





Charophytic Green Algae Encode Ancestral Polymerase IV/Polymerase V Subunits and a CLSY/DRD1 Homolog

Tania Chakraborty ^{1,†}, Joshua T. Trujillo ^{2,3,†}, Timmy Kendall ¹, Rebecca A. Mosher ^{1,*}

¹The School of Plant Sciences, University of Arizona, Tucson, USA

²Department of Molecular and Cellular Biology, University of Arizona, Tucson, USA

³Department of Biochemistry, Purdue University, West Lafayette, USA

*Corresponding author: E-mail: rmosher@arizona.edu.

[†]These authors contributed equally.

Accepted: June 03, 2024

Abstract

In flowering plants, euchromatic transposons are transcriptionally silenced by RNA-directed DNA Methylation, a small RNA-guided de novo methylation pathway. RNA-directed DNA Methylation requires the activity of the RNA Polymerases IV and V, which produce small RNA precursors and noncoding targets of small RNAs, respectively. These polymerases are distinguished from Polymerase II by multiple plant-specific paralogous subunits. Most RNA-directed DNA Methylation components are present in all land plants, and some have been found in the charophytic green algae, a paraphyletic group that is sister to land plants. However, the evolutionary origin of key RNA-directed DNA Methylation components, including the two largest subunits of Polymerase IV and Polymerase V, remains unclear. Here, we show that multiple lineages of charophytic green algae encode a single-copy precursor of the largest subunits of Polymerase IV and Polymerase V, resolving the two presumed duplications in this gene family. We further demonstrate the presence of a Polymerase V-like C-terminal domain, suggesting that the earliest form of RNA-directed DNA Methylation utilized a single Polymerase V-like polymerase. Finally, we reveal that charophytic green algae encode a single CLSY/DRD1-type chromatin remodeling protein, further supporting the presence of a single specialized polymerase in charophytic green algae.

Key words: RNA-directed DNA methylation, charophytic algae, RNA Pol IV, RNA Pol V.

Significance

All land plants contain duplicate subunits of RNA polymerase II that assemble into functionally unique polymerases and participate in small RNA-mediated methylation of DNA. However, the origin and order of these duplications remain unclear. This study demonstrates that duplicate polymerase subunits first appeared in algal ancestors of land plants, suggesting that small RNA-mediated DNA methylation might predate the evolution of land plants.

Introduction

Eukaryotic genomes encode three highly conserved nuclear DNA-dependent RNA Polymerases (Pol), each responsible for transcription of different subsets of RNA within the cell (Cramer et al. 2008; Werner and Grohmann 2011).

Pol I, Pol II, and Pol III synthesize transcripts that mature into rRNA, mRNA, and tRNA, respectively. These RNA Pols are holoenzyme complexes consisting of two large catalytic subunits and additional smaller noncatalytic subunits that regulate transcriptional activity and RNA processing.

© The Author(s) 2024. Published by Oxford University Press on behalf of Society for Molecular Biology and Evolution.

This is an Open Access article distributed under the terms of the Creative Commons Attribution-NonCommercial License (<https://creativecommons.org/licenses/by-nc/4.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is properly cited. For commercial re-use, please contact reprints@oup.com for reprints and translation rights for reprints. All other permissions can be obtained through our RightsLink service via the Permissions link on the article page on our site—for further information please contact journals.permissions@oup.com.

The largest two subunits of the complex possess highly conserved amino acid residues that constitute the active site, template DNA binding, and DNA–RNA hybrid-binding regions. Each RNA Pol complex is composed of a unique pair of these catalytic subunits, but many of the smaller subunits are shared between complexes. The evolution of distinct catalytic and accessory subunits has driven divergence in function between eukaryotic RNA Pol complexes.

Land plants are unique among eukaryotes in encoding two additional RNA polymerases, Pol IV and V, which are specialized for the production of small interfering RNA (siRNA) and longer noncoding transcripts in the RNA-directed DNA Methylation (RdDM) pathway, respectively (Ream et al. 2014). RdDM is functionally similar to small RNA-mediated transcriptional gene silencing in fission yeast, utilizing an RNA-dependent RNA polymerase, a Dicer endonuclease, and an Argonaute protein to establish DNA methylation de novo (Matzke and Mosher 2014). However, transcriptional silencing in fission yeast uses Pol II for both production of small RNAs and synthesis of noncoding transcripts that recruit the silencing machinery, while these activities are performed separately by Pol IV and Pol V in plants.

Like Pol II, Pol IV and Pol V have 12 subunits, with the largest two forming the catalytic region (Ream et al. 2014). Although most subunits are shared between the three complexes, some are shared between two of the three, while others are unique to a single complex. Pol IV and V are primarily distinguished from Pol II by unique largest subunits, NRPD1 and NRPE1, respectively. Furthermore, Pol IV and V share a second-largest subunit (NRPD/E2), which is distinct from the Pol II second-largest subunit (NRPB2). The catalytic domains of NRPD1, NRPE1, and NRPD/E2 have diverged from NRPB1, but key amino acids and domains associated with transcription are conserved (Ream et al. 2014).

Although clearly homologous in the catalytic region, the largest subunits are significantly diverged at the carboxyl terminal domain that serves as a binding platform for interacting proteins (Ream et al. 2014). NRPB1 contains well-characterized seven amino acid (heptad) repeats, while NRPD1 and NRPE1 instead possess a domain of unknown function (DUF3223 or defective chloroplasts and leaves [DeCL] domain) (Ream et al. 2014). NRPE1 also contains an intrinsically disordered, repetitive, and evolutionarily labile region between the catalytic and DeCL domains (Trujillo et al. 2016). This region is enriched in AGO hooks, peptide motifs that mediate interaction with Argonaute proteins (El-Shami et al. 2007).

Pol IV and Pol V also differ in how they are recruited to DNA. Pol IV requires CLSY proteins, which are part of the SNF2 ATPase family of chromatin remodelers. *Arabidopsis* possesses four CLSY proteins, which recruit Pol IV in a tissue and locus-specific fashion (Zhou et al. 2018). Pol V instead

requires DRD1, the only other member of the CLSY subclade of SNF2 ATPases, for association with DNA (Zhong et al. 2012). While evolutionary histories of other RdDM components have been elucidated, we lack information on the evolutionary origin of CLSY/DRD1 proteins.

Previous studies indicate that Pol IV and V arose prior to the emergence of land plants through the retention of paralogous subunits following the duplication of RNA Pol II subunits (Luo and Hall 2007). Land plants are sister to green algae and together form the Viridiplantae (green plants) lineage, all of which share a common ancestor (Becker and Marin 2009). Green algae consist of several different lineages that are divided among the monophyletic clade of chlorophytic green (Chlorophyta) and a paraphyletic sister group called charophytic green algae (CGA). Together, CGA and land plants form a monophyletic clade known as Streptophyta (de Vries and Archibald 2018) (supplementary fig. S1, Supplementary Material online). Although most Pol IV and V subunits are absent in CGA, a putative NRPD1 homolog was identified in a single CGA order, suggesting a Pol IV-like complex evolved immediately prior to the emergence of land plants (Luo and Hall 2007). However, the origin of this homolog, other Pol IV and V subunits, and Pol-associated RdDM components remains unclear.

Here, we report that a single NRPD1/NRPE1-like sequence is present in multiple CGA lineages. This sequence has many of the structural hallmarks of NRPD1 and NRPE1 but has a carboxy-terminal domain (CTD) similar to Pol V. In addition, we demonstrate that a single CLSY/DRD1 protein is present in the CGAs, raising the possibility that a primitive RdDM utilizing a single specialized polymerase evolved in plants prior to terrestrialization.

Results

CGA Taxa Encode a Single NRPD1/NRPE1 Homolog, and It Is Expressed in *Penium margaritaceum*

To examine whether CGAs possess homologs of NRPD1 and NRPE1, we queried the protein and transcript databases of selected CGA taxa by using the *Arabidopsis thaliana* NRPB1, NRPD1, and NRPE1 sequences. We assessed genomes of both early- and late-diverging CGA taxa, including Klebsormidiophyceae, Charophyceae, and Zygnematomphyceae. Protein BLAST identified NRPB1 homologs in multiple lineages of CGA. These genes were generally well annotated and supported by transcriptome data. In contrast, similar searches for NRPD1 and NRPE1 obtained partial gene predictions that matched mostly to the catalytic region of the protein, where conservation between NRPB1, NRPD1, and NRPE1 is highest (Herr 2005; Kanno et al. 2005; Pontier et al. 2005; Luo and Hall 2007). We identified partial transcripts that share homology with *Arabidopsis* NRPD1 and NRPE1 in *Chlorokybus*

atmophyticus, *Klebsormidium nitens*, *Entransia fimbriata*, *Chara braunii*, *Coleochaete irregularis*, *Spirogyra* sp., *Mougeotia scalaris*, *Mesotaenium kramstae*, *Closterium*, and *Penium margaritaceum* (supplementary table S1, Supplementary Material online). To further investigate these predicted transcripts, we amplified partial cDNA sequences from *P. margaritaceum* and sequenced the resulting fragments. The sequenced cDNA fragments have numerous differences with respect to the predicted transcript but show high similarity to *Arabidopsis* NRPD1 and NRPE1 (supplementary fig. S2 and Data Set S1, Supplementary Material online).

To examine the evolutionary relationship between known Pol IV and Pol V subunits and the predicted first subunit homologs in CGA taxa, we constructed a multisequence alignment and inferred phylogeny using the maximum likelihood method (Fig. 1). Two major clades consist of NRPB1 and NRPD1/NRPE1 homologs from across green plants, with the newly identified CGA homologs diverging prior to the division of land plant NRPD1 and NRPE1. This position indicates that the predicted CGA sequences represent the single-copy ancestor of NRPD1 and NRPE1. Our result agrees with the 2007 Luo and Hall study, and their predicted sequences fall in the same intermediate position as our predicted CGA first subunits. To confirm that the predicted CGA first subunit homologs are not homologs of Pol I or Pol III largest subunits, we also inferred phylogeny from a multisequence alignment containing NRPA1 and NRPC1 homologs. We observed that the NRPA1 and NRPC1 sequences form distinct clades, and the CGA NRPD/E1 sequences remain in the intermediate position between the NRPB1 clade and the land plant NRPD1 and NRPE1 clades (supplementary fig. S3, Supplementary Material online).

CGA First Subunit Homologs Resemble NRPD1 and NRPE1 at Key Catalytic Motifs

NRPD1 and NRPE1 sequences have a higher rate of substitution than NRPB1 sequences, resulting in longer branch lengths (Luo and Hall 2007). Despite this sequence divergence, the metal A site in domain D remains conserved across NRPB1, NRPD1, and NRPE1 (Rymen et al. 2020). The metal A site, which binds a magnesium ion near the polymerase active site, is essential for Pol IV and Pol V function, as substituting the aspartic acid residues in *Arabidopsis* NRPD1 or NRPE1 causes depletion of siRNA production and loss of DNA methylation (Haag et al. 2009). To examine whether the CGA first subunit homologs also possess the metal A site, we searched the multisequence alignment and found the “DFDGD” sequence that represents the core of the metal A site in all the CGA first subunits (Fig. 2). The strict conservation of these five residues within the CGA first subunit homologs suggests that

they could assemble into a catalytically active polymerase holoenzyme. Haag et al. noted that the broader conservation pattern around the DFDGD residues is not maintained in land plant NRPD1 and NRPE1 sequences, a pattern that extends into the CGA first subunits we evaluated (Fig. 2).

Unlike metal A site, there is significant divergence between NRPB1 and NRPD1/NRPE1 sequences in other functionally important motifs. The bridge helix domain is important for NRPB1 catalytic activity, and the 35-residue region is highly conserved in NRPB1 sequences across the different kingdoms of life (Hein and Landick 2010; Weinzierl 2010). However, conservation in this region is lost in NRPD1 and NRPE1 (Herr 2005). We examined the bridge helix region in the CGA first subunit homologs and observed substantial divergence from the NRPB1 sequence (Fig. 2).

NRPD1 and NRPE1 homologs also possess substitutions and deletions relative to NRPB1 within the trigger loop region (Ferrafiat et al. 2019; Rymen et al. 2020). Loss of the trigger loop causes errors in Pol II transcription, and the loss of conservation of this region might explain Pol IV’s error-prone transcription (Kaplan et al. 2008; Erhard et al. 2009; Haag et al. 2012; Marasco et al. 2017; Rymen et al. 2020). We examined the trigger loop in CGA largest subunit homologs and observed that it is not conserved in the CGA sequences, making them look more like NRPD1 and NRPE1 than NRPB1 (Fig. 2). Immediately upstream of the trigger loop, NRPD1 and NRPE1 also have a ~190 amino acid deletion of the foot domain (Luo and Hall 2007; Matzke et al. 2015). We confirmed that this deletion is conserved across land plant lineages and also exists in the CGA homologs (Fig. 2). Together, analysis of domains and motifs in the CGA largest subunit homologs shows that these sequences contain the sequence features that are characteristic of Pol IV and Pol V largest subunits, further supporting the phylogenetic placement of these CGA sequences as single-copy ancestors of NRPD1 and NRPE1.

The CGA NRPD1/NRPE1 Homologs Have C-terminal Domains with AGO Hooks

Although NRPD1 and NRPE1 are clearly homologous to NRPB1 throughout their catalytic regions, their CTDs are not homologous and likely arose through a gene fusion event (Ream et al. 2014). The NRPB1 CTD consists of multiple repeats of seven amino acids (YTPTSPS), a motif that is conserved across all eukaryotic NRPB1s. The CTDs of NRPD1 and NRPE1 lack this motif and instead possess a domain related to the DeCL protein. NRPE1 also has an AGO-binding platform, an intrinsically disordered and repetitive region enriched in “AGO hooks” (GW, WG, or GWG peptides; El-Shami et al. 2007; Till et al. 2007; Trujillo et al. 2016). Luo and Hall identified a single additional polymerase subunit in members of the Charales,

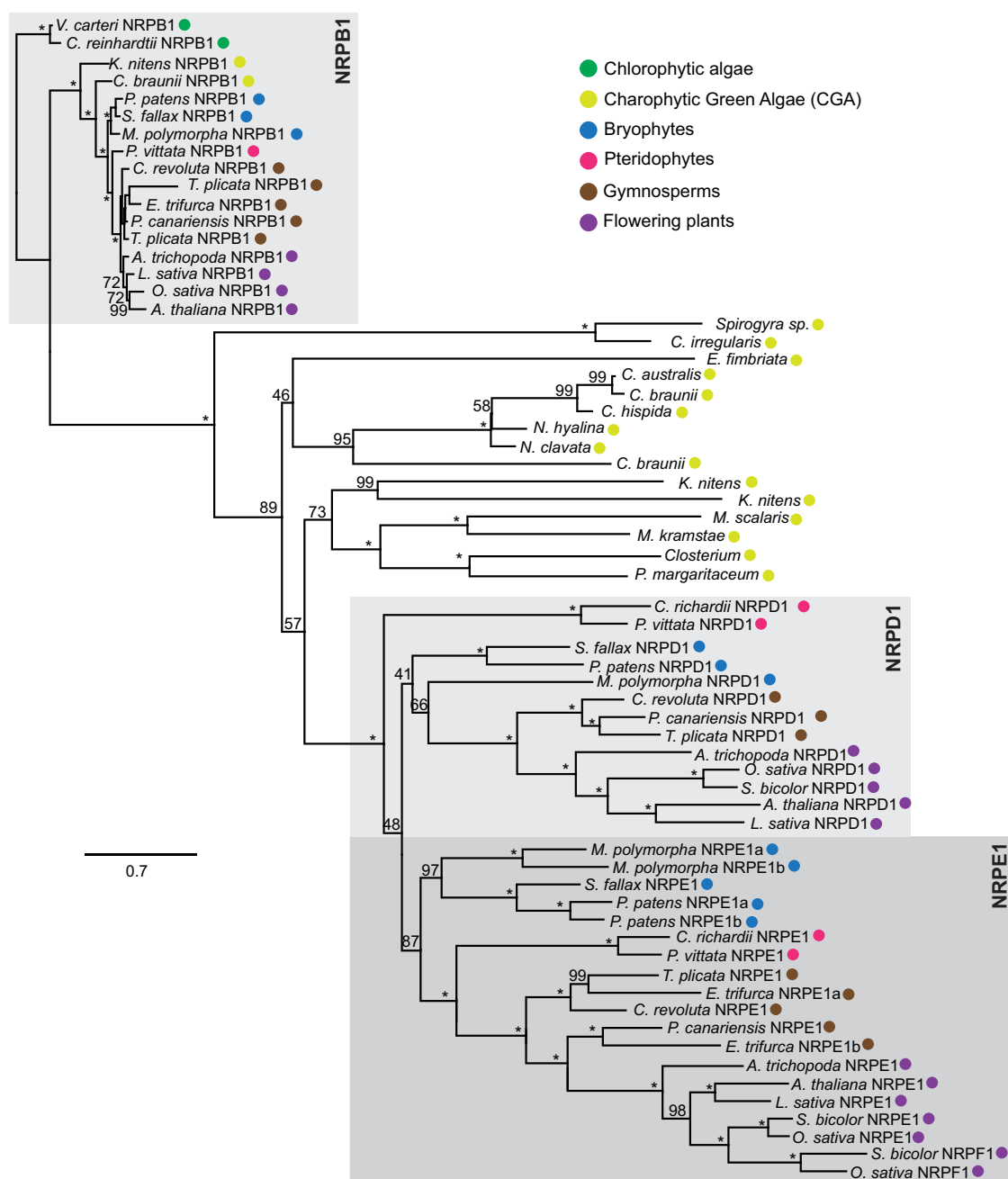


Fig. 1. Phylogenetic analysis of NRPB1, NRPE1, and NRPD1 homologs shows a Pol IV/Pol V-like first subunit in multiple CGA. Amino acid sequences from conserved catalytic regions (C to G domains) were aligned with MAFFT v7.450 and stripped in positions where 50% of the taxa contained a gap. The tree was inferred by maximum likelihood and rooted on chlorophyte NRPB1 sequences. Bootstrap support is listed on each branch (*, 100% support).

and they designated these sequences as (N)RPD1, in part due to their conclusion that NRPE1 sequences were absent from nonflowering plants (Luo and Hall 2007). Subsequently, NRPE1 sequences were identified in all lineages of land plants, raising the possibility that the single additional subunit in CGA lineages might be more similar to NRPE1 in the C-terminal region, or even retain the heptad repeats found in NRPB1. To examine whether the CGA largest

subunit homologs more closely resemble NRPB1, NRPD1, or NRPE1 in the CTDs, we searched for AGO hook motifs and the DeCL domain in regions downstream of the final catalytic domain. We identified numerous AGO hooks in the *K. nitens*, *C. braunii*, *Closterium*, and *P. margaritaceum* homologs, as well as predicted DeCL domains in *K. nitens*, *C. braunii*, and *Closterium* (Fig. 3a). We confirmed the presence of AGO hooks in

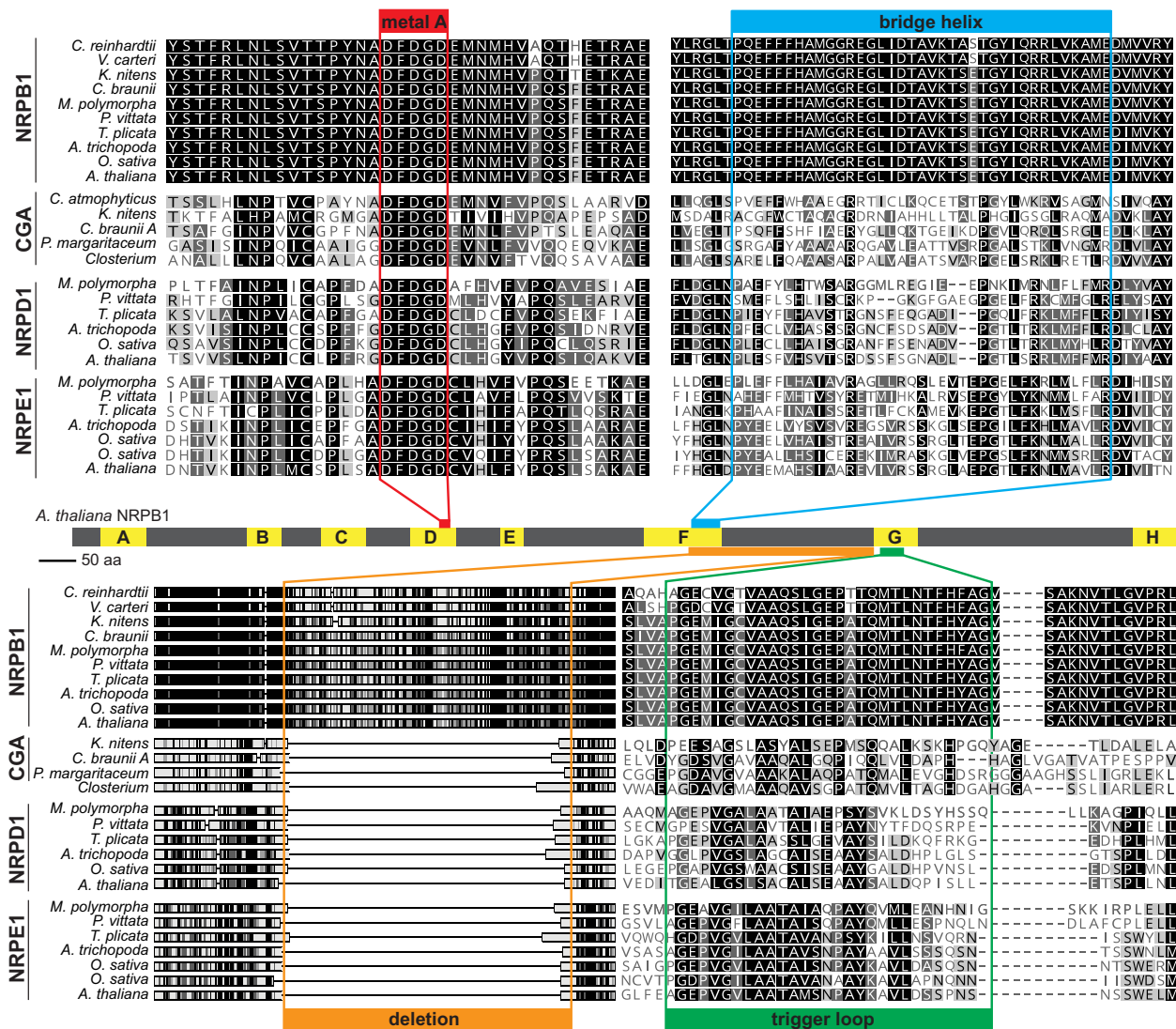


Fig. 2. CGA homologs show hallmarks of NRPD1 and NRPE1 in the catalytic region. Alignment of NRBP1, NRPD1, and NRPE1 sequences from land plants and CGAs demonstrates conservation of the metal A site, loss of conservation in the bridge helix and trigger loop, and a large deletion between domains F and G. All amino acid sequences were aligned with MUSCLE; sequences from a single gene were extracted from the alignment and shaded by similarity in Geneious Prime. The *C. atmophyticus* sequence does not contain the trigger loop region because the recovered sequence is not full-length.

the *P. margaritaceum* transcript through reverse transcription–polymerase chain reaction (RT–PCR) and recovered multiple alleles at the C-terminus (supplementary Data Set S1, Supplementary Material online). However, our *P. margaritaceum* cDNA fragments do not contain stop codons and therefore may not contain the extreme C-terminus. The presence of AGO hooks and a DeCL domain in these homologs indicates that the additional NRBP1 homologs in CGAs more closely resemble NRPE1 than NRPD1 and probably interact with AGO proteins.

Mutations in NRPD1-specific motif (C[KR]YC) reduce siRNA accumulation, with mutations in the cysteine 118

residue being particularly detrimental (Ferrafiat et al. 2019). The NRPD1-specific motif is conserved in angiosperms and at least two gymnosperms (Ferrafiat et al. 2019). To examine whether nonseed plant NRPD1 and the CGA NRPD/E1 homologs possess this motif, we examined the alignment of conserved A to H domains across selected land plant and CGA species. We observed that the C[KR]YC box is conserved across land plants, including bryophytes, but this sequence is not found in the CGA homologs (Fig. 3b). This observation further indicates that the additional largest subunit homologs in CGA more closely resemble NRPE1 than NRPD1.

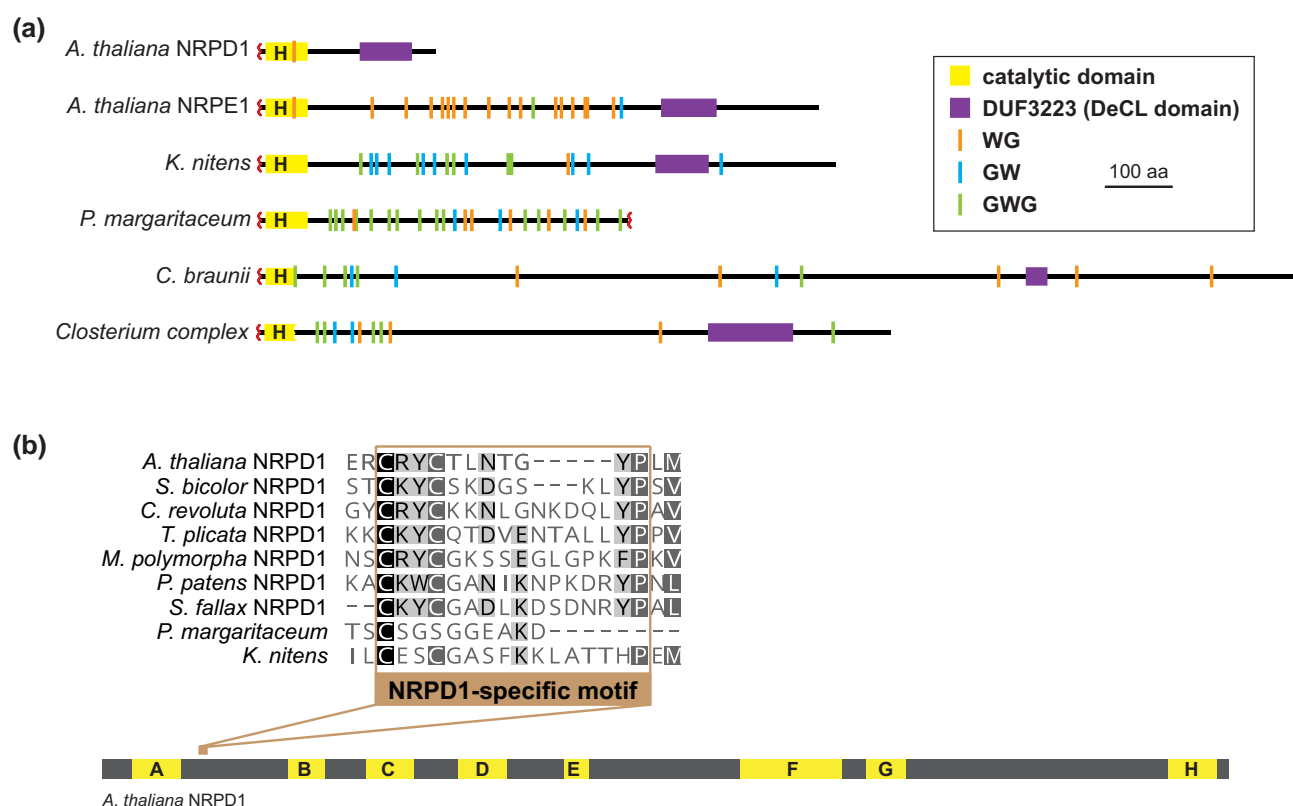


Fig. 3. CGA homologs resemble NRPE1 in the C-terminal domain. a) Diagram of the CTD of *Arabidopsis* NRPD1, *Arabidopsis* NRPE1, and homologous sequences from *K. nitens*, *P. margaritaceum*, *C. braunii*, and *Closterium* demonstrating the presence of numerous AGO hooks (WG, GW, or GWG peptides) between the catalytic region and the DeCL domain. b) Alignment of the NRPD1-specific motif (Ferrafiat et al. 2019) from land plant NRPD1 and homologous CGA sequences demonstrates that conservation of this motif across land plants is absent in the CGA sequences.

Later-diverging CGA Taxa Encode an NRPD/E2 Homolog

Pol IV and Pol V share a second subunit, NRPD/E2, which is conserved across all land plants but was not previously identified in CGA species (Luo and Hall 2007; Huang et al. 2015; Wang and Ma 2015). To examine whether the CGAs possess noncanonical second subunit homologs, we utilized BLAST to identify second subunit homologs in CGA genomes. We identified NRPB2 homologs across all the land plant lineages as well as CGAs and chlorophytic algae, and we identified NRPD/E2 homologs in the tested land plant genomes. We also observed additional second subunit sequences in *E. fimbriata*, *C. braunii*, *M. kramstae*, and *P. margaritaceum*, CGAs from the Zygnematales and Charales orders, suggesting that these later-diverging CGAs might contain a noncanonical second subunit (supplementary table S2, Supplementary Material online). To identify if these sequences are NRPD/E2 homologs, we first inferred phylogeny from maximum likelihood analyses. The topology of these phylogenetic trees indicates that the partial sequences we identified are NRPD/E2-like (Fig. 4). To examine whether the CGA homologous sequences are truly intermediate between the NRPB1 clade and the land plant

NRPD/E2, we gathered NRPA2 and NRPC2 sequences across the land plant lineage as well as from CGA taxa and constructed a multisequence alignment with all the second subunit homologs. We observed that the CGA NRPD/E2 sequences remain outside the well-supported NRPA2 and NRPC2 clades (supplementary fig. S4, Supplementary Material online).

To further examine the identity of these sequences, we assessed characteristic functional domains in the second subunit. Similar to the first subunit homologs, the second subunits of RNA polymerases contain a metal binding site that is responsible for complexing magnesium in the active site. A multisequence alignment of second subunit sequences demonstrates high conservation of the metal B site across all sequences, indicating that the additional CGA second subunit homologs may be catalytically active polymerase subunits (Fig. 5). The hybrid-binding site of NRPB2 binds the nascent DNA–RNA hybrid and interacts with the bridge helix (Cramer et al. 2001). Like the bridge helix, the hybrid-binding site in NRPD/E2 is less conserved than in NRPB2 across the land plant lineage (Herr 2005). We observed that CGA second subunit homologs also have low conservation in this region (Fig. 5).

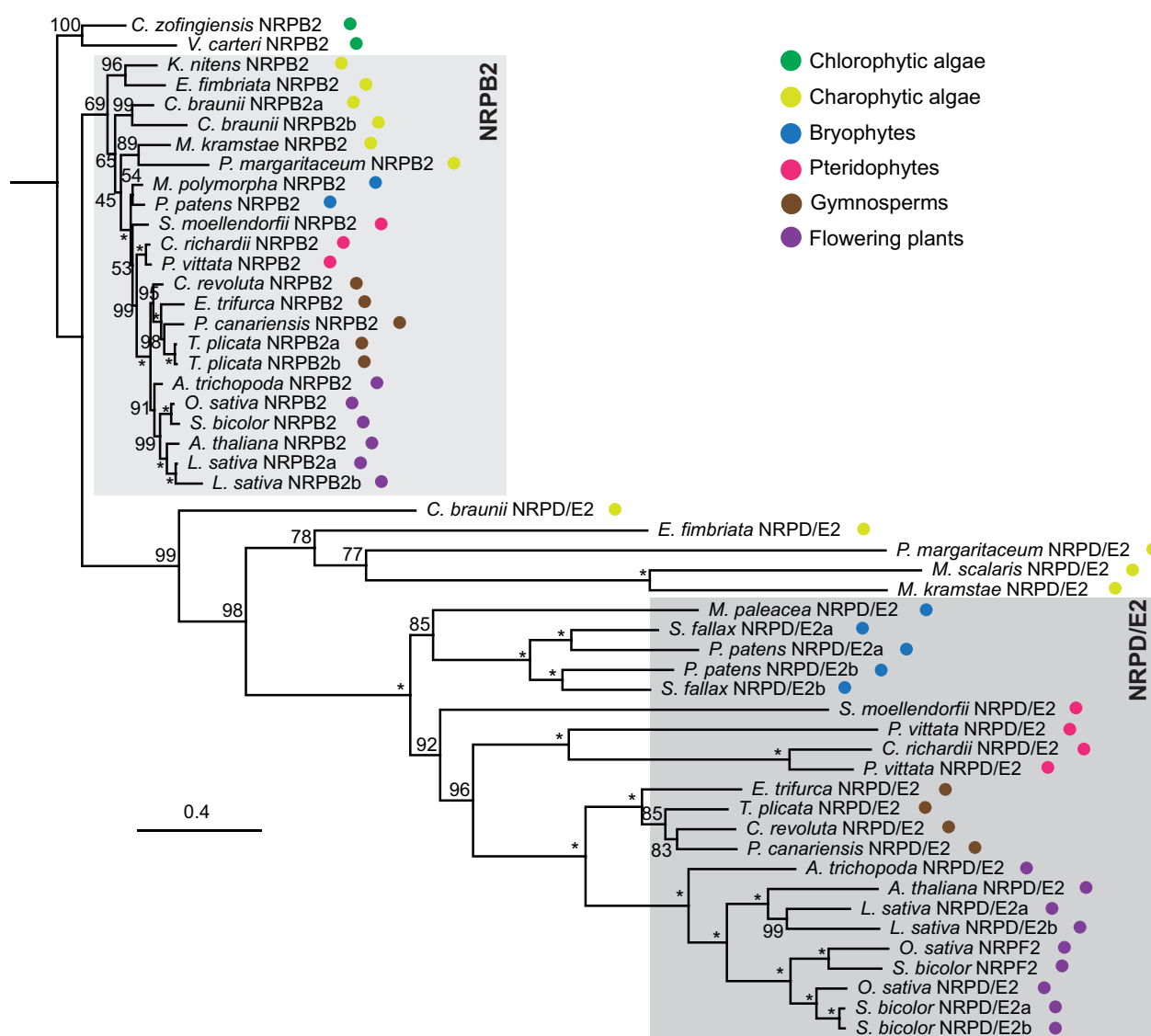


Fig. 4. Phylogenetic analysis of NRPB2 and NRPD/E2 homologs shows a Pol IV/Pol V-like second subunit in multiple CGA. Amino acid sequences were aligned with MAFFT v7.450 and stripped in positions where 80% of the taxa contained a gap. The tree was inferred by maximum likelihood and rooted on chlorophyte NRPB2 sequences. Bootstrap support is listed on each branch (*, 100% support).

Together, these observations support the presence of an NRPD/E2 subunit in at least some CGA lineages. Taken together, our data indicate the first plant-specific homolog of Pol II arose in the CGAs, prior to its duplication into separate Pol IV and Pol V enzymes in land plants.

A Single CLSY/DRD1 Sequence Is Also Present in CGAs and Is Expressed in *Penium*

Pol IV and V require SNF2 ATPase homologs of the CLSY/DRD1 family for their activity. *Arabidopsis* contains four CLSY proteins and the quadruple mutant phenocopies Pol IV mutants with respect to siRNA production

(Smith et al. 2007; Zhou et al. 2022). There are two additional homologs in this small gene family, DRD1 and CHR34 (Kanno et al. 2004; Knizewski et al. 2008). DRD1 is required for Pol V association to DNA (Zhong et al. 2012), and CHR34 has no known role in RdDM. Since CGAs encode Pol IV/Pol V-like first and second subunits, we hypothesized that CGAs would also encode homologs of the CLSY/DRD1 proteins. We used BLAST to gather peptide sequences sharing homology with the *Arabidopsis* CLSY/DRD1 family proteins across land plants and in the CGAs. We observed CLSY and DRD1 homologs across land plants in many CGA taxa. To examine whether this single CLSY/DRD1 protein is expressed, we sequenced

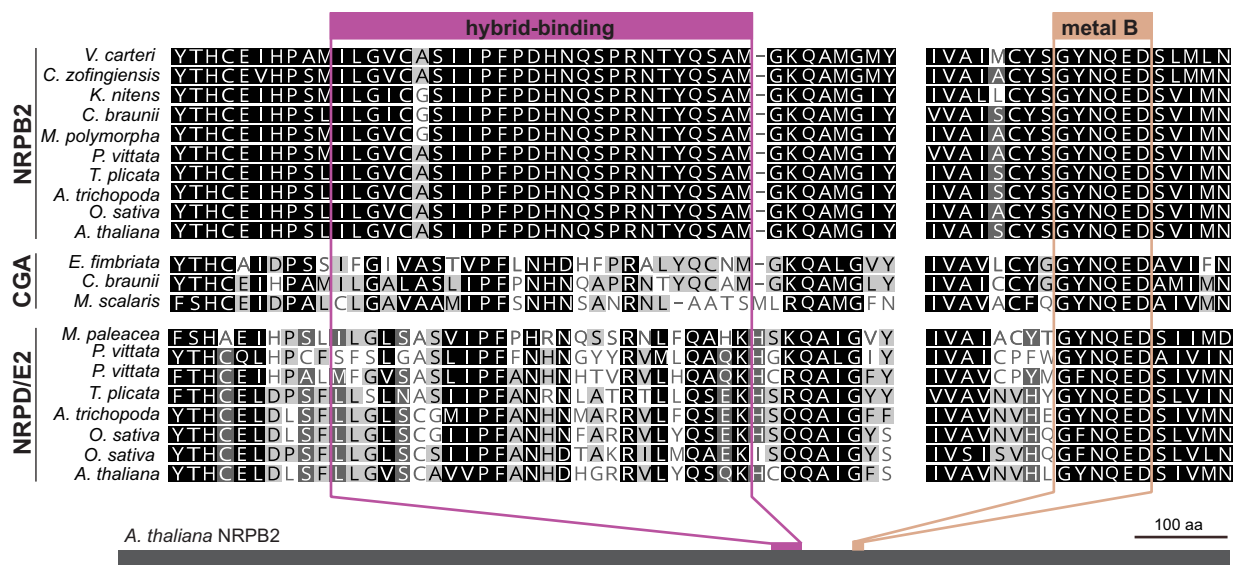


Fig. 5. CGA homologs resemble NRPD/E2 at key motifs. Alignment of conserved domains of Pol II, Pol IV, and Pol V second subunit sequences from land plants and CGAs demonstrates conservation of the metal B site and lack of conservation in the hybrid-binding domain. All amino acid sequences were aligned with MUSCLE; sequences from a single gene were extracted from the alignment and shaded by similarity in Geneious Prime. The partial sequence from *P. margaritaceum* did not contain these regions.

cDNA fragments from *K. nitens* and *P. margaritaceum* (supplementary Data Set S1, Supplementary Material online). The single CLSY/DRD1 was expressed in both species.

To determine the evolutionary relationship between the CGA CLSY/DRD1 homologs, we inferred a phylogeny that included land plant CLSY and DRD1 homologs as well as homologs of RAD54 and ATRX, the next most closely related SNF2 family members (Knizewski et al. 2008). We observed three pairs of related proteins in *Arabidopsis*, representing DRD1/CHR34, CLSY1/CLSY2, and CLSY3/CLSY4. In most vascular plants, these protein pairs exist as single homologs (two CLSY-type and one DRD1), while bryophytes encode a single CLSY homolog and a single DRD1 (Fig. 6 and supplementary fig. S5, Supplementary Material online). CGA genomes reduce this complexity further, encoding only one CLSY/DRD1-like protein, which is positioned sister to the land plant DRD1 and CLSY clades (Fig. 6). This pattern suggests that CGAs, which encode a single Pol IV/V-like polymerase, also encode a single CLSY/DRD1 homolog. The CLSY and DRD1 proteins diverged in the ancestor of bryophytes coincident with the divergence of Pol IV and Pol V.

Discussion

Duplication and diversification of multiple Pol II subunits have led to the presence of Pol IV and Pol V in all land plant lineages (Huang et al. 2015; Wang and Ma 2015). These specialized polymerases produce RNAs responsible for RdDM (Wendte and Pikaard 2017). Here, we demonstrate

that CGA, the sister lineages to land plants, contain the largest and second-largest subunits related to Pol IV and Pol V, as well as a homolog of CLSY/DRD1, which is required for Pol IV and Pol V activity. Together with observation of other RdDM components in CGAs, our results suggest that CGAs might perform an ancient form of RdDM using only a single additional polymerase.

An NRPD1-like transcript was previously identified in two genera within Charales, which at the time was considered the lineage most closely related to land plants (Luo and Hall 2007). However, this study failed to identify NRPD1-like sequences in other CGA orders and found no evidence for NRPD/E2-like sequences in CGAs. Our study identifies first subunit homologs in the four CGA orders most closely related to land plants (Klebsormidiales, Charales, Coleochaetales, and Zygnematales; Fig. 1). We also identified a partial transcript in *C. atrophyticus* (Chlorokybales order) with similarity to Pol IV and Pol V largest subunits, hinting that all CGA orders might contain a specialized polymerase (supplementary table S1 and fig. S3, Supplementary Material online). We identified a second subunit homolog in Charales, Zygnematales, and Klebsormidiales, but not in earlier diverging orders (Fig. 4 and supplementary fig. S4, Supplementary Material online), suggesting asynchronous evolution of the largest two subunits. Similarly, we were unable to identify paralogs of the Pol IV and Pol V seventh subunit in any CGA genome, despite the presence of this subunit in all land plants (Huang et al. 2015). These observations are in line with the presumed increasing elaboration of Pol IV and Pol V



Fig. 6. Phylogenetic analysis of CLSY and DRD1 homologs identifies a single homolog in multiple CGA. Amino acid sequences were aligned with MAFFT v7.450, the conserved SNF2 and helicase domains and intervening sequences were extracted, realigned, and stripped in positions where 50% of the taxa contained a gap. The tree was inferred by maximum likelihood and rooted on RAD54 sequences. Bootstrap support is listed on each branch (*, 100% support).

holoenzyme assemblies over evolutionary time (Luo and Hall 2007; Huang et al. 2015).

Analysis of the conserved functional motifs in the first and second subunits reveals how closely the CGA subunits resemble land plant Pol IV and Pol V, suggesting that these

subunits might function together in a single polymerase (Figs. 2 and 5). Biochemical analysis is necessary to demonstrate such an interaction and to determine whether the resulting polymerase has the enzymatic characteristics of RNA Pol IV or Pol V (Marasco et al. 2017). However, the

C-terminal domain reveals that CGA first subunit homologs are more similar to NRPE1, suggesting that a polymerase formed from this subunit might function like Pol V (Fig. 3).

Our analyses also show that a single CLSY/DRD1 homolog is present in CGA lineages (Fig. 6), paralleling the presence of a single specialized largest polymerase subunit. Other components of RdDM, including RDR2, DCL3, AGO4, and DRM2, also have paralogous sequences in CGA genomes (You et al. 2017; de Mendoza et al. 2018; Wang et al. 2021; Bélanger et al. 2023), raising the possibility that CGAs perform some form of RdDM. The presence of a single specialized polymerase with a Pol V-like tail suggests two models for RdDM in CGAs. First, it is possible that the specialized CGA polymerase produces both siRNA precursors (like Pol IV) and the noncoding RNA scaffold responsible for recruiting siRNA-containing Argonaute complexes (like Pol V). These functions might then have been subfunctionalized into NRPD1 and NRPE1 in land plants. Alternatively, CGA RdDM might utilize Pol II for siRNA production while its specialized polymerase functions like Pol V in land plants. Further research into this intriguing group of algae is needed to reveal the earliest forms of RdDM and expand our understanding of the diverse ways epigenomes are maintained.

Materials and Methods

Identification of Homologs Across Green Plant Lineage

Homologs for all the proteins involved in the study were obtained in *Oryza sativa*, *Sorghum bicolor*, *Lactuca sativa*, and *Thuja plicata* by protein BLAST searches across land plants in phytozome 12 and 13 using *A. thaliana* peptide sequences as query. Sequences in other land plants were obtained either via protein BLAST in phytozome 12 and 13 using *A. thaliana* peptide sequences as query or from submitted data for Huang et al. (2015) (accessed from TreeBASE, study 16473, last accessed on 2023 January 3). For *Cycas revoluta* first subunit sequences, we accessed the shotgun transcriptome data submitted by Huang et al. (EMBL GenBank ID GBJU000000000) and used the transcriptome for search using *A. thaliana* peptides as queries.

For CGA, we utilized a broad range of search techniques. We obtained genomic, transcriptomic, and peptide databases for *K. nitens*, *C. braunii*, *M. kramstae*, *Spirogloea muscicola*, *Zygnema circumcarinatum* (SAG 698-1b, UTEX 1559, and UTEX 1560), and *Zygnema cf. cylindricum* from Phycocosm (<https://phycocosm.jgi.doe.gov/>; Hori et al. 2014; Nishiyama et al. 2018; Grigoriev et al. 2021). Sequences for *P. margaritaceum*, *C. braunii*, and *Closterium peracerosum-strigosum-littorale* were downloaded from respective publications or websites (Nishiyama et al. 2018; Jiao et al. 2020; Sekimoto et al. 2023). These resources

were searched by pBLAST using *Arabidopsis* peptide sequences as query, and resultant peptides were checked by reciprocal BLAST to *Arabidopsis* peptide library (Araport11 protein sequences) on TAIR or to *Arabidopsis* nonredundant protein sequences (nr) on NCBI. Peptide matches were evaluated by checking whether they contained the expected conserved domains and whether they reciprocally matched *Arabidopsis* NRPD1 or NRPE1. Peptide matches that did not meet these criteria are not reported.

Transcriptome searches for CGAs were conducted by downloading 1kP transcriptomes from CyVerse (Carpenter et al. 2019; One Thousand Plant Transcriptomes Initiative 2019) and using TBLASTN with *Arabidopsis* peptide sequences. BLAST hits to genomic sequences were examined by taking 10 to 15 kb of nucleotide sequence around the BLAST hits and finding predicted exons with fgenesh and fgenesh+ as an unbiased approach to exon finding or by using the *Arabidopsis* peptide as reference, respectively. Open reading frames underlying the predicted exons were defined and translated in their frame of reference for generating peptide sequences.

Culturing of *P. margaritaceum* and *K. nitens*

P. margaritaceum (supplied by Professor Jocelyn Rose at Cornell University) and *K. nitens* (University of Texas at Austin Culture Collection, UTEX 623, *Klebsormidium flaccidum*) were cultured in agar slants containing Bristol medium (2.94 mM NaNO₃, 0.17 mM CaCl₂·2H₂O, 0.3 mM MgSO₄·7H₂O, 0.43 mM K₂HPO₄, 1.29 mM KH₂PO₄, and 0.43 mM NaCl) supplemented with UTEX Soilwater: GR+ Medium under constant lighting at 25 °C. The cultures were maintained by transferring once a month.

RT-PCR and Sequencing of *Penium* and *Klebsormidium* Transcripts

Algal culture was scraped from the slant surface, and total nucleic acid was extracted following a protocol adapted from White and Kaper (1989), with the modification of performing the phenol–chloroform–isoamyl alcohol extraction step twice. About 1.1 to 1.3 µg of total nucleic acid was subjected to DNase treatment using Invitrogen DNA-free kit (catalog no. AM1906). 2.5 µg of DNase-treated RNA was used in a 20-µL reaction to convert into cDNA using Thermo Fisher SuperScript IV First-Strand Synthesis System (catalog number 18091050) using random hexamer primers. One microliter of the reaction was then used to amplify predicted transcripts using gene-specific primers. The amplicon was allowed to run out on agarose gel, and the fragment was excised and purified using GeneJET Gel Extraction and DNA Cleanup Micro Kit (catalog number K0831). The fragments were T-cloned using the pGEM-T Easy Vector System I (cat A1360, Promega) and plated on X-gal/IPTG plates. Three white colonies were selected to

undergo restriction digestion with EcoR I to confirm transgene insertion, and upon confirmation, plasmids were sequenced at Plasmidsaurus to generate the cDNA sequence.

Phylogenetic Analysis

Multisubunit alignments were constructed using MAFFT v7.450 on the Geneious Prime (version 2021.2, DotMatics) platform. Alignments were curated manually, and gap trimming was performed using Mask Alignment option in Geneious on sites with gaps in at least 80% of the taxa. Phylogenetic inference was made using the iq-TREE web server using options LG for substitution model, +R free-rate heterogeneity, and Ultrafast Branch Support Analysis using 1,000 bootstrap alignments and default search parameters (Nguyen et al. 2015; Hoang et al. 2018). The same topology was recovered when trees were inferred with ModelFinder and Felsenstein nonparametric bootstrapping (100 bootstraps).

Supplementary Material

Supplementary material is available at *Genome Biology and Evolution* online.

Author Contributions

R.A.M. conceived the research. T.C., J.T.T., and T.K. conducted the research. T.C., J.T.T., and R.A.M. wrote the manuscript.

Funding

Research in the Mosher Lab is supported by the National Science Foundation (IOS-1546825 to R.A.M.); the USDA National Institute of Food and Agriculture (AFRI 2021-67013-33797 to R.A.M.); and the USDA Hatch funding (ARZT 1361510-H25-249 to R.A.M.).

Data Availability

All of the data analyzed in this manuscript are available in public repositories, as described in [supplementary tables S1 to S3, Supplementary Material](#) online. Additional sequences are provided in [supplementary Data Set S1, Supplementary Material](#) online.

Literature Cited

- Becker B, Marin B. Streptophyte algae and the origin of embryophytes. *Ann Bot.* 2009;103(7):999–1004. <https://doi.org/10.1093/aob/mcp044>.
- Bélanger S, Zhan J, Meyers BC. Phylogenetic analyses of seven protein families refine the evolution of small RNA pathways in green plants. *Plant Physiol.* 2023;192(2):1183–1203. <https://doi.org/10.1093/plphys/kiad141>.
- Carpenter EJ, Matasci N, Ayyampalayam S, Wu S, Sun J, Yu J, Jimenez Vieira FR, Bowler C, Dorrell RG, Gitzendanner MA, et al. Access to RNA-sequencing data from 1,173 plant species: the 1000 plant transcriptomes initiative (1KP). *Gigascience.* 2019;8(10):giz126. <https://doi.org/10.1093/gigascience/giz126>.
- Cramer P, Armache KJ, Baumli S, Benkert S, Brueckner F, Buchen C, Damsma GE, Dengl S, Geiger SR, Jasiak AJ, et al. Structure of eukaryotic RNA polymerases. *Annu Rev Biophys.* 2008;37(1):337–352. <https://doi.org/10.1146/annurev.biophys.37.032807.130008>.
- Cramer P, Bushnell DA, Kornberg RD. Structural basis of transcription: RNA polymerase II at 2.8 angstrom resolution. *Science.* 2001;292(5523):1863–1876. <https://doi.org/10.1126/science.1059493>.
- de Mendoza A, Bonnet A, Vargas-Landin DB, Ji N, Li H, Yang F, Li L, Hori K, Pflueger J, Buckberry S, et al. Recurrent acquisition of cytosine methyltransferases into eukaryotic retrotransposons. *Nat Commun.* 2018;9(1):1341. <https://doi.org/10.1038/s41467-018-03724-9>.
- de Vries J, Archibald JM. Plant evolution: landmarks on the path to terrestrial life. *New Phytol.* 2018;217(4):1428–1434. <https://doi.org/10.1111/nph.14975>.
- El-Shami M, Pontier D, Lahmy S, Braun L, Picart C, Vega D, Hakimi MA, Jacobsen SE, Cooke R, Lagrange T. Reiterated WG/GW motifs form functionally and evolutionarily conserved ARGONAUTE-binding platforms in RNAi-related components. *Genes Dev.* 2007;21(20):2539–2544. <https://doi.org/10.1101/gad.451207>.
- Erhard KF Jr, Talbot JE, Deans NC, McClish AE, Hollick JB. RNA polymerase IV functions in paramutation in *Zea mays*. *Science.* 2009;323(5918):1201–1205. <https://doi.org/10.1126/science.1164508>.
- Ferrafiat L, Pflieger D, Singh J, Thieme M, Böhner M, Himber C, Gerbaud A, Bucher E, Pikaard CS, Blevins T. The NRPD1 N-terminus contains a Pol IV-specific motif that is critical for genome surveillance in *Arabidopsis*. *Nucleic Acids Res.* 2019;47(17):9037–9052. <https://doi.org/10.1093/nar/gkz618>.
- Grigoriev IV, Hayes RD, Calhoun S, Kamel B, Wang A, Ahrendt S, Dusheyko S, Nikitin R, Mondo SJ, Salamov A, et al. PhycoCosm, a comparative algal genomics resource. *Nucleic Acids Res.* 2021;49(D1):D1004–D1011. <https://doi.org/10.1093/nar/gkaa898>.
- Haag JR, Pontes O, Pikaard CS. Metal A and metal B sites of nuclear RNA polymerases Pol IV and Pol V are required for siRNA-dependent DNA methylation and gene silencing. *PLoS One.* 2009;4(1):e4110. <https://doi.org/10.1371/journal.pone.0004110>.
- Haag JR, Ream TS, Marasco M, Nicora CD, Norbeck AD, Pasa-Tolic L, Pikaard CS. 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.* 2012;48(5):811–818. <https://doi.org/10.1016/j.molcel.2012.09.027>.
- Hein PP, Landick R. The bridge helix coordinates movements of modules in RNA polymerase. *BMC Biol.* 2010;8(1):141. <https://doi.org/10.1186/1741-7007-8-141>.
- Herr AJ. Pathways through the small RNA world of plants. *FEBS Lett.* 2005;579(26):5879–5888. <https://doi.org/10.1016/j.febslet.2005.08.040>.
- Hoang DT, Chernomor O, von Haeseler A, Minh BQ, Vinh LS. UFBoot2: improving the ultrafast bootstrap approximation. *Mol Biol Evol.* 2018;35(2):518–522. <https://doi.org/10.1093/molbev/msx281>.
- Hori K, Maruyama F, Fujisawa T, Togashi T, Yamamoto N, Seo M, Sato S, Yamada T, Mori H, Tajima N, et al. *Klebsormidium flaccidum* genome reveals primary factors for plant terrestrial adaptation. *Nat Commun.* 2014;5(1):3978. <https://doi.org/10.1038/ncomms4978>.
- Huang Y, Kendall T, Forsythe ES, Dorantes-Acosta A, Li S, Caballero-Pérez J, Chen X, Arteaga-Vázquez M, Beilstein MA, Mosher RA. Ancient origin and recent innovations of RNA polymerase IV and

- V. *Mol Biol Evol.* 2015;32(7):1788–1799. <https://doi.org/10.1093/molbev/msv060>.
- Jiao C, Sørensen I, Sun X, Sun H, Behar H, Alseekh S, Philippe G, Palacio Lopez K, Sun L, Reed R, et al. The *Penium margaritaceum* genome: hallmarks of the origins of land plants. *Cell.* 2020;181(5):1097–1111.e12. <https://doi.org/10.1016/j.cell.2020.04.019>.
- Kanno T, Huettel B, Mette MF, Aufsatz W, Jaligot E, Daxinger L, Kreil DP, Matzke M, Matzke AJ. Atypical RNA polymerase subunits required for RNA-directed DNA methylation. *Nat Genet.* 2005;37(7):761–765. <https://doi.org/10.1038/ng1580>.
- Kanno T, Mette MF, Kreil DP, Aufsatz W, Matzke M, Matzke AJ. Involvement of putative SNF2 chromatin remodeling protein DRD1 in RNA-directed DNA methylation. *Curr Biol.* 2004;14(9):801–805. <https://doi.org/10.1016/j.cub.2004.04.037>.
- Kaplan CD, Larsson K-M, Kornberg RD. The RNA polymerase II trigger loop functions in substrate selection and is directly targeted by alpha-amanitin. *Mol Cell.* 2008;30(5):547–556. <https://doi.org/10.1016/j.molcel.2008.04.023>.
- Knizewski L, Ginalski K, Jerzmanowski A. Snf2 proteins in plants: gene silencing and beyond. *Trends Plant Sci.* 2008;13(10):557–565. <https://doi.org/10.1016/j.tplants.2008.08.004>.
- Luo J, Hall BD. A multistep process gave rise to RNA polymerase IV of land plants. *J Mol Evol.* 2007;64(1):101–112. <https://doi.org/10.1007/s00239-006-0093-z>.
- Marasco M, Li W, Lynch M, Pikaard CS. Catalytic properties of RNA polymerases IV and V: accuracy, nucleotide incorporation and rNTP/dNTP discrimination. *Nucleic Acids Res.* 2017;45(19):11315–11326. <https://doi.org/10.1093/nar/gkx794>.
- Matzke MA, Kanno T, Matzke AJM. RNA-directed DNA methylation: the evolution of a complex epigenetic pathway in flowering plants. *Annu Rev Plant Biol.* 2015;66(1):243–267. <https://doi.org/10.1146/annurev-arplant-043014-114633>.
- Matzke MA, Mosher RA. RNA-directed DNA methylation: an epigenetic pathway of increasing complexity. *Nat Rev Genet.* 2014;15(6):394–408. <https://doi.org/10.1038/nrg3683>.
- Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. IQ-TREE: a fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol Biol Evol.* 2015;32(1):268–274. <https://doi.org/10.1093/molbev/msu300>.
- Nishiyama T, Sakayama H, de Vries J, Buschmann H, Saint-Marcoux D, Ullrich KK, Haas FB, Vanderstraeten L, Becker D, Lang D, et al. The *Chara* genome: secondary complexity and implications for plant terrestrialization. *Cell.* 2018;174(2):448–464.e24. <https://doi.org/10.1016/j.cell.2018.06.033>.
- One Thousand Plant Transcriptomes Initiative. One thousand plant transcriptomes and the phylogenomics of green plants. *Nature.* 2019;574(7780):679–685. <https://doi.org/10.1038/s41586-019-1693-2>.
- Pontier D, Yahubyan G, Vega D, Bulski A, Saez-Vasquez J, Hakimi MA, Lerbs-Mache S, Colot V, Lagrange T. Reinforcement of silencing at transposons and highly repeated sequences requires the concerted action of two distinct RNA polymerases IV in *Arabidopsis*. *Genes Dev.* 2005;19(17):2030–2040. <https://doi.org/10.1101/gad.348405>.
- Ream TS, Haag JR, Pikaard CS. Plant multisubunit RNA polymerases IV and V. In: Murakami KS, Trakselis MA, editors. *Nucleic acid polymerases*. Berlin, Heidelberg: Springer Berlin Heidelberg; 2014. p. 289–308.
- Rymen B, Ferrafiat L, Blevins T. Non-coding RNA polymerases that silence transposable elements and reprogram gene expression in plants. *Transcription.* 2020;11(3-4):172–191. <https://doi.org/10.1080/21541264.2020.1825906>.
- Sekimoto H, Komiya A, Tsuyuki N, Kawai J, Kanda N, Ootsuki R, Suzuki Y, Toyoda A, Fujiyama A, Kasahara M, et al. A divergent RWP-RK transcription factor determines mating type in heterothallic *Closterium*. *New Phytol.* 2023;237(5):1636–1651. <https://doi.org/10.1111/nph.18662>.
- Smith LM, Pontes O, Searle I, Yelina N, Yousafzai FK, Herr AJ, Pikaard CS, Baulcombe DC. An SNF2 protein associated with nuclear RNA silencing and the spread of a silencing signal between cells in *Arabidopsis*. *Plant Cell.* 2007;19(5):1507–1521. <https://doi.org/10.1105/tpc.107.051540>.
- Till S, Lejeune E, Thermann R, Bortfeld M, Hothorn M, Enderle D, Heinrich C, Hentze MW, Ladurner AG. A conserved motif in Argonaute-interacting proteins mediates functional interactions through the Argonaute PIWI domain. *Nat Struct Mol Biol.* 2007;14(10):897–903. <https://doi.org/10.1038/nsmb1302>.
- Trujillo JT, Beilstein MA, Mosher RA. The Argonaute-binding platform of NRPE1 evolves through modulation of intrinsically disordered repeats. *New Phytol.* 2016;212(4):1094–1105. <https://doi.org/10.1111/nph.14089>.
- Wang S, Liang H, Xu Y, Li L, Wang H, Sahu DN, Petersen M, Melkonian M, Sahu SK, Liu H. Genome-wide analyses across Viridiplantae reveal the origin and diversification of small RNA pathway-related genes. *Commun Biol.* 2021;4(1):412. <https://doi.org/10.1038/s42003-021-01933-5>.
- Wang Y, Ma H. Step-wise and lineage-specific diversification of plant RNA polymerase genes and origin of the largest plant-specific subunits. *New Phytol.* 2015;207(4):1198–1212. <https://doi.org/10.1111/nph.13432>.
- Weinzierl ROJ. The nucleotide addition cycle of RNA polymerase is controlled by two molecular hinges in the bridge helix domain. *BMC Biol.* 2010;8(1):134. <https://doi.org/10.1186/1741-7007-8-134>.
- Wendte JM, Pikaard CS. The RNAs of RNA-directed DNA methylation. *Biochim Biophys Acta Gene Regul Mech.* 2017;1860(1):140–148. <https://doi.org/10.1016/j.bbagr.2016.08.004>.
- Werner F, Grohmann D. Evolution of multisubunit RNA polymerases in the three domains of life. *Nat Rev Microbiol.* 2011;9(2):85–98. <https://doi.org/10.1038/nrmicro2507>.
- White JL, Kaper JM. A simple method for detection of viral satellite RNAs in small plant tissue samples. *J Virol Methods.* 1989;23(2):83–93. [https://doi.org/10.1016/0166-0934\(89\)90122-5](https://doi.org/10.1016/0166-0934(89)90122-5).
- You C, Cui J, Wang H, Qi X, Kuo LY, Ma H, Gao L, Mo B, Chen X. Conservation and divergence of small RNA pathways and microRNAs in land plants. *Genome Biol.* 2017;18(1):158. <https://doi.org/10.1186/s13059-017-1291-2>.
- Zhong X, Hale CJ, Law JA, Johnson LM, Feng S, Tu A, Jacobsen SE. DDR complex facilitates global association of RNA polymerase v to promoters and evolutionarily young transposons. *Nat Struct Mol Biol.* 2012;19(9):870–875. <https://doi.org/10.1038/nsmb.2354>.
- Zhou M, Coruh C, Xu G, Martins LM, Bourbonnise C, Lambolez A, Law JA. The CLASSY family controls tissue-specific DNA methylation patterns in *Arabidopsis*. *Nat Commun.* 2022;13(1):244. <https://doi.org/10.1038/s41467-021-27690-x>.
- Zhou M, Palanca AMS, Law JA. Locus-specific control of the de novo DNA methylation pathway in *Arabidopsis* by the CLASSY family. *Nat Genet.* 2018;50(6):865–873. <https://doi.org/10.1038/s41588-018-0115-y>.

Associate editor: Yves Van De Peer