

Exploring Structural Variation and Gene Family Architecture with De Novo Assemblies of 15 Medicago Genomes

Running head: ***De novo assemblies of Medicago genomes***

**Peng Zhou¹, Kevin A. T. Silverstein², Thiruvarangan Ramaraj³, Joseph Guhlin⁴,
Roxanne Denny¹, Junqi Liu⁵, Andrew D. Farmer³, Kelly P. Steele⁶, Robert M.
Stupar⁵, Jason R. Miller⁷, Peter Tiffin⁴, Joann Mudge³, Nevin D. Young^{1,4}***

Affiliations

1 Department of Plant Pathology, University of Minnesota, St. Paul, MN, USA

2 Supercomputing Institute for Advanced Computational Research, University of Minnesota, Minneapolis, MN, USA

3 National Center for Genome Resources, Santa Fe, NM, USA

4 Department of Plant Biology, University of Minnesota, St. Paul, MN, USA

5 Department of Agronomy and Plant Genetics, University of Minnesota, St. Paul, MN, USA

6 Science and Mathematics Faculty, Arizona State University, Mesa, AZ, USA

7 J. Craig Venter Institute, Rockville, MD, USA

* Corresponding Author

Email Addresses

Peng Zhou: zhoux379@umn.edu

Kevin A. T. Silverstein: kats@umn.edu

Thiruvarangan Ramaraj: tr@ncgr.org

Joseph Guhlin: guhli007@umn.edu

Roxanne Denny: denny004@umn.edu

Junqi Liu: liuqx162@umn.edu

Andrew D. Farmer: adf@ncgr.org

Kelly P. Steele: Kelly.Steele@asu.edu

Jason R. Miller: jmiller@jcvi.org

Peter Tiffin: ptiffin@umn.edu

Joann Mudge: jm@ncgr.org

Nevin D. Young: nevin@umn.edu

1 **Abstract**

2 **Background**

3 Previous studies exploring sequence variation in the model legume, *Medicago*
4 *truncatula*, relied on mapping short reads to a single reference. However, read-
5 mapping approaches are inadequate to examine large, diverse gene families or to
6 probe variation in repeat-rich or highly divergent genome regions. *De novo*
7 sequencing and assembly of *M. truncatula* genomes enables near-comprehensive
8 discovery of structural variants (SVs), analysis of rapidly evolving gene families, and
9 ultimately, construction of a pan-genome.

10 **Results**

11 Genome-wide synteny based on 15 *de novo M. truncatula* assemblies effectively
12 detected different types of SVs indicating that as much as 22% of the genome is
13 involved in large structural changes, altogether affecting 28% of gene models. A
14 total of 63 million base pairs (Mbp) of novel sequence was discovered, expanding
15 the reference genome space for *Medicago* by 16%. Pan-genome analysis revealed
16 that 42% (180 Mbp) of genomic sequences is missing in one or more accession,
17 while examination of *de novo* annotated genes identified 67% (50,700) of all
18 ortholog groups as dispensable – estimates comparable to recent studies in rice,
19 maize and soybean. Rapidly evolving gene families typically associated with biotic
20 interactions and stress response were found to be enriched in the accession-specific

1 gene pool. The nucleotide-binding site leucine-rich repeat (NBS-LRR) family, in
2 particular, harbors the highest level of nucleotide diversity, large effect single
3 nucleotide change, protein diversity, and presence/absence variation. However, the
4 leucine-rich repeat (LRR) and heat shock gene families are disproportionately
5 affected by large effect single nucleotide changes and even higher levels of copy
6 number variation.

7 **Conclusions**

8 Analysis of multiple *M. truncatula* genomes illustrates the value of *de novo*
9 assemblies to discover and describe structural variation, something that is often
10 under-estimated when using read-mapping approaches. Comparisons among the *de*
11 *novo* assemblies also indicate that different large gene families differ in the
12 architecture of their structural variation.

1 **Background**

2 Legumes comprise a diverse and ecologically significant plant family that serves
3 as the second most important crop family in the world [1]. As a cool season legume,
4 *Medicago truncatula* is closely related to important crops such as alfalfa (*Medicago*
5 *sativa*), clover (*Trifolium pratense* and *T. repens*), pea (*Pisum sativum*), chickpea
6 (*Cicer arietinum*), and *Lotus japonicas* [2,3]. *M. truncatula* was chosen as a model for
7 studying legume biology due to its small genome size, simple diploid genetics, self-
8 fertility, short generation time, amenability to genetic transformation and large
9 collections of diverse ecotypes [3–5]. *M. truncatula* research has focused especially
10 on its symbiotic relationship with rhizobia and arbuscular mycorrhizae, root
11 development, secondary metabolism and disease resistance [3,6]. A high quality,
12 BAC-based sequence has served as the original “reference genome” for the
13 *Medicago* research community [7] while re-sequencing of additional accessions has
14 enriched the pool of sequence data available [8,9].

15 In plants, large gene families play a crucial role in both biotic interactions and
16 abiotic response. Some of these families are encoded by hundreds of members [10–
17 12] organized in clusters of varying size and thought to evolve through gene
18 duplication and birth-and-death processes [13–17]. Widely studied examples include
19 the nucleotide-binding site, leucine-rich repeat proteins (NBS-LRRs), receptor-like
20 kinases (RLKs), F-box proteins, leucine-rich repeat proteins (LRRs), heat shock

1 proteins (HSPs), and protein kinases [16–20]. In *M. truncatula* and close taxonomic
2 relatives, an additional gene family is important in symbiotic nitrogen fixation, the
3 nodule-specific cysteine-rich peptides (NCRs), a sub-family within the larger
4 cysteine-rich peptide (CRP) superfamily [21–24]. Legume NCRs are highly expressed
5 in rhizobial nodules [22,24,25] where they act as plant effectors directing bacteroid
6 differentiation [26]. NCR genes are abundant, diverse, and frequently clustered
7 [23,24].

8 Previous studies of plant genomes highlighted the important role that gene
9 families play in the architecture of structural variation (SV) (reviewed in [27]). Array-
10 based re-sequencing of 20 *Arabidopsis* accessions indicated that 60% of NBS-LRRs,
11 25% of F-box, and 16% of RLKs exhibited some type of major-effect polymorphism
12 compared with less than 10% for all expressed sequences [28]. In *Arabidopsis*, 33.3%
13 of the NBS-LRR genes in the Columbia reference are deleted in at least one of 80
14 accessions compared with just 12.5% of genes in the *Arabidopsis* genome as a whole
15 [29]. In rice, Schatz et al [30] re-sequenced three divergent genomes and found that
16 genes containing the NB-ARC domain (signature motif of NBS-LRRs) constituted 12%
17 of lineage-specific genes compared with just 0.35% of genes shared among all three
18 genomes.

19 In contrast to earlier alignment-based (read-mapping) studies of sequence
20 diversity, *de novo* sequencing and assembly of genomes from multiple accessions

1 enables near-comprehensive discovery of SVs, gene family membership, and
2 ultimately, construction of a pan-genome. Here, we describe *de novo* genome
3 assemblies for 15 *M. truncatula* accessions, which we analyze together with the *M.*
4 *truncatula* reference. We were especially interested in the level and type of SVs
5 found in different gene families, with a focus on families associated with biotic
6 interactions and abiotic stress. Our results illustrate how different gene families
7 exhibit distinctly different variant architectures, including differing representation
8 within the dispensable portion of the pan-genome.

9

10 **Results**

11 ***De novo* assemblies have scaffold N50s > 250 kb, capturing > 90% of the *M.***

12 ***truncatula* gene space**

13 Fifteen *M. truncatula* accessions were sequenced with Illumina HiSeq2000 using
14 a combination of short and long insert paired-end libraries to an average of 120-fold
15 coverage, then assembled using ALLPATHS-LG [31] (Figure S1, Table S1). Between 80
16 and 94% of each genome could be assembled into scaffolds >100 kbp, with scaffold
17 N50s ranging from 268 kbp to 1,653 kbp and contig N50 sizes averaging around 20
18 kbp (Table S2). Assembled genome sizes ranged from 388 Mbp to 428 Mbp (Table
19 S2), correlating well with cytologically derived genome size estimates ($r = 0.83$, $P =$

1 0.005, Figure S2). Genomes were repeat-masked with a *Medicago*-specific repeat
2 database [32]. About 20% of each assembly was annotated as repeat, which is
3 slightly lower than the 23% repetitive content in *Medicago* reference Mt4.0, (based
4 on accession HM101, also known as A17) (Table S2). The *de novo* assemblies also
5 capture 87 - 96% of unique content in the reference genome, including 90 - 96% of
6 all Mt4.0 gene coding regions.

7

8 **Genic features in *de novo* assemblies largely resemble those of the reference**

9 All 15 genome assemblies were annotated using Augustus [33] incorporating *ab*
10 *initio* gene prediction results, RNA-Seq expression evidence from a subset of
11 accessions as well as protein homolog support from Mt4.0 reference gene models
12 (See Methods). Evidence-guided annotation yielded comparable numbers of coding
13 genes (60,000 to 67,000) for each of the 15 assemblies (Table S3). On average 80-
14 90% of predicted gene models receive support from either RNA-Seq expression or
15 Mt4.0 syntenic homologs. The number of TE-related genes in different accessions
16 (15,000 to 20,000, Table S3) was up to 25% lower than in the Mt4.0 reference,
17 indicating that some *de novo* assemblies missed or collapsed repetitive sequences. A
18 closer look at the number of TE categories suggests certain families were more likely
19 to be missed or collapsed than others (Data file S1). Median protein length (TEs

1 excluded) ranged from 245 to 254 amino acids – nearly equal to the estimate of 255
2 AAs in Mt4.0.

3

4 **Structural variants span as much as 22% of the *M. truncatula* genome**

5 Between 92 and 96% of each assembly could be aligned with the Mt4.0
6 reference typically leading to ~300 Mbp of sequences in syntenic blocks where single
7 nucleotide polymorphisms (SNPs), short InDels, and large SVs could be confidently
8 predicted (Tables S4-S6). Global comparisons revealed long syntenic blocks
9 intermixed with shorter, poorly aligned regions that harbor numerous structural
10 changes (Figures 1 and 2). The pattern of synteny alignment generally reflects
11 across-accession relationships inferred from SNP data (Figure S1), including three
12 “outgroup” accessions (HM022, HM340 and HM324) that are typically considered
13 separate sub-species with distinct diversity patterns compared with the remaining
14 accessions.

15 Within aligned genomic regions, extensive variation including SNPs, short
16 InDels, and large SVs were observed. Between 1.7 million (HM058) and 5.1 million
17 (HM340) SNPs were identified in comparisons with HM101 (Mt4.0) (Table S6). As
18 expected, SNP density correlates well with divergence from HM101 – with SNP bp⁻¹
19 ranging from 0.63% in HM058 (closest to HM101) to 2.37% in HM340 (most distant
20 from HM101). Estimates of nucleotide diversity ($\theta_{\pi} = 0.0073 \text{ bp}^{-1}$) are nearly 70%

1 higher than previous reports ($\theta_{\pi} = 0.0043 \text{ bp}^{-1}$ based on a broader 26 accession
2 panel) (Table S4, see Discussion) [8] . Approximately 70% of *Medicago* SNPs were
3 found in intergenic regions, which are also distinguished by the highest level of
4 nucleotide diversity ($\theta_{\pi} = 0.0089 \text{ bp}^{-1}$) (Table S4). Diversity was much higher for
5 synonymous than replacement polymorphisms in coding regions (Table S4). These
6 findings are consistent with the expectation of stronger purifying selection acting at
7 replacement sites, especially large-effect polymorphisms that significantly alter the
8 protein product [34].

9 Beyond SNPs, we identified 500,000 to 1,500,000 short InDels (<50 bp), 27,000 -
10 110,000 large InDels, 49,000 - 169,000 copy number variants (CNVs), and 2,700 -
11 12,700 translocations. SVs were identified through a rigorous syntenic anchoring
12 approach with each SV receiving support from synteny alignments of both flanking
13 sequences and being free from any intra- or inter- scaffold gaps (see Methods).
14 Nevertheless, these number may still underestimate the true level of variation given
15 that 4 to 8% of each genome could not be covered by our synteny alignment and
16 therefore likely to involve additional complex changes (Table S5). In count, SVs are
17 far less numerous than single-base variants, yet each of these SV classes affects
18 more total base pairs. Small InDels affect 3 – 10 Mbp, large insertions and deletions
19 affect 7.5 – 30 Mbp, CNVs affect 26 – 85 Mbp, and translocations affect 3.5 – 14
20 Mbp (Table S6). Altogether between 7% (HM058) and 22% (HM022) of genome

1 content is affected by at least one type of structural change (Table S6). This is
2 consistent with findings in other systems where large variants typically affect more
3 bases than SNPs [35,36]. Nearly equivalent numbers of small insertions versus
4 deletions were observed in contrast to traditional read mapping-based approaches
5 (which incorrectly predict more deletions than insertions relative to the reference
6 sequence [37,38]). Nonetheless, large deletions and copy number losses were still
7 30-50% higher, even with our use of synteny-based variant discovery, indicating
8 reduced power in detecting large insertions and copy number gains (Table S6).

9 To estimate the accuracy of our SV prediction, we performed PacBio sequencing
10 on three accessions (HM034, HM056 and HM340). For each SV, the number of
11 PacBio reads fully spanning ± 500 bp of the breakpoints was counted and scored as
12 valid only if each of its breakpoints received at least five supporting PacBio reads.
13 Based on these criteria, between 88 and 94% of all synteny-based SV calls could be
14 validated using long read technology (Table S7). Insertion and deletion of unique
15 (single-copy) genomic contents tended to have higher validation rates than gain or
16 loss of repetitive genomic contents (i.e., copy number gain or loss). This is consistent
17 with assembly quality in repetitive regions generally being lower than in unique
18 regions. Also, SVs involving genic regions tend to have the highest validation rates
19 compared with other genomic contexts (TEs, unknown genes, intergenic). Some of
20 the genic SVs provide good candidates in studying gene birth-and-death processes.

1 As an example, we identified a tandem duplication of a NBS-LRR gene in HM034 (or
2 gene deletion in HM101) which is supported by long PacBio reads (Figure S3)
3 Interestingly, the altered gene copy doesn't have RNA-Seq expression, whereas all
4 the neighboring copies do, a possible indication of pseudogene removal.

5 Global comparisons revealed long, conserved syntenic blocks intermixed with
6 shorter, poorly aligned regions that harbor numerous structural changes (Figures 1
7 and 2). The global pattern of synteny alignment generally reflect the *Medicago*
8 phylogeny – with three “outgroup” accessions (HM022, HM340 and HM324) that are
9 typically considered separate sub-species showing distinct diversity pattern from the
10 remaining accessions (Figure 1, Figure 2A). Nevertheless, peri-centromeric locations
11 generally display increased levels of diversity (and reduced levels of synteny) due to
12 enrichment of transposable elements (TEs) (Figure 1). In genomic regions where
13 synteny disappears altogether, our ability to identify different variant types (i.e.
14 SNPs, short InDels, or structural variants) also disappears. This is illustrated in Figure
15 2 (panels B-E) where high densities of TEs and selected gene families (RLKs, NBS-
16 LRRs, LRRs) are associated with reduced synteny coverage and loss of power in
17 detecting all variant types (grey areas). Non-centromeric regions with higher TE
18 density show high level of diversity and reduced synteny (e.g., Figure 1B, Figure 2B).
19 Like TEs, large clusters of NBS-LRRs, RLKs and LRRs lead to fragile genome
20 architecture and higher level of diversity (Figure 2 C-E). Genomic locations of these

1 gene family clusters are generally uncorrelated with one another, but there are
2 notable examples they co-localize (Figure 2 C-E). In these highlighted regions,
3 substantial clusters of NBS-LRRs, RLKs, NCRs, LRRs and F-box genes are all found
4 within a single 1Mb segment.

5

6 **180 Mbp is dispensable sequence out of a total pan-genome content of 430 Mbp**

7 Sequences that could not be aligned to the Mt4.0 reference even at relaxed
8 stringency (~80% sequence identity) were extensive across the 15 *de novo*
9 assemblies. These sequences often exist in the form of novel insertions or complex
10 substitutions, sometimes as separate scaffolds. After filtering potential contaminant
11 sequences, we identified between 9 and 22 Mbp of novel segments (1.3 to 2.4 Mbp
12 in coding regions) longer than 50 bp among the 15 *de novo* assemblies (Table S5).
13 All-against-all alignments were made among these novel segments (See Method)
14 and a total of 63 Mbp non-redundant novel sequences were identified, with 47% (30
15 Mbp) present in two or more accessions and 53% (33 Mbp) being specific to a single
16 accession (Figure 3A).

17 Size curves for both pan- and core-genomes were obtained by adding one
18 genome to the population pool at a time (Figure 3B). For this analysis, only the 13
19 “ingroup” accessions out of the total 16 were used, excluding the three distinct sub-
20 species accessions (HM340, HM324, HM022). The core-genome size curve drops

1 quickly at first, flattening once 5 accessions are added, though still slightly negative
2 in slope even at the point where all 13 have been added. Approximately 250 Mbp
3 sequences are shared among the 13 “ingroup” accessions representing conserved
4 regions that presumably play core functions in all *M. truncatula* (Figure 3A). Another
5 ~180 Mbp is missing from at least one accession (i.e., “dispensable”), reflecting the
6 dynamic nature of genome content and prevalence of InDels and other SVs (Figure
7 3B). The corresponding pan-genome size curve sees steady increases each time a
8 new genome is added, approaching 430 Mbp when all 13 accessions have been
9 added. Indeed, fitting the observed pan-genome curve using a asymptotic regression
10 model led to estimates for the total pan-genome size of 431 Mbp and a core-
11 genome of 256 Mbp for *M. truncatula*.

12 To understand the effect of sequence variation on gene families, we annotated
13 all *de novo* assemblies and systematically identified orthologous relationships for
14 each gene among the 13 ingroup accessions – *i.e.*, the entire collection of ortholog
15 groups in the population. We placed a total of 607k non-TE genes (44k to 47k per
16 accession) into 75k ortholog groups based on sequence similarity. On average each
17 ortholog group contained 8.1 protein sequences coming from six different
18 accessions (see Methods, Figure 4). In addition to the 37k reference (Mt4.0 /
19 HM101) ortholog groups, this analysis resulted in another 38k ortholog groups with
20 no HM101 members. We identified a substantial number (25k) of accession-specific

1 genes that were only observed in a single accession, 25.7k ortholog groups shared
2 by 2-12 accessions, and 24k more shared among all 13 (Figure 4). Accession-specific
3 ortholog groups numbered as few as 1,500 specific to accession HM060 and as many
4 as 3,000 specific to HM101.

5

6 **Variation in different gene families results from differing mechanisms**

7 Several different diversity measures were estimated for different gene families
8 (Figure 5; Figure S4 A-D). The θ_π statistic, large effect SNP change, and mean protein
9 pairwise distance are metrics that provide insights into the rates of evolution for
10 different gene families, while the coefficient of variation (C.V.) of ortholog groups
11 tracks the level of copy number variation (orthology vs paralogy). The gene families
12 we examined exhibit distinctly different patterns of variation compared with the
13 genome as a whole and among themselves (Figure 5; Figure S4). NBS-LRRs are in
14 every aspect like TEs, showing the highest SNP diversity (θ_π), most frequent large-
15 effect SNP changes (premature stop codon, start codon lost, stop codon lost and
16 splice site changes), highest mean pairwise protein distance (a proxy for all protein
17 structural variants), enrichment in accession-specific gene content, and highest
18 ortholog group size coefficient of variation (CNV) (Figure 5; Figure S4). LRRs and
19 HSPs show intermediate levels of SNP diversity and pairwise protein distance, but
20 are frequently affected by large effect SNP changes and even higher CNV (Figure 5;

1 Figure S4). RLKs, F-box proteins and NCRs all show elevated levels of certain diversity
2 measures, but are much less diverse than NBS-LRRs, LRRs or HSPs. Interestingly,
3 protein kinases show high CNV despite low levels of SNP diversity and pairwise
4 protein distance. Differences in variant architecture among gene families are
5 illustrated in Figure 6, where the percent sequence similarity between the reference
6 gene model and its syntenic orthologs in the other 15 accessions is shown for three
7 example protein families (Zinc-Finger, NCRs and NBS-LRRs). Both the NCR and NBS-
8 LRR protein families are clearly more variable than Zinc-Fingers, but NBS-LRRs
9 exhibit more orthologs with significant sequence dissimilarities (structural variants,
10 red color) as well as higher numbers of CNVs (white regions corresponding to
11 missing orthologs).

12 We further examined these gene families to estimate their contribution to
13 accession-specific ortholog groups (Figure S5). Most striking were TEs, 49.2% of
14 which were accession-specific compared with just 8.3% in the core set of ortholog
15 groups (6.0x). Likewise, LRRs (50.2% accession-specific, 10.4% core; 4.8x), NBS-LRRs
16 (45.3% accession-specific versus 10.7% core; 4.3x), HSP70s (41.2% accession-specific
17 versus 19.3% core; 2.1x) and protein kinases (43.6% accession-specific versus 23.4%
18 core; 1.9x) were all over-represented in terms of accession-specific ortholog groups.
19 By contrast, NCRs (23.8% accession-specific versus 34.1% core; 0.7x), F-box proteins
20 (17.6% accession-specific versus 44.5% core; 0.4x) and RLKs (23.4% accession-

1 specific versus 60% core; 0.4x) (Figure S5) all showed lower rates of representation
2 in the accession-specific portion of the genome.

3

4 **Discussion**

5

6 **Synteny analysis based on *de novo* assemblies effectively discovers SNPs, small**

7 **InDels and large SVs**

8 Exploring plant genome variation increasingly involves the sequencing of
9 multiple accessions within a species. Early efforts simply aligned short reads against
10 a reference to discover SNPs and short indels (so-called “read-mapping approach”).
11 This includes our own earlier surveys of *M. truncatula* variation [8,9] as well as
12 similar studies in *Arabidopsis*, maize, soybean, rice and others [39–45]. In these
13 previous analyses, variation in very divergent or repetitive regions, as well as larger
14 and more complex types of variation would typically have been overlooked. Recent
15 studies have turned to *de novo* genome assembly combined with synteny
16 comparison as a basis for exploring genome variation. In *Arabidopsis*, sequencing
17 and assembling multiple genomes led to the discovery of 14.9 Mb Col-0 sequences
18 missing in at least one other accession along with unprecedented proteome diversity
19 [46]. In soybean, comparison of multiple wild relatives against the reference found

1 that 20% of the genome and 51.4% of gene families were dispensable and also
2 identified hundreds of lineage-specific genes as well as genes exhibiting CNVs as
3 potential targets of selection [47]. Sequencing three divergent rice strains revealed
4 several megabases of novel sequences specific to one strain [30]. In the present
5 study, we deeply re-sequenced 15 *M. truncatula* accessions and used the ALLPATHS-
6 LG algorithm to create high quality assemblies followed by synteny comparison as a
7 basis for global variant discovery. The resulting genome assemblies had scaffold
8 N50s >250 kb and synteny coverage >92% of the *M. truncatula* reference Mt4.0.
9 Synteny-based estimates of θ_w (Watterson's estimator of population mutation rate)
10 suggests the level of diversity is 30% higher than original read-mapping published
11 estimates (Table S4) [8]. Looking at θ_π (i.e., average number of nucleotide
12 differences per site between two accessions), the underestimate is 70%, though this
13 could be due, in part, to a more complete reference, deeper sequencing of the
14 accessions used in this study, and/or population structure among the selected
15 accessions. Examination of the syntenic blocks enabled extensive, high confidence
16 discovery of SVs, including most large indels, CNVs and translocations. These SVs
17 affect 7-22% of the alignable genome space for each *Medicago* accession, with large
18 indels spanning as much as 30 Mbp per accession and CNVs affecting as much as 85
19 Mbp (out of a genome ~450 Mbp in total size). The values reported here provide a
20 better estimate of genomic diversity within *M. truncatula*, allowing for divergent

1 genomic regions to be assessed accurately and helping to resolve repetitive and
2 variable genomic regions and gene families.

3

4 **The *Medicago* pan-genome largely resembles that of other analyzed plant species**

5 *De novo* sequencing of multiple accessions enabled us to construct a draft pan-
6 genome for *M. truncatula*, indicating a core genome of ~250 Mbp and a dispensable
7 genome of ~180 Mbp (Figure 3B). Annotation of the *Medicago de novo* genomes
8 followed by clustering using OrthoMCL resulted in a core set of 24,000 (non-TE)
9 ortholog groups present in all *M. truncatula* accessions sequenced and another
10 50,700 (67% of the total) that are dispensable (Figure 4). As *de novo* genomes were
11 added during the pan-genome analysis, the rate of increase declined quickly, with
12 both the pan-genome and core-genome curves nearly flat with the last genome
13 added. Limited novel sequence discovery would therefore be expected with the
14 addition of further accession genomes. Indeed, our estimation suggests an
15 asymptotic pan-genome size of 431 Mbp and core-genome of 256 Mbp (Figure 3).
16 Similar trends have been observed in pan-genomic analyses of seven *de novo*
17 *Glycine soja* genome [47], ten *Brassica oleracea* genomes [48], as well as a pan-
18 transcriptome analysis 503 maize accessions [49], results that together suggest
19 higher plant pan-genomes may generally be restricted in size. The finding that 67%
20 of *Medicago* ortholog groups are dispensable is likewise comparable to earlier

1 estimates of 51% in the *G. soja* analysis mentioned above [47], 73% in a study of five
2 *Oryza* AA genomes [50], and 83% of the representative transcript assemblies (RTAs)
3 in the pan-transcriptome analysis of maize [49]. All these values are higher,
4 however, than an estimate of just ~20% dispensable gene families observed in the
5 study of the *B. oleracea* pan-genome, an observation that might be attributable to
6 their focus on cultivated genotypes [48].

7 Important caveats should be kept in mind when interpreting these results. Due
8 to the incompleteness of the *de novo Medicago* assemblies (*i.e.*, certain portions of
9 genome were difficult to assemble), sequences present in one assembly but absent
10 in others could have been due to technical artifact. This would have resulted in
11 overestimates of dispensable genome size. By contrast, the pan-genome size
12 estimate should be more robust since it surveys novel sequences across all
13 accessions – and it is much less likely that a given genome region would be missed in
14 all assemblies.

15

16 **Differences in variant architecture among different gene families**

17 Genome regions high in SVs often coincide with genome regions rich in
18 either TEs or one of the biotic interaction and stress related gene families examined
19 in this study (Figures 1 and 2). This is a relationship that has frequently been
20 observed in plant genomes [30,46–48,50], but in our study, we were especially

1 interested in the range and type of SVs found in different gene families (Figure 5,
2 Supplemental Figure S4A-D). NBS-LRRs are the most variable and the most like TEs in
3 their variant structure. Both NBS-LRRs and TEs exhibit frequent large-effect SNP
4 changes, very high levels of protein diversity (mean protein distance), enrichment in
5 the accession-specific gene content, and high levels of CNVs (C.V. of gene copy
6 number). While LRRs and HSPs only exhibit intermediate levels of SNP diversity and
7 protein diversity, they are frequently affected by large effect SNP changes and even
8 higher levels of CNV. Like NBS-LRRs, these two gene families are over-represented in
9 accession-specific gene content. By contrast, protein kinases show notably low SNP
10 and protein diversity together with high levels of CNVs and over-representation in
11 accession-specific content. Finally, RLKs, F-box proteins, and NCRs are all much less
12 diverse than the other families studied here. Not surprisingly, they are also under-
13 represented in terms of accession-specific gene content. Some of these differences
14 make sense when considering the genome features of different gene families. For
15 example, NBS-LRRs have long been known to include a large proportion of
16 pseudogenes [51], a feature thought to result from the value of maintaining a
17 reservoir of genetic diversity against future pathogen pressure. Consequently, very
18 high levels of large-effect SNPs are to be expected. Likewise, NBS-LRRs are large,
19 multi-module proteins, so high levels of protein diversity, often involving domain
20 swapping, should be common [10,13–15]. By contrast, NCR genes, which are just as

1 numerous and comparably clustered in the *M. truncatula* genome, code for
2 expressed, short, single peptide, modular proteins [24,25,51]. Not surprisingly, NCRs
3 are quite low in large effect SNPs.

4

5 **Limitations remain in *de novo* assemblies based on short read sequencing**
6 **technology**

7 Even with very deep re-sequencing and *de novo* assembly using the ALLPATHS-
8 LG algorithm, important limitations remain. The contig N50 for most assemblies was
9 only 20 kb and any of the thousands of sequencing gap potentially represents a
10 missing SV. We also lacked the ability to discover SVs in regions without synteny to
11 the Mt4.0 reference. Altogether, these missing regions account for 4-8% of the
12 genome space for each *Medicago* accession. Moreover, gaps remaining in the Mt4.0
13 reference reduce its effectiveness as a framework for SV discovery. These factors all
14 presumably result in missed SV calls. Nevertheless, the SVs we did predict could
15 largely be validated. By comparing SVs discovered in the ALLPATHS assemblies of
16 three *M. truncatula* accessions to (a minimum of five) long uninterrupted reads
17 coming from PacBio sequencing, we confirmed 88-94% of SV predictions from our
18 synteny analysis. As more PacBio and other long read technologies are used to
19 resequence and assemble genomes, fewer gaps will remain and analyses of SVs,
20 dynamic gene families, and pan-genomes will become more complete and accurate.

1 **Conclusions**

2 Analysis of multiple *M. truncatula* genomes illustrates the value of *de novo*
3 assemblies to discover and describe structural variation, something that is often
4 under-estimated when using read-mapping approaches. Comparisons among the *de*
5 *nov*o assemblies also indicate that different large gene families differ in the
6 architecture of their structural variation.

7

8

9 **Methods**

10

11 **Plant material**

12 Fifteen *M. truncatula* accessions from geographically distinct populations
13 (Figure S1) broadly spanning the entire *Medicago* range were chosen for deep
14 sequencing and *de novo* assembly. These accessions were chosen for both biological
15 interest and to facilitate evaluation of assemblies. In particular, three accessions
16 were selected from the A17 clade, nine were selected from the France-Italy clade,
17 and three were selected from more distantly related clades [52]. While most
18 analyses were done on all 16 accessions including the reference HM101, some
19 statistics sensitive to population structure were derived from a subset of 13

1 accessions (three distant accessions were excluded), which we refer to as “ingroup”
2 accessions. Each accession was self-fertilized for three or more generations before
3 growing seedlings for DNA extraction. Cloning and sequencing grade DNA was
4 extracted from a pool of ~30 day old dark-grown seedlings by Amplicon Express
5 (Pullman, WA) through Ultra Clean BAC Clone Preparation followed by a CTAB liquid
6 DNA preparation [53].

7

8 **Sequencing and genome assembly**

9 Library preparation, sequencing and assembly were performed at the National
10 Center for Genome Resources (NCGR) in Santa Fe, NM. DNA sequencing was
11 performed using Illumina HiSeq 2000 instruments. For each accession, one Short
12 Insert Paired End (SIPE) library and 1 - 2 Long Insert Paired End (LIPE) libraries were
13 created following the ALLPATHS-LG assembler [31]. The SIPE library consisted of
14 fragments of ~300 nucleotides (180 nucleotides plus adapters) while LIPE libraries
15 consisted of either a 5 kb Illumina or 9 kb Nextera library. The ALLPATHS-LG
16 assembly algorithm (version 49962) [31] was run on a linux server with default
17 parameters to complete the assemblies.

18

1 **Functional annotation**

2 AUGUSTUS [33] was used to make *ab initio* gene predictions for each assembly
3 using both RNA-Seq expression evidence and *M. truncatula* HM101 reference
4 sequence (Mt4.0) [7] homology evidence. RNA-Seq data came from transcript
5 sequencing of four diverse accessions, HM034, HM056, HM101 and HM340. Reads
6 from HM034, HM056 and HM340 were directly mapped to their *de novo* assemblies
7 using Tophat [54] to generate intron hints for AUGUSTUS. For the remaining 12
8 accessions, RNA-Seq reads from the closest available accession were mapped to the
9 corresponding assembly to generate intron hints. Predicted protein sequences were
10 scanned for PFAM domains (Pfam-A.hmm) [55] using HMMER [56] and processed
11 using custom scripts. Domain categories were then assigned according to the most
12 significant Pfam hits. Among the resulting Pfam domains, 160 were associated with
13 transposable elements and grouped into a large "TE" category. NBS-LRR and RLK
14 genes were scanned using sub-family alignments from previous work [57] with 37
15 NBS-LRR sub-family identifiers (TNL0100-TNL0850, CNL0100-CNL1600) and 35 RLK
16 sub-family identifiers (LRR_I-LRR_XIII, RLCK_I-RLCK_XI) created in consistent with
17 previous research. NCRs and the broader CRP super-family were annotated by
18 running the SPADA pipeline [58] with group identifiers exactly following previous
19 literature [23]: sub-family CRP0000-CRP1030 representing defending-like genes

1 (DEFLs), CRP1040-CRP1530 representing NCRs, and CRP1600-CRP6250 representing
2 other types of CRPs.

3

4 **Flow cytometry genome size estimates for *Medicago* accessions**

5 Nine accessions (HM004, HM005, HM006, HM029, HM030, HM034, HM056,
6 HM101 and HM324) were examined for cytological genome size. Seeds of known
7 size standards were also obtained from Dolezel [59]. Seedlings were grown in
8 chambers under identical light and humidity conditions, then leaf nuclei were
9 prepared following the procedure of [59] and analyzed on a BD FACS-Calibur flow
10 cytometer at the Bio-Design Institute, Arizona State University. Mean DNA content
11 was based on 15,000 nuclei, with peak means identified using Cell-Quest software
12 (Becton Dickson). Each plant accession was sampled 3 or more times on different
13 days. Correlation analysis was then done between these cytological estimates of
14 genome size and assembled genome sizes to make Figure S2.

15

16 **Comparative genomics analysis**

17 Each *de novo* assembly was first aligned to the HM101 reference (i.e., Mt4.0)
18 using BLAT [60]. Unaligned sequences (query sequences with no hit to the
19 reference) were extracted and aligned a second time because BLAT tended to over-

1 extend gap length when it encountered stretches of 'N's (i.e., assembly gap) in the
2 target sequence. The resulting alignments were merged, fixed (removing non-
3 syntenic or overlapping alignment blocks), and cleaned (removing alignment blocks
4 containing assembly gaps). BLAT Chain/Net tools were then used to obtain a single
5 coverage best alignment net in the target genome (HM101) as well as a reciprocal-
6 best alignment net between genomes. Finally, genome-wide synteny blocks were
7 built for each *de novo* assembly (against HM101), enabling downstream analyses
8 including variant calling, novel sequence identification, and ortholog detection.

9 Based on synteny blocks generated, we identified SNPs, short InDels (alignment
10 gaps \leq 50 bases), and different types of SVs including large deletions, insertions,
11 translocations and copy number gains and losses. SVs were identified in a rigorous
12 syntenic anchoring approach: scaffolds were first aligned to and anchored on the
13 HM101 reference genome, genome-wide synteny blocks were then built for each *de*
14 *novo* assembly (against HM101). SVs were then called only in these well-built
15 synteny blocks, with each SV (insertion, deletion or translocation) receiving support
16 from both flanking sequence alignments. Variants, including large SVs, from the 15
17 accessions were merged to a single VCF file using Bcftools [61]. Since variants were
18 called independently in different accessions, the merging process resulted in missing
19 data for any variant/accession combinations where the variant was not called in that
20 accession. Custom scripts were run to impute "reference genotype" for these

1 variant/accession combinations whenever the underlying synteny alignment
2 supports the non-variant (i.e., reference) allele call. We then partitioned the
3 reference genome into 1-Mbp sliding windows to calculate gene density, TE density,
4 selected gene family density, as well as pairwise nucleotide diversity (θ_{π}) for SNPs,
5 short InDels and SVs within each window.

6

7 **Pan-genome construction and identification of accession-specific genes**

8 Based on pairwise genome comparison of each *de novo* assembly against the
9 reference (HM101), we obtained a raw set of novel sequences (present in *de novo*
10 assembly but absent in HM101) by subtracting all aligned regions from the gap-
11 removed assembly. Low-complexity sequences and short tandem repeats were
12 scanned and removed using Dustmasker and Tandem Repeat Finder [62,63].
13 Potential contaminant sequences (best hit in non-plant species) were filtered by
14 BLASTing [64] against NCBI Nucleotide (nr/nt) database. Contamination removal was
15 done after pairwise comparison with the HM101 reference based on the logic that
16 everything that aligns to HM101 should be of plant origin and free of contaminant,
17 so it was only necessary to scan the sequences that do not align to HM101 - i.e.,
18 novel sequences. Novel sequences (longer than 50 bp) from 12 accessions (13
19 “ingroup” accessions excluding HM101) were pooled and aligned using Para-Mugsy
20 [65]. The resulting alignments were parsed to determine how each segment was

1 shared among accessions – private to one accession or shared by multiple. We then
2 constructed a pan-genome that included the HM101 reference as backbone plus all
3 non-redundant novel segments identified in the other accessions. We further
4 derived genome size curves by adding one *de novo* assembly to the pool at a time
5 and calculating the size of shared genomic regions (core-genome) and the size of
6 total non-redundant sequences (pan-genome). The pan- and core-genome size size
7 curves were fitted using the asymptotic regression model $y = b_0 + b_1*(1-\exp(-$
8 $\exp(\text{Irc}) * x))$ [66]. The model was fitted using means.

9 Accession-specific genomic segments were extracted from Para-Mugsy
10 alignments mentioned above. Genes with more than 50% CDS locating in these
11 regions were selected to make the accession-specific gene set. Pfam analysis and
12 functional enrichment were then performed on this accession-specific gene list.

13

14 **Protein ortholog group analysis and comparisons**

15 Protein sequences from all 16 accessions (1,028,566 total genes) were pooled to
16 construct ortholog groups using OrthoMCL [67]. This resulted in 150k ortholog
17 groups with an average of 6 genes per group. Further analysis only focused on non-
18 TE genes in 13 “ingroup” accessions since the three distant accessions (HM340,
19 HM324, HM022) tend to introduce extra ortholog group due to high divergence.
20 Ortholog groups could contain from 0 to any number of protein sequences from any

1 one accession. A total of 607k non-TE genes from 13 ingroup accessions were
2 grouped into 75k ortholog groups. Grouping of protein sequences was based on
3 BlastP significance so the actual sequence similarities within groups vary – but
4 typically above 70% identity threshold (i.e., pairwise protein distance less than 0.3).
5 On average, each ortholog group contains 8.1 protein sequences, but from only 6.7
6 different accessions. For each group a functional category was assigned based on
7 Pfam annotation of all group members. Ortholog groups were also binned based on
8 the number of accessions contributing to them: from 1 (accession-specific) to 13
9 (present in all ingroup accessions, i.e., “core” ortholog groups).

10

11 **Diversity of different gene families**

12 SNPs were called based on pairwise genome comparisons of each accession
13 against HM101. SNP-based nucleotide diversity (θ_π) was estimated for coding
14 regions of each gene and the distribution of θ_π for different gene families was
15 obtained. To account for poorly covered regions, only genes where $\geq 80\%$ of the CDS
16 regions were covered in at least 10 out of the 13 accessions were retained.
17 Functional effects of SNPs in genic regions were determined using snpEff [68], and
18 the proportion of genes with large effect SNP changes (e.g., gain or loss of stop
19 codon) in each gene family was calculated.

1 In addition to SNPs, we identified a large number of small InDels and large SVs
2 inside/overlapping genic regions. Since these types of variants often lead to frame-
3 shift, splice-site change, exon skipping, domain swapping or other gene structural
4 changes, we decided to use protein sequence distance as a measure to quantify the
5 functional impact of SVs. Since the OrthoMCL-defined ortholog groups do not
6 explicitly define one-to-one orthologous relationship among accessions, we used
7 synteny alignment information and derived a smaller set of syntenic ortholog groups
8 with one-to-one relationship among accessions. Filtering was done requiring
9 syntenic orthologs be present in ≥ 10 accessions (i.e., missing data in less than 3
10 accessions) for each group. We then did multiple-sequence alignment for each
11 syntenic ortholog group, calculated mean pairwise protein distance (MPPD), and
12 characterized the distribution of MPPDs for different gene family categories (Pfam
13 domains).

14 To assess the level of copy number variation (CNV) for different gene families,
15 we grouped protein sequences from 13 accessions into ortholog groups using
16 OrthoMCL (see previous section). Pfam category of each ortholog group was
17 assigned by the most abundant category among group members. Members in each
18 ortholog group were treated as copies of a common ancestor, thus enabling
19 quantification of gene copy number variation among accessions. In practice, we

1 calculated the coefficient of variation (C.V.) of gene copy number among accessions
2 for each ortholog group and summarized its distribution for different gene families.

3

4 **Validation of SVs using PacBio long reads**

5 We performed PacBio sequencing on three accessions (HM034, HM056 and
6 HM340) to validate the breakpoints of identified structural variants. Each accession
7 was sequenced to 14-20 fold coverage using either P4C2 or P5C3 chemistry. The
8 average read length was 4-7 Kbp. PacBio reads were first mapped to the
9 corresponding ALLPATHS assembly using BLASR [69]. For each SV, the number of
10 PacBio reads fully spanning ± 500 bp of the breakpoints were counted. We consider
11 an SV to be “validated” only if each of its breakpoints received at least five such
12 PacBio reads support.

13

14

1 **Declarations**

2 **Acknowledgments**

3 We thank the University of Minnesota Supercomputing Institute for computational
4 infrastructure, storage, and systems administrative support.

5

6 **Funding**

7 This research was funded by National Science Foundation Grant 1237993 to NDY, PT,
8 KATS, RMS, JRM and JM. The funding source was not involved in the design of the
9 study, collection, analysis, interpretation of data, or in the writing of the manuscript.

10

11 **Availability of data and materials**

12 Illumina and PacBio reads data from this article can be found in the NCBI Sequence
13 Read Archive (SRA) under accession number PRJNA256006. RNA-Seq reads can be
14 found under SRA accession number SRP077692. Genome assembly sequences, SNP
15 genotype files are available for download from the *Medicago* Hapmap project
16 website (<http://www.medicagohapmap.org/downloads/assemblies>).

1

2 **Author Contributions**

3 Conceived and designed experiments: KATS, RMS, PT, JRM, JM, NDY. Performed
4 experiments: PZ, TR, RD, JL, KPS. Analyzed data: PZ, KATS, TR, JG, JL, ADF, KPS, RMS,
5 PT, JRM, JM, NDY. Wrote paper: PZ, KATS, NDY. Collected and processed sequence
6 data: PZ, KATS, TR, JG, JL, ADF, RMS, JRM, JM.

7

8 **Competing interests**

9 The authors declare that they have no competing interests.

10

11 **Consent for publication**

12 Not applicable.

13

14 **Ethics approval and consent to participate**

15 *Medicago* germplasm resources (seed) were obtained and used, with permission,
16 from Jean-Marie Prospero at Unité mixte de recherche / Amélioration génétique et
17 adaptation des plantes méditerranéennes et tropicales (UMR-AGAP) at INRA-
18 Montpellier, France.

19

20

1 **Additional Files**

2 **Additional file 1:** Supplementary figures (Figure S1-S5) described in the manuscript.

3

4 **Additional file 2:** Supplementary tables (Table S1-S7) described in the manuscript.

5

6 **Additional file 3:** Supporting data file S1 (Excel spreadsheet listing the member
7 counts of different gene families including all NBS-LRR, NCR, RLK and TE subfamilies,
8 that are predicted in 15 *de novo* assemblies).

9

10 **References**

11 1. Graham PH. Legumes: Importance and Constraints to Greater Use. *Plant Physiol.*
12 2003;131:872–7.

13 2. Lavin M, Herendeen P, Wojciechowski M. Evolutionary Rates Analysis of
14 Leguminosae Implicates a Rapid Diversification of Lineages during the Tertiary. *Syst.*
15 *Biol.* 2005;54:575–94.

16 3. Young ND, Udvardi M. Translating *Medicago truncatula* genomics to crop
17 legumes. *Curr. Opin. Plant Biol.* 2009;12:193–201.

18 4. Ronfort J, Bataillon T, Santoni S, Delalande M, David JL, Prospero J-M.
19 Microsatellite diversity and broad scale geographic structure in a model legume:
20 building a set of nested core collection for studying naturally occurring variation in

- 1 Medicago truncatula. BMC Plant Biol. 2006;6:28.
- 2 5. Tadege M, Wen J, He J, Tu H, Kwak Y, Eschstruth A, et al. Large-scale insertional
3 mutagenesis using the Tnt1 retrotransposon in the model legume Medicago
4 truncatula. Plant J. 2008;54:335–47.
- 5 6. Oldroyd GED, Downie JA. Coordinating Nodule Morphogenesis with Rhizobial
6 Infection in Legumes. Annu. Rev. Plant Biol. 2008;59:519–46.
- 7 7. Tang H, Krishnakumar V, Bidwell S, Rosen B, Chan A, Zhou S, et al. An improved
8 genome release (version Mt4.0) for the model legume Medicago truncatula. BMC
9 Genomics. 2014;15:312.
- 10 8. Branca A, Paape TD, Zhou P, Briskine R, Farmer a. D, Mudge J, et al. Whole-
11 genome nucleotide diversity, recombination, and linkage disequilibrium in the
12 model legume Medicago truncatula. Proc. Natl. Acad. Sci. 2011;108:E864–70.
- 13 9. Stanton-Geddes J, Paape T, Epstein B, Briskine R, Yoder J, Mudge J, et al.
14 Candidate Genes and Genetic Architecture of Symbiotic and Agronomic Traits
15 Revealed by Whole-Genome, Sequence-Based Association Genetics in Medicago
16 truncatula. PLoS One. 2013;8:e65688.
- 17 10. Meyers BC. Genome-Wide Analysis of NBS-LRR-Encoding Genes in Arabidopsis.
18 Plant Cell. 2003;15:809–34.
- 19 11. Shiu S-H, Bleecker a B. Receptor-like kinases from Arabidopsis form a
20 monophyletic gene family related to animal receptor kinases. Proc. Natl. Acad. Sci.

- 1 2001;98:10763–8.
- 2 12. Kuroda H. Classification and Expression Analysis of Arabidopsis F-Box-Containing
3 Protein Genes. *Plant Cell Physiol.* 2002;43:1073–85.
- 4 13. Michelmore RW, Meyers BC. Clusters of resistance genes in plants evolve by
5 divergent selection and a birth-and-death process. *Genome Res.* 1998;8:1113–30.
- 6 14. Richly E, Kurth J, Leister D. Mode of Amplification and Reorganization of
7 Resistance Genes During Recent Arabidopsis thaliana Evolution. *Mol. Biol. Evol.*
8 2002;19:76–84.
- 9 15. Leister D. Tandem and segmental gene duplication and recombination in the
10 evolution of plant disease resistance genes. *Trends Genet.* 2004;20:116–22.
- 11 16. Shiu S-H, Bleecker AB. Expansion of the Receptor-Like Kinase/Pelle Gene Family
12 and Receptor-Like Proteins in Arabidopsis. *Plant Physiol.* 2003;132:530–43.
- 13 17. Xu G, Ma H, Nei M, Kong H. Evolution of F-box genes in plants: Different modes
14 of sequence divergence and their relationships with functional diversification. *Proc.*
15 *Natl. Acad. Sci.* 2009;106:835–40.
- 16 18. Meyers BC, Dickerman AW, Michelmore RW, Sivaramakrishnan S, Sobral BW,
17 Young ND. Plant disease resistance genes encode members of an ancient and
18 diverse protein family within the nucleotide-binding superfamily. *Plant J.*
19 1999;20:317–32.
- 20 19. Padmanabhan M, Cournoyer P, Dinesh-Kumar SP. The leucine-rich repeat

- 1 domain in plant innate immunity: A wealth of possibilities. *Cell. Microbiol.* 2009. p.
2 191–8.
- 3 20. Sung D, Kaplan F, Guy CL. Plant Hsp70 molecular chaperones: protein structure ,
4 gene family , expression and function. *Physiol. Plant.* 2001;113:443–51.
- 5 21. Graham M a. Computational Identification and Characterization of Novel Genes
6 from Legumes. *Plant Physiol.* 2004;135:1179–97.
- 7 22. Silverstein KAT. Genome Organization of More Than 300 Defensin-Like Genes in
8 *Arabidopsis*. *Plant Physiol.* 2005;138:600–10.
- 9 23. Silverstein K a. T, Moskal W a., Wu HC, Underwood B a., Graham M a., Town CD,
10 et al. Small cysteine-rich peptides resembling antimicrobial peptides have been
11 under-predicted in plants. *Plant J.* 2007;51:262–80.
- 12 24. Alunni B, Kevei Z, Redondo-Nieto M, Kondorosi A, Mergaert P, Kondorosi E.
13 Genomic Organization and Evolutionary Insights on GRP and NCR Genes, Two Large
14 Nodule-Specific Gene Families in *Medicago truncatula*. *Mol. Plant-Microbe Interact.*
15 2007;20:1138–48.
- 16 25. Mergaert P. A Novel Family in *Medicago truncatula* Consisting of More Than 300
17 Nodule-Specific Genes Coding for Small, Secreted Polypeptides with Conserved
18 Cysteine Motifs. *Plant Physiol.* 2003;132:161–73.
- 19 26. Farkas A, Maroti G, Durg H, Gyorgypal Z, Lima RM, Medzihradzsky KF, et al.
20 *Medicago truncatula* symbiotic peptide NCR247 contributes to bacteroid

- 1 differentiation through multiple mechanisms. Proc. Natl. Acad. Sci. 2014;111:5183–
2 8.
- 3 27. Young ND, Zhou P, Silverstein KAT. Exploring structural variants in
4 environmentally sensitive gene families. Curr. Opin. Plant Biol. 2016. p. 19–24.
- 5 28. Clark RM, Schweikert G, Toomajian C, Ossowski S, Zeller G, Shinn P, et al.
6 Common Sequence Polymorphisms Shaping Genetic Diversity in *Arabidopsis*
7 *thaliana*. Science. 2007;317:338–42.
- 8 29. Cao J, Schneeberger K, Ossowski S, Günther T, Bender S, Fitz J, et al. Whole-
9 genome sequencing of multiple *Arabidopsis thaliana* populations. Nat. Genet.
10 2011;43:956–63.
- 11 30. Schatz MC, Maron LG, Stein JC, Wences A, Gurtowski J, Biggers E, et al. Whole
12 genome de novo assemblies of three divergent strains of rice, *Oryza sativa*,
13 document novel gene space of aus and indica. Genome Biol. 2014;15:506.
- 14 31. Gnerre S, MacCallum I, Przybylski D, Ribeiro FJ, Burton JN, Walker BJ, et al. High-
15 quality draft assemblies of mammalian genomes from massively parallel sequence
16 data. Proc. Natl. Acad. Sci. 2011;108:1513–8.
- 17 32. Smit A, Hubley R, Green P. RepeatMasker Open-4.0. 2013-2015 .
18 <http://www.repeatmasker.org>. 2013.
- 19 33. Stanke M, Stanke M, Waack S, Waack S. Gene prediction with a hidden Markov
20 model and a new intron submodel. Bioinformatics. 2003. p. 0.

- 1 34. Nielsen R. Molecular Signatures of Natural Selection. *Annu. Rev. Genet.*
2 2005;39:197–218.
- 3 35. Conrad DF, Pinto D, Redon R, Feuk L, Gokcumen O, Zhang Y, et al. Origins and
4 functional impact of copy number variation in the human genome. *Nature.*
5 2010;464:704–12.
- 6 36. Redon R, Ishikawa S, Fitch KR, Feuk L, Perry GH, Andrews TD, et al. Global
7 variation in copy number in the human genome. *Nature.* 2006;444:444–54.
- 8 37. Ye K, Schulz MH, Long Q, Apweiler R, Ning Z. Pindel: A pattern growth approach
9 to detect break points of large deletions and medium sized insertions from paired-
10 end short reads. *Bioinformatics.* 2009;25:2865–71.
- 11 38. Li S, Li R, Li H, Lu J, Li Y, Bolund L, et al. SOAPindel: Efficient identification of
12 indels from short paired reads. *Genome Res.* 2013;23:195–200.
- 13 39. Gore M a., Chia J-M, Elshire RJ, Sun Q, Ersoz ES, Hurwitz BL, et al. A First-
14 Generation Haplotype Map of Maize. *Science.* 2009;326:1115–7.
- 15 40. Lam H-M, Xu X, Liu X, Chen W, Yang G, Wong F-L, et al. Resequencing of 31 wild
16 and cultivated soybean genomes identifies patterns of genetic diversity and
17 selection. *Nat. Genet.* 2010;42:1053–9.
- 18 41. Zheng L-Y, Guo X-S, He B, Sun L-J, Peng Y, Dong S-S, et al. Genome-wide patterns
19 of genetic variation in sweet and grain sorghum (*Sorghum bicolor*). *Genome Biol.*
20 2011;12:R114.

- 1 42. Chia J-M, Song C, Bradbury PJ, Costich D, de Leon N, Doebley J, et al. Maize
2 HapMap2 identifies extant variation from a genome in flux. *Nat. Genet.*
3 2012;44:803–7.
- 4 43. Huang X, Kurata N, Wei X, Wang Z-XX, Wang A, Zhao Q, et al. A map of rice
5 genome variation reveals the origin of cultivated rice. *Nature.* 2012;490:497–501.
- 6 44. Gordon SP, Priest H, Des Marais DL, Schackwitz W, Figueroa M, Martin J, et al.
7 Genome diversity in *Brachypodium distachyon*: deep sequencing of highly diverse
8 inbred lines. *Plant J.* 2014;79:361–74.
- 9 45. Xu X, Liu X, Ge S, Jensen JDDJJ, Hu F, Li X, et al. Resequencing 50 accessions of
10 cultivated and wild rice yields markers for identifying agronomically important
11 genes. *Nat. Biotechnol.* 2012;30:105–11.
- 12 46. Gan X, Stegle O, Behr J, Steffen JG, Drewe P, Hildebrand KL, et al. Multiple
13 reference genomes and transcriptomes for *Arabidopsis thaliana*. *Nature.*
14 2011;477:419–23.
- 15 47. Li Y, Zhou G, Ma J, Jiang W, Jin L, Zhang Z, et al. De novo assembly of soybean
16 wild relatives for pan-genome analysis of diversity and agronomic traits. *Nat.*
17 *Biotechnol.* 2014;32:1045–52.
- 18 48. Golicz AA, Bayer PE, Barker GC, Edger PP, Kim H, Martinez PA, et al. The
19 pangenome of an agronomically important crop plant *Brassica oleracea*. *Nat.*
20 *Commun.* 2016;7:13390.

- 1 49. Hirsch CN, Foerster JM, Johnson JM, Sekhon RS, Muttoni G, Vaillancourt B, et al.
2 Insights into the Maize Pan-Genome and Pan-Transcriptome. *Plant Cell*.
3 2014;26:121–35.
- 4 50. Zhang Q-J, Zhu T, Xia E-H, Shi C, Liu Y-L, Zhang Y, et al. Rapid diversification of five
5 *Oryza* AA genomes associated with rice adaptation. *Proc. Natl. Acad. Sci*.
6 2014;111:E4954–62.
- 7 51. Marone D, Russo M, Laidò G, De Leonardis A, Mastrangelo A. Plant Nucleotide
8 Binding Site–Leucine-Rich Repeat (NBS-LRR) Genes: Active Guardians in Host
9 Defense Responses. *Int. J. Mol. Sci*. 2013;14:7302–26.
- 10 52. Yoder JB, Briskine R, Mudge J, Farmer a., Paape T, Steele K, et al. Phylogenetic
11 Signal Variation in the Genomes of *Medicago* (Fabaceae). *Syst. Biol*. 2013;62:424–38.
- 12 53. Murray MG, Thompson WF. Rapid isolation of high molecular weight plant DNA.
13 *Nucleic Acids Res*. 1980;8:4321–6.
- 14 54. Trapnell C, Pachter L, Salzberg SL. TopHat: discovering splice junctions with RNA-
15 Seq. *Bioinformatics*. 2009;25:1105–11.
- 16 55. Finn RD, Bateman A, Clements J, Coggill P, Eberhardt RY, Eddy SR, et al. Pfam: the
17 protein families database. *Nucleic Acids Res*. 2014;42:D222–30.
- 18 56. Eddy SR. Accelerated Profile HMM Searches. *PLoS Comput. Biol*.
19 2011;7:e1002195.
- 20 57. Ameline-Torregrosa C, Wang B-B, O’Bleness MS, Deshpande S, Zhu H, Roe B, et

- 1 al. Identification and Characterization of Nucleotide-Binding Site-Leucine-Rich
2 Repeat Genes in the Model Plant *Medicago truncatula*. *Plant Physiol.* 2007;146:5–
3 21.
- 4 58. Zhou P, Silverstein KA, Gao L, Walton JD, Nallu S, Guhlin J, et al. Detecting small
5 plant peptides using SPADA (Small Peptide Alignment Discovery Application). *BMC*
6 *Bioinformatics.* 2013;14:335.
- 7 59. DOLEZEL J. Plant DNA Flow Cytometry and Estimation of Nuclear Genome Size.
8 *Ann. Bot.* 2005;95:99–110.
- 9 60. Kent WJ. BLAT - The BLAST-like alignment tool. *Genome Res.* 2002;12:656–64.
- 10 61. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence
11 Alignment/Map format and SAMtools. *Bioinformatics.* 2009;25:2078–9.
- 12 62. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST
13 plus : architecture and applications. *BMC Bioinformatics.* 2009;10.
- 14 63. Benson G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic*
15 *Acids Res.* 1999;27:573–80.
- 16 64. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search
17 tool. *J. Mol. Biol.* 1990;215:403–10.
- 18 65. Angiuoli S V., Salzberg SL. Mugsy: Fast multiple alignment of closely related
19 whole genomes. *Bioinformatics.* 2011;27:334–42.
- 20 66. R Development Core Team. R: A Language and Environment for Statistical

1 Computing. R Found. Stat. Comput. Vienna Austria. 2016;0:{ISBN} 3-900051-07-0.

2 67. Li L. OrthoMCL: Identification of Ortholog Groups for Eukaryotic Genomes.

3 Genome Res. 2003;13:2178–89.

4 68. Cingolani P, Platts A, Wang LL, Coon M, Nguyen T, Wang L, et al. A program for

5 annotating and predicting the effects of single nucleotide polymorphisms, SnpEff:

6 SNPs in the genome of *Drosophila melanogaster* strain w 1118; iso-2; iso-3. Fly.

7 2012;6:80–92.

8 69. Chaisson MJ, Tesler G. Mapping single molecule sequencing reads using basic

9 local alignment with successive refinement (BLASR): application and theory. BMC

10 Bioinformatics. 2012;13:238.

11

12 **List of Abbreviations**

13 **AAs:** amino acids **CDS:** Coding sequence **CNVs:** Copy number variants **CRPs:**

14 Cysteine-rich peptides **HSPs:** Heat shock proteins **LIPE:** Long insert paired end **LRR:**

15 Leucine-rich repeat **Mbp:** Million base pairs **NBS-LRR:** Nucleotide-binding site

16 leucine-rich repeat **NCRs:** nodule-specific cysteine-rich peptides **RLKs:** Receptor-like

17 kinases **SIPE:** Short insert paired end **SNPs:** Single nucleotide polymorphisms **SVs:**

18 Structural variants **TEs:** Transposable elements **VCF:** Variant call format

19

1 **Figures**

2 **Figure 1. Heatmap showing percent covered by synteny alignment for each 1Mb**
3 **window in 15 *de novo* *M. truncatula* assemblies (Upper 15 tracks), reference gap**
4 **position ('Gaps'), percent bases covered by synteny blocks in at least 10 out 13**
5 **accessions ('Coverage'), nucleotide diversity ($\theta\pi$) for SNPs ('Pi_SNP'), short InDels**
6 **(< 50bp, 'Pi_InDel') and large SVs (\geq 50bp, 'Pi_SV'), as well as gene density of**
7 **different categories (TE, NBS-LRR, RLK, NCR, LRR and F-boxes).**

8 Nucleotide diversity ($\theta\pi$) estimates were calculated using only 13 "ingroup" *M.*
9 *truncatula* accessions.

10

11 **Figure 2. Zoom-in view of five 1-Mb regions (A-E) selected from Figure 1.**

12 Upper 15 tracks show percentage covered by synteny alignment for each 50kb
13 window (column) in 15 *M. truncatula* assemblies. Bottom tracks show reference gap
14 position ('Gaps'), percent bases covered by synteny blocks in at least 10 out 13
15 accessions ('Coverage'), nucleotide diversity ($\theta\pi$) for SNPs ('Pi_SNP'), short InDels (<
16 50bp, 'Pi_InDel') and large SVs (\geq 50bp, 'Pi_SV'), as well as gene density of different
17 categories (TE, NBS-LRR, RLK, NCR, LRR and F-boxes) in relative scale (minimum to
18 maximum spaced equally in grayscale within each panel) with grey columns
19 representing missing data due to lack of synteny coverage. Starting position for

1 each region is provided at the bottom (*e.g.*, chr7:28Mb, indicating that a 1 Mb
2 region beginning at position 28,000,001 on chromosome 7 is displayed).

3

4

5 **Figure 3. Sharing status of the *Medicago* pan-genome (A) and the pan-genome size
6 curve (B).**

7

8 **Figure 4. Sharing status of *Medicago* protein ortholog groups.**

9

10 **Figure 5. Diversity estimates of different gene families: (A) SNP-based nucleotide
11 diversity (i.e., θ_π), (B) proportion members affected by different types of large-
12 effect SNPs, (C) mean pairwise protein distance for syntenic ortholog groups and
13 (D) coefficient of variation (CV) of gene copy number in each ortholog group (i.e.,
14 an estimate of copy number variation) among accessions.**

15 Numbers in parenthesis reflect: (A) & (B) number of genes where $\geq 80\%$ of the CDS
16 regions were covered in at least 10 out of the 13 accessions; (C) number of syntenic
17 ortholog groups where syntenic orthologs were present in ≥ 10 accessions (i.e.,
18 missing data in less than 3 accessions); (D) number of OrthoMCL-defined ortholog
19 groups based entirely on protein sequence similarity.

20

1 **Figure 6. Sequence similarity of selected gene families in 15 *Medicago* accessions:**
2 **(A) Zinc-Finger domain, (B) NCRs and (C) NBS-LRRs.**

3 Each cells in the score matrix indicates percent sequence similarity (1-100) between
4 an HM101 gene and its syntenic ortholog from one of the 15 accessions. Blank
5 (white) cells indicate missing data.

6

7 **Additional files**

8 Additional file 1: Supplementary figures (Figure S1-S5) described in the manuscript.

9

10 Additional file 2: Supplementary tables (Table S1-S7) described in the manuscript.

11

12 Additional file 3: Excel spreadsheet listing the member counts of different gene
13 families including all NBS-LRR, NCR, RLK and TE subfamilies, that are predicted in 15
14 *de novo* assemblies.

15