



**HAL**  
open science

# Systematic sequencing of chloroplast transcript termini from *Arabidopsis thaliana* reveals >200 transcription initiation sites and the extensive imprints of RNA-binding proteins and secondary structures

Benoît Castandet, Arnaud Germain, Amber Hotto, David Stern

## ► To cite this version:

Benoît Castandet, Arnaud Germain, Amber Hotto, David Stern. Systematic sequencing of chloroplast transcript termini from *Arabidopsis thaliana* reveals >200 transcription initiation sites and the extensive imprints of RNA-binding proteins and secondary structures. *Nucleic Acids Research*, 2019, 47 (22), pp.11889-11905. 10.1093/nar/gkz1059 . hal-02946563

**HAL Id: hal-02946563**

**<https://hal.inrae.fr/hal-02946563>**

Submitted on 23 Sep 2020

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

**Systematic sequencing of chloroplast transcript termini from *Arabidopsis thaliana* reveals >200 transcription initiation sites and the extensive imprints of RNA-binding proteins and secondary structures**

Benoît Castandet, Arnaud Germain, Amber Hotto, David Stern

► **To cite this version:**

Benoît Castandet, Arnaud Germain, Amber Hotto, David Stern. Systematic sequencing of chloroplast transcript termini from *Arabidopsis thaliana* reveals >200 transcription initiation sites and the extensive imprints of RNA-binding proteins and secondary structures. *Nucleic Acids Research*, Oxford University Press, 2019, 47, pp.11889 - 11905. 10.1093/nar/gkz1059 . hal-02946563

**HAL Id: hal-02946563**

**<https://hal.inrae.fr/hal-02946563>**

Submitted on 23 Sep 2020

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Systematic sequencing of chloroplast transcript termini from *Arabidopsis thaliana* reveals >200 transcription initiation sites and the extensive imprints of RNA-binding proteins and secondary structures

Benoît Castandet<sup>1,2,3,†</sup>, Arnaud Germain<sup>1,†</sup>, Amber M. Hotto<sup>1</sup> and David B. Stern<sup>1,\*</sup>

<sup>1</sup>Boyce Thompson Institute, Ithaca, NY 14853, USA, <sup>2</sup>Institut des Sciences des Plantes de Paris Saclay (IPS2), UEVE, INRA, CNRS, Univ. Paris Sud, Université Paris-Saclay, F-91192 Gif sur Yvette, France and <sup>3</sup>Université de Paris, IPS2, F-91192 Gif sur Yvette, France

Received April 30, 2019; Revised October 02, 2019; Editorial Decision October 25, 2019; Accepted November 05, 2019

## ABSTRACT

Chloroplast transcription requires numerous quality control steps to generate the complex but selective mixture of accumulating RNAs. To gain insight into how this RNA diversity is achieved and regulated, we systematically mapped transcript ends by developing a protocol called Terminome-seq. Using *Arabidopsis thaliana* as a model, we catalogued >215 primary 5' ends corresponding to transcription start sites (TSS), as well as 1628 processed 5' ends and 1299 3' ends. While most termini were found in intergenic regions, numerous abundant termini were also found within coding regions and introns, including several major TSS at unexpected locations. A consistent feature was the clustering of both 5' and 3' ends, contrasting with the prevailing description of discrete 5' termini, suggesting an imprecision of the transcription and/or RNA processing machinery. Numerous termini correlated with the extremities of small RNA footprints or predicted stem-loop structures, in agreement with the model of passive RNA protection. Terminome-seq was also implemented for *pnp1-1*, a mutant lacking the processing enzyme polynucleotide phosphorylase. Nearly 2000 termini were altered in *pnp1-1*, revealing a dominant role in shaping the transcriptome. In summary, Terminome-seq permits precise delineation of the roles and regulation of the many factors involved in organellar transcriptome quality control.

## INTRODUCTION

The sophisticated interplay of factors regulating chloroplast gene expression results from over a billion year symbiosis between the plastid and nucleus. Transcription of the entire plastome (1,2) combined with inefficient termination (3,4) leads to a complex primary transcriptome that undergoes numerous maturation steps. These may include 5' and 3' end processing, intergenic cleavage of polycistronic transcripts, intron removal through RNA splicing and RNA editing to convert specific cytosines into uracils (5–7). Among the protein factors involved in maturation, a variety of endo- and exoribonucleases (RNases) are responsible for processing and cleavage (8), while RNA-binding proteins (RBPs) counteract their activities to stabilize some transcript ends (9) or participate in splicing or editing protein complexes.

Legacy chloroplast RNA maturation analyses have used tedious gene by gene molecular techniques, restricting most detailed RNA analyses to several mono- or polycistronic transcripts including *psbA*, *rbcL*, *atpI-atpH-atpF-atpA* and *psbB-psbT-psbN-psbH-petB-petD* (10), along with the ribosomal RNA operon, which together may not be fully representative of plastid RNA maturation pathways. Gene-by-gene analyses of processing factors may also have limited value. For example, the RNase CSP41a was shown to possess strong endoribonuclease activity *in vitro* (11), however *in vivo* mutant analysis suggests regulatory roles that may have little to do with RNase activity *per se* (12–15). Additionally, while mutants for the endoribonuclease RNase E, its specificity partner RHON1 and the 5'→3' exoribonuclease and endoribonuclease RNase J accumulate novel and apparently unprocessed transcripts, pleiotropic effects are prone to masking their precise sites of cleavage or interaction (16–18).

\*To whom correspondence should be addressed. Tel: +1 607 254 1300; Fax: +1 607 254 1242; Email: ds28@cornell.edu  
Present address: Arnaud Germain, Department of Molecular Biology and Genetics, Cornell University, Ithaca, NY 14853, USA  
†The authors wish it to be known that, in their opinion, the first two authors should be regarded as Joint First Authors.

To overcome some of these limitations, we and others have increasingly employed RNA-seq-based approaches that yield genome-wide cataloging and mechanistic insights into chloroplast transcription, editing, splicing and translation (19–28). Our own results, for example, revealed a large number of non-coding RNAs, which additional research suggests may include a class that exerts its functions through sense-antisense RNA pairing (19,29,30). It is difficult to fully understand transcript function, however, without knowledge of the 5' and 3' termini which together help define promoter sequences, regulatory UTRs and the potential for sense-antisense pairing. On a genome-wide scale, we refer to these 5' and 3' ends as the (RNA) terminome, and the associated technique as Terminome-seq.

Efforts to define the chloroplast RNA terminome have been limited to date, with the most comprehensive study focused on 5' ends in barley, a model that has been effectively used to dissect the respective roles of nucleus- and chloroplast-encoded RNA polymerases (31,32). On a genome-wide level, barley chloroplasts were found to possess larger than expected numbers of both primary and processed 5' termini, consistent with a highly complex transcriptional landscape (20). We chose *Arabidopsis* for our analysis, because of its broad use to dissect post-transcriptional RNA events in the chloroplast, including the analysis of RNase mutants, pentatricopeptide repeat (PPR) and other helical repeat proteins, and RNA editing and splicing factors. We found that the *Arabidopsis* chloroplast terminome is complex, and in some cases surprising. For example, both known and new transcription start sites (TSS) were identified, sometimes internal or antisense to known transcripts and a general imprecision of both processed 5' and 3' ends was observed. To highlight the comparative potential of Terminome-seq, we examined the *pnp1-1* mutant which lacks the major 3' processing enzyme polynucleotide phosphorylase (33,34), revealing a largely reshaped terminome. Overall, our results showcase Terminome-seq as a valuable addition to the organelle gene expression analysis toolkit.

## MATERIALS AND METHODS

### Plant material

*Arabidopsis thaliana* Col-0 and *pnp1-1* seeds were germinated on MS medium with 16 h of light per day at 23°C. Three-week old leaf material was flash-frozen in liquid nitrogen, and total RNA was isolated using TRI Reagent according to the manufacturer's instructions ([www.sigmaaldrich.com](http://www.sigmaaldrich.com)).

### Terminome library synthesis and analysis

All libraries were produced from 1 µg of DNase I-treated RNA ([www.neb.com](http://www.neb.com)), and for TAP-treated samples, tobacco acid phosphorylase (TAP; [www.epibio.com](http://www.epibio.com)) was used according to the manufacturer's instructions with heat inactivation at the end of the incubation period. Library synthesis was carried out using the Illumina TruSeq Small RNA library preparation kit ([www.illumina.com](http://www.illumina.com)) intended to capture the RNA population containing a 5' phosphate and 3' hydroxyl group. Minor modifications were made to

the protocol depending on whether a native 5' or 3' end was the target (Supplementary Figure S1). Libraries intended for native 3' end capture, followed the protocol with initial 3' adapter ligation using T4 RNA ligase 2, a deletion mutant that can only ligate a 3' hydroxyl group to a 5' adenylated RNA, consistent with the 3' RNA adapter chemistry. After ligation, the RNA was fragmented using a Covaris sonicator ([www.covaris.com](http://www.covaris.com)), with a target size of 200 nt, followed by ethanol precipitation for concentration and 5' adapter ligation with T4 RNA ligase. Libraries intended for native 5' end capture required further adjustments. The order of adapter ligation was reversed: 5' adapter ligation (with T4 RNA ligase)—sonication—ethanol precipitation—3' adapter ligation (with T4 RNA ligase 2). In this case, excess 5' adapter remaining following the sonication and ethanol precipitation could ligate to added 3' adapter, but not to any new 5' ends created through sonication as the new 5' ends would not be adenylated. This resulted in unwanted adapter dimers that were preferentially amplified during library amplification (PCR1) due to their small size (~133 bp). Therefore, size selection was performed on the products from PCR1, retaining only products over 200 bp using Pippin Prep ([www.sagescience.com](http://www.sagescience.com)). A second polymerase chain reaction amplification (PCR2) was executed on these products. Quality control was performed after Pippin size selection and before library submission for sequencing using an Agilent BioAnalyzer ([www.agilent.com](http://www.agilent.com)). Details about the procedure and the Pippin Prep are available at <https://github.com/BenoitCastandet/Terminome.Seq>. The final cDNA libraries were purified using magnetic AMPure beads ([www.beckman.com](http://www.beckman.com)) following the manufacturer's protocol. Multiple steps in the above protocol, including fragmentation, ethanol precipitation, Pippin size selection (5' libraries only) and AMPure purification of cDNA libraries, resulted in a bias toward the retention, and therefore sequencing, of fragments >67 nt. As a consequence, ends of small RNAs (smRNAs) and tRNAs would be expected to be underrepresented in the results. Additionally, a minor bias was introduced because the RT primer is fully complementary to the 3' adapter, ending in a sequence complementary to TGG. Illegitimate priming by the adapter resulted in an estimated 52 additional 3' ends terminating in TGG which are included in our data as they cannot be distinguished from legitimate 3' termini ending in T(U)GG. While this bias was inevitable, the overall interpretation of the results was not affected.

Libraries were pooled and sequenced on a NextSeq500 Sequencer ([www.illumina.com](http://www.illumina.com)) using the v3 kit, with paired-end reads generating 40 bp long R1 reads and 35 bp long R2 reads for all libraries. R1 reads are only of use for libraries generated to obtain 5' related data, while R2 reads contain data related to 3' ends and therefore are only relevant for libraries generated to obtain 3' data. Raw sequences have been deposited on the SRA database with the number PRJNA533962 and can be accessed here <https://www.ncbi.nlm.nih.gov/sra/PRJNA533962>. The detailed pipeline used to analyze the relevant reads is available at <https://github.com/BenoitCastandet/Terminome.Seq>. Briefly, the quality of relevant reads was checked using fastq-mcf (<https://github.com/ExpressionAnalysis/ea-utils/blob/wiki/FastqMcf.md>) followed by align-

ment to the chloroplast reference genome, *Arabidopsis* TAIR10 version modified to add the first exon of the chloroplast gene *ycf3*, using tophat2 (<https://ccb.jhu.edu/software/tophat/index.shtml>). Two customized scripts allowed us to extract the positions of the 5' and 3' termini and the results were normalized according to the numbers of reads aligned to the chloroplast genome. Normalized data are available in Supplementary Table S1.

## 5' RACE

Total RNA was isolated from mature leaf tissue using TRI Reagent<sup>®</sup> and treated with DNase I (Ambion; <http://www.thermofisher.com>). 5' Rapid Amplification of cDNA Ends (RACE) reactions were completed as described (35) with some modifications. For analysis of primary transcripts, 4  $\mu$ g of RNA was treated with TAP (0.5U/4  $\mu$ g RNA; Epicenter, <http://www.epibio.com>) for 1 h at 37°C with RNase-OUT (40U; Invitrogen, <http://www.thermofisher.com>), followed by phenol/chloroform extraction and ethanol precipitation. The 5' RACE adapter (Supplementary Table S2) was then ligated to the 5' ends of the + or -TAP treated RNA with T4 RNA ligase (Ambion). For this, 10  $\mu$ M of 5' RACE adapter and 4  $\mu$ g of  $\pm$ TAP-treated RNA were incubated for 5 min at 65°C. Samples were chilled on ice and then 5U of T4 RNA ligase, 1 $\times$  ligase buffer, 1 mM adenosine triphosphate (ATP) and 40U of RNaseOUT were added and the reactions were incubated for 1 h at 37°C. Ligated RNA was phenol/chloroform purified and precipitated with ethanol, and cDNA was generated using SuperScript III (Invitrogen) with random hexamers according to the manufacturer's protocol. Transcripts were amplified by PCR using the 3' gene-specific primers indicated in the figure legends and the 5' RACE primer (Supplementary Table S2). Amplicons were visualized after separation through a 1% agarose gel and gel purified. Purified bands were used for a second, nested PCR reaction with the indicated 3' RACE primer and the RACE 5' nested primer to ensure amplified sequences were specific to the intended target and/or increase amplicon intensity as the bands were directly sequenced, or cloned and sequenced. The nested PCR reactions were visualized in 1% agarose gels, and specific bands were gel purified, cloned and sequenced.

## RESULTS

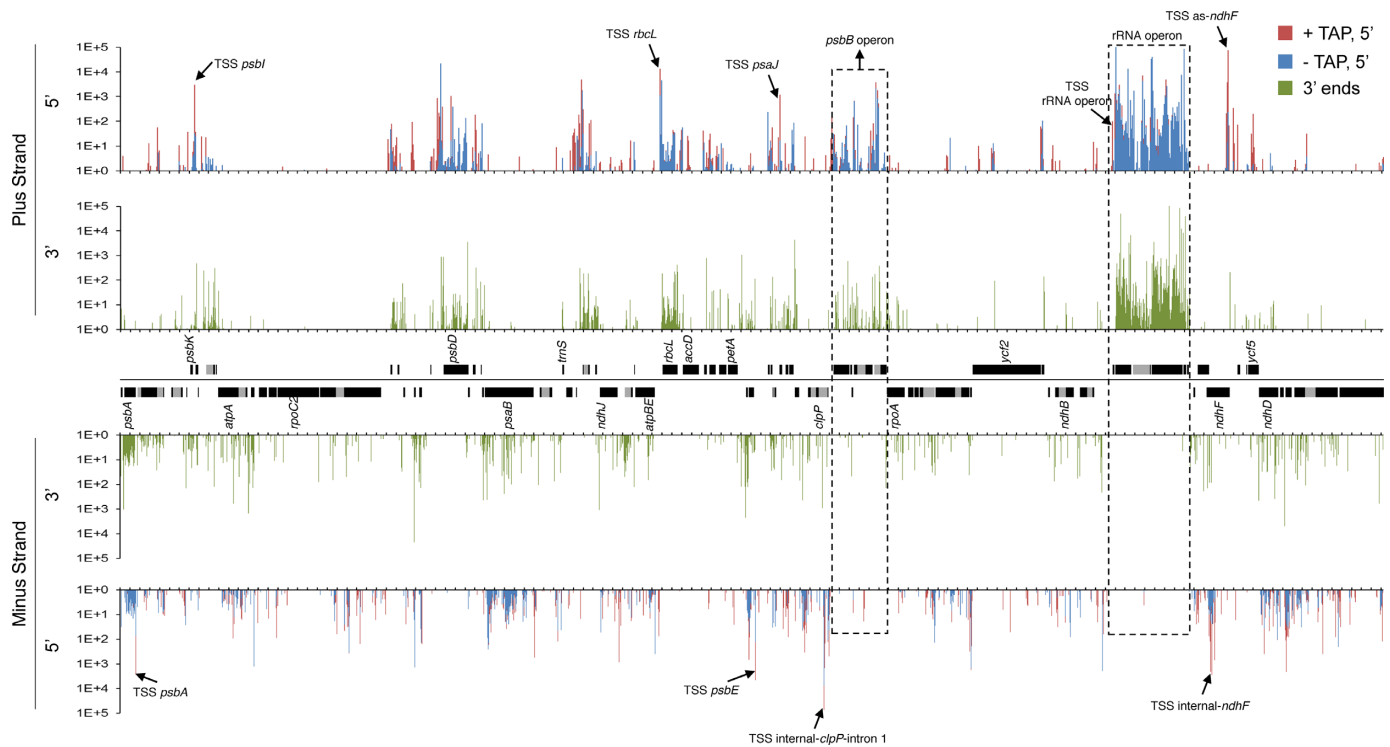
### Terminome sequencing strategy and overview

Identification of RNA termini has traditionally been performed by sequencing individual clones from RACE products, primer extension, or nuclease protection assays. A recent improvement was the substitution of high-throughput sequencing for product-by-product analysis (36–40). Because of their focus on nucleus-encoded RNAs, these protocols rely on the presence of 3' poly(A) tails, rendering them unsuitable for studying organellar RNAs where such tails mark transcripts for degradation (27,41,42). Our experimental design included the sequencing of three distinct libraries devoted to the identification of both primary and processed 5' ends, or 3' ends. The protocol (Supplementary Figure S1) was modified from the Illumina TruSeq Small RNA kit, which relies on sequential ligation

of adaptors prior to cDNA synthesis and amplification. For Terminome-seq, the initial ligation of 5' or 3' adapter captures the native 5' and 3' termini, respectively, then the RNA is fragmented followed by a second adapter ligation step before amplification. To differentiate and identify primary transcripts that correspond to TSS, RNA was pre-treated with TAP prior to 5' adaptor ligation. TAP converts triphosphorylated 5' ends, uniquely found in primary transcripts, to monophosphates, making them amenable to adapter ligation. Thus, libraries created with or without TAP pretreatment (+TAP and -TAP) can be compared to identify TSS. Duplicate libraries were constructed from *Arabidopsis thaliana* Col-0 biological replicates and the results show high correlation (Supplementary Figure S2). Native transcript ends were defined as the first nucleotide of each read, following suitable manipulation of the primary sequencing. The full coverage for Terminome-seq at single nucleotide resolution for 5' ( $\pm$ TAP) and 3' ends generated for this manuscript can be accessed in Supplementary Table S1. It is important to note that transcripts with post-transcriptional modifications, such as tRNAs with a post-transcriptionally added CCA tail, or RNA degradation intermediates bearing polynucleotide tails, would not fully align to the genome and are therefore excluded from our curated dataset.

A genome-wide view of end distribution is depicted in Figure 1. Clusters of 5' and 3' ends, covering a remarkable 13 and 16% of the genome, respectively, are clearly discernable, but are largely absent from areas antisense to known transcripts. In contrast, the full chloroplast genome is covered by reads from total RNA-seq experiments (2,19,20,23). However, when considering only positions representing at least 0.001% of total reads (10 reads per million, or RPM), <0.7% of the genome (1628 processed 5' ends and 1299 3' ends, respectively) was represented (Figure 2A). Thus, the vast majority of termini are stochastically found at low levels, probably representing RNA degradation products or processing intermediates.

To compare the abundance of termini across regions of the chloroplast genome, we informatically divided them into six types: rRNAs, tRNAs, exons (corresponding to protein-coding regions), introns, intergenic regions (anything else) and positions on the antisense strand of exons. When categorized this way, Figure 2B shows that ~85% of the 5' and 3' ends were found in the rRNAs. The high degree of rRNA transcript ends was expected due to the elevated expression level of the rRNA operon and lack of rRNA depletion during the library preparation. Addition of an rRNA depletion step could help to map low abundance transcript ends that may have been overlooked in this study. Ends corresponding to intergenic regions were over-represented compared to exons (Figure 2B). Indeed, even though intergenic regions cover around 25% of the genome (43), they contained 10.2 and 8.8% of the -TAP 5' and 3' ends, respectively, whereas only 1.7 and 1.5% of the ends mapped within exons. TAP treatment, which allowed TSS to be represented, significantly altered these proportions, with the proportion of 5' ends mapping to rRNAs decreasing to 65%, mainly in favor of an increase in reads mapping to introns (12.4%) and antisense to exons (8.3%; Figure 2B). Remarkably, an overabundance of TSS is found in



**Figure 1.** Plastome-scale view of Terminome-seq results. End coverages are the average of two Col-0 biological replicates and given in RPM. 5' ends obtained with or without TAP treatment are red and blue, respectively, and 3' ends are displayed in green. Gene models are indicated between the tracks corresponding to the plus and minus strands of the plastome. One copy of the large inverted repeat is omitted for clarity. Selected TSS described in more detail in the main text are marked by arrows and labeled, as are the *psbB* and rRNA operons. Tick marks are every 1000 nt.

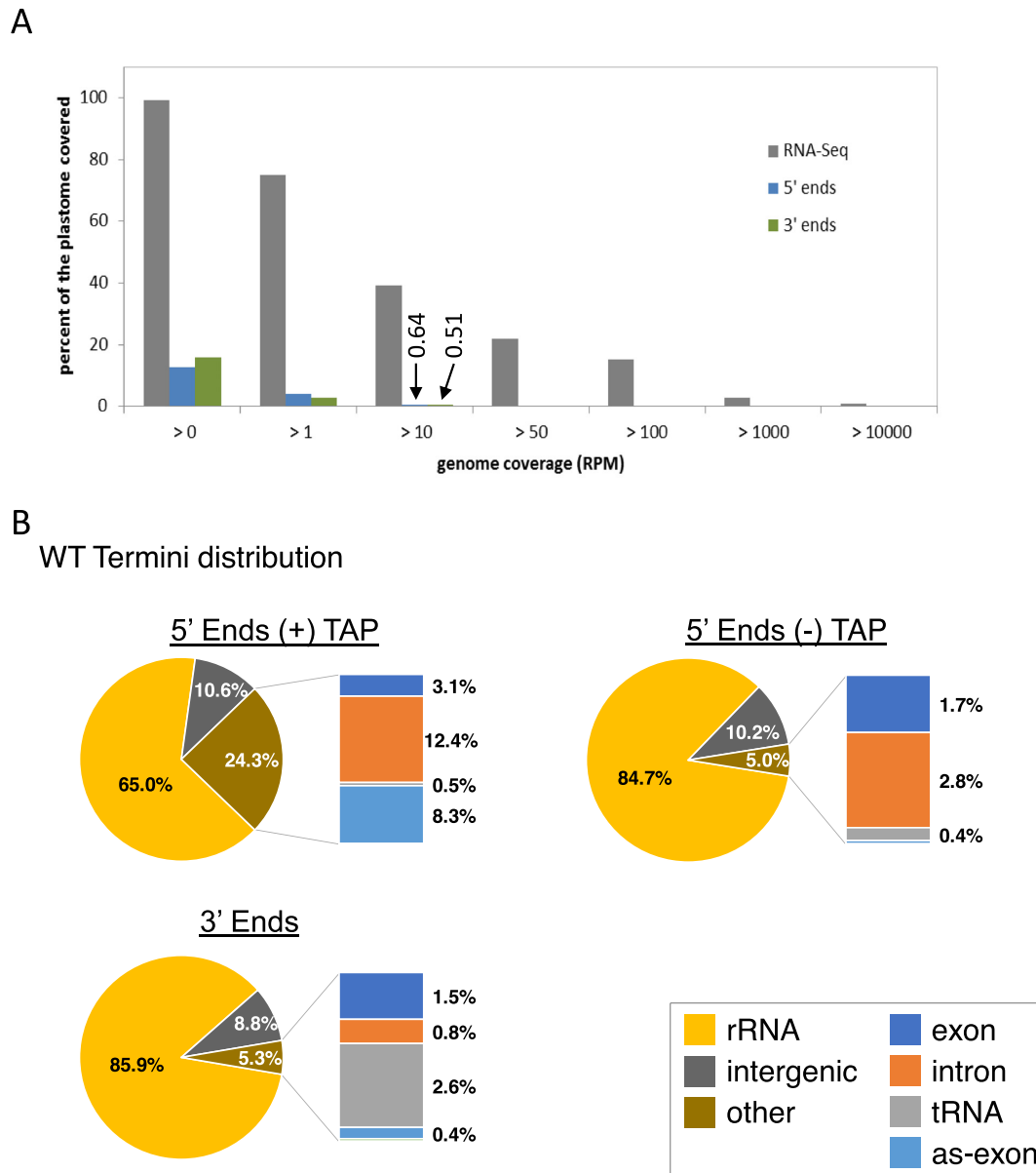
the first *clpP* intron (TSS internal-*clpP*-intron 1) and antisense to the *ndhF* transcript (*as-ndhF*; highlighted in Figure 1).

### The chloroplast genome contains at least 215 TSS

Because the 5' ends of most primary transcripts are marked by a 5' triphosphate, they should be highly over-represented in +TAP-treated versus –TAP libraries. We therefore defined a TSS as any 5' end that was at least 10 times more abundant in +TAP libraries, a calculation that generated 352 positions (Figure 3A). Because of the generally precise mechanism of transcription initiation, a true TSS should be a single predominant nucleotide rather than a cluster. We therefore filtered out putative TSS where the first nucleotide (the TSS itself) represented <50% of the coverage over a 5-nt stretch where it was the first nucleotide, as was done for barley (20). This reduced the number of TSS being considered to 215, with the others potentially representing stochastic initiation (Supplementary Table S3). Among the 215 defined TSS, 81% of the initiating nucleotides were purines (119 A's and 55 G's; Figure 3A). This trend is consistent with barley, where purines defined >80% of TSS (20). Our data identified 45 previously described *Arabidopsis* TSS, 16 of which are in orthologous positions in barley (20), meeting our expectations considering that more closely related species like *Arabidopsis thaliana* and *Nicotiana tabacum* differ in their promoter usage for some genes (44).

Most plastid genes belong to polycistronic units and have classically been described as being co-transcribed. A few genes, such as *psbA*, *rbcL* and *ndhF*, are considered to be monocistronic, which is also the tendency for tRNAs outside the rRNA operon (45,46). TSS could be identified at the 5' ends of all but one of the 20 main transcriptional units, with *petL-petG-psaJ-rpl33-rps18* being the exception. The processed *petL* 5' end is easily identifiable, however (see below). This suggests rapid maturation of the primary transcript, a phenomenon that is more prevalent in *Chlamydomonas reinhardtii* chloroplasts, where a recent transcriptome analysis revealed only 23 TSS, albeit using an entirely different method (47). Terminome-seq also confirmed our earlier discovery of TSS upstream of *trnC*, *trnF* and *trnN* (22). All transcriptional units also contained internal TSS, similar to barley (20), which could reflect mechanisms to allow for differential expression of genes within a larger cluster as is particularly well-documented for the *psbD-C* segment of the barley *psbK* gene cluster (48).

TAP treatment revealed three genomic areas with unexpectedly massive numbers of initiation events. As mentioned above, these lie in the first *clpP* intron and on the sense and antisense strands of *ndhF* (Figure 1 and Supplementary Table S3). TSS antisense to *ndhF* are distributed over nearly 300 nt and some of them were also characterized in barley and tobacco (20,49). The TSS internal to *clpP* intron 1 and on the sense strand of *ndhF* are distributed over ~160 and 600 nt, respectively, and to our knowledge have not been previously described. As a validation step, we confirmed that at least some of these ends could be amplified



**Figure 2.** Coverage and distribution of transcript termini (A) Comparison of genome coverage between RNA-seq and Terminome-seq. While the plastome is almost fully covered by at least one read in RNA-seq, only 12.7 and 15.8% is covered by 5' and 3' ends, respectively. Data for RNA-seq correspond to the average of two previously published WT replicates (19). Termini coverage at >10 RPM is marked and discussed in the text. (B) Terminome-seq read distribution in the WT. The results are the average of two biological replicates. Reads antisense to exons (as-exon) refers to reads mapping to the antisense strand of a known coding region.

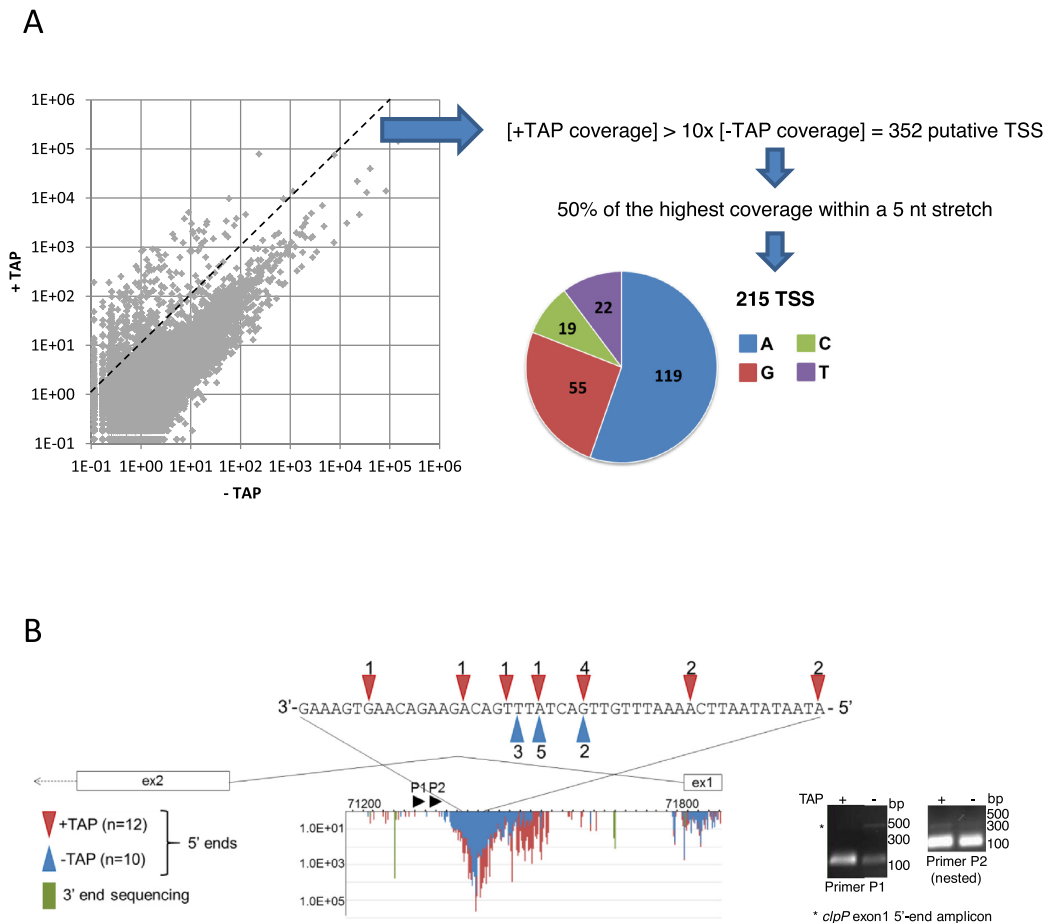
by 5' RACE following TAP treatment (Figure 3B and Supplementary Figure S3A).

Although we attempted to obtain an exhaustive catalog of TSS, the data filters we used to reduce the false discovery rate eliminated a few known initiation sites (Supplementary Figure S3B). For instance, the *PpsbN* -32 promoter (50) and *PatpE* -430 (51,52) are absent from our list because they have a +TAP:-TAP ratio of <10 (9.9 for *PatpE* and 2.5 for *PpsbN*). On the other hand, our data shed light on the uncertainty of whether *PpsbD* -186 (position 32525) is a genuine promoter or processed end (53–55), with our results being strongly diagnostic of a post-transcriptional processing site. The opposite is true

for the *ndhA* -66 (position 122076) 5' end that has been described as being created post-transcriptionally in maize through the action of PPR53 (56), but is a TSS based on our data.

### The case of the *psbB* operon

The *psbB* operon is particularly well studied, containing five genes on the plus strand (*psbB-psbT-psbH-petB-petD*, the last two containing introns) and one on the minus strand (*psbN*). Altogether, over 20 accumulating transcripts generated from this operon have been characterized (5,10,57), making it an attractive testbed for the abil-



**Figure 3.** TSS analysis (A) The abundance of 5' ends at each position for both genome strands was compared between +TAP and –TAP libraries. The dashed line separates the 352 ends that have a +TAP/–TAP ratio >10 from those with a lower ratio. Putative TSS were filtered to remove any ends that did not reach 50% of the coverage of the most represented read within a 5 nt stretch. The pie chart graphs the initiating nucleotide of the remaining 215 TSS. (B) The novel TSS detected within *clpP* intron 1 were mapped using 5' RACE. 5' RACE was completed with (red) and without (blue) prior TAP treatment, and sequenced clones are represented by colored arrowheads above/below the nucleotide sequence. The stained gel of the corresponding PCR reactions is shown at right. The gene model between exons 1 and 2 is shown with Terminome-seq results. +TAP 5' ends are in red, –TAP ends in blue and 3' ends are in green. The X-axis is genomic position and the Y-axis is RPM coverage. Black arrowheads P1 and P2 represent the 3' primers used for 5' RACE.

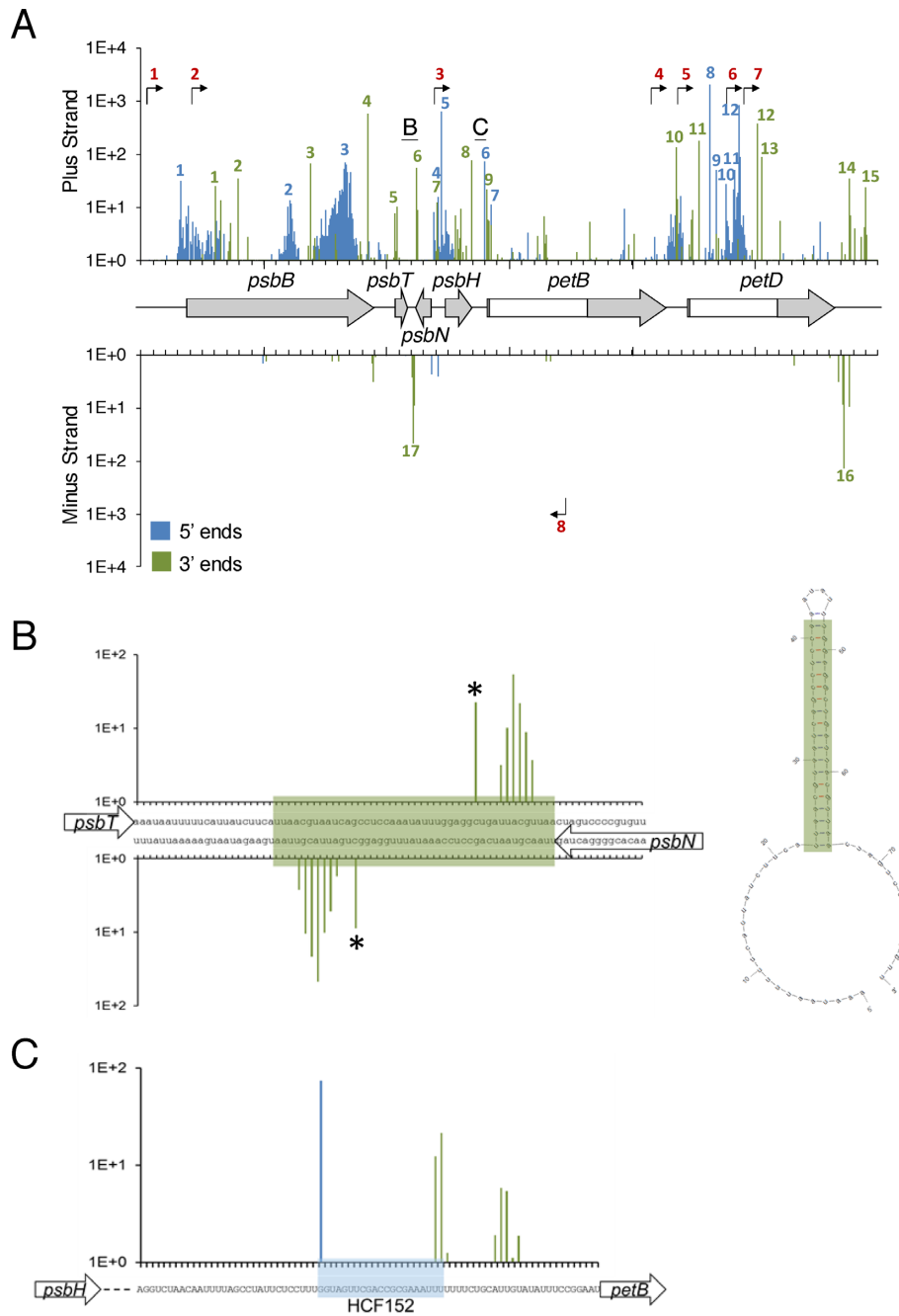
ity of Terminome-seq to reflect such a complex landscape accurately.

Figure 4A shows the positions of the eight TSS (numbered in red), along with at least 12 processed 5' ends (numbered in blue) and 17 3' ends (numbered in green), which are annotated in Table 1. A consistent feature of 5' ends is their clustered organization around a dominant peak, whereas the 3' termini are more discrete. Heterodisperse 5' ends are reminiscent of degradation intermediates, perhaps created by an enzyme such as RNase J that would progressively stall as it encountered secondary structures or bound proteins. 3' ends that were more discrete could be identified in both coding and intergenic regions, probably representing a mix of degradation intermediates and mature 3' ends. Among the processed 5' ends and 3' termini found by Terminome-seq, several had been previously described. These include the 5' ends *psbB* -51 (57), *psbH* -44 (58,59) and -67 (60) and *petB* -47 (61,62), as well as the 3' termini *psbT* +60 (30), +223 (60), *psbH* +109 (61,62), *petB* +67 (63), *petD* +94 (64) and *psbN* +39 (30). Thus, there was ex-

cellent overlap between our data and previously published work.

We conducted a similar analysis for the *atpI-atpH-atpF-atpA* and *ndhH-ndhA-ndhI-ndhG-ndhE-psaC-ndhD* gene clusters (Supplementary Figure S4 and Table S4). Five TSS were identified in the *atpI* cluster, including what might be specific promoters for *atpH* and *atpA*, and 15 TSS in the *ndhH* cluster, which were located predominantly toward the distal end. Most of these 3' and 5' termini were previously identified in a thorough investigation of the *atpI* cluster processing (52). The existence of an accumulating monocistronic *atpI* transcript in Arabidopsis has been debated (52,65) and our data suggest that the main *atpI* 3' end is 584 nt downstream of the stop codon, inside *atpH*. This 3' end is more abundant than the one at position +493, whose processing depends on PPR10 (52,66). Twenty seven termini in addition to TSS were found for the *ndhH* cluster. No detailed information was previously published on the transcript population, however it is quite complex based on gel blot analysis (67), in keeping with our results.





**Figure 4.** Transcript termini of the *psbB* operon highlighting the role of a secondary structure and RNA binding protein in defining ends (A) Terminome-seq coverage of the *psbB* operon with the corresponding gene models, with exons in gray and introns in white. –TAP 5' ends are in blue and 3' ends are in green; bent arrows represent TSS inferred from +TAP data. Underlined letters mark the locations that are expanded in panels B and C; numbered peaks and promoters refer to features listed in Table 1. (B) 3' end coverage for a stem-loop structure between *psbT* and *psbN*. The stem is highlighted in green in the nucleotide sequence and the Mfold (119) predicted secondary structure is at right. Asterisks highlight the previously described ends (30). (C) The gene model, nucleotide sequence and end coverage for the HCF152 binding site. Reads accumulate at both the 5' and 3' ends of the binding site on the plus strand, indicative of a protected RNA fragment. The color code is the same as in panel A.

### The roles of RNA secondary structures and RNA-binding proteins in shaping the terminome

Chloroplast RNA termini are known to be stabilized by stem-loop structures, as well as by sequence-specific and general RBPs (6,7,68). Both of these mechanisms act on transcripts from the *psbB* operon, for which termini were described above. While the strategy used here to make the

terminome libraries would be biased against smRNAs, the processing sites from longer, precursor RNAs could be detected. The *psbT-psbN* intergenic region, for example, forms a stem-loop that defines the 3' ends of transcripts encoded on opposite strands (30). Terminome-seq confirmed the previously identified 3' ends inside the stem and additionally showed staggered 3' ends closer to the base of the stem,

**Table 1.** Description of transcript ends originating from the *psbB* operon

End number	Genome position	Notes
TSS, +TAP 5' ends		
1	72 200	<i>PpsbB</i> -171, described in (57). Seen in barley
2	72 409	Internal to <i>psbB</i> . Distal promoter for <i>psbT</i> ?
3	74 393	<i>PpsbH</i> -92, described in (60)
4	76 153	Internal to <i>petB</i> exon 2. Distal promoter for <i>petD</i> ?
5	76 375–76 376	Upstream <i>petD</i>
6	76 391	Internal to <i>petD</i> intron
7	76 780	Internal to <i>petD</i> intron
8	75 482	Antisense to <i>petB</i> intron. Distal promoter for <i>psbN</i> ?
Processed 5' ends, –TAP 5' ends		
1	72 320	<i>psbB</i> -51 mature end, described in (57)
2	73 211	Internal to <i>psbB</i> . Highest peak from a region with multiple 5' ends. Degradation intermediate?
3	73 658	Internal to <i>psbB</i> . Highest peak from a region with multiple 5' ends. Degradation intermediate?
4	74 418	<i>psbT-psbH</i> intergenic region, <i>psbH</i> -67 mature end, described as a precise endoribonuclease cleavage in (60). See also 3' end #7
5	74 441	<i>psbT-psbH</i> intergenic region, <i>psbH</i> -44 mature end. Main <i>psbH</i> 5' end, processing depends on HCF107, described in (58,59)
6	74 794	<i>psbH-petB</i> intergenic region, <i>petB</i> -47 mature end. Processing depends on HCF152, described in (61,62). See 3' end #9
7	74 847	first nucleotide of <i>petB</i> intron. Sign of a hydrolytic splicing?
8	76 627	internal to <i>petD</i> intron. Degradation intermediates?
9	76 679	internal to <i>petD</i> intron. Degradation intermediates?
10	76 760	internal to <i>petD</i> intron. Degradation intermediates?
11	76 830	internal to <i>petD</i> intron. Degradation intermediates?
12	76 863	internal to <i>petD</i> intron. Highest peak from a region with multiple 5' ends. Degradation intermediates?
3' ends		
1	72 601	Internal to <i>psbB</i> . Degradation intermediate?
2	72 786	Internal to <i>psbB</i> . Degradation intermediate?
3	73 371	Internal to <i>psbB</i> . Degradation intermediate? Downstream of numerous 5' ends, see 5' end #2, maybe the signature of an endo-ribonuclease cleavage?
4	73 838	Internal to <i>psbB</i> . Degradation intermediate? Downstream of numerous 5' ends, see 5' end #3, maybe the signature of an endo-ribonuclease cleavage?
5	74 082	First nucleotide of <i>psbT</i> , might be the 3' end of a cDNA identified in (30)
6	74 242	3' end of <i>psbT</i> , <i>psbT</i> +60, defined by a stem loop that also defines the 3' end of the antisense <i>psbN</i> transcript, see 3' end #17. Described in (30)
7	74 405	<i>psbT-psbH</i> intergenic region, 3' end of <i>psbT</i> , <i>psbT</i> +223, described as a precise endo-ribonuclease cleavage in (60). See also 5' end #4
8	74 687	Internal to <i>psbH</i> , 20 nt upstream of the stop codon. Degradation intermediate?
9	74 814	<i>psbH-petB</i> intergenic region, <i>psbH</i> +109 mature end. Processing depends on HCF152, described in (61,62). See 5' end #6
10	76 358	<i>petB-petD</i> intergenic region, <i>petB</i> +67 mature end. Processing depends on CRP1, described in (63)
11	76 543	internal to <i>petD</i> intron. Degradation intermediates?
12	77 014	internal to <i>petD</i> intron. Degradation intermediates? 5 nt downstream of a smRNA footprint.
13	77 047	internal to <i>petD</i> intron. Degradation intermediates?
14	77 765	3' end of <i>petD</i> , <i>petD</i> +94, defined by a stem loop that also defines the 3' end of the antisense <i>rpoA</i> transcript, see 3' end #16. Processing requires mTERF6, described in (64)
15	77 892	3' end of <i>petD</i> , <i>petD</i> +221.
16	77 716	3' end of <i>rpoA</i> , <i>rpoA</i> +185, defined by a stem loop that also defines the 3' end of the antisense <i>petD</i> transcript, see 3' end #14. Processing requires mTERF6, described by (64)
17	74 211	3' end of <i>psbN</i> , <i>psbN</i> +39, defined by a stem loop that also defines the 3' end of the antisense <i>psbT</i> transcript, see 3' end #6. Described in (30)

an expected phenomenon given the tendency of exoribonucleases to stall at such positions (Figure 4B). Another stem-loop, originally described as a 'twin terminator' in spinach (69), defines the 3' ends of the *petD* and *rpoA* transcripts encoded on opposite strands, and a similar structure exists for *petA* and *psbJ* (Supplementary Figure S5). The co-existence of two sets of 3' ends associated with the *psbT* - *psbN* stem-loop is also true for the transcripts of *rbcL* (ends in positions 56 485 and 56 488; see below) and *psbA* (ends in positions 293 and 285; Supplementary Figure S5), probably reflecting breathing in AU-rich regions of the secondary structures.

The *psbH-petB* intergenic region is known to be bound by HCF152 (61,62), a PPR protein that defines the *psbH* 3' and *petB* 5' ends (70). In agreement with these results, Terminome-seq data analysis identified a single major 5' end for *petB* correlating with the 5' end of a smRNA established as the footprint protected by HCF152 binding (9,66,71). The Terminome-seq *psbH* 3' end also correlates with HCF152 binding, with an additional cluster of 3' ends found ~10 nt downstream, possibly reflecting the different stalling characteristics of the exoribonucleases PNPase and RNase II (Figure 4C).

The correlation between RBPs, their footprints found in smRNAs and Terminome-seq ends can be extended beyond HCF152 to at least 13 RBPs that have been described as involved in RNA maturation (Supplementary Figure S6). These include PPR10, which protects the two adjacent transcripts *atpI* and *atpH* and assists in the processing of both their 5' and 3' ends (66,70,72); and HCF107 and CRP1 which, like HCF152, target the *psbB* operon, in particular the *psbH* -44 5' end and the *petB* +67 3' end. We can further generalize this phenomenon, since, on a plastome-wide level, termini are enriched in areas containing smRNAs, which are in many cases likely to be marks of RBPs that still await identification (Supplementary Table S5).

### PNPase deficiency has broad impacts on both 5' and 3' termini

The roles of PNPase in chloroplast transcript 3' maturation, intron degradation, and tRNA processing, as well as an ancillary role in phosphorus metabolism, have been well documented (22,33,34,73–75). Although RNA-seq shows that most chloroplast RNAs contain 3' extensions in the *pnp1-1* null mutant (22), the individual termini were not systematically compared between WT and mutant. Such a comparison could reveal more precisely the types of sequences and structures that impede PNPase activity, and highlight its overall impact on the plastid terminome. We proceeded to analyze duplicate *pnp1-1* Terminome-seq libraries, and decided against performing TAP treatment to capture primary 5' ends because PNPase is not known or expected to affect transcription initiation.

A plastome-wide view of *pnp1-1* termini is presented in Figure 5A. While the WT and *pnp1-1* patterns show substantial overlap, there are numerous areas where reads are specific to *pnp1-1*. This is quantified as genome coverage of reads (Figure 5B) where, regardless of the RPM threshold applied, coverage is higher in *pnp1-1* than in WT. In aggregate, 5' end coverage increases from 12.7 to 26.1%, and 3' coverage from 26.1 to 39.5%. When a threshold of >10 RPM is applied to remove low abundance termini, coverage is ~1% in the mutant, with 2420 5' termini and 2744 3' termini exceeding the threshold. This, total of 5164 >10 RPM termini in *pnp1-1* represents an increase of 76% over the combined 2927 termini found in the WT. The locations of termini are also dramatically altered between genotypes (compare Figures 2B and 5C). Another noticeable difference is the decline in the proportion of termini found in rRNA, which can be rationalized as a relative increase in non-rRNA termini in the mutant. In the -TAP 5' termini population, the mutant differs by having an increase in intronic ends and a decrease in tRNA ends, the latter of which holds true for 3' termini as well.

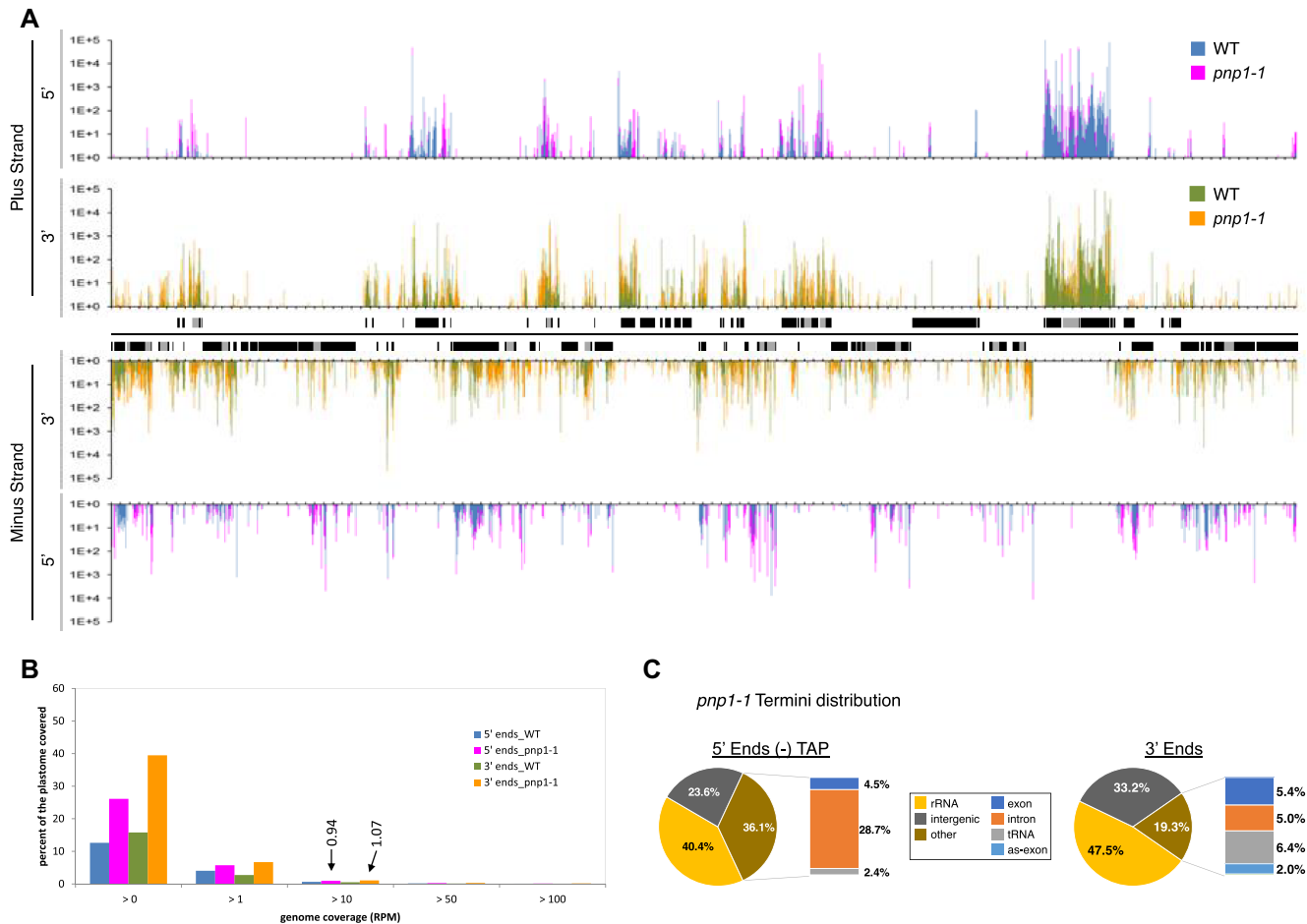
A more pronounced effect on 3' versus 5' ends in *pnp1-1* can be gleaned from Figure 6A, which shows that 349 5' termini and 1348 3' termini are at least 10 times more abundant in the mutant compared to the WT, whereas only 96 and 150, respectively, decrease at least 10-fold in the mutant. The tendency of termini to accumulate in the mutant is in keeping with the degradative function of PNPase, with the expected bias toward 3' termini. The strong effect of PNPase on the 5' terminome was unexpected and

somewhat counterintuitive, given that it is a 3'→5' exonuclease. We have previously shown, however, that PNPase degrades tRNA 5' leader sequences following their liberation by RNase P cleavage (22), and Terminome-seq-based evidence of this phenomenon is shown in Supplementary Figure S7. We speculate that other 5' termini that hyperaccumulate in *pnp1-1* represent the upstream termini of other endonucleolytic cleavage products that are usually removed by the 3'→5' activity of PNPase.

To illustrate the utility of Terminome-seq for characterizing a ribonuclease mutant at the individual gene level, results are shown for the monocistronic *rbcl* transcript (Figure 6B and C), which in plants accumulates as two species with primary and processed 5' ends. The processed 5' end is protected by the PPR protein MRL1, which prevents degradation by RNase J (76,77). MRL1 leaves a smRNA footprint (Figure 6B, blue shading), whose 5' end is represented as similarly abundant termini in WT and *pnp1-1* at genome position 54 889. There is, however, a cluster of 3' ends located ~40 nt downstream of the MRL1 binding site (with a peak at position 54 956) that is present only in the mutant, even though there is a less abundant 3' end slightly downstream (position 54 989) in both genetic backgrounds. This indicates that PNPase likely degrades *rbcl* mRNA until it is stalled by MRL1 and/or nearby RNA secondary structures, as such termini are absent in the WT. The 3' end at position 54 901 may be an RNase II stall site, given the known cooperation of RNase II and PNPase in 3'→5' RNA decay (78). A comparison of WT and *pnp1-1* ends in the vicinity of all described RBP sites is provided in Supplementary Figure S6.

At its 3' end, the *rbcl* transcript is defined by a highly conserved stem-loop (79,80) represented as a smRNA (66,71). RNA-seq showed higher coverage in *pnp1-1* compared to the WT beginning about 10 nt downstream of this structure (22). Using Terminome-seq results, we identified three different 3' ends in WT plants, two of them directly at the 3' base of the stem (positions 56 485 and 56 488) and one 36 nt downstream (position 56 524; Figure 6C). The 3' end at position 56 485 could also be identified in *pnp1-1*, although it was less abundant than in WT, suggesting that its production is not fully dependent on PNPase. On the contrary, ends at positions 56 488 and 56 524 are certainly produced through the action of PNPase because they are missing in the mutant. The most abundant 3' ends in the mutant cluster around position 56 608, 120 nt downstream of the stem-loop (Figure 6C), accounting for the 3' extension previously noted in RNA gel blots, which represents the stall point of RNase II (78). Terminome-seq evidence for several other putative 3' extensions is presented in Supplementary Figure S8.

At last, we present an overview of rRNA operon ends for the WT and *pnp1-1* (Supplementary Figure S9). This analysis revealed termini corresponding to known processing sites (6) for all four rRNAs, and also showed clear peaks in *pnp1-1* for the well-known 23S rRNA extension in this mutant (73). Of note, the first hidden break of the 23S rRNA was not clearly delineated compared to the second 23S rRNA hidden break, in keeping with previously mapped ends in this region (19). Surprisingly, in the 5' part of 16S rRNA there were numerous, staggered 5' termini in-



**Figure 5.** Distribution, coverage and location of transcript termini in *pnp1-1* (A) Plastome-scale view of end coverages from the average of two Col-0 and *pnp1-1* biological replicates in RPM. 5' ends obtained without TAP treatment are blue (WT) and pink (*pnp1-1*), and 3' ends are displayed in green (WT) and orange (*pnp1-1*). Gene models are indicated between the tracks corresponding to the plus and minus strands of the plastome. One copy of the large inverted repeat is omitted for clarity. Tick marks are every 1000 nt. (B) Comparison of Terminome-seq coverage for WT and *pnp1-1*. While 12.7 and 15.8% of the WT plastome is represented by 5' and 3' ends, respectively, these numbers increase to 26.1 and 39.5%, respectively, in *pnp1-1*. Termini coverage at >10 RPM (0.94 and 1.07% for 5' and 3' ends, respectively, in *pnp1-1*) is marked and discussed in the text. (C) Terminome-seq read distribution in *pnp1-1*. The results are the average of two biological replicates. as-exon refers to reads mapping to the antisense strand of known coding regions.

dicative of a processing or degradation event, albeit at 1–2 logs lower abundance than the mature 5' and 3' ends. RNA processing at this position, however, has not been previously described and further analysis of this region is needed to evaluate the origins of these termini.

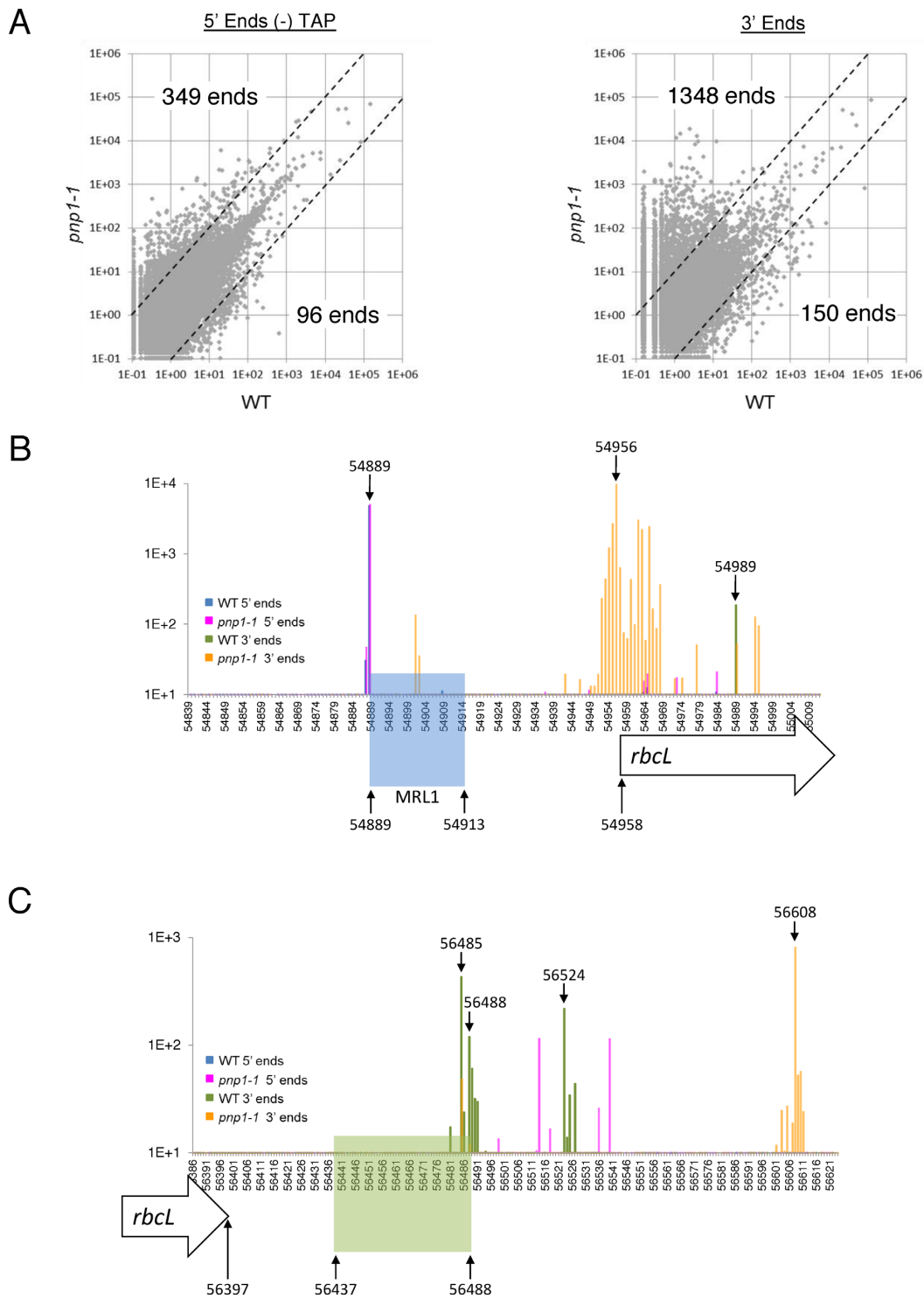
## DISCUSSION

### The chloroplast TSS landscape

In this work, we systematically sequenced *Arabidopsis thaliana* chloroplast RNA 3' termini as well as primary and processed 5' ends, using TAP treatment to discriminate between them. Such strategies are under continual improvement, as transcriptome analyses generally endeavor to be as comprehensive and quantitative as possible. In our case, transcript chemistry, size and secondary structure will all introduce biases to the dataset ultimately used to draw the main conclusions. In addition, factors such as ligation bias—the preferential ligation of adapters to certain sequences (81,82), undoubtedly impacted the selection and

ratios of the transcript ends we have reported. Therefore, our conclusions—particularly quantitative ones—should be evaluated with those caveats in mind.

The strategy as implemented led to the description of 215 TSS meeting defined expression thresholds, which are widely distributed in the plastome (Figure 1) and is consistent in number with the 176 TSS mapped in mature barley leaves (20). These TSS are created by two RNA polymerase types, a bacterial-like, plastid-encoded RNA polymerase (PEP) and two phage-like, nucleus-encoded, RNA polymerases (NEP). NEP and PEP operate simultaneously, however our samples were taken from tissue populated by mature chloroplasts, where PEP activity predominates (83). Since NEP transcription primarily occurs during early differentiation, the TSS we describe likely underestimate the number of promoters utilized over the course of chloroplast differentiation. In barley, the use of the *albostrians* mutant with sectors lacking PEP allowed the discovery of 254 additional NEP-dependent TSS, giving a total of 398 unique TSS when overlap between PEP and NEP TSS was consid-



**Figure 6.** Terminome-seq coverage in *pnp1-1* (A) The RPM abundance of –TAP 5' ends and 3' ends was compared between WT and *pnp1-1* at a genomic level. The dashed lines separate ends that are at least 10-fold more abundant in a given genotype. For example, 349 5' ends are more abundant in the PNase mutant. (B) Terminome-seq coverage upstream of the *rbcL* gene. Color coding of ends is provided in an inset. Genome position 54 958 is the *rbcL* coding region 5' end according to the TAIR10 annotation. The *rbcL* processed 5' end (position 54 889) correlates with the 5' end of the smRNA footprint of MRL1 (highlighted in blue). (C) Terminome-seq coverage downstream of the *rbcL* gene, with labeling as in Panel B. Genome position 56 397 is the 3' end of the coding region. The stem-loop downstream of the gene (positions 56 437–56 488) is highlighted in green and matches a smRNA (66). Other genome positions discussed in the text are also labeled.

ered. At the same time, NEP is known to become promiscuous when PEP has been eliminated genetically (1,84). The proportion of NEP promoters active under these circumstances that are also utilized in WT plants remains to be determined.

If we assume a total of ~400 TSS in Arabidopsis when taking PEP and 'developmentally hidden' NEP promoters into account, this would average to a TSS every ~600 nt when both strands, but only one of the large inverted repeats, are used as a basis. This frequency is not surprising, and such phenomenon occurs in bacteria as well (85,86). Given the AT-rich content of the plastid genome, functional PEP promoter -10 elements are likely to be present by chance, along with the short and highly variable elements that seem to constitute NEP promoters. The average 600 nt spacing we calculate would be reduced if we changed the 10-fold ratio threshold used to define TAP versus non-TAP-dependent transcripts. Some known TSS such as *PatpE* -430, *PpsbN* -32 and *PpsbD* -948 (87), have a  $\pm$ TAP ratio <10 and were therefore not formally considered to be TSS (Supplementary Figure S3B). Perhaps the biggest surprise from this Terminome-seq dataset is the evidence for massive transcription initiation activity in the 3' part of the *ndhF* gene, antisense to *ndhF* and within the first intron of *clpP*. The former two areas represent a staggering ~20% of the reads derived from primary transcripts, and could have functional implications in the chloroplast.

### Do some TSS mark transcripts with novel coding potential?

Chloroplasts are known to contain plastid non-coding RNAs (pncRNAs), with >100 in Arabidopsis (88), whose functions remain mostly unclear. Although many pncRNAs appear to be generated by read-through from adjacent genes, some are known to be primary transcripts (20,23), and the high number of TSS identified here potentially increases the count of pncRNA promoters. For example, the internal *ndhF* TSS (Supplementary Figure S3A) probably initiates a 300–600 nt RNA that ends at the *ndhF* termini 110 331 and 109 924, the latter position being protected by an RBP or stem-loop since a smRNA from this position was identified (66). This pncRNA overlaps *ycf1* on the antisense strand and was previously thought to originate from *ndhF* readthrough (23). Initiation antisense to *ndhF* was also described in tobacco (89) and more recently in barley (20), and is likely responsible for additional Arabidopsis pncRNAs (nc89 and nc90) which were validated by RT-PCR and gel blots (23). The tobacco *ndhF* antisense TSS was proposed to be a proximal TSS for the downstream *rpl32* gene (89), which might also be the case in Arabidopsis.

Many pncRNAs are antisense to coding sequences, and for a few there is evidence they exert a regulatory function at the RNA level (19,29,60). On the other hand, many pncRNAs contain small open reading frames (ORFs) and could therefore encode unknown chloroplast proteins. Such ORFs were discarded from the chloroplast genome sequences first obtained from tobacco and *Marchantia* (90,91), specifically ORFs shorter than 70 nt unless the products were known. In contrast, the functions of the

longer, conserved, hypothetical coding frames (*ycfs*) were still discussed. Whether pncRNA-encoded ORFs are represented in the proteome remains to be determined, however data from bacterial and animal systems suggests that at least some antisense or intergenic RNAs harbor hidden genetic functions (92–94). In this case, the retention of TSS would be expected.

Among the better-studied YCFs, our data shows that *ycf15* (also annotated as ORF77) contains its own TSS (position 93 369, Supplementary Table S3) and shares an abundant 3' end with the upstream *ycf2* transcript at position 93 750, probably defined by an RBP (Supplementary Table S5). *ycf15* is a short ORF downstream of *ycf2* whose functionality as a protein-coding gene has been debated (95–98); however the presence of discrete 5' and 3' ends would be consistent with functionality.

### Terminome-seq corroborates the production of smRNA footprints from post-transcriptional processing

Despite their high number, TSS only account for a small fraction of 5' end diversity. Instead, most 5' and 3' ends are the result of a maturation and winnowing process where RBPs and secondary structures selectively protect RNAs from degradation by RNases with low sequence specificity (6,8). Many of these RBPs are members of plant-specific or plant-amplified helical repeat protein families (99), most prominently the PPR family (100). The correlation between smRNA footprints and transcript termini is a key argument which posits that target sequences and secondary structures largely exist to define the ends of functional transcripts (66,71,101). It follows that the presence of termini correlating with an smRNA represents an effective way to distinguish true footprints from other smRNAs that are more scattered or whose termini are much more ragged.

A good example of using Terminome-seq to distinguish smRNA footprints is the *psbB* transcriptional unit (Figure 4), which gives rise to eight accumulating smRNAs (66). These include footprints for three RBPs: HCF107, HCF152 and CRP1; two derived from stem-loops downstream of *psbT/psbN* and *petD*, and three smRNAs complementary to the *petD* intron. All but the three intron antisense smRNAs, which may be too unstable to be found in the longer transcripts we sequenced, correlate with termini (Table 1). Another example is the recently described PPR10-mediated protection of the maize *psaI* 3' end (72), which correlates with a smRNA from an analogous position in Arabidopsis (genomic position 59 475, see Supplementary Figure S6 and Table S5). Across the transcriptome, we were able to identify (using the 10 RPM threshold) 45 5' ends and 44 3' ends which correlate with the termini of 81 smRNAs (8 footprints correlate with both 5' and 3' ends, 37 with 5' ends only and 36 with 3' ends only), accounting for 33.5% of the 242 smRNAs previously described (66,71). Because of the way in which our libraries were made, these ends belong to the longer transcripts from which the smRNAs were ultimately derived. The ability to look for *bona fide* footprints, along with the elucidation of the PPR code explaining their binding specificity (102–105), will allow the prediction and

discovery of new PPR proteins involved in chloroplast RNA stability. Not all of the hundreds of chloroplast PPRs generate accumulating footprints, as demonstrated by PGR3, which is required for *rpl14* 3' end formation (72).

According to the footprint model, a single RBP can protect the 5' and 3' ends of overlapping transcripts. The original example of this phenomenon is maize PPR10, which in binding the *atpI-atpH* intergenic region protects the *atpI* 3' and *atpH* 5' ends (70,106). Although Terminome-seq could identify the expected *atpH* 5' end, only a minor 3' end mapped to the PPR10-dependent *atpI* 3' end (Supplementary Figure S6). The major Terminome-seq 3' end is located further downstream, inside the *atpH* coding region and does not correlate with a smRNA. Systematic mapping of the *atp* operon termini already revealed that transcripts with overlapping ends are a minority (52) and co-immunoprecipitation showed that PPR10 is preferentially associated with processed *atpH* mRNA rather than processed *atpI* mRNA (66). At a genome-wide level, only ~10 footprints (including PPR38 and HCF152) can be linked to simultaneous stabilization of 5' and 3' ends, suggesting that this is the exception rather than the rule. Secondary structures, on the other hand, can stabilize the 3' ends of two adjacent transcripts, for example in the cases of *psbT-psbN*, *petD-rpoA* and *petA-psbJ* (Figure 4B and Supplementary Figure S5). Currently, the degradation pathway of RBP footprints is unknown. Evidence derived from Terminome-seq is consistent with PNPase participating in this mechanism, as there are more termini correlating with smRNAs in the PNPase mutant than in WT (42.5% versus 33.5% of all termini; Supplementary Table S5), however a direct analysis of smRNAs in both genotypes would be required to substantiate this possibility.

### Terminome-seq of the PNPase mutant points to its potential use in reverse genetics

mRNA processing up to the border of the region protected by RBPs or secondary structures is principally performed by three enzymes, RNase J for 5' ends and PNPase and RNase II for 3' ends (76,78,106). As expected, Terminome-seq of the PNPase mutant showed a more pronounced effect on 3' versus 5' termini, although many termini of both types were affected (Figures 5 and Figure 6A). We could confirm that PNPase is required to remove tRNA precursor 5' extensions subsequent to RNase P processing (Supplementary Figures S7,(22)). Similarly, we confirmed that the 3' extensions observed in the PNPase mutant (22,34,73) correlate with 3' ends distal to secondary structures or RBP binding sites (Supplementary Figures S5, 6 and 8), two elements known to stall PNPase (106,107).

Using *pnp1-1* as a proof of concept opens the door to understanding the roles of other chloroplast gene regulators using Terminome-seq. Concerning 5' termini, this might include the roles of the sigma factors that specify PEP initiation sites (108) and other PEP-associated proteins (PAPs; 109), a closer view at RNase J specificity, or a better understanding of some of the numerous factors that have been implicated in rRNA maturation (6). Terminome-seq would be an excellent choice to gain additional insight into poorly-

characterized RNases whose targets have only begun to be described (110–112).

### Transcription termination analysis with Terminome-seq

Transcription termination in plastids has recently been reviewed (113). The endosymbiont hypothesis would predict that chloroplast termination would resemble that of bacteria, which use both Rho-dependent and Rho-independent mechanisms. Rho-independent termination occurs in AT-rich sequences downstream of GC-rich stem-loops, however *in vitro* and *in vivo* assays found that PEP terminates inefficiently at chloroplast stem-loops, although it does recognize certain bacterial terminators. This is in keeping with the longstanding hypothesis that chloroplast stem-loops are involved in transcript stability rather than transcription termination (3), a finding well illustrated by the correlation between 3' ends and secondary structures (Figure 4B and Supplementary Figure S5). On the other hand, Rho-dependent termination, which works through destabilization of the elongating polymerase (114), would have to rely on an alternative cofactor since Rho is not found in organelles. Such a factor was recently identified by its partial similarity to Rho, RHON1, which assists in transcription termination downstream of *rbcL* (115), binding the transcript between two termini identified here (Figure 6C). The 3' end in the PNPase mutant is further downstream, suggesting a model similar to *Escherichia coli* where Rho factors initiate termination and exoribonucleases like PNPase produce the precise, mature 3' ends directly downstream of stem-loops or RBP sites (116). Similarly, MTERF6, a member of a gene family related to a human mitochondrial transcription termination factor, acts in the *trnI* and *petD-rpoA* regions (64,113,117). As additional factors involved in termination are discovered, Terminome-seq offers a tool to decipher their sites of action.

### CONCLUSION

Several high-throughput RNA-seq-based strategies have recently been used to gain an unprecedented genome-level understanding of chloroplast RNA biology. It is now possible to study RNA abundance, splicing and editing (26–28), to monitor translation rates with ribosome profiling (24,118) and to infer protein binding sites from smRNA sequencing (66,71,101). Terminome-seq now creates access to a single-nucleotide resolution map of the full set of chloroplast RNA termini. We anticipate that this strategy, combined with the other plastome-wide approaches, will be instrumental in deciphering mechanisms such as transcription (from initiation to termination), as well as the roles of RNases and RBPs in shaping the chloroplast transcriptome.

### DATA AVAILABILITY

Raw sequences have been deposited on the SRA database with the number PRJNA533962 and can be accessed here <https://www.ncbi.nlm.nih.gov/sra/PRJNA533962>.

### SUPPLEMENTARY DATA

Supplementary Data are available at NAR Online.

## FUNDING

Division of Chemical Sciences, Geosciences, and Biosciences, Office of Basic Energy Sciences, of the U.S. Department of Energy [DE-FG02-10ER20015, in part]; LabEx Saclay Plant Sciences-SPS [ANR-10-LABX-0040-SPS to B.C.]. Funding for open access charge: Internal Funds (to D.B.S).

*Conflict of interest statement.* None declared.

## REFERENCES

- Legen, J., Kemp, S., Krause, K., Profanter, B., Herrmann, R.G. and Maier, R.M. (2002) Comparative analysis of plastid transcription profiles of entire plastid chromosomes from tobacco attributed to wild-type and PEP-deficient transcription machineries. *Plant J.*, **31**, 171–188.
- Sanitá Lima, M. and Smith, D.R. (2017) Pervasive, genome-wide transcription in the organelle genomes of diverse plastid-bearing protists. *G3*, **7**, 3789–3796.
- Stern, D.B. and Gruissem, W. (1987) Control of plastid gene expression: 3' inverted repeats act as mRNA processing and stabilizing elements, but do not terminate transcription. *Cell*, **51**, 1145–1157.
- Mullet, J.E. and Klein, R.R. (1987) Transcription and RNA stability are important determinants of higher plant chloroplast RNA levels. *EMBO J.*, **6**, 1571–1579.
- Barkan, A. (2011) Expression of plastid genes: organelle-specific elaborations on a prokaryotic scaffold. *Plant Physiol.*, **155**, 1520–1532.
- Germain, A., Hotto, A.M., Barkan, A. and Stern, D.B. (2013) RNA processing and decay in plastids. *Wiley Interdiscip. Rev. RNA*, **4**, 295–316.
- Stern, D.B., Goldschmidt-Clermont, M. and Hanson, M.R. (2010) Chloroplast RNA metabolism. *Annu. Rev. Plant Biol.*, **61**, 125–155.
- Stoppel, R. and Meurer, J. (2012) The cutting crew - ribonucleases are key players in the control of plastid gene expression. *J. Exp. Bot.*, **63**, 1663–1673.
- Manavski, N., Schmid, L.-M. and Meurer, J. (2018) RNA-stabilization factors in chloroplasts of vascular plants. *Essays Biochem.*, **62**, 51–64.
- Stoppel, R. and Meurer, J. (2013) Complex RNA metabolism in the chloroplast: an update on the psbB operon. *Planta*, **237**, 441–449.
- Yang, J., Schuster, G. and Stern, D.B. (1996) CSP41, a sequence-specific chloroplast mRNA binding protein, is an endoribonuclease. *Plant Cell*, **8**, 1409–1420.
- Beligni, M.V. and Mayfield, S.P. (2008) Arabidopsis thaliana mutants reveal a role for CSP41a and CSP41b, two ribosome-associated endonucleases, in chloroplast ribosomal RNA metabolism. *Plant Mol. Biol.*, **67**, 389–401.
- Bollenbach, T.J., Sharwood, R.E., Gutierrez, R., Lerbs-Mache, S. and Stern, D.B. (2009) The RNA-binding proteins CSP41a and CSP41b may regulate transcription and translation of chloroplast-encoded RNAs in Arabidopsis. *Plant Mol. Biol.*, **69**, 541–552.
- Bollenbach, T.J., Tatman, D.A. and Stern, D.B. (2003) CSP41a, a multifunctional RNA-binding protein, initiates mRNA turnover in tobacco chloroplasts. *Plant J.*, **36**, 842–852.
- Qi, Y., Armbruster, U., Schmitz-Linneweber, C., Delannoy, E., de Longevialle, A.F., Rühle, T., Small, I., Jahns, P. and Leister, D. (2012) Arabidopsis CSP41 proteins form multimeric complexes that bind and stabilize distinct plastid transcripts. *J. Exp. Bot.*, **63**, 1251–1270.
- Stoppel, R., Manavski, N., Schein, A., Schuster, G., Teubner, M., Schmitz-Linneweber, C. and Meurer, J. (2012) RHON1 is a novel ribonucleic acid-binding protein that supports RNase E function in the Arabidopsis chloroplast. *Nucleic Acids Res.*, **40**, 8593–8606.
- Walter, M., Piepenburg, K., Schöttler, M.A., Petersen, K., Kahlau, S., Tiller, N., Drechsel, O., Weingartner, M., Kudla, J. and Bock, R. (2010) Knockout of the plastid RNase E leads to defective RNA processing and chloroplast ribosome deficiency. *Plant J.*, **64**, 851–863.
- Sharwood, R.E., Halpert, M., Luro, S., Schuster, G. and Stern, D.B. (2011) Chloroplast RNase J compensates for inefficient transcription termination by removal of antisense RNA. *RNA*, **17**, 2165–2176.
- Hotto, A.M., Castandet, B., Gilet, L., Higdon, A., Condon, C. and Stern, D.B. (2015) Arabidopsis chloroplast Mini-Ribonuclease III participates in rRNA maturation and intron recycling. *Plant Cell*, **27**, 724–740.
- Zhelyazkova, P., Sharma, C.M., Forstner, K.U., Liere, K., Vogel, J. and Börner, T. (2012) The primary transcriptome of barley chloroplasts: numerous noncoding RNAs and the dominating role of the plastid-encoded RNA polymerase. *Plant Cell*, **24**, 123–136.
- Ruwe, H., Castandet, B., Schmitz-Linneweber, C. and Stern, D.B. (2013) Arabidopsis chloroplast quantitative editotype. *FEBS Lett.*, **587**, 1429–1433.
- Castandet, B., Hotto, A.M., Fei, Z. and Stern, D.B. (2013) Strand-specific RNA sequencing uncovers chloroplast ribonuclease functions. *FEBS Lett.*, **587**, 3096–3101.
- Hotto, A.M., Schmitz, R.J., Fei, Z., Ecker, J.R. and Stern, D.B. (2011) Unexpected diversity of chloroplast noncoding RNAs as revealed by deep sequencing of the Arabidopsis transcriptome. *G3*, **1**, 559–570.
- Chotewutmontri, P. and Barkan, A. (2016) Dynamics of chloroplast translation during chloroplast differentiation in maize. *PLoS Genet.*, **12**, e1006106.
- Guillaumot, D., Lopez-Obando, M., Baudry, K., Avon, A., Rigai, G., Falcon de Longevialle, A., Broche, B., Takenaka, M., Berthomé, R., De Jaeger, G. et al. (2017) Two interacting PPR proteins are major Arabidopsis editing factors in plastid and mitochondria. *Proc. Natl. Acad. Sci. U.S.A.*, **114**, 8877–8882.
- Michel, E.J.S., Hotto, A.M., Strickler, S.R., Stern, D.B. and Castandet, B. (2018) A guide to the chloroplast transcriptome analysis using RNA-seq. *Methods Mol Biol.*, **1829**, 295–313.
- Castandet, B., Hotto, A.M., Strickler, S.R. and Stern, D.B. (2016) ChloroSeq, an optimized chloroplast rna-seq bioinformatic pipeline, reveals remodeling of the organellar transcriptome under heat stress. *G3*, **6**, 2817–2827.
- Malbert, B., Rigai, G., Brunaud, V., Lurin, C. and Delannoy, E. (2018) Bioinformatic analysis of chloroplast gene expression and RNA posttranscriptional maturations using RNA sequencing. *Methods Mol Biol.*, **1829**, 279–294.
- Hotto, A.M., Huston, Z.E. and Stern, D.B. (2010) Overexpression of a natural chloroplast-encoded antisense RNA in tobacco destabilizes 5S rRNA and retards plant growth. *BMC Plant Biol.*, **10**, 213.
- Zghidi-Abouzid, O., Merendino, L., Buhr, F., Malik Ghulam, M. and Lerbs-Mache, S. (2011) Characterization of plastid psbT sense and antisense RNAs. *Nucleic Acids Res.*, **39**, 5379–5387.
- Hübschmann, T. and Börner, T. (1998) Characterisation of transcript initiation sites in ribosome-deficient barley plastids. *Plant Mol. Biol.*, **36**, 493–496.
- Hess, W.R., Prombona, A., Fieder, B., Subramanian, A.R. and Börner, T. (1993) Chloroplast rps15 and the rpoB/C1/C2 gene cluster are strongly transcribed in ribosome-deficient plastids: evidence for a functioning non-chloroplast-encoded RNA polymerase. *EMBO J.*, **12**, 563–571.
- Marche, C., Yehudai-Resheff, S., Germain, A., Fei, Z., Jiang, X., Judkins, J., Wu, H., Fernie, A.R., Fait, A. and Stern, D.B. (2009) Abnormal physiological and molecular mutant phenotypes link chloroplast polynucleotide phosphorylase to the phosphorus deprivation response in Arabidopsis. *Plant Physiol.*, **151**, 905–924.
- Germain, A., Herlich, S., Larom, S., Kim, S.H., Schuster, G. and Stern, D.B. (2011) Mutational analysis of Arabidopsis chloroplast polynucleotide phosphorylase reveals roles for both RNase PH core domains in polyadenylation, RNA 3'-end maturation and intron degradation. *Plant J.*, **67**, 381–394.
- Steglich, C., Futschik, M.E., Lindell, D., Voss, B., Chisholm, S.W. and Hess, W.R. (2008) The challenge of regulation in a minimal photoautotroph: non-coding RNAs in Prochlorococcus. *PLoS Genet.*, **4**, e1000173.
- Olivarius, S., Plessy, C. and Carninci, P. (2009) High-throughput verification of transcriptional starting sites by Deep-RACE. *Biotechniques*, **46**, 130–132.
- Denise, H., Moschos, S.A., Sidders, B., Burden, F., Perkins, H., Carter, N., Stroud, T., Kennedy, M., Fancy, S.-A., Laphorn, C. et al. (2014) Deep Sequencing Insights in therapeutic shRNA processing and siRNA target cleavage precision. *Mol. Ther. Nucleic Acids*, **3**, e145.
- Hoff, A.M., Johannessen, B., Alagaratnam, S., Zhao, S., Nome, T., Lövf, M., Bakken, A.C., Hektoen, M., Sveen, A., Lothe, R.A. et al.



- (2015) Novel RNA variants in colorectal cancers. *Oncotarget*, **6**, 36587–36602.
39. Lagarde, J., Uszczynska-Ratajczak, B., Santoyo-Lopez, J., Gonzalez, J.M., Tapanari, E., Mudge, J.M., Steward, C.A., Wilming, L., Tanzer, A., Howald, C. *et al.* (2016) Extension of human lncRNA transcripts by RACE coupled with long-read high-throughput sequencing (RACE-Seq). *Nat. Commun.*, **7**, 12339.
  40. Park, D., Morris, A.R., Battenhouse, A. and Iyer, V.R. (2014) Simultaneous mapping of transcript ends at single-nucleotide resolution and identification of widespread promoter-associated non-coding RNA governed by TATA elements. *Nucleic Acids Res.*, **42**, 3736–3749.
  41. Rorbach, J., Bobrowicz, A., Pearce, S. and Minczuk, M. (2014) Polyadenylation in bacteria and organelles. In: Rorbach, J and Bobrowicz, J.A (eds). *Polyadenylation: Methods and Protocols*. Humana Press, Totowa, NJ, pp. 211–227.
  42. Slomovic, S. and Schuster, G. (2011) Exonucleases and endonucleases involved in polyadenylation-assisted RNA decay. *Wiley Interdiscip. Rev. RNA*, **2**, 106–123.
  43. Sato, S., Nakamura, Y., Kaneko, T., Asamizu, E. and Tabata, S. (1999) Complete structure of the chloroplast genome of *Arabidopsis thaliana*. *DNA Res.*, **6**, 283–290.
  44. Swiatecka-Hagenbruch, M., Liere, K. and Börner, T. (2007) High diversity of plastidial promoters in *Arabidopsis thaliana*. *Mol. Genet. Genomics*, **277**, 725–734.
  45. Sugita, M. and Sugiura, M. (1996) Regulation of gene expression in chloroplasts of higher plants. *Plant Mol. Biol.*, **32**, 315–326.
  46. Shahar, N., Weiner, I., Stotsky, L., Tuller, T. and Yacoby, I. (2019) Prediction and large-scale analysis of primary operons in plastids reveals unique genetic features in the evolution of chloroplasts. *Nucleic Acids Res.*, **47**, 3344–3352.
  47. Cavaiuolo, M., Kuras, R., Wollman, F., Choquet, Y. and Vallon, O. (2017) Small RNA profiling in *Chlamydomonas*: insights into chloroplast RNA metabolism. *Nucleic Acids Res.*, **45**, 10783–10799.
  48. Christopher, D.A., Kim, M. and Mullet, J.E. (1992) A novel light-regulated promoter is conserved in cereal and dicot chloroplasts. *Plant Cell*, **4**, 785–798.
  49. Vera, A., Hirose, T. and Sugiura, M. (1996) A ribosomal protein gene (rpl32) from tobacco chloroplast DNA is transcribed from alternative promoters: Similarities in promoter region organization in plastid housekeeping genes. *Mol. Gen. Genet.*, **251**, 518–525.
  50. Zghidi, W., Merendino, L., Cottet, A., Mache, R. and Lerbs-Mache, S. (2007) Nucleus-encoded plastid sigma factor SIG3 transcribes specifically the psbN gene in plastids. *Nucleic Acids Res.*, **35**, 455–464.
  51. Kapoor, S., Wakasugi, T., Deno, H. and Sugiura, M. (1994) An atpE-specific promoter within the coding region of the atpB gene in tobacco chloroplast DNA. *Curr. Genet.*, **26**, 263–268.
  52. Ghulam, M.M., Courtois, F., Lerbs-Mache, S., Merendino, L. and Merendino, L. (2013) Complex processing patterns of mRNAs of the large ATP synthase operon in *Arabidopsis* chloroplasts. *PLoS One*, **8**, e78265.
  53. Hanaoka, M., Kanamaru, K., Takahashi, H. and Tanaka, K. (2003) Molecular genetic analysis of chloroplast gene promoters dependent on SIG2, a nucleus-encoded sigma factor for the plastid-encoded RNA polymerase, in *Arabidopsis thaliana*. *Nucleic Acids Res.*, **31**, 7090–7098.
  54. Hoffer, P.H. and Christopher, D.A. (1997) Structure and blue-light-responsive transcription of a chloroplast psbD promoter from *Arabidopsis thaliana*. *Plant Physiol.*, **115**, 213–222.
  55. Nagashima, A., Hanaoka, M., Shikanai, T., Fujiwara, M., Kanamaru, K., Takahashi, H. and Tanaka, K. (2004) The multiple-stress responsive plastid sigma factor, SIG5, directs activation of the psbD blue light-responsive promoter (BLRP) in *Arabidopsis thaliana*. *Plant Cell Physiol.*, **45**, 357–368.
  56. Zoschke, R., Watkins, K.P., Miranda, R.G. and Barkan, A. (2016) The PPR-SMR protein PPR53 enhances the stability and translation of specific chloroplast RNAs in maize. *Plant J.*, **85**, 594–606.
  57. Westhoff, P. and Herrmann, R.G. (1988) Complex RNA maturation in chloroplasts: the psbB operon from spinach. *Eur. J. Biochem.*, **171**, 551–564.
  58. Felder, S., Meierhoff, K., Sane, A.P., Meurer, J., Driemel, C., Plucken, H., Klaff, P., Stein, B., Bechtold, N. and Westhoff, P. (2001) The nucleus-encoded HCF107 gene of *Arabidopsis* provides a link between intergenic RNA processing and the accumulation of translation-competent psbH transcripts in chloroplasts. *Plant Cell*, **13**, 2127–2141.
  59. Hammani, K., Cook, W.B. and Barkan, A. (2012) RNA binding and RNA remodeling activities of the half-a-tetratricopeptide (HAT) protein HCF107 underlie its effects on gene expression. *Proc. Natl. Acad. Sci. U.S.A.*, **109**, 5651–5656.
  60. Chevalier, F., Ghulam, M.M., Rondet, D., Pfanschmidt, T., Merendino, L. and Lerbs-Mache, S. (2015) Characterization of the psbH precursor RNAs reveals a precise endoribonuclease cleavage site in the psbT/psbH intergenic region that is dependent on psbN gene expression. *Plant Mol. Biol.*, **88**, 357–367.
  61. Nakamura, T., Meierhoff, K., Westhoff, P. and Schuster, G. (2003) RNA-binding properties of HCF152, an *Arabidopsis* PPR protein involved in the processing of chloroplast RNA. *Eur. J. Biochem.*, **270**, 4070–4081.
  62. Meierhoff, K., Felder, S., Nakamura, T., Bechtold, N. and Schuster, G. (2003) HCF152, an *Arabidopsis* RNA Binding pentatricopeptide repeat protein involved in the processing of chloroplast psbB-psbT-psbH-petB-petD RNAs. *Plant Cell*, **15**, 1480–1495.
  63. Ferrari, R., Tadini, L., Moratti, F., Lehniger, M.-K., Costa, A., Rossi, F., Colombo, M., Masiero, S., Schmitz-Linneweber, C. and Pesaresi, P. (2017) CRP1 Protein: (dis)similarities between *Arabidopsis thaliana* and *Zea mays*. *Front. Plant Sci.*, **8**, 163.
  64. Zhang, Y., Cui, Y.L., Zhang, X.L., Yu, Q.B., Wang, X., Yuan, X.B., Qin, X.M., He, X.F., Huang, C. and Yang, Z.N. (2018) A nuclear-encoded protein, mTERF6, mediates transcription termination of rpoA polycistron for plastid-encoded RNA polymerase-dependent chloroplast gene expression and chloroplast development. *Sci. Rep.*, **8**, 11929.
  65. Malik Ghulam, M., Zghidi-Abouzid, O., Lambert, E., Lerbs-Mache, S. and Merendino, L. (2012) Transcriptional organization of the large and the small ATP synthase operons, atpI/H/F/A and atpB/E, in *Arabidopsis thaliana* chloroplasts. *Plant Mol. Biol.*, **79**, 259–272.
  66. Ruwe, H., Wang, G., Gusewski, S. and Schmitz-Linneweber, C. (2016) Systematic analysis of plant mitochondrial and chloroplast small RNAs suggests organelle-specific mRNA stabilization mechanisms. *Nucleic Acids Res.*, **44**, 7406–7417.
  67. Meurer, J., Berger, A. and Westhoff, P. (1996) A nuclear mutant of *Arabidopsis* with impaired stability on distinct transcripts of the plastid psbB, psbD/C, ndhH, and ndhC operons. *Plant Cell*, **8**, 1193–1207.
  68. Shi, X., Hanson, M.R. and Bentolila, S. (2017) Functional diversity of *Arabidopsis* organelle-localized RNA-recognition motif-containing proteins. *Wiley Interdiscip. Rev. RNA*, **8**, e1420.
  69. Sijben-Müller, G., Hallick, R.B., Alt, J., Westhoff, P. and Herrmann, R.G. (1986) Spinach plastid genes coding for initiation factor IF-1, ribosomal protein S11 and RNA polymerase alpha-subunit. *Nucleic Acids Res.*, **14**, 1029–1044.
  70. Pfalz, J., Bayraktar, O.A., Prikryl, J. and Barkan, A. (2009) Site-specific binding of a PPR protein defines and stabilizes 5' and 3' mRNA termini in chloroplasts. *EMBO J.*, **28**, 2042–2052.
  71. Ruwe, H. and Schmitz-Linneweber, C. (2012) Short non-coding RNA fragments accumulating in chloroplasts: footprints of RNA binding proteins? *Nucleic Acids Res.*, **40**, 3106–3116.
  72. Rojas, M., Ruwe, H., Miranda, R.G., Zoschke, R., Hase, N., Schmitz-Linneweber, C. and Barkan, A. (2018) Unexpected functional versatility of the pentatricopeptide repeat proteins PGR3, PPR5 and PPR10. *Nucleic Acids Res.*, **46**, 10448–10459.
  73. Walter, M., Kilian, J. and Kudla, J. (2002) PNPase activity determines the efficiency of mRNA 3'-end processing, the degradation of tRNA and the extent of polyadenylation in chloroplasts. *EMBO J.*, **21**, 6905–6914.
  74. Yehudai-Resheff, S., Zimmer, S.L., Komine, Y. and Stern, D.B. (2007) Integration of chloroplast nucleic acid metabolism into the phosphate deprivation response in *Chlamydomonas reinhardtii*. *Plant Cell*, **19**, 1023–1038.
  75. Zimmer, S.L., Schein, A., Zipor, G., Stern, D.B. and Schuster, G. (2009) Polyadenylation in *Arabidopsis* and *Chlamydomonas* organelles: the input of nucleotidyltransferases, poly(A) polymerases and polynucleotide phosphorylase. *Plant J.*, **59**, 88–99.
  76. Luro, S., Germain, A., Sharwood, R.E. and Stern, D.B. (2013) RNase J participates in a pentatricopeptide repeat protein-mediated 5' end

- maturation of chloroplast mRNAs. *Nucleic Acids Res.*, **41**, 9141–9151.
77. Johnson, X., Wostrikoff, K., Finazzi, G., Kurat, R., Schwarz, C., Bujaldon, S., Nickelsen, J., Stern, D.B., Wollman, F.-A. and Vallon, O. (2010) MRL1, a conserved pentatricopeptide repeat protein, is required for stabilization of rbcL mRNA in *Chlamydomonas* and *Arabidopsis*. *Plant Cell*, **22**, 234–248.
  78. Germain, A., Kim, S.H., Gutierrez, R. and Stern, D.B. (2012) Ribonuclease II preserves chloroplast RNA homeostasis by increasing mRNA decay rates, and cooperates with polynucleotide phosphorylase in 3' end maturation. *Plant J.*, **72**, 960–971.
  79. Zurawski, G., Perrot, B., Bottomley, W. and Whitfield, P.R. (1981) The structure of the gene for the large subunit of ribulose 1,5-bisphosphate carboxylase from spinach chloroplast DNA. *Nucleic Acids Res.*, **9**, 3251–3270.
  80. Zurawski, G. and Clegg, M.T. (1987) Evolution of higher-plant chloroplast DNA-encoded genes: implications for structure-function and phylogenetic studies. *Ann. Rev. Plant Physiol.*, **38**, 391–418.
  81. Lama, L., Cobo, J., Buenaventura, D. and Ryan, K. (2019) Small RNA-seq: The RNA 5'-end adapter ligation problem and how to circumvent it. *J. Biol. Methods*, **6**, e108.
  82. Xu, H., Yao, J., Wu, D.C. and Lambowitz, A.M. (2019) Improved TGIRT-seq methods for comprehensive transcriptome profiling with decreased adapter dimer formation and bias correction. *Sci. Rep.*, **9**, 7953.
  83. Liebers, M., Grübler, B., Chevalier, F., Lerbs-Mache, S., Merendino, L., Blanvillain, R. and Pfannschmidt, T. (2017) Regulatory shifts in plastid transcription play a key role in morphological conversions of plastids during plant development. *Front. Plant Sci.*, **8**, 23.
  84. Silhavy, D. and Maliga, P. (1998) Mapping of promoters for the nucleus-encoded plastid RNA polymerase (NEP) in the *iojap* maize mutant. *Curr. Genet.*, **33**, 340–344.
  85. Lloréns-Rico, V., Cano, J., Kamminga, T., Gil, R., Latorre, A., Chen, W.-H., Bork, P., Glass, J.I., Serrano, L. and Lluch-Senar, M. (2016) Bacterial antisense RNAs are mainly the product of transcriptional noise. *Sci. Adv.*, **2**, e1501363.
  86. Georg, J. and Hess, W.R. (2018) Widespread antisense transcription in prokaryotes. *Microbiol. Spectr.*, **6**, doi:10.1128/microbiolspec.RWR-0029-2018.
  87. Shimmura, S., Nozoe, M., Kitora, S., Kin, S., Matsutani, S., Ishizaki, Y., Nakahira, Y. and Shiina, T. (2017) Comparative analysis of chloroplast psbD promoters in terrestrial plants. *Front. Plant Sci.*, **8**, 1186.
  88. Hotto, A.M., Germain, A. and Stern, D.B. (2012) Plastid non-coding RNAs: emerging candidates for gene regulation. *Trends Plant Sci.*, **17**, 737–744.
  89. Vera, A., Matsubayashi, T. and Sugiura, M. (1992) Active transcription from a promoter positioned within the coding region of a divergently oriented gene: the tobacco chloroplast rpl32 gene. *Mol. Gen. Genet.*, **233**, 151–156.
  90. Ohyama, K., Fukuzawa, H., Kohchi, T., Shirai, H., Sano, T., Sano, S., Umesono, K., Shiki, Y., Takeuchi, M., Chang, Z. *et al.* (1986) Chloroplast gene organization deduced from complete sequence of liverwort *Marchantia polymorpha* chloroplast DNA. *Nature*, **322**, 572–574.
  91. Shinozaki, K., Ohme, M., Tanaka, M., Wakasugi, T., Hayashida, N., Matsubayashi, T., Zaita, N., Chunwongse, J., Obokata, J., Yamaguchi-Shinozaki, K. *et al.* (1986) The complete nucleotide sequence of the tobacco chloroplast genome: Its gene organization and expression. *EMBO J.*, **5**, 2043–2050.
  92. Bazzini, A.A., Johnstone, T.G., Christiano, R., Mackowiak, S.D., Obermayer, B., Fleming, E.S., Vejnar, C.E., Lee, M.T., Rajewsky, N., Walther, T.C. *et al.* (2014) Identification of small ORFs in vertebrates using ribosome footprinting and evolutionary conservation. *EMBO J.*, **33**, 981–993.
  93. Washietl, S., Findeiss, S., Müller, S.A., Kalkhof, S., von Bergen, M., Hofacker, I.L., Stadler, P.F. and Goldman, N. (2011) RNCODE: robust discrimination of coding and noncoding regions in comparative sequence data. *RNA*, **17**, 578–594.
  94. Slavoff, S.A., Mitchell, A.J., Schwaid, A.G., Cabili, M.N., Ma, J., Levin, J.Z., Karger, A.D., Budnik, B.A., Rinn, J.L. and Saghatelian, A. (2013) Peptidomic discovery of short open reading frame-encoded peptides in human cells. *Nat. Chem. Biol.*, **9**, 59–64.
  95. Wicke, S., Schneeweiss, G.M., dePamphilis, C.W., Müller, K.F. and Quandt, D. (2011) The evolution of the plastid chromosome in land plants: gene content, gene order, gene function. *Plant Mol. Biol.*, **76**, 273–297.
  96. Shi, C., Liu, Y., Huang, H., Xia, E.-H., Zhang, H.-B. and Gao, L.-Z. (2013) Contradiction between plastid gene transcription and function due to complex posttranscriptional splicing: an exemplary study of ycf15 function and evolution in angiosperms. *PLoS One*, **8**, e59620.
  97. Schmitz-Linneweber, C., Maier, R.M., Alcaraz, J.P., Cottet, A., Herrmann, R.G. and Mache, R. (2001) The plastid chromosome of spinach (*Spinacia oleracea*): complete nucleotide sequence and gene organization. *Plant Mol. Biol.*, **45**, 307–315.
  98. Raubeson, L.A., Peery, R., Chumley, T.W., Dziubek, C., Fourcade, H.M., Boore, J.L. and Jansen, R.K. (2007) Comparative chloroplast genomics: analyses including new sequences from the angiosperms *Nuphar advena* and *Ranunculus macranthus*. *BMC Genomics*, **8**, 174.
  99. Hammani, K., Bonnard, G., Bouchoucha, A., Gobert, A., Pinker, F., Salinas, T. and Giegé, P. (2014) Helical repeats modular proteins are major players for organelle gene expression. *Biochimie*, **100**, 141–150.
  100. Barkan, A. and Small, I. (2014) Pentatricopeptide repeat proteins in plants. *Annu. Rev. Plant Biol.*, **65**, 415–442.
  101. Zhelyazkova, P., Hammani, K., Rojas, M., Voelker, R., Vargas-Suarez, M., Borner, T. and Barkan, A. (2012) Protein-mediated protection as the predominant mechanism for defining processed mRNA termini in land plant chloroplasts. *Nucleic Acids Res.*, **40**, 3092–3105.
  102. Yagi, Y., Hayashi, S., Kobayashi, K., Hirayama, T. and Nakamura, T. (2013) Elucidation of the RNA recognition code for pentatricopeptide repeat proteins involved in organelle RNA editing in plants. *PLoS One*, **8**, e57286.
  103. Takenaka, M., Zehrmann, A., Brennicke, A. and Graichen, K. (2013) Improved Computational target site prediction for pentatricopeptide repeat RNA editing factors. *PLoS One*, **8**, e65343.
  104. Yan, J., Yao, Y., Hong, S., Yang, Y., Shen, C., Zhang, Q., Zhang, D., Zou, T. and Yin, P. (2019) Delineation of pentatricopeptide repeat codes for target RNA prediction. *Nucleic Acids Res.*, **47**, 3728–3738.
  105. Barkan, A., Rojas, M., Fujii, S., Yap, A., Chong, Y.S., Bond, C.S. and Small, I. (2012) A combinatorial amino acid code for RNA recognition by pentatricopeptide repeat proteins. *PLoS Genet.*, **8**, e1002910.
  106. Prikrýl, J., Rojas, M., Schuster, G. and Barkan, A. (2011) Mechanism of RNA stabilization and translational activation by a pentatricopeptide repeat protein. *Proc. Natl. Acad. Sci. U.S.A.*, **108**, 415–420.
  107. Yehudai-Resheff, S., Portnoy, V., Yogeve, S., Adir, N. and Schuster, G. (2003) Domain analysis of the chloroplast polynucleotide phosphorylase reveals discrete functions in RNA degradation, polyadenylation, and sequence homology with exosome proteins. *Plant Cell*, **15**, 2003–2019.
  108. Chi, W., He, B., Mao, J., Jiang, J. and Zhang, L. (2015) Plastid sigma factors: their individual functions and regulation in transcription. *Biochim. Biophys. Acta Bioenerg.*, **1847**, 770–778.
  109. Pfannschmidt, T., Blanvillain, R., Merendino, L., Courtois, F., Chevalier, F., Liebers, M., Grubler, B., Hommel, E. and Lerbs-Mache, S. (2015) Plastid RNA polymerases: orchestration of enzymes with different evolutionary origins controls chloroplast biogenesis during the plant life cycle. *J. Exp. Bot.*, **66**, 6957–6973.
  110. Zhou, W., Lu, Q., Li, Q., Wang, L., Ding, S., Zhang, A., Wen, X., Zhang, L. and Lu, C. (2017) PPR-SMR protein SOT1 has RNA endonuclease activity. *Proc. Natl. Acad. Sci. U.S.A.*, **114**, E1554–E1563.
  111. Yang, Z., Hou, Q., Cheng, L., Xu, W., Hong, Y., Li, S. and Sun, Q. (2017) RNase H1 cooperates with DNA Gyrase to restrict R-Loops and maintain genome integrity in *Arabidopsis* chloroplasts. *Plant Cell*, **29**, 2478–2497.
  112. Condon, C., Pilon, J. and Braun, F. (2018) Distribution of the ribosome associated endonuclease Rael and the potential role of conserved amino acids in codon recognition. *RNA Biol.*, **15**, 1–6.
  113. Ji, D., Manavski, N., Meurer, J., Zhang, L. and Chi, W. (2019) Regulated chloroplast transcription termination. *Biochim. Biophys. Acta Bioenerg.*, **1860**, 69–77.

114. Porrua, O., Boudvillain, M. and Libri, D. (2016) Transcription termination: variations on common themes. *Trends Genet.*, **32**, 508–522.
115. Chi, W., He, B., Manavski, N., Mao, J., Ji, D., Lu, C., Rochaix, J.D., Meurer, J. and Zhang, L. (2014) RHON1 mediates a Rho-Like activity for transcription termination in plastids of *Arabidopsis thaliana*. *Plant Cell*, **26**, 4918–4932.
116. Dar, D. and Sorek, R. (2018) High-resolution RNA 3'-ends mapping of bacterial Rho-dependent transcripts. *Nucleic Acids Res.*, **46**, 6797–6805.
117. Romani, I., Manavski, N., Morosetti, A., Tadini, L., Maier, S., Kühn, K., Ruwe, H., Schmitz-Linneweber, C., Wanner, G., Leister, D. *et al.* (2015) A member of the Arabidopsis mitochondrial transcription termination factor family is required for maturation of chloroplast transfer RNA<sup>Ile</sup> (GAU). *Plant Physiol.*, **169**, 627–646.
118. Zoschke, R., Watkins, K.P. and Barkan, A. (2013) A rapid ribosome profiling method elucidates chloroplast ribosome behavior in vivo. *Plant Cell*, **25**, 2265–2275.
119. Zuker, M. (2003) Mfold web server for nucleic acid folding and hybridization prediction. *Nucleic Acids Res.*, **31**, 3406–3415.