UNIVERSITY *of* York

This is a repository copy of *A recombineering pipeline to clone large and complex genes in Chlamydomonas*.

White Rose Research Online URL for this paper:
https://eprints.whiterose.ac.uk/169738/

Version: Accepted Version

**Article:**

White Rose
university consortium
Universities of Leeds, Sheffield & York

eprints@whiterose.ac.uk
https://eprints.whiterose.ac.uk/

1    **LARGE-SCALE BIOLOGY ARTICLE**
2
3    **A Recombineering Pipeline to Clone Large and Complex Genes in Chlamydomonas**
4
5    Tom Z. Emrich-Mills[a,b,1], Gary Yates[a,1], James Barrett[a], Philipp Girr[a], Irina Grouneva[a], Chun Sing Lau[a],
6    Charlotte E Walker[a], Tsz Kam Kwok[a], John W. Davey[a], Matthew P. Johnson[b], Luke C.M. Mackinder[a,2]
7
8    [a]University of York, Department of Biology, York YO10 5DD, UK
9    [b]University of Sheffield Department Molecular Biology and Biotechnology, Sheffield S10 2TN, UK
10   [1]These authors contributed equally to this work
11   [2]Address correspondence to luke.mackinder@york.ac.uk
12
13   **Short title:** A Recombineering Pipeline for Chlamydomonas
14
15   **One-sentence summary:** We have developed a high-throughput, gene size and gene complexity
16   independent recombineering pipeline in *Chlamydomonas reinhardtii* and applied it to clone 157 $CO_2$
17   concentrating mechanism genes.
18
22
23
24   **Abstract**
25
26   The ability to clone genes has driven fundamental advances in cell and molecular biology, enabling
27   researchers to introduce precise mutations, generate fluorescent protein fusions for localization and to
28   confirm genetic causation by mutant complementation. Most gene cloning is PCR or DNA synthesis
29   dependent, which can become costly and technically challenging as genes increase in size and particularly
30   if they contain complex regions. This has been a long-standing challenge for the *Chlamydomonas*
31   *reinhardtii* research community, with a high percentage of genes containing complex sequence structures,
32   an average genomic GC content of 64% and gene expression requiring regular introns for stable
33   transcription. Here we overcome these challenges via the development of a recombineering pipeline that
34   enables the rapid parallel cloning of genes from a Chlamydomonas BAC collection. We show the method
35   can successfully retrieve large and complex genes that PCR-based methods have previously failed to
36   clone, including genes as large as 23 kilobases, thus making previously technically challenging genes to
37   study now amenable to cloning. We applied the pipeline at both batch and high-throughput scales to 203
38   genes relating to the Chlamydomonas $CO_2$ concentrating mechanism (CCM) with an overall cloning
39   success rate of 77% that is independent of gene size. Localization of a subset of CCM targets has confirmed
40   previous mass spectrometry data and identified new pyrenoid components. To expand the functionality of
41   our system, we developed a series of localization vectors that enable complementation of mutants (e.g.
42   Chlamydomonas Library Project and CRISPR/Cas generated mutants) and enable protein tagging with a
43   range of fluorophores. Vectors and detailed protocols are available to facilitate the easy adoption of this
44   method by the Chlamydomonas research community and to enable the development of recombineering
45   pipelines in other algal and plant species. We envision that this technology will open up new possibilities in
46   algal and plant research and be complementary to the Chlamydomonas mutant library.
47
48   Keywords: *Chlamydomonas reinhardtii*, CCM, $CO_2$ concentrating mechanism, cloning, recombineering,
49   recombination-mediated cloning, genetic engineering, photosynthesis.
50

**Introduction**

The unicellular alga *Chlamydomonas reinhardtii* (hereafter Chlamydomonas) is a widely used model organism for studying photosynthesis, biofuel production, ciliopathies, flagella-powered motility and cell cycle control (Salomé and Merchant, 2019). Its nuclear, chloroplast and mitochondrial genomes are sequenced, well annotated and transformable, and a variety of genetic resources are available to any institution including a close-to-genome-saturating mutant library (Li et al., 2019), extensive -omics based data and a wealth of molecular tools developed over decades by a dedicated research community (Salomé and Merchant, 2019). These collections, data and tools are a vital resource for studies that aim to understand fundamental biological processes, to guide engineering efforts such as improved photosynthetic efficiency and to enable efficient biomolecule production.

Reverse genetic approaches in Chlamydomonas often depend on localizing target proteins to understand spatial distribution and the complementation of mutants to link genotype to phenotype. Both of these methods generally rely on cloning a gene of interest into a plasmid from genomic DNA (gDNA) by PCR, followed by amplification in *Escherichia coli* and reintroduction to Chlamydomonas cells. PCR-based cloning from gDNA presents its own challenges and limitations that are particularly problematic when working with Chlamydomonas nuclear genes, which generally have a high GC content (68% in coding regions), contain one or more introns and can include complex repeating regions (Merchant et al., 2007). On the other hand, cloning from complementary DNA can result in low or no expression of target genes most likely due to lack of introns and lack of regulatory elements (Lumbreras et al., 1998; Schroda, 2019). Some of the challenges associated with PCR-based cloning can be circumvented via whole or partial gene synthesis followed by re-assembly using cloning strategies such as Golden Gate. Although the falling costs of gene synthesis make this a viable option for some genes, for many others the need to include introns, high GC content and high gene complexity, typical of the Chlamydomonas nuclear genome, results in synthesis failure or is prohibitively expensive. For example, SAGA1 (StArch Granules Abnormal 1), a 16.7 kilo base pair (kbp) gene target, required over 12 months of work, included multiple gene synthesis failures and ultimately had to be assembled from three synthesised fragments with 14 introns removed due to repetitive regions (Itakura et al., 2019).

Improved Chlamydomonas target gene and foreign gene (collectively transgenes) expression (e.g., GFP) has been achieved through strain optimization (Neupert et al., 2009), the development of systems with linked transgene and antibiotic resistance gene expression (Rasala et al., 2012; Onishi and Pringle, 2016) and an advanced understanding of transgene silencing (reviewed in Schroda, 2019). Furthermore, release of the Chlamydomonas Golden Gate based Modular Cloning kit has provided a cloning framework and selection of genetic elements to enable labs to rapidly assemble and test transgene constructs (Crozet et al., 2018). Independent of background strain and expression system, it is now clear that inserting or maintaining introns, correct codon usage and promoter sequence are all critical for robust transgene expression (Barahimipour et al., 2015; López-Paz et al., 2017; Baier et al., 2018; Weiner et al., 2018; Schroda, 2019). These considerations have made the cloning of Chlamydomonas target genes directly from gDNA the community standard for mutant complementation and fluorescent protein tagging. However, there are considerable technical hurdles to overcome when working with the expression of large Chlamydomonas genes, predominantly caused by inefficient amplification of gDNA due to gene size, GC content and complexity of target genes (Sahdev et al., 2007). Though modern polymerases have been engineered to overcome sequence challenges (Hommelsheim et al., 2014) they may still suffer from replication slippage events, which are exacerbated by repetitive regions (Levinson and Gutman, 1987; Clarke et al., 2001). In addition to considerations of size and complexity, cloning native genes based on current genome annotations can be complicated by the abundance of upstream transcription start sites corresponding to possible alternative open reading frames (Cross, 2015) and hence potentially resulting in incorrect target gene cloning.

The results of a recent high-throughput localization study illustrate the challenges of PCR-based cloning of Chlamydomonas nuclear genes (Mackinder et al., 2017). In Mackinder et al. (2017) genes were PCR

103 amplified from start site to stop site using gDNA as the template. Amplicons were then cloned in-frame via
104 Gibson assembly with a fluorescent protein and a constitutive promoter and terminator, resulting in the
105 successful cloning of 298 genes out of an attempted 624 (48% success rate), with most failures at the PCR
106 amplification step. This relatively low success rate led us to develop a cloning platform based on
107 recombination-mediated genetic engineering (recombineering) to enable size and sequence independent
108 cloning of Chlamydomonas genes. Recombineering enables gene cloning by homologous recombination
109 in *E. coli* without PCR amplification of the template, and so is predominantly independent of the target
110 region size. Large-scale recombineering pipelines have been developed for bacterial artificial chromosome
111 (BAC) and fosmid libraries from a broad range of organisms including *Caenorhabditis elegans* (Sarov et
112 al., 2006), *Drosophila melanogaster* (Sarov et al., 2016), human and mice (Poser et al., 2008) and
113 *Arabidopsis thaliana* (Brumos et al., 2020) but are lacking in algae. Our developed pipeline involves making
114 BAC-containing *E. coli* homologous recombination competent by introducing the recombinogenic viral
115 proteins Red α, β and γ from the bacteriophage lambda virus (Yu et al., 2000; Copeland et al., 2001), then
116 retrieving a target sequence via introduction of 50 bp homology regions flanking a linearized plasmid.
117
118 We decided to apply our recombineering pipeline to an extended list of putative $CO_2$ concentrating
119 mechanism (CCM) genes. The CCM functions to enhance photosynthesis by increasing the concentration
120 of $CO_2$ around Rubisco. To achieve this Chlamydomonas actively accumulates inorganic carbon in the
121 chloroplast and delivers it as $CO_2$ to tightly packed Rubisco within the pyrenoid (Wang et al., 2015). The
122 pyrenoid is essential for CCM function in Chlamydomonas (Meyer et al., 2012; Mackinder et al., 2016) and
123 due to the photosynthetic turbocharging properties of pyrenoid based CCMs there is growing interest in
124 engineering them into crop plants to boost yields (Mackinder, 2017; Rae et al., 2017). Recent studies have
125 identified a large number of potential pyrenoid and CCM components (Mackinder et al., 2017; Zhan et al.,
126 2018) that require functional characterization to understand their priority for future synthetic CCM
127 engineering efforts. However, many of these are proving challenging to clone due to size and sequence
128 complexity, making localization and mutant complementation studies difficult.
129
130 By applying our pipeline, we have successfully cloned 157 CCM related genes with their native promoters.
131 Cloning appears independent of target gene size and many target genes had multiple complex features
132 that would typically result in PCR failure. The average cloned region was 7.3 kbp and target regions up to
133 22.7 kbp in size were successfully cloned. The inclusion of the native promoters ensures any upstream
134 open reading frames have been incorporated. The localization of a subset of the proteins encoded by these
135 genes has enabled identification of diverse cellular locations, confirming interaction data (Mackinder et al.,
136 2017) and pyrenoid proteomic data (Mackinder et al., 2016; Zhan et al., 2018). We go on to develop a
137 series of recombineering vectors to enable protein tagging with a range of fluorescent proteins and selection
138 markers for localization, complementation and relative protein abundance studies. The method takes four
139 days to implement, is accessible for any lab equipped for molecular biology and requires no specialized
140 reagents or equipment. The BAC library used in this work and all developed plasmids are available from
141 the Chlamydomonas Resource Center and a detailed protocol is provided to enable the rapid adoption of
142 this method by research labs to clone nuclear Chlamydomonas genes.
143
144 **Results**
145
146 **Analysis of the Chlamydomonas genome highlights the challenges affecting PCR-based cloning**
147 Cloning Chlamydomonas genes for successful localization and complementation often requires the
148 amplification of complete open reading frames from gDNA, spanning from their start site to their stop site
149 including any introns (ATG-Stop). To gain a better understanding of the challenges involved in cloning
150 Chlamydomonas genes we performed a whole genome analysis of gene size, complexity, intron
151 prevalence, splice variants, and ATG-Stop primer suitability, including comparisons to available datasets
152 and other organisms.
153

154    *Gene size* - A major limitation of PCR-based cloning is the target amplicon size. ATG-Stop cloning data
155    from Mackinder et al. (2017) for 624 genes using gDNA as a template and Phusion Hot Start II DNA
156    polymerase (ThermoFisher Scientific) shows an association between cloning success and gene size; the
157    average cloned ATG-Stop region was ~2.3 kbp long while the average uncloned region was ~4.5 kbp
158    (Mann-Whitney $U$ = 16306, $P$ < 0.001, two-tailed). Extrapolation of PCR efficiency relative to target size
159    from Mackinder et al. (2017) to whole genes in the Chlamydomonas genome (version 5.5) indicates that
160    68% of genes would be technically challenging to clone via PCR-based methods (Figure 1A), predominantly
161    due to a severe drop off in amplification efficiency for genes >3 kbp long. The largest amplified target in
162    Mackinder et al. (2017) was 8 kbp, and genes at least as large as 9.7 kbp have been cloned before
163    (Kobayashi et al., 2015), but this appears to be highly gene specific. Alternative approaches exist to clone
164    larger genes, such as testing a broad range of PCR conditions and DNA polymerases, amplification in
165    fragments and re-stitching together, cloning from cDNA, and gene synthesis. While some of these
166    approaches avoid the challenges presented here, they can be time consuming, costly, have low success
167    rates and may still result in no or poor expression.

168

169    *Gene complexity* - High GC content and the presence of numerous repetitive regions can make PCR-based
170    cloning challenging. Data from Mackinder et al. (2017) shows that the average GC content for successfully
171    cloned targets by ATG-Stop PCR cloning was 61.4%, while the average for unsuccessful targets was 64.3%
172    – a value exceeded by over 41% of Chlamydomonas nuclear genes. To analyse the genome for repetitive
173    regions, we determined the frequency of simple tandem repeats, inverted repeats, and larger, interspersed
174    repeats between the start of the 5'UTR and the end of the 3'UTR of each gene. Tandem repeats were
175    assessed by counting individual regions that consist of consecutive mono-, di- or trinucleotide repeats.
176    Mononucleotide repeats shorter than 10 bp and regions of di- and trinucleotide repeats shorter than 20 bp
177    were excluded. Some slight imperfections in the repeating pattern of a region were allowed, with regions
178    that showed ≥90% identity included such as GGGGGTGGGG. Of the 17,741 coding genes in the nuclear
179    genome 8,810 contain one or more mono-, di- or trinucleotide repeats (Figure 1B). In terms of prevalence
180    per kilobase, the average Chlamydomonas gene contains 0.21 tandem repeats whereas Arabidopsis
181    contains 0.16 and *Saccharomyces cerevisiae* contains 0.10. Interestingly, if polynucleotide repeats with
182    higher period numbers are counted as well (from tetranucleotide repeats to tandem repeating units of
183    hundreds of base pairs), these values increase 5 fold for Chlamydomonas (1.07 per kbp), 2.5 fold for
184    Arabidopsis (0.39 per kbp) and 3 fold for yeast (0.3 per kbp), highlighting the repetitive nature of the
185    Chlamydomonas genome. Inverted repeats were assessed by counting regions over 10 bp long that are
186    followed closely downstream by their reverse complement, with some mismatches allowed so that regions
187    with ≥90% identity were included. 14,454 genes contain one or more inverted repeats of this kind (Figure
188    1B), with an average of 0.93 repeats per kbp. To further validate these findings we analysed nuclear gene
189    sequences for repeats using WindowMasker, a program for detecting global repeats that include larger
190    non-adjacent sequences as well as a diverse range of tandem repeats and inverted repeats (Morgulis et
191    al., 2006). With this expanded detection range, Chlamydomonas genes contain an average of 38.9 repeats
192    (6.8 per kbp) whereas Arabidopsis contains 13.7 (5.5 per kbp) and yeast contains 6.0 (4.2 per kbp). On
193    average, Chlamydomonas genes are more repetitive between their start and stop codons than in their
194    untranslated regions (Figure 1B), although at least one repeat was detected by WindowMasker in 36.6%
195    of 5'UTRs and 87.6% of 3'UTRs. Crucially, analysis of sequence data from Mackinder et al. (2017) for 624
196    Chlamydomonas genes indicates an association between ATG-Stop PCR cloning success and repeat
197    frequency; the average cloned ATG-Stop region contained 6.1 repeats per kbp whereas the average
198    uncloned region contained 7.5 repeats per kbp (Mann-Whitney $U$ = 24110, $P$ < 0.001, two-tailed).

199

200    *Mis-annotation of start sites* - Another challenge associated with PCR-based and gene synthesis-based
201    cloning is incorrectly annotated gene models that lead to cloning of a non-biologically relevant sequence.
202    The analysis of transcript models in the Chlamydomonas genome shows that additional ATGs upstream of
203    the annotated start site are highly prevalent (Cross, 2015; Figure 1C top 4 bars). Cross (2015) categorized
204    these potential upstream open reading frames (uORFs) into three classes: class 1 uORFs initiate in-frame
205    with the annotated start site, potentially producing an N-terminal extension relative to the annotated gene

4

206    model; class 2 uORFs initiate out-of-frame with the annotated start site and terminate within the coding
207    sequence; and class 3 uORFs initiate and terminate within the 5'UTR. Data from Cross (2015) on the
208    presence of Kozak sequences preceding class 1 uORFs suggests that approximately half are the correct
209    translation initiation site *in vivo*. In a PCR-based approach where a constitutive promoter is used, cloning
210    from the wrong ATG may result in an out-of-frame or truncated product, potentially removing essential
211    signal sequences for correct targeting. 57 of the 298 successfully cloned genes from Mackinder et al. (2017)
212    contained a class 1 in-frame ATG upstream of the cloned region, therefore ~10% of cloned regions may
213    have encoded truncated protein products.

215    *Introns, UTRs and splice variants* - Chlamydomonas genes have a relatively high intron frequency,
216    providing a further challenge for PCR-based cloning. The average gene contains 7.3 introns with an
217    average intron length of 373 bp compared to an average exon length of 190 bp. 94% of genes contain
218    introns between their start and stop codons, 13% of genes contain one or more introns in their 5'UTRs and
219    3.4% have introns in their 3'UTRs. ATG-Stop cloning would omit introns in UTR regions, potentially missing
220    critical regulatory information. Furthermore, approximately 9% of genes are annotated with two or more
221    transcript models that result from alternative splicing (Figure 1C). This variation would be missed through
222    cloning from cDNA or through gene synthesis that excludes native introns.

224    *Unsuitable primers* - ATG-Stop PCR cloning of either gDNA or cDNA results in limited flexibility of primer
225    design. Sequence analysis of a set of genome-wide primer pairs for ATG-Stop cloning (Mackinder et al.,
226    2017) indicates that primers are frequently of poor quality and unsuitable for efficient PCR. The average
227    primer in the dataset had a predicted melting temperature (Tm) of 69.2°C and an average GC content of
228    64.2%. Primer Tm and GC content are expected to be high in comparison to other organisms with less GC-
229    rich genomes, however, many primers also breached recommended thresholds pertaining to length,
230    secondary structure formation, repetitive sequences and 3' GC content. Primers are shown in Figure 1D
231    (blue bars) as having breached these four thresholds if, (1) they were longer than 30 bp; (2) the free energy
232    ($\Delta G$) required to disrupt secondary structure formation (self-dimers, cross-dimers or hairpins) was less than
233    -9 kcal mol$^{-1}$ at PCR-relevant annealing temperatures (66-72°C); (3) they contained mono- or dinucleotide
234    repeats of 5 or more; or (4) their 3' end contained five or more consecutive G/C bases. A stricter set of
235    thresholds is utilized by the Primer3 check_primers module (Rozen and Skaletsky, 2000), which results in
236    the rejection of over 60% of individual primers in the dataset, even when the program is set to ignore
237    predicted annealing temperatures (Figure 1D, orange bar). Under these settings, only 13% of pairs are free
238    from detectable issues in both primers. Interestingly, there is a high GC content mismatch between forward
239    and reverse primers with a considerably higher GC content of reverse primers (Figure 1D, inset).

241    Many individual genes contain a range of the above features that result in challenges faced during PCR
242    cloning or gene synthesis. Figure 1E shows a gene from chromosome 8 that exhibits several examples and
243    was a target for recombineering. Cre08.g379800 is >16 kbp with 40 introns, contains mono-, di-, tri- and
244    pentanucleotide repeat regions of ≥9 repeats. It also contains a potential misannotated upstream ATG in
245    the 5'UTR that could initiate a class 1 uORF, as well as seven class 3 uORFs (Cross, 2015). Cre08.g379800
246    structural information was obtained from the version 5.5 gene model currently available on Phytozome.

248    To further compare whether the challenges faced in Chlamydomonas were similar in other organisms we
249    analysed gene size and gene complexity relative to gene size for the model eukaryote *S. cerevisiae,* the
250    model plant Arabidopsis and the ~17 Gb hexaploid genome of *Triticum aestivum* (bread wheat). Figure 1F
251    shows that Chlamydomonas has a higher proportion of long genes and fewer short genes than the three
252    other genomes tested, along with a considerably higher average gene size for Chlamydomonas (5322 bp
253    versus 1430 bp for yeast, 2187 bp for Arabidopsis and 3521 bp for chromosome-assigned genes in wheat).
254    Unlike wheat, Arabidopsis and yeast, Chlamydomonas genes show a trend of increasing complexity per
255    kilobase for longer genes (Figure 1F), potentially in line with an increase in average UTR length as gene
256    length increases (Salomé and Merchant, 2019).

**Recombineering pipeline development**

To overcome the challenges associated with PCR-based cloning we developed a high-throughput recombineering pipeline for large-scale parallel cloning of Chlamydomonas nuclear genes from BACs with their native promoter regions intact. During pipeline development we decided to pursue a simplified 1-step DNA retrieval recombineering approach rather than a BAC editing approach (i.e. Poser et al., 2008; Brumos et al., 2020) for several reasons: (1) Using a gene retrieval method enables all cloning to be performed in the BAC host *E. coli* strain, thereby avoiding BAC purification, which can be timely and low yielding; (2) assembled constructs contain only the gene of interest making them considerably smaller than the original BAC, this allows a medium copy origin of replication to be used that improves ease of handling, and the smaller constructs minimize DNA fragmentation during Chlamydomonas transformation (Zhang et al., 2014); (3) BACs contain many genes, with additional copies of adjacent genes to the gene of interest potentially having an unwanted phenotypic impact on transformed Chlamydomonas lines; (4) the backbone of the available BAC collection lacks a suitable Chlamydomonas selection marker, therefore additional BAC editing to insert a suitable selection marker (Aksoy and Forest, 2019) or inefficient and poorly understood plasmid co-transformation strategies would be required for selection; and (5) a typical BAC engineering approach would require two recombination steps, which would increase pipeline time, decrease pipeline efficiency and add further challenges due to the repetitive nature of the Chlamydomonas genome.

The simplicity of our pipeline enables completion in four days using only generic reagents. The final recombineered construct is a vector containing the target region (typically including the native promoter, 5'UTR and open reading frame) recombined in-frame with a downstream fluorescent protein followed by the *PSAD* terminator (see Figure 2 for a pipeline schematic and Supplemental Method 1 for a detailed protocol). Our pipeline has four key steps: (1) *E. coli* harbouring a BAC containing the gene of interest is made recombination competent by transformation with the pRed vector containing the lambda viral *exo*, *beta* and *gam* genes (Redαβγ) and *recA* (Sarov et al., 2006) (Figure 2A); (2) Redαβγ and *recA* induction by arabinose followed by transformation with a linear tagging cassette including 50 bp homology arms to the target gene (Figure 2B); (3) kanamycin selection for successful recombination events and temperature inhibition of the pRed pSC101 replication origin to minimise further undesired recombination (Figure 2C); and (4) plasmid isolation and verification via restriction digest and junction sequencing (Figure 2D).

The original tagging cassette consists of the codon optimized YFP CrVenus, a 3xFLAG tag, the *PSAD* terminator, the paromomycin selection marker (*AphVIII)*, the p15A medium-copy-number origin of replication and the kanamycin resistance gene (*Kan^R*). Amplification of the tagging cassette from pLM099 is performed using primers containing 50 bp homology arms corresponding to regions flanking the target gene; the forward primer at least 2,000 bp upstream of the start codon to encompass the native 5' promoter and UTR region and the reverse primer at the 3' end of the coding region (immediately upstream of the stop codon). The annealing site of the reverse primer can easily be altered to amplify a cassette from pLM099 that can clone genes without a fluorescent tag or with only the 3xFLAG tag (see Supplemental Method 1). To minimise false positives due to pLM099 carryover, pLM099 contains the *ccdB* counter selection gene (Bernard and Couturier, 1992). In addition, the cassette includes an I-SceI restriction site. I-SceI has an 18 bp recognition site not found within the reference Chlamydomonas genome (strain CC-503) and allows cassette linearization prior to transformation into Chlamydomonas.

We initially tested our pipeline on 12 targets. To ensure that the BAC library (available from the Chlamydomonas Resource Center; https://www.chlamycollection.org/) was correctly mapped we performed PCR to check for the presence of the 5' and 3' ends of our target genes (Figure S1A). We next implemented the pipeline according to a small-scale batch protocol (Supplemental Method 1A). For all targets except one, plasmids isolated from most picked colonies gave a correct banding pattern after restriction digest (Figure S1B). After sequence confirmation we successfully cloned 11 out of our 12 targets, resulting in a 92% success rate (Figure S1C). To further expand the capabilities of our pipeline we tested whether we could successfully recombineer a large and complex gene from a fosmid (available from the Chlamydomonas Resource Center). We targeted SAGA1 (Cre11.g467712; fosmid VTP41289), that had previously been highly challenging to gene synthesize (see above; Itakura et al., 2019) and was not available in the BAC library. Restriction digest of recombineered plasmids purified from three colonies all showed the correct digestion pattern (Figure S1D). Sequencing confirmed that the 19,601 bp target region,

310 that included 2,913 bp upstream of the predicted SAGA1 start codon, was successfully cloned. Confident
311 that our recombineering method was robust we pursued the development of a large-scale pipeline that
312 would allow the parallel tagging of genes with most steps achievable in 96-well format.
313
314 **Successful large-scale application of the recombineering pipeline**
315 To test the efficiency of the pipeline we shortlisted 191 genes which could be mapped to a clone from the
316 Chlamydomonas BAC library. To more easily identify BACs within the library that contain a target gene we
317 designed a Python script (BACSearcher; Supplemental Code) and have outputted the five smallest BACs
318 for all targets in the genome in Supplemental Data Set 1, revealing that 86% of nuclear genes are covered
319 by at least one BAC (87% if BACs are included that terminate within 3'UTRs). BACSearcher also enables
320 automated design of primers containing 50 bp homology regions to target genes in optimal positions; the
321 script reports suitable 5' homology regions 2000-3000 bp upstream of the annotated start codon and takes
322 into account local DNA complexity features, including mono- and dinucleotide repeating runs and GC
323 content. This feature can be easily modified to design 5' homology regions further upstream of the target
324 (see Supplemental Method 2A). The length of 50 bp is short enough to design into an oligonucleotide but
325 long enough to be unlikely to share homology with more than one site within a BAC. Supplemental Data
326 Set 1 includes sequences for the top five optimal 5' homology regions for each target, all >2000 bp upstream
327 of the start codon, along with the corresponding 50 bp 3' homology region. In addition, four pairs of primer
328 sequences are included that can be used to check for the presence of each target in a BAC.
329 Our 191 targets were primarily chosen based on our 2017 association study for CCM components
330 (Mackinder et al., 2017), transcriptomics (Brueggeman et al., 2012; Fang et al., 2012) and pyrenoid
331 proteomics (Mackinder et al., 2016; Zhan et al., 2018). 81 genes previously targeted in 2017 were retried
332 here by recombineering, this time with >2000 bp upstream sequence included. 41 of these were previously
333 unsuccessful by PCR and 40 were previously successful but included here in order to compare the effect
334 of retaining the native promoter. These included five targets that contain a class 1 uORF (Cross, 2015) and
335 so may have previously produced misleading localization data due to expression of a truncated protein.
336 Selection of the remaining 110 targets was guided by new pyrenoid proteome (Zhan et al., 2018) and CCM
337 interactome data (Mackinder et al., 2017). *E. coli* strains containing the correct BAC as identified by
338 BACSearcher were recovered from the BAC library and processed in parallel using 96-format culturing
339 plates. To optimise the efficiency of our high-throughput pipeline, we successively ran the pipeline three
340 times removing successful targets once confirmed. Supplemental Method 1B provides a detailed protocol
341 for the optimized high-throughput pipeline. In summary, 100% of our 191 target BAC lines were made
342 recombination competent (Figure 2A) and out of the 191 target genes, one gene-specific tagging cassette
343 failed to amplify (Figure 2B), likely due to the formation of secondary structure(s) within the 50 bp homology
344 regions of the primers. Of the 190 that amplified successfully, 187 yielded colonies after selection with
345 kanamycin (Figure 2C). Validation by enzymatic digestion confirmed that 146 of these lines contained
346 correct recombineering plasmid products (Figure 2D). Recombineering plasmid products from the 146
347 successful lines were extracted and their junctions confirmed by Sanger sequencing. Our high-throughput
348 pipeline had an overall efficiency of 76%, an average recombineered region length of 7259 bp and a
349 maximum cloned length of 22,773 bp corresponding to gene Cre10.g427850 (Supplemental Data Set 2).
350 26 target genes that were unsuccessful by PCR in 2017 were successfully cloned here by recombineering,
351 and all five previously successful targets containing class 1 uORFs retried here were successful.
352 During pipeline development, we found that optimising bacterial growth prior to transformation with
353 the recombineering cassette was critical (see protocol notes in Supplemental Method 1). In addition, for 14
354 out of the 146 correctly recombineered lines in our high-throughput pipeline, use of an alternative BAC from
355 the library yielded success after an initial failure. We found that for approximately half of the target genes it
356 was necessary to validate multiple colonies by enzymatic digest in order to rule out false positives;
357 beginning with the 187 colony-producing lines from our high-throughput pipeline, picking just a single colony
358 gave a 49% success rate, screening a second colony increased the success rate to 66% and a third colony
359 gave a 76% success rate. For a small proportion of targets screening >3 colonies led to identification of a
360 correctly recombineered construct (Figure 2E). Restriction digest analysis of plasmids isolated from
361 incorrectly assembled recombineering events suggested that cloning could fail due to a broad range of

362 reasons including cassette recircularization, cassette duplication, cassette insertion into the BAC or
363 retrieval of incorrect target regions. Increasing homology arm length, using alternative homology arms,
364 using alternative BACs and using fosmids are potential solutions to overcome incorrect recombineering for
365 specific targets. Supplemental Data Set 1 provides up to five options for homology arms and up to five
366 available BACs per gene, and BACSearcher can be easily modified to increase homology arm length (see
367 Supplemental Method 2A). Taken together with our 12 initial targets, we successfully cloned 157 out of 203
368 target regions from BACs using our recombineering pipeline, achieving an efficiency of 77%.
369
370 **Cloning success is size independent and tolerant of sequence complexity**
371 To investigate if our developed recombineering approach was gene size and complexity independent, we
372 compared our successful targets against unsuccessful targets (Figure 3). Here we define a target region to
373 mean the ATG-stop ORF for PCR-based cloning and the ATG-stop ORF plus an upstream region of >2000
374 bp designed to encompass the 5'UTR and native promoter for recombineering. The results show that there
375 is no significant difference in the region lengths between cloned and uncloned targets for recombineering
376 (Figure 3A; Mann-Whitney $U$ = 3303, $P$ = 0.38, two-tailed), indicating that our method is target size
377 independent. This contrasts to the clear effect of target size on cloning success for our previous PCR-based
378 cloning data (Figure 3A; Mackinder et al., 2017). We then compared our cloning success to the number of
379 simple and global repeats per kilobase in target regions. Our method appears far more tolerant of repetitive
380 sequences than PCR-based cloning, both in the per-kilobase prevalence of simple and global repeats and
381 in the number of repeats per target region (Figure 3B and 3C). For our recombineering pipeline there is no
382 significant difference detectable in the average repeat prevalence per kilobase between cloned and
383 uncloned regions (Mann-Whitney $U$ = 3129, $P$ = 0.17, two-tailed), while there is a clear negative effect on
384 PCR-based cloning success for targets with over ~4.8 repeats per kbp (Figure 3B). For the most repetitive
385 targets involved in our analysis (>9 repeats per kbp), our recombineering cloning efficiency remained above
386 60%; an efficiency over three times higher than PCR-based cloning (Figure 3B). Extrapolation of these data
387 overlaid with the genome wide distribution of repeat frequencies indicates that a large proportion of genes
388 that are technically challenging for PCR-based cloning due to high repeat frequencies may be cloned by
389 recombineering (Figure 3B).
390
391 **Localization of Venus-tagged proteins**
392 To assess the validity of the pipeline for localization studies we transformed wild type Chlamydomonas cells
393 with a subset of linearized recombineering plasmid products tagged at the C-terminus with CrVenus (Figure
394 4A). Paromomycin resistant colonies were directly screened for YFP fluorescence on transformation plates,
395 picked, grown in TP media at air-levels of $CO_2$ (~0.04%), imaged by fluorescence microscopy to examine
396 the localization pattern (Figure 4B and Figure S2) and immunoblotted against the C-terminal 3xFLAG
397 epitope to confirm fusion protein size (Figure S2A). Transformed genes were selected based on previous
398 affinity purification mass spectrometry data (Mackinder et al., 2017) and pyrenoid proteomics data
399 (Mackinder et al., 2016; Zhan et al., 2018). The localization data supports the proteomics data with PSAF
400 (Photosystem I subunit F; Cre09.g412100), ISA1 (Isoamylase 1; Cre03.g155001) and CSP41B
401 (Chloroplast Stem-loop Binding Protein of 41 kDa B; Cre10.g435800) present in the pyrenoid. PSAF is a
402 core transmembrane subunit of photosystem I. As expected PSAF shows strong colocalization with
403 chlorophyll outside of the pyrenoid, however in addition it clearly localizes to the thylakoid tubules traversing
404 the pyrenoid. Interestingly, in the pyrenoid tubules the chlorophyll signal is minimal, particularly at the
405 "pyrenoid tubule knot" where the tubules converge (Engel et al., 2015). These data along with the
406 localization of other PSI and PSII components to the pyrenoid tubules (Mackinder et al., 2017) suggest that
407 the tubules contain both PSI and PSII but that chlorophyll-containing light harvesting complexes found
408 within the pyrenoid may be quenched or at low abundance. Tagged Cre17.g702500 (TAB2), a protein linked
409 to early PSI assembly (Dauvillée et al., 2003) and which was identified as an interactor with PSBP4 found
410 within and at the periphery of the pyrenoid (Mackinder et al., 2017), was also enriched at the pyrenoid.
411 Interestingly, the location of TAB2 is not just restricted to the pyrenoid periphery but is also found within the
412 pyrenoid forming distinct small foci (Figure 4B). This may indicate that early PSI assembly could be
413 occurring within the pyrenoid as well as at the pyrenoid periphery (Uniacke and Zerges, 2009).

414      CSP41B localized to the pyrenoid matrix, and analysis of the translated product of CSP41B shows
415    that it belongs to a family of NAD-dependent epimerase/dehydratases (IPR001509) and contains a UDP-
416    galactose 4-epimerase domain that may be involved in galactose metabolism. Its role in pyrenoid function
417    is unclear. Localization of ISA1 shows it was enriched in the pyrenoid with an uneven distribution. ISA1 is
418    a starch debranching enzyme that is essential for starch synthesis with *ISA1* deletion lines lacking both
419    chloroplast and pyrenoid starch (Mouille et al., 1996). The presence of pyrenoid starch and its correct
420    organization is critical for correct CCM function (Itakura et al., 2019; Toyokawa et al., 2020), with the
421    absence of starch in an *ISA1* knock out (4-D1) having incorrect LCIB localization (see below), retarded
422    growth at very low $CO_2$ (0.01% v/v) and reduced inorganic carbon affinity (Toyokawa et al., 2020).
423    Interestingly in Toyokawa et al. (2020) they failed to attain localization data for an ISA1-mCherry fusion
424    driven by the HSP70A/RBCS2 hybrid promoter.
425      Cre14.g613950 encodes a protein belonging to the ABC transporter family identified as an
426    interactor of HLA3 (high light activated gene 3) (Mackinder et al., 2017), a putative $HCO_3^-$ transporter
427    located in the plasma membrane (Duanmu et al., 2009; Gao et al., 2015). Like HLA3, Cre14.g613950 shows
428    a typical plasma membrane localization pattern with YFP signal at the cell periphery and signal typical of
429    the Golgi network. However, immunoblotting against the C-terminal 3xFLAG tag of Cre14.g613950 in two
430    independent transformants shows a smaller molecular weight band than predicted (Figure S2). This
431    potentially indicates that the gene model for Cre14.g613950 is incorrect or that the protein undergoes post-
432    translation cleavage as seen for other CCM related proteins that transit via the secretory pathway
433    (Fukuzawa et al., 1990; Tachiki et al., 1992).
434

435    **Development of backbones with additional tags and markers**
436    To further expand the functional application of our recombineering pipeline we designed additional
437    backbone vectors that enable protein tagging with the fluorophores mScarlet-i (Bindels et al., 2017),
438    mNeonGreen (Shaner et al., 2013) and mTurquoise2 (Goedhart et al., 2012) and that allow selection with
439    hygromycin or zeocin (Figure 5A and 5B). This enables complementation of Chlamydomonas Library
440    Project (CLiP) mutants that have been generated using the *AphVIII* marker conferring paromomycin
441    resistance (Li et al., 2016; Li et al., 2019) and also enables expression of two or three differently tagged
442    proteins within the same cell. For comparison, we tested these vectors on *LCI9* (Cre09.g394473), which
443    encodes the low-$CO_2$ inducible protein LCI9 that, via PCR-based cloning, we previously showed to localize
444    to the pyrenoid periphery (Mackinder et al., 2017). Recombineered *LCI9* was 7160 bp long including the
445    native promoter region. All fluorophores displayed the same pyrenoid periphery localization pattern (Figure
446    5C) and agree with the localization information obtained when LCI9 expression was driven from the *PSAD*
447    promoter (Figure 5C bottom image; the *PSAD* promoter is here defined as the sequence spanning from 3-
448    763 bp upstream of the *PSAD* start codon (Cre05.g238332), encompassing both the 5'UTR and promoter
449    region), thus further supporting the use of ~2000 bp upstream regions as promoters for fusion protein
450    expression.
451      To further confirm that localization of proteins driven by their native promoter does not differ from
452    those driven by the constitutive *PSAD* promoter we compared localization between *native*-LCIB-Venus and
453    *PSAD*-LCIB-Venus. LCIB is an essential CCM component that shows dynamic relocalization to the pyrenoid
454    periphery at $CO_2$ levels <0.04% (Yamano et al., 2010). LCIB expressed from its endogenous promoter was
455    localized to the pyrenoid periphery at very low $CO_2$ (0.01% v/v), in full agreement with localization data
456    when LCIB expression is driven by the constitutive *PSAD* promoter (Figure 5D).
457      Finally, we tested that our recombineering pipeline could be used to successfully complement a
458    CLiP mutant. We transformed *native-LCIB-Venus* (cloned into pLM161 that contains the *APHVII* gene
459    conferring hygromycin resistance) into a CLiP *lcib* mutant (LMJ.RY0402.215132). Four transformants
460    showing Venus fluorescence were selected for microscopy and growth phenotyping. All showed a typical
461    pyrenoid peripheral localization pattern when grown at very low $CO_2$ and all rescued the *lcib* mutant
462    phenotype to varying degrees, with *lcib::LCIB-Venus-1* showing complete rescue (Figure S3).
463

464    **Maintaining the native promoter enables relative protein abundances to be monitored**

465   As our pipeline retains the native promoter of the target gene we hypothesized that fluorescence output
466   would be representative of relative changes in protein abundance in response to environmental conditions.
467   To test this we grew lines with LCIB driven from either the constitutive *PSAD* promoter (*PSAD*-LCIB-Venus)
468   or its native promoter (*Native*-LCIB-Venus). LCIB-Venus signal stayed relatively constant between high (3%
469   v/v) and low (0.04% v/v) $CO_2$ when LCIB was expressed from the *PSAD* promoter (*PSAD*-LCIB-Venus),
470   but showed an approximate 8-fold increase between these conditions when the native promoter was used,
471   with this change consistent across three independently transformed lines (Figure 5E). This agrees with
472   previous immunoblotting data, in which a comparable fold increase was seen in LCIB abundance when
473   cells were transferred from high $CO_2$ to low $CO_2$ (Yamano et al., 2010). This indicates that our
474   recombineering lines can be used to monitor relative protein abundance across different growth conditions.
475
476   **Discussion**
477   We have established a rapid recombineering based method to clone large and complex Chlamydomonas
478   genes from BACs. Our approach circumvents the challenges associated with cloning large, GC-rich and
479   complex genes that are prevalent in Chlamydomonas. We demonstrate that the method can be applied for
480   small batch cloning as well as 96-well high-throughput cloning. Our overall cloning success rate (combined
481   batch and high-throughput results) was 77%, considerably higher than our previous PCR-based high-
482   throughput cloning pipeline (48%), which was inflated due to an enrichment of small target genes. Our
483   overall success rate is slightly lower when compared to recombineering pipelines in other organisms, with
484   success rates of 89% achieved in *C. elegans* (Sarov et al., 2012) and ~93% for Arabidopsis (Brumos et al.,
485   2020). This reduced overall efficiency is likely due to the complexity of the Chlamydomonas genome (Figure
486   1), with DNA secondary structure having been previously linked to recombineering failure (Nelms and
487   Labosky, 2011). We expect a higher success rate when the pipeline is applied to small sample numbers
488   since it is easier to optimise bacterial growth prior to electrotransformation on a per-sample basis if there
489   are fewer samples to manage. This may be evidenced by our successful cloning of 11 out of 12 targets in
490   an initial batch-scale pipeline attempt (Figure S1), although the sample size is insufficient to generalize from
491   with confidence.
492           To enable expression of multiple fluorophores simultaneously and for the complementation of CLiP
493   mutants we designed a series of vectors with modern fluorophores and varying selection markers and
494   demonstrated their performance in Chlamydomonas (Figure 5). The presence of either 3xFLAG or 3xHA
495   tag enables use of the vectors for affinity purification to explore interacting partners of tagged proteins.
496   Different fluorophore pairs (i.e. mNeonGreen and mScarlet-i) could also be used for FRET based studies
497   to explore protein-protein interactions. In addition, all vectors can be used for cloning genes without
498   fluorescence tags or with only short affinity tags (3xFLAG and 3xHA).
499           Due to the size independence of our method we could maintain the native promoter of target genes.
500   For two genes, LCI9 and LCIB, the comparison between native promoter-driven expression and *PSAD*
501   promoter-driven expression showed no noticeable differences in localization. Interestingly, using a native
502   promoter allows relative protein abundance to be tracked between conditions (Figure 5E). Once validated,
503   acquiring relative abundance data is straightforward and can be easily parallelized. This enables relative
504   protein abundance to be tracked in real-time across a broad range of conditions. Future experiments could
505   include tracking relative protein abundance in 96-well libraries of tagged proteins in response to a
506   perturbation (i.e. high to low $CO_2$ transition). This would be highly supportive of available transcriptomic and
507   proteomic data sets and provide novel insights into cellular processes (Mettler et al., 2014; Zones et al.,
508   2015; Strenkert et al., 2019). Although our relative abundance data for LCIB appears to closely reflect
509   immunoblotting data, it should be noted that using a native promoter may not always fully reflect native
510   changes. This discrepancy can be due to insertional effects caused by integration into transcriptionally
511   unfavourable regions of the genome and absence of cis-regulatory regions in the recombineered construct,
512   or transcriptional silencing (Schroda, 2019). At a protein level, fluorescent protein folding time could affect
513   protein stability and turnover and the presence of the fused fluorescence protein could affect function or
514   multi-subunit assembly.
515

516 Whilst our approach allows the native promoter, 5'UTR region and open reading frame to be cloned, the
517 native 3'UTR is not maintained. This could be addressed through a two-step recombineering pipeline where
518 the tag is first inserted into the BAC at the desired location, markers could then be removed via a Flp-*FRT*
519 recombinase system (Sarov et al., 2006; Brumos et al., 2020), and the edited target gene can then be
520 retrieved into a final Chlamydomonas expression vector. When establishing our pipeline, we decided not
521 to pursue this strategy in order to maximise the success rate by limiting the number of steps, with a focus
522 on developing a simple, easy to apply approach. In addition, whilst we have focused on C-terminal tagging
523 as this allows conservation of N-terminal transit peptides required for organelle targeting, our
524 recombineering pipeline could be applied for N-terminal tagging by modification of our cloning vectors with
525 a constitutive promoter and N-terminal tag.
526        The simplicity of our framework and vector design could be adopted for other organisms with
527 relative ease provided a BAC or fosmid library and efficient transformation protocols are available. Multiple
528 features of our recombineering cassette could make adaptation to different organisms relatively
529 straightforward, such as the use of ccdB counter-selection and the rare I-SceI recognition site used for
530 linearization of the recombineering cassette prior to transformation. For organisms in which selection with
531 paromomycin, hygromycin or zeocin is ineffective, or which cannot utilise the *AphVII*, *AphVIII* or *BLE* genes
532 included in the pLM099-derived cassettes, alternative selection genes can be quickly incorporated by
533 restriction-ligation using flanks containing KpnI and I-SceI recognition sites at the 5' and 3' respectively.
534
535 One limitation we encountered was that only 86% of nuclear genes are covered by the BAC library.
536 However, this value only takes into account ~73% of BACs, with the remaining BACs potentially incorrectly
537 mapped to the current version of the Chlamydomonas genome (see Supplemental Method 2B). Our
538 analysis suggests the true percentage of genes covered could be higher than 86% but confirming this may
539 require a careful re-mapping of the library. A promising solution is cloning from fosmids. We demonstrated
540 that our pipeline can be successfully applied for cloning from fosmids and a Chlamydomonas fosmid library
541 is now available (released July 2020; Chlamydomonas Resource Center). The use of fosmids, with smaller
542 DNA fragments compared to BACs, could help improve efficiency by reducing off-target recombination
543 between the PCR-amplified cassette and the BAC or by reducing recombination between two repetitive
544 regions of the BAC. In addition, the fosmid library is expected to have close to 100% genome coverage.
545
546 Our recombineering approach has enabled the efficient cloning of large and complex genes that could not
547 be achieved via PCR-based cloning. It opens the door to a better understanding of the functional role of a
548 large fraction of the Chlamydomonas genome though protein localization, protein-protein interaction
549 studies, real-time monitoring of relative protein abundance and complementation of mutants (e.g. random
550 insertion and CRISPR/Cas generated mutants). In addition, it provides a highly complementary method to
551 the recently released CLiP mutant collection.
552

**Methods**

**Availability of materials, data and software**

All plasmid sequences are available in Supplemental Data Set 4 and have been deposited in GenBank with the following IDs: pLM099, MT737960; pLM160, MT737961; pLM161, MT737962; pLM162, MT737963; pLM459, MT737964. Plasmids are available from the Chlamydomonas Resource Center (https://www.chlamycollection.org/), as are the BAC and fosmid libraries. Full protocols for batch and high-throughput recombineering are available in Supplemental Method 1. Data used for the genome analyses presented in Figure 1 are available on request. The python computer code used for identifying BACs, fosmids and suitable homology regions for recombineering is supplied as Supplemental Code and is available at https://github.com/TZEmrichMills/Chlamydomonas_recombineering.

**Plasmid and cassette construction**

Fragments for pLM099 were amplified by PCR (Phusion Hotstart II polymerase, ThermoFisher Scientific) from the following plasmids: Venus-3xFLAG, PSAD terminator and *AphVIII* from pLM005 (Mackinder et al., 2017); the p15A origin of replication from pNPC2; the *Kan^R* resistance gene from pLM007; the counter-selection *ccdB* gene from Gateway pDONR221 Vector (ThermoFisher Scientific). The resulting amplicons were gel purified (MinElute Gel Extraction Kit, QIAGEN) and assembled by Gibson assembly (see Figure 5A for detailed map). pLM160 was constructed from pLM099 to replace CrVenus with mNeonGreen (Shaner et al., 2013), and pLM161 was constructed from pLM099 to replace the paromomycin resistance gene (*AphVIII)* with the hygromycin resistance gene (*AphVII)*. pLM162 was constructed from pLM161 with the synthetic fluorophore mScarlet-i (Bindels et al., 2017) replacing CrVenus. pLM459 was constructed from pLM161 to replace CrVenus with mTurquoise2 (Goedhart et al., 2012), the 3xFLAG with the 3xHA haemagglutinin tag, and *AphVII* with the zeocin resistance gene (*Sh ble)*. Gene-specific cloning primers were designed to amplify a ~4.6 kbp cassette from the recombineering vectors pLM099, 160, 161, 162 and 459 (Figure 5), excluding *ccdB*, and providing 50 bp of sequence homology to the target gene an average of ~2500 bp upstream of the 5'UTR and directly upstream of the stop codon. This enables the retrieval of each target gene into the cassette in frame with a fluorescent tag and with the native promoter region intact. All oligonucleotide and plasmid sequences can be found in Supplemental Data Sets 3 and 4**.**

**Culturing**

*E. coli* cells were cultured in lysogeny broth (LB) or yeast extract nutrient broth (YENB) at 37°C unless they contained the temperature sensitive pSC101-BAD-gbaA-tet (pRed), in which case 30°C was used. All DNA for transformation was introduced by electroporation and transformants were recovered in super optimal broth with catabolite repression (SOC). DH10B cells containing fragments of the Chlamydomonas genome in the form of BACs were obtained from the Clemson University Genomics Institute (now distributed by the Chlamydomonas Resource Center, University of Minnesota, USA). DB3.1 cells expressing the *ccdB* antidote gene, *ccdA*, were obtained from ThermoFisher Scientific and used for maintenance of the recombineering vectors.

Chlamydomonas wild type cells (strain CC-4533) were cultured in Tris-acetate-phosphate media (TAP) with revised Hutner's trace elements (Kropat et al., 2011) and illuminated by white fluorescent light. Assembled recombineering vectors were prepared for transformation into Chlamydomonas by restriction digest with I-SceI endonuclease (NEB). Transformation and selection of fluorescence lines was performed in accordance with Mackinder et al. (2017) using a Typhoon Trio fluorescence scanner (GE Healthcare). Viable Chlamydomonas transformants were screened for CrVenus and mNeonGreen expression at 555/20 nm, and for mScarlet-i at 615/12 nm. Several strains emitting the strongest fluorescence for each line were picked. The average number of fluorescent colonies for recombineered Venus fusion proteins with their native promoter was ~10%, however this varied considerably between constructs (PSAF (10/134) 7%, TAB2 (6/44) 13.6%, CSP41B (6/43) 13.9%, ISA1 (25/297) 8%, Cre14.g613950 (2/22) 9%, LCI9 (6/25) 24%, LCIB (6/19) 31.5%). Picked fluorescent strains were cultured in Tris-phosphate minimal media (TP) under ambient $CO_2$ (~0.04%) conditions then imaged by fluorescent microscopy to visualise protein localization. To ensure that determined localizations were not due to in-frame integration of a fluorophore-containing-

605  fragment of the cassette with another gene we confirmed localization in at least two independent
606  transformants and performed immunoblotting against the 3xFLAG epitope to confirm expected fusion
607  protein size.
608          For spot tests cells were grown to ~8 x $10^6$ cells/ml in TAP at ~50 µmol photons/m$^2$/s, washed with
609  TP and then serial diluted in TP prior to spotting 1000, 100 and 10 cells on TP 1.5% agar plates. Replica
610  plates were incubated in 0.04% or 3% $CO_2$ chambers for 24 hours at 50 µmol photons/m$^2$/s, then 24 hours
611  at 150 µmol photons/m$^2$/s followed by 48 hours at 300 µmol photons/m$^2$/s prior to imaging.
612
613  **Protein extraction and immunoblotting**
614  Lines expressing recombineered fusion proteins were cultured in 50 ml TAP media containing 5 µg/mL
615  paromomycin to a cell density of ~2 x $10^6$ cells/ml. Cells were harvested by centrifugation at 5,000 x g for
616  10 min at room temperature. The supernatant was discarded, and the pellet was resuspended in 500 µL
617  of protein extraction buffer (20 mM Tris-HCl pH 7.5, 5 mM $MgCl_2$, 300 mM NaCl, 5 mM DTT, 0.1% Triton
618  X100, Roche protease inhibitor) then flash frozen in liquid nitrogen in 100 µl aliquots. Cells were thawed
619  on ice and flash frozen again before a final thaw on ice. Samples were then centrifuged at 17,000 x g for
620  15 min at 4°C to separate the soluble and insoluble fractions. The soluble supernatant was transferred to
621  a new tube and mixed 1:1 with 2x Laemmli buffer containing β-mercaptoethanol, then heated at 80°C for
622  10 min prior to SDS-PAGE.
623          15-30 µl of each sample was loaded onto a 10% mini-protean TGX gel (Bio-Rad) then transferred
624  to a polyvinylidene difluoride (PVDF) membrane via semi-dry transfer (10V, 60 min). Fusion proteins were
625  immuno-detected using the monoclonal anti-flag M2 antibody (1:1000; Sigma-Aldrich; catalog # F1804)
626  followed by Alexa-Fluor 555 goat anti-mouse secondary antibody (1:10 000; Invitrogen; catalog # A-
627  21422). The membrane was imaged using a Typhoon 5 Scanner.
628
629  **Microscopy**
630  Sample preparation for microscopy was performed as per (Mackinder et al., 2017). Images were acquired
631  using a Zeiss LSM880 confocal microscope on an Axio Observer Z1 invert, equipped with a 63x 1.40 NA
632  oil planapochromat lens. Images were analysed using ZEN 2.1 software (Zeiss) and FIJI. Excitation and
633  emission filter settings were as follows: Venus and mNeonGreen, 514 nm excitation, 525-550 nm emission;
634  mScarlet-i, 561 nm excitation, 580-600 nm emission; and chlorophyll, 561 nm excitation, 665-705 nm
635  emission.
636
637  **Plate reader assay**
638  To monitor fluorescence changes in response to $CO_2$, three independent *native*-LCIB-Venus lines, a single
639  *PSAD*-LCIB-Venus line and WT were grown in TP bubbled at low $CO_2$ (0.04%) or high $CO_2$ (3%) at 300
640  µmol photons/m$^2$/s. Four samples per line were aliquoted into a 96-well plate and chlorophyll (excitation
641  625/34, emission 692/50) and Venus (excitation 504/10, emission 540/12)  fluorescence was immediately
642  measured using a BMG Labtech Clariostar Plate Reader. Venus fluorescence was normalised by
643  chlorophyll then WT background subtracted. The average low $CO_2$ fluorescence was divided by the average
644  high $CO_2$ fluorescence for each line. Error was calculated by the propagation of variance across both low
645  and high $CO_2$ values and is shown as the standard error of the mean.
646
647  **Recombineering procedure for 1-step subcloning and tagging**
648  The following outlines the batch-scale recombineering protocol. Extended batch and multi-well plate-scale
649  recombineering protocols are supplied in Supplemental Method 1.
650          For each target, a recombineering cassette was amplified from plasmid pLM099 (Phusion
651  Hotstart II polymerase, ThermoFisher Scientific) using primers containing 50 bp homology arms, one
652  homologous to a region upstream of the annotated start codon of the target gene, and one homologous to
653  the 3' end of the coding sequence (excluding the stop codon). The resulting PCR product was purified
654  (MinElute Gel Extraction Kit, QIAGEN) and its concentration measured using a nanodrop
655  spectrophotometer. Upstream region lengths ranged from 1000-4000 bp from the start codon, with an

656 average of ~2500 bp. For two genes, Cre04.g220200 and Cre16.g678661, the first 50 bp of the 5'UTR
657 was used as the upstream homology region due to BAC coverage limitations.

658        The pRed plasmid, pSC101-BAD-gbaA-tet, was extracted from *E. coli* cells grown overnight at
659 30°C (Plasmid Mini Kit, QIAGEN), and its concentration measured by nanodrop. *E. coli* cells harbouring a
660 BAC containing the target gene were recovered from the Chlamydomonas BAC library and used to
661 inoculate 20 ml of YENB media containing 12.5 µg/ml chloramphenicol, followed by overnight growth in a
662 50 ml conical flask at 37°C with vigorous shaking. After 16 h of growth, 120 µl of the culture was used to
663 inoculate 4 ml of fresh YENB containing 12.5 µg/ml chloramphenicol. This was grown for ~2 h at 37°C until
664 an optical density ($OD_{600}$) of 2 was reached. 2 ml of the culture was then incubated on ice for 2 min, followed
665 by centrifugation at 5000 x g for 10 min at 4°C. After removing the supernatant, the pellet was placed back
666 on ice and washed by resuspension in 1 ml of chilled 10% glycerol, followed immediately by centrifugation
667 at 5000 x g for 10 min at 4°C. The resulting supernatant was removed, and the pellet was placed back on
668 ice and resuspended in 100 µl of 0.1 ng/µl pRed. This mixture was transferred to a pre-chilled 2 mm gap
669 electroporation cuvette and electroporated at 2500 V, 400 Ω and 25 µF using a Gene Pulser II (Bio-Rad).
670 The electroporated cells were immediately recovered in 800 µl SOC and incubated at 30°C for 90 min with
671 vigorous shaking. The whole outgrowth was added to 20 ml of YENB containing 12.5 µg/ml chloramphenicol
672 and 5 µg/ml tetracycline and grown overnight at 30°C with vigorous shaking.

673        After 16 h of growth, 600 µl of culture was used to inoculate 4 ml of fresh YENB containing 12.5
674 µg/ml chloramphenicol and 5 µg/ml tetracycline. This was grown for 3 h at 30°C, or until reaching an $OD_{600}$
675 >2, at which point 80 µl of 10% L-arabinose was added to induce pRed expression and growth was shifted
676 to 37°C for 1 h with vigorous shaking. 2 ml of the induced culture was incubated on ice for 2 min, then
677 centrifuged at 5000 x g for 10 min at 4°C, the supernatant removed, and the pellet placed back on ice. Cells
678 were then washed in 10% glycerol, centrifuged at 5000 x g for 10 min at 4°C, the supernatant removed,
679 and the pellet placed back on ice. The pellet was resuspended in 100 µl of 5 ng/µl PCR product and
680 transferred to a pre-chilled 2 mm gap electroporation cuvette, followed by electroporation as before.
681 Electroporated cells were immediately added to 800 µl of SOC and recovered at 37°C for 90 min with
682 vigorous shaking. 450 µl of outgrowth was spread onto 1.5% LB-agar containing 25 µg/ml kanamycin, air-
683 dried and incubated overnight at 37°C. Selected colonies were used to inoculate 4 ml of LB containing 25
684 µg/ml kanamycin and grown for 16-18 h at 37°C with shaking. Recombineering products were extracted
685 and validated by restriction digest using appropriate enzymes, followed by Sanger sequencing using
686 primers designed to amplify the junctions between the pLM099-derived cassette and the target region.

687

688 **Statistics**
689 Confidence intervals for Figure 1A were calculated using the Wilson score interval method based on the
690 number of attempted and successfully cloned ATG-Stop amplicons per size category in Mackinder et al.
691 (2017). Statistical differences in the distribution of sizes and repeat frequencies between successful and
692 unsuccessful PCR and recombineering targets (presented in Figure 3) were assessed using the Mann-
693 Whitney U test. A non-parametric test was chosen based on results of the Kolmogorov-Smirnov test for
694 normality for recombineering targets. Test statistics are detailed in Supplemental Table 1.

695

696 **Genome analysis**
697 Chlamydomonas, Arabidopsis, yeast and wheat nuclear genes were analysed for gene size and sequence
698 complexity. Gene sizes are defined from the start of the 5'UTR to the end of the 3'UTR. Note that in Figure
699 1A the predicted clonable proportion of genes in each size category is based on cloning success for ATG-
700 Stop regions not full genes. Sequence complexity is defined in relation to intron prevalence, GC content,
701 and the prevalence of various repeat regions. We designate regions containing a high frequency of repeats
702 as being more complex than regions with a low frequency. This reflects the increased potential for cloning
703 complications presented by sequences with large numbers of repetitive regions, though it differs from
704 descriptions given by Morgulis et al. (2006). Sequences were analysed for complexity using the freely
705 available bioinformatics software detailed below (see Supplemental Method 3 for settings), and outputs
706 were processed using custom python scripts (Supplemental Code; see Supplemental Method 4 for usage
707 information). GC content was calculated using annotated bases only.

708

*Sequence data sources* – Unspliced Chlamydomonas nuclear gene sequences used for the analyses were generated using a custom python script (see Supplemental Code) to extract whole-gene, 5'UTR, ATG-Stop and 3'UTR sequences from the genome based on their start and end positions in the current gene models (Phytozome version 5.5). Chlamydomonas gene models are based on predictions using Augustus (annotation version u11.6) and refined using a range of RNA-seq datasets. Files containing the whole genome nucleotide sequence (version 5.0) and the annotation information for each of the 17,741 nuclear genes (version 5.5) were downloaded from Phytozome 12 and are provided as precursor files for running the BACSearcher script (see Supplemental Code and Supplemental Method 2). Sequence data for *Arabidopsis thaliana* (TAIR10 assembly) and *Triticum aestivum* nuclear genes (International Wheat Genome Consortium assembly) were obtained from EnsemblPlants BioMart. Analysis was limited to the 105,200 chromosome-assigned wheat genes. Sequence data for *Saccharomyces cerevisiae* (S288C reference genome, 2015 release) were obtained from the Saccharomyces Genome Database. Gene sequences were appended to include all annotated UTRs and introns, resulting in a dataset that is more closely comparable to the unspliced gene data used for Chlamydomonas, Arabidopsis and wheat.

*Analysis of repeats* – Repetitive regions in the nucleotide sequences analysed in this work are categorized into simple and global repeats. We use the term simple repeats to refer to relatively short (tens to hundreds of bases) repetitive regions in a nucleotide sequence that display regular or semi-regular repeating patterns. We include consecutive repeating motifs of varying unit lengths, known as tandem repeats, as well as inverted patterns in which a short region is followed closely (or immediately, if palindromic) by its reverse complement sequence. Chlamydomonas genes were analysed for tandem repeats using Tandem Repeats Finder (Benson, 1999). The default settings were modified to provide a cut-off for detection such that no repeats under 10 bp in length were reported (see Supplemental Method 3A). All Tandem Repeats Finder outputs were processed using a custom python script and analysed in spreadsheet format to generate mean values for the number of genes with either, (1) at least one mononucleotide repeat ≥10 bp in length and with ≥90% identity; (2) at least one di- or trinucleotide repeat ≥20 bp in length with ≥90% identity; (3) at least one tandem repeat ≥20 bp in length, with a period length of four or more (tetra+), with ≥90% identity; and (4) the mean number of repeats of these types per kilobase of sequence.

Chlamydomonas genes were analysed for inverted repeats using the Palindrome Analyser webtool (Brázda et al., 2016), available at http://bioinformatics.ibp.cz:9999/#/en/palindrome. The default settings were modified to report repeats with a maximum of 1 mismatch for every 10 bp of stem sequence, a maximum spacer length of 10 bp and a maximum total length of 210 bp (see Supplemental Method 3B for settings). All Palindrome Analyser outputs were downloaded and analysed in spreadsheet format to generate mean values for the number of genes containing one or more inverted repeats over 20 bp long with ≥90% identity and the mean number of inverted repeats of this type per kilobase.

All nuclear genes from Chlamydomonas (Figure 1B), Arabidopsis, yeast and wheat (Figure 1F), and recombineering target regions (Figure 3B and C) were analysed for global repeats using the NCBI WindowMasker program (Morgulis et al., 2006). We use the term global repeats to denote the combined number of individual masked regions detected by the WindowMasker modules DUST and WinMask. DUST detects and masks shorter repetitive regions including tandem and inverted repeats, overlapping with and providing support for the Tandem Repeats Finder and Palindrome Analyser outputs. WinMask detects and masks families of longer repetitive regions that do not necessarily occur adjacently in the genome. Default settings were used throughout (see Supplemental Method 3C). These modules mask repetitive regions using only the supplied sequence as a template.

Chlamydomonas repeats localized to the 5'UTRs, ATG-Stop regions and 3'UTRs were distinguished using positional information from Phytozome (genome annotation version 5.5). Repeats that spanned from a 5'UTR across the start codon or across the stop codon into the 3'UTR were not counted, though were included in the whole-gene repeat analyses described above.

*uORFs, transcripts and intron analysis* – Data on the presence of uORFs in Chlamydomonas transcripts were obtained from the results of a BLASTP analysis performed by Cross (2015) and adapted to provide

760     the per-gene values. A list of Chlamydomonas transcripts was downloaded from Phytozome Biomart and
761     used to identify the number of genes with more than one transcript model. Genomic data detailing the
762     number and order of exons within each gene were also downloaded from Phytozome Biomart; this
763     information was used to ascertain the number of genes containing introns in their translated and
764     untranslated regions.
765

766     *Primer analysis* – To assess the impact of inefficient priming on PCR-based cloning, analysis was
767     performed on a dataset of PCR primers designed to clone every gene in the Chlamydomonas genome from
768     start to stop codon using gDNA as the template and generated such that the predicted Tm difference for
769     each pair was not more than 5°C where possible. Primer sequences were then assessed against four
770     thresholds pertaining to efficient priming, set in accordance with advice found in the Primer3 manual,
771     support pages provided by IDT, and the Premier Biosoft technical notes. These thresholds relate to primer
772     length, propensity for secondary structure formation, the presence of repeats and the GC content of the 3'
773     end. Long primers can have a reduced amplification efficiency, secondary structure formation can reduce
774     the number of primers available to bind to the intended template during a PCR, multiple repeats can
775     increase the risk of mispriming, and a high 3' end GC content can increase the risk of primer-dimer
776     formation. Thresholds for each were set as follows: (1) primer length should not be more than 30 bp, (2)
777     the $\Delta G$ required to disrupt predicted secondary structures should be above -9 kcal/mol at 66 or 72°C, (3)
778     tandem single nucleotides or dinucleotide motifs should repeat no more than 4 times, and (4) the 3' end
779     should consist of no more than 4 G/C bases in a row. The number of primers in breach of each of these
780     thresholds is shown in Figure 1D as a percentage of the dataset. The percentage of unsuitable primer pairs
781     was calculated by counting pairs for which one or both primers breached one or more of these thresholds.
782     Tm considerations were omitted from analysis since Chlamydomonas genes have an unusually high GC
783     content, so primers designed to amplify gDNA are expected to have higher than recommended Tms
784     according to generic primer design guidelines. GC content was calculated using annotated bases only.
785         To complement these results, primers were analysed using the check_primers algorithm from
786     Primer3 (Rozen and Skaletsky, 2000). Settings used were as default for Primer3Plus (Untergasser et al.,
787     2007) – an updated, online version of the Primer3 package – with minimal modifications that included
788     removing the Tm constraints (see Supplemental Method 3D for full settings used). The output was analysed
789     with a custom python script that reported the primary reason for rejection of individual primers (see
790     Supplemental Method 4C). Tm was removed as a constraint to allow for more detailed analysis of primer
791     sequence parameters, since the default maximum allowable Tm for Primer3Plus is 63°C, which results in
792     rejection of almost 90% of primers for this reason alone if used. 1.6% of primers were too long to be
793     considered for analysis (>36 bp); these were included in Figure 1D (orange bar) as having been rejected
794     for breaching the length constraint. The majority of rejected primers produced one of the following three
795     reasons for rejection: (1) 'high end complementarity' for primer pairs, which implies a high likelihood that
796     the 3' ends of the forward and reverse primers will anneal, enabling amplification of a short, heterogeneous
797     primer-dimer (cross-dimer); (2) 'high end complementarity' for single primers, which implies a high likelihood
798     that a primer's 3' end will bind to that of another identical copy, self-priming to form a homogenous primer-
799     dimer (self-dimer); and (3) 'high any complementarity' for single primers, which implies a high likelihood of
800     self-annealing without necessarily self-priming, relevant to both the inter-molecular annealing of identical
801     copies and to instances of hairpin formation resulting from intra-molecular annealing. Primers rejected for
802     these three reasons are labelled in Figure 1D (orange bar) as cross-dimers, self-dimers and hairpins,
803     respectively.
804

805     *Note on differences between Chlamydomonas BAC library strain and CLiP mutant strain* – The
806     Chlamydomonas BAC library was constructed using the genome reference strain CC-503, so researchers
807     working with alternative strains need to take into account potential genomic divergence. For example, here
808     we transformed recombineered DNA from the BAC library into CC-4533, the wild type strain used for the
809     CLiP mutant collection and a popular strain for studying the CCM. Genomic analysis of CC-4533 relative
810     to CC-503 has revealed 653 instances of variation that may be disruptive to protein function, although only
811     three of these are unique to CC-4533 when compared to other common lab strains (Li et al., 2016). Two

812  genes affected by this variation were successfully cloned using our recombineering pipeline;
813  Cre06.g250650 in CC-4533 contains three short deletions relative to CC-503 with an uncertain impact on
814  the protein, while Cre06.g249750 in CC-4533 contains a predicted inversion affecting the final three exons
815  and part of the 3'UTR.
816
817  **BACSearcher python resource**
818  Suitable BACs containing the target genes were identified using a python script that also identifies 50 bp
819  binding sites for recombineering cloning primers and provides sequences for primers that can be used to
820  check for the presence of a target gene within a BAC (see Supplemental Method 2). BACSearcher output
821  is available for all 17,741 genes in the genome in Supplemental Data Set 1. For individual targets in our
822  recombineering pipeline that were not covered by a BAC in the BACSearcher output, an alternative method
823  was employed to search for BAC coverage. This method is detailed in Supplemental Method 2, along with
824  usage and modification instructions for BACSearcher, including instructions to output suitable fosmids for
825  all genes in the genome. BACSearcher resources can also be found in the associated GitHub repository at
826  https://github.com/TZEmrichMills/Chlamydomonas_recombineering.
827
828  **Accession numbers**
829  Cre11.g467712: SAGA1
830  Cre09.g412100: PSAF
831  Cre03.g155001: ISA1
832  Cre10.g435800: CSP41B
833  Cre17.g702500: TAB2
834  Cre10.g452800: LCIB
835  Cre09.g394473: LCI9
836
837  **Author contributions**
838  TZEM developed the initial recombineering pipeline, designed and assembled the original pLM099
839  recombineering plasmid and performed the genome wide analysis. GY, TZEM and TKK assembled
840  additional recombineering plasmids. TZEM and GY optimized and performed the large-scale
841  recombineering pipeline. GY performed the microscopy and Venus quantification data. PG validated the
842  pipeline using fosmids. JB performed the complementation experiments. JB, IG, CSL, CEW and TKK
843  supported the development and implementation of the recombineering pipeline. JWD wrote the
844  BACSearcher code and provided bioinformatics support to TZEM for the remaining code. LCMM conceived
845  the idea and led the research. LCMM and MPJ received funding to support the work. LCMM, TZEM and
846  GY wrote the manuscript.
847
848  **Supplemental data**
849  Supplemental Figure 1. Batch-scale recombineering results
850  Supplemental Figure 2. Validation of fluorescently localized lines
851  Supplemental Figure 3. Complementation of the *lcib* CLiP mutant
852  Supplemental Table 1. Mann-Whitney U test statistics
853  Supplemental Method 1. Protocols for batch and large-scale recombineering
854  Supplemental Method 2. BACSearcher usage
855  Supplemental Method 3. Bioinformatics software usage
856  Supplemental Method 4. Bioinformatics python analysis
857  Supplemental Data Set 1. BACSearcher output
858  Supplemental Data Set 2. Large-scale pipeline results summary
859  Supplemental Data Set 3. Oligonucleotide sequences
860  Supplemental Data Set 4. Plasmid sequences
861  Supplemental Code. BACSearcher python code, BACSearcher precursor files and python codes for
862  processing outputs from bioinformatics programs and generating unspliced gene sequences
863

874     **References**
875

876     **Aksoy, M., and Forest, C.** (2019). One step modification of Chlamydomonas reinhardtii BACs
877             using the RED/ET system. Mediterranean Agricultural Sciences **32,** 49-55.
878     **Baier, T., Wichmann, J., Kruse, O., and Lauersen, K.J.** (2018). Intron-containing algal
879             transgenes mediate efficient recombinant gene expression in the green microalga
880             *Chlamydomonas reinhardtii*. Nucleic Acids Research **46,** 6909-6919.
881     **Barahimipour, R., Strenkert, D., Neupert, J., Schroda, M., Merchant, S.S., and Bock, R.**
882             (2015). Dissecting the contributions of GC content and codon usage to gene expression
883             in the model alga *Chlamydomonas reinhardtii*. The Plant Journal **84,** 704-717.
884     **Benson, G.** (1999). Tandem repeats finder: a program to analyze DNA sequences. Nucleic Acids
885             Research **27,** 573-580.
886     **Bernard, P., and Couturier, M.** (1992). Cell killing by the F plasmid CcdB protein involves
887             poisoning of DNA-topoisomerase II complexes. Journal of Molecular Biology **226,** 735-
888             745.
889     **Bindels, D.S., Haarbosch, L., Van Weeren, L., Postma, M., Wiese, K.E., Mastop, M.,**
890             **Aumonier, S., Gotthard, G., Royant, A., and Hink, M.A.** (2017). mScarlet: a bright
891             monomeric red fluorescent protein for cellular imaging. Nature Methods **14,** 53.
892     **Brázda, V., Kolomazník, J., Lýsek, J., Hároníková, L., Coufal, J., and Šťastný, J.** (2016).
893             Palindrome analyser–a new web-based server for predicting and evaluating inverted
894             repeats in nucleotide sequences. Biochemical and biophysical research communications
895             **478,** 1739-1745.
896     **Brueggeman, A.J., Gangadharaiah, D.S., Cserhati, M.F., Casero, D., Weeks, D.P., and**
897             **Ladunga, I.** (2012). Activation of the carbon concentrating mechanism by $CO_2$ deprivation
898             coincides with massive transcriptional restructuring in *Chlamydomonas reinhardtii*. The
899             Plant Cell **24,** 1860-1875.
900     **Brumos, J., Zhao, C., Gong, Y., Soriano, D., Patel, A.P., Perez-Amador, M.A., Stepanova,**
901             **A.N., and Alonso, J.M.** (2020). An Improved Recombineering Toolset for Plants. The
902             Plant Cell **32,** 100.
903     **Clarke, L., Rebelo, C., Goncalves, J., Boavida, M., and Jordan, P.** (2001). PCR amplification
904             introduces errors into mononucleotide and dinucleotide repeat sequences. Molecular
905             Pathology **54,** 351.
906     **Copeland, N.G., Jenkins, N.A., and Court, D.L.** (2001). Recombineering: a powerful new tool
907             for mouse functional genomics. Nature Reviews Genetics **2,** 769-779.
908     **Cross, F.R.** (2015). Tying Down Loose Ends in the Chlamydomonas Genome: Functional
909             Significance of Abundant Upstream Open Reading Frames. G3 **6,** 435-446.
910     **Crozet, P., Navarro, F.J., Willmund, F., Mehrshahi, P., Bakowski, K., Lauersen, K.J., Pérez-**
911             **Pérez, M.-E., Auroy, P., Gorchs Rovira, A., Sauret-Gueto, S., Niemeyer, J., Spaniol,**
912             **B., Theis, J., Trösch, R., Westrich, L.-D., Vavitsas, K., Baier, T., Hübner, W., de**
913             **Carpentier, F., Cassarini, M., Danon, A., Henri, J., Marchand, C.H., de Mia, M.,**
914             **Sarkissian, K., Baulcombe, D.C., Peltier, G., Crespo, J.-L., Kruse, O., Jensen, P.-E.,**
915             **Schroda, M., Smith, A.G., and Lemaire, S.D.** (2018). Birth of a Photosynthetic Chassis:
916             A MoClo Toolkit Enabling Synthetic Biology in the Microalga *Chlamydomonas reinhardtii*.
917             ACS Synthetic Biology **7,** 2074-2086.
918     **Dauvillée, D., Stampacchia, O., Girard-Bascou, J., and Rochaix, J.D.** (2003). Tab2 is a novel
919             conserved RNA binding protein required for translation of the chloroplast psaB mRNA.
920             The EMBO journal **22,** 6378-6388.
921     **Duanmu, D., Miller, A.R., Horken, K.M., Weeks, D.P., and Spalding, M.H.** (2009). Knockdown
922             of limiting-$CO_2$-induced gene HLA3 decreases $HCO_3^-$ transport and photosynthetic Ci
923             affinity in *Chlamydomonas reinhardtii*. Proceedings of the National Academy of Sciences
924             **106,** 5990-5995.
925     **Engel, B.D., Schaffer, M., Kuhn Cuellar, L., Villa, E., Plitzko, J.M., and Baumeister, W.** (2015).
926             Native architecture of the *Chlamydomonas* chloroplast revealed by in situ cryo-electron
927             tomography. Elife **4,** e04889.

928 **Fang, W., Si, Y., Douglass, S., Casero, D., Merchant, S.S., Pellegrini, M., Ladunga, I., Liu, P.,**
929     **and Spalding, M.H.** (2012). Transcriptome-wide changes in *Chlamydomonas reinhardtii*
930     gene expression regulated by carbon dioxide and the $CO_2$-concentrating mechanism
931     regulator CIA5/CCM1. The Plant Cell **24,** 1876-1893.
932 **Fukuzawa, H., Fujiwara, S., Yamamoto, Y., Dionisio-Sese, M.L., and Miyachi, S.** (1990).
933     cDNA cloning, sequence, and expression of carbonic anhydrase in *Chlamydomonas*
934     *reinhardtii*: regulation by environmental $CO_2$ concentration. Proceedings of the National
935     Academy of Sciences **87,** 4383.
936 **Gao, H., Wang, Y., Fei, X., Wright, D.A., and Spalding, M.H.** (2015). Expression activation and
937     functional analysis of HLA 3, a putative inorganic carbon transporter in *Chlamydomonas*
938     *reinhardtii*. The Plant Journal **82,** 1-11.
939 **Goedhart, J., Von Stetten, D., Noirclerc-Savoye, M., Lelimousin, M., Joosen, L., Hink, M.A.,**
940     **Van Weeren, L., Gadella, T.W., and Royant, A.** (2012). Structure-guided evolution of
941     cyan fluorescent proteins towards a quantum yield of 93%. Nature Communications **3,** 1-
942     9.
943 **Hommelsheim, C.M., Frantzeskakis, L., Huang, M., and Ülker, B.** (2014). PCR amplification of
944     repetitive DNA: a limitation to genome editing technologies and many other applications.
945     Scientific Reports **4,** 5052.
946 **Itakura, A.K., Chan, K.X., Atkinson, N., Pallesen, L., Wang, L., Reeves, G., Patena, W.,**
947     **Caspari, O., Roth, R., and Goodenough, U.** (2019). A Rubisco-binding protein is
948     required for normal pyrenoid number and starch sheath morphology in *Chlamydomonas*
949     *reinhardtii*. Proceedings of the National Academy of Sciences **116,** 18445-18454.
950 **Kobayashi, Y., Takusagawa, M., Harada, N., Fukao, Y., Yamaoka, S., Kohchi, T., Hori, K.,**
951     **Ohta, H., Shikanai, T., and Nishimura, Y.** (2015). Eukaryotic Components Remodeled
952     Chloroplast Nucleoid Organization during the Green Plant Evolution. Genome Biology and
953     Evolution **8,** 1-16.
954 **Kropat, J., Hong-Hermesdorf, A., Casero, D., Ent, P., Castruita, M., Pellegrini, M., Merchant,**
955     **S.S., and Malasarn, D.** (2011). A revised mineral nutrient supplement increases biomass
956     and growth rate in *Chlamydomonas reinhardtii*. The Plant Journal **66,** 770-780.
957 **Levinson, G., and Gutman, G.A.** (1987). Slipped-strand mispairing: a major mechanism for DNA
958     sequence evolution. Molecular biology and evolution **4,** 203-221.
959 **Li, X., Zhang, R., Patena, W., Gang, S.S., Blum, S.R., Ivanova, N., Yue, R., Robertson, J.M.,**
960     **Lefebvre, P.A., Fitz-Gibbon, S.T., Grossman, A.R., and Jonikas, M.C.** (2016). An
961     Indexed, Mapped Mutant Library Enables Reverse Genetics Studies of Biological
962     Processes in *Chlamydomonas reinhardtii*. The Plant Cell **28,** 367.
963 **Li, X., Patena, W., Fauser, F., Jinkerson, R.E., Saroussi, S., Meyer, M.T., Ivanova, N.,**
964     **Robertson, J.M., Yue, R., Zhang, R., Vilarrasa-Blasi, J., Wittkopp, T.M., Ramundo, S.,**
965     **Blum, S.R., Goh, A., Laudon, M., Srikumar, T., Lefebvre, P.A., Grossman, A.R., and**
966     **Jonikas, M.C.** (2019). A genome-wide algal mutant library and functional screen identifies
967     genes required for eukaryotic photosynthesis. Nature Genetics.
968 **López-Paz, C., Liu, D., Geng, S., and Umen, J.G.** (2017). Identification of *Chlamydomonas*
969     *reinhardtii* endogenous genic flanking sequences for improved transgene expression. The
970     Plant Journal **92,** 1232-1244.
971 **Lumbreras, V., Stevens, D.R., and Purton, S.** (1998). Efficient foreign gene expression in
972     *Chlamydomonas reinhardtii* mediated by an endogenous intron. The Plant Journal **14,**
973     441-447.
974 **Mackinder, L.C., Meyer, M.T., Mettler-Altmann, T., Chen, V.K., Mitchell, M.C., Caspari, O.,**
975     **Freeman Rosenzweig, E.S., Pallesen, L., Reeves, G., Itakura, A., Roth, R., Sommer,**
976     **F., Geimer, S., Muhlhaus, T., Schroda, M., Goodenough, U., Stitt, M., Griffiths, H.,**
977     **and Jonikas, M.C.** (2016). A repeat protein links Rubisco to form the eukaryotic carbon-
978     concentrating organelle. Proceedings of the National Academy of Sciences **113,** 5958-
979     5963.
980 **Mackinder, L.C.M.** (2017). The *Chlamydomonas* $CO_2$-concentrating mechanism and its potential
981     for engineering photosynthesis in plants. New Phytologist **217,** 54-61.

Mackinder, L.C.M., Chen, C., Leib, R.D., Patena, W., Blum, S.R., Rodman, M., Ramundo, S., Adams, C.M., and Jonikas, M.C. (2017). A Spatial Interactome Reveals the Protein Organization of the Algal $CO_2$-Concentrating Mechanism. Cell **171,** 133-147.e114.

Merchant, S.S., Prochnik, S.E., Vallon, O., Harris, E.H., Karpowicz, S.J., Witman, G.B., Terry, A., Salamov, A., Fritz-Laylin, L.K., Marechal-Drouard, L., Marshall, W.F., Qu, L.H., Nelson, D.R., Sanderfoot, A.A., Spalding, M.H., Kapitonov, V.V., Ren, Q., Ferris, P., Lindquist, E., Shapiro, H., Lucas, S.M., Grimwood, J., Schmutz, J., Cardol, P., Cerutti, H., Chanfreau, G., Chen, C.L., Cognat, V., Croft, M.T., Dent, R., Dutcher, S., Fernandez, E., Fukuzawa, H., Gonzalez-Ballester, D., Gonzalez-Halphen, D., Hallmann, A., Hanikenne, M., Hippler, M., Inwood, W., Jabbari, K., Kalanon, M., Kuras, R., Lefebvre, P.A., Lemaire, S.D., Lobanov, A.V., Lohr, M., Manuell, A., Meier, I., Mets, L., Mittag, M., Mittelmeier, T., Moroney, J.V., Moseley, J., Napoli, C., Nedelcu, A.M., Niyogi, K., Novoselov, S.V., Paulsen, I.T., Pazour, G., Purton, S., Ral, J.P., Riano-Pachon, D.M., Riekhof, W., Rymarquis, L., Schroda, M., Stern, D., Umen, J., Willows, R., Wilson, N., Zimmer, S.L., Allmer, J., Balk, J., Bisova, K., Chen, C.J., Elias, M., Gendler, K., Hauser, C., Lamb, M.R., Ledford, H., Long, J.C., Minagawa, J., Page, M.D., Pan, J., Pootakham, W., Roje, S., Rose, A., Stahlberg, E., Terauchi, A.M., Yang, P., Ball, S., Bowler, C., Dieckmann, C.L., Gladyshev, V.N., Green, P., Jorgensen, R., Mayfield, S., Mueller-Roeber, B., Rajamani, S., Sayre, R.T., Brokstein, P., Dubchak, I., Goodstein, D., Hornick, L., Huang, Y.W., Jhaveri, J., Luo, Y., Martinez, D., Ngau, W.C., Otillar, B., Poliakov, A., Porter, A., Szajkowski, L., Werner, G., Zhou, K., Grigoriev, I.V., Rokhsar, D.S., and Grossman, A.R. (2007). The Chlamydomonas genome reveals the evolution of key animal and plant functions. Science **318,** 245-250.

Mettler, T., Mühlhaus, T., Hemme, D., Schöttler, M.-A., Rupprecht, J., Idoine, A., Veyel, D., Pal, S.K., Yaneva-Roder, L., Winck, F.V., Sommer, F., Vosloh, D., Seiwert, B., Erban, A., Burgos, A., Arvidsson, S., Schönfelder, S., Arnold, A., Günther, M., Krause, U., Lohse, M., Kopka, J., Nikoloski, Z., Mueller-Roeber, B., Willmitzer, L., Bock, R., Schroda, M., and Stitt, M. (2014). Systems Analysis of the Response of Photosynthesis, Metabolism, and Growth to an Increase in Irradiance in the Photosynthetic Model Organism *Chlamydomonas reinhardtii*. The Plant Cell **26,** 2310.

Meyer, M.T., Genkov, T., Skepper, J.N., Jouhet, J., Mitchell, M.C., Spreitzer, R.J., and Griffiths, H. (2012). Rubisco small-subunit α-helices control pyrenoid formation in Chlamydomonas. Proceedings of the National Academy of Sciences **109,** 19474-19479.

Morgulis, A., Gertz, E.M., Schäffer, A.A., and Agarwala, R. (2006). WindowMasker: window-based masker for sequenced genomes. Bioinformatics **22,** 134-141.

Mouille, G., Maddelein, M.L., Libessart, N., Talaga, P., Decq, A., Delrue, B., and Ball, S. (1996). Preamylopectin Processing: A Mandatory Step for Starch Biosynthesis in Plants. The Plant Cell **8,** 1353.

Nelms, B.L., and Labosky, P.A. (2011). A predicted hairpin cluster correlates with barriers to PCR, sequencing and possibly BAC recombineering. Scientific Reports **1,** 106.

Neupert, J., Karcher, D., and Bock, R. (2009). Generation of Chlamydomonas strains that efficiently express nuclear transgenes. The Plant Journal **57,** 1140-1150.

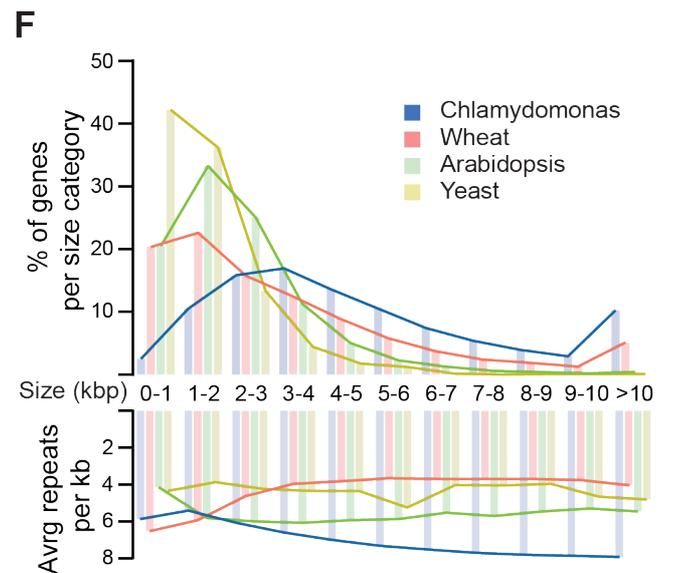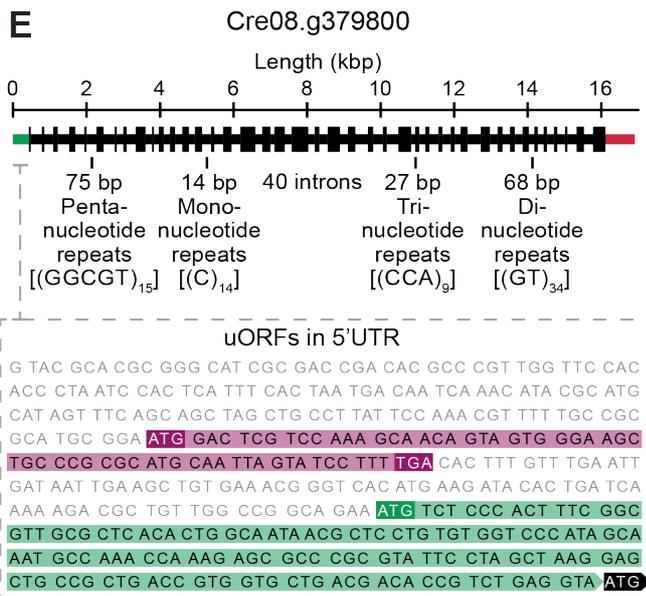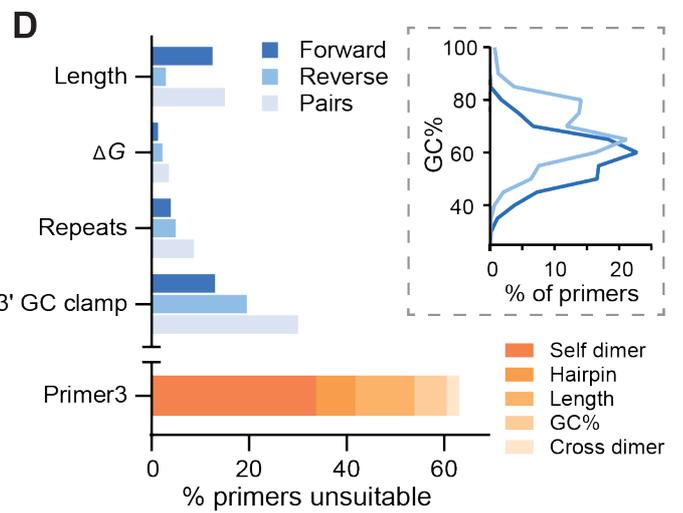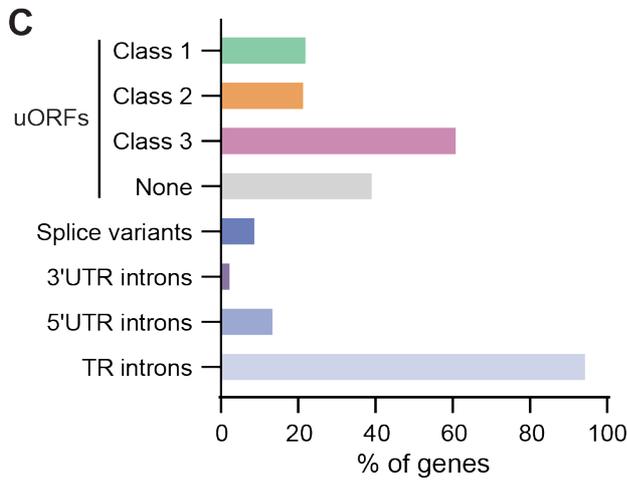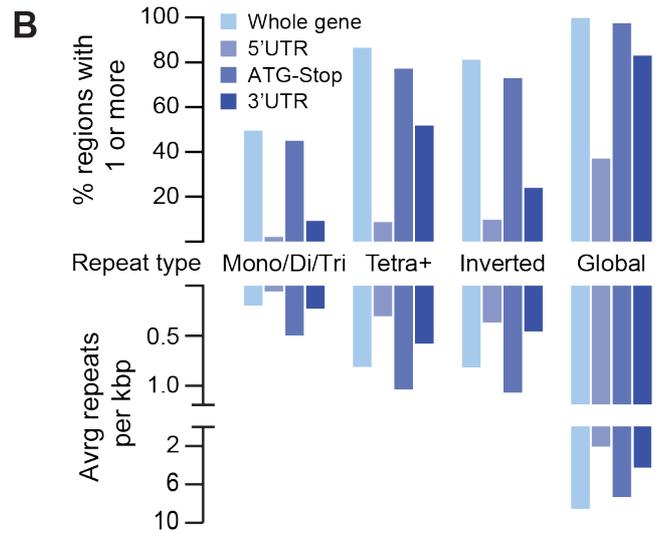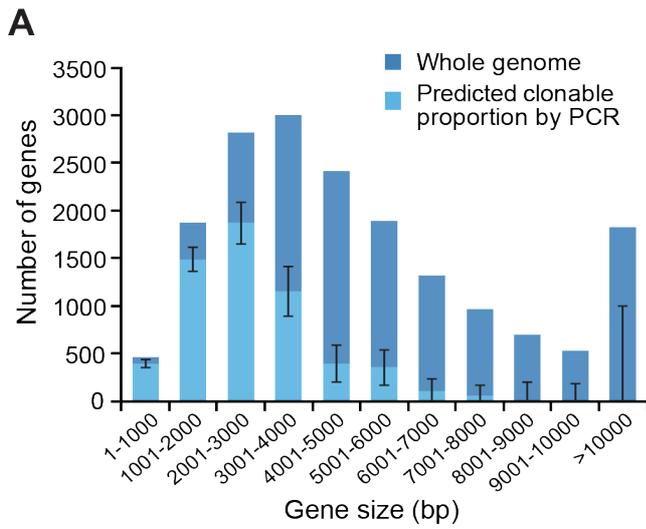Onishi, M., and Pringle, J.R. (2016). Robust Transgene Expression from Bicistronic mRNA in the Green Alga *Chlamydomonas reinhardtii*. G3 **6,** 4115-4125.

Poser, I., Sarov, M., Hutchins, J.R., Heriche, J.K., Toyoda, Y., Pozniakovsky, A., Weigl, D., Nitzsche, A., Hegemann, B., Bird, A.W., Pelletier, L., Kittler, R., Hua, S., Naumann, R., Augsburg, M., Sykora, M.M., Hofemeister, H., Zhang, Y., Nasmyth, K., White, K.P., Dietzel, S., Mechtler, K., Durbin, R., Stewart, A.F., Peters, J.M., Buchholz, F., and Hyman, A.A. (2008). BAC TransgeneOmics: a high-throughput method for exploration of protein function in mammals. Nature Methods **5,** 409-415.

Rae, B.D., Long, B.M., Förster, B., Nguyen, N.D., Velanis, C.N., Atkinson, N., Hee, W.Y., Mukherjee, B., Price, G.D., and McCormick, A.J. (2017). Progress and challenges of engineering a biophysical carbon dioxide-concentrating mechanism into higher plants. Journal of Experimental Botany **68,** 3717–3737.

1037  **Rasala, B.A., Lee, P.A., Shen, Z., Briggs, S.P., Mendez, M., and Mayfield, S.P.** (2012). Robust
1038      expression and secretion of Xylanase1 in *Chlamydomonas reinhardtii* by fusion to a
1039      selection gene and processing with the FMDV 2A peptide. PLoS One **7,** e43349.
1040  **Rozen, S., and Skaletsky, H.** (2000). Primer3 on the WWW for general users and for biologist
1041      programmers. In Bioinformatics methods and protocols (Springer), pp. 365-386.
1042  **Sahdev, S., Saini, S., Tiwari, P., Saxena, S., and Saini, K.S.** (2007). Amplification of GC-rich
1043      genes by following a combination strategy of primer design, enhancers and modified PCR
1044      cycle conditions. Molecular and cellular probes **21,** 303-307.
1045  **Salomé, P.A., and Merchant, S.S.** (2019). A Series of Fortunate Events: Introducing
1046      Chlamydomonas as a Reference Organism. The Plant Cell **31,** 1682.
1047  **Sarov, M., Schneider, S., Pozniakovski, A., Roguev, A., Ernst, S., Zhang, Y., Hyman, A.A.,**
1048      **and Stewart, A.F.** (2006). A recombineering pipeline for functional genomics applied to
1049      *Caenorhabditis elegans*. Nature Methods **3,** 839-844.
1050  **Sarov, M., Murray, J.I., Schanze, K., Pozniakovski, A., Niu, W., Angermann, K., Hasse, S.,**
1051      **Rupprecht, M., Vinis, E., and Tinney, M.** (2012). A genome-scale resource for in vivo
1052      tag-based protein function exploration in *C. elegans*. Cell **150,** 855-866.
1053  **Sarov, M., Barz, C., Jambor, H., Hein, M.Y., Schmied, C., Suchold, D., Stender, B., Janosch,**
1054      **S., K, J.V., Krishnan, R.T., Krishnamoorthy, A., Ferreira, I.R., Ejsmont, R.K., Finkl,**
1055      **K., Hasse, S., Kampfer, P., Plewka, N., Vinis, E., Schloissnig, S., Knust, E.,**
1056      **Hartenstein, V., Mann, M., Ramaswami, M., VijayRaghavan, K., Tomancak, P., and**
1057      **Schnorrer, F.** (2016). A genome-wide resource for the analysis of protein localisation in
1058      Drosophila. Elife **5,** e12068.
1059  **Schroda, M.** (2019). Good News for Nuclear Transgene Expression in Chlamydomonas. Cells **8**.
1060  **Shaner, N.C., Lambert, G.G., Chammas, A., Ni, Y., Cranfill, P.J., Baird, M.A., Sell, B.R., Allen,**
1061      **J.R., Day, R.N., and Israelsson, M.** (2013). A bright monomeric green fluorescent protein
1062      derived from Branchiostoma lanceolatum. Nature Methods **10,** 407.
1063  **Strenkert, D., Schmollinger, S., Gallaher, S.D., Salomé, P.A., Purvine, S.O., Nicora, C.D.,**
1064      **Mettler-Altmann, T., Soubeyrand, E., Weber, A.P., and Lipton, M.S.** (2019). Multiomics
1065      resolution of molecular events during a day in the life of Chlamydomonas. Proceedings of
1066      the National Academy of Sciences **116,** 2374-2383.
1067  **Tachiki, A., Fukuzawa, H., and Miyachi, S.** (1992). Characterization of Carbonic Anhydrase
1068      Isozyme CA2, Which Is the CAH2 Gene Product, in *Chlamydomonas reinhardtii*.
1069      Bioscience, Biotechnology, and Biochemistry **56,** 794-798.
1070  **Toyokawa, C., Yamano, T., and Fukuzawa, H.** (2020). Pyrenoid Starch Sheath Is Required for
1071      LCIB Localization and the $CO_2$-Concentrating Mechanism in Green Algae. Plant
1072      Physiology.
1073  **Uniacke, J., and Zerges, W.** (2009). Chloroplast protein targeting involves localized translation
1074      in *Chlamydomonas*. Proceedings of the National Academy of Sciences **106,** 1439-1444.
1075  **Untergasser, A., Nijveen, H., Rao, X., Bisseling, T., Geurts, R., and Leunissen, J.A.** (2007).
1076      Primer3Plus, an enhanced web interface to Primer3. Nucleic Acids Research **35,** W71-
1077      W74.
1078  **Wang, Y., Stessman, D.J., and Spalding, M.H.** (2015). The $CO_2$ concentrating mechanism and
1079      photosynthetic carbon assimilation in limiting $CO_2$: how *Chlamydomonas* works against
1080      the gradient. Plant Journal **82,** 429-448.
1081  **Weiner, I., Atar, S., Schweitzer, S., Eilenberg, H., Feldman, Y., Avitan, M., Blau, M., Danon,**
1082      **A., Tuller, T., and Yacoby, I.** (2018). Enhancing heterologous expression in
1083      *Chlamydomonas reinhardtii* by transcript sequence optimization. The Plant Journal **94,**
1084      22-31.
1085  **Yamano, T., Tsujikawa, T., Hatano, K., Ozawa, S.-i., Takahashi, Y., and Fukuzawa, H.** (2010).
1086      Light and low-$CO_2$-dependent LCIB–LCIC complex localization in the chloroplast supports
1087      the carbon-concentrating mechanism in *Chlamydomonas reinhardtii*. Plant and Cell
1088      Physiology **51,** 1453-1468.
1089  **Yu, D., Ellis, H.M., Lee, E.C., Jenkins, N.A., and Copeland, N.G.** (2000). An efficient
1090      recombination system for chromosome engineering in *Escherichia coli*. Proceedings of
1091      the National Academy of Sciences **97,** 5978-5983.

1092 **Zhan, Y., Marchand, C.H., Maes, A., Mauries, A., Sun, Y., Dhaliwal, J.S., Uniacke, J.,**
1093     **Arragain, S., Jiang, H., and Gold, N.D.** (2018). Pyrenoid functions revealed by
1094     proteomics in *Chlamydomonas reinhardtii*. PloS One **13**.
1095 **Zhang, R., Patena, W., Armbruster, U., Gang, S.S., Blum, S.R., and Jonikas, M.C.** (2014).
1096     High-throughput genotyping of green algal mutants reveals random distribution of
1097     mutagenic insertion sites and endonucleolytic cleavage of transforming DNA. The Plant
1098     Cell **26**, 1398-1409.
1099 **Zones, J.M., Blaby, I.K., Merchant, S.S., and Umen, J.G.** (2015). High-resolution profiling of a
1100     synchronized diurnal transcriptome from *Chlamydomonas reinhardtii* reveals continuous
1101     cell and metabolic differentiation. The Plant Cell **27**, 2743-2769.
1102
1103

**Figure 1**

1106

1107 **Figure 1.** Chlamydomonas nuclear genes are often large, complex, or misannotated, affecting PCR-based
1108 cloning attempts and transgene expression success.
1109 **A** The distribution of gene sizes for the 17,741 genes in the Chlamydomonas nuclear genome (dark blue).
1110 Gene sizes are measured from the start of the 5'UTR to the end of the 3'UTR. Within each size category,
1111 the predicted proportion amenable to PCR-based cloning is shown in light blue. These proportions were
1112 extrapolated from cloning success for 624 CCM-related genes from Mackinder et al. (2017) in which PCR-
1113 based cloning was used to amplify the ATG-Stop region of each gene, excluding any UTRs. The strong
1114 size-dependence of ATG-Stop cloning efficiency seen in 2017 indicates that 68% of the genome would be
1115 challenging to clone. 95% confidence intervals for the predicted clonable proportions of each size category
1116 were calculated using the Wilson score interval method. No genes over 8000 bp are predicted to be clonable
1117 by PCR although only a handful of regions of these sizes were tested in 2017 giving rise to the large
1118 confidence intervals for these categories.
1119 **B** Genome wide sequence complexity as indicated by the presence of one or more repetitive sequences
1120 and frequency of repeats per kilobase (kbp) in each gene (pale blue). Values are also given for repeats
1121 localised to the 5'UTR (light indigo), ATG-Stop (indigo) and 3'UTR (dark indigo) within each gene. Note that
1122 while all 17,741 genes contain a start-to-stop region, not all genes contain a 5'UTR and/or 3'UTR so the
1123 percentages presented for these are relative to totals of 17,721 and 17,717 respectively. Simple repeats
1124 are shown in the left three categories. Mono/di/tri refers to tandem repeats with a period length of one, two
1125 or three; tetra+ refers to all oligonucleotide tandem repeats with a period length of 4 or more and a total
1126 length ≥20 bp. Combining whole-gene counts for mono-, di-, tri- and tetra+ produces an average value of
1127 1.07 tandem repeats per kbp. Inverted repeats refer to short (20-210 bp) sequences that have the potential
1128 to form secondary structures by self-complementary base pairing. 836 genes were free from detectable
1129 tandem and inverted repeats under our criteria, most of which are small, with an average length of 1766
1130 bp. Global repeats refer to repetitive sequences masked by the NCBI WindowMasker program (Morgulis et
1131 al., 2006), which includes both longer, non-adjacent sequences and shorter, simple repeats (see Methods).
1132 All genes contained detectable repetitive regions using the default WindowMasker settings, with an average
1133 of 40.07 per gene. UTR data are based on gene models from Phytozome (version 5.5).
1134 **C** Gene features that complicate correct transgene expression. Top four bars illustrate potential
1135 misannotation of functional start sites in the genome shown by the percentage of genes containing one or
1136 more uORFs of each class (see text). Note that some genes contain multiple classes of uORF. Shown
1137 below this is the percentage of Chlamydomonas genes with multiple transcript models (splice variants), and
1138 those containing introns in the UTRs and translated regions (TR; between start and stop codons). uORF
1139 data is from Cross (2015). Splice variant and intron data are based on gene models from Phytozome
1140 (version 5.5).
1141 **D** Analysis of a set of ATG-Stop PCR primers designed to clone every gene in the genome from start to
1142 stop codon using gDNA as the template (Mackinder et al., 2017). Many primers are predicted to be
1143 unsuitable for efficient PCR, as shown by the percentage of forward (dark blue) and reverse (light blue)
1144 primers that breach various recommended thresholds associated with good primer design. Pairs (pale blue)
1145 are shown for which one or both primers breach the respective thresholds. Thresholds shown pertain to
1146 length, secondary structure stability, tandem repeats and 3' GC content. The inset shows the distribution
1147 of GC content of primers in the dataset, illustrating a clear trend in higher GC content at the 3' end of coding
1148 sequences. Below this, the given reason for rejection of primers by the Primer3 check_primers module is
1149 shown in orange. Dimer and hairpin values refer to primers rejected for 'high end complementarity' and
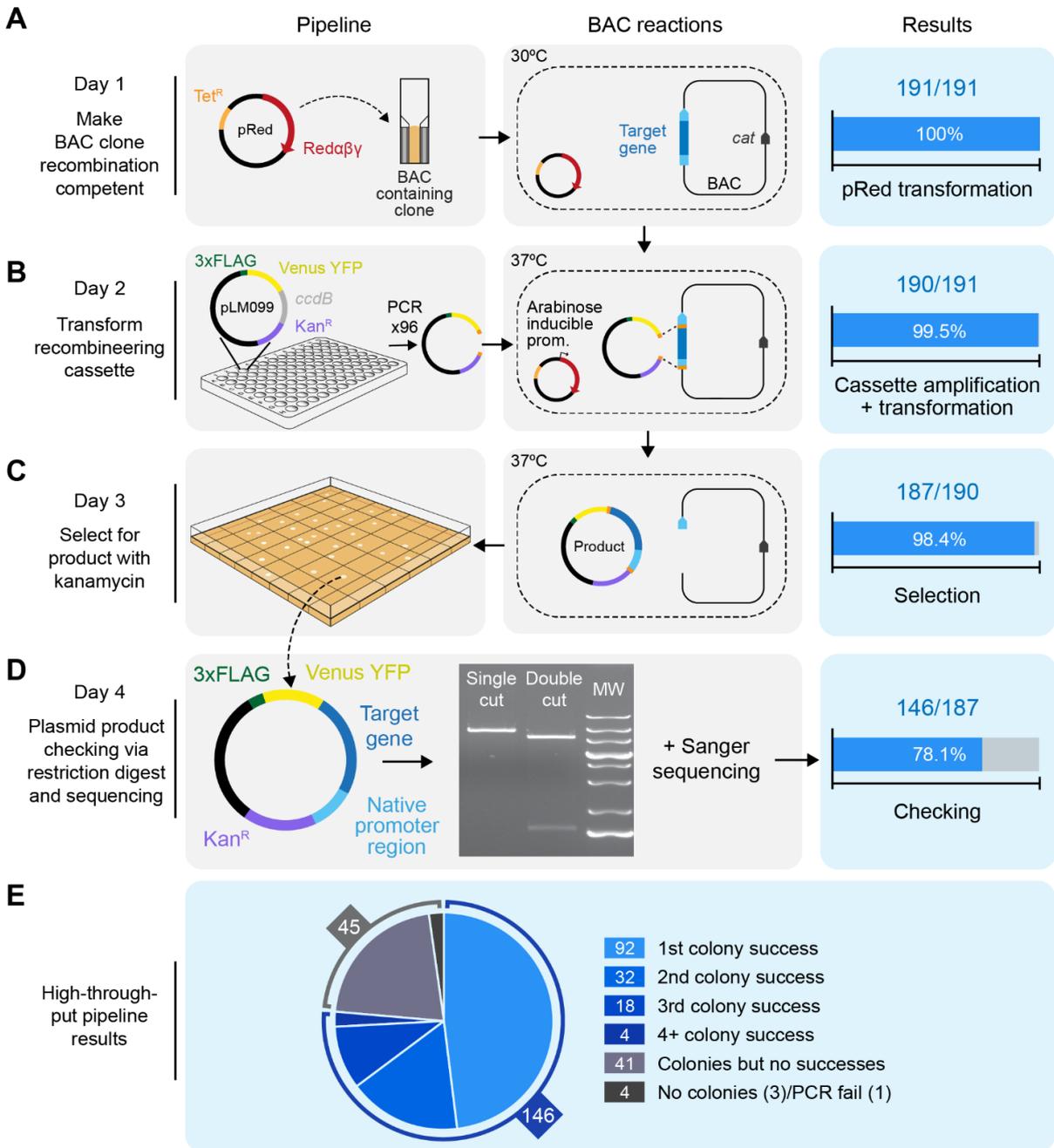1150 'high any complementarity' errors, respectively.
1151 **E** Annotated gene structure of Cre08.g379800. The gene encodes a predicted protein of unknown function
1152 but shows examples of several sequence features that contribute to sequence complexity. The unspliced
1153 sequence is 16,892 bases long with a GC content of 64.3%. The 41 exons are shown as regions of
1154 increased thickness, with 40 introns between them, the annotated 5'UTR in green and the 3'UTR in red.
1155 Labels denote selected examples of simple repeats throughout the gene. The inset shows the 5'UTR
1156 sequence, displaying examples of two classes of uORFs (see text); class 3 is highlighted in magenta and
1157 class 1 in green. For simplicity only one of the seven class 3 uORFs are shown in full. Cre08.g379800 was
1158 successfully cloned and tagged using recombineering.
1159 **F** A comparison of gene size and complexity between Chlamydomonas, bread wheat (*Triticum aestivum*),
1160 *Arabidopsis thaliana* and *Saccharomyces cerevisiae*. Gene sizes were binned as in **A**, and the average
1161 number of global repeats per kilobase (kbp) masked by the NCBI WindowMasker program was counted for
1162 genes in each size category (Morgulis et al., 2006). Genes were measured from the start of the 5'UTR to
1163 the end of the 3'UTR.
1164

**Figure 2**

1167 **Figure 2.** We developed a high-throughput recombineering pipeline for generating Venus-tagged fusion
1168 proteins with native promoter regions intact.
1169 **A** On day 1, BAC clones containing target genes are made recombineering competent by transformation
1170 with the pRed plasmid, which encodes the viral recombinogenic Redαβγ genes and *recA* under the control
1171 of an arabinose inducible promoter. Transformation efficiency shown on the right-hand side relates to BAC
1172 clones that yielded colonies after selection with tetracycline and chloramphenicol. *Cat*: the chloramphenicol
1173 resistance gene in the backbone of every BAC clone in the BAC library.
1174 **B** On or before day 2, the recombineering cassette is amplified from pLM099 using primers that contain 50
1175 bp homology arms complementary to regions flanking the target gene (shown in orange); one >2000 bp
1176 upstream of the annotated ATG and one at the 3' end of the coding sequence. On day 2, BAC-containing
1177 cells are electrotransformed with the recombineering cassette after induction with L-arabinose.
1178 Recombination between the BAC and the cassette results in a plasmid product containing the target gene
1179 in frame with CrVenus-3xFLAG and under its native promoter. Efficiency shown at this stage relates to PCR
1180 reactions that yielded efficient amplification of the recombineering cassette.
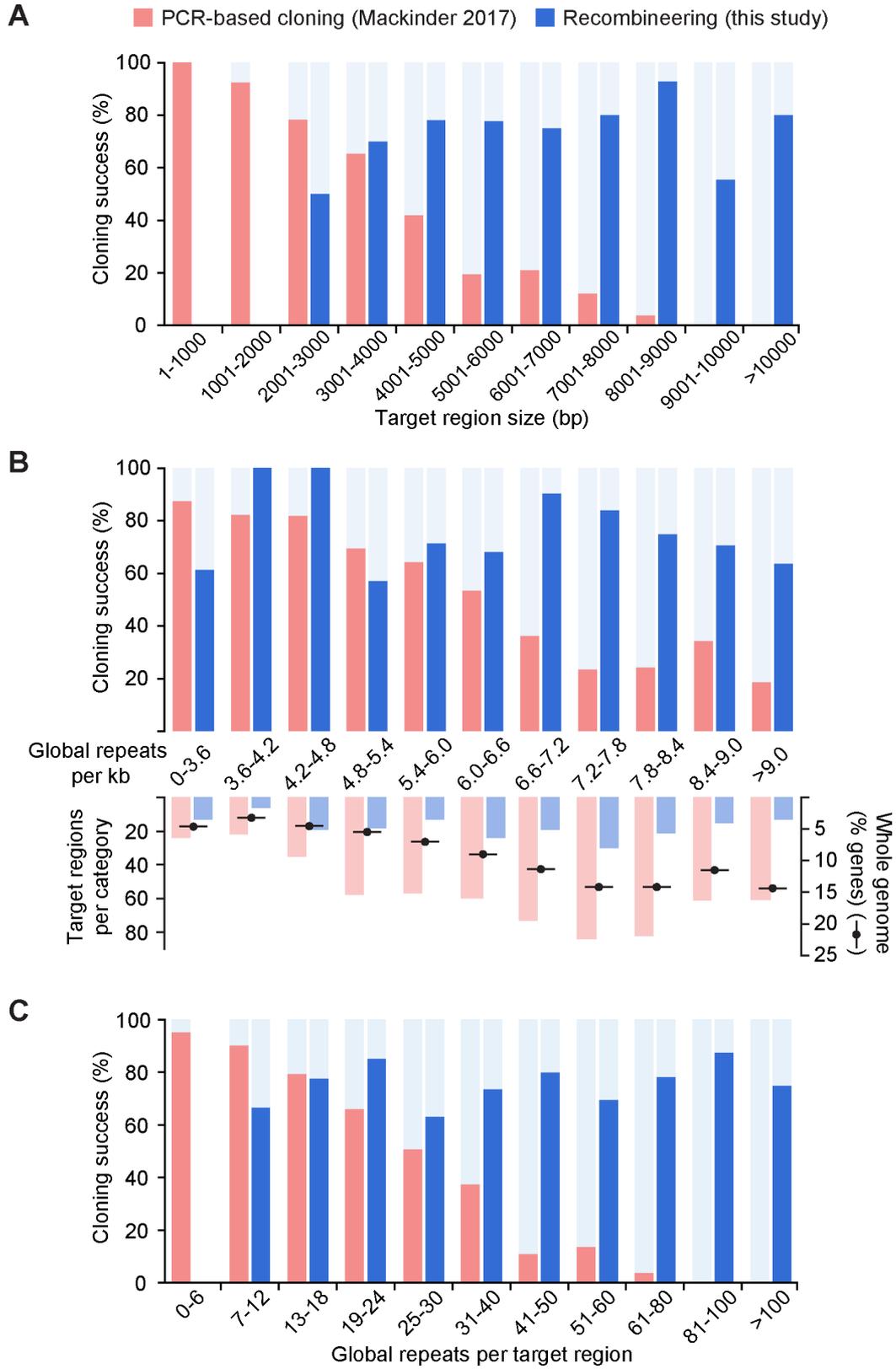1181 **C** On day 3, colonies containing plasmid products are isolated. Efficiency at this stage relates to the number
1182 of transformations that yielded colonies after selection with kanamycin.
1183 **D** On day 4, plasmid products are extracted from cells, screened by enzymatic digestion and confirmed by
1184 sequencing. Efficiency shown at this stage relates to correct digest patterns with single and double cutting
1185 restriction enzymes. MW: molecular weight marker.
1186 **E** Overall efficiency split into number of colonies screened via restriction digest. For 74% of target regions,
1187 the correct digest pattern was observed from plasmids isolated from the first, second or third colony picked
1188 per target. For 3% of targets, analysing >3 colonies yielded the correct product.
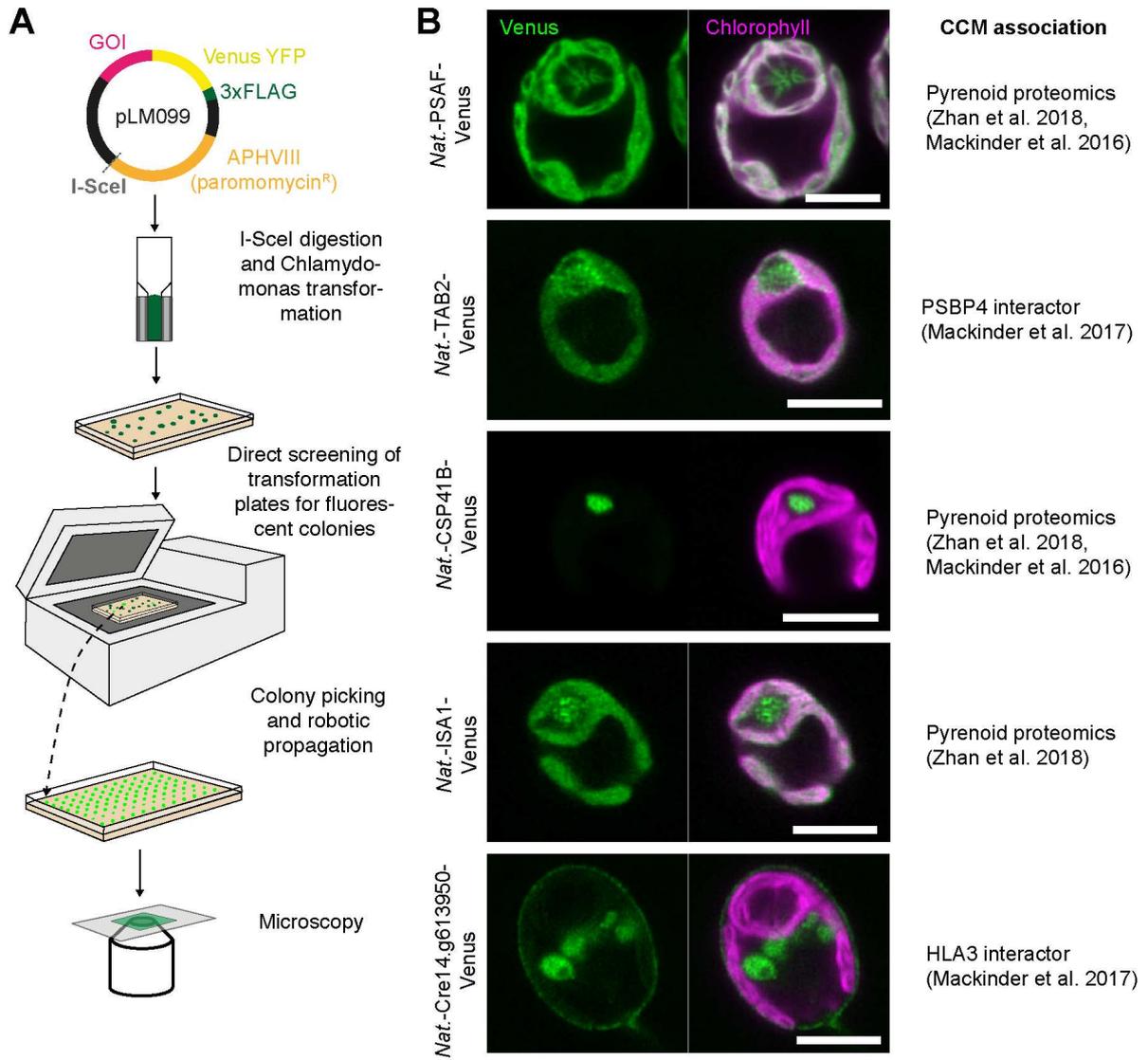1189

**Figure 3**

1193 **Figure 3.** Our recombineering pipeline is target gene size independent and tolerant of sequence complexity
1194 **A** The size distribution of successfully PCR-cloned coding sequences (Mackinder et al., 2017; red) or
1195 recombineered regions (this study; blue) are shown. Regions cloned by recombineering include ~2 kbp of
1196 flanking DNA upstream of the annotated start codon to incorporate native 5' promoter sequences. A severe
1197 drop in PCR-based cloning efficiency can be seen for templates >3 kbp long, whereas recombineering
1198 cloning efficiency does not show size dependency. No recombineering target regions were less than 2000
1199 bp long due to inclusion of native 5' promotor sequences.
1200 **B** As above but showing the dependence of cloning success on the per-kilobase frequency of repeats
1201 masked by the NCBI WindowMasker program with default settings (Morgulis et al., 2006). The number of
1202 target regions per repeat category is shown beneath this, overlaid with the percentage of Chlamydomonas
1203 genes in each category. The distribution of targets for this study and our previous PCR-based cloning
1204 attempt (Mackinder et al., 2017) gives a reasonably close representation of the whole genome distribution.
1205 Almost a third of nuclear genes contain 7.2-8.4 repeats per kbp; this peak corresponds to a clear drop in
1206 PCR-based cloning efficiency, but to a high recombineering efficiency of 75-85%. Data for repeats per kbp
1207 was continuous and there are no values present in more than one category.
1208 **C** As above but showing the number of simple and global repeats masked by WindowMasker per template.
1209 Data are binned to provide a higher resolution for the lower value categories, since the targets for PCR-
1210 based cloning were enriched in targets with low numbers of repeats. As in **A**, a severe negative trend in
1211 PCR-based cloning efficiency can be seen, reflecting a strong positive correlation between repeat number
1212 and region size. No negative association is present for recombineering cloning efficiency, likely illustrating
1213 the benefit of avoiding size- and complexity-associated polymerase limitations. No recombineering target
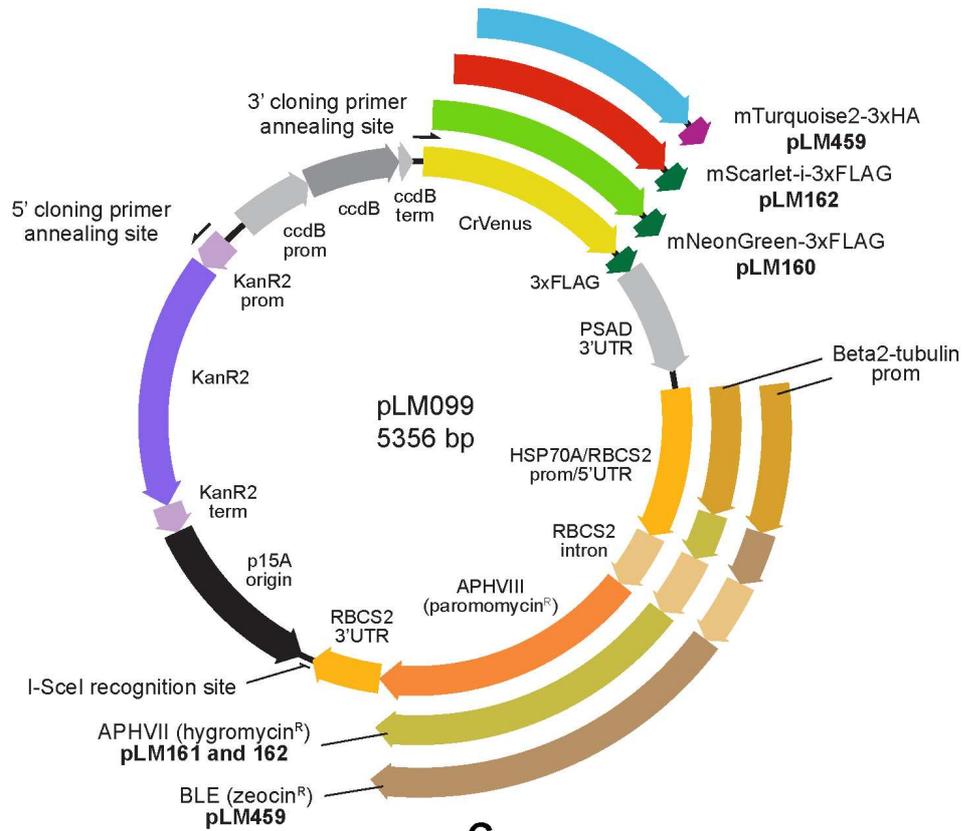1214 regions contained fewer than 6 repeats.
1215

1216    **Figure 4**
1217



A

pLM099

GOI
Venus YFP
3xFLAG
APHVIII (paromomycin^R)
I-SceI

I-SceI digestion and Chlamydomonas transformation

Direct screening of transformation plates for fluorescent colonies

Colony picking and robotic propagation

Microscopy

B

| Venus | Chlorophyll | **CCM association** |
|---|---|---|

*Nat.*-PSAF-Venus — Pyrenoid proteomics (Zhan et al. 2018, Mackinder et al. 2016)

*Nat.*-TAB2-Venus — PSBP4 interactor (Mackinder et al. 2017)

*Nat.*-CSP41B-Venus — Pyrenoid proteomics (Zhan et al. 2018, Mackinder et al. 2016)

*Nat.*-ISA1-Venus — Pyrenoid proteomics (Zhan et al. 2018)

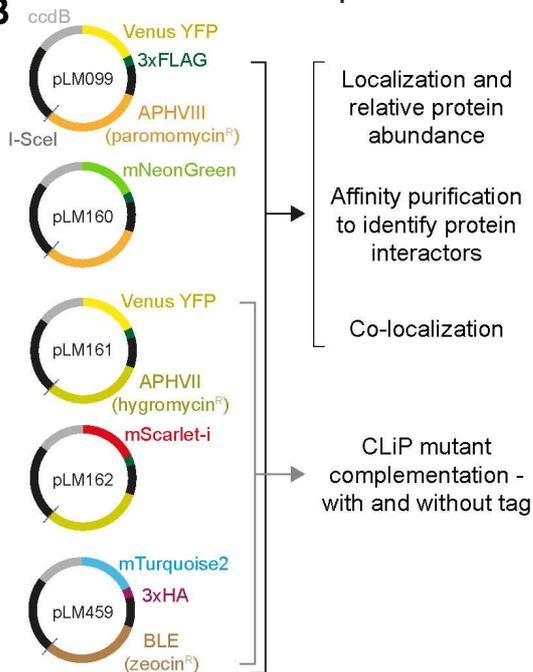*Nat.*-Cre14.g613950-Venus — HLA3 interactor (Mackinder et al. 2017)

1218   **Figure 4.** Transformation and localization of a subset of recombineered targets.
1219   **A** Chlamydomonas transformation pipeline. The I-SceI cut site allows vector linearization prior to
1220   Chlamydomonas transformation via electroporation. Transformants are directly screened for fluorescence
1221   using a Typhoon scanner (GE Healthcare) and then picked and propagated prior to imaging**.**
1222   **B** The localization for a subset of the recombineered target genes. Localizations agree with data from an
1223   affinity-purification followed by mass spectrometry study (Mackinder et al. 2017) or pyrenoid proteomics
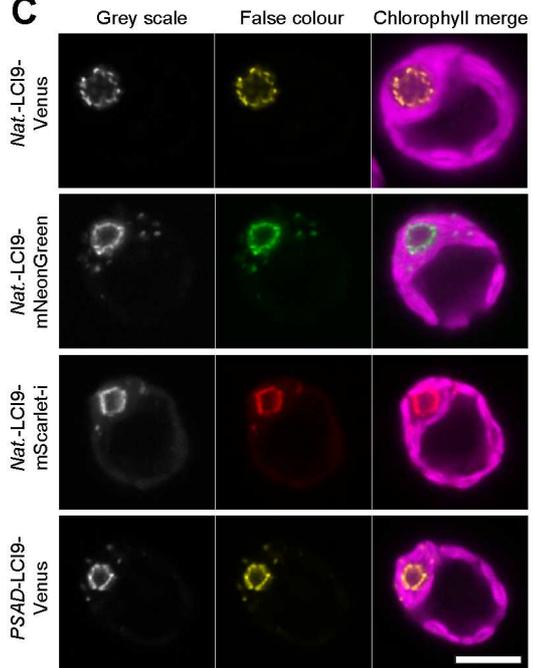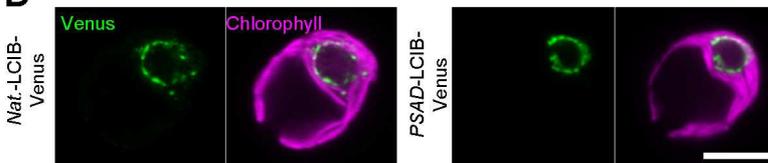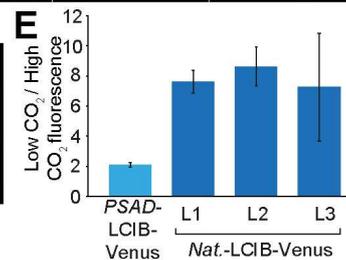1224   (Zhan et al. 2018 and/or Mackinder et al. 2016). Scale bars: 5 μm.
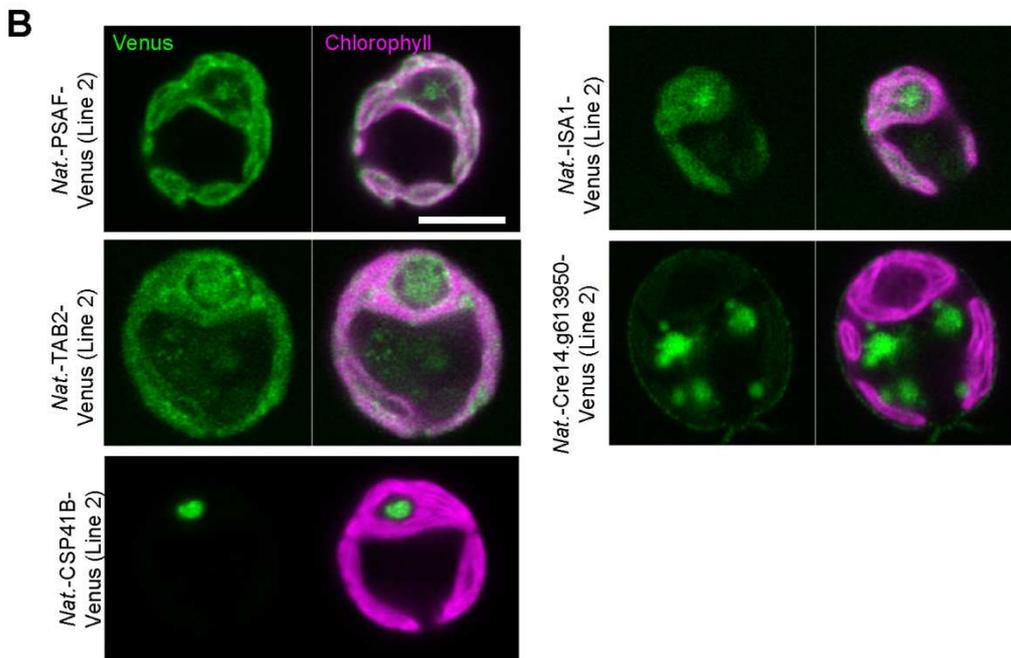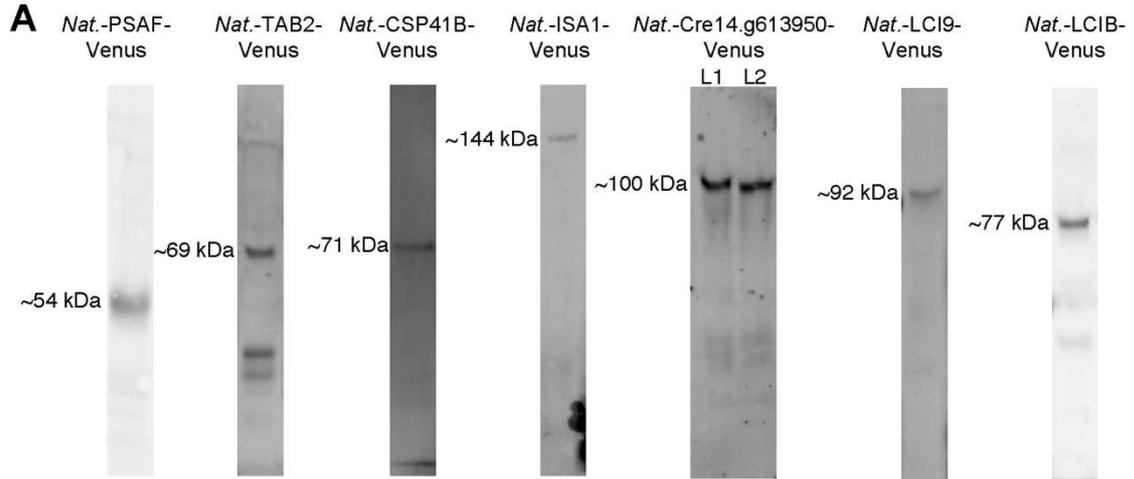1225

**Figure 5**

1227    **Figure 5.** Development and application of different recombineering vectors to enable novel biological
1228    insights into Chlamydomonas biology.
1229    **A** Plasmid map for pLM099 and derivative recombineering vectors. PCR amplification with 5' and 3' cloning
1230    primers at the annealing sites shown results in a ~4.6 kbp linear cassette for recombineering target genes
1231    in-frame with a fluorescent protein and affinity tag. For each recombineering vector, the fluorescent protein
1232    sequence is preceded by a flexible linker (GGLGGSGGR) and followed by a tri-glycine linker prior to the
1233    affinity tag. The PSAD 3'UTR terminates all four fluorescent protein-affinity tag cassettes. The RBCS2
1234    3'UTR terminates all three Chlamydomonas selection cassettes. The same RBCS2 intron is present in all
1235    three Chlamydomonas selection cassettes but is only inter-exonic in the hygromycin and zeocin resistance
1236    cassettes.
1237    **B** Additional vectors for tagging with different fluorophores and for complementation of Chlamydomonas
1238    library mutants generated using insertion of the *AphVIII* paromomycin resistant gene.
1239    **C** Localization of LCI9 with different fluorescence protein tags. *LCI9* was recombineered with its native
1240    promoter (*Nat.*) using pLM099, pLM160 and pLM162. A previously developed line cloned by PCR and using
1241    the constitutively expressed promotor *PSAD* is shown for comparison (*PSAD*-LCI9-Venus). Scale bar: 5
1242    µm.
1243    **D** A comparison of the low $CO_2$ upregulated gene LCIB cloned with its native promoter via recombineering
1244    vs LCIB under the constitutive *PSAD* promoter. Cells were grown and imaged at atmospheric $CO_2$ levels.
1245    Scale bar: 5 µm.
1246    **E** Relative change in LCIB-Venus fluorescence between high (3% v/v) and low (0.04% v/v) $CO_2$ when
1247    expressed from the constitutive *PSAD* promoter vs expression from the native LCIB promoter. Data is
1248    shown for three independent *native LCIB* promoter lines (L1-L3). Error bars are standard error of the mean.
1249

1250 **Figure S1**



**A**

Cre09.g394621    Cre13.g571700    Cre16.g663150

**B**

**C**

**D** Cre11.g467712 (SAGA1)

1251
1252

1253 **Figure S1.** Batch-scale recombineering results.
1254 **A** Three examples of colony PCRs to check for presence of target genes in BACs. Primer pairs were
1255 designed to the 5' and 3' end of each target gene. All amplicons were of the expected size.
1256 **B** Restriction digest checks for isolated recombineered plasmids from two colonies per gene, corresponding
1257 to the same genes as in **A**. Expected sizes are shown in bp. Note that colonies 1 and 2 for Cre09.g394621
1258 produced low-abundance bands in addition to the expected banding pattern that potentially correspond to
1259 incomplete digestion products. Colony 1 for Cre13.g571700 gave the incorrect size and banding patterns
1260 after digestion indicating incorrect recombination. L: GeneRuler 1 kb DNA Ladder (ThermoFisher Scientific).
1261 U: undigested.
1262 **C** Overall batch-scale recombineering success for 12 target genes.
1263 **D** Restriction digest checks of plasmids isolated from three colonies for *SAGA1* recombineered from fosmid
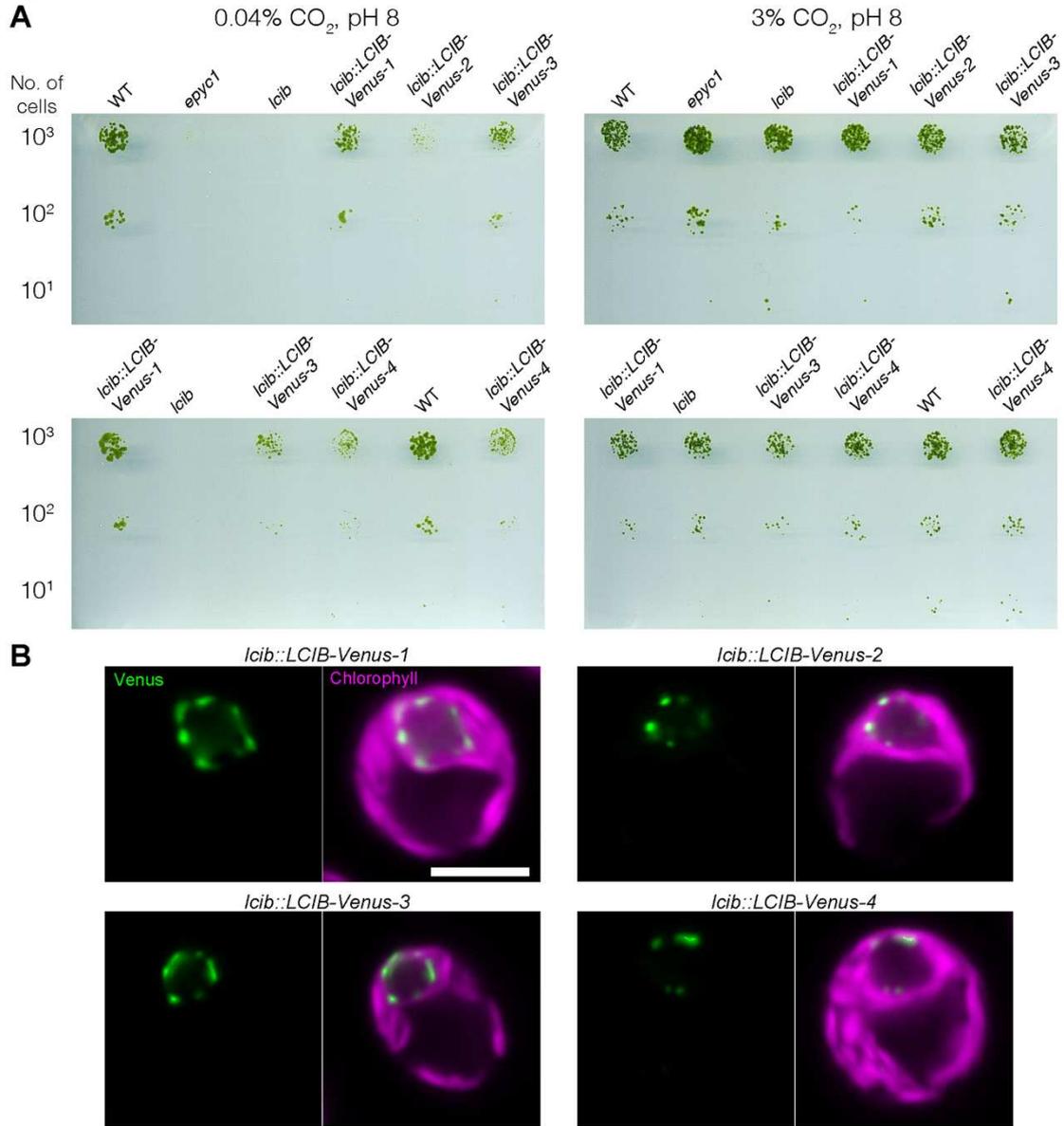1264 VTP41289 using pLM160.
1265 Supports Figure 2.
1266

1267    **Figure S2**

1270 **Figure S2.** Validation of fluorescently localized lines.
1271 **A** Immunoblots against the 3xFLAG epitope for recombineered targets. Molecular weights indicate the
1272 approximate band size. All cloned targets except Cre14.g613950 (expected molecular weight of 141 kDa)
1273 showed the expected molecular weight. Two independent transformants were tested for Cre14.g613950 to
1274 confirm that the observed lower molecular weight was consistent between transformants.
1275 Contrast/brightness were adjusted separately for each image.
1276 **B** Localization of target proteins in additional independent transformants (line 2). All localizations are
1277 consistent with line 1 localizations shown in Figure 4. Scale bar: 5 µm.
1278 Supports Figure 4 and Figure 5.
1279

1280    **Figure S3**



A

0.04% CO$_2$, pH 8                    3% CO$_2$, pH 8

B

*lcib::LCIB-Venus-1*              *lcib::LCIB-Venus-2*

Venus    Chlorophyll

*lcib::LCIB-Venus-3*              *lcib::LCIB-Venus-4*

1281
1282

1283    **Figure S3.** Complementation of the *lcib* CLiP mutant.

1284    **A** Spot tests of *lcib* CLiP mutant LMJ.RY0402.215132 complemented with recombineered *LCIB-Venus*

1285    driven by its native promoter. Four independent transformants we spotted onto pH 8 TP minimal media

1286    plates and grown at 0.04% and 3% $CO_2$. The *epyc1* mutant that has a severe CCM phenotype due to

1287    incorrect pyrenoid assembly was included as a CCM growth phenotype control. Note varying degrees of

1288    complementation between lines. Top and bottom images for each $CO_2$ condition are from the same plate

1289    but split for labelling clarity.

1290    **B** Corresponding confocal microscope images of complemented lines all showing the typical localization of

1291    LCIB at the pyrenoid periphery. Scale bar: 5 µm.

1292    Supports Figure 5.

1293