



Published in final edited form as:

*J Evol Biol.* 2017 June ; 30(6): 1205–1218. doi:10.1111/jeb.13100.

## Genetic Basis for Soma is Present in Undifferentiated Volvocine Green Algae

Zachariah I. Grochau-Wright<sup>\*</sup>, Erik R. Hanschen<sup>\*</sup>, Patrick J. Ferris<sup>\*</sup>, Takashi Hamaji<sup>†</sup>, Hisayoshi Nozaki<sup>†</sup>, Bradley J. S. C. Olson<sup>‡</sup>, and Richard E. Michod<sup>\*</sup>

<sup>\*</sup>Ecology and Evolutionary Biology, University of Arizona, Tucson, AZ, USA

<sup>†</sup>Department of Biological Sciences, Graduate School of Science, University of Tokyo, Tokyo, Japan

<sup>‡</sup>Division of Biology, Kansas State University, Manhattan, KS, USA

### Abstract

Somatic cellular differentiation plays a critical role in the transition from unicellular to multicellular life, but the evolution of its genetic basis remains poorly understood. By definition somatic cells do not reproduce to pass on genes and so constitute an extreme form of altruistic behavior. The volvocine green algae provide an excellent model system to study the evolution of multicellularity and somatic differentiation. In *Volvox carteri*, somatic cell differentiation is controlled by the *regA* gene, which is part of a tandem duplication of genes known as the *reg* cluster. While previous work found the *reg* cluster in divergent *Volvox* species, its origin and distribution in the broader group of volvocine algae has not been known. Here we show that the *reg* cluster is present in many species without somatic cells, and determine that the genetic basis for soma arose before the phenotype at the origin of the family Volvocaceae approximately 200 million years ago. We hypothesize the ancestral function was involved in regulating reproduction in response to stress and that this function was later co-opted to produce soma. Determining that the *reg* cluster was co-opted to control somatic cell development provides insight into how cellular differentiation, and with it greater levels of complexity and individuality, evolves.

### Introduction

The evolutionary transition from unicellular to multicellular life involves an increase in organismal complexity and a shift in individuality from the level of the cell to the level of the multicellular organism. A key step in this transition is the evolution of cellular differentiation, in particular the division of labor between non-reproductive somatic cells and reproductive germ cells (Buss, 1987; Michod, 1999; Grosberg & Strathmann, 2007; Simpson, 2012). The evolution of altruistic somatic cells is a major step in transferring the level of selection, as well as individuality, from the level of the cell to the level of the multicellular organism (Queller, 2000; Michod, 2005; Folse III & Roughgarden, 2013). This transition also represents an increase in complexity as measured by the hierarchical nestedness of the organism (from single-celled to multicellular) and the number of cell types present (Maynard Smith & Szathmáry, 1995; Bell & Mooers, 1997; Marcot & McShea, 2007; Niklas *et al.*, 2014). In organisms with somatic differentiation, somatic cells give up reproductive capacity to specialize on maintenance and survival-related functions of the

multicellular organism, whereas germ cells specialize on reproduction and have reduced viability at the cell level (were they to exist outside the context of the group). Consequently, germ and somatic cells must work together as a team to ensure high fitness for the multicellular organism (Michod, 2005).

The volvocine green algae, along with their unicellular relative *Chlamydomonas reinhardtii*, span a gradient in complexity from single-celled organisms, to undifferentiated multicellular groups, to species with thousands of cells and germ-soma differentiation. They are photosynthetic, facultatively sexual, haploid eukaryotes in the chlorophycean order Volvocales which includes three families: the Tetrabaenaceae, the Goniaceae, and the Volvocaceae (Figure 1). The Tetrabaenaceae contain two genera, *Tetrabaena* and *Basichlamys*, which are made up of four undifferentiated cells held together by extracellular matrix. The Goniaceae contain the genera *Astrephomene*, 32 – 64 celled spheroidal colonies with 2 – 4 sterile somatic cells in the posterior pole of the colony; and *Gonium*, colonies shaped like a slightly curved plate made up of 8 – 32 undifferentiated cells. The Volvocaceae are a diverse family containing genera with undifferentiated 8 – 32 celled colonies (such as *Pandorina*, *Volvulina*, *Platydorina*, *Yamagishiella*, and *Eudorina*), 64 – 128 celled *Pleodorina* with specialized somatic cells in the anterior portion of the colony, and *Volvox* with thousands of cells arranged in a sphere and complete germ-soma cellular differentiation (Figure 1). Many of these genera were originally described morphologically, but molecular phylogenies have revealed many to be polyphyletic, including *Volvox* and *Eudorina* (Figure 1) (Nozaki *et al.*, 2002; Herron & Michod, 2008; Coleman, 2012).

Multicellularity arose relatively recently in the volvocine green algae (~240 million years ago) compared to embryophytes (~748 – 872 million years ago) and metazoa (~574 – 852 million years ago) (Herron *et al.*, 2009; Sharpe *et al.*, 2014). In addition, ancestral character state reconstructions based on molecular phylogenies have inferred that numerous multicellular traits, including the evolution of somatic cells, have been gained multiple times in this group (Figure 5) (Herron and Michod 2008). Thus, these algae provide a uniquely detailed and recent timeline of the stages involved in the evolutionary transition from unicellular to multicellular individuals (Kirk, 2005; Herron & Michod, 2008). In addition, multiple genomes (Merchant *et al.*, 2007; Prochnik *et al.*, 2010; Hanschen *et al.*, 2016) and available molecular techniques (Kirk, 1998; Umen & Olson, 2012), make this group particularly well suited for studying the evolution of the genetic basis for cellular differentiation and individuality (Kirk, 2005; Michod, 2007; Coleman, 2012).

Mutants of the differentiated multicellular alga *Volvox carteri* lacking terminally differentiated somatic cells were first described by Starr (1970). These mutants are known as “regenerator” mutants because colonies first develop seemingly normally, but later the somatic cells develop into reproductive germ cells. The locus found to be responsible for the regenerator phenotype, *regA*, encodes a putative transcription factor (Huskey & Griffin, 1979; Kirk *et al.*, 1999). *regA* is part of a tandem duplication of several paralogous genes known as the *reg* cluster (Figure 2) (Duncan *et al.*, 2007; Hanschen *et al.*, 2014). All *reg* cluster genes encode a DNA-binding SAND domain, also known as the VARL (volvocine algae *regA*-like) domain. During development in *V. carteri* the *regA* gene is turned on in soma-progenitor cells where it is thought to down regulate chloroplast biogenesis. This is

thought to prevent somatic cells from being able to grow large enough for cell division, thus keeping them in a non-reproductive somatic state (Kirk, 2001). Preliminary work has shown that two other *reg* cluster genes, *rlsB* and *rlsC*, are co-expressed with *regA* during development suggesting that the other members of the *reg* cluster also play a role in somatic cell differentiation (Harryman, 2012).

The *reg* cluster is absent from the genomes of *Chlamydomonas reinhardtii* (Duncan *et al.*, 2007; Merchant *et al.*, 2007) and *Gonium pectorale* (Hanschen *et al.*, 2016), but present in divergent *Volvox* species including *V. carteri*, *V. ferrisii*, and *V. gigas* suggesting that the *reg* cluster evolved at the origin of the Volvocaceae (Figures 1 and 3) (Hanschen *et al.*, 2014). The most closely related homolog to the *reg* cluster found in *C. reinhardtii* and *G. pectorale* is *RLS1*, which is orthologous to *rlsD* from *V. carteri* and *V. ferrisii* (Duncan *et al.*, 2007; Hanschen *et al.*, 2014). It is thought that the *reg* cluster arose via the duplication of the *RLS1* gene after the Goniaceae diverged from their last common ancestor with the Volvocaceae (Hanschen *et al.*, 2014, 2016). However, the possibility that the *V. carteri* and *V. ferrisii* lineages gained the *reg* cluster through independent duplication events could not be rejected (Hanschen *et al.*, 2014). The single duplication hypothesis predicts that the *reg* cluster is ancestral to all descendants of the last common ancestor of *V. carteri* and *V. ferrisii* including genera without somatic cells such as *Pandorina*, *Eudorina*, *Yamagishiella*, and *Platydorina*. Alternatively, independent evolutions of the *reg* cluster may correlate with the independent evolutions of somatic cells in *V. ferrisii* and the *V. carteri/V. gigas* group. If so, the *reg* cluster would not be found in undifferentiated genera such as *Pandorina* and *Yamagishiella* (Figure 1).

In order to understand the origin and evolution of *regA*, a key component of the genetic basis for somatic differentiation in *V. carteri*, we investigated for the first time the presence of the *reg* cluster in non-*Volvox* volvocacean algae species representing multiple clades and levels of morphological complexity. We report the discovery of the *reg* cluster in *Pandorina morum*, *Platydorina caudata*, *Yamagishiella unicocca*, *Eudorina elegans* UTEX 1212, and *Pleodorina californica* (Figure 1). The phylogenetic position and morphology of these species indicate that the *reg* cluster arose in an undifferentiated ancestor early in the diversification of the volvocine lineage. Thus, there is a disparity between the known function of the *reg* cluster (*i.e.*, somatic differentiation in *V. carteri*) and the phenotype observed in species with the *reg* cluster but lacking somatic cells, suggesting the cluster originally served a different function and was later co-opted to produce somatic cells.

## Methods

### Cultures and Cosmid Construction

We grew *Pan. morum* (UTEX 1727), *Pla. caudata* (UTEX 1658), *E. elegans* UTEX 1212 (NEIS 721), and a wild isolate of *Ple. californica* (see Supplemental methods for identification details) in standard *Volvox* medium (SVM) at 25°C with a 16:8 hour light:dark cycle (~35  $\mu\text{mol photons/m}^2/\text{s}$ ). Genomic DNA was extracted from cultures using a previously described protocol (Miller & Kirk, 1999). Degenerate PCR was then performed as described in Hanschen *et al.* (2014), using two forward (regF2 and regF3) and two reverse primers (regR and DregR) designed to amplify the VARL domains of *regA* and related *reg*

cluster genes (Supplemental Table 1). All reactions were performed using 2X Phusion HF Master Mix (Thermo Scientific, Waltham, MA) with 3% DMSO and 25 ng of template DNA. Cycling conditions used for reactions were as follows: an initial melting step of 98°C for 3 minutes; followed by 35 cycles of 98°C for 10 seconds, 63–68°C for 30 seconds, and 72°C for 60 seconds; then a final 5 minute extension step at 72°C. For sequences that required extension beyond the degenerate PCR amplified region TAIL-PCR was used based on methods described in Dent *et al.* (2005) using 2× Taq Master Mix (New England BioLabs). Cosmid library construction was performed as in Hanschen *et al.* (2014) using Epicentre's pWEB-TNC Cosmid Cloning Kit. Libraries were screened as described in Hanschen *et al.* (2014) using probes designed from sequences obtained through degenerate PCR (Supplemental Table 1). Cosmids found to contain genes of interest were then isolated for further analysis. See Supplemental Methods for detailed description of cosmid isolation in each species.

### Sequencing and Assembly

DNA sequencing of PCR products, plasmids, and cosmid segments was performed by the University of Arizona Genetics Core service using Applied Biosystems 3730 DNA Analyzers (Waltham, MA). All PCR product sequences were derived from at least two reactions pooled together. Preparation of plasmid and cosmid DNA was performed using the Qiagen Plasmid Midi Kit or Qiagen Spin Miniprep Kit according to manufacturer's instructions. DNA sequences were assembled into contigs using the CLC Main Workbench or Geneious (version 6.1.7) software.

Full cosmids were sequenced using Illumina® sequencing. Illumina libraries were prepared by shearing cosmids with a Covaris® S220 using settings optimized for cosmids then generating libraries with the Illumina TruSeq® HT kit. All size selection steps during library preparation were for an average insert size of ~550 base pairs. Each cosmid library was individually barcoded and paired-end sequenced with a MiSeq instrument using version 3 reagents for 600 cycles. Library sequences were de-multiplexed on the BaseSpace platform, resultant fastq files were then quality trimmed with Sickle, and assembled with Abyss 1.5.2 (Simpson *et al.*, 2009) using a sequential *k*-mer optimization. Following assembly, gaps in cosmid sequences were manually edited to close gaps in Geneious (version 6.1.7). After masking the cosmid vector, genes were annotated with Augustus version 2.7 (Stanke *et al.*, 2008) with settings optimized for *C. reinhardtii*, followed by manual curation.

We used the *regA* VARL domain from *V. carteri f. nagariensis* to search the preliminary genome assembly of a *Y. unicocca plus* strain, and found that scaffold 42 contains the *reg* cluster. This region also contained several gaps. Forward and reverse PCR primers were designed flanking the gaps and used for PCR with genomic DNA, followed by sequencing the PCR product of successful reactions. Three gaps could be filled: one in the intergenic region between *ackB* and *rlsD*, and one near the C-terminal end of both *rlsB* and *rlsO*. Sequencing of genomic PCR products was used to confirm the nucleotide sequence of all *reg* cluster VARL domains. Due to the extensive sequence homology present in the non-VARL portions of *regA*, *rlsB*, and *rlsO*, these sequences were also confirmed by genomic PCR.

## Annotation

Gene predictions for all reported genomic sequences, including *Y. unicocca*, were created using Augustus version 3.1 (Stanke & Morgenstern, 2005). Augustus was selected as its algorithm has been tuned to predict high GC content genomes such as *C. reinhardtii* and *V. carteri*. Annotations were created using partial gene models based on *C. reinhardtii*. Based on previous volvocine genome annotations (Hanschen *et al.*, 2014, 2016), the UTR option was turned off. Models were manually modified to increase interspecific homology when possible. Gene identification was determined by phylogenetic relationship (Figure 3), syntenic position (Figure 2), characteristic protein homologies (Figure 4), and intron positions. Supplemental Table 3 summarizes support for all VARL gene annotations and additional annotation details can be found in the supplemental material.

The presence of SAND domains was predicted using both Pfam (version 28, (Finn *et al.*, 2014)) and SMART (version 7, (Letunic *et al.*, 2012)) databases. Annotations for both Pfam and SMART were obtained using direct submission via Perl script. E-values of SAND domains (PF01342, SM00258) are reported in Supplemental Table 2. We also searched the conserved motifs highlighted in Figure 4 against the PROSITE motif database but found no significant hits.

## Phylogenetic and Protein Similarity Analyses

VARL domain protein sequences were defined following the N-terminal extension and core VARL domain structure of Duncan *et al.* (2007). Sequences were aligned with MAFFT (version 7.213) using the accurate L-INS-i option (Katoh *et al.*, 2005). Maximum likelihood gene phylogenies were generated using RAxML (version 8.1.12) with a rapid bootstrapping analysis with an automatically selected 1,000 bootstrap replicates (Stamatakis, 2014). The LG+G substitution model was automatically determined as the best substitution model by RAxML using the corrected Akaike information criterion. Stabilizing selection was tested for by calculating dN/dS ratios using codeml in PAML (Yang, 2007).

## Species tree construction

We generated a species tree (Figure 1, Figure 5) using Bayesian Markov chain Monte Carlo (MCMC) implemented in MrBayes version 3.2.2 (Ronquist *et al.*, 2012) using default parameters except as described below. Sequences for 80 volvocine operational taxonomic units (OTUs) and six non-volvocine green algae, consisting of five chloroplast genes (ATP synthase beta-subunit, *atpB*; P700 chlorophyll a-apoprotein A1, *psaA*; P700 chlorophyll a-apoprotein A2, *psaB*; photosystem II CP43 apoprotein, *psbC*; and the large subunit of Rubisco, *rbcL*; Supplemental Table 6) were included (excess OTUs were trimmed for Figure 1 and outgroup taxa were trimmed for Figure 5) and a codon partitioning scheme was used following Herron and Michod (2008). Four independent Bayesian analyses of four chains (three heated chains and one cold chain) were run for  $2.5 \times 10^7$  generations with a burn-in of  $5 \times 10^6$  generations. Trees were sampled every 100 generations and assembled to construct a majority rule consensus phylogram. Posterior probabilities for nodes were calculated using all post burn-in trees. We considered the run to have adequately sampled the solution space as the standard deviation of split frequencies was below  $5 \times 10^{-3}$ . Statistical support for this species tree was further estimated using maximum likelihood bootstrapping. jModelTest

version 2.1.10 (Darriba *et al.*, 2012) was used to select GTR+ $\Gamma$ +I as the best-fit model (AICc=135.26). RAxML version 7.2.8 (Stamatakis, 2006) with a GTR+ $\Gamma$ +I model was used to draw 200 topologies onto the Bayesian consensus tree. The resulting species tree (Figure 1, Figure 5) is consistent with recently published phylogenetic analyses (Herron & Michod, 2008; Nozaki *et al.*, 2014). Specifically, there are no statistically supported differences between our topology and other published tree topologies (Herron & Michod, 2008).

### Ancestral character state reconstruction

The presence of obligate somatic cells in each species and strain was compiled from published reports (Supplemental Table 5). Ancestral character states were reconstructed using maximum likelihood and Bayesian methods.

For the maximum likelihood analysis, the consensus species tree was analyzed in R (R Core Team, 2013) using the *diversitree* package (Fitzjohn, 2012). An ultrametric tree was calculated using a penalized likelihood function in the *ape* package (Paradis *et al.*, 2004). Two evolutionary models were evaluated, one constraining the rate of gain and loss of somatic cells to be equal and one with these two parameters unequal. Comparisons of these models using AICc (Akaike, 1974) with a small sample size correction (Burnham & Anderson, 2002), revealed the constrained, equal transition rates model as the best compromise between parameter number and model fit (AICc-equal=49.32, AICc-unequal=50.35, AICc=1.03). Given this low value (AICc < 2), the equal rates model is presented in Figure 5 and the unequal rates model is presented in Supplemental Figure 13. All nodes in the phylogeny were set to the state (presence of somatic cells or not) with the highest probability. A node was determined to be significantly supported if it was at least 7.39 times (if the natural logarithm of the ratio of two likelihoods is greater than 2) more likely than the alternative state (Pagel, 1999).

For the Bayesian analysis, phylogenetic uncertainty was explicitly taken into account by analyzing a sample of trees from the MCMC runs. The systematic subsample included every 1,000<sup>th</sup> tree from the four independent MrBayes runs (Ronquist *et al.*, 2012) for a total of 800 trees. Unicellular outgroup taxa were not pruned from these trees for the Bayesian hypothesis tests. As with the maximum likelihood analysis, both equal and unequal transition rate models were analyzed. A Bayes Factor (BF) was estimated based on twice the difference between the highest harmonic mean log likelihood from five independent MCMC runs for each model. As with the maximum likelihood analysis, the equal rates model is weakly favored (Bayes Factor of -0.3) (Kass & Raftery, 1995). For specific nodes of interest, the presence of somatic cells was tested using hypothesis tests, estimating a Bayes Factor from three independent MCMC runs in which the node in question was constrained to one state or the other. Uniform priors, maximum likelihood priors, and gamma-distributed hyperpriors seeded from a uniform distribution were used. All MCMC runs included 5,500,000 generations with a burn-in period of 500,000 generations.



## Results

We used degenerate PCR followed by cosmid cloning and sequencing to determine the sequence of the *reg* cluster from undifferentiated *Pan. morum*, *Pla. caudata*, and *E. elegans* UTEX 1212, as well as soma-differentiated *Ple. californica*. We also obtained sequences of the *reg* cluster and other VARL genes from a preliminary genome assembly of the undifferentiated *Y. unicocca*. Inferring the presence of the *reg* cluster in these species is based on several results. First, the syntenic structure of the *reg* cluster and nearby markers is conserved among all species studied here except for a translocation of *rlsD* and *ackB* in *V. carteri* (Duncan *et al.*, 2007; Hanschen *et al.*, 2014) and an inversion in *Pla. caudata* (Figure 2). Second, all *reg* cluster VARL domains fall into a single phylogenetic clade, though most relationships show little statistical support (Figure 3). Third, we identified several conserved amino acid motifs outside of the VARL domain that are shared between genes in the same syntenic positions (Figure 4). The most notable of these motifs (“*Pandorina's Box*”) is always found in the first *reg* cluster gene (*rlsA*) and is relatively long (~40 amino acids) and highly conserved (Figure 4). Finally, all *reg* cluster genes contain the same intron within the VARL domain classified as intron 4 by Duncan *et al.* (2007) (Supplemental Table 3).

### Phylogenetic relationships of VARL genes

A maximum likelihood phylogeny (Figure 3) was constructed for all known VARL domains (Duncan *et al.*, 2007; Merchant *et al.*, 2007; Prochnik *et al.*, 2010; Hanschen *et al.*, 2014, 2016) using the core VARL and N-terminal extension boundaries described by Duncan *et al.* (2007). Due to the short sequence length of the VARL domain (86 amino acids) most nodes have very low bootstrap support, although several features are worth noting.

First, all *reg* cluster genes fall within a single clade with *RLS1/rlsD* as an outgroup, consistent with the hypothesis that the *reg* cluster arose through the duplication of *RLS1/rlsD* (Hanschen *et al.*, 2014, 2016). Statistical support is especially low among relationships within the *reg* cluster clade compared to relationships between non-*reg* cluster VARL domains. This may be due to the presence of many more domains within the *reg* cluster clade, which can lower bootstrap support by oversampling (Sanderson & Wojciechowski, 2000). Rebuilding VARL domain trees with sequences from different numbers of species is consistent with this interpretation (data not shown).

Second, several non-*reg* cluster VARL domains form clades with highly supported nodes such as the *rlsE*, *rlsF*, *rlsG*, *rlsH*, and *rlsL* clades (Figure 3). This is consistent with previous analyses comparing the VARL gene content of *C. reinhardtii* and *V. carteri* but with an additional clade; *rlsF*, which has no ortholog from *C. reinhardtii*. Also of note is that the *rlsG* and *rlsE* clades lack an ortholog from *G. pectorale*.

Third, a gene modeled on scaffold 300 of the *Y. unicocca* draft genome (*sc300.g10*) falls within the *reg* cluster clade, sister to the second VARL domain in the *rlsN* gene of *V. ferrisii*, though with poor bootstrap support (Figure 3). This gene contains a VARL domain with a single intron at position 4 as seen in *reg* cluster genes and other VARL genes including *rlsE*, *rlsF*, *rlsG*, *rlsI*, and *rlsM*; but, it does not lie in the *reg* cluster itself (Figure 2).

## Non-VARL domain homology

Previous analyses of the VARL gene family have found that, except for *rlsN* in *V. ferrisii*, VARL genes are characterized as having a single SAND domain with generally low protein homology outside of this region (Duncan *et al.*, 2007; Hanschen *et al.*, 2014). Our Pfam and SMART gene annotations support the presence of only a single SAND domain in the genes reported here (Supplemental Table 2). However, we identified several smaller regions of homology in the predicted protein sequences of *reg* cluster genes and *RLS1//rlsD* that were previously unknown or underappreciated. Almost all *reg* cluster VARL domains are followed by an acidic 5 amino acid sequence (labeled DSGDE in Supplemental Figures 1-5) which is lacking in *RLS1//rlsD* (Duncan *et al.*, 2007). One particularly notable region is found upstream of the VARL domain in *rlsA*, which we've named “*Pandorina's Box*” (Figure 4, Supplemental Figure 1). Other homologous regions of note include a conserved “PRL” motif downstream of the VARL domain in *rlsA*; two conserved motifs found in *regA*, *rlsB*, and *rlsO* (the “LALRP” motif found upstream of the VARL domain, and the “FLQ” motif found downstream of the VARL domain); and an “EQ” motif downstream of the VARL domain in *RLS1//rlsD* (Figure 4, Supplemental Figures 1-6). We searched the Pfam, SMART and PROSITE databases for matches to these conserved regions but no significant matches were found. The three central genes (*regA*, *rlsB*, and *rlsO*) of the *reg* cluster in *Y. unicocca* have a higher degree of protein similarity to each other than to their putative counterparts in other species (Supplemental Figure 7). We are unable to computationally predict functions for these conserved regions.

Intriguingly, the “FLQ” motif we identified lies within a ~1,200 base-pair region known to contain a *cis*-regulatory repressor in the *regA* gene of *V. carteri* (Stark *et al.*, 2001). It is possible then that the conservation of the “FLQ” motif is due to selection to maintain a shared *cis*-regulatory nucleotide sequence rather than selection operating on protein function. However, inspection of nucleotide alignments of the “FLQ” region show saturation of synonymous site mutations (average dS = 1.29, 1.35, and 0.95 for *regA*, *rlsB*, and *rlsO*, respectively) suggesting that purifying selection is acting at the amino acid level rather than the nucleotide level for this region.

## Synteny

We analyzed the synteny of the *reg* cluster from *Pan. morum*, *Ple. californica*, *Pla. caudata*, *Y. unicocca*, and *E. elegans* UTEX 1212 (Figure 2, Supplemental Figure 8). In *V. carteri* the *reg* cluster is bordered by gene model 91984 upstream of *rlsA*, a relationship we observe in *Ple. californica*, *Y. unicocca*, and *Pan. morum* (Figure 2, Supplemental Figure 8). Likewise, the *reg* cluster is bordered by *ackB* followed by *rlsD* in *Pan. morum*, *Y. unicocca*, and *V. ferrisii* (Figure 2) (Hanschen *et al.*, 2014). The close syntenic linkage between *ackB* and *rlsD* is also observed in *V. carteri* and *Pla. caudata* though in *V. carteri* it appears that *rlsD* and *ackB* have been translocated away from the *reg* cluster (Figure 2) (Duncan *et al.*, 2007; Hanschen *et al.*, 2014). The *reg* cluster of *Pla. caudata* has undergone an inversion, as its *reg* cluster is bordered by 91984 and *ackB*, as observed in *Y. unicocca* and *Pan. morum*, but in the reverse order (Figure 2, Supplemental Figure 8). We were unable to clone syntenic markers for *E. elegans* UTEX 1212 but the relative order of *reg* cluster genes is the same as in *Ple. californica*, *V. gigas*, and *V. carteri*.



## **reg cluster nomenclature**

The namesake of the *reg* cluster is the *regA* gene which was named for its role in the regenerator mutant phenotype of *V. carteri* that results in spheroids whose somatic cells regenerate into gonidial cells (Huskey & Griffin, 1979). *Pandorina*, *Platydorina*, *Yamagishiella*, and *Eudorina* do not have somatic cells, however, and thus are incapable of becoming regenerator mutants. Thus, the *regA* gene terminology is somewhat misleading with respect to the effect of a knockout mutation in undifferentiated species. Nevertheless, we choose to use the *Volvox* nomenclature for species with and without soma due to the absence of the *reg* cluster in *C. reinhardtii* and *G. pectorale* and to keep the nomenclature consistent with previous work (Duncan *et al.*, 2007; Hanschen *et al.*, 2014), with one modification. The *reg* clusters of *Pan. morum*, *Y. unicocca* and *Pla. caudata* all contain five VARL domain genes. Hanschen *et al.* (2014) found that *V. ferrisii* also has a five gene *reg* cluster but the fourth gene in the *reg* cluster, *rlsN*, is unique due to having two VARL domains. The fourth *reg* cluster genes from *Pan. morum*, *Pla. caudata*, and *Y. unicocca* only have a single VARL domain and possess the LALRP and FLQ motifs which *rlsN* lacks, suggesting that they are not orthologous to *rlsN* in *V. ferrisii*. For these reasons, we decided to name the fourth gene of the *reg* cluster in *Pan. morum*, *Pla. caudata*, and *Y. unicocca* as *rlsO*. Naming conventions and justifications used for all VARL genes described in the work can be found in Supplemental Table 3.

## **Purifying selection suggests function**

We calculated dN, dS, and dN/dS for interspecific comparisons of *rlsA*, *regA*, *rlsB*, *rlsO*, *rlsC*, and *RLS1/rlsD* VARL domains from *C. reinhardtii*, *G. pectorale*, *Pan. morum*, *Pla. caudata*, *Ple. californica*, *Y. unicocca*, *E. elegans* UTEX 1212, *V. ferrisii*, *V. gigas*, and *V. carteri*. Average dN/dS ratios ranged between 0.16-0.27 for each *reg* cluster gene and *RLS1/rlsD* (Supplemental Figure 9). While relatively high for between species comparisons, these dN/dS values are reflective of the long divergence times between the species compared (Herron *et al.*, 2009) and are consistent with the genome-wide dN/dS values found between *C. reinhardtii*, *G. pectorale*, and *V. carteri* (Hanschen *et al.*, 2016). Thus, our results suggest purifying selection is operating on the VARL domain implying functional conservation. Given their high level of protein similarity, we also calculated dN, dS, and dN/dS values for pairwise comparisons of *regA*, *rlsB*, and *rlsO* coding regions in *Y. unicocca* and found dS values of 0.43-0.74 but dN values of 0.17-0.31, resulting in dN/dS values of 0.40-0.46 (Supplemental Table 4).

## **Ancestral character state reconstruction**

We constructed a Bayesian species tree following methods from Herron & Michod (2008) and inferred ancestral character states using maximum likelihood and Bayesian methods (Figure 5). Both methods provide novel statistical support for inferring that the common ancestor of all species with a *reg* cluster, the ancestor of the *Volvocaceae*, was undifferentiated (Figure 5). Furthermore, both methods statistically significantly infer that the presence of somatic cells in *Volvox* section *Volvox* (i.e., *V. ferrisii*) and the group including *Pleodorina* and the remaining *Volvox* (i.e., *V. carteri*, *V. gigas*, and *Ple. californica*)

represent independent evolutions. In addition, there is a statistically supported third independent evolution of somatic cells in *Astrephomene* (Figure 5).

This analysis was repeated with two models of character evolution, equal (Figure 5) and unequal transition rates (Supplemental Figure 13). The equal transition rate model was found to be marginally statistically preferred ( $AICc = 1.03$ ). Both models of evolution support the results above, but differ in their inference regarding the ancestral state of the lineage of *E. elegans* UTEX 1212. Under the marginally-unfavored unequal rates model, *E. elegans* UTEX 1212 has a soma differentiated ancestor (Supplemental Figure 13), but under the equal rates model all ancestors of *E. elegans* UTEX 1212 lacked cellular differentiation (Figure 5). Yet the *reg* cluster is intact in *E. elegans* UTEX 1212 and presumably serves a different function (Figure 6). If the ancestor did have soma, then *E. elegans* UTEX 1212 would represent a loss of soma indicating that loss of somatic cells may occur without loss of the *reg* cluster.

## Discussion

We show that the *reg* cluster, a key component of the genetic basis for soma in *V. carteri*, is present in many species without somatic cells. All *reg* cluster VARL domains form a single phylogenetic clade (Figure 3), are found in the same relative order in the genome (Figure 2), and most *reg* cluster proteins contain regions of conservation outside of the VARL domain across species (Figure 4). These results and the phylogenetic position of the species examined (Figure 1) indicate that the *reg* cluster arose in the common ancestor of the Volvocaceae following the group's divergence from the Goniaceae. In addition, we performed ancestral character state reconstruction analyses and inferred that the Volvocaceae ancestor lacked terminally differentiated somatic cells, and that soma has evolved at least twice independently in the Volvocaceae (Figure 5, Supplemental Figure 13). Thus, we conclude that the *reg* cluster evolved in an undifferentiated organism and was later co-opted to control somatic differentiation in *V. carteri*.

### Origin of the *reg* cluster

Previous work has shown that the *reg* cluster, a key component of the genetic basis for somatic cell differentiation in *V. carteri*, arose through the duplication of the *RLS1/rlsD* gene, and is present in divergent *Volvox* species that evolved soma independently (Figure 5) (Hanschen *et al.*, 2014, 2016). This evidence indicated that the *reg* cluster likely arose at the origin of the Volvocaceae; however, the possibility that the *reg* cluster arose twice, coinciding with the two independent evolutions of soma in the Volvocaceae, could not be ruled out (Hanschen *et al.*, 2014). The presence of the *reg* cluster in *Pan. morum*, *Pla. caudata*, and *Y. unicocca*, but absence in *G. pectorale* (Hanschen *et al.*, 2016), demonstrates the *reg* cluster arose once, at the origin of the Volvocaceae, following this group's separation from the Goniaceae.

### VARL gene evolution in the volvocine green algae

The *reg* clusters of *Y. unicocca* and *V. ferrisii* indicate that *reg* cluster genes may not be strictly orthologous between species. The fourth gene in *V. ferrisii*'s *reg* cluster, *rlsN*, has a

unique two domain structure which may have arisen via domain duplication (Li, 1997). This gene structure has only been found in one other volvocine species, *Volvox rousseletii*, a close relative of *V. farrisii* (Hanschen *et al.*, 2014). Furthermore, the *reg* cluster of *Y. unicocca* is also unique because the predicted protein sequences of *regA*, *rlsB*, and *rlsO* have high similarity (Supplemental Figure 7, 10 - 12). Pairwise dN and dS analyses of these genes show dS values averaged across the entirety of the alignment are much higher than corresponding dN values (Supplemental Table 4) indicating that gene-birth death with purifying selection is the most likely explanation for the high protein similarity (Nei & Rooney, 2005). Thus, our results indicate that *reg* cluster evolution may be more complicated than previously thought with genes evolving via a birth-death process and at least one instance of domain duplication.

We also investigated the evolution of non-*reg* cluster VARL genes. Previous analyses of VARL genes in the volvocine green algae were limited to the genomes of *C. reinhardtii* and *V. carteri*, restricting their inferential power (Duncan *et al.*, 2007). In this analysis, we included VARL genes from two additional genomes, *G. pectorale* (Hanschen *et al.*, 2016) and *Y. unicocca*. We found that *RLS10/rlsL* and *RLS2/rlsI* form distinct clades containing orthologs from all four genomes (Figure 3, Supplemental Table 3). In contrast, the *RLS11/rlsG* and *RLS8/rlsE* clades lack orthologs from *G. pectorale* which suggests lineage specific gene loss. Similarly, there is no *rlsF* homolog from *C. reinhardtii* suggesting this gene arose following the divergence of the multicellular volvocine green algae from *C. reinhardtii* or gene loss along the *C. reinhardtii* lineage. Finally, the clade formed by *RLS6* from *C. reinhardtii* and *rlsK* from *V. carteri* shows very similar gene content between *V. carteri* and *Y. unicocca* but also many genes unique to *C. reinhardtii* and *G. pectorale* (Figure 3). Taken together, these results suggest little change in non-*reg* cluster VARL gene content within the Volvocaceae but lineage specific evolution in more distantly related groups.

### Evolution of soma in the volvocine green algae

We performed ancestral character state reconstructions of soma to infer when the *reg* cluster arose relative to the evolution of somatic cells. Previous ancestral character state analyses in the volvocine green algae used parsimony and Bayesian methods (Herron and Michod 2008). We built upon this work by incorporating additional taxa, including novel species described since the previous analyses were conducted (Nozaki *et al.*, 2006, 2014; Nozaki & Coleman, 2011; Isaka *et al.*, 2012), and using both maximum likelihood and Bayesian approaches. These new analyses infer, with statistical significance, that the ancestor of the Volvocaceae, and hence the ancestor in which the *reg* cluster arose, was undifferentiated (Figure 5, Supplemental Figure 13).

The ancestral character state analyses were performed under two different models of trait evolution, one where the transition rates between differentiated and undifferentiated were constrained to be equal (Figure 5) and the other where the rates were unequal (Supplemental Figure 13). The equal transition rates model was found to be statistically preferred, but only marginally. Both models infer, with statistical significance, that the ancestor of the Volvocaceae was undifferentiated but differ in how many times somatic cells have been gained or lost. Each model infers six character state changes total, but in the equal transition

rates model (Figure 5) all six changes are independent gains of soma whereas the unequal rates model infers three gains and three losses (Supplemental Figure 13). Additionally, both models infer independent gains of somatic cells in *Astrephomene*; the clade of *Volvox* species containing *V. ferrisii* (often referred to as the “Euvolvox” or section *Volvox*); and at least one gain of soma in the clade containing *V. carteri*, *V. gigas*, and species in the genera *Eudorina* and *Pleodorina* (often referred to as the “*Eudorina* group”) (Figure 5, Supplemental Figure 13). Thus, the *reg* cluster arose at the origin of the Volvocaceae about 200 million years ago (Herron *et al.*, 2009) in an undifferentiated organism prior to the evolution of somatic cells. Somatic cell differentiation appears to be an evolutionarily dynamic trait in the volvocine green algae with multiple gains and potentially multiple losses with the *reg* cluster present in at least two lineages that evolved somatic cells independently.

### Hypotheses for *reg* cluster functions in undifferentiated species

The presence of the *reg* cluster, the genetic basis for soma in *V. carteri*, in organisms lacking somatic cells presents a puzzle. What are these genes doing there? Here we propose three hypotheses for the function of the *reg* cluster in species without somatic cells (Figure 6). It should be noted that the three hypotheses below are not mutually exclusive and different members of the *reg* cluster may have different functions within or between undifferentiated species.

First, all organisms, whether unicellular, undifferentiated multicellular or differentiated multicellular, may be expected to have genes that down regulate growth and reproduction in stressful environments. The nearest homolog of the *reg* cluster in *C. reinhardtii*, *RLS1/rlsD*, is expressed in response to environmental stressors that disfavor photosynthesis and cell growth, such as low light, and sulfur or phosphorus deprivation, when expression of chloroplast proteins is down-regulated (Figure 6A) (Nedelcu & Michod, 2006; Nedelcu, 2009). In the somatic cells of *V. carteri*, *regA* is thought to down-regulate the expression of chloroplast proteins, keeping somatic cells in a starved state and thereby preventing their growth and reproduction (Kirk, 2001). These observations suggest the hypothesis that, following duplication, the ancestral function of *RLS1/rlsD* (down regulating cell growth in response to environmental stress) was co-opted by *regA* to down regulate somatic cell growth in a developmental context (Figure 6E) (Nedelcu & Michod, 2006; Nedelcu, 2009). Therefore, we hypothesize that the function of down-regulating cell growth, and likely reproduction, in response to stress is maintained by the *reg* cluster in undifferentiated species (Figure 6B).

Second, volvocine species differ in their commitment to living in a group. Colonies of *G. pectorale* are known to dissociate under stressful conditions (Graves Jr *et al.*, 1961), but to our knowledge this has never been reported in any members of the Volvocaceae. This may be because *G. pectorale* cells are held together via intercellular wall connections between cells and an outer ECM capsule while the Volvocaceae possess a colony boundary wall (Nozaki, 1990). Given that the *reg* cluster is considered to be a key regulator of cell growth during development in *V. carteri*, the *reg* cluster may regulate cell growth in species that are undifferentiated but committed to group living to prevent premature or spurious cleavage (Figure 6C).

Lastly, the Volvocaceae exhibit a gradient in eyespot and cell size along their anterior-posterior axis (AP axis) (Coleman, 2012). Such a gradient is not seen in the Tetrabaenaceae but *G. pectorale* (Goniaceae) does have a central-to-peripheral axis in flagella basal body rotation (Kirk, 2005; Coleman, 2012). Since the *reg* cluster is thought to control cell growth, we hypothesize that it may be involved in forming the AP axis seen in the Volvocaceae by regulating cell growth with respect to the cell's position within the colony (Figure 6D).

### Implications for evolution of multicellularity and individuality

The evolution of multicellularity in the volvocine algae is thought to have arisen through three major phases: the evolution of cell cycle regulation, the evolution of increased body size, and the evolution of cellular differentiation (Hanschen *et al.*, 2016). The genes underlying the first phase, the evolution of cell cycle regulation, evolved early in the evolution of the volvocine algae, probably to establish a life cycle at the group level (Hanschen *et al.*, 2016). Similarly, our results show that the *reg* cluster, the genetic basis for somatic cells in *V. carteri*, evolved relatively early. There appears to be a common theme of early evolution and co-option of the genetic basis for traits important for multicellularity in the volvocine algae (Hanschen *et al.*, 2014, 2016; Olson & Nedelcu, 2016), though when the genetic basis for the remaining phase, increased body size, evolved is currently unknown.

We have demonstrated that the *reg* cluster is present in undifferentiated species (*Pan. morum*, *Pla. caudata*, and *Y. unicocca*) and arose at the origin of the Volvocaceae (Hanschen *et al.*, 2014) in an ancestor inferred to be undifferentiated (Figure 5, Supplemental Figure 13). We proposed three hypotheses for the function of the *reg* cluster in species without soma including response to environmental stress, adapting to a commitment to group living, and forming the anterior-posterior gradient (Figure 6). Whatever the function of the *reg* cluster is in undifferentiated species, the origin of the *reg* cluster predating the evolution of somatic cells demonstrates that the *reg* cluster must have been secondarily co-opted to control somatic cell development in *V. carteri* (Figure 6). Understanding the function of the *reg* cluster in undifferentiated species and what changes underlie its co-option to control somatic cell development will provide insight into how cellular differentiation, and with it greater levels of complexity and individuality, evolves.

### Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

### Acknowledgments

The authors would like to thank T. Mittelmeier, C. Dieckmann, and D. Shahnooshi for technical help; D. E. Shelton for isolating the novel *Plse. californica* strain; A. Toyoda and A. Fujiyama of the National Institute of Genetics in Mishima, Shizuoka, Japan for access to the *Y. unicocca* genome; as well as D. R. Davison, A. M. Nedelcu, and S. Miller for discussions and comments. This work was supported by the National Aeronautics and Space Administration (grant number NNX13AH41G), the National Science Foundation (grant number MCB-1412395), the National Institute of Health (grant number GM084905), Grants-in-Aid for Scientific Research on Innovative Areas "Genome Science" (grant number 221S0002), and Scientific Research (A) (grant number 24247042) from MEXT/JSPS KAKENHI.

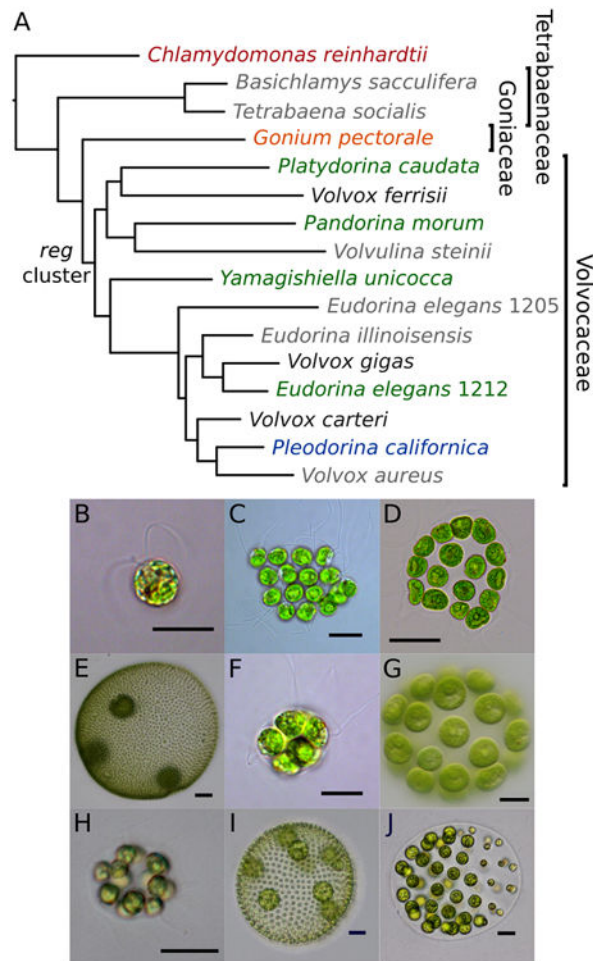
## Literature Cited

- Akaike H. A New Look at the Statistical Model Identification. *IEEE Trans Automat Contr.* 1974; 19:716–723.
- Bell G, Mooers AO. Size and complexity among multicellular organisms. *Biol J Linn Soc.* 1997; 60:345–363.
- Burnham, KP., Anderson, DR. Model selection and multi-model inference: a practical information-theoretic approach. Springer-Verlag; New York: 2002.
- Buss, L. The Evolution of Individuality. Princeton University Press; Princeton, NJ: 1987.
- Coleman AW. A Comparative Analysis of the Volvocaceae (Chlorophyta). *J Phycol.* 2012; 48:491–513. [PubMed: 27011065]
- Darriba, D., Taboada, GL., Doallo, R., Posada, D. *Nat Methods.* Vol. 9. Nature Publishing Group; 2012. jModelTest 2: more models, new heuristics and parallel computing; p. 772-772.
- Dent RM, Haglund CM, Chin BL, Kobayashi MC, Niyogi KK. Functional genomics of eukaryotic photosynthesis using insertional mutagenesis of *Chlamydomonas reinhardtii*. *Plant Physiol.* 2005; 137:545–556. [PubMed: 15653810]
- Duncan L, Nishii I, Harryman A, Buckley S, Howard A, Friedman NR, et al. The VARL gene family and the evolutionary origins of the master cell-type regulatory gene, *regA*, in *Volvox carteri*. *J Mol Evol.* 2007; 65:1–11. [PubMed: 17646893]
- Finn RD, Bateman A, Clements J, Coghill P, Eberhardt RY, Eddy SR, et al. Pfam: the protein families database. *Nucleic Acids Res.* 2014; 42:D222–D230. [PubMed: 24288371]
- Fitzjohn RG. Diversitree: Comparative phylogenetic analyses of diversification in R. *Methods Ecol Evol.* 2012; 3:1084–1092.
- Folse HJ III, Roughgarden J. What is an individual organism? A multilevel selection perspective. *Q Rev Biol.* 2013; 85:447–472.
- Graves LB Jr, Kostir WJ, Lynn B, Jerome W. Some factors affecting the formation of colonies in *Gonium pectorale*. *Ohio J Sci.* 1961; 61:321.
- Grosberg RK, Strathmann RR. The Evolution of Multicellularity: A Minor Major Transition? *Annu Rev Ecol Evol Syst.* 2007; 38:621–654.
- Hanschen ER, Ferris PJ, Michod RE. Early Evolution of the Genetic Basis for Soma in the Volvocaceae. *Evolution (N Y).* 2014; 68:2014–2025.
- Hanschen ER, Marriage TN, Ferris PJ, Hamaji T, Toyoda A, Fujiyama A, et al. The *Gonium pectorale* genome demonstrates co-option of cell cycle regulation during the evolution of multicellularity. *Nat Commun.* 2016; 7:11370. [PubMed: 27102219]
- Harryman, A. Investigating the roles of *regA* and related genes in the evolution of multicellularity in the volvocine green algae. University of Maryland; Baltimore County: 2012.
- Herron MD, Hackett JD, Aylward FO, Michod RE. Triassic origin and early radiation of multicellular volvocine algae. *Proc Natl Acad Sci U S A.* 2009; 106:3254–8. [PubMed: 19223580]
- Herron MD, Michod RE. Evolution of complexity in the volvocine algae: transitions in individuality through Darwin's eye. *Evolution.* 2008; 62:436–51. [PubMed: 18031303]
- Huskey RJ, Griffin BE. Genetic Control of Somatic Cell Differentiation in *Volvox*. *Dev Biol.* 1979; 72:226–235. [PubMed: 510785]
- Isaka N, Kawai-Toyooka H, Matsuzaki R, Nakada T, Nozaki H. Description of two new monoecious species of *volvox* sect. *volvox* (volvocaceae, chlorophyceae), based on comparative morphology and molecular phylogeny of cultured material. *J Phycol.* 2012; 48:759–767. [PubMed: 27011093]
- Kass RE, Raftery AE. Bayes Factors. *J Am Stat Assoc.* 1995; 90:773–795.
- Katoh K, Kuma KI, Toh H, Miyata T. MAFFT version 5: Improvement in accuracy of multiple sequence alignment. *Nucleic Acids Res.* 2005; 33:511–518. [PubMed: 15661851]
- Kirk DL. A twelve-step program for evolving multicellularity and a division of labor. *Bioessays.* 2005; 27:299–310. [PubMed: 15714559]
- Kirk DL. Germ-soma differentiation in *Volvox*. *Dev Biol.* 2001; 238:213–23. [PubMed: 11784005]
- Kirk, DL. *Volvox: Molecular-Genetic Origins of Multicellularity and Cellular Differentiation.* Cambridge University Press; Cambridge, UK: 1998.



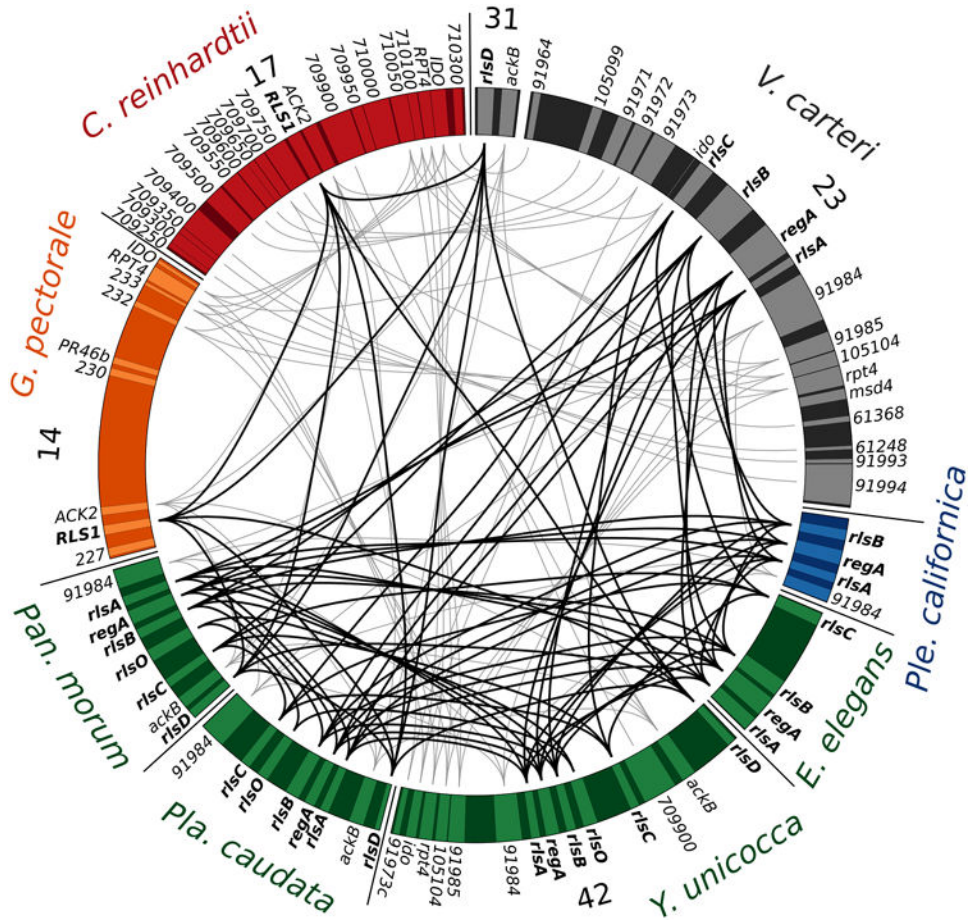
- Kirk MM, Stark K, Miller SM, Müller W, Taillon BE, Gruber H, et al. *regA*, a *Volvox* gene that plays a central role in germ-soma differentiation, encodes a novel regulatory protein. *Development*. 1999; 126:639–47. [PubMed: 9895312]
- Letunic I, Doerks T, Bork P. SMART 7: recent updates to the protein domain annotation resource. *Nucleic Acids Res*. 2012; 40:D302–D305. [PubMed: 22053084]
- Li, WH. *Molecular Evolution*. Sinauer Associates; Sunderland, MA: 1997.
- Marcot JD, Mcshea DW. Increasing hierarchical complexity throughout the history of life: phylogenetic tests of trend mechanisms. *Paleo*. 2007; 33:182–200.
- Maynard, Smith J., Szathmáry, E. *The Major Transitions in Evolution*. Oxford University Press; 1995.
- Merchant SS, Prochnik SE, Vallon O, Harris EH, Karpowicz SJ, Witman GB, et al. The *Chlamydomonas* genome reveals the evolution of key animal and plant functions. *Science*. 2007; 318:245–50. [PubMed: 17932292]
- Michod, RE. *Darwinian Dynamics*. Princeton University Press; Princeton, NJ: 1999.
- Michod RE. Evolution of individuality during the transition from unicellular to multicellular life. *Proc Natl Acad Sci U S A*. 2007; 104:8613–8. [PubMed: 17494748]
- Michod RE. On the transfer of fitness from the cell to the multicellular organism. *Biol Philos*. 2005; 20:967–987.
- Miller SM, Kirk DL. *glsA*, a *Volvox* gene required for asymmetric division and germ cell specification, encodes a chaperone-like protein. *Development*. 1999; 126:649–658. [PubMed: 9895313]
- Nedelcu AM. Environmentally induced responses co-opted for reproductive altruism. *Biol Lett*. 2009; 5:805–8. [PubMed: 19578098]
- Nedelcu AM, Michod RE. The evolutionary origin of an altruistic gene. *Mol Biol Evol*. 2006; 23:1460–4. [PubMed: 16720695]
- Nei M, Rooney AP. Concerted and Birth-and-Death Evolution of Multigene Families. *Annu Rev Genet*. 2005; 39:121–152. [PubMed: 16285855]
- Niklas KJ, Cobb ED, Dunker AK. The number of cell types, information content, and the evolution of complex multicellularity. *Acta Soc Bot Pol*. 2014; 83:337–347.
- Nozaki H. Ultrastructure of the extracellular matrix of *Gonium* (Volvocales, Chlorophyta). *Phycologia*. 1990; 29:1–8.
- Nozaki H, Coleman AW. A New Species of *Volvox* Sect. *Merrillosphaera* (Volvocaceae, Chlorophyceae) From Texas. *J Phycol*. 2011; 47:673–679. [PubMed: 27021996]
- Nozaki H, Ott FD, Coleman AW. Morphology, Molecular Phylogeny and Taxonomy of Two New Species of *Pleodorina* (Volvocaceae, Chlorophyceae). *J Phycol*. 2006; 42:1072–1080.
- Nozaki H, Takahara M, Nakazawa A, Kita Y, Yamada T, Takano H, et al. Evolution of *rbcL* group IA introns and intron open reading frames within the colonial Volvocales (Chlorophyceae). *Mol Phylogenet Evol*. 2002; 23:326–38. [PubMed: 12099791]
- Nozaki H, Yamada TK, Takahashi F, Matsuzaki R, Nakada T. New “missing link” genus of the colonial volvocine green algae gives insights into the evolution of oogamy. *BMC Evol Biol*. 2014; 14:37–47. [PubMed: 24589311]
- Olson, BJSC., Nedelcu, AM. *Curr Opin Genet Dev*. Vol. 39. Elsevier Ltd; 2016. Co-option during the evolution of multicellularity and developmental complexity in the volvocine green algae; p. 107-115.
- Pagel M. The maximum likelihood approach to reconstructing ancestral character states of discrete characters on phylogenies. *Syst Biol*. 1999; 48:612–622.
- Paradis E, Claude J, Strimmer K. APE: Analyses of phylogenetics and evolution in R language. *Bioinformatics*. 2004; 20:289–290. [PubMed: 14734327]
- Prochnik SE, Umen J, Nedelcu AM, Hallmann A, Miller SM, Nishii I, et al. Genomic analysis of organismal complexity in the multicellular green alga *Volvox carteri*. *Science*. 2010; 329:223–6. [PubMed: 20616280]
- Queller DC. Relatedness and the fraternal major transitions. *Philos Trans R Soc Lond B Biol Sci*. 2000; 355:1647–55. [PubMed: 11127911]
- R Core Team. *R: A Language and Environment for Statistical Computing*. Vienna; 2013.

- Ronquist F, Teslenko M, van der Mark P, Ayres DL, Darling A, Hohna S, et al. MrBayes 3.2: Efficient Bayesian Phylogenetic Inference and Model Choice Across a Large Model Space. *Syst Biol.* 2012; 61:539–542. [PubMed: 22357727]
- Sanderson MJ, Wojciechowski MF. Improved Bootstrap Confidence Limits in Large-Scale Phylogenies, with an Example from Neo-Astragalus (Leguminosae). *Syst Biol.* 2000; 49:671–685. [PubMed: 12116433]
- Sharpe, SC., Eme, L., Brown, MW., Roger, AJ. Timing the Origins of Multicellular Eukaryotes Through Phylogenomics and Relaxed Molecular Clock Analyses. In: Nedelcu, AM., Ruiz-Trillo, I., editors. *Evolutionary Transitions to Multicellular Life*. Springer; 2014. p. 3-30.
- Simpson C. The evolutionary history of division of labour. *Proc Biol Sci.* 2012; 279:116–21. [PubMed: 21561969]
- Simpson JT, Wong K, Jackman SD, Schein JE, Jones SJM, Birol I. ABySS: a parallel assembler for short read sequence data. *Genome Res.* 2009; 19:1117–1123. [PubMed: 19251739]
- Stamatakis A. RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses with thousands of taxa and mixed models. *Bioinformatics.* 2006; 22:2688–90. [PubMed: 16928733]
- Stamatakis A. RAxML version 8: A tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics.* 2014; 30:1312–1313. [PubMed: 24451623]
- Stanke M, Diekhans M, Baertsch R, Haussler D. Using native and syntenically mapped cDNA alignments to improve de novo gene finding. *Bioinformatics.* 2008; 24:637–644. [PubMed: 18218656]
- Stanke M, Morgenstern B. AUGUSTUS: a web server for gene prediction in eukaryotes that allows user-defined constraints. *Nucleic Acids Res.* 2005; 33:W465–W467. [PubMed: 15980513]
- Stark K, Kirk DL, Schmitt R. Two enhancers and one silencer located in the introns of *regA* control somatic cell differentiation in *Volvox carteri*. *Genes Dev.* 2001; 15:1449–1460. [PubMed: 11390364]
- Starr RC. Control of differentiation in *Volvox*. *Dev Biol.* 1970; 4:59–100.
- Umen, JG., Olson, BJSC. Genomics of Volvocine Algae. In: Piganeau, G., editor. *Genomic Insights into the Biology of Algae*. Elsevier Ltd: Academic Press; 2012. p. 185-243.
- Yang Z. PAML 4: Phylogenetic Analysis by Maximum Likelihood. *Mol Biol Evol.* 2007; 24:1586–1591. [PubMed: 17483113]



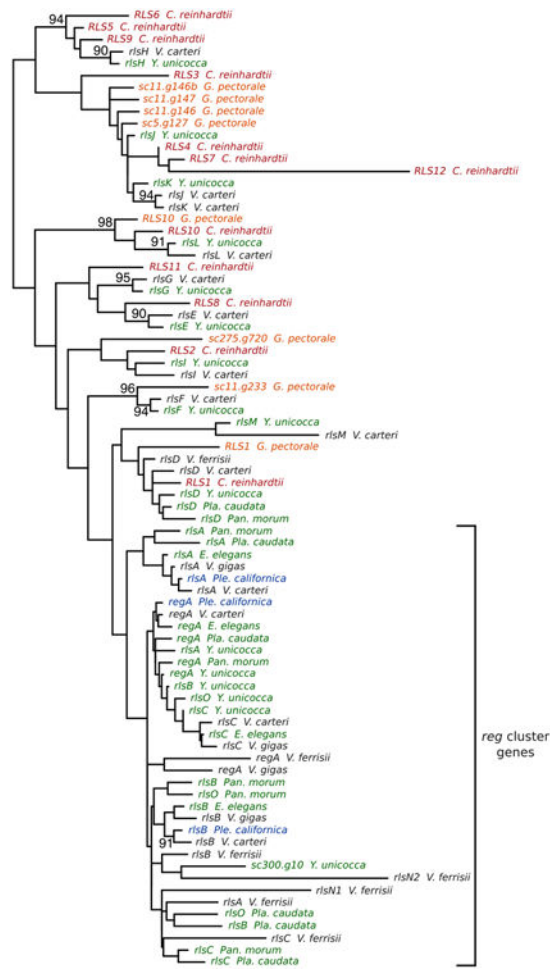
**Figure 1.**

Species phylogeny and micrographs of exemplar species of volvocine algae. A. Bayesian species tree, consistent with previously published species trees. Color of species without (*Chlamydomonas reinhardtii*, red; *Gonium pectorale*, orange) and with the *reg* cluster (undifferentiated *Pandorina morum*, *Platydorina caudata*, *Yamagishiella unicocca*, *Eudorina elegans* UTEX 1212, green; soma differentiated *Pleodorina californica*, blue) correspond to other figures, *Volvox* (germ and soma differentiated) species for which the *reg* cluster has been previously sequenced are shown in black species in grey are not included in this analysis. Note that numbers following *E. elegans* species refer to UTEX strain numbers. Inferred origin of the *reg* cluster is denoted. See Figure 5 for maximum likelihood and Bayesian support values. B. *C. reinhardtii* (scale bar, 10  $\mu\text{m}$ ); C. *G. pectorale* (10  $\mu\text{m}$ ); D. *Pla. caudata* (25 $\mu\text{m}$ ); E. *V. ferrisii* (50  $\mu\text{m}$ ); F. *Pan. morum* (10  $\mu\text{m}$ ); G. *Y. unicocca* (20  $\mu\text{m}$ ); H. *E. elegans* UTEX 1212 (10  $\mu\text{m}$ ); I. *V. carteri* f. *nagariensis* (50  $\mu\text{m}$ ); J. *Ple. californica* (25  $\mu\text{m}$ ).

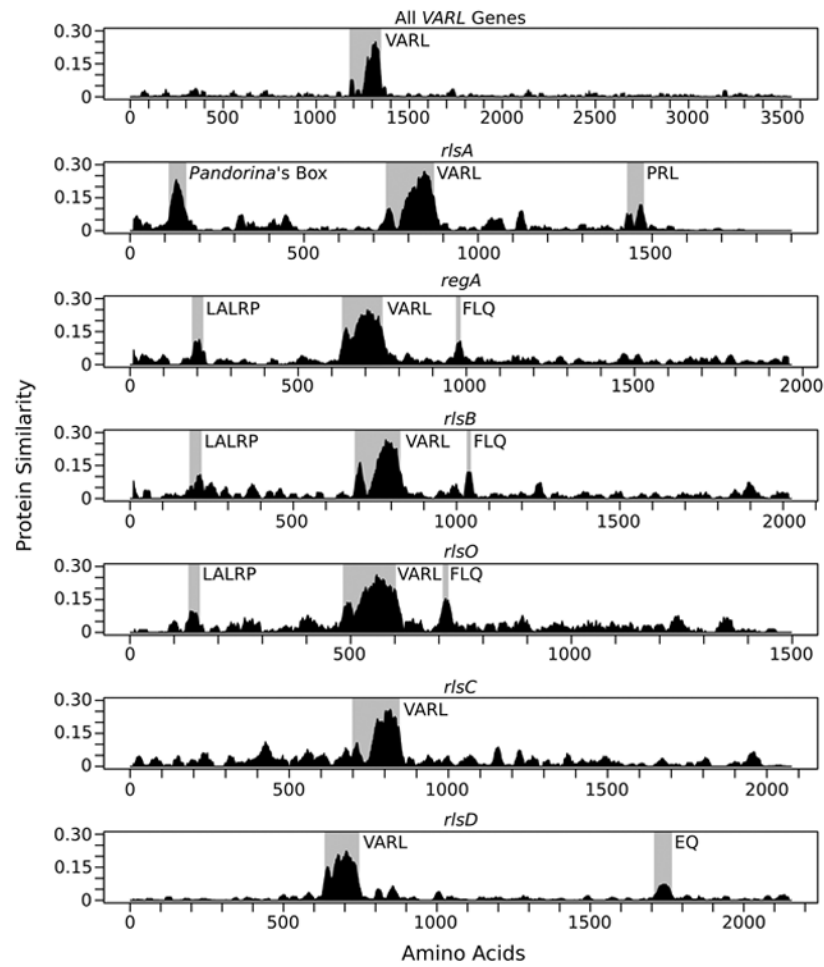


**Figure 2.**

Gene synteny near the *reg* cluster and closely related *regA*-like genes (bold). Synteny of *C. reinhardtii* (red), *G. pectorale* (orange), *Pan. morum* (green), *Pla. caudata* (green), *Y. unicocca* (green), *E. elegans* UTEX 1212 (green), *Ple. californica* (blue), and *V. carteri* (black) is shown. All available data from *Pan. morum*, *Pla. caudata*, *E. elegans* UTEX 1212, and *Ple. californica* are shown, while representative genomic regions from *C. reinhardtii*, *G. pectorale*, *Y. unicocca*, and *V. carteri* are shown. Scaffold or chromosome numbers are indicated for *C. reinhardtii*, *G. pectorale*, *Y. unicocca*, and *V. carteri*. Putative *reg* cluster and *RLS1/rlsD* orthologs are connected with black lines, syntenic genes are connected with gray lines.

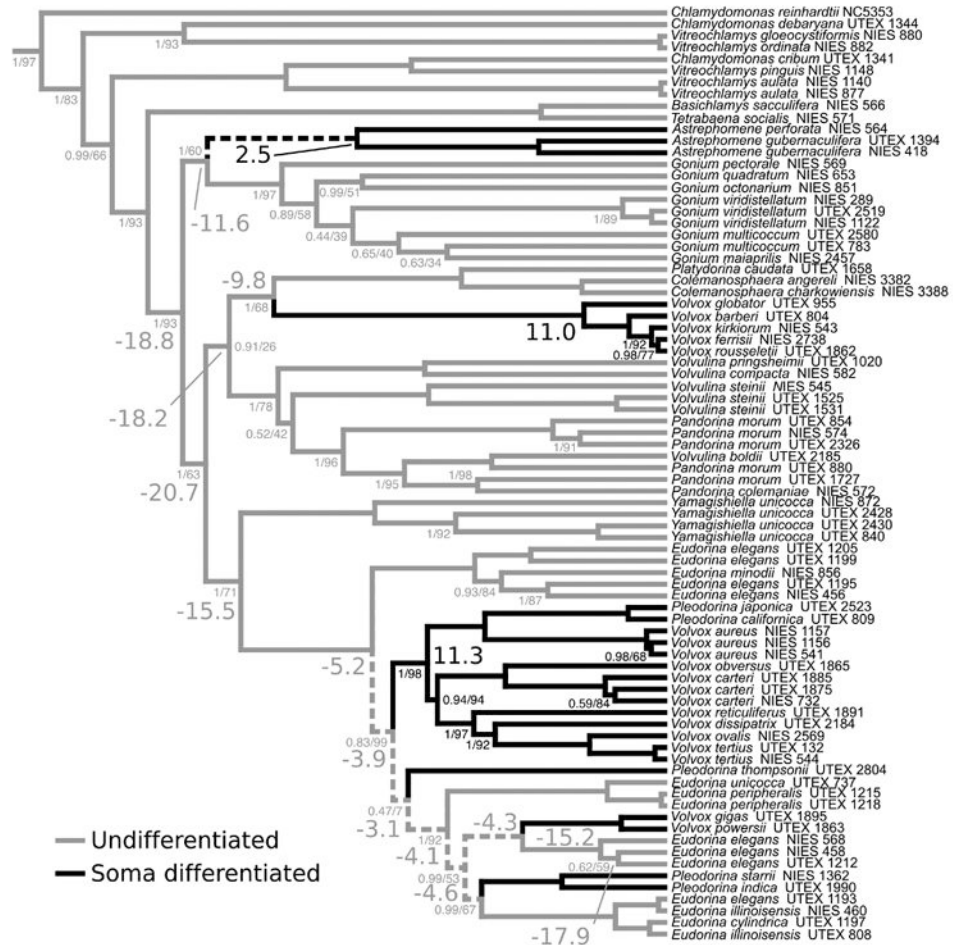


**Figure 3.** Maximum likelihood VARL domain tree. Color of species without (*C. reinhardtii*, red; *G. pectorale*, orange) and with the *reg* cluster (*Pan. morum*, *Pla. caudata*, *Y. unicocca*, and *E. elegans* UTEX 1212, green; *Ple. californica*, blue; *V. carteri*, *V. ferrisii*, and *V. gigas*, black) correspond to other figures. Nodes with 80% or higher bootstrap support values are labeled with support values; unlabeled nodes have less than 80% bootstrap support.

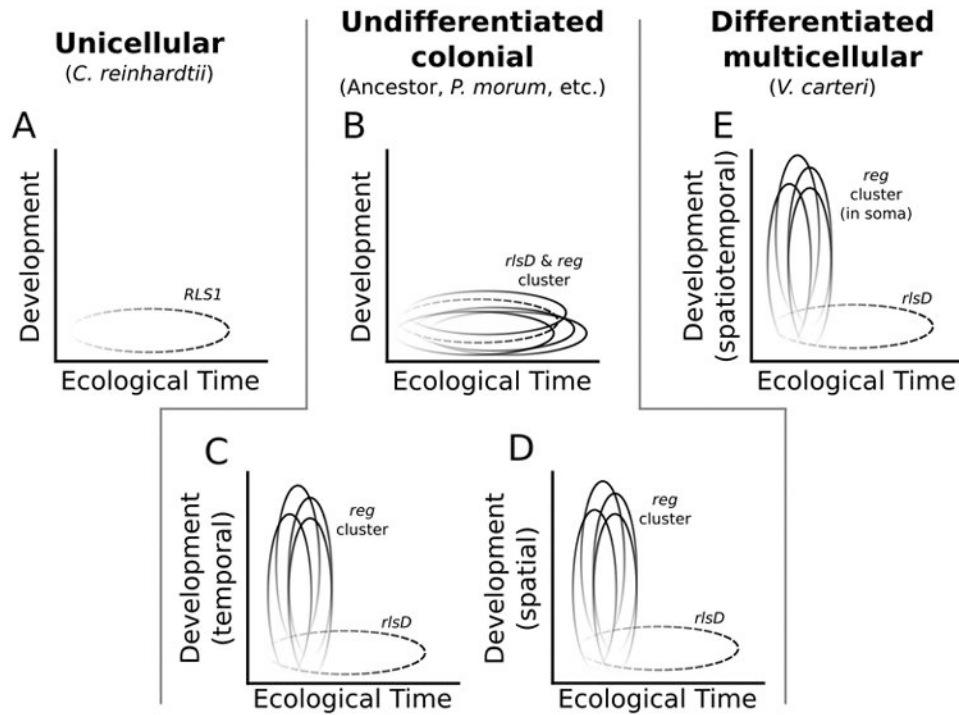


**Figure 4.** Protein similarity plots for all *reg* cluster and *RLS1/rlsD* proteins based on syntenic position (*rlsA*, *regA*, *rlsB*, *rlsO*, *rlsC*, and *rlsD*). Regions showing high similarity are highlighted with gray boxes. The two peaks in the shaded VARL region represent the N-terminal extension and core VARL domain separated by the less conserved linker region (Duncan *et al.*, 2007).





**Figure 5.** Ancestral character state reconstruction of somatic differentiation. Branch color refers to undifferentiated (gray) or somatic differentiation (black) inferred by maximum likelihood methods using the equal transition rates model. Dashed branches indicate an ambiguous maximum likelihood reconstruction. Large font numbers at selected nodes indicate Bayes Factors using the equal rates model; negative, support for undifferentiated; positive, support for somatic differentiated. Bayes Factors are interpreted following Kass and Raftery (1995): 0-2 weak evidence, 2-6 positive evidence, 6-10 strong evidence, >10 very strong evidence. Small font numbers along branches indicate Bayesian posterior probabilities (left of slash) and maximum likelihood bootstrap values (right of slash). Unlabeled nodes are supported with 1.00 PP and 100% bootstrap values. Bayes Factors and support values are colored consistent with the reconstructed state at that node.



**Figure 6.**

Conceptual schematic of hypothesized VARL gene expression patterns in unicellular, undifferentiated, and differentiated species. A. Temporal expression of *RLS1* (dashed) in *C. reinhardtii* in response to environmental change. B. The *reg* cluster (shown as four genes, though some species have five *reg* cluster genes) maintains expression in response to environmental change following its origin from the duplication of *rlsD* (dashed). C. The *reg* cluster is developmentally co-opted to control cell division throughout the life cycle. D. The *reg* cluster is developmentally co-opted to control the AP axis. Panels B, C, and D are not mutually exclusive as different hypotheses may apply to different undifferentiated species and a single hypothesis may not uniformly apply to all genes within the *reg* cluster of a given species. E. The *reg* cluster is co-opted to regulate somatic differentiation in *V. carteri*. For all panels, darker shading represents higher gene expression.