

# ***PHASIS*: A computational suite for *de novo* discovery and characterization of phased, siRNA-generating loci and their miRNA triggers**

Atul Kakrana<sup>1,2</sup>, Pingchuan Li<sup>2,3</sup>, Parth Patel<sup>1,2</sup>, Reza Hammond<sup>1,2</sup>, Deepti Anand<sup>4</sup>, Sandra M. Mathioni<sup>5</sup>, Blake C. Meyers<sup>5,6\*</sup>

<sup>1</sup>Center for Bioinformatics and Computational Biology, University of Delaware, Newark, DE 19714, USA

<sup>2</sup>Delaware Biotechnology Institute, University of Delaware Newark, DE 19714, USA

<sup>3</sup>Morden Research and Development Centre, Agriculture and Agri-Food Canada, Morden, Canada

<sup>4</sup>Department of Biological Sciences, University of Delaware, Newark, DE, 19716, USA

<sup>5</sup>Donald Danforth Plant Science Center, St. Louis, MO 63132, USA

<sup>6</sup>University of Missouri – Columbia, Division of Plant Sciences, 52 Agriculture Lab, Columbia, MO 65211, USA.

\*To whom correspondence should be addressed; Email: [bmeyers@danforthcenter.org](mailto:bmeyers@danforthcenter.org)

**Keywords:** phased siRNAs, phasiRNA, tasiRNA, piRNA, small RNA.

## 1 **Abstract**

2 Phased, secondary siRNAs (phasiRNAs) are found widely in plants, from protein-coding  
3 transcripts and long, non-coding RNAs; animal piRNAs are also phased. Integrated methods  
4 characterizing “*PHAS*” loci are unavailable, and existing methods are quite limited and  
5 inefficient in handling large volumes of sequencing data. The *PHASIS* suite described here  
6 provides complete tools for the computational characterization of *PHAS* loci, with an emphasis  
7 on plants, in which these loci are numerous. Benchmarked comparisons demonstrate that  
8 *PHASIS* is sensitive, highly scalable and fast. Importantly, *PHASIS* eliminates the requirement of  
9 a sequenced genome and PARE/degradome data for discovery of phasiRNAs and their miRNA  
10 triggers.

11

## 12 **Background**

13 Phased siRNAs (phasiRNAs) are a major subclass of secondary siRNAs, found extensively in  
14 plants [1]. The defining characteristic of phasiRNAs is the DCL-catalyzed processing of double-  
15 stranded RNA (dsRNA) precursors, starting from a precisely delimited 5' terminus and  
16 generating regularly-spaced 21- or 24-nt populations of siRNAs [2]. PhasiRNAs can be further  
17 subdivided into three main categories based on their precursor mRNAs and spatiotemporal  
18 patterns of accumulation: i) The first phasiRNAs identified, so-called *trans*-acting siRNAs  
19 (tasiRNAs) generated from a small set of long, non-coding mRNAs (lncRNAs) referred to as *TAS*  
20 genes [3–5]; ii) phasiRNAs from protein-coding transcripts, such as *NB-LRRs* or *PPRs* [6]; and iii)  
21 two classes, 21-nt premeiotic or 24-nt meiotic phasiRNAs, highly enriched in reproductive  
22 tissues and also produced from lncRNAs, reported in grasses, but with no as-yet reported  
23 targets [2,7]. Thus, the umbrella name of “phasiRNAs” refers simply to their biogenesis and not  
24 their function (unlike the subset of tasiRNAs) because many phasiRNAs lack validated targets,  
25 either in *cis* or *trans* [6,8].

26

27 The biogenesis of phasiRNAs in plants is dependent on a triggering mechanism that sets the  
28 phase of the resulting secondary siRNAs, generated from a specific nucleotide in the mRNA  
29 precursor. To date, the only described type of trigger is a miRNA, and a breakthrough in our

30 understanding of plant miRNA function came with the observation that all or nearly all 22-nt  
31 miRNAs trigger phasiRNA biogenesis from their targets [9,10]. These miRNA triggers function  
32 via the ARGONAUTE (AGO) proteins into which they are loaded, and since phasiRNA biogenesis  
33 requires both SGS3 and RDR6 [3,4], there may be interactions between these proteins,  
34 ultimately recruiting DCL4 or DCL5. SGS3 and RDR6 proteins function in the cytoplasm, forming  
35 siRNA bodies [11]. Recent work has identified membrane-bound polysomes in the rough ER as  
36 the site where miRNA triggers of phasiRNAs accumulate, leading to phasiRNA biogenesis [12].  
37 miRNA triggers are thus an important component in the analysis of plant phasiRNAs, and the  
38 identification of specific triggers with specific *PHAS* targets is an integral part of phasiRNA  
39 analysis.

40  
41 Since the discovery of phasiRNAs in 2005, *TAS* genes have been characterized in detail,  
42 especially the eight loci in Arabidopsis but these represent only a small fraction of the *PHAS*  
43 repertoire found in many plant genomes. In other eudicot genomes, there are hundreds of  
44 protein-coding genes that are targeted by diverse miRNAs, many of which are lineage-specific  
45 [8,13–15]. Grass genomes contain even more *PHAS* loci. For example, loci yielding reproductive  
46 phasiRNAs number in the hundreds to thousands in maize [7] and rice [16], and have yet to be  
47 characterized broadly in monocots or other lineages outside of the grasses. These include the  
48 premeiotic *PHAS* loci that are targeted by miR2118 family members, triggering production of  
49 21-nt phasiRNAs accumulating in early anther development, and the 24-*PHAS* loci that are  
50 targeted by miR2275 family members, triggering production of 24-nt phasiRNAs, accumulating  
51 in anthers during meiosis [7]. Analysis of the spruce genome, a gymnosperm that speciated  
52 ~325 million years before the evolution of monocots and eudicots, identified over 2000 *PHAS*  
53 loci mostly from protein-coding genes, including over 750 *NB-LRRs* [17]. Thus, plant *PHAS* loci  
54 are widely prevalent and highly variable from genome to genome both in the total number and  
55 in terms of the types of loci that generate them. Characterization of *PHAS* loci from each  
56 sequenced plant genome will provide insights into this unusual type of post-transcriptional  
57 control, its evolution, and diversification.

58

59 Tools for the *de novo* identification of *PHAS* genes (or loci) to date have required an assembled  
60 genome for their discovery and additional experimental data such as PARE [18], degradome  
61 [19] or GMUCT [20] libraries to identify their miRNA triggers. Integrated tools for discovery and  
62 in-depth characterization of *PHAS* genes have not yet been developed, and the existing options  
63 are both limited in number and function. These algorithmic limitations and bioinformatic gaps  
64 along with the increasing depth and volume of sequencing data necessitates scalable, fast and  
65 advanced methods to study this relatively new class of secondary siRNAs for which parallels  
66 exist between plants and animals [2,13,21,22]. Motivated by this need for software, the  
67 prospect of discovering novel phasiRNA modules, the emerging importance of phasiRNAs, and  
68 the explosion in the number of plant species that are being investigated for small RNAs (sRNAs),  
69 we developed a new computational suite that we call “*PHASIS*”. The name “*PHASIS*” is from the  
70 ancient Greek city of Phasis, a destination for Jason and the Argonauts according to Greek  
71 mythology; we selected the name as it links the colloquialism “phasis” as short for phasiRNAs,  
72 with the Argonaut proteins that bind them. This set of tools facilitates the discovery,  
73 quantification, annotation, comparison of *PHAS* loci (and precursors) and identification of their  
74 miRNA triggers, from a few to hundreds of sRNA libraries in a single run. *PHASIS* not only  
75 addresses crucial bioinformatic gaps while providing an integrated and flexible workflow for the  
76 comprehensive study of *PHAS* loci, but it is also fast and sensitive.

77

## 78 **Results**

### 79 *Assessment and benchmarking*

80 We first sought to assess the sensitivity and specificity for *PHASIS*; ideally, this would be done  
81 with a gold-standard reference set of experimentally-validated *PHAS* loci in plants. While the  
82 definition of “gold standard” is as-yet unclear for *PHAS* loci, the recently-described maize loci  
83 are among the most exhaustively characterized [7], and thus we used these data below. We  
84 also compared *PHASIS* predictions and performance with PhaseTank [23]. Currently, two  
85 computational tools are capable of *de novo* discovery of *PHAS* loci – PhaseTank [23] and  
86 ShortStack [24]. *PhaseTank* is exclusively built for predicting *PHAS* loci in plants, while  
87 *ShortStack* aims to annotate and quantify diverse sRNA-associated genes (or clusters), and it’s

88 typically deployed for characterizing miRNAs in plants and animals [24]. A direct comparison  
89 between *PHASIS* and *ShortStack* is not possible due to significant differences in their scope,  
90 utility and workflow (**Table 1**). So, for comparative benchmarking, we chose *PhaseTank*, mainly  
91 because of matching objectives and its published superiority over *ShortStack* in predicting *PHAS*  
92 loci [23]. Benchmarking was performed across five plant species – *Arabidopsis thaliana*  
93 (*Arabidopsis*), *Brachypodium distachyon* (*Brachypodium*), *Oryza sativa* (rice), *Zea mays* (maize)  
94 and *Lilium maculatum* (*Lilium*). These species were selected based on availability of high-quality  
95 nuclear genome assemblies or anther transcriptomes (in case of *Lilium* – generated for a  
96 different study but included here), and deep sRNA libraries from premeiotic and meiotic anther  
97 or from at least one of these two stages that should contain many reproductive phasiRNAs  
98 (**Supplementary Table 1**). *Arabidopsis* was included because it was originally used in *PhaseTank*  
99 benchmarking [23]. For *PhaseTank*, the reference genome, transcriptome and sRNA libraries  
100 were converted to the appropriate formats, and the time for file conversion process, although  
101 complex and lengthy, was not added in the *PhaseTank* runtimes. *PHASIS* and *PhaseTank* use  
102 inherently different scoring schemas; because of this difference, we used a conservative p-value  
103 (1e-05) for *PHASIS* and the recommended score (i.e. 15) for *PhaseTank*. All benchmarks were  
104 performed on a 28 core, 2.42 GHz machine with 512 GB of RAM, running CentOS 6.6.

105

#### 106 *PHAS prediction and runtime performance*

107 We first compared *PHAS* loci and transcript predictions from *PHASIS* and *PhaseTank*. Since  
108 *Arabidopsis* lacks 24-*PHAS* loci (none have ever been published, nor have we found any), and  
109 the genome encodes just eight *TAS* genes, these were excluded from quantification of  
110 prediction and speed comparisons. *PHASIS* demonstrated an edge over *PhaseTank* in *PHAS*  
111 predictions: in genomic analyses, it predicted up to 2.5 times more *PHAS* loci, ranging from 73  
112 24-*PHAS* (145% gain) to 380 21-*PHAS* (24% gain) loci in *Brachypodium* and rice respectively  
113 (**Table 2**). The biggest gain was observed in an analysis of the *Lilium* transcriptome, in which  
114 *PHASIS* predicted ~10 times (n=408) more 21-*PHAS* and 18 times (n=9065) more 24-*PHAS*  
115 precursor transcripts compared to *PhaseTank* (**Figure 2**). The specific data format requirements  
116 of *PhaseTank* made it difficult to accurately determine the set of common *PHAS* predictions

117 (the ‘common *PHAS* pool’, hereafter) for transcriptome level analysis, however, by matching  
118 the sequences we determined that *PHASIS* captured at least 66% of 21-*PHAS* and 99% of 24-  
119 *PHAS* predictions from *PhaseTank*. For genomic analyses, *PHASIS* captured >80% of *PhaseTank*  
120 predictions, except in rice and Arabidopsis in which *PhaseTank* predicted additional 24-*PHAS*  
121 loci (**Table 2**).

122  
123 The additional 24-*PHAS* loci predicted by *PhaseTank* in rice and Arabidopsis all had significantly  
124 lower quality scores (from *PhaseTank*) compared to the common *PHAS* pool, as did the  
125 *PhaseTank*-exclusive 21- and 24-*PHAS* predictions from other species. The average quality  
126 scores computed for each species were 1.7 to 7.8 times lower compared to the common *PHAS*  
127 pool (p-value < 0.001, t-test) (**Supplementary Table 2**); thus, the predictions exclusive to  
128 *PhaseTank* are likely unphased and a misinterpretation of loci yielding profuse heterochromatic  
129 siRNAs (hc-siRNAs). This may explain the 24-*PHAS* predictions in Arabidopsis by *PhaseTank*  
130 (**Figure 2b and Table 2**), as 24-nt phasiRNAs have not been reported in Arabidopsis despite  
131 exhaustive analyses [1]. Nonetheless, considering that these *PhaseTank* predictions could  
132 represent weak *PHAS* loci, we attempted to capture them by running *PHASIS* at lower p-value  
133 cutoff ( $1e-03$ ) but still failed to detect >96% of them. Manual investigation of a portion of these  
134 *PHAS* loci using our custom sRNA browser, which uses a slightly different *PHAS* scoring schema  
135 [25], revealed that these are indeed either unphased or show typical characteristics of hc-siRNA  
136 loci such as similarity to transposons, and we concluded that these are false positives predicted  
137 by *PhaseTank* (**Supplementary Figure 1A**). However, we could detect 70% (n=67) of the total  
138 24-*PHAS* *PhaseTank* predictions in rice at the lower p-value cutoff ( $1e-03$ ) of *PHASIS*, and a  
139 majority of these showed weak phasing patterns (**Supplementary Figure 1B**), suggesting that  
140 *PHASIS* missed these at the selected cutoff. However, the count of 24-*PHAS* loci predicted in  
141 rice by both tools in these libraries from a recent study [16], was lower than earlier estimates  
142 [2], indicating that the libraries likely missed meiotic peak of accumulation. These contrasting  
143 observations – Arabidopsis, in which *PHASIS* correctly excluded 24-*PHAS* predictions even at  
144 relaxed cutoff, versus rice, in which it correctly captured 70% of weakly phased 24-*PHAS* loci –  
145 highlights differences in scoring in the two tools, with the default *PHASIS* p-value cutoff ( $1e-05$ )

146 more stringent than that of *PhaseTank* (score=15). Using a lower p-value cutoff for *PHASIS*  
147 could further increase the gain in *PHAS* predictions over *PhaseTank* without adding much noise.

148

149 We manually investigated 21- and 24-*PHAS* predictions that are exclusive to *PHASIS*, using our  
150 public, custom genome browser (<https://mpss.danforthcenter.org/>). The majority of these  
151 displayed characteristics matching those of the canonical 21- and 24-*PHAS* loci reported in  
152 maize [7] (**Supplementary Figure 2**). Moreover, a major proportion of these *PHASIS*-exclusive  
153 predictions had PARE-validated miRNA triggers, matching to the earlier reports from maize, rice  
154 and *Brachypodium* [2,7,13]. Next, we compared prediction runtimes of *PHASIS* and *PhaseTank*  
155 from genome- and transcriptome-level experiments. To get the correct runtimes for both tools,  
156 we excluded the execution time for a common step performed by an external tool (Bowtie,  
157 version 1) that prepares the index for the reference genome or transcriptome. For genome-  
158 level experiments, *PHASIS* displayed a minimum speed gain of 3x in *Arabidopsis* and rice and a  
159 maximum speed gain of 7x in maize (**Figure 3**). In transcriptome-level experiments, both tools  
160 took almost equal time (**Figure 3**). However, *PHASIS* yielded 10x (n=408) to 17x (n=9065) more  
161 *PHAS* predictions for 21- and 24-*PHAS* loci, respectively (**Table 2 and Supplementary Figure 3**),  
162 compared to *PhaseTank*, which means that *PHASIS* processed a high number of *PHAS*  
163 transcripts in the same runtime. Moreover, the time and effort required to convert the  
164 reference genome as well as the sRNA libraries to meet *PhaseTank* input requirements were  
165 not included in these runtime comparisons. Lastly, it should be noted that *PHASIS* takes  
166 significantly less time for any subsequent analyses in these species because of its unique ability  
167 to systematically store ancillary data in the first run, check data integrity and compatibility with  
168 parameters for subsequent runs, and avoid redoing the slowest steps, such as reference pre-  
169 processing, index preparation, etc.

170

#### 171 *Comparison of PHASIS predictions with manually-curated data*

172 We next wanted to address how well the predictions from *PHASIS* compare with a set of  
173 manually-curated *PHAS* loci. We and collaborators curated a set of 21- and 24-*PHAS* (n= 463  
174 and 176 loci from precisely-staged, premeiotic and meiotic maize anthers [7]. This curated set

175 was prepared by first combining all libraries from the sampled premeiotic and meiotic stages  
176 into a single file, followed by genome wide scans to identify phasiRNA generating loci using a  
177 score-based approach [5] and finally curating each *PHAS* locus to exclude those that overlap  
178 with repeat-associated regions or display sRNA distribution atypical of hc-siRNA generating loci  
179 [7]. *PHASIS* processes each library separately mainly to a) detect phased patterns independently  
180 in at least one of the input sRNA libraries, b) minimize any noise that could be added by  
181 combining sRNAs from multiple stages, tissues or treatments, and c) infer the correct 5'-end of  
182 *PHAS* loci by collating data from different libraries. Therefore, unlike the original analysis, we  
183 did not combine the 32 libraries (see **Supplementary Table 1**) for predictions by *PHASIS*.  
184 Furthermore, to emulate 'real world' conditions in which *PHASIS* would be used by non-experts,  
185 we did not provide a confidence cutoff - i.e. *PHASIS* was run in the default mode. Of the  
186 manually-curated 463 21-*PHAS* and 176 24-*PHAS* loci, *PHASIS* captured 89.0% (n=411) and  
187 85.79% (n=151) (**Supplementary Table 5**). The majority of those missed either lacked  
188 continuous phased positions or had a very low abundance across all sRNA libraries, and some  
189 had a single sRNA read accounting for the major proportion (>90%) of the abundance at the  
190 *PHAS* locus. The average abundance of siRNAs in the 'missed' 21- and 24-*PHAS* set was ~12- and  
191 252-times lower compared to the common pool ( $p < 1.02e-09$ ), supporting the observation that  
192 those missed by *PHASIS* were weakly phased loci; a portion of these could be captured with a  
193 relaxed cutoff. Nonetheless, these results demonstrate that *PHASIS* predictions are largely  
194 consistent with the manually-curated data, and for most studies, the use of *PHASIS* may  
195 ameliorate the need to manually curate *PHAS* locus predictions, an otherwise complex and  
196 cumbersome task especially when *PHAS* loci number in the hundreds to thousands, as reported  
197 in many plant genomes [2,7,8,13,14,17].

198

#### 199 *Trigger prediction and runtime performance*

200 The identification of the miRNA triggers of *PHAS* loci is important for understanding their  
201 potential roles, classification and for discovery of secondary siRNA cascades. In addition, a set  
202 of *PHAS* loci or transcripts when combined with the trigger identity may serve as a gold-  
203 standard reference set for downstream experimental and bioinformatics studies. Given the

204 importance of trigger identification, we compared the trigger prediction performance of *PHASIS*  
205 in ‘validation’ mode with *PhaseTank*. The *PHASIS* ‘validation’ mode will identify triggers for  
206 *PHAS* loci or transcripts using experimental data such as PARE, degradome or GMUCT libraries  
207 (‘PARE’, henceforth) [18–20]. *PhaseTank* by default predicts triggers in ‘validation’ mode, i.e.  
208 experimental data is required. Since *PHASIS* predicted more *PHAS* loci compared to *PhaseTank*,  
209 the number of *PHAS* loci (and transcripts) with the predicted triggers by *PHASIS* was higher too.  
210 So, for a fair comparison, we used only the common pool of *PHAS* loci to evaluate the trigger  
211 prediction performances. *PHASIS* displayed a gain of up to 76.0% in predicted triggers, except  
212 for 21-*PHAS* loci in Arabidopsis (**Figure 2A and B**), with a minimum accuracy of 96.0% for 24-  
213 *PHAS* maize loci and maximum accuracy of 99.5% in Brachypodium 21-*PHAS* loci  
214 (**Supplementary Table 3**). This accuracy was computed as the proportion of triggers (out of the  
215 total) that match to known triggers of phasiRNAs and tasiRNAs described in earlier studies  
216 [2,5,6,8,25–27]. These estimates of accuracy are likely conservative, given that there might be a  
217 few new and unknown triggers that we counted as false positives in our accuracy  
218 computations. We excluded rice 24-*PHAS* loci from our comparisons because both tools failed  
219 to report triggers for these loci, likely due to sRNA libraries that were not precisely staged  
220 relative to the accumulation of 24-nt phasiRNAs and thereby making it difficult to capture the 5’  
221 and 3’ ends of *PHAS* loci – information crucial to the identification of correct triggers. Liliun 21-  
222 and 24-*PHAS* transcripts were also excluded from the comparisons because of a lack of PARE  
223 data from the corresponding anther stages, data required by *PhaseTank* to predict triggers.  
224 Likewise, Arabidopsis 24-*PHAS* couldn’t be included in our comparison as *PhaseTank* predicted  
225 loci (n=146) were false positives, and there were no overlapping loci with *PHASIS*.

226  
227 We noticed a decline in number of predicted triggers by *PHASIS* for 21-*PHAS* loci in Arabidopsis,  
228 compared to those predicted by *PhaseTank* (**Figure 2A**). This decline in predicted triggers was  
229 traced to seven phased loci corresponding to the pentatricopeptide repeat (PPR) gene family,  
230 with phasiRNAs triggered by miR161. We found trigger sites predicted by *PhaseTank* for five of  
231 these loci, located 214 nt to 310 nt from the first or last phased cycle of the *PHAS* loci, towards  
232 their middle (**Supplementary Table 4**). Since *phastrigs*, the trigger discovery tool of *PHASIS*, is

233 built with the aim to eliminate the need for experimental data and because trigger sites are  
234 expected to overlap with 5' or 3' ends of the phased region, it uses a narrow search space at  
235 the 5' and 3' ends to search for triggers. Hence, these miR161 target sites were missed by  
236 *PHASIS*. In *phastrigs*, the search space to identify triggers is defined by the number of phased  
237 positions (*PHAS*-index) on either side of 5' and 3' ends of phased regions, and by default *PHAS*-  
238 index is set to  $\pm 3$  positions for both ends. The *PHAS*-index setting to expand or the narrow  
239 search space for triggers is user tunable and can be adjusted to capture such cases.  
240 Nonetheless, these 21-*PHAS* loci from Arabidopsis support our estimates that trigger  
241 identification by *phastrigs* is conservative, and relaxing the *phastrigs* search parameters could  
242 further increase the gain in predicted triggers compared to *PhaseTank*.

243

#### 244 *Identifying PHAS triggers without additional experimental data*

245 We next evaluated the performance of *PHASIS* in trigger 'prediction' mode by comparing it with  
246 *PhaseTank* and *PHASIS* in the 'validation' mode. We define *PHASIS* 'prediction' mode as an  
247 analysis to predict triggers for *PHAS* loci or transcripts without any supporting experimental  
248 data such as PARE, degradome or GMUCT libraries. Liliium was excluded from the comparison of  
249 predicted triggers due to the lack of PARE data, which is compulsory for *PhaseTank* to predict  
250 triggers and required by *PHASIS* in 'validation' mode. Also, for reasons mentioned above, 24-  
251 *PHAS* loci from Arabidopsis and rice were excluded from the comparisons. *PHASIS* displayed a  
252 minimum gain of 40.3% and maximum gain of 178.3% over *PhaseTank* in predicting triggers for  
253 21-*PHAS* and 24-*PHAS* loci from Brachypodium, respectively (**Table 2** and **Figure 2**). The gain in  
254 the number of triggers ranged from a minimum of 35 for maize 24-*PHAS* loci to a maximum of  
255 611 for rice 21-*PHAS* loci. In addition to the gain in trigger prediction, *PHASIS* also displayed  
256 significant accuracy in prediction mode, with a minimum accuracy of 89.9% in predicting  
257 triggers for 24-*PHAS* loci from maize and maximum accuracy of 99.9% in the case of Liliium 24-  
258 *PHAS* precursor transcripts, however, with an exception for Liliium 21-*PHAS* triggers. The  
259 accuracy of predicted triggers of Liliium 21-*PHAS* loci was significantly lower (43.9%) compared  
260 to the other species (**Supplementary Table 2**). For Liliium, we used miRNAs from well-  
261 characterized monocots like rice and maize because a complete set of miRNAs were not

262 available due to the absence of a sequenced genome. Surprisingly, we found that for *Lilium* 21-  
263 *PHAS* transcripts a majority of triggers corresponded to miR2275 instead of miR2118; this  
264 observation was puzzling because miR2275 is known to trigger 24-nt phasiRNAs in the grasses  
265 [2,7,13], and this was the basis for the low recorded accuracy in predicting *Lilium* 21-*PHAS*  
266 triggers. We did not further investigate the miR2275-triggered 21-*PHAS* transcripts. We also  
267 noticed that the proportion of 21- and 24-*PHAS* precursors for which triggers could be  
268 identified in *Lilium*, 18.1% and 25.9% respectively (**Table 2 and Figure 2**), was substantially  
269 lower compared to the overall average of 73.8% in other species for which genomic analysis  
270 was performed. Plant *PHAS* precursor transcripts are typically cleaved by the miRNA trigger,  
271 converted to dsRNA by an RNA-dependent RNA polymerase, and then successively diced by a  
272 Dicer enzyme. Since no data on transcriptional rate, stability and half-life of phasiRNA  
273 precursors are available, we speculated that a portion of the *Lilium* *PHAS* precursor transcripts  
274 were shortened by processing from the 5' end, removing the trigger target sites. Identifying  
275 triggers from such "processed" precursor transcripts is not possible because the P1 site  
276 corresponding to the first phasiRNA (at the 5' terminus) could be missing from the transcript. In  
277 addition, the presence of already-processed mRNAs will confound the *de novo* assembly of  
278 precursor transcripts from short-reads.

279  
280 To test whether the low yield of triggers by *phastrigs* resulted from our use of processed  
281 precursor transcripts and not a technical shortcoming of *PHASIS*, we generated Single Molecule  
282 Real Time (SMRT) PacBio sequencing data from *Lilium* anthers 4 mm to 6 mm in length. These  
283 sizes represented premeiotic and meiotic stages of anther development (**see supplementary**  
284 **methods**) and were selected based on the availability of the samples. Capturing *PHAS*  
285 precursors is complex, not just because these are targets of miRNAs presumably rapidly  
286 processed by a Dicer, but reproductive phasiRNAs are ephemeral in development and thus not  
287 easily captured [7]. SMRT-seq produced 425,897 full-length transcripts for 176,373 unique  
288 isoforms, which were pre-processed to generate 122,779 high quality (polished) transcripts.  
289 This set had 5,131 unique proteins covered by more than 80% protein length, relative to the  
290 Uniprot protein-sequence resource, thereby suggesting a reasonable assembly of the anther

291 transcriptome. *PHASIS* identified 87 21-*PHAS* and 175 24-*PHAS* precursor transcripts. This low  
292 yield of *PHAS* transcripts was expected, though not to such a degree, because of the  
293 combination of the following: a) low read counts for SMRT-seq compared to the deep RNA-seq  
294 data, b) the coverage-based error correction algorithm - ‘Quiver’ implemented in the IsoSeq  
295 protocol (SMRT Analysis software version 2.3, Pacific Biosciences) which filters out transcripts  
296 with insufficient coverage, i.e. those that cannot be confidently corrected, and c) the  
297 aforementioned processive cleavage of *PHAS* precursors by Dicer. *phastrigs* could identify  
298 triggers for only 21.8% (n=19) of 21-*PHAS* precursors, a slight increase compared to 18.1% in  
299 the RNA-seq assembly, and these triggers included miR2275, miR2118 and miR390. This low  
300 proportion of triggers detected for 21-*PHAS* could result from missing the precise stage at  
301 which 21-*PHAS* precursors accumulate in the *Lilium* samples. However, *phastrigs* could identify  
302 triggers for 54.2% of the 24-*PHAS* precursors, a significant increase over the 25.9% in the RNA-  
303 seq assembly, supporting our premise about the completeness of the *PHAS* precursor  
304 transcripts. The processed precursors were likely collapsed into the full-length or the longest  
305 transcript in SMRT-seq assembly, thereby enriching the proportion of uncleaved precursor  
306 transcripts. Hence, it should be noted that neither the precursors from neither RNA-seq nor  
307 SMRT-seq may accurately represent the true total count of *PHAS* loci in *Lilium*.

308  
309 Lastly, we compared runtimes for both tools for miRNA trigger prediction of *PHAS* loci and  
310 transcripts. *PHASIS* showed a minimum speed gain of 3.3x and a maximum speed gain of 12.6x  
311 over *PhaseTank* in ‘validation’ mode (**Figure 3**). In ‘prediction’ mode, *PHASIS* was at least 5.0x  
312 and at most 31.2x faster compared to its own ‘validation’ mode without any significant loss in  
313 accuracy (**Supplementary Table 3**). *PhaseTank* requires PARE data to predict triggers, and lacks  
314 a function equivalent to *PHASIS* ‘prediction’ mode, but since *PHASIS*, even without the  
315 additional experimental data (like PARE) displays >89.9% accuracy in trigger prediction, we  
316 decided to compare runtimes for both. *PHASIS* in ‘prediction’ mode displayed a minimum speed  
317 gain of 33.3x and a maximum gain of 104.3x for Arabidopsis 21-*PHAS* loci (**Figure 3**). The trigger  
318 predictions for 24-*PHAS* loci from Arabidopsis and rice, which displayed even higher speed  
319 gains, were excluded from the runtime comparisons due to the reasons described above. This

320 gain in *PHAS* trigger identification demonstrates the capacity of *PHASIS* to predict triggers  
321 without experimental data. This functionality will save time and the cost of preparing PARE  
322 libraries; it will also reduce the amount of sample required for phasiRNA analysis. Protocols for  
323 preparing PARE libraries require comparatively more input RNA relative to RNA-seq or sRNA-  
324 seq [28].

325

## 326 **Conclusions**

327 Loci generating 21- and 24-nt phasiRNAs are widely prevalent across land plants  
328 [2,5,8,14,16,17,29], varying in numbers per genome from tens to thousands, displaying diverse  
329 spatial and temporal expression patterns, and participating in an array of different functions  
330 [5,7,8,29,30]. Recently, piRNAs in *Drosophila* too were reported to be phased, generating  
331 ‘trailer’ piRNAs in 27-nt intervals after cleavage by secondary siRNA and Zucchini-dependent  
332 processing of cleaved transcript [21,22]. Given the wide prevalence of phasiRNAs and the rate  
333 of genome sequencing, it is likely that they will be better characterized and studied in the  
334 coming years. The existing tools for computational characterization of *PHAS* loci or transcript  
335 are limited both in number and functionality.

336

337 The *PHASIS* suite provides an integrated solution for the large-scale survey of tens to hundreds  
338 of sRNA libraries for the following applications: a) *de novo* discovery of *PHAS* loci and precursor  
339 transcripts, b) a summarization of *PHAS* loci from specific groups of sRNA libraries, c) a  
340 comparison of *PHAS* summaries between groups corresponding to samples from different  
341 stages, tissues and treatments, d) quantification and annotations of *PHAS* loci, and e) discovery  
342 of their miRNA triggers. *PHASIS* generates easily parsed output files for downstream  
343 bioinformatics analysis, formatted result files for immediate consumption and organized  
344 ancillary data to facilitate optimizations like a re-summarization to exclude or include libraries.

345

346 More complete characterization of phasiRNAs in evolutionarily diverse plant genomes will  
347 advance our understanding of phasiRNA function and the adaptation of the pathway, and it  
348 may yet discover new classes of *PHAS* genes. *PHASIS* will thus facilitate the discovery of

349 phasiRNAs and their precursors, and the identification of their triggers by eliminating the  
350 requirement of a genome assembly and experimental PARE/degradome data. *PHASIS* offers  
351 flexibility to users to tailor analyses for their own goals and it integrates an array for functions  
352 in one package.

353

## 354 **Methods**

355 *PHASIS* comprises three components that together perform *de novo* discovery, annotation,  
356 quantification, comparison and trigger identification for *PHAS* loci or precursor transcripts. We  
357 chose a modular approach over the single ‘one-command’ style for the following reasons: a) to  
358 maximize the flexibility for specific data or study requirements; b) to integrate multiple,  
359 connected analyses; and, c) to reduce overall runtime by maximizing phase- and step-specific  
360 parallelization. A description of these tools – *phasdetect*, *phasmerge*, *phastrigs* – in order of  
361 their utility or phases of study is provided below (see also **Figure 1**). *PHASIS* leverages the  
362 Python (v3) process-based “threading” interface to achieve efficient scalability and significantly  
363 reduce runtimes through parallel computing.

364

365 *phasdetect* performs *de novo* prediction of *PHAS* loci or precursor transcripts using user-  
366 supplied sRNA libraries along with a reference genome or transcriptome. It can efficiently  
367 process tens to hundreds of sRNA libraries in parallel, reducing runtimes. *phasdetect* operates  
368 via three main steps: a) first, sRNA libraries are normalized and mapped to the reference; b)  
369 second, mapped sRNA reads are scanned to identify regions rich for specific size classes, such as  
370 those generated by Dicer activity (typically 21, 22, or 24 nt in plants); and, c) finally these  
371 regions are stitched into clusters and the phasing of the small RNAs is computed as a *p-value*.  
372 We adopted a standard approach to compute *p-values* [9]. Parameters controlling these steps  
373 can be modified by users via the setting file “*phasis.set*”, including values for *phase*, *mindepth*  
374 and *clustbuffer*; these refer to the phasing periodicity, minimum sRNA abundance to be  
375 included for *p-value* computation, and the minimum distance separating two clusters. These  
376 parameters are explained in detail on the *PHASIS* wiki page  
377 (<https://github.com/atulkakrana/PHASIS/wiki/>). The output for *phasdetect* includes library-

378 specific list of *PHAS* loci (or transcripts) at several different confidence levels plus ancillary data,  
379 used to reduce runtime for subsequent analyses. For example, in case of a reanalysis after  
380 adding new libraries, *phasdetect* checks for any changes in parameters from the earlier analysis,  
381 assesses the integrity and compatibility of the ancillary data for, and reuses existing data to  
382 avoid repetition. This ancillary data also enables an array for downstream analyses and analysis-  
383 specific optimizations directly through *phasdetect*.

384  
385 *phasmerge* generates a summary, matches *PHAS* loci to annotations and performs a  
386 comparison between the *PHAS* summaries using the library-specific *PHAS* lists and ancillary  
387 data generated by *phasdetect*. These operations are selected by using the *-mode* option with  
388 the *merge* (default) or *compare* values. The *merge* mode prepares a *PHAS* summary for the  
389 libraries of interest, or for libraries that belong to different groups based on sample stages,  
390 tissues or treatments. The analysis can be tailored to meet the study requirements. For  
391 example, to maximize discovery, a user might set a lower confidence level (*p-value*) for  
392 summarization and consider all loci with a trigger predicted without the PARE data (identified  
393 through *phastrigs*) for downstream analyses. In contrast, a user motivated to maximize the  
394 quality might identify *PHAS* loci with the highest confidence level, followed by pruning of results  
395 with stringent quality parameters (described on the *phasmerge* wiki), and use *PHAS* loci that  
396 have PARE-supported triggers. *PHAS* summaries from different groups of libraries can be  
397 compared using *compare* mode. This is particularly useful to identify intersecting and exclusive  
398 *PHAS* loci between different groups of stages, tissues or treatments. In *merge* mode, if an  
399 additional annotation file is provided, then merged *PHAS* loci are matched to genome  
400 annotations so as to identify coding *PHAS* loci or other available annotations. This function also  
401 supports quick discovery of precursor transcripts for summarized *PHAS* loci when provided with  
402 a GTF file generated from mapping the transcriptome assembly to genome. Furthermore,  
403 *phasmerge* attempts to determine the correct 5' terminus of *PHAS* loci by optimizing for the  
404 best 5' or 3' coordinates based on the user's sRNA data – a crucial functionality for  
405 determination of the correct miRNA trigger. *phasmerge* benefits from the modular *PHASIS*

406 workflow, allowing users to optimize their results for the study which may vary in purpose, and  
407 making *phasmerge* independent from other tools.

408  
409 The *phasmerge* workflow has three mandatory and two optional steps: a) via *merge* mode,  
410 *phasmerge* first generates a unique list of *PHAS* loci (or transcripts) for each user-specified  
411 library, by selecting predictions with the highest available confidence score (lowest *p-value*)  
412 that pass a user-supplied *p-value* cutoff, after comparing predictions from all available  
413 confidence levels; b) *phasmerge* clusters the “best” candidate loci from specified libraries  
414 specific by the user, based on the degree of overlap in phased positions (or ‘cycles’) to select a  
415 representative locus for each cluster; finally, c) *phasmerge* computes library-specific  
416 abundances, a size-class ratio, the maximum to total phasiRNAs abundance ratio, and other  
417 quality information. Optional steps include d) *compare* mode, which first reads *PHAS* loci (or  
418 transcripts) from user-supplied summaries ( $n=2$ ) and then identifies matching *PHAS* pairs based  
419 on the overlap in phased positions, to report a combined matrix including both shared and  
420 unique loci in each *PHAS* summary file, and e) *merge* mode; when supplied with annotations,  
421 as described above, *phasmerge* matches a merged set of *PHAS* loci with genome annotations or  
422 with a genome-matched transcriptome assembly, both provided as GTF file, to report exonic or  
423 complete overlaps with annotated transcripts. This step requires prior installation of *SQLite* on  
424 user’s machine. *phasmerge* generates several reports as output, most importantly, *PHAS*  
425 summary for libraries of interest which includes quality parameters (see online wiki for more  
426 information), FASTA files for size-specific siRNAs and all the siRNAs from phased positions along  
427 with detailed information on phased clusters with phasiRNAs, positions, associated *p-values*,  
428 etc.

429  
430 *phastrigs* identifies sRNA triggers for *PHAS* loci and precursor transcripts using the *phasmerge*  
431 summaries and a user-provided list of miRNAs (or any other small RNA). It was developed with  
432 the idea to minimize the requirement of experimental PARE libraries [18–20]. However, if such  
433 data (‘PARE’, henceforth) are provided, then *phastrigs* reports sRNA triggers with experimental  
434 support; these may be of higher confidence for some downstream experimental analyses.

435 *phastrigs* uses an algorithm designed to be both fast and exhaustive. It uses *miRferno*, an  
436 exhaustive target prediction algorithm that we developed [31] to predict target sites for user-  
437 supplied miRNAs. The speed and precision of *phastrigs* is enhanced by a scan focused on the 5'  
438 terminus of each *PHAS* locus (5'-end of the first cycle, the P1 position) for the trigger site, which  
439 reduces the search space and chance of reporting false triggers. This 5' terminus is inferred at  
440 the summarization step by *phasmerge* while collating data from different sRNA libraries. In the  
441 case of *PHAS* transcripts, only the 5' terminus of the phased precursor is scanned, while in case  
442 of genomic *PHAS* loci, either the 5' or 3' end of the phased region is chosen, based on the  
443 strand targeted by a specific miRNA. *phastrigs* analysis is divided into two main steps: a) *PHAS*  
444 transcripts or genomic sequences are extracted, and targets for user-supplied miRNAs are  
445 predicted; b) next, a scan of phased positions located at the 5' or 3' termini of precursor for a  
446 target site that corresponds with the production of phasiRNAs is performed; this scan looks for  
447 target sites within  $\pm 3$  nt of the '*PHAS* index', defined as theoretical phased positions upstream  
448 from the 5' terminus of P1. If PARE data is supplied, then PARE-validated cleavage sites are used  
449 for trigger identification. The *phastrigs* report includes detailed information on miRNA-target  
450 interactions, PARE abundances at the predicted cleavage site, and the *PHAS* index of the  
451 predicted trigger site relative to the P1 position.

452

### 453 *Software*

454 The methods and algorithm described in this article, implemented as *PHASIS* suite of tools for  
455 *PHAS* discovery, are freely available from <https://github.com/atulkakrana/PHASIS>. *PHASIS* is  
456 released under the OSI Artistic License 2.0. Tools and Perl libraries required to use *PHASIS* along  
457 with the instructions to install and usage of individual tools is provided in detail in the *PHASIS*  
458 wiki (<https://github.com/atulkakrana/PHASIS/wiki/>).

459

### 460 **Abbreviations**

461 *PHASIS*: *PHAS* Inspection Suite; tasiRNAs: *trans*-acting siRNAs; PMC: pollen mother cells; PARE:  
462 Parallel Analysis of RNA Ends; SMRT: Single Molecule Real Time Sequencing; GMUCT: genome-  
463 wide mapping of uncapped and cleaved transcripts

464

## 465 **Additional files**

466 **Additional file 1:** Supplementary methods in Microsoft Word (.doc) format.

467 **Additional file 2:** Supplementary figures from S1 to S3 in Portable Document Format (PDF).

468 **Additional file 3:** Supplementary tables from S1 to S4 in Office Open XML Spreadsheet format  
469 (.xlsx)

470

## 471 **Acknowledgments**

472 We thank Bruce Kingham and Olga Shevchenko from Delaware Biotechnology Institute  
473 (Newark, DE, USA) for their help with PacBio sequencing. We thank Karol Miaskiewicz from  
474 Delaware Biotechnology Institute (Newark, DE, USA) for work on the SMRT portal and the *PB-*  
475 *Tofu* command line environment. We also thank Kun Huang from University of Delaware  
476 (Newark, DE, USA) for helpful discussions on the developmental stage and anther size  
477 correlations.

478

## 479 **Funding**

480 This work was supported by U.S. National Science Foundation Plant Genome Research Program  
481 (NSF-PGRP) grant (1649424) and University Competitive Fellow Award (2015-2016) from  
482 University of Delaware.

483

## 484 **Author contributions**

485 AK conceived the project. AK and PL designed and implemented the method, individual  
486 contributions are marked on scripts. PP, RH, DA and AK tested tools and compiled  
487 benchmarking data. AK and SM collected the *Lilium* samples. SM prepared the SMRT  
488 sequencing libraries. AK and BCM wrote the manuscript. All authors read and approved the  
489 final manuscript.

490

491 **Competing interests**

492 The authors declare that they have no competing interests.

493

494 **References**

495 1. Axtell MJ. Classification and comparison of small RNAs from plants. *Annu. Rev. Plant Biol.*  
496 2013;64:137–59.

497 2. Johnson C, Kasprzewska A, Tennessen K, Fernandes J, Nan G-L, Walbot V, et al. Clusters and  
498 superclusters of phased small RNAs in the developing inflorescence of rice. *Genome Res.*  
499 2009;19:1429–40.

500 3. Vazquez F, Vaucheret H, Rajagopalan R, Lepers C, Gascioli V, Mallory AC, et al. Endogenous  
501 trans-Acting siRNAs Regulate the Accumulation of Arabidopsis mRNAs. *Mol. Cell.* 2004;16:69–  
502 79.

503 4. Peragine A, Yoshikawa M, Wu G, Albrecht HL, Poethig RS. SGS3 and SGS2/SDE1/RDR6 are  
504 required for juvenile development and the production of trans-acting siRNAs in Arabidopsis.  
505 *Genes Dev.* 2004;18:2368–79.

506 5. Allen E, Xie Z, Gustafson AM, Carrington JC. microRNA-directed phasing during trans-acting  
507 siRNA biogenesis in plants. *Cell.* 2005;121:207–221.

508 6. Fei Q, Xia R, Meyers BC. Phased, secondary, small interfering RNAs in posttranscriptional  
509 regulatory networks. *Plant Cell.* 2013;25:2400–15.

510 7. Zhai J, Zhang H, Arikait S, Huang K, Nan G-L, Walbot V, et al. Spatiotemporally dynamic, cell-  
511 type-dependent premeiotic and meiotic phasiRNAs in maize anthers. *Proc. Natl. Acad. Sci. U. S.*  
512 *A.* 2015;112:3146–51.

513 8. Zhai J, Jeong D-H, Paoli ED, Park S, Rosen BD, Li Y, et al. MicroRNAs as master regulators of  
514 the plant NB-LRR defense gene family via the production of phased, trans-acting siRNAs. *Genes*  
515 *Dev.* 2011;25:2540–53.

516 9. Chen H-M, Chen L-T, Patel K, Li Y-H, Baulcombe DC, Wu S-H. 22-Nucleotide RNAs trigger  
517 secondary siRNA biogenesis in plants. *Proc. Natl. Acad. Sci. U. S. A.* 2010;107:15269–74.

518 10. Cuperus JT, Carbonell A, Fahlgren N, Garcia-Ruiz H, Burke RT, Takeda A, et al. Unique  
519 functionality of 22-nt miRNAs in triggering RDR6-dependent siRNA biogenesis from target  
520 transcripts in Arabidopsis. *Nat. Struct. Mol. Biol.* 2010;17:997–1003.

521 11. Jouannet V, Moreno AB, Elmayan T, Vaucheret H, Crespi MD, Maizel A. Cytoplasmic  
522 Arabidopsis AGO7 accumulates in membrane-associated siRNA bodies and is required for ta-  
523 siRNA biogenesis. *EMBO J.* 2012;31:1704–13.

- 524 12. Li S, Le B, Ma X, Li S, You C, Yu Y, et al. Biogenesis of phased siRNAs on membrane-bound  
525 polysomes in Arabidopsis. *eLife*. 2016;5:e22750.
- 526 13. Jeong DH, Schmidt SA, Rymarquis LA, Park S, Ganssmann M, German MA, et al. Parallel  
527 analysis of RNA ends enhances global investigation of microRNAs and target RNAs of  
528 *Brachypodium distachyon*. *Genome Biol*. 2013;14:R145.
- 529 14. Arikiti S, Xia R, Kakrana A, Huang K, Zhai J, Yan Z, et al. An atlas of soybean small RNAs  
530 identifies phased siRNAs from hundreds of coding genes. *Plant Cell*. 2014;26:4584–601.
- 531 15. Xia R, Ye S, Liu Z, Meyers BC, Liu Z. Novel and recently evolved microRNA clusters regulate  
532 expansive F-BOX gene networks through phased small interfering RNAs in wild diploid  
533 strawberry. *Plant Physiol*. 2015;169:594–610.
- 534 16. Fei Q, Yang L, Liang W, Zhang D, Meyers BC. Dynamic changes of small RNAs in rice spikelet  
535 development reveal specialized reproductive phasiRNA pathways. *J. Exp. Bot*. 2016;67:6037–49.
- 536 17. Xia R, Xu J, Arikiti S, Meyers BC. Extensive families of miRNAs and *PHAS* loci in Norway  
537 Spruce demonstrate the origins of complex phasiRNA networks in seed plants. *Mol. Biol. Evol*.  
538 2015; 32: 2905–2918.
- 539 18. German MA, Luo S, Schroth G, Meyers BC, Green PJ. Construction of Parallel Analysis of RNA  
540 Ends (PARE) libraries for the study of cleaved miRNA targets and the RNA degradome. *Nat*.  
541 *Protoc*. 2009;4:356–62.
- 542 19. Addo-Quaye C, Eshoo TW, Bartel DP, Axtell MJ. Endogenous siRNA and miRNA targets  
543 identified by sequencing of the Arabidopsis degradome. *Curr. Biol*. 2008;18:758–762.
- 544 20. Gregory BD, O'Malley RC, Lister R, Urich MA, Tonti-Filippini J, Chen H, et al. A link between  
545 RNA metabolism and silencing affecting Arabidopsis development. *Dev. Cell*. 2008;14:854–866.
- 546 21. Mohn F, Handler D, Brennecke J. piRNA-guided slicing specifies transcripts for Zucchini-  
547 dependent, phased piRNA biogenesis. *Science*. 2015;348:812–7.
- 548 22. Han BW, Wang W, Li C, Weng Z, Zamore PD. piRNA-guided transposon cleavage initiates  
549 Zucchini-dependent, phased piRNA production. *Science*. 2015;348:817–21.
- 550 23. Guo Q, Qu X, Jin W. PhaseTank: genome-wide computational identification of phasiRNAs  
551 and their regulatory cascades. *Bioinformatics*. 2015;31:284–6.
- 552 24. Axtell MJ. ShortStack: Comprehensive annotation and quantification of small RNA genes.  
553 *RNA*. 2013;19:740–51.
- 554 25. Allen E, Howell MD. miRNAs in the biogenesis of trans-acting siRNAs in higher plants. *Semin*.  
555 *Cell Dev. Biol*. 2010;21:798–804.

- 556 26. Axtell MJ, Jan C, Rajagopalan R, Bartel DP. A two-hit trigger for siRNA biogenesis in plants.  
557 Cell. 2006;127:565–77.
- 558 27. Zheng Y, Li Y-F, Sunkar R, Zhang W. SeqTar: an effective method for identifying microRNA  
559 guided cleavage sites from degradome of polyadenylated transcripts in plants. Nucleic Acids  
560 Res. 2012;40:e28.
- 561 28. Zhai J, Arikiti S, Simon SA, Kingham BF, Meyers BC. Rapid construction of parallel analysis of  
562 RNA end (PARE) libraries for Illumina sequencing. Methods. 2014;67:84–90.
- 563 29. Shivaprasad PV, Chen H-M, Patel K, Bond DM, Santos BACM, Baulcombe DC. A microRNA  
564 superfamily regulates nucleotide binding site–leucine-rich repeats and other mRNAs. Plant Cell.  
565 2012;24:859–74.
- 566 30. Dukowic-Schulze S, Sundararajan A, Ramaraj T, Kianian S, Pawlowski WP, Mudge J, et al.  
567 Novel meiotic miRNAs and indications for a role of phasiRNAs in meiosis. Plant Genet.  
568 Genomics. 2016;762.
- 569 31. Kakrana A, Hammond R, Patel P, Nakano M, Meyers BC. sPARTA: a parallelized pipeline for  
570 integrated analysis of plant miRNA and cleaved mRNA data sets, including new miRNA target-  
571 identification software. Nucleic Acids Res. 2014;42:e139.

572

## 573 **Figure legends**

### 574 **Figure 1.** *PHASIS* workflow.

575 *PHAS* loci or precursors transcripts are predicted through *phasedetect* in the first step. The  
576 library-specific list of *PHAS* predictions can be summarized and annotated through *phasmerge*  
577 for libraries of interest into a *PHAS* summary. These summaries from two different groups can  
578 also be compared using “compare” mode of *phasmerge*. Triggers for *PHAS* summaries are  
579 identified through *phastrigs* either with PARE data in “validation” mode or without any  
580 experimental data in “prediction” mode. Selection between these two modes is made  
581 automatically based on a PARE library input or the lack of it. All analysis steps are independent  
582 and their execution depends upon the requirements of the user.

583

### 584 **Figure 2.** Number of *PHAS* loci or transcripts and their triggers, predicted by *PHASIS*.

585 *PHASIS* is labelled as ‘*PS*’ and it is compared to *PhaseTank* for benchmarking. **A)** 21-*PHAS* and  
586 **B)** 24-*PHAS* loci identified by both tools along with their triggers in *Arabidopsis* (*ath*),  
587 *Brachypodium* (*bdi*), *Lilium* (*lma*), rice (*osa*) and maize (*zma*). For *PHASIS* trigger prediction,  
588 results from both “validation” and “prediction” mode were included. The bars for *Lilium* 24-  
589 *PHAS* loci are split at two different points for display purposes. Triggers assigned to *PHAS* loci  
590 that do not match with known or published miRNA triggers were represented as ‘unknown’  
591 triggers.

592

### 593 **Figure 3.** Runtime comparisons between *PHASIS* and *PhaseTank*.

594 **A)** Time taken by both tools in prediction of 21- and 24-*PHAS* loci or precursors transcripts.  
595 Speed gain displayed by *PHASIS* over *PhaseTank*, approximated for both size classes, is  
596 individually marked for each species. **B)** and **C)** Time taken by both tools in predicting 21- and  
597 24-*PHAS* triggers, respectively. Speed gain displayed by *PHASIS* in “validation” and “prediction”  
598 mode over *PhaseTank* is displayed in blue and orange colors respectively. In all comparisons,  
599 *Arabidopsis* is marked as “*ath*”, *Brachypodium* as “*bdi*”, rice as “*osa*”, maize as “*zma*” and *Lilium*  
600 as “*lma*”.

601

602 **Tables**

603

604 **Table 1. Comparison of features from existing tools for phasiRNA characterization.**

	Software & citation:	<i>ShortStack</i> , Axtell MJ (2013)	<i>PhaseTank</i> , Guo et al. (2015)	<i>PHASIS</i>
PHAS prediction-related features	Tool-specific data format requirement?	no	yes	no
	Library-specific results?	no	no	yes
	<i>PHAS</i> prediction in w/o genome assembly?	no	no	yes
	Results grouped based on stage, tissue or treatments?	no	no	yes
	<i>PHAS</i> comparison between groups?	no	no	yes
	<i>PHAS</i> locus annotation?	yes (from genome GFF only)	no	yes
	<i>PHAS</i> trigger prediction w/o PARE data?	no	no	yes
Features unrelated to <i>PHAS</i> prediction	miRNA/hairpin locus prediction?	yes	no	no
	Whole-genome report of sRNA cluster characteristics?	yes	no	no

605

606 **Table 2. Comparison of predictions for *PHAS* loci, precursor transcripts, and their miRNA**  
 607 **triggers.**  
 608

Species	Type	<i>PHAS</i> locus gain with <i>PHASIS</i> over <i>PhaseTank</i>	<i>PhaseTank</i> <i>PHAS</i> loci captured by <i>PHASIS</i>	Gain in miRNA triggers: <i>PHASIS</i> (PARE supported) vs. <i>PhaseTank</i> (PARE supported)	Gain in miRNA triggers, (predicted) <i>PhaseTank</i> support
Arabidopsis	21- <i>PHAS</i>	21%	84%	-54%	-18%
Brachypodium	21- <i>PHAS</i>	145%	79%	76%	178%
	24- <i>PHAS</i>	49%	85%	36%	69%
Rice	21- <i>PHAS</i>	24%	97%	5%	54%
	24- <i>PHAS</i>	-33%	29%	No predictions	N.D.
Maize	21- <i>PHAS</i>	81%	97%	4%	55%
	24- <i>PHAS</i>	59%	86%	9%	64%
Lilium	21- <i>PHAS</i>	907%	67%*	No PARE data	No PARE data
	24- <i>PHAS</i>	1694%	94%*	No PARE data	No PARE data

For Lilium, no PARE data were available for trigger predictions. N.D. = not determined, capture trigger miRNAs.

\* Estimates, as processing for PhaseTank data made it difficult to accurately assess the proportion of

Figure 1

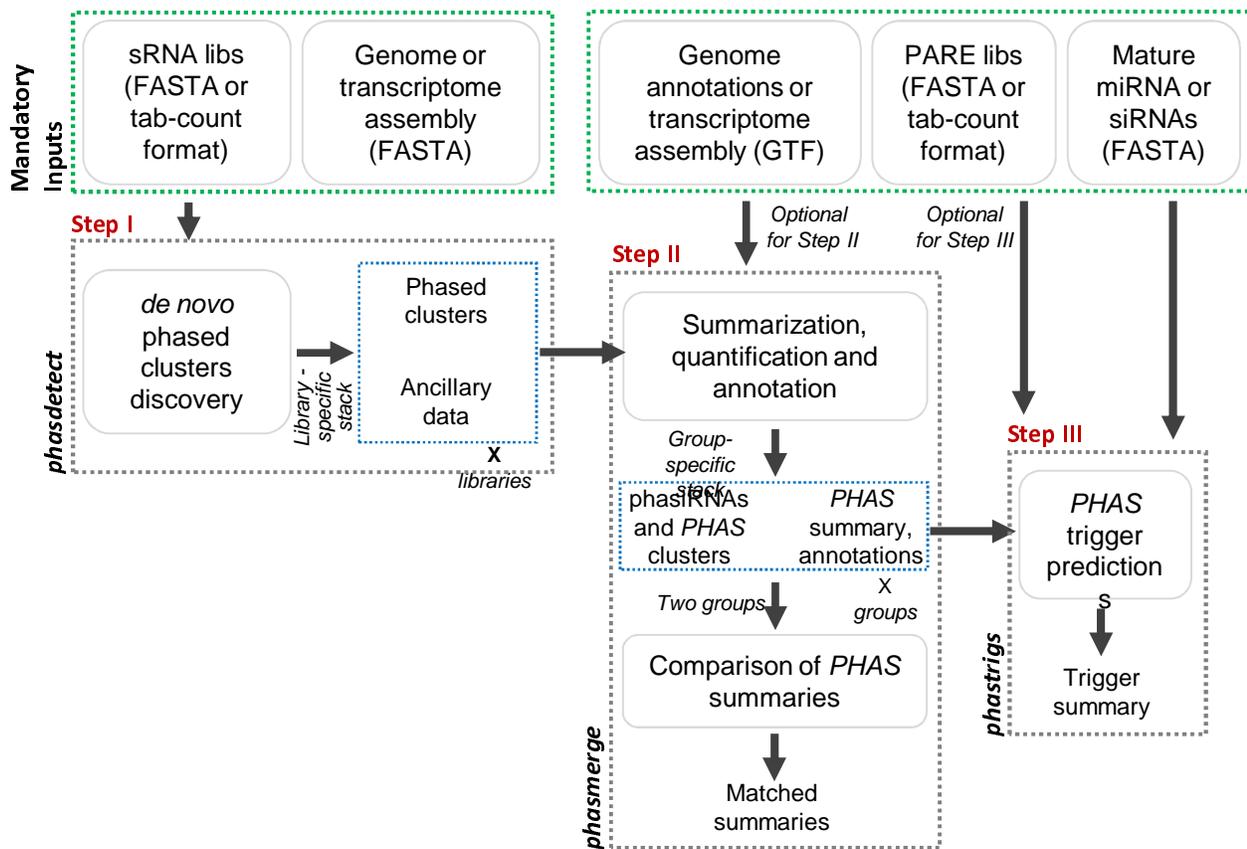


Figure 2

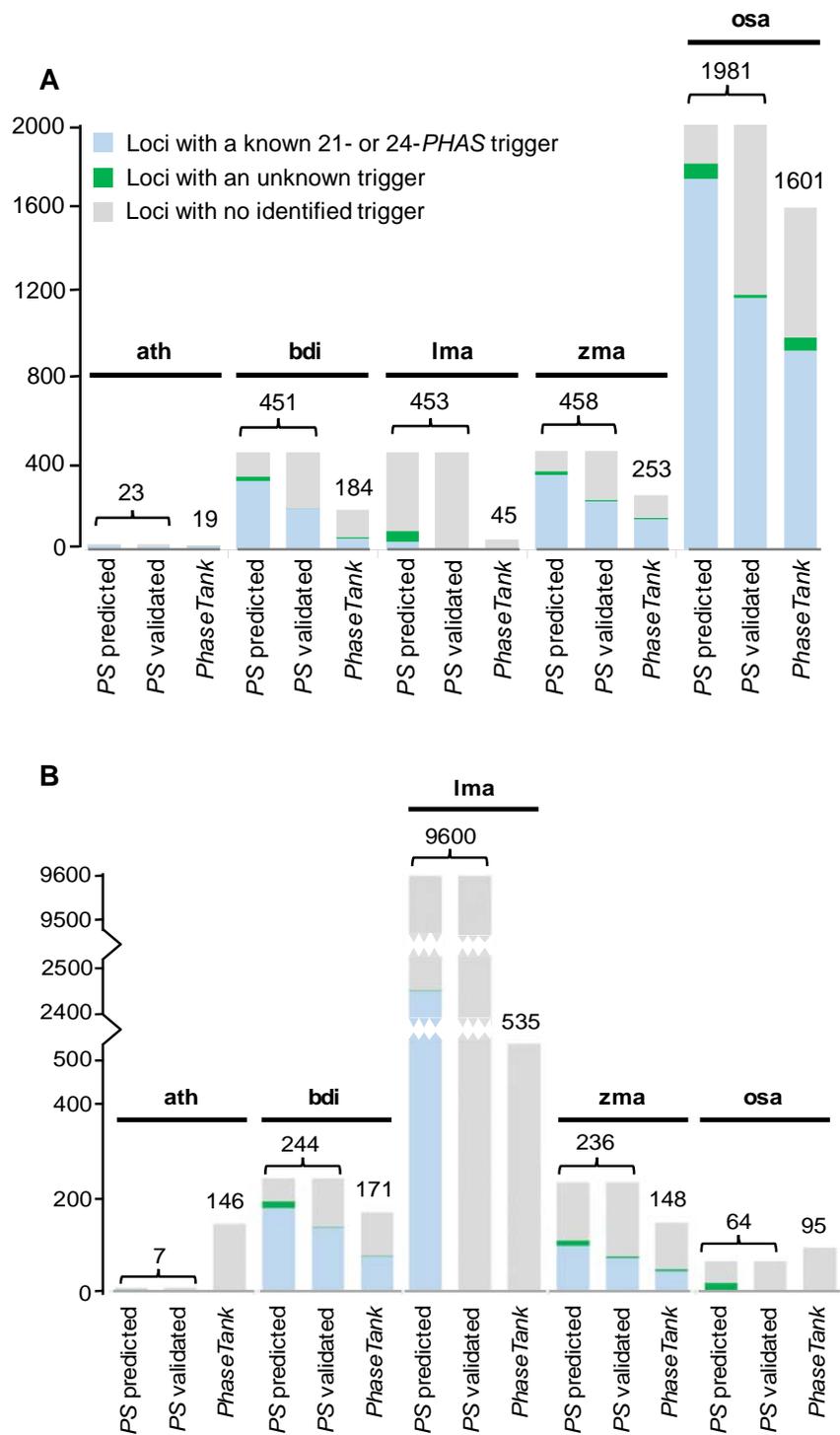


Figure 3

