueye B (2013) Cowpea. In: Genetic and Genomic Resources of Grain Legume Improvement (Singh M, Upadhyaya HD, Bisht IS eds.), Chapter 6, pp. 137-156, Elsevier, London, UK. ISBN 978-0-12-397935-3

Cao J, Schneeberger K, Ossowski S, Günther T, Bender S, Fitz, J, Koenig D, Lanz C, Stegle O, Lippert C, Wang X, Ott F, Müller J, Alonso-Blanco C, Borgwardt K, Schmid KJ, Weigel D (2011) Whole-genome sequencing of multiple *Arabidopsis thaliana* populations. Nature Genetics 43:956–963, doi.org/10.1038/ng.911

Chapman JA, Ho I, Sunkara S, Luo S, Schroth GP, Rokhsar DS (2011) Meraculous: *de novo* genome assembly with short paired-end reads. PLOS ONE 6:e23501

Cleary A, Farmer A (2018) Genome Context Viewer: visual exploration of multiple annotated genomes using microsynteny. Bioinformatics 34:1562-1564, doi.org/10.1093/bioinformatics/btx757

Contreras-Moreira B, Cantalapiedra CP, Garcia Pereira MJ, Gordon S, Vogel JP, Igartua E, Casas AM, Vinuesa P (2017) Analysis of plant pan-genomes and transcriptomes with GET_HOMOLOGUES-EST, a clustering solution for sequences of the same species. Frontiers in Plant Science 8:184, doi.org/10.3389/fpls.2017.00184

Corbett-Detig RB, Hartl DL, Sackton TB (2015) Natural selection constrains neutral diversity across a wide range of species. PLOS Biology 13:e1002112, doi: 10.1371/journal.pbio.1002112

**Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)**

1  Darling AE, Mau B, Perna NT (2010) progressiveMauve: multiple genome alignment with gene
2    gain, loss and rearrangement. PLOS ONE 5:e11147, doi.org/10.1371/journal.pone.0011147

3  Dash S, Campbell JD, Cannon EK, Cleary AM, Huang W, Kalberer SR, Karingula V, Rice AG,
4    Singh J, Umale PE, Weeks NT, Wilkey AP, Farmer AD, Cannon SB. (2016) Legume
5    information system (LegumeInfo.org): a key component of a set of federated data resources for
6    the legume family. Nucleic Acids Research 44:D1181-1188, doi: 10.1093/nar/gkv1159

7  de Mooy BE (1985) Germplasm evaluation of Botswana cowpea (*Vigna unguiculata* [L.] Walp.)
8    landraces. M.S. Thesis, Michigan State University, Lansing, Michigan, USA.

9  Eddy SR (2011) Accelerated Profile HMM Searches. PLoS Computational Biology 7:e1002195,
10   doi: 10.1371/journal.pcbi.1002195.

11  Ehlers JD, Fery RL, Hall AE (2002) Cowpea breeding in the USA: new varieties and improved
12   germplasm. *In* Challenges and Opportunities for Enhancing Sustainable Cowpea Production.
13   Proceedings of the World Cowpea Conference III held at the International Institute of Tropical
14   Agriculture (IITA), Ibadan, Nigeria, 4-8 September 2000, (Eds. Fatokun CS, Tarawali SA,
15   Singh BB, Kormawa PM, Tamò M), Chapter 1.6, pp. 62-77, IITA, Ibadan, Nigeria, ISBN 978-
16   131-190-8

17  Goel M, Schneeberger K (2022) plotsr: visualizing structural similarities and rearrangements
18   between multiple genomes. Bioinformatics 38:2922-2926,
19   doi.org/10.1093/bioinformatics/btac196

20  Goel M, Sun H, Jiao W-B, Schneeberger K (2019) SyRI: finding genomic rearrangements and
21   local sequence differences from whole-genome assemblies. Genome Biology 20:277,
22   doi:10.1186/s13059-019-1911-0

23  Golicz AA, Bayer PE, Bhalla PL, Batley J, Edwards D (2020) Pangenomics comes of age: from
24   bacteria to plant and animal applications. Trends in Genetics 36:132-145

25  Golicz AA, Bayer PE, Barker GC, Edger PP, Kim HR, Martinez PA, Chan CKK, Severn-Ellis A,
26   McCombie R, Parkin IAP, Paterson AH, Pires JC, Sharpe AG, Tang H, Teakle GR, Town CD,
27   Bately J, Edwards D (2016) The pangenome of an agronomically important crop plant
28   *Brassica oleracea*. Nature Communications 7:13390, doi.org/10.1038/ncomms13390

29  Gordon SP, Contreras-Moreira B, Woods DP, Des Marais DL, Burgess D, Shu S, Stritt C, Roulin
30   AC, Schackwitz W, Tyler L, Martin J, Lipzen A, Dochy N, Phillips J, Barry K, Geuten K,

**Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)**

1   Budak H, Juenger TE, Amasino R, Caicedo AL, Goodstein D, Davidson P, Mur LAJ, Figueroa
2   M, Freeling M, Catalan P, Vogel JP (2017) Extensive gene content variation in the
3   *Brachypodium distachyon* pan-genome correlates with population structure. Nature
4   Communications 8:2184, doi.org/10.1038/s41467-017-02292-8

5   Haas BJ, Delcher AL, Mount SM, Wortman JR, Smith Jr RK, Hannick LI, Maiti R, Ronning
6   CM, Rusch DB, Town CD, Salzberg SL, White O (2003) Improving the *Arabidopsis* genome
7   annotation using maximal transcript alignment assemblies. Nucleic Acids Research 31:5654-
8   5666, http://nar.oupjournals.org/cgi/content/full/31/19/5654

9   Hastings PJ, Lupski JR, Rosenberg SM, Ira G (2009) Mechanisms of change in gene copy
10  number. Nature Review Genetics 10:551-564, doi: 10.1038/nrg2593. PMID: 19597530;
11  PMCID: PMC2864001.

12  Hall AE (2004) Breeding for adaptation to drought and heat in cowpea. European Journal of
13  Agronomy 21:447-454

14  Herniter IA, Muñoz-Amatriaín M, Close TJ (2020) Genetic, textual, and archeological evidence
15  of the historical global spread of cowpea (*Vigna unguiculata* [L.] Walp.). Legume Science 3:57

16  Herniter IA, Muñoz-Amatriaín M, Lo S, Guo Y-N, Close TJ (2018) Identification of candidate
17  genes controlling black seed coat and pod tip color in G3 8:3347-3355

18  Hirsch CN, Foerster JM, Johnson JM, Sekhon RS, Muttoni G, Vaillancourt B, Peñagaricano F,
19  Lindquist E, Pedraza MA, Barry K, de Leon N, Kaeppler SM, Buell CR (2014) Insights into
20  the maize pan-genome and pan-transcriptome. The Plant Cell 26:121–135

21  Iwata-Otsubo A, Lin J-Y, Gill N, Jackson SA (2016) Highly distinct chromosome structures in
22  cowpea (*Vigna unguiculata*), as revealed by molecular cytogenetic analysis. Chromosome
23  Research 24:197-216

24  Kalderimis A, Lyne R, Butano D, Contrino S, Lyne M, Heimbach J, Hu F, Smith R, Stěpán R,
25  Sullivan J, Micklem G (2014) InterMine: extensive web services for modern biology. Nucleic
26  Acids Research 42:W468-472, doi: 10.1093/nar/gku301

27  Kang YJ, Kim SK, Kim MY, Lestari P, Kim KH, Ha BK, Jun TH, Hwang WJ, Lee T, Lee J,
28  Shim S, Yoon MY, Jang YE, Han KS, Taeprayoon P, Yoon N, Somta P, Tanya P, Kim KS,
29  Gwag JG, Moon JK, Lee YH, Park BS, Bombarely A, Doyle JJ, Jackson SA, Schafleitner R,

**Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)**

1  Srinives P, Varshney RK, Lee SH (2014) Genome sequence of mungbean and insights into
2  evolution within *Vigna* species. Nature Communications 5:5443, doi: 10.1038/ncomms6443

3  Kent WJ (2002) BLAT - the BLAST-like alignment tool. Genome Research 12:656-64

4  Kirkpatrick M (2010) How and why chromosome inversions evolve. PLOS Biology 8:e1000501

5  Kirkpatrick M, Barton N (2006) Chromosome inversions, local adaptation and speciation.
6  Genetics 173:419-434

7  Korunes KL, Samuk K (2021) PIXY: Unbiased estimation of nucleotide diversity and
8  divergence in the presence of missing data. Molecular Ecology Resources 21:1359-1368

9  Leffler EM, Bullaughey K, Matute DR, Meyer WK, Ségurel L, Venkat A, Andolfatto P,
10  Przeworski M (2012) Revisiting an old riddle: what determines genetic diversity levels within
11  species. PLOS Biology 10:e1001388

12  Lei L, Goltsman E, Goodstein D, Wu GA, Rokhsar DS, Vogel JP (2021) Plant pan-genomics
13  comes of age. Annual Reviews of Plant Biology 72:411-435

14  Li H (2018) Minimap2: pairwise alignment for nucleotide sequences. Bioinformatics 34:3094-
15  3100

16  Li H, Durbin R (2009) Fast and accurate short read alignment with Burrows-Wheeler transform.
17  Bioinformatics 25:1754–1760

18  Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R
19  (2009) 1000 Genome Project Data Processing Subgroup, The Sequence Alignment/Map format
20  and SAMtools. Bioinformatics 25:2078–2079

21  Li YH, Zhou G, Ma J, Jiang W, Jin LG, Zhang Z, Guo Y, Zhang J, Sui Y, Zheng L, Zhang SS,
22  Zuo Q, Shi XH, Li YF, Zhang WK, Hu Y, Kong G, Hong HL, Tan B, Song J, Liu ZX, Wang
23  Y, Ruan H, Yeung CKL, Liu J, Wang H, Zhang LJ, Guan RX, Wang KJ, Li WB, Chen SY,
24  Chang RZ, Jiang Z, Jackson SA, Li R, Qiu LJ (2014) *De novo* assembly of soybean wild
25  relatives for pan-genome analysis of diversity and agronomic traits. Nature Biotechnology
26  32:1045–1052, doi.org/10.1038/nbt.2979

27  Liang Q, Lonardi S (2021) Reference-agnostic representation and visualization of pan-genomes.
28  BMC Bioinformatics 22:502

**Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)**

1   Lonardi S, Muñoz-Amatriaín M, Liang Q, Shu S, Wanamaker SI, Lo S, Tanskanen J, Schulman
2      AH, Zhu T, Luo MC, Alhakami H, Ounit R, Hasan AM, Verdier J, Roberts PA, Santos JRP,
3      Ndeve A, Doležel J, Vrána J, Hokin SA, Farmer AD, Cannon SB, Close TJ (2019) The
4      genome of cowpea (*Vigna unguiculata* [L.] Walp.). The Plant Journal 98:767-782

5   Mackie WW (1946) Blackeye beans in California. Bulletin 696, University of California
6      Agricultural Experiment Station

7   McLaren W, Gil L, Hunt SE, Riat HS, Ritchie GRS, Thormann A, Flicek P, Cunningham F
8      (2016) The Ensembl variant effect predictor. Genome Biology 17:122

9   Miller AJ, Gross BL (2011) From forest to field: perennial fruit crop domestication. American
10     Journal of Botany 98:1389-1414

11  Moghaddam SM, Oladzad A, Koh C, Ramsay L, Hart JP, Mamidi S, Hoopes G, Sreedasyam A,
12     Wiersma A, Zhao D, Grimwood J, Hamilton JP, Jenkins J, Vaillancourt B, Wood JC, Schmutz
13     J, Kagale S, Porch T, Bett KE, Buell CR, McClean PE (2021) The tepary bean genome
14     provides insight into evolution and domestication under heat stress. Nature Communications
15     12:2638, doi: 10.1038/s41467-021-22858-x

16  Montenegro JD, Golicz AA, Bayer PE, Hurgobin B, Lee H, Chan C-KK., Visendi P, Lai K,
17     Doležel J, Batley J, Edwards D (2017) The pangenome of hexaploid bread wheat. The Plant
18     Journal 90:1007-1013, doi.org/10.1111/tpj.13515

19  Morgante M, De Paoli E, Radovic S (2007) Transposable elements and the plant pan-genomes.
20     Current Opinions in Plant Biology 10:149–155

21  Morrell PL, Gonzales AM, Meyer KKT, Clegg MT (2014) Resequencing data indicate
22     domestication's modest effect on barley diversity: a cultigen with multiple origins. Journal of
23     Heredity 105:253-264, doi: 10.1093/jhered/est083

24  Muñoz-Amatriaín M, Eichten SR, Wicker T, Richmond TA, Mascher M, Steuernagel B, Scholz
25     U, Ariyadasa R, Spannagl M, Nussbaumer T, Mayer KF, Taudien S, Platzer M, Jeddeloh JA,
26     Springer NM, Muehlbauer GJ, Stein N (2013) Distribution, functional impact, and origin
27     mechanisms of copy number variation in the barley genome. Genome Biology 14:R58

28  Muñoz-Amatriaín M, Lo S, Herniter IA, Boukar O, Fatokun C, Carvalho M, Castro I, Guo Y-N,
29     Huynh B-L, Roberts PA, Carnide V, Close TJ (2021) The UCR Minicore: a resource for
30     cowpea research and breeding. Legume Science 3:e95, doi.org/10.1002/leg3

**Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)**

1  Muñoz-Amatriaín M, Mirebrahim H, Xu P, Wanamaker SI, Luo M, Alhakami H, Alpert M,
2  Atokple I, Batieno BJ, Boukar O, Bozdag S, Cissé N, Drabo I, Ehlers JD, Farmer A, Fatokun
3  C, Gu YQ, Guo Y, Huynh B, Jackson SA, Kusi F, Lawley CT, Lucas MR, Ma Y, Timko MP,
4  Wu J, You F, Barkley NA, Roberts PA, Lonardi S, Close TJ (2017) Genome resources for
5  climate-resilient cowpea, an essential crop for food security. The Plant Journal 89:1042-1054

6  Onsongo G, Xie H, Griffin TJ, Carlis (2008) Generating GO slim using relational database
7  management systems to support proteomics analysis. 21$^{st}$ IEEE International Symposium on
8  Computer-Based Medical Systems, doi.org/10.1109/CBMS.2008.77

9  Pedersen BS, Quinlan AR (2018) Mosdepth: quick coverage calculation for genomes and
10  exomes. Bioinformatics 34:867–868, doi.org/10.1093/bioinformatics/btx699

11  Quinlan AR, Hall IM (2010) BEDtools: a flexible suite of utilities for comparing genomic
12  features. Bioinformatics 26:841-841

13  Sakai H, Naito K, Ogiso-Tanaka E, Takahashi Y, Iseki K, Muto C, Satou K, Teruya K, Shiroma
14  A, Shimoji M, Hirano T, Itoh T, Kaga A, Tomooka N (2015) The power of single molecule
15  real-time sequencing technology in the de novo assembly of a eukaryotic genome. Scientific
16  Reports 5:16780, doi:10.1038/srep16780

17  Salamov AA, Solovyev VV (2000) *Ab initio* gene finding in *Drosophila* genomic DNA. Genome
18  Research 10:516-522

19  Schmid K, Kilian B, Russell J (2018) Barley Domestication, Adaptation and Population
20  Genomics. editors. In: Stein N, Muehlbauer G (eds.) The Barley Genome. Compendium of
21  Plant Genomes. Springer, Cham. pp. 317-336.

22  Schmutz J, McClean PE, Mamidi S, Wu GA, Cannon SB, Grimwood J, Jenkins J, Shu S, Song
23  Q, Chavarro C, Torres-Torres M, Geffroy V, Moghaddam SM, Gao D, Abernathy B, Barry K,
24  Blair M, Brick MA, Chovatia M, Gepts P, Goodstein DM, Gonzales M, Hellsten U, Hyten DL,
25  Jia G, Kelly JD, Kudrna D, Lee R, Richard MM, Miklas PN, Osorno JM, Rodrigues J, Thareau
26  V, Urrea CA, Wang M, Yu Y, Zhang M, Wing RA, Cregan PB, Rokhsar DS, Jackson SA
27  (2014) A reference genome for common bean and genome-wide analysis of dual
28  domestications. Nature Genetics 46:707-713, doi:10.1038/ng.3008

29  Shu S, Goodstein D, Rokhsar D (2013) PERTRAN: genome-guided RNA-seq read assembler.
30  https://www.osti.gov/biblio/1241180

**Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)**

1    Shu S, Rokhsar D, Goodstein D, Hayes D, Mitros T (2014) JGI plant genomics gene annotation
2    pipeline. Lawrence Berkeley National Laboratory, Berkeley, CA, USA

3    Singh BB, Olufajo OO, Ishiyaku MF, Adeleke RA, Ajeigbe HA, Mohammed SG (2006)
4    Registration of six improved germplasm lines of cowpea with combined resistance to *Striga*
5    *gesnerioides* and *Alectra vogelii*. Crop Science 46:2332-2333

6    Springer NM, Ying K, Fu Y, Ji T, Yeh CT, Jia Y, Wu W, Richmond T, Kitzman J, Rosenbaum
7    H, Iniguez AL, Barbazuk WB, Jeddeloh JA, Nettleton D, Schnable PS (2009) Maize inbreds
8    exhibit high levels of copy number variation (CNV) and presence/absence variation (PAV) in
9    genome content. PLoS Genetics 5:e1000734

10   Stai JS, Yadav A, Sinou C, Bruneau A, Doyle JJ, Fernández-Baca D, Cannon SB (2019) Cercis:
11   a non-polyploid genomic relic within the generally polyploid legume family. Frontiers in Plant
12   Science 10:345, doi:10.3389/fpls.2019.00345

13   Tajima F (1983) Evolutionary relationship of DNA sequences in finite populations. Genetics
14   105:437-460

15   Tang H, Zhang X, Miao C, Zhang J, Ming R, Schnable JC, Schnable PS, Lyons E, Lu J (2015)
16   ALLMAPS: robust scaffold ordering based on multiple maps. Genome Biology 16:3,
17   doi:10.1186/s13059-014-0573-1

18   Tettelin H, Masignani V, Cieslewicz MJ, Donati C, Medini D, Ward NL, Angiuoli SV, Crabtree
19   J, Jones AL, Durkin AS, DeBoy RT, Davidsen TM, Mora M, Scarselli M, Immaculada
20   Margarit y Ros I, Peterson JD, Hauser CR, Sundaram JP, Nelson WC, Madupu R, Brinkac LM,
21   Dodson RJ, Rosovitz MJ, Sullivan SA, Daugherty SC, Haft DH, Selengut J, Gwinn ML, Zhou
22   L, Zafar N, Khouri H, Radune D, Dimitrov G, Watkins K, O'Connor KJB, Smith S, Utterback
23   TR, White O, Rubens CE, Grandi G, Madoff LC, Kasper DL, Telford JL, Wessels MR,
24   Rappuoli R, Fraser CM (2005) Genome analysis of multiple pathogenic isolates of
25   *Streptococcus agalactiae*: implications for the microbial "pan-genome". Proc. Natl. Acad.
26   Sci. USA 102:13950-13955, doi:10.1073/pnas.0506758102

27   Tian T, Liu Y, Yan H, You Q, Yi X, Du Z, Xu W, Su Z (2017) agriGO v2.0: a GO analysis
28   toolkit for the agricultural community, 2017 update. Nucleic Acids Research 45:W122-W129,
29   doi: 10.1093/nar/gkx382

**Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)**

1  Tittes S, Lorant A, McGinty S, Doebley JF, Holland JB, de Jesus Sánchez-González, Seetharam
2    A, Tenaillon M, Ross-Ibarra J (2021) Not so local: the population genetics of convergent
3    adaptation in maize and teosinte. bioRxiv, doi: 10.1101/2021.09.09.459637

4  Torkamaneh D, Lemay M-A, Belzile F (2021) The pan-genome of the cultivated soybean
5    (PanSoy) reveals an extraordinarily conserved gene content. Plant Biotechnology Journal
6    19:1852-1862, https://doi.org/10.1111/pbi.13600

7  Valliyodan B, Cannon SB, Bayer PE, Shu S, Brown AV, Ren L, Jenkins J, Chung CY, Chan TF,
8    Daum CG, Plott C, Hastie A, Baruch K, Barry KW, Huang W, Patil G, Varshney RK, Hu H,
9    Batley J, Yuan Y, Song Q, Stupar RM, Goodstein DM, Stacey G, Lam HM, Jackson SA,
10   Schmutz J, Grimwood J, Edwards D, Nguyen HT (2019) Construction and comparison of three
11   reference-quality genome assemblies for soybean. The Plant Journal 100:1066-1082, doi:
12   10.1111/tpj.14500

13  Wang Y, Tang H, Debarry JD, Tan X, Li J, Wang X, Lee TH, Jin H, Marler B, Guo H, Kissinger
14   JC, Paterson AH (2012) MCScanX: a toolkit for detection and evolutionary analysis of gene
15   synteny and collinearity. Nucleic Acids Research 40:e49, doi: 10.1093/nar/gkr1293

16  Xu P, Wu X, Muñoz-Amatriaín M, Wang B, Wu X, Hu Y, Huynh BL, Close TJ, Roberts PA,
17   Zhou W, Lu Z, Li G (2017) Genomic regions, cellular components and gene regulatory basis
18   underlying pod length variations in cowpea (*V. unguiculata* L. Walp). Plant Biotechnology
19   Journal 15:547-557, doi:10.1111/pbi.12639

20  Yeh R-F, Lim LP, Burge CB (2001) Computational inference of homologous gene structures in
21   the human genome. Genome Research 11:803-816

22  Yu G, Wang L, Han Y, He Q (2012). clusterProfiler: an R package for comparing biological
23   themes among gene clusters. OMICS: A Journal of Integrative Biology 16:284-287,
24   doi:10.1089/omi.2011.0118

25

**Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)**

1 **Figure and Table Captions**

2 **Figure 1. Principal component analysis of the UCR Minicore, indicating the accessions**
3 **selected for sequencing and the subpopulation they belong to.** Accessions in the plot are
4 colored by the result of STRUCTURE for *K=6*, as shown in Muñoz-Amatriaín *et al.* (2021).

5 **Figure 2. The number of genes identified in the pan-genome (pan genes) and core genome**
6 **(core genes) as new accessions are added.** Green curves are fitted Tettelin functions.

7 **Figure 3. Gene Ontology (GO) term enrichment analysis.** Significantly enriched GO terms for
8 core (A) and noncore genes (B) are shown for GO-Slim categories belonging to Biological
9 Process, Cellular Component, and Molecular Function aspects (in different colors). -log10 of
10 FDR-adjusted p-values (q-values) are shown on the right of each bar.

11 **Figure 4. Representation of structural variations (of any size) detected by SyRI from the**
12 **output of whole-genome pairwise alignments between the seven cowpea accessions.** The
13 black track indicates gene density in the reference genome IT97K-499-35, while the blue track
14 indicates SNP density in the reference genome IT97K-499-35.

15 **Figure 5. Variant effect predictor (VeP) annotations for SNPs and indels found in the core**
16 **and noncore genes present in IT97K-499-35.** Values on the y-axis are the absolute number of
17 variants in each variant class.

18 **Figure 6. Conservation of gene content within and across species.** (A) A region depicting
19 gene content conservation and variability among cowpea genomes and other representative
20 Phaseoleae species. Triangular glyphs represent order and orientation of genes, with color
21 representing gene family memberships. (https://vigna.legumeinfo.org/tools/gcv) (B) All cowpea
22 proteins assigned to the family whose members exhibit copy number variation in (A) are shown
23 augmenting a dynamically recomputed gene tree at the Legume Information System, with genes
24 from unanchored contigs not present in the chromosomes aligned in (A) indicated with arrows
25 (https://mines.legumeinfo.org/cowpeamine).

26 **Table 1. Summary of assembly statistics, repetitive content, gene content, and BUSCO**
27 **completeness for the seven genomes.**

28 **Table 2. Genomic coordinates of all inversions of size > 1 Mbp detected by comparing the**
29 **seven cowpea genomes pairwise.** IT97K- 499-35 is abbreviated as IT97K.

**Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)**

1 **Supplemental Figure S1. Gene and repeat density.** Accessions are represented by different
2 shades of red (genes) and blue (repeats). The reference genome of IT97K-499-35 is the longest
3 curve, i.e., the one that extends furthest to the right of each graph.

4 **Supplemental Figure S2. SNP density (number of SNPs per Mb) using each genome as the**
5 **"reference."** (A) IT97K-499-35 (IT97K), centromeres are marked with orange in the innermost
6 circle, (B) CB5-2, (C) Suvita-2, (D) Sanzi, (E) UCR779, (F) ZN016, (G) TZ30.

7 **Supplemental Figure S3. Output of ALLMAPS for chromosome Vu06 for Suvita-2.** Ten
8 genetic maps were used to orient the five Dovetail contigs in Vu06 (two of which are larger than
9 1Mb – see Supplemental Table S14). The first four were arbitrarily oriented by ALLMAPS due
10 to low recombination in that region, as shown on the graphs on the right, which plot cM position
11 (y-axis) as a function of physical position (x-axis). In particular, the 8.2 Mb contig represented in
12 gray in the bottom left figure is a region of very low recombination frequency and was likely
13 oriented incorrectly.

14 **Supplemental Figure S4. Structural variants (of any size) detected by SyRI between any**
15 **pairs of genomes in this study.** (A) CB5-2 vs Sanzi, (B) CB5-2 vs Suvita-2, (C) CB5-2 vs
16 TZ30, (D) CB5-2 vs UCR779, (E) CB5-2 vs ZN016, (F) IT97K-499-35 vs CB5-2, (G) IT97K-
17 499-35 vs Sanzi, (H) IT97K-499-35 vs Suvita-2, (I) IT97K-499-35 vs TZ30, (J) IT97K-499-35
18 vs UCR779, (K) IT97K-499-35 vs ZN016, (L) Sanzi vs TZ30, (M) Sanzi vs UCR779, (N) Sanzi
19 vs ZN016, (O) Suvita-2 vs Sanzi, (P) Suvita-2 vs TZ30, (Q) Suvita-2 vs UCR779, (R) Suvita-2
20 vs ZN016, (S) UCR779 vs TZ30, (T) UCR779 vs ZN016, (U) ZN016 vs TZ30.

21 **Supplemental Figure S5**. **Macrosynteny views.** (A) Macrosynteny view with blocks
22 representing regions in the IT97K-499-35 reference cowpea genome with conserved gene order
23 relative to each of the genomes shown as tracks below. The region from the microsynteny view
24 of Figure 6A is shown with a vertical gray bar, and the set of chromosomes displayed is
25 restricted to those showing synteny in that region (i.e., the non-cowpea chromosomes have an
26 apparent lack of synteny downstream because of genomic rearrangements that have moved
27 corresponding content to other chromosomes than those shown). Various inversions are seen as
28 blocks with orientations opposing those of their neighboring blocks. Gaps in otherwise syntenic
29 regions indicate regions where gene content diversity outweighs conserved content through
30 presence-absence and copy-number variation. (B) Counts of genes participating in conserved
31 collinear blocks for all pairwise genome comparisons among the cowpea pangenome members
32 and across representative genomes from several genera in the Phaseoleae tribe. Self-comparisons
33 are included to illustrate within-species conservation of duplicated content from ancient whole

27

**Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)**

1   genome duplication (WGD) events shared by subfamily Faboideae species and the more recent

2   WGD in the *Glycine max* genome.

3   **Supplemental Table S01. Cowpea accessions used in this work.**

4   **Supplemental Table S02. Cowpea iSelect SNP positions on each of the seven genome**

5   **assemblies.** The allele, chromosome and position in the assembled genome are indicated for

6   each accession (columns D-X). For IT97K-499-35 (97K) the orientation of the sequence used for

7   the cowpea iSelect array (Muñoz-Amatriaín et al. 2017) "forward" strand is indicated in column

8   B and the two possible alleles for iSelect assay "forward" strand are in column C.

9   **Supplemental Table S03. Cowpea iSelect arrray SNP positions and alleles relative to the**

10  **IT97K-499-35 sequence of Lonardi et al. 2019 (columns B-F,I,J) and to the iSelect**

11  **"Forward Strand" of Muñoz-Amatriaín et al. 2017 and Muñoz-Amatriaín et al. 2021**

12  **(columns G&H).** Other columns indicate reasons for exclusion of data from 2,316 SNPs: blastn

13  alignment ambiguities (columns K-M,S), poor technical performance on array (column N),

14  monomorphic across all DNA samples (column O), excess heterozygote and or no-call (columns

15  P-R).

16  **Supplemental Table S04. Statistics of the six new assemblies at each step of the Dovetail**

17  **assembly pipeline.**

18  **Supplemental Table S05. Putative centromeric region coordinates (all numbers are bp).**

19  **Supplemental Table S06. Gene annotation statistics.**

20  **Supplemental Table S07. BUSCO v4 completeness results.**

21  **Supplemental Table S08. Core and noncore genes identified from sequencing the seven**

22  **cowpea genomes, tabulated by gene cluster.**

23  **Supplemental Table S09. Enrichment analysis of GO Terms for core and noncore genes**

24  **performed in AgriGO v2.** Only significantly enriched GO terms (FDR < 0.05) are shown for

25  the three different ontology aspects.

26  **Supplemental Table S10. Average diversity at the chromosome (pseudomolecule) level**

27  **relative to the IT97K-4899-35 assembly.** Values reported are $\theta\pi$ and the standard deviation for

28  "callable regions."

**Title: A view of the pan-genome of domesticated cowpea (*Vigna unguiculata* [L.] Walp.)**

1    **Supplemental Table S11. Number of SNPs when considering each accession as the**
2    **"reference" genome and the resulting union of unique SNPs (merged GVCF) for each**
3    **accession.**

4    **Supplemental Table S12. Number of indels of size 1 to 300 bp when considering each**
5    **accession as the "reference" genome and the union set of all indels (merged GVCF) for**
6    **each accession.**

7    **Supplemental Table S13. Genomic coordinates of all structural variants detected via SyRI**
8    **by comparing the seven cowpea genomes pairwise.** IT97K-499-35 is abbreviated as IT97K.

9    **Supplemental Table S14. Largest inversions, using each of the Dovetail assemblies as**
10    **reference.** Left table: ALLMAPS' orientation of assembled contigs based on markers' position
11    on the genetic maps ("?" indicates a contig that was arbitrarily oriented). Right tables: Large
12    (>1Mb) inversions detected by SyRI, and whether they are within an oriented ALLMAPS contig.

13    **Supplemental Table S15. Summary of nucleotide sequence variants in core and noncore**
14    **genes with potential consequences on coding sequence as identified by Variant Effect**
15    **Predictor (VeP).** SNPs and indels were analyzed separately. These values are shown in Figure
16    5. Predictions are based on annotations from the IT97K-499-35 genome assembly.

**Figure 1. Principal component analysis of the UCR Minicore, indicating the accessions selected for sequencing and the subpopulation they belong to.** Accessions in the plot are colored by the result of STRUCTURE for *K*=6, as shown in Muñoz-Amatriaín et al. (2021).

**Figure 2. The number of genes identified in the pan-genome (left) and core genome (right) as new accessions are added.** Green curves are fitted Tettelin functions.

**Figure 3. Gene Ontology (GO) term enrichment analysis.** Significantly enriched GO terms for core (A) and noncore genes (B) are shown for GO-Slim categories belonging to Biological Process, Cellular Component, and Molecular Function aspects (in different colors). -log$_{10}$ of FDR-adjusted p-values (q-values) are shown on the right of each bar.

**Figure 4. Representation of structural variations (of any size) detected by SyRI from the output of whole-genome pairwise alignments between the seven cowpea accessions.** The black track indicates gene density in the reference genome IT97K-499-35, while the blue track indicates SNP density in the reference genome IT97K-499-35.

**Figure 5. Variant effect predictor (VeP) annotations for SNPs and indels found in the core and noncore genes present in IT97K-499-35.** Values on the y-axis are the absolute number of variants in each variant class.

**Figure 6A. Conservation of gene content within and across species.** A region depicting gene content conservation and variability among cowpea genomes and other representative Phaseoleae species. Triangular glyphs represent order and orientation of genes, with color representing gene family memberships. (https://vigna.legumeinfo.org/tools/gcv).

**Figure 6B. Conservation of gene content within and across species.** All cowpea proteins assigned to the family whose members exhibit copy number variation in Figure 6A are shown augmenting a dynamically recomputed gene tree at the Legume Information System, with genes from unanchored contigs not present in the chromosomes aligned in 6A indicated with arrows (https://mines.legumeinfo.org/cowpeamine).

**Table 1. Summary of assembly statistics, repetitive content, gene content and BUSCO4 completeness for the seven genomes**.

| | IT97K-499-35 | CB5-2 | Suvita-2 | Sanzi | UCR779 | ZN016 | TZ30 |
|---|---|---|---|---|---|---|---|
| Assembly size (bp) | 519,435,864 | 448,043,751 | 447,585,192 | 447,277,261 | 453,970,486 | 451,130,807 | 451,468,680 |
| N50 (bp) | 41,684,185 | 36,897,245 | 36,142,647 | 34,759,918 | 35,700,653 | 37,764,243 | 36,906,789 |
| #Contigs/scaffolds | 686 | 6,534 | 9,123 | 11,268 | 12,939 | 7,032 | 6,771 |
| #Contigs/scaffolds ≥ 100kbp | 103 | 28 | 28 | 17 | 13 | 28 | 48 |
| #Contigs/scaffolds ≥ 1Mbp | 13 | 11 | 11 | 11 | 11 | 11 | 11 |
| #Contigs/scaffolds ≥ 10Mp | 11 | 11 | 11 | 11 | 11 | 11 | 11 |
| Longest contig (bp) | 65,292,630 | 60,086,998 | 58,539,223 | 58,655,738 | 58,369,212 | 60,653,587 | 59,481,915 |
| Repetitive content | 47.25% | 45.52% | 45.43% | 45.50% | 45.89% | 45.68% | 45.76% |
| Annotated genes (#) | 31,948 | 28,297 | 28,545 | 28,461 | 28,562 | 27,723 | 27,742 |
| *BUSCO completeness* | | | | | | | |
| Genome | 1595   98.8% | 1574   97.5% | 1580   97.8% | 1581   97.9% | 1574   97.6% | 1589   98.5% | 1583   98.1% |
| Transcripts | 1594   98.8% | 1570   97.2% | 1582   98.0% | 1585   98.2% | 1581   97.9% | 1584   98.1% | 1580   97.8% |
| Proteins | 1595   98.8% | 1569   97.3% | 1584   98.2% | 1587   98.3% | 1585   98.2% | 1584   98.1% | 1582   98.0% |

**Table 2. Genomic coordinates of all inversions of size > 1 Mbp detected by comparing the seven cowpea genomes pairwise.**
IT97K- 499-35 is abbreviated as IT97K.

**CBS-2 vs IT97K**

| IT97K | | | CBS-2 | | |
|---|---|---|---|---|---|
| chr03 | 36,118,990 | 40,333,678 | chr03 | 32,390,474 | 36,391,036 |
| chr04 | 17,622,506 | 20,917,095 | chr04 | 14,956,441 | 17,919,052 |
| chr05 | 25,746,455 | 27,269,915 | chr05 | 23,493,794 | 24,846,894 |
| chr10 | 17,517,032 | 18,768,535 | chr10 | 14,803,744 | 15,920,434 |
| chr11 | 30,575,657 | 33,619,557 | chr11 | 26,382,727 | 29,314,033 |

**Sanzi comparisons**

| IT97K | | | Sanzi | | | | CBS-2 | | | Sanzi | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| chr03 | 36,118,990 | 40,333,678 | chr03 | 31,262,282 | 35,207,418 | | chr05 | 23,469,613 | 24,641,440 | chr05 | 22,486,336 | 23,465,592 |
| chr04 | 17,944,012 | 20,826,275 | chr04 | 13,863,317 | 16,451,333 | | chr11 | 13,755,783 | 16,506,588 | chr11 | 13,607,535 | 14,719,408 |
| chr11 | 15,536,600 | 19,182,958 | chr11 | 12,656,484 | 13,536,503 | | chr11 | 26,389,390 | 29,309,356 | chr11 | 24,170,618 | 27,116,888 |

**Suvita2 comparisons**

| IT97K | | | Suvita2 | | | | CBS-2 | | | Suvita2 | | | | Sanzi | | | Suvita2 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| chr01 | 7,882,945 | 11,829,903 | chr01 | 6,880,499 | 10,585,839 | | chr01 | 6,882,777 | 10,598,583 | chr01 | 6,880,226 | 10,585,839 | | chr01 | 6,841,294 | 10,522,595 | chr01 | 6,890,381 | 10,585,839 |
| chr03 | 36,118,990 | 40,333,678 | chr03 | 31,052,542 | 35,034,290 | | chr06 | 6,407 | 8,165,285 | chr06 | 58,197 | 8,215,985 | | chr06 | 384,940 | 7,854,307 | chr06 | 12,864 | 8,215,985 |
| chr04 | 17,620,322 | 21,043,166 | chr04 | 14,040,335 | 16,532,009 | | chr11 | 26,392,061 | 29,309,384 | chr11 | 25,612,417 | 28,506,945 | | | | | | | |
| chr05 | 25,755,748 | 27,141,671 | chr05 | 23,368,354 | 24,577,329 | | | | | | | | | | | | | | |
| chr06 | 15,232 | 8,928,750 | chr06 | 105,004 | 8,215,985 | | | | | | | | | | | | | | |
| chr10 | 17,518,369 | 18,714,514 | chr10 | 14,746,438 | 15,817,466 | | | | | | | | | | | | | | |

**TZ30 comparisons**

| IT97K | | | TZ30 | | | | CBS-2 | | | TZ30 | | | | Sanzi | | | TZ30 | | | | Suvita2 | | | TZ30 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| chr01 | 8,078,984 | 11,774,887 | chr01 | 7,037,437 | 12,547,938 | | chr01 | 7,073,895 | 10,598,703 | chr01 | 6,952,274 | 12,533,400 | | chr01 | 7,006,278 | 10,479,491 | chr01 | 7,037,894 | 12,533,400 | | chr04 | 14,053,478 | 15,710,525 | chr04 | 14,784,086 | 16,411,900 |
| chr03 | 36,118,990 | 40,333,678 | chr03 | 31,657,822 | 35,692,568 | | chr04 | 14,858,370 | 16,730,160 | chr04 | 14,787,251 | 16,411,900 | | chr06 | 25,185,327 | 26,275,534 | chr06 | 26,599,616 | 27,728,577 | | chr06 | 58,197 | 8,156,866 | chr06 | 48,056 | 8,235,862 |
| chr04 | 21,097,140 | 22,121,274 | chr04 | 17,439,719 | 18,187,637 | | chr06 | 26,597,860 | 27,660,154 | chr06 | 26,590,042 | 27,686,349 | | | | | | | | | chr06 | 26,021,210 | 27,076,895 | chr06 | 26,597,598 | 27,696,602 |
| chr05 | 25,746,455 | 27,141,671 | chr05 | 23,450,412 | 24,709,628 | | chr11 | 26,383,812 | 29,292,784 | chr11 | 26,278,507 | 29,207,281 | | | | | | | | | | | | | | |
| chr06 | 28,973,468 | 30,151,695 | chr06 | 26,590,042 | 27,712,655 | | | | | | | | | | | | | | | | | | | | |
| chr10 | 17,517,032 | 18,768,535 | chr10 | 14,761,180 | 15,883,623 | | | | | | | | | | | | | | | | | | | | |

**UCR779 comparisons**

| IT97K | | | UCR779 | | | | CBS-2 | | | UCR779 | | | | Suvita2 | | | UCR779 | | | | TZ30 | | | UCR779 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| chr03 | 36,119,072 | 40,333,678 | chr03 | 30,810,075 | 34,783,649 | | chr11 | 26,392,061 | 29,312,944 | chr11 | 25,119,570 | 28,030,631 | | chr06 | 110,090 | 8,055,043 | chr06 | 25,279 | 8,037,739 | | chr04 | 14,846,045 | 16,411,900 | chr04 | 13,944,667 | 15,345,835 |
| chr04 | 17,502,725 | 21,030,423 | chr04 | 13,944,667 | 16,474,544 | | | | | | | | | | | | | | | | | | | | |
| chr05 | 25,746,455 | 27,075,873 | chr05 | 22,985,675 | 24,003,974 | | | | | | | | | | | | | | | | | | | | |

**ZN016 comparisons**

| IT97K | | | ZN016 | | | | CBS-2 | | | ZN016 | | | | Sanzi | | | ZN016 | | | | Suvita2 | | | ZN016 | | | | TZ30 | | | ZN016 | | | | UCR779 | | | ZN016 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| chr01 | 7,839,391 | 11,820,835 | chr01 | 7,019,074 | 10,785,201 | | chr01 | 6,857,029 | 10,589,386 | chr01 | 7,020,634 | 10,774,128 | | chr01 | 6,813,588 | 10,522,595 | chr01 | 7,020,980 | 10,775,930 | | chr04 | 15,405,491 | 17,105,526 | chr04 | 16,118,615 | 17,254,093 | | chr04 | 15,234,860 | 16,411,900 | chr04 | 14,817,248 | 15,900,761 | | chr01 | 6,946,955 | 10,601,930 | chr01 | 7,020,634 | 10,728,684 |
| chr03 | 36,118,990 | 40,333,678 | chr03 | 32,843,668 | 36,858,590 | | chr04 | 16,349,146 | 17,910,419 | chr04 | 16,191,362 | 17,334,982 | | chr06 | 25,198,638 | 26,275,534 | chr06 | 26,369,920 | 27,460,156 | | chr06 | 60,001 | 8,215,985 | chr06 | 745 | 8,344,949 | | chr10 | 12,029,208 | 15,883,623 | chr10 | 12,026,103 | 15,749,232 | | chr06 | 25,789,928 | 26,848,252 | chr06 | 26,369,920 | 27,427,414 |
| chr04 | 18,707,087 | 20,943,046 | chr04 | 14,895,114 | 17,690,481 | | chr06 | 26,590,252 | 27,653,879 | chr06 | 26,366,480 | 27,439,152 | | chr10 | 11,712,476 | 14,725,119 | chr10 | 12,070,488 | 15,606,154 | | chr06 | 26,003,233 | 27,074,234 | chr06 | 26,369,920 | 27,458,526 | | | | | | | | | chr10 | 12,451,287 | 15,886,883 | chr10 | 12,177,925 | 15,742,142 |
| chr05 | 25,746,455 | 27,223,043 | chr05 | 24,265,392 | 25,383,012 | | chr10 | 12,100,642 | 15,920,434 | chr10 | 12,026,103 | 15,751,269 | | | | | | | | | chr10 | 12,037,878 | 15,817,466 | chr10 | 12,027,428 | 15,751,264 | | | | | | | | | | | | | | |
| chr06 | 28,999,740 | 30,127,069 | chr06 | 26,369,920 | 27,439,171 | | chr11 | 26,357,807 | 29,309,414 | chr11 | 27,061,318 | 30,061,182 | | | | | | | | | | | | | | | | | | | | | | | | | | |
| chr10 | 14,534,255 | 18,390,632 | chr10 | 12,350,844 | 15,751,259 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |