

Differences in firing efficiency, chromatin and transcription underlie the developmental plasticity of Arabidopsis originome

Sequeira-Mendes, J.¹, Vergara, Z.¹, Peiró, R.¹, Morata, J.², Aragüez, I.¹, Costas, C.^{1,3}, Mendez-Giraldez, R.^{1,4}, Casacuberta, J.M.², Bastolla, U.^{1,*}, Gutierrez, C.^{1,*}

1 Centro de Biología Molecular Severo Ochoa, CSIC-UAM, Nicolas Cabrera 1, Cantoblanco, 28049 Madrid, Spain

2 Center for Research in Agricultural Genomics, CRAG (CSIC-IRTA-UAB-UB), Campus Universitat Autònoma de Barcelona, Bellaterra, Cerdanyola del Valles, 08193 Barcelona, Spain

3 Current address: Departamento de Química Física, Universidad de Vigo, 36310 Vigo, Spain

4 Current address: Lineberger Comprehensive Cancer Center, Department of Genetics, University of North Carolina at Chapel Hill, Chapel Hill, North Carolina, USA

* Correspondence to: CG (cgutierrez@cbm.csic.es) or UB (ubastolla@cbm.csic.es)

SUMMARY

Full genome replication in eukaryotes depends on the function of thousands of DNA replication origins (ORIs) that constitute the originome. The identification of ORIs in cultured cells has served to define their basic DNA and chromatin features. However, a major challenge is to learn the biology of ORIs in adult organisms to understand their developmental plasticity. Here, we have determined the originome and chromatin landscape in *Arabidopsis thaliana* at two stages of vegetative development. We found that ORIs associate with multiple chromatin signatures including the most frequent at TSS but also at proximal and distal gene regulatory regions or repressed Polycomb (PcG) domains. In constitutive heterochromatin, a high fraction of ORIs colocalize with GC-rich retrotransposons. Quantitative analysis of ORI activity led us to conclude that strong ORIs possess high scores of local GC content and clusters of GGN trinucleotides that may form G quadruplexes and other G-rich structures. We also found that development primarily influences ORI firing strength rather than ORI location. Moreover, ORIs that preferentially fire at early vegetative stages colocalize with GC-rich heterochromatin whereas those at later stages associate with transcribed genes. Our study provides the first originome of an adult organism, which is a developmentally and genetically tractable, thus opening new avenues of studying the biology of ORIs in response to developmental cues, transcriptional programs, environmental challenges and mutant backgrounds.

INTRODUCTION

Replication of the large and complex genomes of multicellular organisms occurs during S-phase, once per cell cycle. Full genome replication is not achieved from a single DNA replication origin (ORI) but instead it depends on the coordinated function of thousands of ORIs scattered across the genome (Mechali, 2010; Mechali et al., 2013; Sanchez et al., 2012). The association of pre-replication complexes (pre-RCs) with DNA determines the specification of genomic sites that can potentially act as ORIs. Some clues on the contribution of DNA sequence, chromatin marks, transcription factor binding sites and GC content to ORI specification and activity have been obtained (Costas et al., 2011; Gutierrez et al., 2016; Leonard and Mechali, 2013; Mechali, 2010; Mechali et al., 2013; Vergara and Gutierrez, 2017). However, in spite of extensive efforts, the molecular determinants that specify the genomic location of ORIs in eukaryotes are still largely unknown. One of the reasons for this scarcity in our knowledge about ORI specification is the relative lack of well-characterized examples.

Genomic approaches carried out in various cultured cell models, including mammals, insects and plants, have helped to gain a picture supporting the view that ORIs frequently colocalize with chromatin marks associated with active chromatin (Cayrou et al., 2015; Cayrou et al., 2011; Comoglio et al., 2015; Costas et al., 2011; Karnani et al., 2010; Macalpine et al., 2010; Sequeira-Mendes et al., 2009). In addition, ORIs in the euchromatin of multicellular organisms are associated with GC-stretches, of which some may form G quadruplexes (Castillo Bosch et al., 2014; Cayrou et al., 2012; Valton et al., 2014). However there are large genomic regions where marks of active chromatin are absent but must be replicated, posing the question of whether there is a single chromatin signature able to define ORI location.

One complication in identifying the mechanisms of ORI specification is the variability introduced by cell culture conditions (Mesner et al., 2013), which most studies have used so far. Unlike *in vitro* culture conditions, cells within the body of a multicellular organism are subjected to hormonal signals and developmental cues that influence cell proliferation, cell fate decisions and differentiation. These intra- and extracellular factors, which are lost in the *in vitro* cell culture studies, can be crucial for the integration of genome replication with cell proliferation during development. Therefore, one major challenge is to identify ORIs in the

cells of an adult organism to assess potential effects of cell fate acquisition and developmental cues on ORI specification.

Several studies have shown that in eukaryotes only a subset of ORIs is activated at each replication round. Thus, besides ORI specification, it is necessary to quantify the variability of ORI activity and investigate the factors that influence it. This has been recently approached in *Caenorhabditis elegans* by analyzing embryos of different ages (Pourkarimi et al., 2016; Rodriguez-Martinez et al., 2017). Here we have taken the challenge of defining the originome, that is the localization of all ORIs, and studying the variability of their activity in a living organism beyond the embryo stage. To identify ORIs in a whole organism and study their plasticity during development we used the model plant *Arabidopsis thaliana* at two stages of vegetative development: 4 day-old seedlings (shortly after germination, when the hormonal and developmental signals necessary for vegetative growth have been established) and 10 day-old seedlings (before the transition to reproductive development). Since the young seedling contains a mixed population of dividing, endoreplicating, differentiating, embryo-derived and stem cells ((Gutierrez, 2005); [Figure 1A](#)), our approach should provide a collection of ORIs active in a wide variety of cell types and developmental stages. With this strategy we sought to obtain an understanding of (i) the molecular determinants of ORI specification and function, and (ii) their relationships with cell proliferation potential and the gene expression programs. We have found that ORIs are located in regions of all known chromatin states, indicating that there is not a single epigenetic signature that defines an ORI. Instead, different chromatin states do appear to affect the ORI relative strength because ORIs in active chromatin, though more prevalent, show lower signal intensity than the less abundant heterochromatin ORIs.

This suggests that once a location is chosen as an ORI in less accessible regions, it is maintained in more cells of the population. Finally, we have identified a subset of ORIs, preferentially active in 4 day-old seedlings, which are located in heterochromatin and another subset preferred in 10 day-old seedlings, associated with transcriptionally active genes. These results are supporting the association of ORI activity with certain DNA sequence, chromatin features and developmentally regulated transcriptional programs.

RESULTS

Identification of ORIs and their replicative strength

Active ORIs are characterized by the presence of newly synthesized single-stranded DNA (ssDNA) molecules, also known as nascent strands (NS). NS purification from whole seedlings is challenging because of the limited amount of NS, even in highly proliferating cultured cells. Here, we have (i) implemented procedures to obtain sufficient amounts of a clean NS sample from whole plant seedlings, (ii) designed protocols to reduce possible biases associated with NS preparation and dsDNA conversion and (iii) developed computational tools to analyze ORIs in a quantitative manner.

NS were isolated from *Arabidopsis* seedlings in two stages of vegetative development: 4 day-old seedlings, soon after germination and 10 day-old seedlings, before the transition to reproductive development, which in both cases contain proliferating and endoreplicating cells of various differentiation stages and types. Our enhanced procedure yielded sufficient amounts of clean NS samples (see [Experimental Procedures](#), [Figure 1A](#) and [Figure S1](#)). Briefly, after nuclei purification, NS were purified from DNA replication bubbles by sucrose gradient centrifugation to isolate DNA fragments of appropriate size in several gradient fractions (300 bp < nascent strands < 2 kb, longer than Okazaki fragments but without compromising resolution). Any contaminating DNA degradation products were removed by λ -exonuclease (λ -Exo) treatment because they are not protected by an RNA primer, as *bona fide* NS fragments are. It has been reported that the λ -Exo treatment produces a bias towards GC-rich DNA sequences (Foulk et al., 2015). However, this is significantly reduced

provided that the treatment is carried out at least twice and under optimal conditions of substrate and λ -Exo concentrations (Cayrou et al., 2015; Cayrou et al., 2011; Picard et al., 2014; Comoglio et al., 2015; Lombrana et al., 2016).

Purified NS were further processed to generate libraries and submitted to sequencing ([Experimental Procedures](#)). Non-redundant sequence reads were aligned to the Arabidopsis genome (TAIR10) and those uniquely mapping were kept for further analysis. To allow a stringent identification of ORIs we carried out three independent experiments for each developmental stage and processed 2-3 consecutive sucrose gradient fractions in each case. We reasoned that ORIs are in principle active in all the experiments, but with different probabilities of firing in the cell population. We identified a reliable set of ORIs active in at least two experiments and quantified their strength in each experiment through the excess of NS reads with respect to the genomic control ([Figure 1B](#) and [Experimental Procedures](#)). In this way, the NSS value of each ORI fully characterizes all samples (two developmental stages and three independent experiments), whereas the locations of the ORIs do not vary since they are obtained by combining all samples. We computed weighted averages over the set of ORIs using the NSS as weights, so that ORIs that are weak in all samples contribute little to the weighted averages, making our method robust to false positives. Conversely, the strong correlation that we found between NSS in different samples makes it unlikely that we missed any strong ORI, as we verified by visual inspection in the genome browser. To generate the originome we identified ORIs using our own peak-calling ZPeaks algorithm because (1) it provides a well-defined profile of the NSS over all of the genome, crucial for our analysis, and (2) it localizes an ORI at the local maximum of the NSS over the ORI box called, which is needed for carefully centering the metaplots ([Figure 1B](#)).

We also carried out experiments to evaluate a possible bias in peak detection introduced during the dsDNA conversion step prior to library preparation. We found that the routine protocol using random primers mixed with heat-denatured genomic DNA and fast cooling to 37 °C rendered a collection of peaks when compared with the untreated genomic DNA. These biased regions were locally enriched in GC and, consequently, they may overlap with ORIs, which are also locally enriched in GC in many eukaryotes (Comoglio et al., 2015; Costas et al., 2011; Mechali et al., 2013). Therefore, we searched for conditions that reduced the dsDNA conversion (dsc) bias and found that carrying the primer-annealing step slowly (see [Experimental Procedures](#)) halved the number of biased regions from 6479 to 3223. Therefore, we used this procedure with our NS samples, which we believe is advisable in all ORI mapping approaches that require a dsDNA conversion step. To remove any residual bias we used as background controls genomic DNA sheared, denatured and converted into dsDNA as for the NS samples. Importantly, the stringent strategy used to control for the dsDNA conversion bias has shown effective since (1) ~80% of ORIs do not overlap with biased regions and possess high NSS and CDC6 (data from Costas et al., 2011) values, as expected for *bona fide* ORIs, (2) likewise, ORIs that overlap with biased regions present even higher CDC6 binding and higher NSS, despite the dsDNA conversion bias contributes negatively to the NSS, and (3) ~84% of biased regions do not overlap with ORIs and show a vanishing NSS score as expected for random genomic regions ([Figure 1C](#)).

The ORI midpoint was identified with a resolution of 25 bp, the bin size used in the ZPeaks algorithm. We found that when an ORI was identified in different samples, the midpoint varied ~120 bp on average, which estimates the precision of our measurements. The ORI midpoint was then computed as the weighted average of the midpoints of the individual samples. The distribution of inter-ORI distance is highly skewed towards small values with a median of 27.6 kb much smaller than the mean of 43.4 kb. We tested ZPeaks both by visual inspection of the overlap between sequencing reads and candidate ORIs ([Figure 1D](#)) and by statistical tests. Together our strategy allowed us to identify a robust set of 2374 *bona fide* ORIs that constitute the seedling originome ([Table S1](#)) and can be confidently used to analyze the determinants of ORI specification.

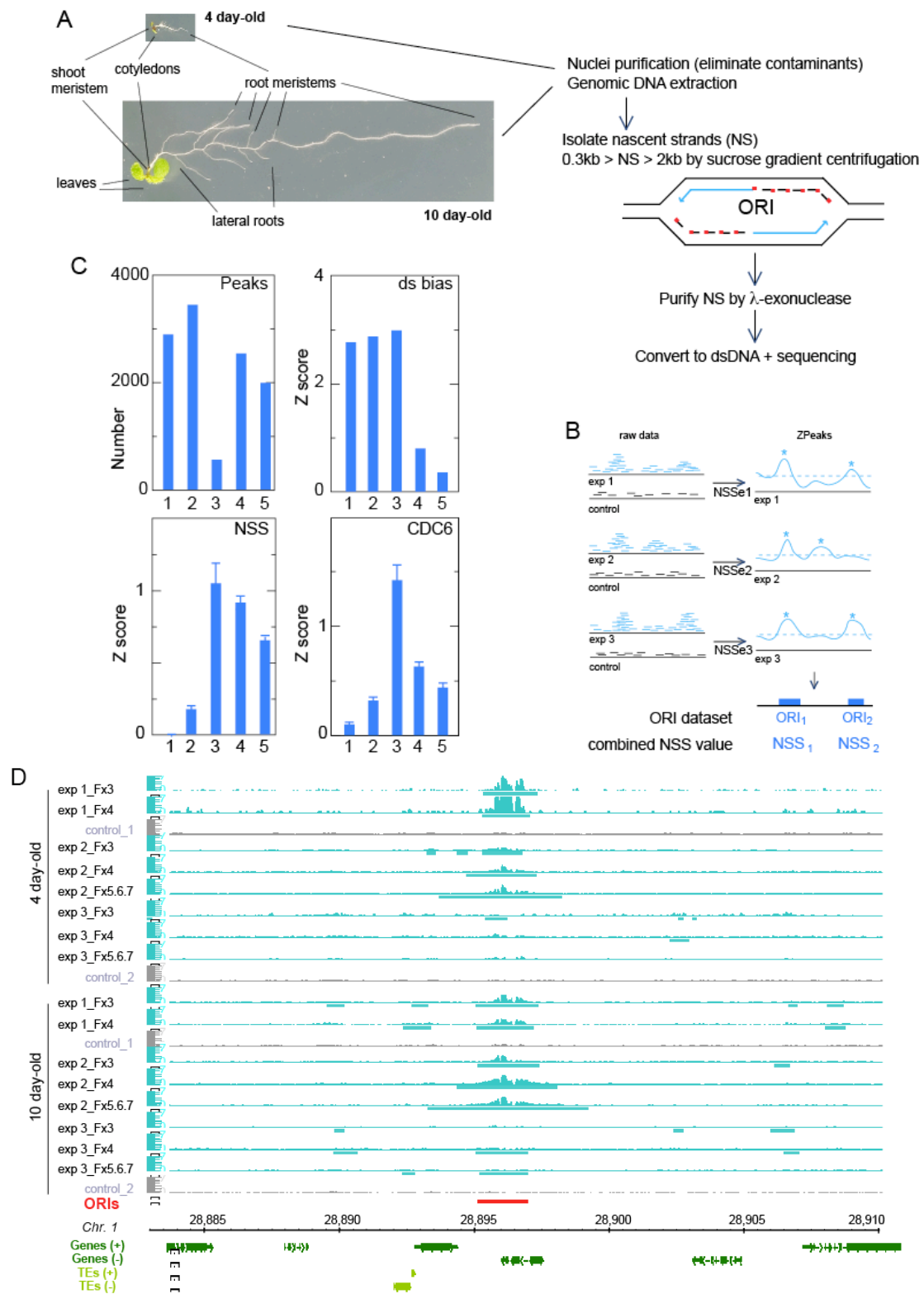


Figure 1A-D. DNA replication origin (ORI) identification in whole Arabidopsis seedlings and evaluation of reproducibility and quality of sequencing datasets.

(A) Summary of basic steps for purification of nascent strands (NS) from seedlings at two developmental stages of Arabidopsis vegetative growth. Seedlings contain cells undergoing cell proliferation and the endocycle in different locations. In 4 day-old seedlings, the shoot and the root apical meristems contain dividing cells whereas the cotyledons and the transition domain above the root apical meristem contain endocycling cells. The rest of the root is made up by different cell types some of them dividing and some differentiated. In addition to all these organs and cell types, 10 day-old seedlings contain growing leaf

primordia and lateral root primordia, with proliferating cells, as well as a longer root with more differentiated cells.

(B) Flowchart summary of the assignment of the nascent strand score (NSS) and the ZPeaks algorithm used to identify ORIs.

(C) Quality controls of the NS purification and dsDNA conversion step. The horizontal axis represents the following datasets: class 1, biased peaks that do not overlap with the ORI set; class 2, all peaks with dsDNA conversion bias; class 3, biased peaks overlapping with ORIs; class 4, all ORIs; class 5, ORIs that do not overlap with the biased set.

(D) Representative genome browser view of a ~35 kb region of chromosome 1, to illustrate ORI identification in various sucrose gradient fractions of the three independent experiments in 4 and 10 day-old seedlings.

We found that the NSS values are broadly distributed and span several orders of magnitude. Importantly, the NSS of different experiments are strongly correlated with each other, with correlation coefficients ranging from 0.65 to 0.92 (Figure S2). This result lends support to our approach and suggests that some intrinsic properties of the ORIs influence their firing rates in all the experiments although with important variations, as discussed below. Despite that the firing probability should saturate at 100% probability, the NSS did not show any sign of saturation, suggesting that they are far from 100% activity. To facilitate the comparison between the two developmental stages, we obtained scores that averaged the three independent experiments performed for each stage, as described in Experimental Procedures.

Genomic landscape of ORI locations

To determine the preferences of ORI localizations, we examined the association of ORIs with various genomic elements. Most ORIs (>78%) are associated with genic regions, including 1 kb upstream regions, much more than expected by chance. Within genes, ORIs locate more frequently in exons (Figure 2A). Intergenic regions and TEs comprise ~5% and ~13% of ORIs, respectively, while these genomic regions represent a much larger fraction of the genome (~15% and ~21%, respectively). We recently discovered that ~5% of ORIs active in cultured *Arabidopsis* cells colocalize with TEs, although in the gene-poor pericentromeric heterochromatin the frequency of colocalization with TEs increases to ~34%. Moreover, these ORIs are much more frequently located within retrotransposons than in DNA transposons, and among them in TEs of the Gypsy and LINE families (Vergara et al., 2017). Thus, we use the same approach to investigate the distribution of ORIs active in seedlings and found that in pericentromeric regions the frequency of ORIs located within TEs increases significantly compared to the overall genome (Figure 2B). Furthermore ORI-TEs show a preference to colocalize with retrotransposons, as it occurs in cultured cells, although in seedlings the fraction of ORIs located in DNA transposons is higher, in particular for the MuDR family (Figure 2B).

We also found that ORIs have a preference to localize ~0.5 kb downstream from the Transcription Start Site (TSS) and tend to avoid Transcription Termination Sites (TTS) at both developmental stages (Figure 2C).

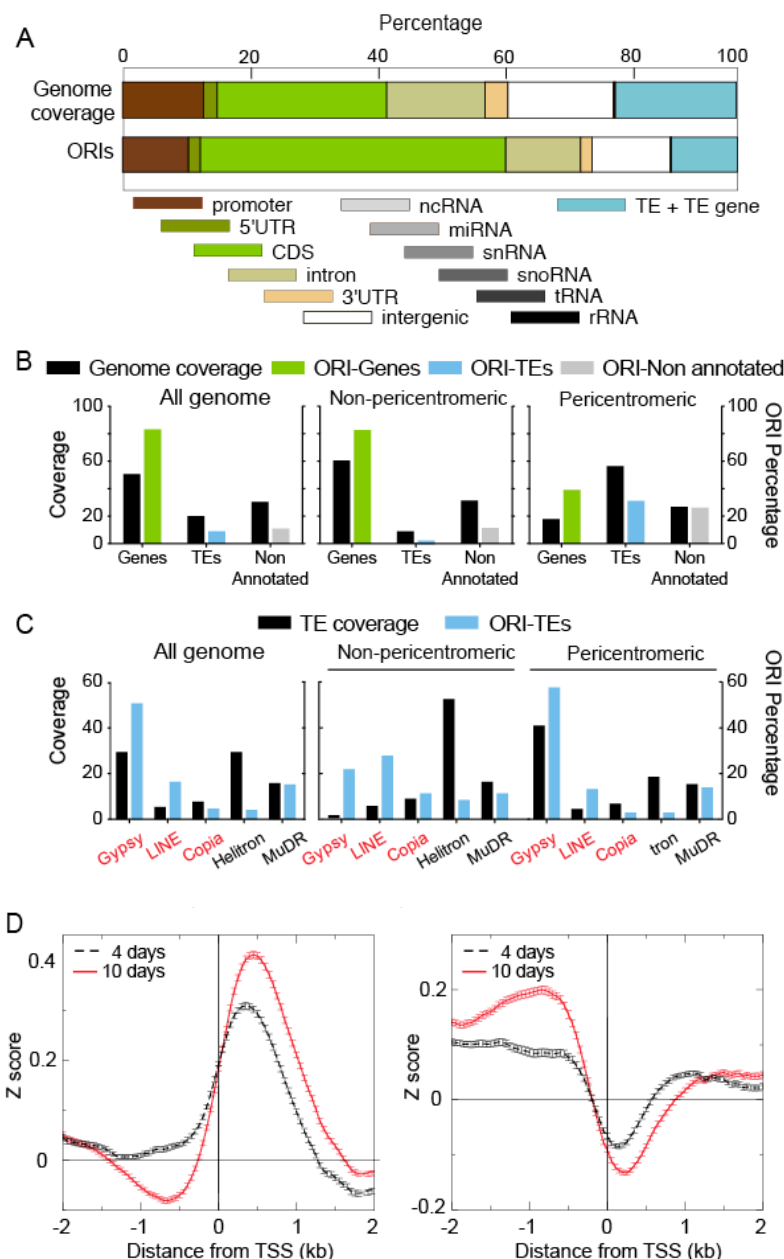


Figure 2. Association of ORIs with genomic elements and TE families.

(A) Relationship between ORI location and genomic elements. The overlap (in base pairs) between the indicated genomic elements and each ORI was computed and expressed as a percentage. A region of 1 kb upstream the coding sequence was considered as the promoter. Note that TEs are large genomic elements that may have one or more TE genes associated with them. Here, the class TE refers to genomic regions that contain TEs but do not overlap with TE genes.

(B) Frequency distribution of ORIs colocalizing with genes (green), TEs (blue) and non-annotated regions (grey) compared with the respective nucleotide coverage.

(C) Frequency distribution of ORI-TEs (blue bars) in TE families in all the Arabidopsis genome, the non-pericentromeric regions

and the pericentromeric regions compared with the respective TE family nucleotide coverage of total TE nucleotides (black bars). In the X-axis, retrotransposon families (red) and DNA transposon families (black).

(D) Metaplots of the combined NSS of the three independent experiments of 4 day-old and 10 day-old seedlings with respect to the transcription start sites (TSS; left panel) or the transcription termination site (TTS; right panel), oriented in both cases with the transcribed RNAs.

Local properties of ORI locations

To study the properties of Arabidopsis ORIs in seedlings, we assessed the average local neighborhood of all ORIs by computing metaplots centered at the ORI midpoint with each ORI being weighted with its own NSS. To account for the asymmetry between the G and C and between the A and T bases, known to be a feature of replication origins, ORIs were oriented along the 5'-3' direction of the strand with positive GC skew. For this analysis each

variable was transformed into a Z score with respect to the entire genome, which allows comparing the strength of different genomic and epigenomic variables.

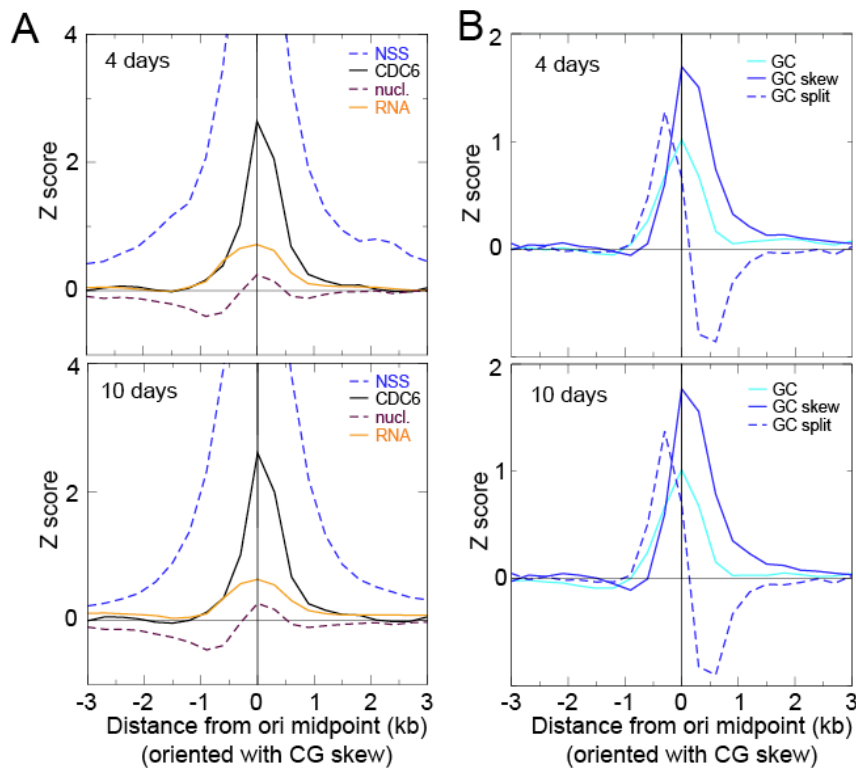


Figure 3. Features of the local neighborhood of ORIs in whole Arabidopsis seedlings.

(A) Metaplots of NSS, CDC6, transcript content (RNA) and nucleosome (nucl.) content weighted with the combined Z score of the three independent experiments of 4 day-old (top) and 10 day-old (bottom) seedlings. The metaplots for individual scores of each experiment are shown in Supplemental Figure S3.

(B) Metaplots of GC, GC skew and GC split weighted with the combined Z score of the three independent experiments of 4 day-old (top) and 10 day-old (bottom) seedlings. The metaplots for individual scores of each experiment are shown in Figure S4.

First of all, we confirmed that the NSS of all experiments had a very prominent peak at the ORI midpoint (Figure 3A and Figure S3). Additionally, we examined the position of the pre-RC protein CDC6 (Costas et al., 2011). As expected, we found that CDC6 has a high peak centered at the midpoint of active ORIs with a width of ± 1 kb (Figure 3A and Figure S3). This provides strong independent support to the peak-calling procedure used here to define ORI location based on nascent-strand mapping. We also found that the location of ORI midpoint coincides with a peak of nucleosome occupancy (Figure 3A and Figure S3), as described for cultured cells in Arabidopsis and mammals (Lombrana et al., 2013; Stroud et al., 2012). Furthermore, the regions around ORIs are more frequently transcribed than the average of the genome, with a broad peak of approximately ± 1 kb centered at the ORI midpoint (Figure 3A and Figure S3).

The Arabidopsis genome is particularly rich in A-T (63.8%). ORIs identified in cultured cells preferentially colocalize with short G+C-rich stretches (Costas et al., 2011). Measuring the G+C content in 100 bp sliding windows across all ORI locations revealed that ORIs in seedlings also colocalize with G+C-rich regions that show a peak of ~ 0.8 kb in width, centered at the ORI midpoint (Figure 3B and Figure S4). The asymmetry between the G and

C nucleotide, called GC skew, is a signature of ORIs both in prokaryotes and metazoa (Arakawa and Tomita, 2012; Cayrou et al., 2011; Comoglio et al., 2015; Macalpine et al., 2010; Xia, 2012). The GC skew presents a strong peak centered at the ORI midpoint but asymmetrically distributed (-0.5 kb through 1.0 kb from ORI midpoint).

The GC skew associated with DNA replication is attributed to the asymmetry of mutation processes on the leading and lagging strand. To test this we defined the GC split as the difference of GC skew downstream and upstream of a genomic point. We observed that the ORI midpoints typically lay in between a strong maximum of the GC split upstream of the ORI midpoint at approximately -0.3 kb, and a slightly less strong minimum at ~0.4 kb (Figure 3B and Figure S4). This result does not support the hypothesis that the different mutation processes at the leading and lagging strand are the sole cause of the GC skew.

ORIs associate with multiple chromatin signatures

The mechanisms responsible for ORI specification along the genome in multicellular eukaryotes remain unknown. Studies in several model organisms, including mammals, *Drosophila* and plants, have revealed a preferential association of ORIs with activating chromatin marks (Cayrou et al., 2015; Cayrou et al., 2011; Comoglio et al., 2015; Costas et al., 2011; Picard et al., 2014; Pourkarimi et al., 2016; Rodriguez-Martinez et al., 2017; Sequeira-Mendes et al., 2009). A simplistic interpretation of these observations may suggest that some combination of chromatin features may be sufficient for ORI specification. However, simple inspection of genomic data clearly shows that there is not a single epigenetic mark or a simple combination of them common to all ORIs. Recent studies demonstrate the existence of three major classes of ORIs with different organization, chromatin environment, and sequence motifs (Cayrou et al., 2015), suggesting that ORIs are associated with different signatures.

To investigate the preferences of ORIs to occur in particular chromatin settings, we assigned each ORI midpoint to one of the chromatin states defined by high-resolution analysis of 16 chromatin and DNA features in the *Arabidopsis* genome (Sequeira-Mendes et al., 2014). These states simplify the combinatorial complexity of DNA and histone marks across the *Arabidopsis* genome into nine chromatin states characterized by unique signatures, in a manner similar to what has been done for *Drosophila* and human cells (Ernst et al., 2011; Kharchenko et al., 2011). Importantly, the chromatin states identified in *Arabidopsis* reflect show a preferred linear sequence defining proximal promoters (state 2) – TSS (state 1) – 5' end of genes (state 3) – long genes (state 7) – 3' end of genes (state 6), followed by repressed states containing Polycomb marks (states 5 and 4) and two types of heterochromatin (states 8 and 9; (Sequeira-Mendes and Gutierrez, 2016)).

The cumulative weight of ORIs in the chromatin states is higher for ORIs colocalizing with state 1 (TSS), which are the most numerous (Figure 4A and Figure S5A). Moreover, ORIs that colocalize with state 2 (proximal promoters and 5'UTRs) and state 3 (5' end of genes) also have a relatively high weight (Figure 4A and Figure S5A). On the contrary, state 5 (PcG-repressed regions) and states 8 and 9 (the two heterochromatin types), state 7 (long coding genes), state 6 (3' end of genes) and in particular state 4 (distal regulatory intergenic regions) contain more moderate amounts of ORIs.

Since not all states have the same frequency in the genome we transformed the normalized weight into propensities by dividing them by the fraction of the genome that belongs to the same state, so that positive values identify states with a NSS weight larger than expected based on its genome coverage. We confirmed that ORIs close to TSS (state 1) show the highest values, followed by ORIs in proximal promoters (state 2; Figure 4B and Figure S5B). Among the rest of ORIs, those in distal regulatory intergenic regions (state 4), long genes (states 7) and 3' end of genes (state 6) showed a negative propensity (Figure 4B and Figure S5B). Interestingly, ORIs in heterochromatin (states 8 and 9) showed a negative propensity in 10 day-old but not in 4 day-old seedlings.

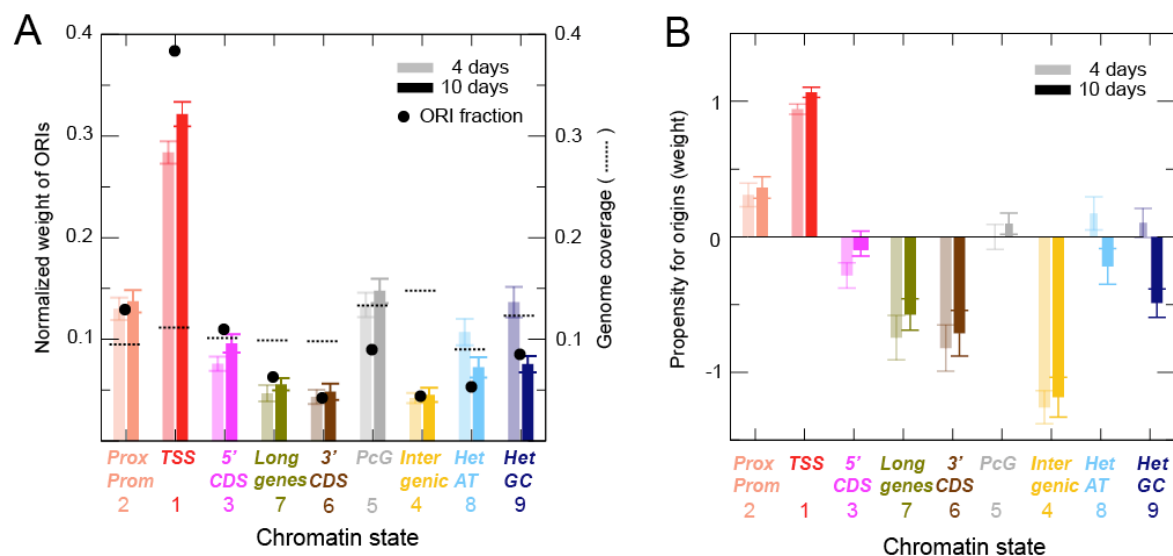


Figure 4. Association of ORIs with chromatin states.

(A) Normalized weight of ORIs belonging to the 9 chromatin states in the combined experiments of 4 (transparent colours) and 10 day-old (solid colours) seedlings. The results of the three independent experiments are shown in the Supplemental Figure S5A. Black circles indicate the fractions of ORIs in each chromatin state. Broken lines indicate the genome coverage of each chromatin state, as indicated in the right Y-axis.

(B) Same as Fig. 4A, showing the propensity (instead of the cumulative weight) for ORIs in the 9 chromatin states is depicted. This reveals ORI types according to the associated chromatin states that have larger NSS than expected by chance based on the fraction of genome that they represent. The results of the three independent experiments are shown in Figure S5B.

Factors associated with ORI specification in different chromatin landscapes

The presence of ORIs across all the different chromatin states clearly demonstrates that ORI activity is not associated with a single chromatin signature but with different signatures depending on ORI localization. To define ORI features in a quantitative manner we calculated five traits over a region of 300 bp around the ORI midpoint for each state and experimental dataset, such as nascent strand score (NSS), CDC6, GC content, GC skew and the number of GGN trinucleotides. All of these traits correlate positively with the NSS, which suggests that they contribute positively to the strength of the ORIs. We then computed weighted averages of these traits (using the NSS as weight) over the ORI of a given state, and transformed them into Z scores, so that a positive value indicates that the trait in that state is larger than the average score over the whole genome.

We found that the GGN score profile across states is very similar to that of CDC6 and, to a lower extent, to those of GC and GC skew, showing that these traits are rather consistent with each other. This is relevant because it might be argued that λ -exo has a lowered activity on G-rich secondary structures (Fouk et al., 2015). Furthermore, as discussed below, the CDC6 binding data were obtained in ChIP experiments and ORIs in Arabidopsis cultured cells are locally enriched in GC using procedures that do not rely on the use of λ -exo (Costas et al., 2011). Moreover, these profiles were consistently similar in all experimental situations tested (Figure 5A and Figure S6). We can distinguish two groups of ORIs with low and high NSS values (Figure 5A). Thus, the average NSS of ORIs in active chromatin states 1 and 3 (TSS and 3'-end of genes) is significantly lower than for ORIs in repressed heterochromatin (states 8 and 9) and Polycomb regions (state 5; the two-tailed t-test $p < 0.0001$ in all cases except for state 3-state 9 ($p < 0.0046$ and $p < 0.0002$ in 10 and 4 day-old seedlings, respectively). ORIs in Polycomb chromatin were consistently high in all other traits, which

Together our results strongly support the following conclusions regarding the definition of different classes of ORIs depending on the chromatin states typical of their neighborhood.

- (1) ORIs located in genic regions (states 1, 3, 6 and 7), associated with active transcription and more open chromatin, possess low or intermediate NSS values, despite their large cumulative weight, suggesting a more variable usage of ORI sites. This is in part explained by the relatively lower values of CDC6 and GGN, and in part suggests interference between replication and transcription.
- (2) The contrary holds for ORIs in heterochromatin, with overall low accessibility for replication proteins, which tend to have a high NSS, despite their low cumulative weight. This suggested to us that once a region is specified as a potential ORI in a disfavored chromatin landscape, it is used more frequently in all cells of the population. This also applies to other compact and poorly transcribed regions such as Polycomb chromatin.

Differences in the average NSS of ORIs of different states may either stem from a global effect that affects all ORIs in the same way, or indicate strong variability of the NSS. To investigate how the chromatin states influence the variability of ORIs, we measured the correlation coefficients of the NSS over the sets of ORIs that belong to a given state (Figure S7). The correlation coefficients are close to one in most states, except state 1 (associated with the TSS), state 8 (AT-rich heterochromatin) and state 4 (intergenic, Polycomb repressed and AT-rich), which are most variable.

Interplay between DNA replication origins and transcriptional programs

The relationship between ORI activity and transcriptional programs during development has been demonstrated (Comoglio et al., 2015; Lubelsky et al., 2014; Muller and Nieduszynski, 2017; Nordman et al., 2011; Siefert et al., 2017). We observed that this is highly dependent on the ORI type according to the chromatin state where it is located, as visualized by plotting the quantity of transcripts in the different chromatin states identified by RNA-seq of 4 and 10 day-old seedlings. To make data comparable we transformed them into Z-scores with respect to the whole genome, so that positive values revealed more transcription than the genomic average.

The first observation is that ORI locations are more transcribed than genomic regions of the same chromatin state at both developmental stages and in all states, except GC-rich heterochromatin (state 9), supporting the strong relationship between DNA replication and transcription. Next, we compared the transcription scores of the two developmental stages. Interestingly, repressed or less transcribed regions (states 5, 4, 8 and 9) showed more transcripts through ORI sites in 4 day-old than in 10 day-old seedlings (Figure 6 and Figure S8). This is particularly striking for ORIs in Polycomb chromatin (state 5) and to a lesser extent in the AT-rich heterochromatin (state 8) where a strong enhancement of transcription in 4 day-old seedlings was observed. The opposite happens for the active regions of long genes (state 7), which are more transcribed in 10 day-old seedlings. In other words, repressed regions in 10 day-old seedlings are more actively transcribed in 4 day-old than, suggesting that there are differences in the accessibility of the different chromatin states, which may affect ORI activity in a developmentally regulated manner.

ORI specification and usage during vegetative development

The systematic differences in transcriptional activity and chromatin organization observed between early (4 day-old) and late (10 day-old) vegetative stages support the idea that chromatin organization, ORI specification and the transcriptional program change during vegetative development. To further investigate these differences, we identified ORIs that

have a higher NSS value in 4 day-old seedlings than in 10 day-old seedlings and vice versa, using the combined NSS of the two stages. We first analyzed the variation of the frequency

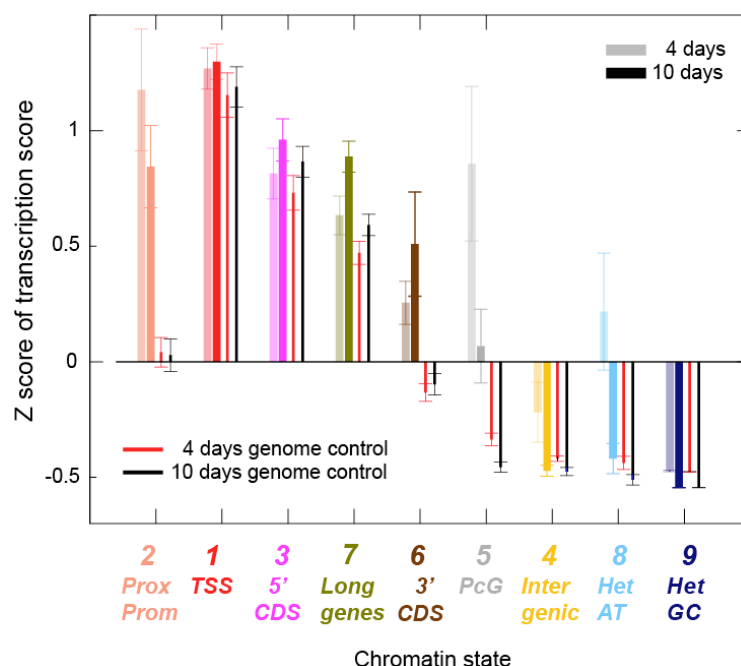


Figure 6. ORIs and transcriptional activity. The average Z score of the transcription score with respect to the average transcription score of the entire genome is shown for ORIs and genomic locations belonging to all chromatin states.

of ORIs in different chromatin states as a function of the threshold used to define the preferred ORIs in each developmental stage (see [Experimental Procedures](#)). To simplify the analysis we grouped ORIs in three classes: genic chromatin (states 2, 1, 3, 7 and 6), Polycomb chromatin (states 5 and 4) and heterochromatin (states 8 and 9). We found that ORIs preferentially used in 4 day-old seedlings have a strong preference for being located in heterochromatin with a very stringent threshold, whereas ORIs preferentially used in 10 day-old seedlings are located in genic states for all thresholds ([Figure 7A](#)). Using the same threshold for both developmental times but strong enough to evidence state preferences, led us to generate a list of 94 ORIs preferentially used in 4 day-old seedlings and 90 ORIs preferentially used in 10 day-old seedlings ([Table S3](#)).

We analyzed several genomic features of these developmentally regulated ORIs. ORIs preferentially active in 10 day-old seedlings possess lower NSS, CDC6, GC and GGN scores than the average over all ORIs, indicating that they are weak ([Figure S9](#)). In contrast, ORIs preferentially used in 4 day-old seedlings have NSS, CDC6, GC and GGN scores comparable to or only slightly lower than the average ORIs. The most distinctive feature between these specific ORI sets is the RNA score that measures the amount of RNA reads across the regions where ORIs are located. The RNA score was comparable to the average over all ORIs for ORIs preferentially expressed in 10-day old seedlings, whereas genomic regions containing ORIs stronger in 4-day old seedlings were significantly less transcribed.

The ORIs preferentially activated in a developmental stage-specific manner were not distributed randomly among the chromatin states. ORIs preferentially activated in 4 day-old plants occur much more frequently in the two types of heterochromatin ([Figure 7B](#); [Table S3](#)). Furthermore, these ORIs in TEs are more frequently located in pericentromeric heterochromatin of 4 day-old seedlings and, among them, more skewed towards Gypsy elements (81.8%), in line with data obtained in Arabidopsis cultured cells (Vergara et al., 2017). In contrast, ORIs preferentially activated in 10 day-old plants colocalize more frequently with genic regions, typically in promoters, TSS and proximal promoters ([Figure 7B](#)). Since ORI activation is highly related to chromatin accessibility, our data strongly

suggest that ORI usage changes significantly in a locus-specific manner during vegetative development, most likely associated with changes in chromatin organization.

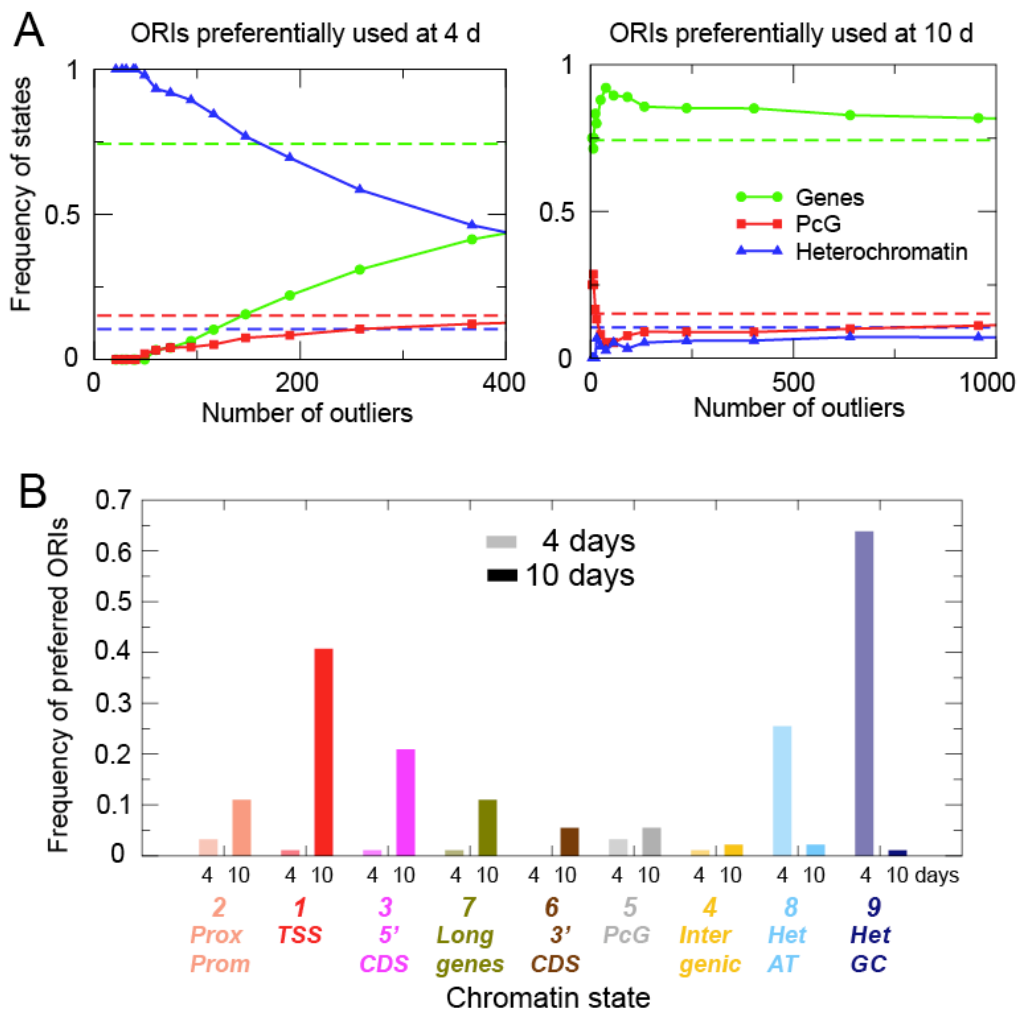


Figure 7. Properties of developmentally regulated ORIs.

(A) Frequency of ORIs preferentially stronger in 4 day-old seedlings (left panel) or in 10 day-old seedlings (right panel), as a function of the number of ORIs obtained by varying the threshold. To simplify the analysis, ORIs have been grouped into those associated with genes (states 2, 1, 3, 7 and 6), with Polycomb chromatin (states 4 and 5) and with heterochromatin (states 8 and 9). Dotted lines represent the average values for the entire genome.

(B) Frequency of ORIs preferentially used at two developmental stages defined as in panel A colocalizing with different chromatin states in 4 day-old (transparent colours) and 10 day-old seedlings (solid colours).

DISCUSSION

Genomic features of Arabidopsis DNA replication origins

In this study, we have generated a whole-body originome map with the location and properties of ORIs in *A. thaliana* plants at two different developmental stages of vegetative growth. ORI activity was assessed quantitatively from the sequencing data by determining a

nascent strand score (NSS) that measures the propensity of a certain genomic location to behave as an ORI.

We have identified 2374 potential ORIs in the Arabidopsis genome that can function during vegetative growth. Our results indicate that Arabidopsis ORIs are organized in discrete sites rather than in large initiation zones, in agreement with ORI mapping in cells with similar genome size (Comoglio et al., 2015; Lombrana et al., 2016; Pourkarimi et al., 2016; Rodriguez-Martinez et al., 2017). A comparison of ORI locations in cultured cells and seedlings revealed a coincidence of ~14-25%, depending on the threshold tolerance. This amount of ORIs common to both sources reveals the existence of technical and biological variables in ORI usage, e.g. the presence of many different cell types in the seedling compared to the cell culture, and that need to be identified in the future using the tools established in this work.

Most Arabidopsis ORIs in seedlings (~78%) also associate with genic elements, in particular the 5' end of genes, reinforcing the strong preference of ORIs for genic regions. Overall this is similar to the situation in metazoan cultured cells and embryos (Cayrou et al., 2015; Comoglio et al., 2015; Macalpine et al., 2010; Rodriguez-Martinez et al., 2017; Sequeira-Mendes et al., 2009). Similar to the situation in cultured cells (Costas et al., 2011), Arabidopsis ORIs colocalize with transposable elements (TEs) less frequently than expected at random. However, as previously found in cultured cells (Vergara et al., 2017), we found that ORIs in pericentromeric regions, which contain a lower gene density, increase their tendency to colocalize within TEs and in particular in retrotransposons of the Gypsy and LINE families. These results reinforce the conclusion that Arabidopsis ORIs have a high preference of being associated with genes in euchromatin and with both genes and transposons in pericentromeric heterochromatin. Moreover, we show that in seedlings the frequency of ORIs colocalizing with DNA transposons, mostly of the MuDR family, is much higher than in cultured cells, perhaps due to a different chromatin landscape in cultured cells and seedlings.

GGN clusters are a strong determinant of ORI strength

The number of GGN trinucleotides within ± 150 nt around the ORI midpoint is the feature that correlates most strongly with the NSS value, although it explains <30% of the variance ($r^2=0.28$). From a structural point of view, four consecutive GGN motifs may form G-rich secondary structures, such as G4 with two tetrads (Chen and Yang, 2012; Sen and Gilbert, 1988). Although some computer programs require at least three tetrads to identify G4 (Todd et al., 2005), examples with two tetrads have been experimentally characterized, such as the thrombin binding aptamer d(G2T2G2TGTG2T2G2) (Macaya et al., 1993) and the *Bombyx mori* telomeric sequence d(AG2T2AG2T2AG2T2AG2) (Sacca et al., 2005). Moreover, the stability of G4 is largest for loops of length 1, such as those in consecutive GGN motifs, and it has been observed that GjNGj sequence motifs form a robust parallel stranded structure motif with 1 nt loop (Chen and Yang, 2012). Stabilizing interactions between distinct G4 structures have been observed (Palumbo et al., 2009), suggesting that there could be synergy between the large number of GGN motifs observed in Arabidopsis ORIs.

Only 35 (1.4%) ORIs contain <4 GGN motifs, while the maximum frequency is between 8 and 16 GGN, finding up to 77 GGN out of the theoretical maximum of 100 in the 300 nt windows. Such nucleotide distribution produces one strand very enriched in Gs, favoring the possibility to form G4 (Cayrou et al., 2015). Our finding that GGN motifs are enriched in ORIs is consistent with reports that the presence of G4 influences ORI activity in animal cells (Besnard et al., 2012; Cayrou et al., 2015; Cayrou et al., 2012; Valton et al., 2014). Moreover, the local GC content is much higher for ORIs than for random genomic locations, which seems to be a common feature of ORIs in all animal and plant cells mapped so far (Besnard et al., 2012; Cayrou et al., 2015; Cayrou et al., 2011; Costas et al., 2011; Macalpine et al., 2010). The enrichment of ORIs in GGN motifs might be related to a reduced

efficiency of I-exo to digest G-mediated secondary structures. However, we have carried out λ -exo treatments under optimal enzyme/substrate conditions. Moreover, it must be kept in mind that (1) we observed a strong correlation between occurrence of GGN motifs at ORIs and the CDC6 binding score, and (2) Arabidopsis ORIs were found to be locally enriched in Gs in cultured cells using procedures that do not rely on λ -exo treated samples (Costas et al., 2011).

The enrichment in GGN stems from two genomic properties of Arabidopsis ORIs, the high GC content and the GC skew, which concur to produce one G-rich strand. As mentioned above, the GC skew associated with replication has been prevalently associated with the mutational asymmetry between the leading and lagging strand (Lobry, 1996). However, we found that the asymmetry between G and C starts approximately 300 nt before the ORI, identified as the local maximum of the NSS. Thus, the replication asymmetry cannot be the sole cause of the GC skew at Arabidopsis ORIs, selection favoring GGN clusters being a likely mechanism that generates GC skew. Moreover, many GGN clusters are formed by tandems of quasi-repeats, so that a likely mutational mechanism is through trinucleotides insertions produced by the slippage of the polymerase. Importantly, the three hypothesized mechanisms, replication asymmetry, insertion and selection, are expected to cooperate for the formation and maintenance of GGN clusters.

Other structural features of Arabidopsis DNA replication origins

Besides GGN clusters we have identified other genomic features associated with the ORI activity. A first important feature is chromatin accessibility because more accessible states, which are transcriptionally active, are more prone to be ORI locations. The second driving factor is the propensity to bind the pre-initiation protein CDC6. Interestingly, the affinity for CDC6 is on average significantly larger in the repressed than in the active chromatin states, probably in order to compensate their reduced accessibility. Finally, despite both transcription and replication activities are affected by chromatin accessibility, we observed an overall negative correlation between the average ORI strength and transcription of the chromatin state. In particular, ORIs colocalizing with TSS are overall among the weakest ones. There are reports of a preferential location of ORIs in actively transcribed genes (Aladjem, 2004; Goren et al., 2008; MacAlpine et al., 2004; Saha et al., 2004). However, in these cases the ORIs are not in close proximity to the TSS. We hypothesize that the negative correlation may originate from possible interference between replication and transcription (Aguilera and Garcia-Muse, 2013).

Our data clearly demonstrate the existence of different types of ORIs according to their features, including primarily their chromatin landscape. There are multiple signatures that can accommodate ORIs to ensure full replication of the entire genome, although ORIs show an overall preference for localizing in the TSS (state 1) and adjacent states containing relatively open chromatin properties. This points to the compatibility of ORI activity with multiple chromatin signatures (Cayrou et al., 2015). Long genes (state 7), the 3'-end of coding sequences (state 6) and the distal regulatory intergenic regions (state 4) are significantly depleted of ORIs. Surprisingly, ORIs located within PcG repressed chromatin (state 5) as well as within heterochromatic domains (states 8 and 9) tend to be stronger than average. Thus, it is conceivable that finding an appropriate local ORI landscape within repressed and compact chromatin may favor that it is used in more cells of the population, leading to a higher NSS value.

Developmental regulation of DNA replication origins and the chromatin landscape

We found remarkable that ORI activity undergoes systematic changes in the course of development. In particular, we observed that ORIs that are stronger in 10 day-old seedlings are particularly frequent in the TSS (state 1) and adjacent chromatin regions (states 2 and 3).

In contrast, ORIs preferentially used in 4 day-old seedlings are particularly recurrent in heterochromatin (state 9). These trends are consistent with our transcriptional analysis that revealed lower repression of typically repressed chromatin states in 4 day-old seedlings compared to 10 day-old seedlings, suggesting that chromatin organization may be different at these two stages of vegetative development. The increased frequency of ORIs in TEs, in particular of the Gypsy family, in 4 day-old seedlings suggests that the role of these TEs as ORIs could be important at early developmental stages where the dynamics of cytosine methylation suggests differences in the repression level of heterochromatin (Bouyer et al., 2017). This is a major difference with recent studies in *C. elegans*, where early pregastrula embryos are depleted of ORIs in heterochromatin whereas ORIs have a preference for non-coding regions and enhancers in postgastrula embryos (Rodriguez-Martinez et al., 2017). These differences between animal and plants suggest different mechanisms of coupling ORI activity, developmental programs and heterochromatin dynamics. It is also worth noting that the spatial organization of the Arabidopsis genome reveals that typical TADs and distal enhancers as in animals are lacking, or very infrequent (Liu et al., 2016; Vergara and Gutierrez, 2017; Wang et al., 2015).

Based on ORI identification in various cell types in culture, a general consensus exists that developmentally regulated ORIs are not very efficient (Besnard et al., 2012). The quantitative parameter of ORI activity (NSS), developed in our study, clearly showed that every genomic location associated with an ORI possesses a certain firing efficiency. In agreement with studies in animal cells in culture (Besnard et al., 2012; Comoglio et al., 2015), it seems that modulation of ORI activity rather than the selection of different genomic locations determines ORI usage at different developmental stages and, most likely, in different cell types.

Polycomb complexes are involved in regulating gene expression associated with developmental phase transitions in Arabidopsis (Kuwabara and Gruissem, 2014). We found ORIs associated with Polycomb-regulated genes (state 5), although they are under-represented in these chromatin regions. Somehow surprisingly, they behave as strong ORIs, suggesting that once a genomic site is chosen as an ORI within a Polycomb-repressed chromatin domain, the position of ORIs is restricted and used in many cells. This is in agreement with the hypothesis developed for mammalian cells in culture, where Polycomb factors are strongly associated with efficiently used ORIs (Cayrou et al., 2015; Cayrou et al., 2011; Picard et al., 2014). Developmentally regulated genes in animal pluripotent stem cells share H3K4me3 and H3K27me3 marks, typical of bivalent chromatin (Bernstein et al., 2006). An attractive possibility is that some of the H3K27me3 regions, colocalizing with ORIs, may actually constitute regions of bivalent chromatin, although this has not been experimentally demonstrated. Using sequential re-ChIP experiments (Sequeira-Mendes et al., 2014) we demonstrated the presence of this bivalent chromatin type in somatic cells of Arabidopsis seedlings, restricted to proximal and distal regulatory regions (states 2 and 4, respectively). We also found that H3K27me3 (and H3K4me3) is enriched in ORIs located in proximal promoters (state 2) but not in distal regulatory regions (state 4). Thus, it is conceivable that the subset of ORIs potentially associated with bivalent chromatin may have a fast response to developmental cues, as it occurs to genes where they colocalize.

Our genome-wide results have defined the main DNA and chromatin properties associated with different ORI classes in a living adult organism and at two developmental stages during vegetative growth. These properties allowed us to show that ORI activity is compatible with a variety of signatures and demonstrate the existence of various classes of ORIs defined by their strength, DNA features and chromatin landscape. The feasibility to study ORI activity in a developmentally and genetically tractable organism opens new avenues to determine how ORI activity is regulated in response to developmental cues, in association with transcriptional programs, in response to environmental challenges and in a variety of mutant backgrounds.

EXPERIMENTAL PROCEDURES

Plant growth

Arabidopsis thaliana seeds (Col-0 ecotype) were stratified for 48h and grown in Murashige and Skoog (MS) medium supplemented with 1% (w/v) sucrose and 1% (w/v) agar in a 16h:8h light/dark regime at 22°C, for either 4 or 10 days.

Purification of short nascent strands (SNS)

Total genomic DNA and SNS preparations were obtained under RNase-free conditions, by an optimization of the protocol described (Sequeira-Mendes et al., 2009), which has been enhanced to obtain sufficiently clean SNS preparations. Nuclei isolation from 4 or 10 days post-sowing (dps) *Arabidopsis* seedlings was performed prior to genomic DNA extraction as described (Chodavarapu et al., 2010), in order to minimize genomic DNA contamination with cytosolic polyphenols and other secondary metabolites. Twelve grams of whole seedlings were collected, frozen and ground in liquid nitrogen in the presence of 10% PVPP (Sigma). The ground material was resuspended in 10 ml per gram of Honda Buffer Modified for 30 min in a rotary shaker at 4 °C (HBM; 2% (p/v) PVP10 (Sigma), 25 mM Tris-HCl, pH 7.6, 440 mM sucrose (Merck), 10 mM magnesium chloride, 0.1% Triton X-100, 10 mM β -mercaptoethanol). To better release the nuclei, the resuspended material was processed in a dounce homogenizer twice with a loose and a tight pestles and filtered through a double miracloth mesh into corex tubes. The nuclei were centrifuged 10 min at 3000xg and 4 °C. The supernatant was discarded and nuclear pellet was resuspended in 5 ml per gram of Nuclei Isolation Buffer (NIB; 2% (p/v) PVP10 (Sigma), 20 mM Tris-HCl, pH 7.6, 250 mM sucrose (Merck), 5 mM magnesium chloride, 5 mM potassium chloride, 0.1% Triton X-100, 10 mM β -mercaptoethanol). The sample was loaded onto a 15/50% gradient of percoll in NIB and centrifuged 20 min at 500xg and 4 °C with slow brake. The green upper layer was discarded and the same volume was replaced with NIB. Nuclei were centrifuged 5 min at 1100xg and 4 °C, washed twice with 10 ml of NIB and 4 °C, and resuspended in 20 ml of lysis buffer per 12 grams of starting material (0.5% (p/v) PVP10, 50 mM Tris-HCl, pH 8.0, 10 mM EDTA pH 8.0, 1% SDS, 10 mM β -mercaptoethanol) by agitation 15 min at 4 °C. To digest the proteins, 100 μ g/ml proteinase K was added and incubated overnight at 37 °C with mild rotation. Total DNA was extracted twice, first using phenol, pH 8.0, then with phenol:chloroform:IAA and the aqueous phase containing genomic DNA was collected into polyallomer tubes (Beckman). DNA was precipitated by adding 1.5M sodium chloride and 2 volumes of absolute ethanol, incubated 1h at -80 °C and pelleted by centrifugation for 45 min at 52,000xg at 4 °C using an AH-627 rotor (Sorvall). DNA was washed twice with 70% ethanol, centrifuged 20 min at 52,000xg and room temperature in an AH-627 rotor (Sorvall), air dried and resuspended in 1 ml of TE (10 mM Tris-HCl, pH 8.0, 1 mM EDTA) containing 160 U of RNase OUT (Invitrogen). DNA was incubated at 4 °C overnight without pipetting or vortexing.

Purified DNA was denatured by heating 10 min at 100 °C and size-fractionated in a seven-step neutral sucrose gradient (5-20% sucrose in TEN buffer (10 mM Tris-HCl, pH 8.0, 1 mM EDTA and 100 mM NaCl), by centrifugation at 102,000xg in a SW-40T1 Beckman rotor for 20 h at 20 °C (Gomez and Antequera, 2008). Fractions (1 ml) were collected from the top and the DNA was ethanol-precipitated. An aliquot of each fraction was analyzed in a 1% alkaline agarose gel (50 mM NaOH, 1 mM EDTA) to monitor size fractionation. Normally, fractions 3 (~100-600 nt), 4 (~300-800 nt) and 5+6+7 (~500-3000 nt) were processed further by treating with 0.67 U/ μ l of polynucleotide kinase (PNK, Fermentas) to phosphorylate 5'-hydroxyl ends in the presence of 1.34 mM dATP for 30 min at 37 °C. After PNK inactivation, phosphorylated DNA was extracted, precipitated and resuspended in water. SNS were distinguished from randomly broken DNA molecules based on the presence of 4-6 nt-long RNA primers at their 5'ends, which made them resistant to λ -exonuclease treatment (Cayrou et al., 2015; Comoglio et al., 2015; Costas et al., 2011; Gerbi and Bielinsky, 1997). The λ -

exonuclease digestion was carried out with 5 U/μl of enzyme (Thermo Scientific) following the manufacture's instructions at 37 °C overnight. The efficiency of the digestion was monitored by adding 40 ng of phosphorylated linearized plasmid to an aliquot of each reaction tube. DNA from each λ-treated fraction was extracted, precipitated and resuspended in TE. The phosphorylation and λ-exonuclease treatments were repeated at least twice. RNA was digested 0.05 μg/ml RNase A (Roche) and 0.16 U/μl RNase I (Thermo Scientific) for 30 min at 37 °C. RNases were digested with 100 μg/ml proteinase K and DNA was extracted, precipitated and resuspended in miliQ water. The ssDNA of purified SNS was converted into dsDNA: first, SNS together with 2 pmol random hexamer primers (Roche) were denatured 5 min at 100 °C, then a slow annealing was achieved by cooling down the samples from 80 °C to room temperature; second, the dsDNA was synthesized by using 5U of Klenow fragment for 1h at 37 °C; third, the fragments were ligated with 40 U of Taq DNA ligase (New England Biolabs) for 45 min at 45 °C; finally, dsDNA was extracted, precipitated, resuspended in miliQ water and quantified before proceeding to the library preparation. The same method of dsDNA conversion was applied to sheared and denatured genomic DNA to be used as sequencing control.

Next-generation sequencing

DNA libraries of both SNS DNA (sucrose gradient fractions 3, 4 and 5+6+7 combined) and gDNAs were first sheared by an S2 focused-ultrasonicator (Covaris) for 2 minutes (Intensity 5, Duty Cycle 10%, Cycles per Burst 200), and then used as inputs to generate sequencing libraries by Ovation Ultralow V1 library prep kits (NuGen). The libraries were subjected to deep sequencing on HiSeq 2000 per manufacturer instructions (Illumina). In two out of the three experiments we used different amplification protocols for library generation with the purpose of estimating possible bias introduced by this crucial but unavoidable step. RNA-Seq libraries were made by TruSeq Stranded mRNA library prep kit and NeoPrep (Illumina), and subjected to deep sequencing on HiSeq 2000 per manufacturer instructions (Illumina). Single-end sequenced reads (51 nt) were aligned to the reference Arabidopsis genome (TAIR10), using the Bowtie alignment tool (Langmead et al., 2009), allowing up to one mismatch and discarding multihit reads. PCR duplicate reads were removed using an in-house script.

Peak-calling

For each sample and each fraction, we call ORIs with our own peak calling algorithm ZPeaks (Bastolla et al., in preparation) that (i) provides a well defined, genome-wide profile of Nascent Strand Score (NSS), instrumental for weighting candidate ORIs and generic genomic locations, and (ii) localizes an ORI at the local maximum of the NSS over the ORI box called, needed for centering the metaplots. We tested ZPeaks by visual inspection of the overlap between experiment and control reads and candidate ORIs as well as by the statistical analysis of the ORIs properties. Furthermore, our procedure was robust with respect to false positive ORIs because (i) it requires that each ORI is detected in several independent experiments and (ii) it weights each ORI with its NSS, so that spurious ORIs have low NSS and contribute little to the average properties.

Thus, ZPeaks computes optimally smoothed profiles of the reads of the experiment and the control, obtains from them a normalized smoothed profile, calls peaks when the profile is above an user-specified threshold, and sets the ORI location at the maximum of the normalized profile. More in detail, the algorithm works as follows, once the sequencing reads have been aligned to the reference Arabidopsis TAIR10 genome: (1) The wig files (normalized read counts) are input to ZPeaks and the number of reads is rescaled so that its mean number over each chromosome is the same both for the experiment *e* and the control *c*. If the control is not available, a constant profile is used. To increase the reliability of bins where the control is low, values of the control below the mean are interpolated between the

current value and the mean: if $c_i < \langle c \rangle$, then we use $c'_i = (c_i + \langle c \rangle) / 2$, where i indicates the genomic location and $\langle C \rangle$ is the mean value of C over the chromosome where i is located; (2) The profiles of the rescaled experiment and control are smoothed as $c'_i = \sum_k c_k w_{ik}$ where the weights w_{ik} are given by $w_{ik} = \exp(-d_{ik}/d_0) / \sum_l \exp(-d_{il}/d_0)$, d_{ik} is the distance between the center of bin i and bin k , d_0 is a parameter that is optimized as described below. A cut-off on distance is used to accelerate the computation, whose value is optimized alongside d_0 ; (3) From the smoothed experiment and control, the difference score $d_i = e'_i - c'_i$ is constructed and it is transformed into the Z score $z_i = (d_i - \langle d \rangle) / s_d$, where $\langle d \rangle$ is the mean value of d over the chromosome of i and s_d is the standard deviation; (4) For the chosen threshold T , the program counts the number of bins with $z_i > T$, $N(T)$. The smoothing parameter d_0 that yields the largest $N(T)$ for the chosen threshold is chosen as the optimal parameter. The rationale is that, if the profiles are smoothed too much, then the experiment and the control will tend to become equal to their mean values and d_i will tend to be zero, thus decreasing $N(T)$, whereas if the profiles are smoothed too little the standard deviations will be large, also decreasing $N(T)$. We can always determine numerically an optimal parameter d_0 for which $N(T)$ is maximum, which justifies our procedure. We defined the nascent strand score (NSS) profile of the experiment e as $NSS_{ei} = z_i$; (5) We then joined together consecutive bins with $NSS_{ei} > T$ separated by less than 200 nucleotides, obtaining boxes that represent candidate origins; (6) Finally, the putative DNA replication origin is set at the bin where z_i is maximum within the box, and the limits of the box are reduced in such a way that the ORI is at the center and the new box is contained into the original one. It must be kept in mind that the NSS showed a continuous distribution without any sign of saturation, perhaps suggesting that some bias introduced by the amplification step prior to sequencing may have some contribution.

One may expect that the threshold parameter T may be objectively determined by clustering all genomic bins in two clusters through some clustering algorithm such as K-means, Expectation Maximization (that assumes that the scores z_i are distributed according to a Gaussian distribution) or Hidden Markov Models (that also exploits the positional order of the bins along the chromosome). We followed such strategies, but the thresholds that we obtained were low, a sizable fraction of the genome satisfied $z_i > T$, and visual inspection showed that most candidate ORIs were not reliable. Thus, we had no better choice than selecting an arbitrary threshold T and determining *bona fide* ORIs by combining different experiments, as explained below.

Combining peaks into consensus boxes (potential ORIs)

Our strategy consisted of determining a robust set of ORIs detected in at least two independent experiments and two fractions for each experiment and weighting each candidate ORI with the NSS value of each experiment in such a way that the results are little dependent of false positives with low score.

We analyzed two developmental stages (4 and 10 day-old seedlings) and 3 experiments for each stage (exp1, exp2, exp3), obtaining six different samples. For each of them, either two (F3 and F4) or three (F3, F4, F5+6+7) consecutive fractions of the sucrose gradients for size selection of nascent strands were sequenced. Fractions were equivalent between 4 and 10 day-old seedlings, matching F3, F4 and F5+6+7. We called candidate ORIs with a tolerant threshold ($z > 1.8$) and for each sample we selected candidate regions, or boxes, that were identified in at least two fractions of the same gradient. Boxes with size smaller than 200 bp were eliminated, and boxes closer than 200 bp were joined. In this way we obtained six datasets of high quality ORIs, which numbers were: 842 (4d_exp1), 1938 (4d_exp2), 3008 (4d_exp3), 3298 (10d_exp1), 1686 (10d_exp2), 3107 (10d_exp3).

To increase the reliability of candidate ORIs, we selected only those boxes that had been found in at least two out of six independent samples, obtaining a total of 2374 highly reliable candidate ORIs. We matched the boxes with non-vanishing overlap and if an ORI had multiple overlaps, we selected the largest overlap. The center of the combined box was

computed as the weighted average of the location with maximum score present in the associated boxes, weighting more the boxes with high NSS and small size. When we matched different fractions, the fraction F5+6+7, which contains larger nascent strands, was used to confirm boxes but not to locate their center, in order to obtain better resolution. The limits of the combined box were set in such a way that all of the bins are above the threshold in all fractions.

Scoring ORIs in different samples

For each ORI, we obtained their score NSS_{ek} in the six samples, where e labels the experiment, k labels the ORI, and NSS_{ek} is the maximum value of the score over all bins included in the box that contains the ORI. For each sample, we used the corresponding scores as weights, and we obtained the average values and the metaplots of genomic and epigenetic marks as the weighted average over the set of ORIs. We also generated combined scores by averaging the scores of all 4 day-old and all 10 day-old seedling samples.

Detection of developmentally regulated ORIs

Despite most ORI scores were strongly correlated, indicating that ORIs that are strong in one sample are strong also in the others, we identified a reduced number of ORIs whose strength is significantly different from one sample to the other. For this purpose, we only considered the combined scores of the two developmental stages (4 and 10 day-old seedlings), we rescaled them in such a way that the average value over the set of origins was the same for both samples, and for each ORI k . We computed the mean and standard deviation of the rescaled NSS_{ek} over the samples. We considered variable those ORIs in which the ratio between the standard deviation and the mean is larger than a threshold and studied the properties of outliers as a function of the threshold.

Analysis of ORIs in heterochromatin and TE families was carried out as described in (Vergara et al., 2017).

ACCESSION NUMBERS

The accession number for the SNS-seq data obtained in this study is GSE109668

SUPPLEMENTAL INFORMATION

Supplemental Information includes nine figures (S1-S9) and three tables (S1-S3).

AUTHOR CONTRIBUTIONS

The work was conceived by CG, JS-M and UB. JS-M, together with ZV, implemented protocols for purification and analysis of nascent strands, with the help of CC and IA in the initial steps of the work. UB developed the computational and statistical analysis, with the initial help of RM-G. JS-M, ZV, UB and RP generated data and contributed to analysis. JM and JMC analyzed ORIs in heterochromatin. CG and UB wrote the manuscript with the input of all authors.

ACKNOWLEDGMENTS

We thank E. Martinez-Salas and M. Gomez for discussions and comments during this project and the critical reading of the manuscript, members of our laboratories for continuous feedback and V. Mora-Gil for technical assistance. We thank S. Jacobsen, C. Hale, S. Feng and the BSCRC BioSequencing Core for Illumina DNA Sequencing and discussions. This research was supported by grants BFU2012-34821, BFU2013-50098-EXP and BFU2015-

68396-R to C.G. and grant AGL2013-43244-R to J.M.C., as well as by institutional grants from Fundacion Ramon Areces and Banco de Santander to the CBMSO.

COMPETING FINANCIAL INTERESTS

Authors declare no competing financial interests.

REFERENCES

- Aguilera, A., and Garcia-Muse, T. (2013). Causes of genome instability. *Annual review of genetics* 47, 1-32.
- Aladjem, M.I. (2004). The mammalian beta globin origin of DNA replication. *Front Biosci* 9, 2540-2547.
- Arakawa, K., and Tomita, M. (2012). Measures of compositional strand bias related to replication machinery and its applications. *Curr Genomics* 13, 4-15.
- Bernstein, B.E., Mikkelsen, T.S., Xie, X., Kamal, M., Huebert, D.J., Cuff, J., Fry, B., Meissner, A., Wernig, M., Plath, K., *et al.* (2006). A bivalent chromatin structure marks key developmental genes in embryonic stem cells. *Cell* 125, 315-326.
- Besnard, E., Babled, A., Lapasset, L., Milhavet, O., Parrinello, H., Dantec, C., Marin, J.M., and Lemaitre, J.M. (2012). Unraveling cell type-specific and reprogrammable human replication origin signatures associated with G-quadruplex consensus motifs. *Nat Struct Mol Biol* 19, 837-844.
- Bouyer, D., Kramdi, A., Kassam, M., Heese, M., Schnittger, A., Roudier, F., and Colot, V. (2017). DNA methylation dynamics during early plant life. *Genome Biol* 18, 179.
- Castillo Bosch, P., Segura-Bayona, S., Koole, W., van Heteren, J.T., Dewar, J.M., Tijsterman, M., and Knipscheer, P. (2014). FANCI promotes DNA synthesis through G-quadruplex structures. *Embo J* 33, 2521-2533.
- Cayrou, C., Ballester, B., Peiffer, I., Fenouil, R., Coulombe, P., Andrau, J.C., van Helden, J., and Mechali, M. (2015). The chromatin environment shapes DNA replication origin organization and defines origin classes. *Genome Res* 25, 1873-1885.
- Cayrou, C., Coulombe, P., Puy, A., Rialle, S., Kaplan, N., Segal, E., and Mechali, M. (2012). New insights into replication origin characteristics in metazoans. *Cell Cycle* 11, 658-667.
- Cayrou, C., Coulombe, P., Vigneron, A., Stanojcic, S., Ganier, O., Peiffer, I., Rivals, E., Puy, A., Laurent-Chabalier, S., Desprat, R., *et al.* (2011). Genome-scale analysis of metazoan replication origins reveals their organization in specific but flexible sites defined by conserved features. *Genome Res* 21, 1438-1449.
- Chen, Y., and Yang, D. (2012). Sequence, stability, structure of G-quadruplexes and their drug interactions. *Curr Protoc Nucleic Acid Chem. Chapter 17:Unit17.5*.
- Chodavarapu, R.K., Feng, S., Bernatavichute, Y.V., Chen, P.Y., Stroud, H., Yu, Y., Hetzel, J.A., Kuo, F., Kim, J., Cokus, S.J., *et al.* (2010). Relationship between nucleosome positioning and DNA methylation. *Nature* 466, 388-392.
- Comoglio, F., Schlumpf, T., Schmid, V., Rohs, R., Beisel, C., and Paro, R. (2015). High-resolution profiling of Drosophila replication start sites reveals a DNA shape and chromatin signature of metazoan origins. *Cell Rep* 11, 821-834.
- Costas, C., de la Paz Sanchez, M., Stroud, H., Yu, Y., Oliveros, J.C., Feng, S., Benguria, A., Lopez-Vidriero, I., Zhang, X., Solano, R., *et al.* (2011). Genome-wide mapping of Arabidopsis thaliana origins of DNA replication and their associated epigenetic marks. *Nature Struc Mol Biol* 18, 395-400.

Ernst, J., Kheradpour, P., Mikkelsen, T.S., Shores, N., Ward, L.D., Epstein, C.B., Zhang, X., Wang, L., Issner, R., Coyne, M., *et al.* (2011). Mapping and analysis of chromatin state dynamics in nine human cell types. *Nature* 473, 43-49.

Fouk, M.S., Urban, J.M., Casella, C., and Gerbi, S.A. (2015). Characterizing and controlling intrinsic biases of lambda exonuclease in nascent strand sequencing reveals phasing between nucleosomes and G-quadruplex motifs around a subset of human replication origins. *Genome Res* 25, 725-735.

Gerbi, S.A., and Bielinsky, A.K. (1997). Replication initiation point mapping. *Methods* 13, 271-280.

Gomez, M., and Antequera, F. (2008). Overreplication of short DNA regions during S phase in human cells. *Genes Dev* 22, 375-385.

Goren, A., Tabib, A., Hecht, M., and Cedar, H. (2008). DNA replication timing of the human beta-globin domain is controlled by histone modification at the origin. *Genes Dev* 22, 1319-1324.

Gutierrez, C. (2005). Coupling cell proliferation and development in plants. *Nat Cell Biol* 7, 535-541.

Gutierrez, C., Desvoves, B., Vergara, Z., Otero, S., and Sequeira-Mendes, J. (2016). Links of genome replication, transcriptional silencing and chromatin dynamics. *Current opinion in plant biology* 34, 92-99.

Karnani, N., Taylor, C.M., Malhotra, A., and Dutta, A. (2010). Genomic study of replication initiation in human chromosomes reveals the influence of transcription regulation and chromatin structure on origin selection. *Mol Biol Cell* 21, 393-404.

Kharchenko, P.V., Alekseyenko, A.A., Schwartz, Y.B., Minoda, A., Riddle, N.C., Ernst, J., Sabo, P.J., Larschan, E., Gorchakov, A.A., Gu, T., *et al.* (2011). Comprehensive analysis of the chromatin landscape in *Drosophila melanogaster*. *Nature* 471, 480-485.

Kuwabara, A., and Grissem, W. (2014). Arabidopsis Retinoblastoma-related and Polycomb group proteins: cooperation during plant cell differentiation and development. *Journal of experimental botany* 65, 2667-2676.

Langmead, B., Trapnell, C., Pop, M., and Salzberg, S.L. (2009). Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol* 10, R25.

Leonard, A.C., and Mechali, M. (2013). DNA replication origins. *Cold Spring Harbor perspectives in biology* 5, a010116.

Liu, C., Wang, C., Wang, G., Becker, C., Zaidem, M., and Weigel, D. (2016). Genome-wide analysis of chromatin packing in *Arabidopsis thaliana* at single-gene resolution. *Genome Res* 26, 1057-1068.

Lobry, J.R. (1996). Asymmetric substitution patterns in the two DNA strands of bacteria. *Molecular biology and evolution* 13, 660-665.

Lombrana, R., Almeida, R., Revuelta, I., Madeira, S., Herranz, G., Saiz, N., Bastolla, U., and Gomez, M. (2013). High-resolution analysis of DNA synthesis start sites and nucleosome architecture at efficient mammalian replication origins. *EMBO J* 32, 2631-2644.

Lombrana, R., Alvarez, A., Fernandez-Justel, J.M., Almeida, R., Poza-Carrion, C., Gomes, F., Calzada, A., Requena, J.M., and Gomez, M. (2016). Transcriptionally Driven DNA Replication Program of the Human Parasite *Leishmania major*. *Cell Rep* 16, 1774-1786.

Lubelsky, Y., Prinz, J.A., DeNapoli, L., Li, Y., Belsky, J.A., and MacAlpine, D.M. (2014). DNA replication and transcription programs respond to the same chromatin cues. *Genome Res* 24, 1102-1114.

MacAlpine, D.M., Rodriguez, H.K., and Bell, S.P. (2004). Coordination of replication and transcription along a *Drosophila* chromosome. *Genes & Dev* 18, 3094-3105.

Macalpine, H.K., Gordan, R., Powell, S.K., Hartemink, A.J., and Macalpine, D.M. (2010). *Drosophila* ORC localizes to open chromatin and marks sites of cohesin complex loading. *Genome Res* 20, 201-211.

Macaya, R.F., Schultze, P., Smith, F.W., Roe, J.A., and Feigon, J. (1993). Thrombin-binding DNA aptamer forms a unimolecular quadruplex structure in solution. *Proc Natl Acad Sci U S A* 90, 3745-3749.

Mechali, M. (2010). Eukaryotic DNA replication origins: many choices for appropriate answers. *Nature reviews* 11, 728-738.

Mechali, M., Yoshida, K., Coulombe, P., and Pasero, P. (2013). Genetic and epigenetic determinants of DNA replication origins, position and activation. *Curr Opin Genet Dev* 23, 124-131.

Mesner, L.D., Valsakumar, V., Cieslik, M., Pickin, R., Hamlin, J.L., and Bekiranov, S. (2013). Bubble-seq analysis of the human genome reveals distinct chromatin-mediated mechanisms for regulating early- and late-firing origins. *Genome Res* 23, 1774-1788.

Muller, C.A., and Nieduszynski, C.A. (2017). DNA replication timing influences gene expression level. *J Cell Biol* 216, 1907-1914.

Nordman, J., Li, S., Eng, T., Macalpine, D., and Orr-Weaver, T.L. (2011). Developmental control of the DNA replication and transcription programs. *Genome Res* 21, 175-181.

Palumbo, S.L., Ebbinghaus, S.W., and Hurley, L.H. (2009). Formation of a unique end-to-end stacked pair of G-quadruplexes in the hTERT core promoter with implications for inhibition of telomerase by G-quadruplex-interactive ligands. *J Am Chem Soc* 131, 10878-10891.

Picard, F., Cadoret, J.C., Audit, B., Arneodo, A., Alberti, A., Battail, C., Duret, L., and Prioleau, M.N. (2014). The spatiotemporal program of DNA replication is associated with specific combinations of chromatin marks in human cells. *PLoS Genet* 10, e1004282.

Pourkarimi, E., Bellush, J.M., and Whitehouse, I. (2016). Spatiotemporal coupling and decoupling of gene transcription with DNA replication origins during embryogenesis in *C. elegans*. *eLife* 5, e21728.

Rodriguez-Martinez, M., Pinzon, N., Ghommidh, C., Beyne, E., Seitz, H., Cayrou, C., and Mechali, M. (2017). The gastrula transition reorganizes replication-origin selection in *Caenorhabditis elegans*. *Nat Struct Mol Biol* 24, 290-299.

Sacca, B., Lacroix, L., and Mergny, J.L. (2005). The effect of chemical modifications on the thermal stability of different G-quadruplex-forming oligonucleotides. *Nucleic Acids Res* 33, 1182-1192.

- Saha, S., Shan, Y., Mesner, L.D., and Hamlin, J.L. (2004). The promoter of the Chinese hamster ovary dihydrofolate reductase gene regulates the activity of the local origin and helps define its boundaries. *Genes Dev* 18, 397-410.
- Sanchez, M.P., Costas, C., Sequeira-Mendes, J., and Gutierrez, C. (2012). DNA replication control in plants. *Cold Spring Harb. Perspect. Biol.* 4, a010140.
- Sen, D., and Gilbert, W. (1988). Formation of parallel four-stranded complexes by guanine-rich motifs in DNA and its implications for meiosis. *Nature* 334, 364-366.
- Sequeira-Mendes, J., Araguez, I., Peiro, R., Mendez-Giraldez, R., Zhang, X., Jacobsen, S.E., Bastolla, U., and Gutierrez, C. (2014). The Functional Topography of the Arabidopsis Genome Is Organized in a Reduced Number of Linear Motifs of Chromatin States. *Plant Cell* 26, 2351-2366
- Sequeira-Mendes, J., Diaz-Uriarte, R., Apedaile, A., Huntley, D., Brockdorff, N., and Gomez, M. (2009). Transcription initiation activity sets replication origin efficiency in mammalian cells. *PLoS Genet* 5, e1000446.
- Sequeira-Mendes, J., and Gutierrez, C. (2016). Genome architecture: from linear organisation of chromatin to the 3D assembly in the nucleus. *Chromosoma* 125, 455-469.
- Siefert, J.C., Georgescu, C., Wren, J.D., Koren, A., and Sansam, C.L. (2017). DNA replication timing during development anticipates transcriptional programs and parallels enhancer activation. *Genome Res.*
- Stroud, H., Otero, S., Desvoyes, B., Ramirez-Parra, E., Jacobsen, S.E., and Gutierrez, C. (2012). Genome-wide analysis of histone H3.1 and H3.3 variants in Arabidopsis thaliana. *Proc Natl Acad Sci U S A* 109, 5370-5375.
- Todd, A.K., Johnston, M., and Neidle, S. (2005). Highly prevalent putative quadruplex sequence motifs in human DNA. *Nucleic Acids Res* 33, 2901-2907.
- Valton, A.L., Hassan-Zadeh, V., Lema, I., Boggetto, N., Alberti, P., Saintome, C., Riou, J.F., and Prioleau, M.N. (2014). G4 motifs affect origin positioning and efficiency in two vertebrate replicators. *EMBO J* 33, 732-746.
- Vergara, Z., and Gutierrez, C. (2017). Emerging roles of chromatin in the maintenance of genome organization and function in plants. *Genome Biol* 18, 96.
- Vergara, Z., Sequeira-Mendes, J., Morata, J., Peiro, R., Henaff, E., Costas, C., Casacuberta, J.M., and Gutierrez, C. (2017). Retrotransposons are specified as DNA replication origins in the gene-poor regions of Arabidopsis heterochromatin. *Nucleic Acids Res* 45, 8358–8368.
- Wang, C., Liu, C., Roqueiro, D., Grimm, D., Schwab, R., Becker, C., Lanz, C., and Weigel, D. (2015). Genome-wide analysis of local chromatin packing in Arabidopsis thaliana. *Genome Res* 25, 246-256.
- Xia, X. (2012). DNA replication and strand asymmetry in prokaryotic and mitochondrial genomes. *Curr Genomics* 13, 16-27.

SUPPLEMENTAL INFORMATION

Supplemental Figures

Figure S1. Enhanced protocol for purification of nascent strands (NS) from whole developing seedlings.

Figure S2. Scatter plots of the nascent strand score (NSS) measured in different experiments. Each point represents one of the 2374 ORIs identified in this study. Only pairs corresponding to the same day or the same experiment are shown. V1, V2 and V3 refer to each of the three experiments

Figure S3. Features of the local neighborhood of ORIs in whole Arabidopsis seedlings. Metaplots of NSS, CDC6, transcript content (RNA) and nucleosome (nucl.) content of each of the three independent experiments in 4 day-old (top panels) and 10 day-old seedlings (bottom panels), as indicated.

Figure S4. Features of the local neighborhood of ORIs in whole Arabidopsis seedlings. Metaplots of GC, GC skew and GC split of each of the three independent experiments of 4 day-old (top panels) and 10 day-old (bottom panels) seedlings.

Figure S5. Association of ORIs with chromatin states. A. Normalized weight of ORIs belonging to the 9 chromatin states in each of the three independent experiments of 4 day-old and 10 day-old seedlings, as indicated. B. Same as in panel S4A, showing the propensity (instead of the cumulative weight) for ORIs in the 9 chromatin states is depicted for each of the three independent experiments of 4 day-old and 10 day-old seedlings, as indicated.

Figure S6. Relevance of several genomic variables for ORI specification. The weighted averages of the Z scores for several variables (NSS, CDC6, GC content, GC skew and GGN trinucleotide) are shown for each of the three independent experiments in 4 day-old and 10 day-old seedlings, as indicated.

Figure S7. A. Correlation coefficients of the combined NSS of all ORIs identified in 4 day-old and 10 day-old seedlings, in each chromatin state. B. The correlations between the 6 individual NSS for the indicated experimental pairs.

Figure S8. Average Z score of the transcription score with respect to the average transcription score of the entire genome is shown for ORIs in each chromatin state for each of the three independent experiments in 4 day-old (top panel) and 10 day-old seedlings (bottom panel), as indicated.

Figure S9. Weighted average of several genomic and epigenomic properties, using as weights the combined NSS, for the set of all ORIs and the ORIs preferred at 4 or 10 day-old seedlings, as indicated. In all cases we have used the stronger weight. Measurements were transformed to Z-score with respect to the whole genome, which produces positive values when they are higher than a generic genomic region.

Supplemental Tables

Table S1. Genome localization of Arabidopsis ORIs. The Table contains the chromosomal location, individual NSS values of each of the 2374 ORIs identified in each of the three independent experiments carried out in 4 day-old and 10 day-old seedlings, as well as the combined NSS of each developmental stage. Also the chromatin state associated with each ORI, according to the classes defined in Sequeira-Mendes et al. (2014) are also provided.

Table S2. Number and sequence of GGN motif in ORIs.

Table S3. Genome localization of Arabidopsis ORIs preferentially active in 4 day-old and in 10 day-old seedlings with indication of their NSS values and the chromatin state associated.