

1 One precursor One siRNA model for Pol IV-dependent RNA
2 directed DNA methylation.

3 Jixian Zhai^{1,#}, Sylvain Bischof^{1,#}, Haifeng Wang^{2,1}, Suhua Feng¹, Tzoo-fen Lee³, Chong
4 Teng³, Xinyuan Chen⁴, Soo Young Park⁵, Linshan Liu⁵, Javier Gallego-Bartolome^{1,6},
5 Wanlu Liu¹, Ian R. Henderson^{1,\$}, Blake C. Meyers³, Israel Ausin^{2,*} and Steven E.
6 Jacobsen^{1,7,*}

7

8 ¹ Department of Molecular, Cell and Developmental Biology, University of California at
9 Los Angeles, Los Angeles, CA 90095, USA;

10 ² Haixia Institute of Science and Technology (HIST), Fujian Agriculture and Forestry
11 University, Fuzhou 350002, China;

12 ³ Department of Plant and Soil Sciences and Delaware Biotechnology Institute,
13 University of Delaware, Newark, DE 19716, USA;

14 ⁴ Department of Microbiology, Immunology, and Molecular Genetics, ⁵ Department of
15 Chemistry and Biochemistry, University of California at Los Angeles, Los Angeles, CA
16 90095, USA;

17 ⁶ Plant Biology Laboratory & Howard Hughes Medical Institute, The Salk Institute for
18 Biological Studies, La Jolla, CA 92037, USA;

19 ⁷ Howard Hughes Medical Institute, University of California at Los Angeles, Los
20 Angeles, CA 90095, USA.

21 # Co-first authors

22 \$ Present address: Department of Plant Sciences, University of Cambridge, Cambridge
23 CB2 3EA, United Kingdom.

24 * Correspondence: israelausin@fafu.edu.cn (I.A.), jacobsen@ucla.edu (S.E.J.)

1 **SUMMARY**

2 RNA-directed DNA methylation in *Arabidopsis thaliana* is driven by the plant-specific
3 RNA Polymerase IV (Pol IV), which is believed to transcribe precursor RNAs that give
4 rise to 24-nt small interfering RNAs (siRNAs) that target DNA methylation. However,
5 very little is known about the mechanisms of Pol IV action or the nature of Pol IV
6 transcripts. Here, we describe Pol IV-dependent RNAs (P4RNAs) from wild-type
7 *Arabidopsis* that are surprisingly short in length (30 to 40 nucleotides) and that mirror 24-
8 nucleotide siRNAs in distribution, abundance, strand bias, and 5'-adenine preference.
9 P4RNAs exhibit transcription-start-sites (TSSs) similar to Pol II products, and have
10 unique features such as 5'-monophosphates and 3'-misincorporated nucleotides. The 3'-
11 misincorporation preferentially occurs at methylated cytosines on the template DNA
12 strand, suggesting a co-transcriptional feedback to siRNA biogenesis by DNA
13 methylation to reinforce silencing locally. These results highlight an unusual mechanism
14 of Pol IV transcription and suggest a “one precursor, one siRNA” model for the
15 biogenesis of 24-nt siRNAs in *Arabidopsis*.

1 INTRODUCTION

2 In Arabidopsis, 24-nucleotide (nt) siRNAs are the triggers for RNA-directed
3 DNA methylation (RdDM), which plays central roles in repressing transposable elements
4 (TEs) and maintaining genome integrity (Law and Jacobsen, 2010; Matzke and Moshier,
5 2014). The current model for 24-nt siRNA biogenesis is composed of several sequential
6 steps. First, Pol IV recognizes heterochromatic regions, in part via SAWADEE
7 HOMEODOMAIN HOMOLOG 1 (SHH1) (Law et al., 2013), and transcribes precursor
8 RNAs. These precursor RNAs are then thought to be processed by RNA-DEPENDENT
9 RNA POLYMERASE 2 (RDR2) to form double-stranded RNAs (dsRNAs). The dsRNAs
10 are primarily cleaved by DICER-LIKE 3 (DCL3) to produce 24-nt siRNAs (Matzke and
11 Moshier, 2014). Although Pol IV is clearly required for the biogenesis of 24-nt siRNAs
12 (Herr et al., 2005; Kanno et al., 2005; Onodera et al., 2005; Pontier et al., 2005) and was
13 shown to be transcriptionally active *in vitro* (Haag et al., 2012), the nature and
14 characteristics of Pol IV-dependent transcripts remain poorly understood. The low
15 transcription level at silent Pol IV loci and the efficient downstream processing by dicers
16 has made Pol IV transcripts difficult to study. Recently, regions containing Pol IV-
17 dependent transcripts were described in a *dcl2/3/4* triple mutant that compromises
18 downstream processing of siRNAs (Li et al., 2015). In that study, with the conventional
19 assumption that Pol IV transcribes long precursors, RNA was fragmented before cloning
20 and sequenced reads were assembled into contiguous regions that corresponded to well-
21 known siRNA producing regions (Li et al., 2015). In the present study, we sought to
22 identify the very low abundance Pol IV-dependent transcripts in a wild-type background,
23 and to precisely characterize their start and end positions. We utilized a new method for

1 RNA sequencing named PATH (Parallel Analysis of Tail and Head) that utilizes RNA
2 adapters to capture both ends of any RNA greater than 27-nt with a 5' monophosphate
3 and a 3' hydroxyl. Using PATH, we efficiently cloned Pol IV-dependent RNAs (P4RNAs)
4 and found that they were only on the order of 30 to 40 nt in length. These P4RNAs mirror
5 24-nt Pol IV-dependent siRNAs in distribution, abundance, strand bias, and 5'-adenine
6 preference, suggesting that they are direct precursors that often determine both the start
7 position and strandedness of siRNAs. We also observed extensive incorporation of 3'-
8 nontemplated nucleotides preferentially at methylated cytosines on the template DNA
9 strand, suggesting that Pol IV misincorporates and terminates at DNA methylated sites.
10 Our results support a “one precursor, one siRNA” model for the biogenesis of Pol IV-
11 dependent 24-nt siRNAs in Arabidopsis that creates a positive feedback loop between
12 siRNA biogenesis and DNA methylation.

13

14 **RESULTS**

15 **Identification of short RNAs derived from Pol IV-dependent 24-nt siRNA loci.**

16 We utilized an RNA cloning scheme called PATH that uses RNA ligation and
17 gel-based size selection to capture RNAs larger than 27 nt that contain a 5'
18 monophosphate and a 3' hydroxyl (Fig 1a, see Methods). From a wild-type Columbia
19 (Col) library, around 162 million PATH reads were obtained, of which ~50 million could
20 be mapped uniquely to the Arabidopsis genome (Table S1), and PATH reads with a
21 minimal length of 27 nt were used for further analysis. Although the majority of PATH
22 reads matched to structural RNAs (t/r/sn/snoRNA) (Table S1), a large number of reads
23 mapped to previously defined Pol IV siRNA loci (Law et al., 2013) (Fig. 1a and Table

1 S2). In addition to their co-localization, the abundances of siRNAs and PATH reads were
2 also highly positively correlated (Fig. 1b). At Pol IV siRNA loci, in addition to the 24-nt
3 siRNA peak that is still present in PATH libraries (Fig. 1c), we observed a secondary
4 peak of PATH reads (P4RNAs) that primarily ranged from 30 to 40 nt (Fig. 1c). The size
5 distribution of P4RNAs was different from other PATH reads in the libraries such as
6 reads matching tRNAs and snoRNAs (Fig. S1a). The short length of P4RNAs fits well
7 with the specific preference of DCL3 for short (30 to 50 base pairs, bp) dsRNA substrates
8 (Nagano et al., 2014), consistent with the hypothesis that these short P4RNAs serve as
9 the precursors for siRNA biogenesis.

10

11 **Biogenesis and processing of Pol IV RNAs**

12 To investigate the biogenesis of P4RNAs, we constructed small RNA and PATH
13 libraries from null mutants of *NRPD1* (encoding the largest subunit of Pol IV), *NRPD/E2*
14 (encoding the shared second-largest subunit of both Pol IV and Pol V), *RDR2*, as well as
15 from the double mutants of *DCL3 NRPD1* and *DCL3 RDR2*. We found that both siRNAs
16 and P4RNAs were eliminated in all these mutant backgrounds, indicating a complete
17 dependency of P4RNA biogenesis on both Pol IV and RDR2 (Fig. 2a). Although the loss
18 of P4RNAs in *rdr2* mutants is somewhat unexpected given the assumed downstream
19 function of RDR2 (Haag et al., 2012), this result is consistent with a previous report (Li
20 et al., 2015), and suggests that *in vivo*, Pol IV and RDR2 activities are tightly coupled.
21 These results are also consistent with the earlier observation that RDR2 is in tight
22 association with Pol IV subunits *in vivo* (Haag et al., 2012; Law et al., 2011).

1 To test if P4RNA biogenesis is dependent on downstream processing events by
2 dicer proteins, we carried out additional siRNA and PATH library analysis from plants
3 triply mutant for *DCL2*, 3, and 4. Consistent with a recent report detecting Pol IV
4 dependent transcripts in a *dcl2/3/4* mutant background (Li et al., 2015), we observed a
5 dramatically increased accumulation of P4RNAs in the triple mutant *dcl2/3/4* (Fig. 2b).
6 In contrast, siRNAs levels were dramatically decreased in *dcl2/3/4* at almost all loci (Fig.
7 2b, c) and their sizes were shifted from predominantly 24 nt to predominantly 21 nt (Fig.
8 S1b). In addition, we observed that the abundances of the remaining siRNAs in *dcl2/3/4*
9 were tightly correlated with the abundances of siRNAs in wild-type (Fig 2c). We also
10 found that the abundances of P4RNAs in *dcl2/3/4* were strongly positively correlated
11 with the abundances of siRNAs in Col (Fig. S1c). As one example, P4RNAs and siRNAs
12 showed a very similar enrichment at the boundaries of long transposable elements (Fig.
13 2d). The accumulation of P4RNAs but reduction of siRNAs in *dcl2/3/4* mutants are
14 inconsistent with the interpretation that P4RNAs are merely longer, misprocessed
15 siRNAs, and suggests that P4RNAs are indeed the precursors of Pol IV-dependent
16 siRNAs.

17 We also analyzed siRNA and PATH libraries from other *DCL* mutant
18 combinations (*dcl2/4*, *dcl3*, *dcl2/3*, and *dcl3/4*) and compared these with *dcl2/3/4*. As
19 previously reported (Henderson et al., 2006), all *dcl3* mutant combinations showed a shift
20 in siRNA size from 24 nt to 21 and/or 22 nt (Fig. 3a, top panel). PATH sequencing
21 indicated that the over-accumulation of P4RNAs was negligible in *dcl2/4*, higher in *dcl3*,
22 *dcl2/3*, and *dcl3/4*, and highest in *dcl2/3/4*, indicating that the P4RNA processing by
23 DCL2 and DCL4 is less efficient than by DCL3 (Fig. 3a, bottom panel). As a secondary

1 confirmation of the accumulation of P4RNAs in *dcl3* mutants, we examined small RNA
2 northern blots for the presence of 30- to 40-nt signals. Using LNA probes to six different
3 siRNA producing loci we indeed observed a RNA smear with larger sizes, along with an
4 accumulation pattern (highest in *dcl2/3/4*) that matched the trend of the P4RNA level in
5 our PATH libraries (Fig. 3b, c) (Henderson et al., 2006). In addition, we used a reverse
6 transcription quantitative PCR (real-time RT-PCR) approach to verify the very short
7 nature of the P4RNAs. Because the peak of P4RNA abundance is 30 - 40 nt, whereas
8 longer P4RNAs are much less abundant (Fig. 1c), we reasoned that primers spaced about
9 40 nt from each other should amplify P4RNAs much more efficiently than those spaced
10 slightly further apart. After normalizing for efficiency of DNA amplification, we
11 observed a 100-fold higher level of PCR product when comparing a 36 bp versus a 58 bp
12 amplicon at one locus, and 50-fold higher level of PCR product when comparing a 39 bp
13 versus an 81 bp amplicon at a second locus (Fig. S2a).

14 To directly test whether the longer P4RNAs might simply be misprocessed
15 siRNAs that are loaded into AGO4, we performed AGO4 RNA-IP (RIP) in the *dcl2/3/4*
16 genetic background using antibodies against endogenous AGO4, and then characterized
17 AGO4-associated RNAs with high-throughput sequencing. The AGO4 RIP result showed
18 that even though there is a massive accumulation of longer P4RNAs and a reduction of
19 24nt siRNAs, AGO4 still selectively binds to the remaining 22-24 nt siRNAs but not the
20 longer P4RNAs (Fig. S2c), strongly supporting our model that P4RNAs are the
21 precursors of the 24-nt siRNAs.

22 Even though double-stranded RNAs are thought to be intermediates of siRNA
23 biogenesis, it is known that at many loci in the genome, siRNAs are predominantly found

1 on one strand of the genome but not the other (Lister et al., 2008; Zhong et al., 2014). In
2 addition, these strand-biased clusters of siRNAs correspond to strand biased DNA
3 methylation (Lister et al., 2008; Zhong et al., 2014). At these strand-biased siRNA
4 clusters, we also observed strongly strand-biased P4RNAs in both Col and *dcl2/3/4* (Fig.
5 1d and S1d). These results suggest that our PATH libraries are primarily composed of
6 single-stranded, Pol IV-derived strands, but not the RDR2-derived second strand. In
7 support of this hypothesis, we used RT-PCR to amplify these strand-biased clusters, and
8 indeed found that Pol IV/RDR2-dependent transcripts could be amplified equally from
9 both strands (Fig. S3). This suggests that RDR2 strands are likely present *in vivo*, but are
10 not cloned by the specific PATH sequencing technique employed in this study. The
11 observation that P4RNAs and siRNAs show the same strandedness suggests that the Pol
12 IV-derived strands, rather than the RDR2-derived strands, are strongly favored to become
13 the final 24-nt siRNA products.

14 In summary, results from siRNAs and PATH libraries of mutants deficient in
15 siRNA biogenesis or processing, together with the shared strandedness of P4RNAs and
16 siRNAs, suggest that the 30 to 40 nt P4RNAs serve as precursors to Pol IV dependent
17 siRNAs.

18

19 **Pol IV transcription initiates at Pol II like TSSs and favors 5'-adenine.**

20 Next we investigated the nature of the 5' ends of P4RNAs. Given that our cloning
21 method only captures RNAs with 5-monophosphates (Fig. 1a), we sought to measure the
22 proportion of P4RNAs containing this type of 5' end. To address this, we used
23 Terminator exonuclease to preferentially digest RNAs with a 5'-monophosphate. We

1 subsequently measured the abundance of the remaining P4RNAs by real-time RT-PCR.
2 Consistent with a recent study (Li et al., 2015), Terminator treatment degraded the
3 majority of P4RNAs at all loci tested (Fig. S2b). Thus, while it is possible that
4 subpopulations of P4RNAs have other end structures such as 5'-triphosphates, 5'-caps, or
5 5'-hydroxyl groups, it appears that the majority of P4RNAs contain 5'-monophosphates.

6 Because our method did not include a fragmentation step that is typical in RNA-
7 seq library protocols (Li et al., 2015), it was possible to detect the 5' nucleotide of
8 P4RNA reads. We observed a strong enrichment of T/C (Y) at the -1 position (the
9 nucleotide immediately upstream of the first nucleotide of the P4RNA read) and A/G (R)
10 at the +1 position (beginning nucleotide of the read) (Fig. 4a). Further, the four possible
11 Y/R dinucleotides at the -1/+1 positions were by far the most enriched dinucleotides at
12 the 5' end of P4RNAs (Fig. 4b). This pattern is very similar to that known for the
13 transcriptional start sites (TSSs) of RNA Polymerase II (Pol II) in plants and other
14 organisms and it is referred to as the “Y/R rule” (Cumbie et al., 2015; Nechaev et al.,
15 2010; Yamamoto et al., 2007). This result suggests that Pol IV has retained this
16 preference from its evolutionary ancestor Pol II (Ream et al., 2009), and that the 5' ends
17 of P4RNAs likely represent Pol IV transcriptional start sites. The short TSS like
18 sequences at the 5' ends of P4RNA, along with their very short nature does not support
19 previous models in which Pol IV initiates transcription solely at the nucleosome-depleted
20 promoter regions near the ends of transposons to produce long transcripts (Li et al., 2015).
21 Instead, our results suggest that Pol IV can initiate transcription at many positions that
22 resemble the Y/R features of Pol II TSSs, transcribing many short P4RNAs along the
23 length of transposons.

1 Because Pol IV transcripts feature Pol II-like TSSs, we also performed genome-
2 wide profiling of Pol II occupancy in Arabidopsis via ChIP-seq in wild-type and different
3 mutant backgrounds (*nrpd1* and *dcl2/3/4*), and compared this with Pol IV ChIP-seq (Law
4 et al., 2013). Our results showed that Pol II does not appear to access Pol IV loci even in
5 the absence of Pol IV (*nrpd1*) (Fig. S4a). Furthermore, whole-genome bisulfite
6 sequencing and small RNA sequencing of floral tissues from the weak Pol II mutant
7 (*nrpb2-3*) did not reveal obvious changes in either DNA methylation or siRNA
8 biogenesis (Fig. S4b, c), which is consistent with our previous analysis of *nrpb2-3* using
9 leaf tissue (Stroud et al., 2013). Therefore despite the proposed crosstalk between Pol II
10 and the RdDM pathway (Zheng et al., 2009), our data suggest that Pol II and Pol IV
11 occupy distinct territories on the genome.

12 Arabidopsis 24-nt siRNAs are primarily loaded into AGO4, and are strongly
13 biased toward having a 5'-adenine, which was previously shown to involve an AGO4
14 loading preference (Havecker et al., 2010; Mi et al., 2008). Interestingly, we found that
15 P4RNAs, like siRNAs, also show a strong enrichment for 5'-adenine (Fig. 4c, and S5a,
16 b). Because the size of P4RNAs is approximately 30 to 40 nt, on average only one 24-nt
17 siRNA duplex could be processed from each of these P4RNA precursors. This fact,
18 coupled with the shared 5' adenine preference and the shared strand preference, suggests
19 that 24-nt siRNAs are preferentially cleaved from the 5' portion of P4RNAs. Consistent
20 with this hypothesis, DCL3 was shown to prefer short double-stranded RNAs (30 to 50
21 bp) that contain a 5' adenine (Nagano et al., 2014). Therefore, our results favor a scenario
22 in which the 5'-adenine preference of P4RNAs likely contributes to the 5'-adenine
23 preference of Pol IV siRNAs. In addition, our results provide a plausible explanation as

1 to why AGO4 evolved to bind siRNAs with 5' adenine. Taken together, the short size
2 and 5'-A feature of Pol IV transcripts may help to channel their processing to DCL3
3 rather than dicer proteins in other silencing pathways, and thus lead to production of
4 predominantly 24-nt siRNAs at Pol IV transcribed loci.

5

6 **Pol IV transcription preferentially terminates at methylated cytosines with**
7 **misincorporated nucleotides.**

8 We analyzed the sequence composition of P4RNA reads with perfect match to
9 genome and found enrichment for A, C, and U at last three positions of the 3' end (Fig.
10 S5c). It is not known whether Pol IV tends to cease transcription at this sequence, or
11 whether it might transcribe a longer RNA that is then processed by an unknown
12 endonuclease. We found a similar compositional bias at the 3' end of Pol IV-dependent
13 siRNAs (Fig. S5d), although the magnitude of the biases were lower for siRNAs,
14 suggesting that P4RNAs are processed at some level at their 3' ends to produce siRNAs.
15 Despite the ACU enrichment at the 3' end of siRNAs, they still show enrichment for 5'
16 adenine, which is likely explained by the AGO4 preference for loading siRNAs with a 5'
17 adenine. These results provides additional evidence that the P4RNAs described here are
18 indeed the precursors of siRNAs, and it is again consistent with (1) the shared
19 strandedness of P4RNAs and siRNAs, and (2) the hypothesis that the Pol IV strand,
20 rather than the RDR2 strand, is favored as the final siRNA.

21 We also performed an analysis in which we allowed multiple mismatches during
22 genome mapping of P4RNAs. Interestingly, we found that more than half of the P4RNAs
23 contained one or two non-templated nucleotides at their 3'-ends (Fig. 5a, b). In contrast,

1 in the same PATH libraries, reads derived from Pol II-transcribed coding regions, or from
2 microRNA processing intermediates, had very few mismatches, and these mismatches
3 were not localized to the 3' ends (Fig. S6a, b). To further rule out the possibility that
4 P4RNA 3' non-templated nucleotides are due to lower quality of sequencing toward the
5 end of the read or result from incomplete trimming of adapter sequence, we analyzed the
6 second read (read2) from the paired-end sequencing. On read2, the beginning nucleotide
7 corresponds to the 3' end nucleotide of RNA, where base quality is high and there is no
8 trimming step involved, and we observed the same high-level of non-templated
9 nucleotides at the 3'-end of P4RNA (Fig. S6c).

10 The 3' end base composition of P4RNAs containing non-templated nucleotides
11 was quite different from those with a perfect genome match (Fig. S5c, e), suggesting that
12 P4RNAs with non-templated nucleotides terminate by a different mechanism than those
13 without. On the other hand, the 5' end base composition of reads with non-templated
14 nucleotides was very similar to those without, suggesting that both classes may share a
15 similar transcription initiation mechanism by Pol IV (Fig. S5f, g). All four nucleotides
16 were found among the non-templated nucleotides, although there was some preference
17 for guanines (Fig. S7a). In addition, we observed different preferences for non-templated
18 nucleotides depending on the sequence that should have been present, suggesting that the
19 preference for a particular non-templated nucleotide is determined by the sequence of the
20 DNA template for Pol IV (Fig. S7b, c). Because all four nucleotides were present, and
21 because different incorrect nucleotides were present depending on the template DNA
22 sequence, it seems most likely that these nucleotides arise from misincorporation during
23 Pol IV transcription rather than from the activity of a terminal transferase.

1 Because misincorporation of nucleotides occurred most frequently at positions
2 corresponding to guanines (which would be cytosines on the DNA template) (Fig. 5a),
3 we hypothesized that misincorporation might be caused by in part by cytosine DNA
4 methylation. The most highly methylated sequences in the Arabidopsis genome are CG
5 dinucleotides (Cokus et al., 2008; Lister et al., 2008), and Pol IV siRNA loci targeted by
6 RdDM are usually heavily methylated at most CG sites (Stroud et al., 2013) (Fig. S4b).
7 Agreeing with the DNA methylation hypothesis, we found that CG dinucleotides (G
8 being the last nucleotide of the RNA) exhibited by far the highest enrichment of
9 misincorporation amongst the 16 possible dinucleotide sequences (Fig. 5c). In addition,
10 CG dinucleotides were strongly enriched at the 3'-end of P4RNAs that exhibited
11 misincorporation (Fig. 5d). The second group of most commonly methylated sequences
12 in the genome and at Pol IV siRNA loci are CHG sites (Fig. S4b) (Cokus et al., 2008;
13 Lister et al., 2008). A trinucleotide analysis showed that all of the trinucleotides showing
14 the strongest tendency for misincorporation were those that contained CG sites (Fig. S7d).
15 In addition, CHG sites showed a strong tendency for misincorporation, which was higher
16 than AHG, THG, or GHG sites (Fig. S7d).

17 To directly test whether the loss of methylation can alter the pattern of 3' end
18 misincorporation, we analyzed the *ddm1* (*decrease in DNA methylation 1*) mutant that
19 exhibits a severe loss of methylation in heterochromatin (Matzke and Mosher, 2014).
20 Despite the significant loss of methylation, many silent loci are still producing 24-nt
21 siRNAs in *ddm1* (Colome-Tatche et al., 2012), suggesting that Pol IV is still largely
22 functional in *ddm1*. Because P4RNAs are elevated in *dcl3* mutants, we utilized a small
23 RNA dataset from *ddm1 dcl3* double mutant (RNAs smaller than 200 nt) (McCue et al.,

1 2015) and compared these with a similar dataset from the *dcl3* single mutant. We first
2 focused our analysis on a set of strong *ddm1* hypomethylated CG DMRs (Differentially
3 Methylated Regions) (Fig. 6a, see methods) from whole genome bisulfite sequencing
4 data of *ddm1* (Creasey et al., 2014). We found that the enrichment of P4RNA 3'-
5 misincorporation at CG dinucleotides was eliminated in *ddm1 dcl3* compared to the *dcl3*
6 single mutant (Fig. 6b). We also examined trinucleotide enrichments in *ddm1 dcl3* at a set
7 of CHG DMRs (Fig. S7e), and observed a significant reduction at CHG but not at AHG,
8 THG, or GHG (Fig. S7f). Moreover, P4RNAs in *ddm1/dcl3* were slightly longer than
9 those in *dcl3* single mutant (Fig. S7g). Take together, these results suggest that DNA
10 methylation itself is contributing to the pattern of 3'-misincorporation.

11 If P4RNAs are the precursors of siRNAs and DCL3 can cleave to some extent
12 from the 3'-end of P4RNA, siRNAs should also contain some level of misincorporated
13 bases at their 3' ends. By allowing for mismatches during genome mapping, we indeed
14 found that siRNAs contain non-templated nucleotides that were enriched at the 3' end
15 (Fig. S6d). However, the proportion of siRNAs with 3'-mismatches (~1%) was far lower
16 than that of P4RNAs (~50%, Fig. 5b), suggesting that far fewer siRNAs are processed
17 from the 3' ends of P4RNAs than from the 5' ends, possibly due to the preference of 5'-
18 adenine by DCL3 (Nagano et al., 2014). We also observed enrichment of CG sites at the
19 3' end of siRNAs that showed misincorporation (Fig. S6e), again consistent with some
20 level of processing of the 3' end of P4RNAs into siRNAs.

21 Since DCL3 appears to process from both ends of the double-stranded Pol
22 IV/RDR2-derived RNA, and since P4RNAs are enriched for adenines at their 5' ends and
23 misincorporated nucleotides at their 3' ends (Fig. 7a), two predictions are that siRNAs

1 with 5' adenines should have lower than average 3' misincorporation, and siRNAs with 3'
2 misincorporation should have lower than average 5' adenine content. Indeed, we found
3 that siRNAs with 5' adenines had 50% lower misincorporation than siRNAs with other 5'
4 nucleotides (Fig. 7b), and siRNAs with misincorporated nucleotides showed a lower 5'
5 adenine content than those with perfect matches to the genome (Fig. 7c). These results
6 further support that P4RNAs containing misincorporated 3' nucleotides are processed
7 into siRNAs. In summary, our results support that three different mechanisms can
8 contribute the formation of the 3' end of P4RNAs - termination or 3' end processing at
9 sequences enriched for ACU at the last three positions of the P4RNA, termination
10 associated with misincorporation at methylated cytosines, and termination associated
11 with misincorporation at other nucleotides.

12

13 **DISCUSSION**

14 The results of this study support the general scheme that P4RNAs are first
15 transcribed by Pol IV, made double stranded by RDR2, and diced by DCL3 and other
16 dicers to make the siRNAs which are loaded into AGO4. DCL3 cleavage produces an
17 siRNA duplex with symmetric structure (Matzke and Mosher, 2014), but only one strand
18 (the guide strand) is retained in AGO4, while the other strand (the passenger strand) is
19 cleaved by AGO4 and degraded (Ye et al., 2012). Little is known about guide-strand
20 selection of 24-nt siRNAs, and it remains unclear why, at many loci, siRNAs
21 predominately match to one strand, and direct strand-specific methylation (Zhong et al.,
22 2014). Our analysis revealed that the strand biases for P4RNAs and siRNAs are highly
23 positively correlated (Fig. 1d and Fig. S1d), suggesting that the siRNAs derived from the

1 Pol IV strand rather than the RDR2-synthesized complementary strand, are favored as the
2 guide-strand siRNA. This bias appears to be partially accomplished by a strong bias for
3 adenine to be present as the first nucleotide of P4RNAs, by the preference of DCL3 for 5'
4 adenines for siRNA processing, and by the preference of AGO4 for loading siRNAs with
5 a 5' adenine (Fig 7a). P4RNAs also appear to be processed to some extent at their 3' ends
6 to yield siRNAs that have similar 3' end signatures as are present in P4RNAs, including
7 either an enrichment for ACU sequences, or misincorporated nucleotides. In addition,
8 Pol IV transcripts are so short that on average only one siRNA will arise from each
9 P4RNA. Finally, the short nature of Pol IV transcripts, coupled with the preference of
10 DCL3 for short dsRNAs, may serve to channel Pol IV/RDR2 products into DCL3 rather
11 than other dicers, such as DCL4 that prefers long dsRNAs (Nagano et al., 2014). These
12 results imply a “one precursor, one siRNA” model for processing of siRNAs from Pol IV
13 to RDR2 to DCL3 to AGO4 (Fig. 7a). Based on the patterns of TSSs, strand-bias and 3'-
14 misincorporation we conclude that the majority of PATH reads at siRNA loci are
15 P4RNA, but it is still possible that some of these PATH reads are RDR2-transcribed as
16 they are genetically dependent on both RDR2 and Pol IV.

17 Because Pol IV transcripts are so short, the sites of Pol IV transcriptional
18 initiation and termination are largely determining the positions of siRNAs in the genome.
19 Production of short transcripts by Pol IV could provide for siRNA biogenesis while
20 avoiding the risk of transcribing full-length transposable elements. In addition, the short
21 length of Pol IV transcripts may help to prevent spreading of RdDM to flanking regions,
22 and allow specific silencing of transposons that are close to genes.

1 Our proposed model for biogenesis of 24-nt siRNAs in *Arabidopsis* share
2 characteristics with that of 21U-RNAs (or piRNAs) in *C. elegans*. 21U-RNAs are 21-nt
3 in length and begin with a 5' uridine, and they are autonomously expressed from
4 thousands of loci dispersed in two broad regions of chromosome IV (Batista et al., 2008;
5 Ruby et al., 2006). Like P4RNAs, precursors of 21U-RNAs, which are ~26 nucleotide
6 Pol II transcripts, are terminated by an unknown mechanism to produce unusually short
7 transcripts (Gu et al., 2012). However, unlike P4RNAs, 21U-RNA precursors are
8 processed at their 5' ends to remove two nucleotides, and transcription units usually
9 contain a conserved motif 42-nt upstream of the mature piRNA (Batista et al., 2008; Gu
10 et al., 2012; Ruby et al., 2006). In contrast we did not find evidence for sequence
11 conservation upstream of P4RNA start sites (Fig. S4d), and instead Pol IV appears to be
12 recruited by epigenetic signals such as the binding of H3K9 methylation through the Pol
13 IV interacting protein SHH1 (Law et al., 2013). Furthermore, while there is clear
14 evidence for transposon-derived secondary siRNA production through RNA-dependent
15 RNA Polymerase (RdRP) activity in *C. elegans* (22G-RNA) (Lee et al., 2012), as well as
16 Zucchini-dependent, secondary phased piRNAs in mammals (Han et al., 2015; Mohn et
17 al., 2015), there is little evidence to support the possibility of secondary siRNA
18 production in the RdDM pathway.

19 The 5' ends of P4RNAs are very similar in sequence composition to that of Pol II
20 transcripts, showing strong enrichment for Y/R dinucleotides at the -1/+1 positions.
21 However, as opposed to Pol II transcripts that begin with trimethylguanosine caps, the
22 majority of P4RNAs have 5' monophosphates. Since RNA polymerases normally start
23 transcribing with a triphosphate containing nucleotide, it is unclear how Pol IV

1 transcripts acquire a 5' monophosphate. It is possible that Pol IV is able to initiate
2 transcription with a nucleoside monophosphate as previously reported for some RNA
3 polymerase (Martin and Coleman, 1989; Ranjith-Kumar et al., 2002), or the P4RNAs
4 may be initially capped after which a de-capping enzyme converts the trimethylguanosine
5 cap to a 5' monophosphate, or an unknown polyphosphatase-like enzyme may directly
6 convert the 5'-triphosphate to a 5'-monophosphate. Additionally, since 5'
7 monophosphate containing RNAs are often a target of 5' to 3' exonucleases, there may
8 be mechanisms to protect P4RNA 5' ends until subsequent processing steps are
9 completed. It is also unclear why P4RNA evolved to have 5' monophosphates rather
10 than cap structures, but one possibility is that this helps the cell avoid inadvertently
11 mistaking a P4RNA for a Pol II transcript in order to avoid translation of transposon
12 RNAs.

13 Our observations suggest that a component of the mechanism by which Pol IV
14 transcription terminates is that DNA methylation on the template strand causes
15 misincorporation of inappropriate bases. The mechanism by which this happens is
16 unclear, but it is known that DNA methylation can cause transcriptional elongation
17 defects in *Neurospora* (Rountree and Selker, 1997). The preferential termination of Pol
18 IV transcription near DNA methylation would promote siRNA generation near sites of
19 preexisting DNA methylation, thereby creating a self-reinforcing loop in which siRNAs
20 direct DNA methylation targeting and DNA methylation helps direct the location of
21 further siRNA production.

1 **Experimental Procedures**

2

3 **Biological materials**

4 The mutant alleles of *nrpd1-4*, *nrpb2-3*, *rdr2-1*, *nrp(d/e)2*, *rdr2/dcl3*, and combinations of
5 *dcl2/3/4* used in this study were in the background of *Arabidopsis thaliana* ecotype
6 Columbia-0 and have been previously described (Henderson et al., 2006; Li et al., 2008;
7 Pontier et al., 2005; Xie et al., 2004; Zheng et al., 2009). The *dcl3/nrpd1* double mutant
8 was obtained by crossing *dcl3-1* with *nrpd1-4*. Plants were grown in a growth chamber
9 with 16 hour of light or greenhouse condition for five weeks. Immature inflorescence
10 tissues including inflorescence meristem and early stages floral buds (up to stage 11/12)
11 were collected.

12

13 **Sequencing of small RNA, PATH mRNA, and BS-seq libraries**

14 Total RNA was first treated with RiboMinus™ Plant Kit for RNA-Seq
15 (Invitrogen A10838-08) to remove rRNA, followed by size selection of RNA on a 15%
16 UREA TBE Polyacrylamide gel (Invitrogen, EC6885BOX). Gels containing RNA with
17 size between 15- to 27-nt were kept for small RNA library, while gels containing 28- to
18 ~300-nt RNA were kept for PATH library. After gel elution, library construction for both
19 sRNA and PATH was done using the Illumina TruSeq Small RNA Sample Preparation
20 Kit (RS-200-0012), except that at the final size selection step, PCR products were
21 separated on a 6% TBE Polyacrylamide gel (Invitrogen, EC6265BOX) and selected for
22 the range from 120- to ~1000-bp for PATH library. Gel-eluted PCR products with
23 different TruSeq index sequences were pooled and sent for Illumina sequencing. The

1 mRNA library was constructed using the Illumina TruSeq RNA Sample Preparation Kit
2 (RS-122-2001) according to the standard manual. PATH libraries were sequenced using
3 either paired-end mode with length of read1 being 120-bp and length of read2 being 30-
4 bp (PE120+30), or single-end 100-bp (SE100); while sRNA and mRNA libraries were
5 sequenced with single-end 50-bp (SE50). For BS-seq, DNA was isolated using the
6 DNeasy Plant Mini Kit (Qiagen #69104) according to manufacturer instructions and
7 quantified using the Qubit dsDNA High Sensitivity Kit (Life Technologies #Q32851).
8 Libraries were constructed with 30ng DNA using the Ovation Ultralow Methyl-Seq
9 Library Systems (NuGEN #0335). Bisulfite conversion was done using the EpiTect
10 Bisulfite Kit (Qiagen # 59104). BS-seq Libraries were sequenced at a length of 50 bp. All
11 sequencing was carried out on Illumina HiSeq machines at the Broad Stem Cell Research
12 Center (BSCRC) sequencing core at University of California, Los Angeles.

13

14 **Realtime RT-PCR and RNA gel blot**

15 Realtime RT-PCR to detect P4RNAs in various genotypes was performed as the
16 following: Total RNA were extracted from 100 mg of flowers using Trizol (Ambion)
17 with an extra step of 24:1 Chloroform : isoamyl alcohol to remove remaining phenol
18 prior to isopropanol precipitation, the resuspended RNAs were cleaned up using Quick-
19 RNA miniprep (Zymo research, USA) for further purification and complete gDNA
20 removal. Then 5ug of the purified RNA from each sample was used for RT reaction with
21 SuperScript III first-strand kit (Invitrogen, USA). 1uL of the RT reaction was used for
22 real time PCR using iQ SYBR green supermix (Biorad, USA). We used terminator
23 exonuclease (Epicenter, USA) for the removal of RNA with 5'-monophosphate: 5ug of

1 total RNA were treated for an hour at 30 degree with Terminator, and the control, non-
2 treated, RNA with the same buffer and 50% glycerol instead of terminator exonuclease
3 enzyme. After Terminator treatment the RNA were used for RT reaction as described
4 above. Primer information can be found in [Table S3](#).

5
6 RNA gel blotting was performed as previously described (Henderson et al., 2006) with
7 LNA probes for detecting transposon regions and regular RNA probes for detecting 5S
8 locus. Probe information can be found in [Table S3](#).

9 10 **Pol II ChIP-seq and AGO4-RIP**

11 ChIP-seq was performed as described previously (Johnson et al., 2014). 5 ug Pol II
12 antibodies (Abcam #ab817) were used for each ChIP. Libraries for Pol II ChIP-seq were
13 generated using the Ovation Ultralow DR Multiplex System (NuGen #0330) and
14 sequenced at a length of 50 bp. ChIP-seq data was visualized using ngsplot (Shen et al.,
15 2014).

16 AGO4-RIP was performed as previously described (Ji et al., 2011) with commercial
17 AGO4 antibody (Agrisera #AS09 617). RNA isolated from RIP was used for library
18 construction with Illumina TruSeq Small RNA Sample Preparation Kit (RS-122-2001)
19 according to the standard manual.

20 21 **Bioinformatic analysis**

22 *Data handling*

1 In general, qseq files received from the sequencing core were demultiplexed with an in-
2 house Perl script and converted to fastq files for downstream analysis. For small RNA
3 and PATH data, original reads were first trimmed using Cutadapt (v1.4), then mapped to
4 the reference TAIR10 genome using Bowtie (Langmead et al., 2009) allowing only one
5 unique hit (-m 1). We allowed zero mismatch for small RNA mapping (-v 0) and up to
6 three mismatches for PATH mapping (-v 3). mRNA data were mapped using Tophat
7 (Trapnell et al., 2009) allowing two mismatches and only one unique hit.

8

9 *Analysis of siRNAs and P4RNAs*

10 We used a list of Pol IV dependent siRNA loci that is previously described (Law et al.,
11 2013) and manually inspected and filtered out a few loci that are tRNA related (listed in
12 [Table S2](#)). In brief, these are 200 base pair bins where siRNAs were significantly reduced
13 in a Pol IV mutant compared to two replicates of wild-type controls ($FRD < 10^{-10}$). More
14 details can be found in the “Identification of siRNA clusters” section of Methods in the
15 Law et al. paper (Law et al., 2013). Our definition of P4RNAs is any 27+ nt PATH reads
16 derived from these ~7000 previously defined Pol IV dependent 24-nt siRNA loci. For the
17 abundance calculation, siRNA reads with the length between 18 and 26 nt, and PATH
18 reads with the length greater or equal to 27 nt were included ([Table S1](#)). 5'-Adenine
19 preference was analyzed by calculating the nucleotide composition at each position
20 (counting from 5' end) for Pol IV siRNAs and P4RNAs. Percentage of 3'-nontemplated
21 nucleotides in P4RNAs was calculated at each position by analyzing the mapping results
22 with a custom Perl script with the focus on the “MD:Z” column in sam format. The
23 length of the 3'-misincorporation was done using the seed mapping option of bowtie to

1 map the first 24 nt of P4RNA perfectly and then allow up to ten mismatches for 3'
2 portion. This analysis confirmed that majority of P4RNAs carry one or two 3'-
3 misincorporated nucleotides. Therefore we chose to allow three mismatches for P4RNA
4 mapping for all analysis.

5

6 *Normalization of P4RNA and sRNA abundance*

7 Abundances of P4RNA in each library were normalized to the total number of sequenced
8 reads, and abundances of siRNAs were normalized to the sum of all miRNAs in each
9 library. For example, in Figure 2b, in Col PATH library we obtained 161,241,523 reads
10 in total and of which 63,739 were classified as sasRNA (27+ nt), while in dcl2/3/4 PATH
11 library we obtained 89,226,121 reads in total and of which 4,283,063 were classified as
12 sasRNA. Therefore the percentage of P4RNA in dcl2/3/4 compared to Col is calculated
13 as percentage = $(4,283,063/89,226,121) / (63,739/161,241,523) = 121.43 \approx 120$ fold.

14

15 *BS-seq analysis and DMR calling*

16 Analysis of the floral BS-seq libraries of wild-type and *ddm1* mutant (Creasey et al., 2014)
17 were performed using BSMAP (Xi and Li, 2009), allowing only uniquely mapped reads
18 and discarded sibling PCR products, with the tolerance of 2 mismatches per 50 bp. DMR
19 calling was performed as previously described (Stroud et al., 2013), with a more stringent
20 criteria: 1) sum of all sequenced cytosines in the 100bp bin need to be at least 100; 2) the
21 difference in CG methylation at each bin needs to be at least 0.5. This filtering allows us
22 to focus on regions that are not only highly methylated in wild-type but also lose
23 methylation dramatically in the *ddm1* mutant. CHG DMRs were filtered with similar

1 criteria, with the minimal loss of CHG in *ddm1* compared to wild-type being 0.4, and the
2 count of covered cytosines in each bin no less than 50.

3

4 *Nucleotide composition and enrichment*

5 P4RNAs a defined set of loci were measured for their 3'-end sequence composition. For
6 those P4RNAs that contain the 3' misincorporation, different compositions of the
7 reference sequence at the first mismatched nucleotide (mono-nucleotide), the last
8 matched plus first mismatched (di-nucleotide), or the last two matched plus first
9 mismatched (tri-nucleotide) were counted as the observed value. The mono-, di- or tri-
10 nucleotide compositions of the DNA sequence (or K-mer) were calculated using Jellyfish
11 (Marcais and Kingsford, 2011) and considered as expected value assuming P4RNAs end
12 equally likely on any nucleotide. The observed value and expected value were both
13 normalized to their own population, and the ratio of observed/expected is used to show
14 the relative enrichment of certain mono-, di- or tri-nucleotide composition.

15

16 **Accession codes**

17 Sequencing data have been deposited at GEO (GSE61439). We used public data of floral
18 *ddm1* and wild-type BS-seq (GSE52346), *ddm1/dcl3* sRNA-seq (GSE57191), and *dcl3*
19 sRNA-seq (GSE62801).

20

21 **Acknowledgments**

22 We thank members of the Jacobsen lab for insightful discussions. We thank Molly
23 Megraw for suggestions on TSS analysis, and Mahnaz Akhavan for technical assistance.

1 J.Z. thanks Dr. Chomdao Chommy and Seksan Sapsubbsakul for encouragement. High
2 throughput sequencing was performed at the UCLA BSCRC BioSequencing Core
3 Facility. This work was supported by NIH grant GM60398 to S.E.J and NSF award
4 1051576 to B.C.M. S.B. is supported by a postdoctoral fellowship of the Swiss National
5 Science Foundation. J.G-B is a Human Frontiers Science Program fellow
6 (LT000425/2012-L). J.Z. is a Life Science Research Foundation postdoctoral fellow,
7 sponsored by the Gordon and Betty Moore Foundation. S.E.J. is an Investigator of the
8 Howard Hughes Medical Institute.

9

10 **Author Contributions**

11 J.Z., S.B., S.F., I.A, I.H., X.C., and C.T. performed experiments. T.L., S.F., and J.G-B.
12 provided materials. S.P., L.L., and X.C. participated in the genetic experiments. S.E.J.
13 and B.C.M. oversaw the study. J.Z., H.W., and W.L. analyzed data. J.Z. and S.E.J.
14 designed the study and wrote the manuscript.

Figure Legends

Fig. 1. Identification of siRNA-loci associated short RNAs (P4RNAs).

- a. Procedure for the construction and analysis of PATH libraries. An example siRNA locus is shown at the bottom with IGV screenshots of P4RNAs and siRNAs matching to that region.
- b. Abundances are highly correlated of PATH reads (27+ nt) and sRNA reads (18 to 26 nt) in wild-type Col from previously defined Pol IV siRNA loci (Law et al., 2013).
- c. Size distribution plots of all PATH reads and sRNA reads in Col from Pol IV siRNA loci. A distinct peak at 30 to 40 nt can be seen in the PATH library; we named these “P4RNAs”.
- d. P4RNAs and siRNAs share the same strand bias at Pol IV siRNA loci. Only Pol IV siRNA loci matched by more than 100 P4RNAs in Col were selected to obtain a robust calculation of strandedness. The plus-strand ratio was calculated as the abundance of reads matching to the plus strand divided by the total number of reads at that locus.

Fig. 2. Biogenesis and processing of P4RNAs.

- a. Both Pol IV siRNAs and P4RNAs are eliminated in *nRPD1*, *nRP(d/e)2*, *RDR2*, *dCL3/nRPD1*, and *dCL3/RDR2*. Small RNA abundances were normalized to the sum of all TAIR10-annotated miRNAs in each sRNA library, and then compared to Col;

- P4RNA abundances were normalized to the total number of reads in each PATH library.
- b. Pol IV siRNAs are reduced in *dcl2/3/4*, whereas P4RNAs accumulate substantially.
 - c. Abundances of P4RNA and siRNA at each locus in Col and *dcl2/3/4*.
 - d. At long TEs, P4RNAs resemble the distribution of siRNAs – enriched at promoters and termini of TEs, and also spread into TE bodies.

Fig. 3. P4RNAs in *dcl* mutant combinations.

- a. Size of Pol IV siRNAs shifted from 24 nt in Col to 21/22 nt in *dcl3*-containing mutants, accompanied by the accumulation of P4RNAs. P4RNA abundances are normalized to that in Col.
- b. RNA blot analysis of siRNAs and P4RNAs in *dcl* mutant backgrounds.
- c. Size distribution of P4RNAs in *dcl3*-containing mutants.

Fig. 4. P4RNAs feature Pol II like TSSs and favor 5'-Adenine.

- a. TSSs of Pol IV exhibit preference for C/T at the -1 position and A/G at the +1 position, resembling the “Y/R rule” of Pol II TSS, calculated using P4RNAs from *dcl2/3/4*.
- b. Dinucleotide enrichment at -1/+1 of P4RNA, calculated using P4RNAs from *dcl2/3/4*.
- c. Both siRNA and P4RNAs from Col have a strong preference for 5' adenine.

Fig. 5. Pol IV transcription preferentially terminates at methylated cytosines with misincorporated nucleotides.

- a. Example of 3'-nontemplated nucleotides on P4RNAs. The black bar represents 50 bp in length.
- b. Length of the 3'-end nontemplated nucleotides, defined by the first mismatched nucleotide to the last nucleotide. If a P4RNA has no mismatch, the length is zero.
- c. Di-nucleotide enrichment at the first mismatched position at the 3' end.
- d. Frequency of CG dinucleotide on reference sequence over the P4RNAs with misincorporation, “-1” marks the last perfectly matched position, and “+1” marks the first mismatched position. The count of each CG is designated to the position of the G, therefore the peak at “+1” represents a peak of CG at “-1/+1”.

Fig. 6. Pol IV transcription preferentially terminates at methylated cytosines with misincorporated nucleotides.

- a. CG hypomethylated DMRs in *ddm1* compared to a wild-type control.
- b. 3'-misincorporated nucleotides of P4RNA at *ddm1* CG DMRs in *ddm1 dcl3* compared to *dcl3*.

Fig. 7. “one precursor, one siRNA” model for the biogenesis of Pol IV dependent 24-nt siRNAs.

- a. Pol IV transcription is initiated at Pol II-like TSSs. A short RNA of ~30 to 40 nt (with 5'-adenine preference) is produced by Pol IV at heterochromatic regions; misincorporation then occurs at the cytosine position (red bar) and terminates Pol

IV transcription. This process typically yields a P4RNA with 5' adenine and 3' misincorporation. P4RNA then goes through processing that involves RDR2 synthesizing the complementary strand (gray), DCL3 cutting from the 5' end (major) or 3' end (minor) of the P4RNA, and loading of the P4RNA-derived siRNA strand into AGO4 as the guide-strand. Eventually one P4RNA precursor gives rise to one siRNA that is derived from either its 5' or 3' end. (“M” underneath the DNA template indicates DNA methylation at the heterochromatic region.)

- b. When DCL3 cuts from the 5' end of P4RNAs, the resulting siRNAs are more likely to carry a 5' adenine and be perfectly matched. Indeed a reduction of 3' mismatches is observed for siRNAs with a 5' adenine.
- c. When DCL3 cuts from the 3' end of P4RNAs, the resulting siRNAs are more likely to carry the 3' misincorporation and less likely the 5' adenine. This is consistent with the observation that siRNAs with mismatch have lower percentage of 5' adenine.

References

- Batista, P.J., Ruby, J.G., Claycomb, J.M., Chiang, R., Fahlgren, N., Kasschau, K.D., Chaves, D.A., Gu, W., Vasale, J.J., Duan, S., *et al.* (2008). PRG-1 and 21U-RNAs interact to form the piRNA complex required for fertility in *C. elegans*. *Mol Cell* *31*, 67-78.
- Cokus, S.J., Feng, S., Zhang, X., Chen, Z., Merriman, B., Haudenschild, C.D., Pradhan, S., Nelson, S.F., Pellegrini, M., and Jacobsen, S.E. (2008). Shotgun bisulphite sequencing of the *Arabidopsis* genome reveals DNA methylation patterning. *Nature* *452*, 215-219.
- Colome-Tatche, M., Cortijo, S., Wardenaar, R., Morgado, L., Lahouze, B., Sarazin, A., Etcheverry, M., Martin, A., Feng, S., Duvernois-Berthet, E., *et al.* (2012). Features of the *Arabidopsis* recombination landscape resulting from the combined loss of sequence variation and DNA methylation. *Proc Natl Acad Sci U S A* *109*, 16240-16245.
- Creasey, K.M., Zhai, J., Borges, F., Van Ex, F., Regulski, M., Meyers, B.C., and Martienssen, R.A. (2014). miRNAs trigger widespread epigenetically activated siRNAs from transposons in *Arabidopsis*. *Nature* *508*, 411-415.
- Cumbie, J., Ivanchenko, M., and Megraw, M. (2015). nanoCAGE-XL and CapFilter: an Approach to Genome Wide Identification of High Confidence Transcription Start Sites. *BMC Genomics* *in press*.
- Gu, W., Lee, H.C., Chaves, D., Youngman, E.M., Pazour, G.J., Conte, D., Jr., and Mello, C.C. (2012). CapSeq and CIP-TAP identify Pol II start sites and reveal capped small RNAs as *C. elegans* piRNA precursors. *Cell* *151*, 1488-1500.
- Haag, J.R., Ream, T.S., Marasco, M., Nicora, C.D., Norbeck, A.D., Pasa-Tolic, L., and Pikaard, C.S. (2012). In vitro transcription activities of Pol IV, Pol V, and RDR2 reveal coupling of Pol IV and RDR2 for dsRNA synthesis in plant RNA silencing. *Mol Cell* *48*, 811-818.
- Han, B.W., Wang, W., Li, C., Weng, Z., and Zamore, P.D. (2015). Noncoding RNA. piRNA-guided transposon cleavage initiates Zucchini-dependent, phased piRNA production. *Science* *348*, 817-821.
- Havecker, E.R., Wallbridge, L.M., Hardcastle, T.J., Bush, M.S., Kelly, K.A., Dunn, R.M., Schwach, F., Doonan, J.H., and Baulcombe, D.C. (2010). The *Arabidopsis* RNA-directed DNA methylation argonautes functionally diverge based on their expression and interaction with target loci. *Plant Cell* *22*, 321-334.
- Henderson, I.R., Zhang, X., Lu, C., Johnson, L., Meyers, B.C., Green, P.J., and Jacobsen, S.E. (2006). Dissecting *Arabidopsis thaliana* DICER function in small RNA processing, gene silencing and DNA methylation patterning. *Nature genetics* *38*, 721-725.
- Herr, A.J., Jensen, M.B., Dalmay, T., and Baulcombe, D.C. (2005). RNA polymerase IV directs silencing of endogenous DNA. *Science* *308*, 118-120.

- Ji, L., Liu, X., Yan, J., Wang, W., Yumul, R.E., Kim, Y.J., Dinh, T.T., Liu, J., Cui, X., Zheng, B., *et al.* (2011). ARGONAUTE10 and ARGONAUTE1 regulate the termination of floral stem cells through two microRNAs in Arabidopsis. *PLoS Genet* 7, e1001358.
- Johnson, L.M., Du, J., Hale, C.J., Bischof, S., Feng, S., Chodavarapu, R.K., Zhong, X., Marson, G., Pellegrini, M., Segal, D.J., *et al.* (2014). SRA- and SET-domain-containing proteins link RNA polymerase V occupancy to DNA methylation. *Nature* 507, 124-128.
- Kanno, T., Huettel, B., Mette, M.F., Aufsatz, W., Jaligot, E., Daxinger, L., Kreil, D.P., Matzke, M., and Matzke, A.J. (2005). Atypical RNA polymerase subunits required for RNA-directed DNA methylation. *Nature genetics* 37, 761-765.
- Langmead, B., Trapnell, C., Pop, M., and Salzberg, S.L. (2009). Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol* 10, R25.
- Law, J.A., Du, J., Hale, C.J., Feng, S., Krajewski, K., Palanca, A.M., Strahl, B.D., Patel, D.J., and Jacobsen, S.E. (2013). Polymerase IV occupancy at RNA-directed DNA methylation sites requires SHH1. *Nature* 498, 385-389.
- Law, J.A., and Jacobsen, S.E. (2010). Establishing, maintaining and modifying DNA methylation patterns in plants and animals. *Nature reviews* 11, 204-220.
- Law, J.A., Vashisht, A.A., Wohlschlegel, J.A., and Jacobsen, S.E. (2011). SHH1, a homeodomain protein required for DNA methylation, as well as RDR2, RDM4, and chromatin remodeling factors, associate with RNA polymerase IV. *PLoS Genet* 7, e1002195.
- Lee, H.C., Gu, W., Shirayama, M., Youngman, E., Conte, D., Jr., and Mello, C.C. (2012). *C. elegans* piRNAs mediate the genome-wide surveillance of germline transcripts. *Cell* 150, 78-87.
- Li, C.F., Henderson, I.R., Song, L., Fedoroff, N., Lagrange, T., and Jacobsen, S.E. (2008). Dynamic regulation of ARGONAUTE4 within multiple nuclear bodies in Arabidopsis thaliana. *PLoS Genet* 4, e27.
- Li, S., Vandivier, L.E., Tu, B., Gao, L., Won, S.Y., Li, S., Zheng, B., Gregory, B.D., and Chen, X. (2015). Detection of Pol IV/RDR2-dependent transcripts at the genomic scale in Arabidopsis reveals features and regulation of siRNA biogenesis. *Genome Res* 25, 235-245.
- Lister, R., O'Malley, R.C., Tonti-Filippini, J., Gregory, B.D., Berry, C.C., Millar, A.H., and Ecker, J.R. (2008). Highly integrated single-base resolution maps of the epigenome in Arabidopsis. *Cell* 133, 523-536.
- Marcais, G., and Kingsford, C. (2011). A fast, lock-free approach for efficient parallel counting of occurrences of k-mers. *Bioinformatics* 27, 764-770.

- Martin, C.T., and Coleman, J.E. (1989). T7 RNA polymerase does not interact with the 5'-phosphate of the initiating nucleotide. *Biochemistry* 28, 2760-2762.
- Matzke, M.A., and Mosher, R.A. (2014). RNA-directed DNA methylation: an epigenetic pathway of increasing complexity. *Nature reviews* 15, 394-408.
- McCue, A.D., Panda, K., Nuthikattu, S., Choudury, S.G., Thomas, E.N., and Slotkin, R.K. (2015). ARGONAUTE 6 bridges transposable element mRNA-derived siRNAs to the establishment of DNA methylation. *EMBO J* 34, 20-35.
- Mi, S., Cai, T., Hu, Y., Chen, Y., Hodges, E., Ni, F., Wu, L., Li, S., Zhou, H., Long, C., *et al.* (2008). Sorting of small RNAs into Arabidopsis argonaute complexes is directed by the 5' terminal nucleotide. *Cell* 133, 116-127.
- Mohn, F., Handler, D., and Brennecke, J. (2015). Noncoding RNA. piRNA-guided slicing specifies transcripts for Zucchini-dependent, phased piRNA biogenesis. *Science* 348, 812-817.
- Nagano, H., Fukudome, A., Hiraguri, A., Moriyama, H., and Fukuhara, T. (2014). Distinct substrate specificities of Arabidopsis DCL3 and DCL4. *Nucleic Acids Res* 42, 1845-1856.
- Nechaev, S., Fargo, D.C., dos Santos, G., Liu, L., Gao, Y., and Adelman, K. (2010). Global analysis of short RNAs reveals widespread promoter-proximal stalling and arrest of Pol II in *Drosophila*. *Science* 327, 335-338.
- Onodera, Y., Haag, J.R., Ream, T., Costa Nunes, P., Pontes, O., and Pikaard, C.S. (2005). Plant nuclear RNA polymerase IV mediates siRNA and DNA methylation-dependent heterochromatin formation. *Cell* 120, 613-622.
- Pontier, D., Yahubyan, G., Vega, D., Bulski, A., Saez-Vasquez, J., Hakimi, M.A., Lerbs-Mache, S., Colot, V., and Lagrange, T. (2005). Reinforcement of silencing at transposons and highly repeated sequences requires the concerted action of two distinct RNA polymerases IV in Arabidopsis. *Genes & development* 19, 2030-2040.
- Ranjith-Kumar, C.T., Gutshall, L., Kim, M.J., Sarisky, R.T., and Kao, C.C. (2002). Requirements for de novo initiation of RNA synthesis by recombinant flaviviral RNA-dependent RNA polymerases. *Journal of virology* 76, 12526-12536.
- Ream, T.S., Haag, J.R., Wierzbicki, A.T., Nicora, C.D., Norbeck, A.D., Zhu, J.K., Hagen, G., Guilfoyle, T.J., Pasa-Tolic, L., and Pikaard, C.S. (2009). Subunit compositions of the RNA-silencing enzymes Pol IV and Pol V reveal their origins as specialized forms of RNA polymerase II. *Mol Cell* 33, 192-203.
- Rountree, M.R., and Selker, E.U. (1997). DNA methylation inhibits elongation but not initiation of transcription in *Neurospora crassa*. *Genes & development* 11, 2383-2395.

Ruby, J.G., Jan, C., Player, C., Axtell, M.J., Lee, W., Nusbaum, C., Ge, H., and Bartel, D.P. (2006). Large-scale sequencing reveals 21U-RNAs and additional microRNAs and endogenous siRNAs in *C. elegans*. *Cell* *127*, 1193-1207.

Shen, L., Shao, N., Liu, X., and Nestler, E. (2014). ngs.plot: Quick mining and visualization of next-generation sequencing data by integrating genomic databases. *BMC Genomics* *15*, 284.

Stroud, H., Greenberg, M.V., Feng, S., Bernatavichute, Y.V., and Jacobsen, S.E. (2013). Comprehensive analysis of silencing mutants reveals complex regulation of the *Arabidopsis* methylome. *Cell* *152*, 352-364.

Trapnell, C., Pachter, L., and Salzberg, S.L. (2009). TopHat: discovering splice junctions with RNA-Seq. *Bioinformatics* *25*, 1105-1111.

Xi, Y., and Li, W. (2009). BSMAP: whole genome bisulfite sequence MAPPING program. *BMC bioinformatics* *10*, 232.

Xie, Z., Johansen, L.K., Gustafson, A.M., Kasschau, K.D., Lellis, A.D., Zilberman, D., Jacobsen, S.E., and Carrington, J.C. (2004). Genetic and functional diversification of small RNA pathways in plants. *PLoS biology* *2*, E104.

Yamamoto, Y.Y., Ichida, H., Matsui, M., Obokata, J., Sakurai, T., Satou, M., Seki, M., Shinozaki, K., and Abe, T. (2007). Identification of plant promoter constituents by analysis of local distribution of short sequences. *BMC Genomics* *8*, 67.

Ye, R., Wang, W., Iki, T., Liu, C., Wu, Y., Ishikawa, M., Zhou, X., and Qi, Y. (2012). Cytoplasmic assembly and selective nuclear import of *Arabidopsis* Argonaute4/siRNA complexes. *Mol Cell* *46*, 859-870.

Zheng, B., Wang, Z., Li, S., Yu, B., Liu, J.Y., and Chen, X. (2009). Intergenic transcription by RNA polymerase II coordinates Pol IV and Pol V in siRNA-directed transcriptional gene silencing in *Arabidopsis*. *Genes & development* *23*, 2850-2860.

Zhong, X., Du, J., Hale, C.J., Gallego-Bartolome, J., Feng, S., Vashisht, A.A., Chory, J., Wohlschlegel, J.A., Patel, D.J., and Jacobsen, S.E. (2014). Molecular mechanism of action of plant DRM de novo DNA methyltransferases. *Cell* *157*, 1050-1060.

Figure 1

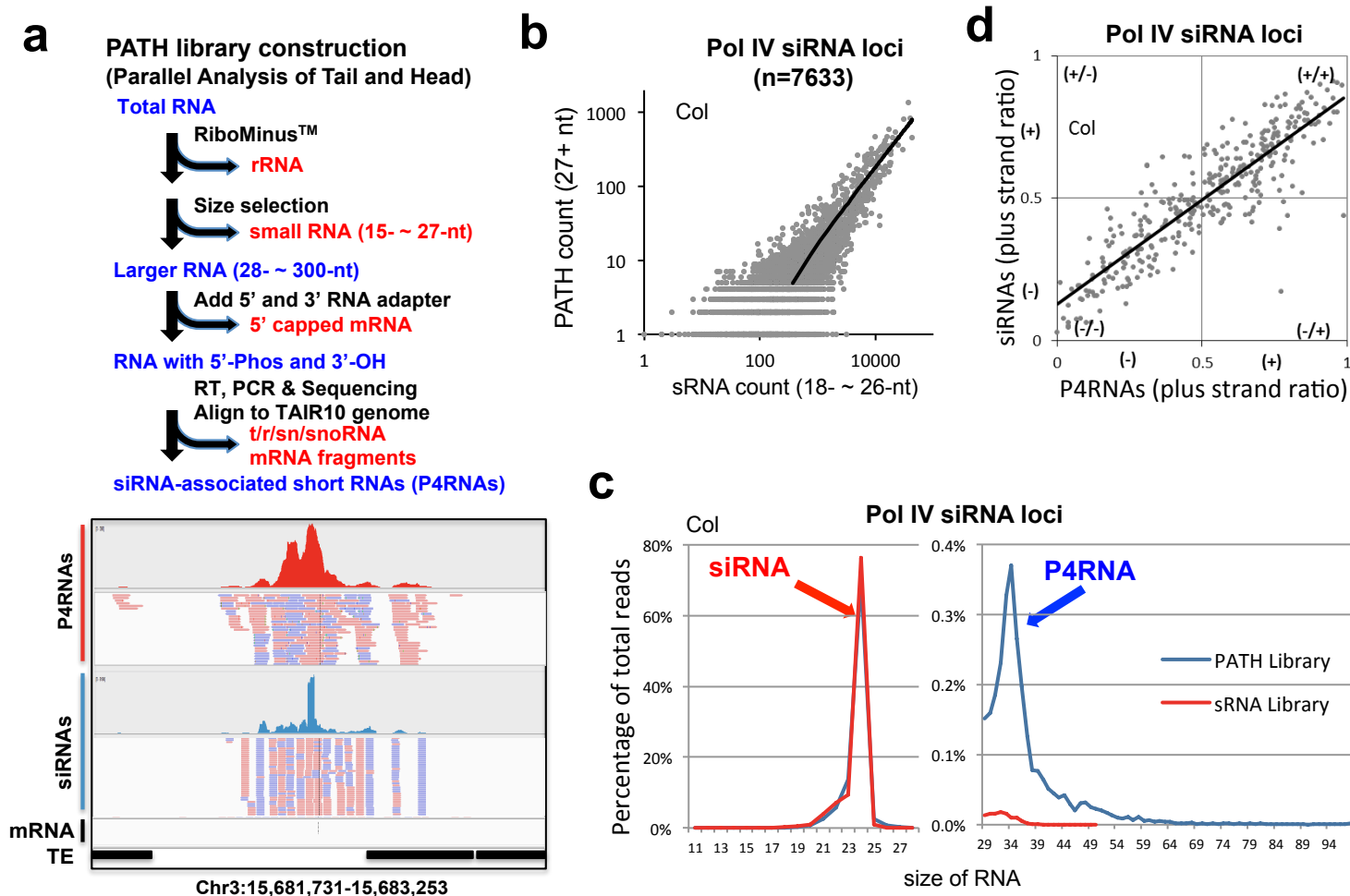


Figure 1. Identification of siRNA-loci associated short RNAs (P4RNAs) .

- Procedure for the construction and analysis of PATH libraries. An example siRNA locus is shown at the bottom with IGV screenshots of P4RNAs and siRNAs matching to that region.
- Abundances are highly correlated of PATH reads (27+ nt) and sRNA reads (18 to 26 nt) in wild-type Col from previously defined Pol IV siRNA loci (Law et al., 2013).
- Size distribution plots of all PATH reads and sRNA reads in Col from Pol IV siRNA loci. A distinct peak at 30 to 40 nt can be seen in the PATH library; we named these “P4RNAs”.
- P4RNAs and siRNAs share the same strand bias at Pol IV siRNA loci. Only Pol IV siRNA loci matched by more than 100 P4RNAs in Col were selected to obtain a robust calculation of strandedness. The plus-strand ratio was calculated as the abundance of reads matching to the plus strand divided by the total number of reads at that locus.

Figure 2

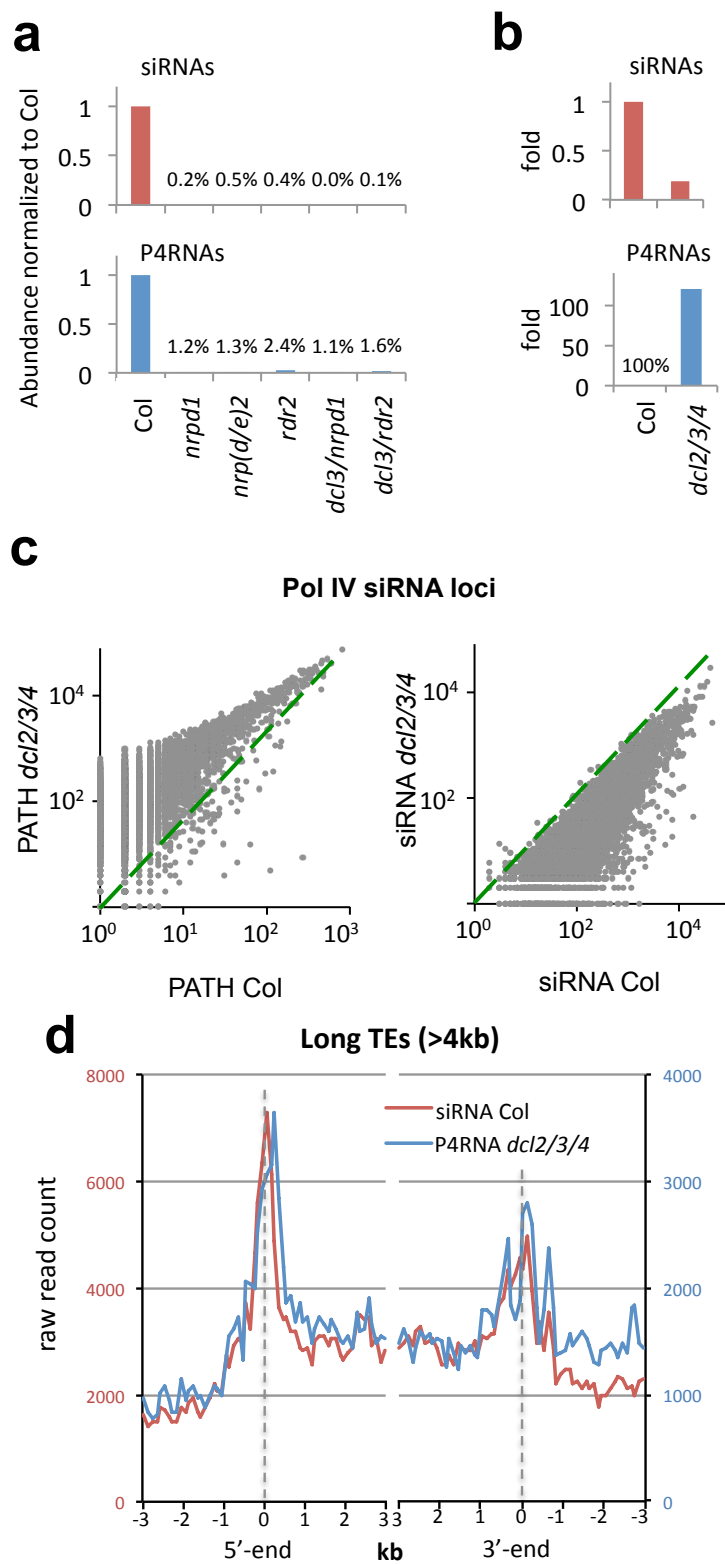


Figure 2. Biogenesis and processing of P4RNAs.

- Both Pol IV siRNAs and P4RNAs are eliminated in *nrpd1*, *nrp(d/e)2*, *rdr2*, *dcl3/nrpd1* and *dcl3/rdr2*. Small RNA abundances were normalized to the sum of all TAIR10-annotated miRNAs in each sRNA library, and then compared to Col; P4RNA abundances were normalized to the total number of reads in each PATH library.
- Pol IV siRNAs are reduced in *dcl2/3/4*, whereas P4RNAs accumulate substantially.
- Correlation of P4RNA and siRNA abundance at each locus in Col and *dcl2/3/4*.
- At long TEs, P4RNAs resemble the distribution of siRNAs – enriched at promoters and termini of TEs, and also spread into TE bodies.

Figure 3

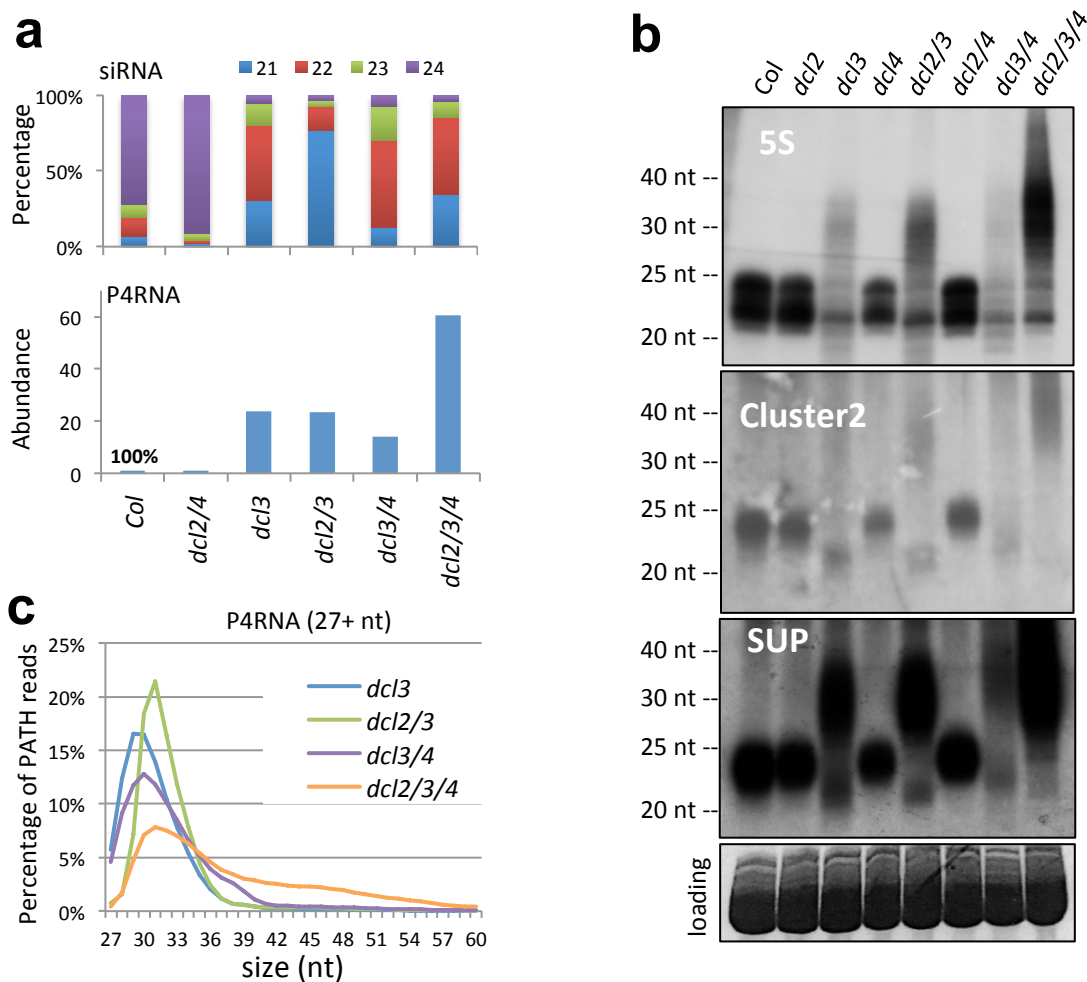


Figure 3. P4RNAs in *dcl* mutant combinations.

- Size of Pol IV siRNAs shifted from 24 nt in Col to 21/22 nt in *dcl3*-containing mutants, accompanied by the accumulation of P4RNAs. P4RNA abundances are normalized to that in Col.
- RNA blot analysis of siRNAs and P4RNAs in *dcl* mutant backgrounds.
- Size distribution of P4RNAs in *dcl3*-containing mutants.

Figure 4

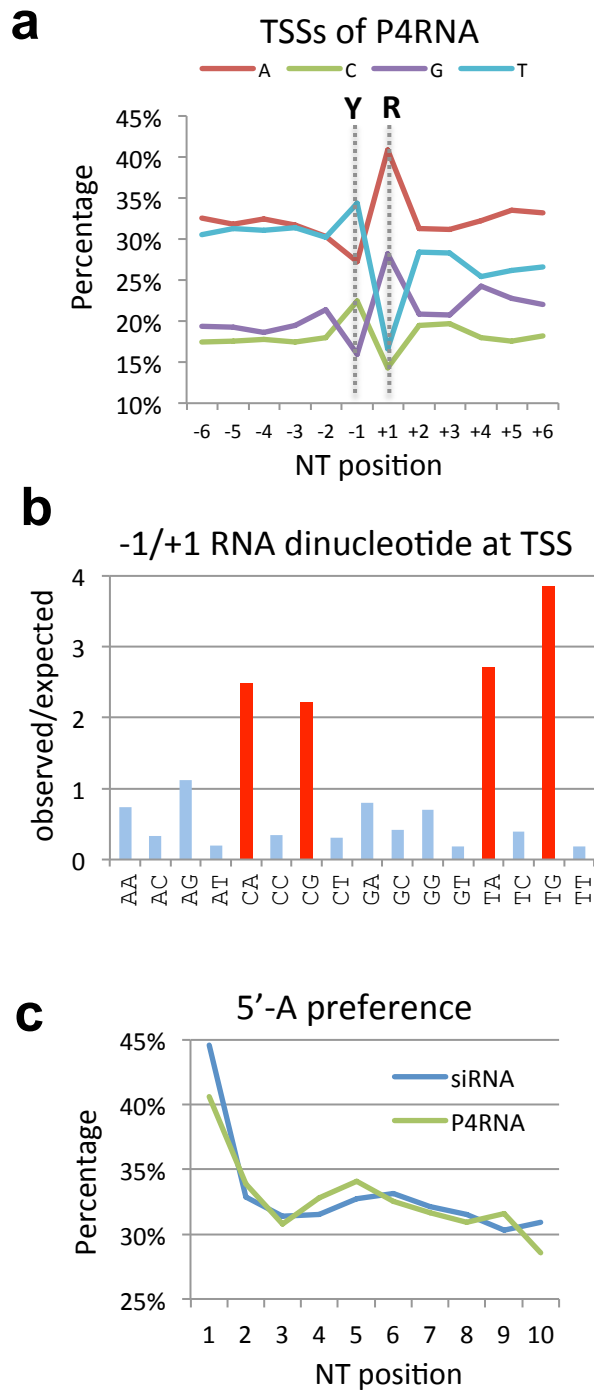


Figure 4. P4RNAs feature Pol II like TSSs and favor 5'-Adenine.

- TSSs of Pol IV exhibit preference for C/T at the -1 position and A/G at the +1 position, resembling the “Y/R rule” of Pol II TSS, calculated using P4RNAs from *dcl2/3/4*.
- Dinucleotide enrichment at -1/+1 of P4RNA, calculated using P4RNAs from *dcl2/3/4*.
- Both siRNA and P4RNAs from Col have a strong preference for 5' adenine.

Figure 5

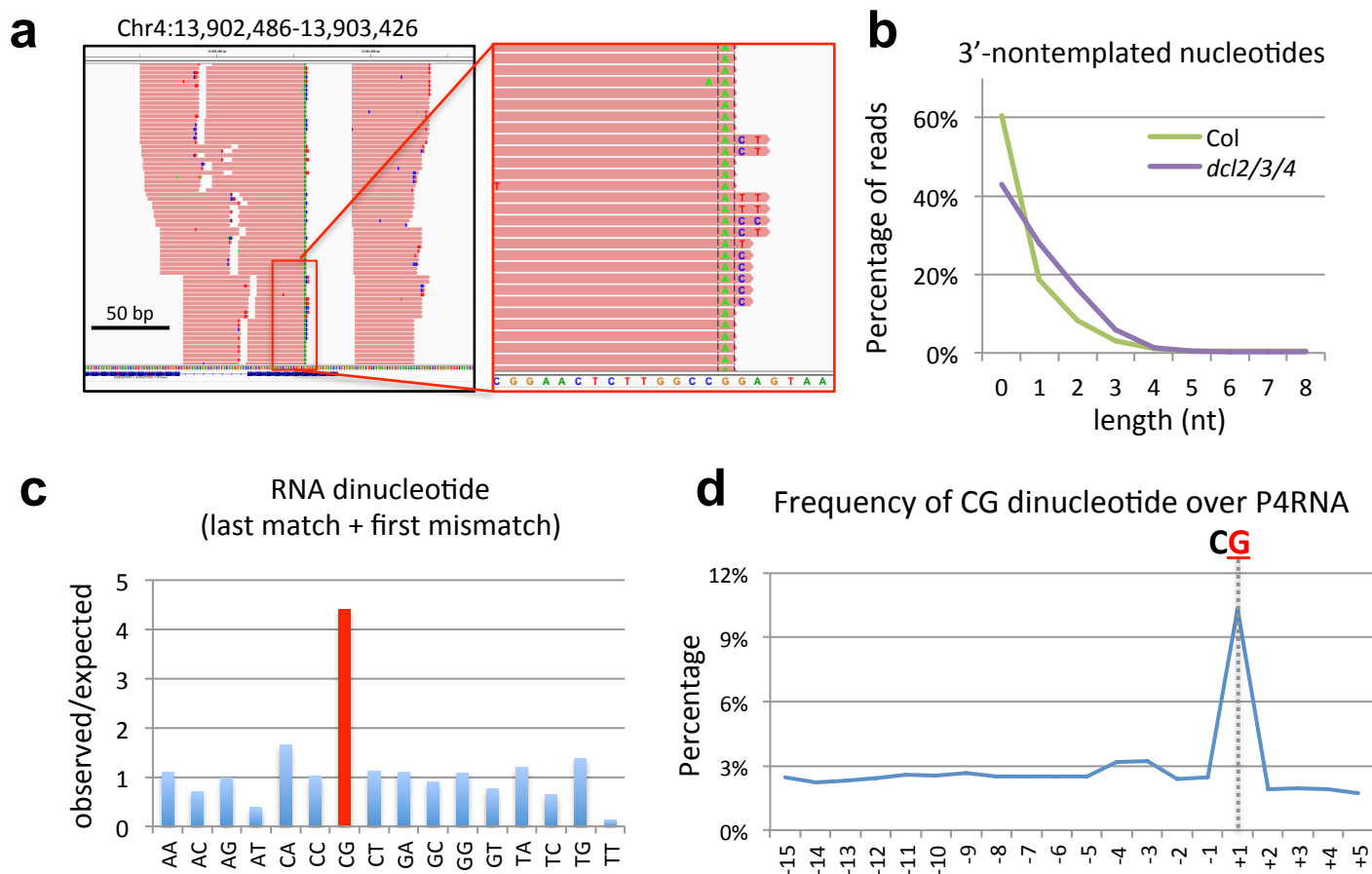


Figure 5. Pol IV transcription preferentially terminates at methylated cytosines with misincorporated nucleotides.

- Example of 3'-nontemplated nucleotides on P4RNAs. The black bar represents 50 bp in length.
- Length of the 3'-end nontemplated nucleotides, defined by the first mismatched nucleotide to the last nucleotide. If a P4RNA has no mismatch, the length is zero.
- Di-nucleotide enrichment at the first mismatched position at the 3' end.
- Frequency of CG dinucleotide on reference sequence over the P4RNAs with misincorporation, "-1" marks the last perfectly matched position, and "+1" marks the first mismatched position. The count of each CG is designated to the position of the G, therefore the peak at "+1" represents a peak of CG at "-1/+1".

Figure 6

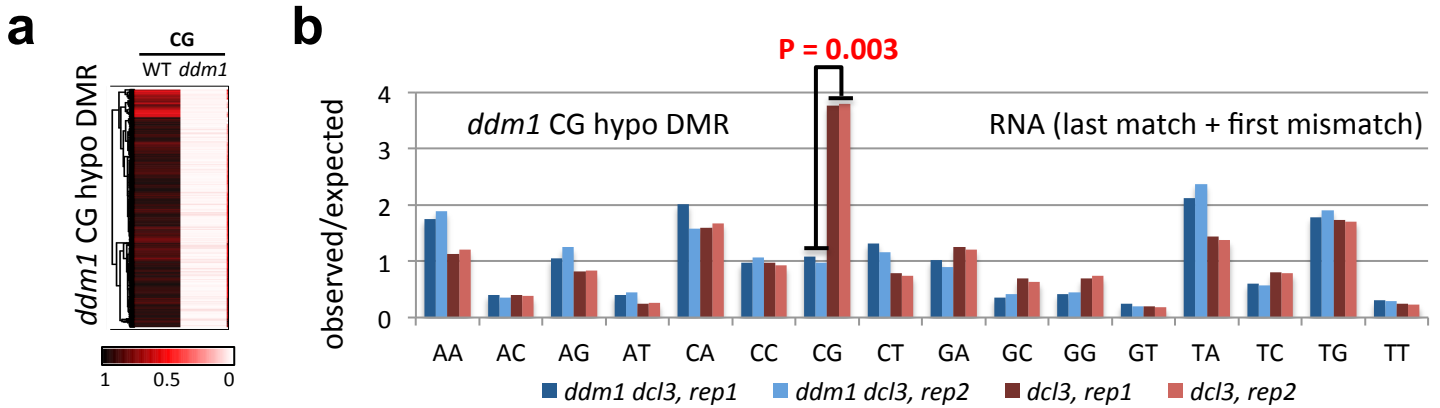


Figure 6. Misincorporation of P4RNAs at CG sites is suppressed in *ddm1/dcl3*.

- CG hypomethylated DMRs in *ddm1* compared to a wild-type control.
- 3'-misincorporated nucleotides of P4RNA at *ddm1* CG DMRs in *ddm1 dcl3* compared to *dcl3*.

Figure 7

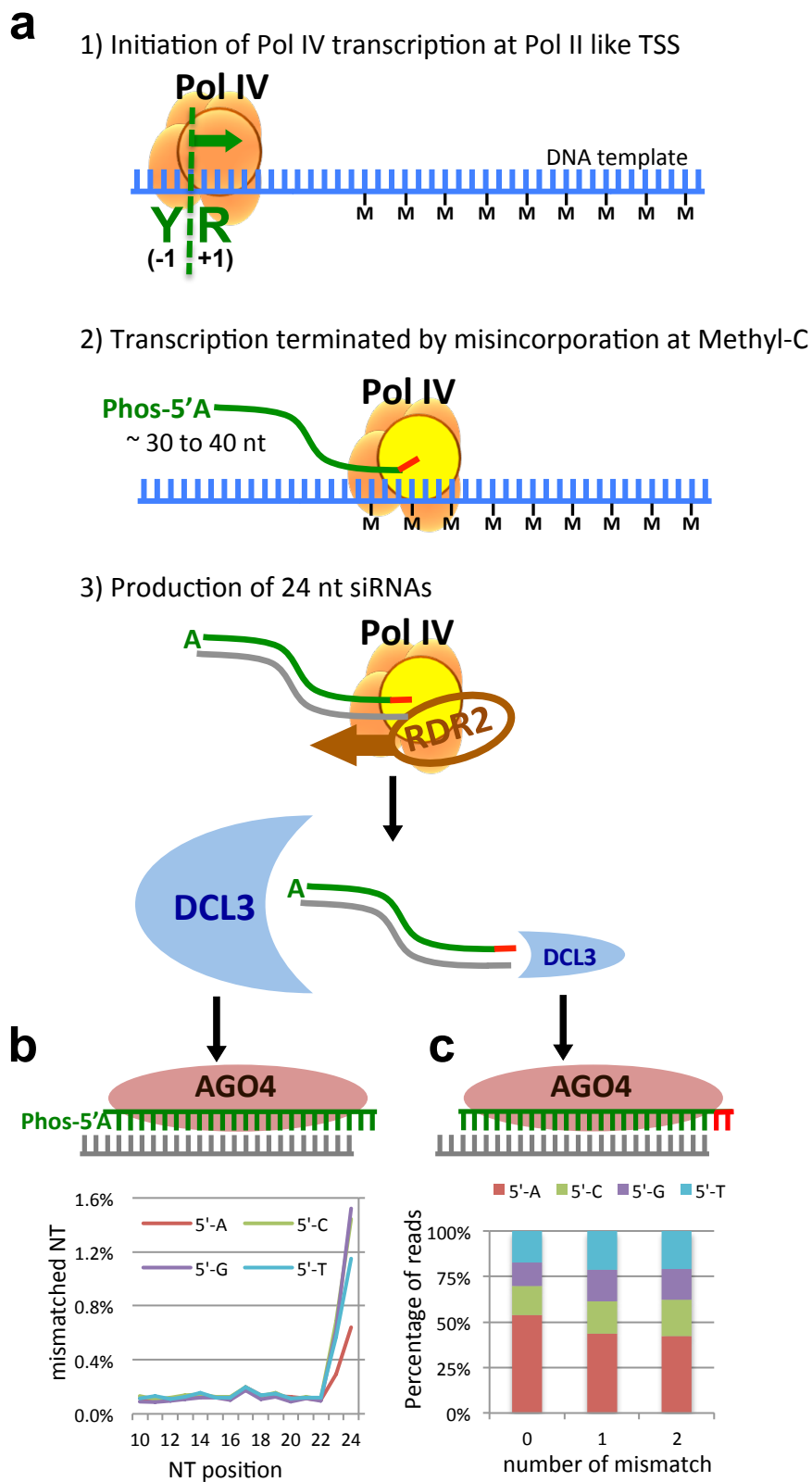


Figure 7. “one precursor, one siRNA” model for the biogenesis of Pol IV dependent 24-nt siRNAs.

- Pol IV transcription is initiated at Pol II-like TSSs. A short RNA of ~30 to 40 nt (with 5'-adenine preference) is produced by Pol IV at heterochromatic regions; misincorporation then occurs at the cytosine position (red bar) and terminates Pol IV transcription. This process typically yields a P4RNA with 5' adenine and 3' misincorporation. P4RNA then goes through processing that involves RDR2 synthesizing the complementary strand (gray), DCL3 cutting from the 5' end (major) or 3' end (minor) of the P4RNA, and loading of the P4RNA-derived siRNA strand into AGO4 as the guide-strand. Eventually one P4RNA precursor gives rise to one siRNA that is derived from either its 5' or 3' end. (“M” underneath the DNA template indicates DNA methylation at the heterochromatic region.)
- When DCL3 cuts from the 5' end of P4RNAs, the resulting siRNAs are more likely to carry a 5' adenine and be perfectly matched. Indeed a reduction of 3' mismatches is observed for siRNAs with a 5' adenine.
- When DCL3 cuts from the 3' end of P4RNAs, the resulting siRNAs are more likely to carry the 3' misincorporation and less likely the 5' adenine. This is consistent with the observation that siRNAs with mismatch have lower percentage of 5' adenine.

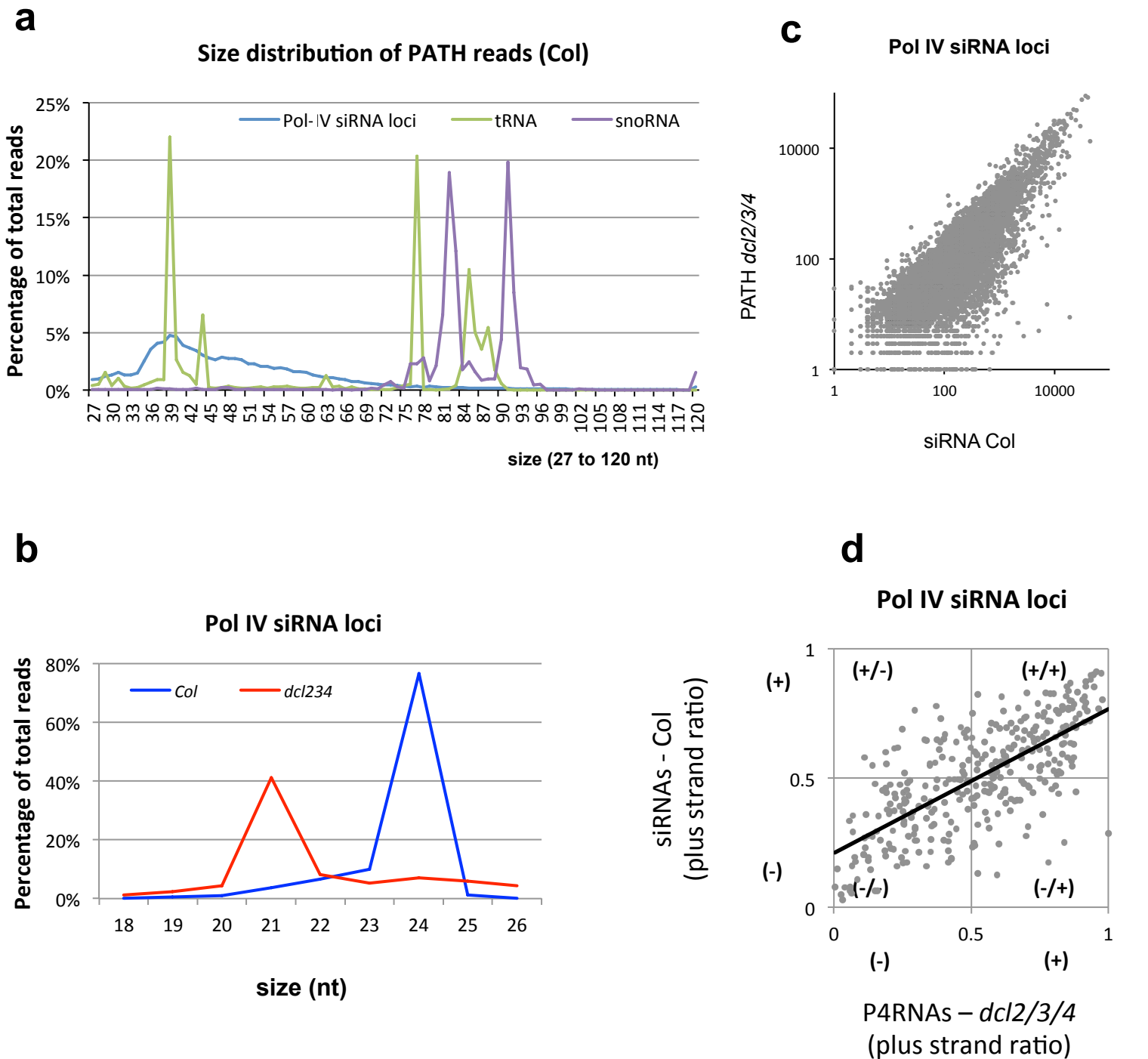


Figure S1. Size distribution of PATH reads and siRNAs

- Size distributions of PATH reads aligned to Pol IV siRNA loci, tRNA loci, and snoRNA loci were shown as the percentage of total read count of that class of loci.
- In Col the siRNAs from Pol IV siRNA loci are predominantly 24-nt in size, while in the *dcl2/3/4* mutant the peak of siRNAs is shifted to 21-nt.
- Correlation of P4RNAs in *dcl2/3/4* with siRNAs in Col at Pol IV siRNA loci.
- P4RNAs and siRNAs share the same strand bias. Shown is the same set of Pol IV siRNA loci that were shown in Figure 1d (bins with no less than 100 P4RNA reads in Col). The plus-strand ratio was calculated as the abundance of reads matching to the plus strand divided by the total number of reads at that locus.

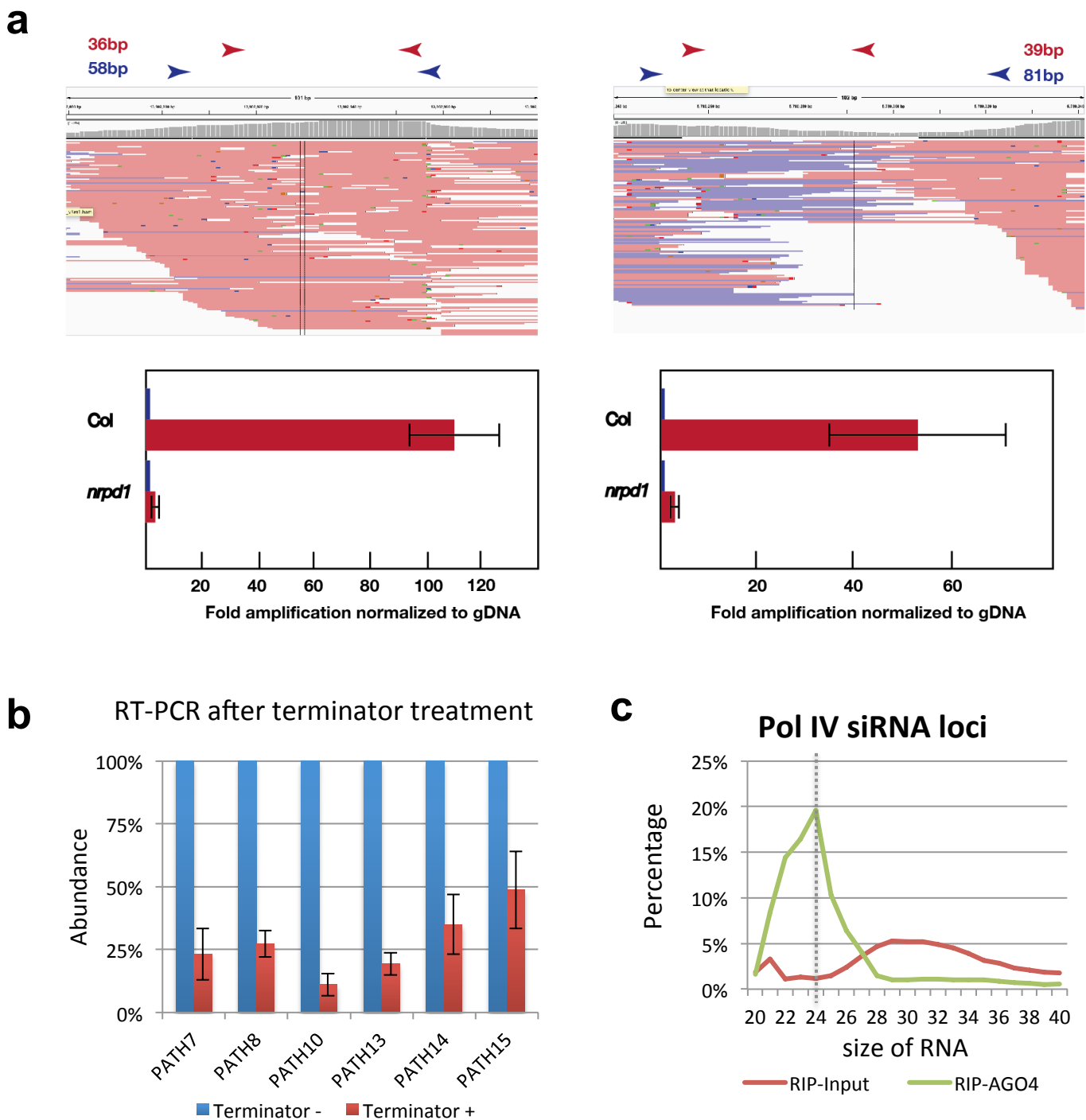


Figure S2. 5' modification and size of Pol IV transcripts.

- Graphical representation of two Pol IV target loci. Red and blue arrows represent primers. Fold change amplification measured by quantitative PCR. Red bars represent the fold change in amplification of the small region (red primers) versus bigger region (blue primers) normalized against amplification of the same primers on genomic DNA. Error bars represent standard error across four replicate. Primer information can be found in Table S3.
- Terminator exonuclease specifically digest RNAs with 5'-monophosphate. Six P4RNA loci were examined by RT-PCR using templates with and without terminator treatment. A large portion of P4RNAs have a 5' monophosphate, evidenced by their degradation via treatment with Terminator.
- Size distribution of RNAs from AGO4 bound compared to input. Experiment was done in the *dcl2/3/4* background.

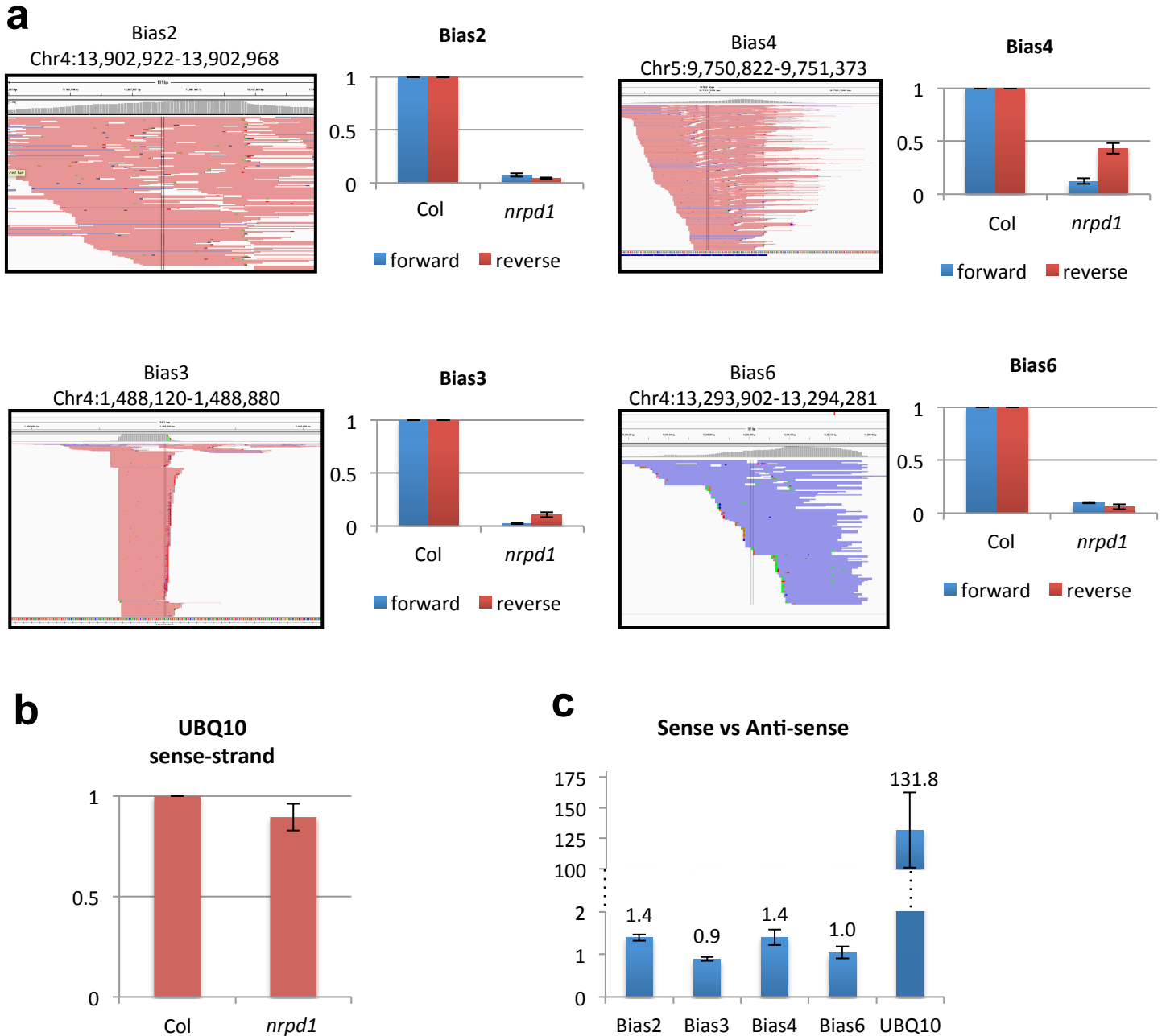


Figure S3. Strand-biased P4RNA loci

- Four loci that exhibit strand-biased accumulation of P4RNAs in the *dcl2/3/4* PATH library were examined by strand-specific RT-PCR. Red bars are reads mapped to Watson strand and blue bars map to Crick strand. Amplifications from both sense and anti-sense strands are Pol IV-dependent.
- Relative abundance of UBQ10 in Col and *nrpd1*. UBQ10 is not Pol IV-dependent and served as a control.
- The relative ratio of sense versus anti-sense amplification is calculated assuming equal reverse transcription efficiency for sense and anti-sense strand at each locus. UBQ10 is used as a control to illustrate clear strand preference.

Primer information can be found in Table S3. Standard errors were calculated from four biological replicates.

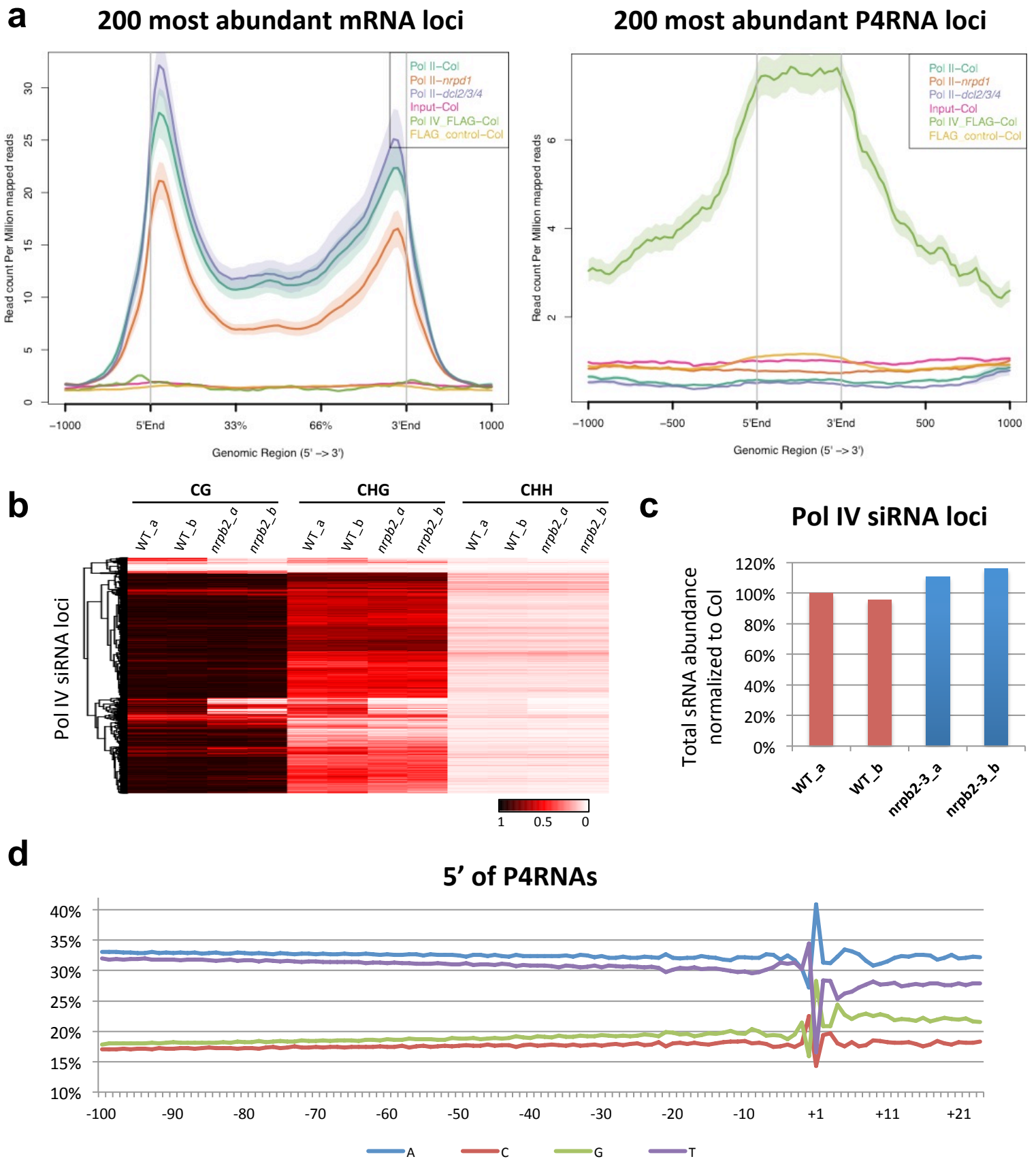


Figure S4. Pol II and Pol IV occupy distinct genomic loci.

- ChIP-seq of Pol II and Pol IV in various genetic backgrounds at either the top expressing mRNA or top expressing P4RNA loci.
- DNA methylation at Pol IV siRNA loci in the weak Pol II mutant (*nrpb2-3*) revealed by BS-seq. Two biological replicates *_a _b* were performed.
- Abundance of 24-nt siRNAs at Pol IV siRNA loci, first normalized to the total number of mapped reads in each library then compared to Col (WT_a). Two biological replicates *_a _b* were performed.
- Nucleotide composition at up to 100 nt upstream of P4RNA TSSs.

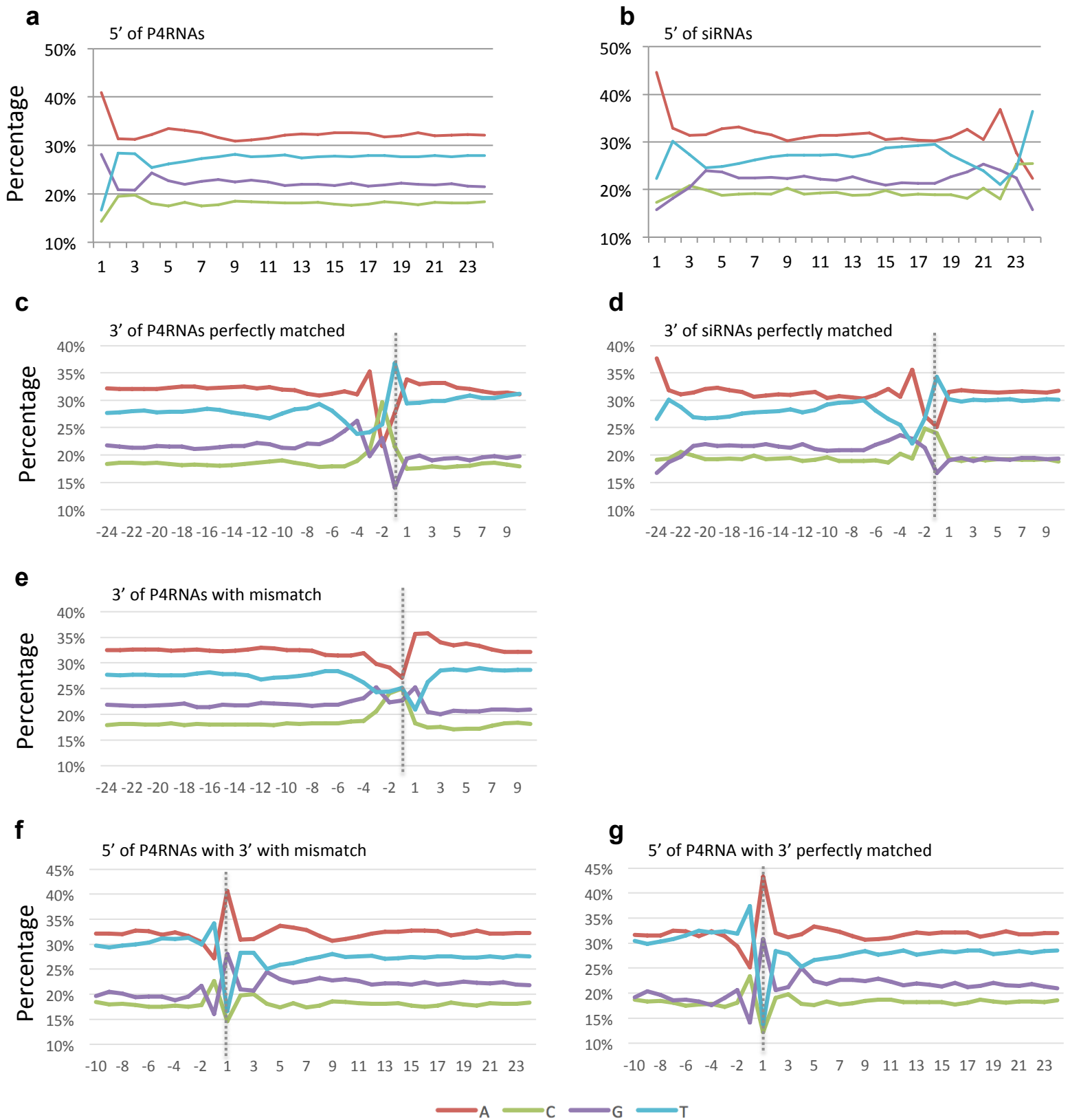


Figure S5. 5' and 3'-end sequence composition of P4RNA and siRNA.

The 5'-end sequence composition of the first 24 nt of P4RNA (a) and the entire 24-nt siRNA (b). The 3'-end sequence composition of perfectly matched P4RNA (c) and 24-nt siRNA (d). The 3'-end sequence composition of mismatched P4RNA (e). P4RNAs with 3' non-templated nucleotides (f) have similar 5' end composition as those without (g). For 3' analysis, “-1” position, indicated by the dashed line, is the last nucleotide of RNA that is perfectly matched, “1” position is the next nucleotide of in genomic DNA sequence; For 5' analysis, “1” position, indicated by the dashed line, is the first nucleotide of RNA.

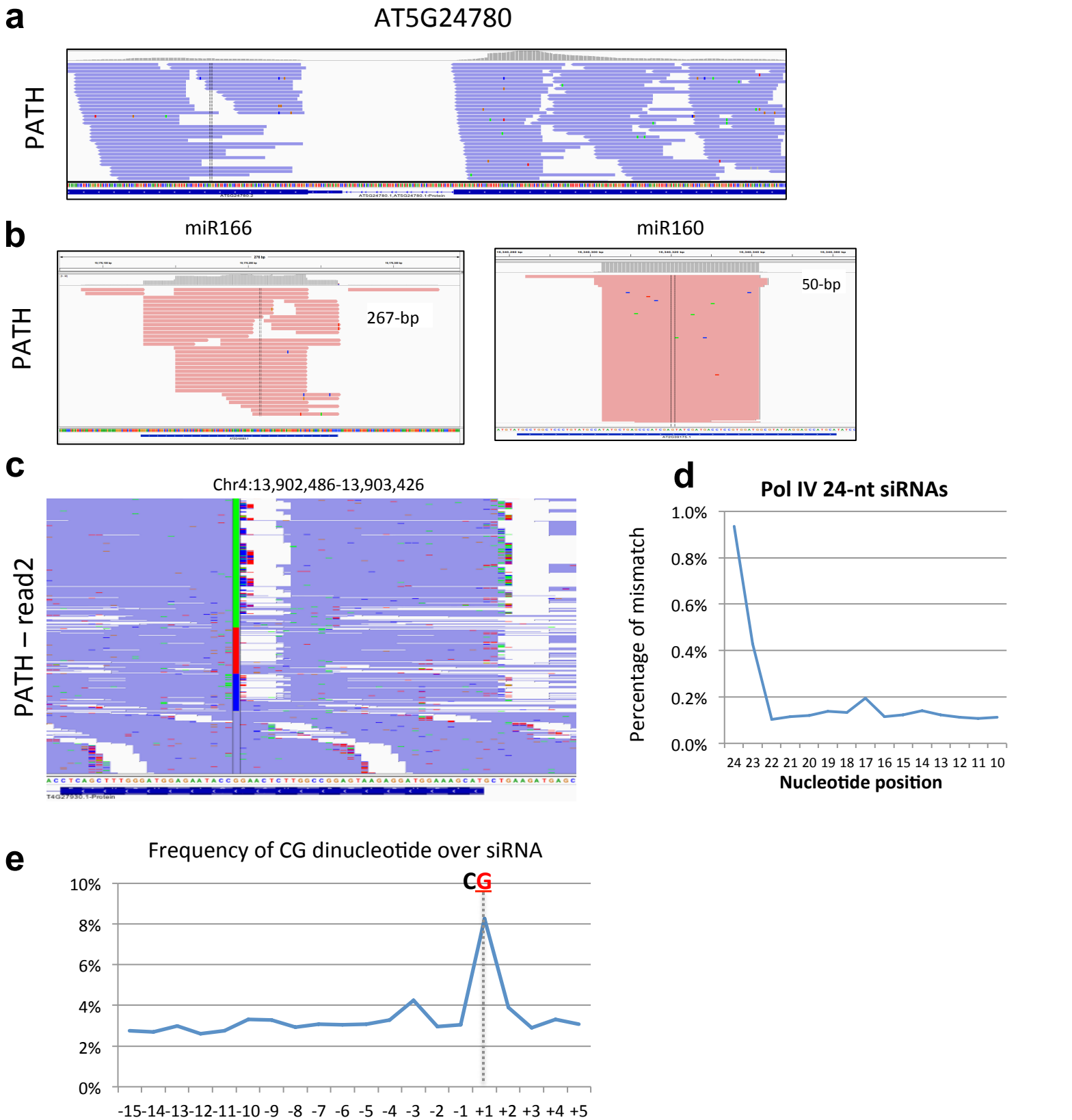


Figure S6. PATH reads from Pol II transcribed regions do not enrich for 3' mismatched nucleotides.

- Example of an mRNA region.
- Examples of miRNA loci.
- Read2 of the paired-end sequencing from PATH library. The beginning nucleotide of read2 corresponds to the 3' end nucleotide of RNA.
- Percentage of mismatched nucleotides on 3'-ends of 24-nt siRNAs. Position 24th represents the last position of the 24 nt siRNA.
- Frequency of CG dinucleotide on reference sequence over the siRNA with misincorporation, "-1" marks the last perfectly matched position, and "+1" marks the first mismatched position. The count of each CG is designated to the position of the G, therefore the peak at "+1" represents a peak of CG at "-1/+1".

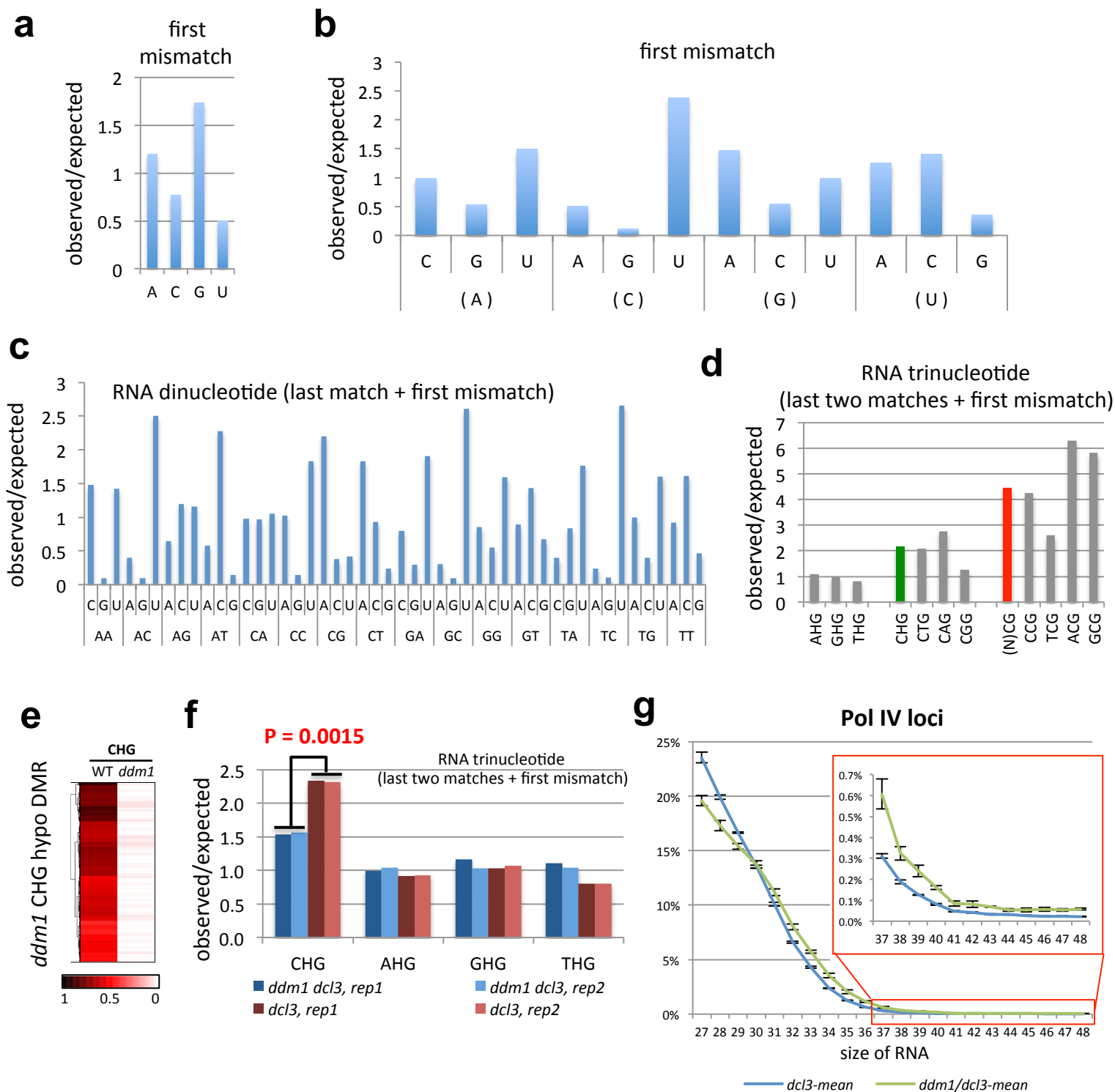


Figure S7. 3'-end mismatch nucleotide composition.

- Mono-nucleotide enrichment at the first mismatched position.
- Relative enrichment for different type of mismatches – nucleotide enclosed in parentheses is the supposed perfectly matched nucleotide, different nucleotide above represents different mismatches.
- Relatively enrichment for different type of mismatches, corresponding to the last nucleotide of the di-nucleotide (last perfectly-matched + first matched).
- Tri-nucleotide enrichment at the first mismatched position.
- CHG hypomethylated DMRs in *ddm1* compared to a wild-type Col control.
- 3'-misincorporated nucleotides of P4RNA at *ddm1* CHG DMRs in *ddm1 dcl3* compared to *dcl3*.
- P4RNAs are slightly longer in *ddm1/dcl3* compared to *dcl3*. Mean and standard error were calculated from three biological replicates of each genotype.