

1 **A genome-wide algal mutant library reveals a global**
2 **view of genes required for eukaryotic photosynthesis**

3

4 **Xiaobo Li^{1,2,†}, Weronika Patena^{1,2}, Friedrich Fauser^{1,2}, Robert E. Jinkerson^{2,†}, Shai**
5 **Saroussi², Moritz T. Meyer¹, Nina Ivanova², Jacob M. Robertson^{1,2}, Rebecca Yue²,**
6 **Ru Zhang^{2,†}, Josep Vilarrasa-Blasi², Tyler M. Wittkopp^{2,3,†}, Silvia Ramundo⁴, Sean**
7 **R. Blum², Audrey Goh¹, Matthew Laudon⁵, Tharan Srikumar¹, Paul A. Lefebvre⁵,**
8 **Arthur R. Grossman², and Martin C. Jonikas^{1,2*}**

9

10 ORCID IDs: 0000-0003-3951-9646 (X.L.); 0000-0001-8516-2591 (M.T.M.); 0000-0002-4860-
11 7800 (R.Z.); 0000-0003-4328-4798 (J.V.-B.); 0000-0001-7061-0611 (T.M.W.); 0000-0002-3747-
12 5881 (A.R.G.); 0000-0002-9519-6055 (M.C.J.)

13 ¹Department of Molecular Biology, Princeton University, Princeton, New Jersey 08544, USA.

14 ²Department of Plant Biology, Carnegie Institution for Science, Stanford, California 94305, USA.

15 ³Department of Biology, Stanford University, Stanford, California 94305, USA. ⁴Department of

16 Biochemistry and Biophysics, University of California, San Francisco, California 94158, USA.

17 ⁵Department of Plant and Microbial Biology, University of Minnesota, St. Paul, Minnesota

18 55108, USA. [†]Present addresses: School of Life Sciences, Westlake University, Hangzhou,

19 Zhejiang Province 310064, China (X.L.); Department of Chemical and Environmental

20 Engineering, University of California, Riverside, California 92521, USA (R.E.J); Donald

21 Danforth Plant Science Center, 975 North Warson Road, St. Louis, Missouri 63132, USA (R.Z.);

22 Salk Institute for Biological Studies, La Jolla, California 92037, USA (T.M.W.). *e-mail:

23 mjonikas@princeton.edu

1
2 **Photosynthetic organisms provide food and energy for nearly all life on Earth, yet**
3 **half of their protein-coding genes remain uncharacterized^{1,2}. Characterization of**
4 **these genes could be greatly accelerated by new genetic resources for unicellular**
5 **organisms that complement the use of multicellular plants by enabling higher-**
6 **throughput studies. Here, we generated a genome-wide, indexed library of mapped**
7 **insertion mutants for the flagship unicellular alga *Chlamydomonas reinhardtii***
8 **(*Chlamydomonas* hereafter). The 62,389 mutants in the library, covering 83% of**
9 **nuclear, protein-coding genes, are available to the community. Each mutant**
10 **contains unique DNA barcodes, allowing the collection to be screened as a pool. We**
11 **leveraged this feature to perform a genome-wide survey of genes required for**
12 **photosynthesis, which identified 303 candidate genes. Characterization of one of**
13 **these genes, the conserved predicted phosphatase *CPL3*, showed it is important for**
14 **accumulation of multiple photosynthetic protein complexes. Strikingly, 21 of the 43**
15 **highest-confidence genes are novel, opening new opportunities for advances in our**
16 **understanding of this biogeochemically fundamental process. This library is the first**
17 **genome-wide mapped mutant resource in any unicellular photosynthetic organism,**
18 **and will accelerate the characterization of thousands of genes in algae, plants and**
19 **animals.**

20 Among unicellular photosynthetic organisms, the green alga *Chlamydomonas* has
21 long been employed for genetic studies of eukaryotic photosynthesis because of its rare
22 ability to grow in the absence of photosynthetic function³. In addition, it has made
23 extensive contributions to our basic understanding of light signaling, stress acclimation,
24 and metabolism of carbohydrates, lipids, and pigments (Fig. 1a)⁴⁻⁶. Moreover,

1 Chlamydomonas retained many genes from the plant-animal common ancestor, which
2 allowed it to reveal fundamental aspects of the structure and function of cilia and basal
3 bodies^{7,8}. Like *Saccharomyces cerevisiae*, Chlamydomonas can grow as a haploid,
4 facilitating genetic studies. However, until now, the value of Chlamydomonas has been
5 limited by the lack of mutants in most of its nuclear genes.

6 In the present study, we sought to generate a genome-wide collection of
7 Chlamydomonas mutants with known gene disruptions to provide mutants in genes of
8 interest for the scientific community, and then to leverage this collection to reveal genes
9 with roles in photosynthesis. To reach the necessary scale, we chose to use random
10 insertional mutagenesis and built on advances in insertion mapping and mutant
11 propagation from our pilot study⁹. To enable mapping of insertion sites and screening
12 pools of mutants on a much larger scale, we developed new tools leveraging unique DNA
13 barcodes in each transforming cassette.

14 We generated mutants by transforming haploid cells with DNA cassettes that
15 randomly insert into the genome and inactivate the genes they insert into. We maintained
16 the mutants as indexed colony arrays on agar media containing acetate as a carbon and
17 energy source to allow recovery of mutants with defects in photosynthesis. Each DNA
18 cassette contained two unique barcodes, one on each side of the cassette (Supplementary
19 Fig. 1a-d). For each mutant, the barcode and genomic flanking sequences on each side of
20 the cassette were initially unknown (Supplementary Fig. 1e). We determined the
21 sequence of the barcode(s) in each mutant colony by combinatorial pooling and deep
22 sequencing (Supplementary Fig. 1f). We then mapped each insertion by pooling all
23 mutants and amplifying all flanking sequences together with their corresponding

1 barcodes followed by deep sequencing (Supplementary Fig. 1g). The combination of
2 these datasets revealed the insertion site(s) in each mutant. This procedure yielded 62,389
3 mutants on 245 plates, with a total of 74,923 insertions that were largely randomly
4 distributed over the chromosomes (Fig. 1, b and c, and Supplementary Table 5).

5 This library provides mutants for ~83% of all nuclear genes (Fig. 2a-d).
6 Approximately 69% of genes are represented by an insertion in a 5' UTR, an exon or an
7 intron – regions most likely to cause an altered phenotype when disrupted. Many gene
8 sets of interest to the research community are well represented, including genes encoding
9 proteins phylogenetically associated with the plant lineage (GreenCut2)¹, proteins that
10 localize to the chloroplast¹⁰, or those associated with the structure and function of flagella
11 or basal bodies^{11,12} (Fig. 2b). Mutants in this collection are available through the website
12 <https://www.chlamylibrary.org/>. Over 1,800 mutants have already been distributed to
13 over 200 laboratories worldwide in the first 18 months of pre-publication distribution
14 (Fig. 2e). These mutants are facilitating genetic investigation of a broad range of
15 processes, ranging from photosynthesis and metabolism to cilia structure and function
16 (Fig. 2f).

17 To identify genes required for photosynthesis, we screened our library for mutants
18 deficient in photosynthetic growth. Rather than phenotyping each strain individually, we
19 pooled the entire library into one culture and leveraged the unique barcodes present in
20 each strain to track its abundance after growth under different conditions. This feature
21 enables genome-wide screens with speed and depth unprecedented in photosynthetic
22 eukaryotes. We grew a pool of mutants photosynthetically in light in minimal Tris-
23 Phosphate (TP) medium with CO₂ as the sole source of carbon, and heterotrophically in

1 the dark in Tris-Acetate-Phosphate (TAP) medium, where acetate provides fixed carbon
2 and energy³ (Fig. 3a). To quantify mutant growth under each condition, we amplified and
3 deep sequenced the barcodes from the final cell populations. We then compared the
4 ability of each mutant to grow under photosynthetic and heterotrophic conditions by
5 comparing the read counts of each barcode from each condition (Supplementary Table
6 10; Methods). Mutant phenotypes were highly reproducible (Fig. 3b and Supplementary
7 Fig. 5, a and b). We identified 3,109 mutants deficient in photosynthetic growth (Fig. 3c
8 and Methods).

9 To identify genes with roles in photosynthesis, we developed a statistical analysis
10 framework that leverages the presence of multiple alleles for many genes. This
11 framework allows us to overcome several sources of false positives that have been
12 difficult to identify with previous methods, including cases where the phenotype is not
13 caused by the mapped disruption. For each gene, we counted the number of mutant
14 alleles with and without a phenotype, and evaluated the likelihood of obtaining these
15 numbers by chance given the total number of mutants in the library that exhibit the
16 phenotype (Supplementary Table 11; Methods).

17 We identified 303 candidate photosynthesis genes based on our statistical analysis
18 above. These genes are enriched for membership in a diurnally regulated photosynthesis-
19 related transcriptional cluster¹³ ($P < 10^{-11}$), are enriched for upregulation upon dark-to-light
20 transitions¹⁴ ($P < 0.003$), and encode proteins enriched for predicted chloroplast
21 localization ($P < 10^{-8}$). As expected¹⁵, the candidate genes also encode a disproportionate
22 number of GreenCut2 proteins ($P < 10^{-8}$), which are conserved among photosynthetic

1 organisms but absent from non-photosynthetic organisms¹: 32 GreenCut2 proteins are
2 encoded by the 303 candidate genes (11%), compared to ~3% in the entire genome.

3 Photosynthesis occurs in two stages: the light reactions and carbon fixation. The
4 light reactions convert solar energy into chemical energy, and require coordinated action
5 of Photosystem II (PSII), Cytochrome *b₆f*, Photosystem I (PSI), ATP synthase
6 complexes, a plastocyanin or cytochrome *c₆* metalloprotein, as well as small molecule
7 cofactors¹⁶. PSII and PSI are each assisted by peripheral light-harvesting complexes
8 (LHCs) known as LHCII and LHCI, respectively. Carbon fixation is performed by
9 enzymes in the Calvin-Benson-Bassham cycle, including the CO₂-fixing enzyme
10 Rubisco. In addition, most eukaryotic algae have a mechanism to concentrate CO₂ around
11 Rubisco to enhance its activity¹⁷.

12 Sixty-five of the genes we identified encode proteins that were previously shown
13 to play a role in photosynthesis or chloroplast function in *Chlamydomonas* or vascular
14 plants (Fig. 3f). These include three PSII-LHCII subunits (PSBP1, PSBP2, and PSB27)
15 and seven PSII-LHCII biogenesis factors (CGL54, CPLD10, HCF136, LPA1, MBB1,
16 TBC2, and Cre02.g105650), two cytochrome *b₆f* complex subunits (PETC and PETM)
17 and six cytochrome *b₆f* biogenesis factors (CCB2, CCS5, CPLD43, CPLD49, MCD1, and
18 MCG1), five PSI-LHCI subunits (LHCA3, LHCA7, PSAD, PSAE, and PSAL) and nine
19 PSI-LHCI biogenesis factors (CGL71, CPLD46, OPR120, RAA1, RAA2, RAA3, RAT2,
20 Cre01.g045902, and Cre09.g389615), one protein required for ATP synthase function
21 (PHT3), plastocyanin (PCY1) and two plastocyanin biogenesis factors (CTP2 and PCC1),
22 12 proteins involved in the metabolism of photosynthesis cofactors or signaling
23 molecules (CHLD, CTH1, CYP745A1, DVR1, HMOX1, HPD2, MTF1, PLAP6,

1 UROD3, Cre08.g358538, Cre13.g581850, and Cre16.g659050), three Calvin-Benson-
2 Bassham Cycle enzymes (FBP1, PRK1, and SEBP1), two Rubisco biogenesis factors
3 (MRL1 and RMT2), three proteins involved in the algal carbon concentrating mechanism
4 (CAH3, CAS1, and LCIB), as well as proteins that play a role in photorespiration
5 (GSF1), CO₂ regulation of photosynthesis (Cre02.g146851), chloroplast morphogenesis
6 (Cre14.g616600), chloroplast protein import (SDR17), and chloroplast DNA, RNA, and
7 protein metabolism (DEG9, MSH1, MSRA1, TSM2, and Cre01.g010864) (Fig. 3h and
8 Supplementary Table 12). We caution that not all genes previously demonstrated to be
9 required for photosynthetic growth are detectable by this approach, especially the ones
10 with paralogous genes in the genome, such as *RBCS1* and *RBCS2* that encode the small
11 subunit of Rubisco¹⁸. Nonetheless, the large number of known factors recovered in our
12 screen is a testament to the power of this approach.

13 In addition to recovering these 65 genes with known roles in photosynthesis, our
14 analysis revealed 238 candidate genes with no previously reported role in photosynthesis
15 (Methods). These 238 genes represent a rich set of targets to better understand
16 photosynthesis. Because our screen likely yielded some false positives, we divided all
17 genes into “higher-confidence” ($P < 0.0011$; $FDR < 0.27$) and “lower-confidence” genes
18 based on the number of alleles that supported each gene’s involvement in photosynthesis
19 (Fig. 3d-f; Tables 1 and 2; Methods). The 21 higher-confidence genes with no previously
20 reported role in photosynthesis are enriched in chloroplast localization (9/21, $P < 0.011$;
21 Fig. 3g) and transcriptional upregulation during dark to light transition (5/21, $P < 0.005$),
22 similar to the known photosynthesis genes. Thus, these 21 higher-confidence genes are
23 particularly high-priority targets for the field to pursue.

1 Functional annotations for 15 of the 21 higher-confidence genes suggest that these
2 genes could play roles in regulation of photosynthesis, photosynthetic metabolism, and
3 biosynthesis of the photosynthetic machinery. Seven of the genes likely play roles in
4 regulation of photosynthesis: *GEF1* encodes a voltage-gated channel, Cre01.g008550 and
5 Cre02.g111550 encode putative protein kinases, *CPL3* encodes a predicted protein
6 phosphatase, *TRX21* contains a thioredoxin domain, Cre12.g542569 encodes a putative
7 glutamate receptor, and Cre13.g586750 contains a predicted nuclear importin domain.
8 Six of the genes are likely involved in photosynthetic metabolism: the Arabidopsis
9 homolog of Cre10.g448950 modulates sucrose and starch accumulation¹⁹,
10 Cre11.g467712 contains a starch-binding domain, Cre02.g073900 encodes a putative
11 carotenoid dioxygenase, *VTE5* encodes a putative phosphatidate cytidyltransferase,
12 Cre10.g429650 encodes a putative alpha/beta hydrolase, and Cre50.g761497 contains a
13 magnesium transporter domain. Finally, two of the genes are likely to play roles in the
14 biogenesis and function of photosynthesis machinery: *EIF2* has a translation initiation
15 factor domain, and *CDJ2* has a chloroplast DnaJ domain. Future characterization of these
16 genes by the community is likely to yield fundamental insights into our understanding of
17 photosynthesis.

18 As an illustration of the value of genes identified in this screen, we sought to
19 explore the specific function of one of the novel higher-confidence hits, *CPL3*
20 (*Conserved in Plant Lineage 3*, Cre03.g185200), which encodes a putative protein
21 phosphatase (Fig. 4a and Supplementary Fig. 6e). Many proteins in the photosynthetic
22 apparatus are phosphorylated, but the role and regulation of these phosphorylations are
23 poorly understood²⁰. In our screen, three mutants in *CPL3* exhibited a deficiency in

1 photosynthetic growth (Fig. 3c and Supplementary Table 13). We chose to examine one
2 allele (LMJ.RY0402.153647, referred to hereafter as *cpl3*; Fig. 4a and Supplementary
3 Fig. 6a) for phenotypic confirmation, genetic complementation, and further studies.

4 Consistent with the pooled growth data, *cpl3* showed a severe defect in
5 photosynthetic growth on agar, which was rescued under heterotrophic conditions (Fig.
6 4b). We confirmed that the *CPL3* gene is disrupted in the *cpl3* mutant and found that
7 complementation with a wild-type copy of the *CPL3* gene rescues the phenotype,
8 demonstrating that the mutation in *CPL3* is the cause of the growth defect of the mutant
9 (Supplementary Note and Supplementary Fig. 6a-d).

10 We then examined the photosynthetic performance, morphology of the
11 chloroplast, and the composition of photosynthetic pigments and proteins in *cpl3*.
12 Photosynthetic electron transport rate was decreased under all light intensities, suggesting
13 a defect in the photosynthetic machinery (Fig. 4c). The chloroplast morphology of *cpl3*
14 appeared similar to the wild type based on chlorophyll fluorescence microscopy
15 (Supplementary Fig. 7a). However, we observed a lower chlorophyll *a/b* ratio in *cpl3*
16 than in the wild type (Supplementary Fig. 7b), which suggests a defect in the
17 accumulation or composition of the protein-pigment complexes involved in the light
18 reactions²¹. Using whole-cell proteomics, we found that *cpl3* was deficient in
19 accumulation of all detectable subunits of the chloroplast ATP synthase (ATPC, ATPD,
20 ATPG, AtpA, AtpB, AtpE, AtpF), some subunits of PSII (D1, D2, CP43, CP47, PsbE,
21 PsbH), and some subunits of PSI (PsaA and PsaB) (FDR<0.31 for each subunit, Fig. 4d,
22 Fig. 4f, and Supplementary Table 14). We confirmed these findings by western blots on
23 CP43, PsaA, and ATPC (Fig. 4e). Our results indicate that *CPL3* is required for normal

1 accumulation of thylakoid protein complexes (PSII, PSI, and ATP synthase) involved in
2 the light reactions of photosynthesis.

3 Our finding that 21/43 of the higher-confidence photosynthesis hit genes were
4 uncharacterized suggests that nearly half of the genes required for photosynthesis remain
5 to be characterized. This finding is remarkable, considering that genetic studies on
6 photosynthesis extend back to the 1950s²². Our validation of *CPL3*'s role in
7 photosynthesis illustrates the value of the uncharacterized hit genes identified in this
8 study as a rich set of candidates for the community to pursue.

9 More broadly, it is our hope that the mutant resource presented here will serve as
10 a powerful complement to newly developed gene editing techniques²³⁻²⁸, and that
11 together these tools will help the research community generate fundamental insights in a
12 wide range of fields, from organelle biogenesis and function to organism-environment
13 interactions.

14

15 **Acknowledgments**

16 We thank Olivier Vallon for helpful discussions; Matthew Cahn and Garret Huntress for developing and
17 improving the CLiP website; Xuhuai Ji at the Stanford Functional Genomics Facility and Ziming Weng at
18 the Stanford Center for Genomics and Personalized Medicine for deep sequencing services; Alan Itakura
19 for help in library pooling; Shriya Ghosh, Kyssia Mendoza, and Matthew LaVoie for technical assistance;
20 Kathryn Barton, Winslow Briggs, and Zhi-Yong Wang for providing lab space; Joseph Ecker, Liz Freeman
21 Rosenzweig and Moshe Kafri for constructive suggestions on the manuscript; and the Princeton Mass
22 Spectrometry Facility for proteomics services. This project was supported by a grant from the National
23 Science Foundation (MCB-1146621) awarded to M.C.J. and A.R.G., grants from the National Institutes of
24 Health (DP2-GM-119137) and the Simons Foundation and HHMI (55108535) awarded to M.C.J., a

1 German Academic Exchange Service (DAAD) research fellowship to F.F., Simons Foundation fellowships
2 of the Life Sciences Research Foundation to R.E.J. and J.V.-B., EMBO long term fellowship (ALTF 1450-
3 2014 and ALTF 563-2013) to J.V.-B and S.R., and a Swiss National Science Foundation Advanced
4 PostDoc Mobility Fellowship (P2GEP3_148531) to S.R.

5

6 **Author contributions**

7 X.L. developed the method for generating barcoded cassettes; R.Y. and S.R.B. optimized the mutant
8 generation protocol; R.Y., N.I., and X.L. generated the library; J.M.R., N.I., A.G., and R.Y. maintained,
9 consolidated, and cryopreserved the library; X.L. developed the barcode sequencing method; N.I., X.L.,
10 R.Y., and W.P. performed combinatorial pooling and super-pool barcode sequencing; X.L. performed
11 LEAP-Seq; W.P. developed mutant mapping data analysis pipeline and performed data analyses of barcode
12 sequencing and LEAP-Seq; W.P. analyzed insertion coverage and hot/cold spots; R.Z. and J.M.R.
13 performed insertion verification PCRs and Southern blots; F.F., R.E.J., and J.V.-B. developed the library
14 screening protocol; F.F., J.V.-B., and X.L. performed the photosynthesis mutant screen and barcode
15 sequencing; R.E.J. and W.P. developed screen data analysis methods and implemented them for the
16 photosynthesis screen; X.L. and T.M.W. annotated the hits from the photosynthesis screen; X.L., J.M.R.,
17 and S.R. performed growth analysis, molecular characterizations, and complementation of *cpl3*; S.S. and
18 T.M.W. performed physiological characterizations of *cpl3*; M.T.M. and S.S. performed western blots on
19 the photosynthetic protein complexes; M.T.M. performed microscopy on *cpl3*; X.L., W.P., and T.S.
20 performed proteomic analyses; M.L. and P.A.L. maintained, cryopreserved, and distributed mutants at the
21 Chlamydomonas Resource Center; X.L., W.P., A.R.G., and M.C.J. wrote the manuscript with input from
22 all authors; M.C.J. and A.R.G. conceived and guided the research and obtained funding.

23

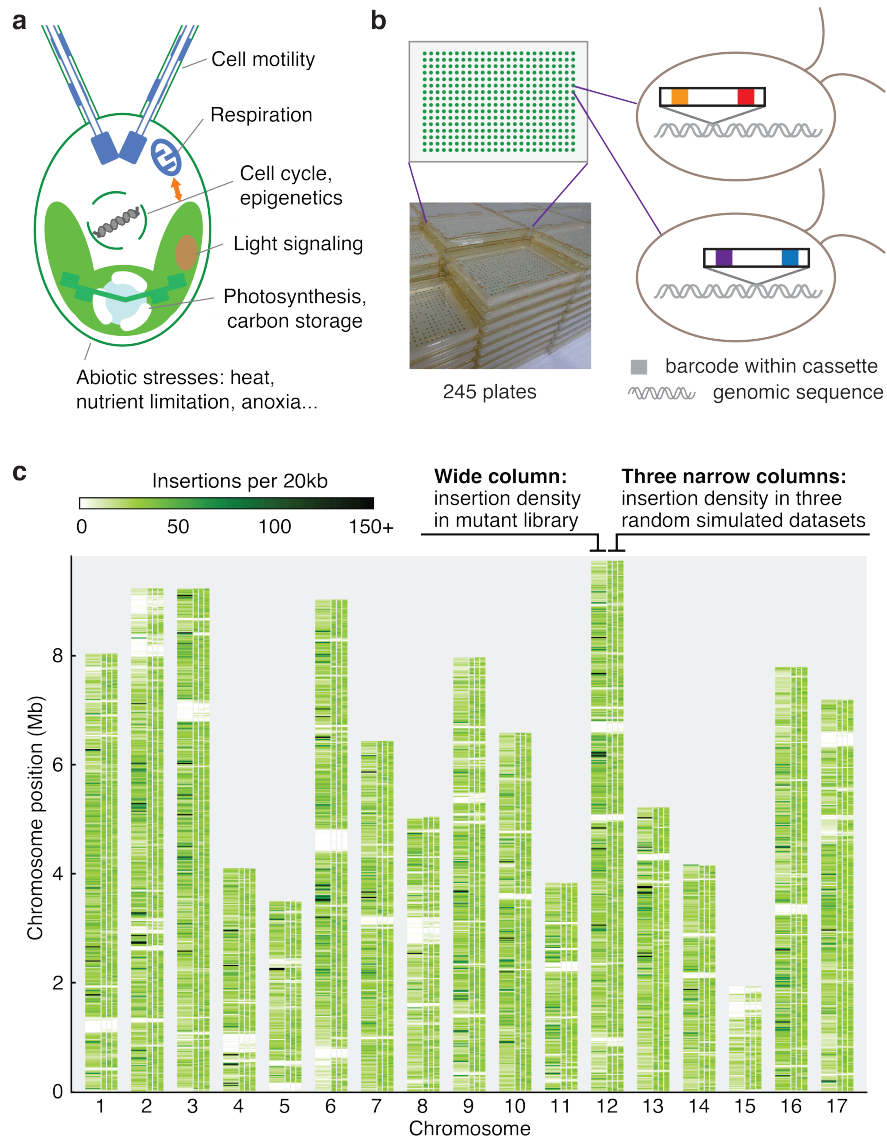
24 **Competing interests**

25 The authors declare no competing interests.

26

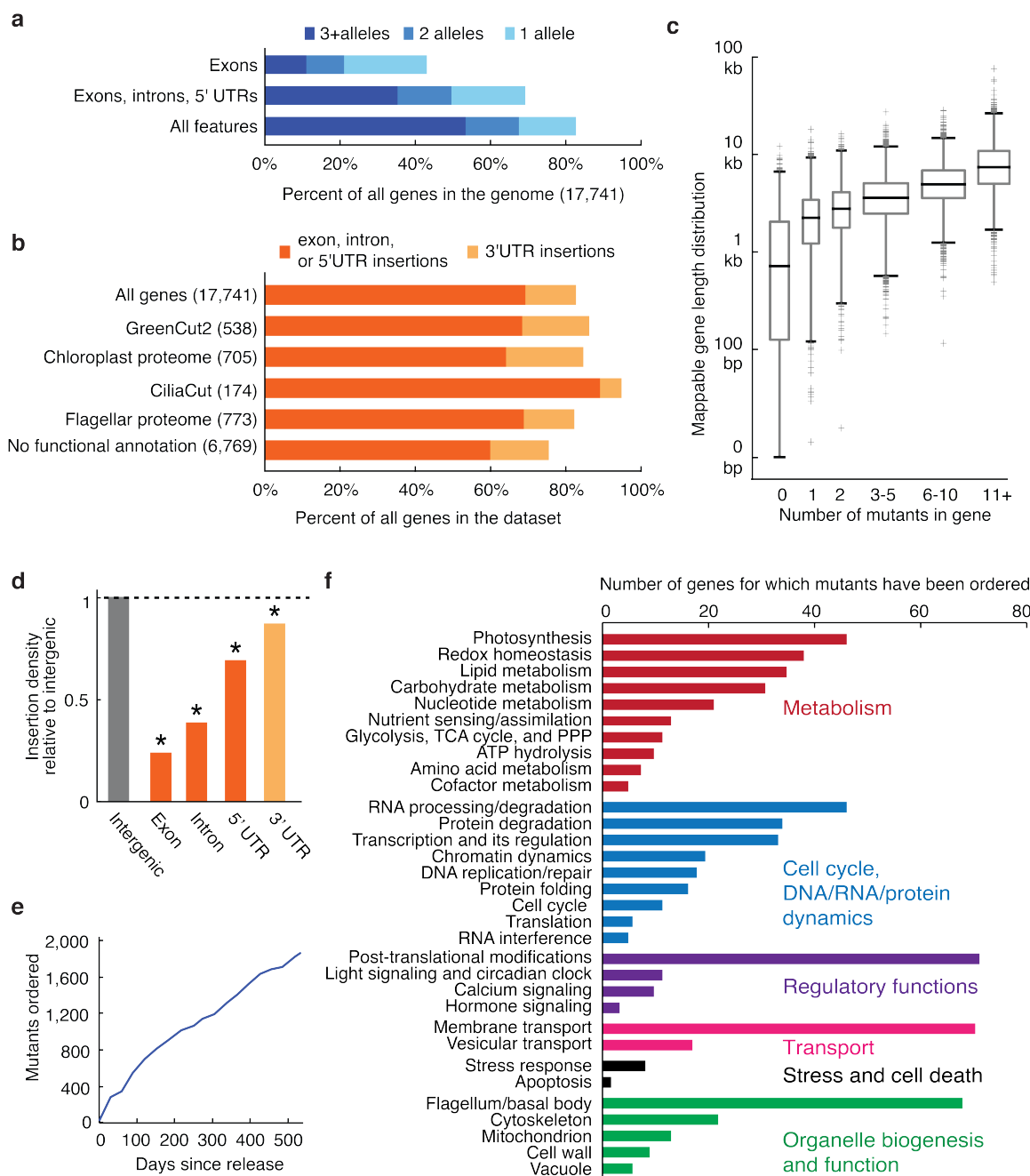
- 1 **Supplementary Information:**
- 2 Methods
- 3 Supplementary Note
- 4 Supplementary Figures 1-7
- 5 Supplementary Tables 1-14 (as separate excel or text files)
- 6

1 Figures



2
3 **Fig. 1 | A genome-wide library of *Chlamydomonas* mutants was generated by**
4 **random insertion of barcoded cassettes and mapping of insertion sites. a,**
5 *Chlamydomonas reinhardtii* is used for studies of various cellular processes and
6 organism-environment interactions. **b,** Our library contains 62,389 insertional mutants
7 maintained as 245 plates of 384-colony arrays. Each mutant contains at least one
8 insertion cassette at a random site in its genome; each insertion cassette contains one

1 unique barcode at each end (Supplementary Fig. 1a-c). **c**, The insertion density is largely
2 random over the majority of the genome. This panel compares the observed insertion
3 density over the genome (the left column above each chromosome number) to three
4 simulations with insertions randomly distributed over all mappable positions in the
5 genome (the three narrow columns to the right for each chromosome). Areas that are
6 white throughout all columns represent regions where insertions cannot be mapped to a
7 unique genomic position due to highly repetitive sequence. See also Supplementary Fig.
8 4.



1

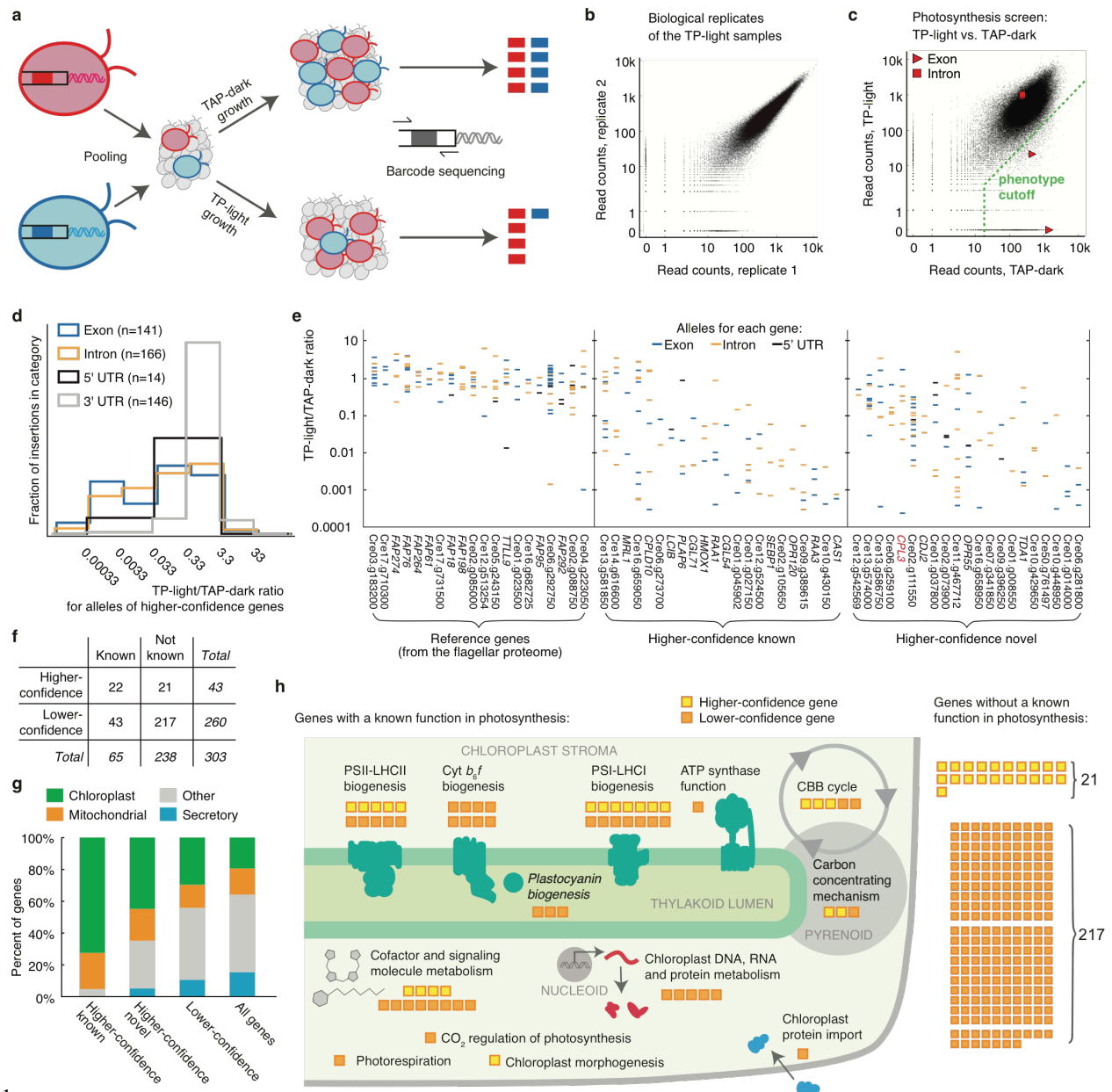
2 **Fig. 2 | The library covers 83% of Chlamydomonas genes. a**, 83% of all

3 Chlamydomonas genes have one or more insertions in the library. **b**, In various functional

4 groups, more than 75% of genes are represented by insertions in the library. **c**, The

5 number of insertions per gene is roughly correlated with gene length. Box heights

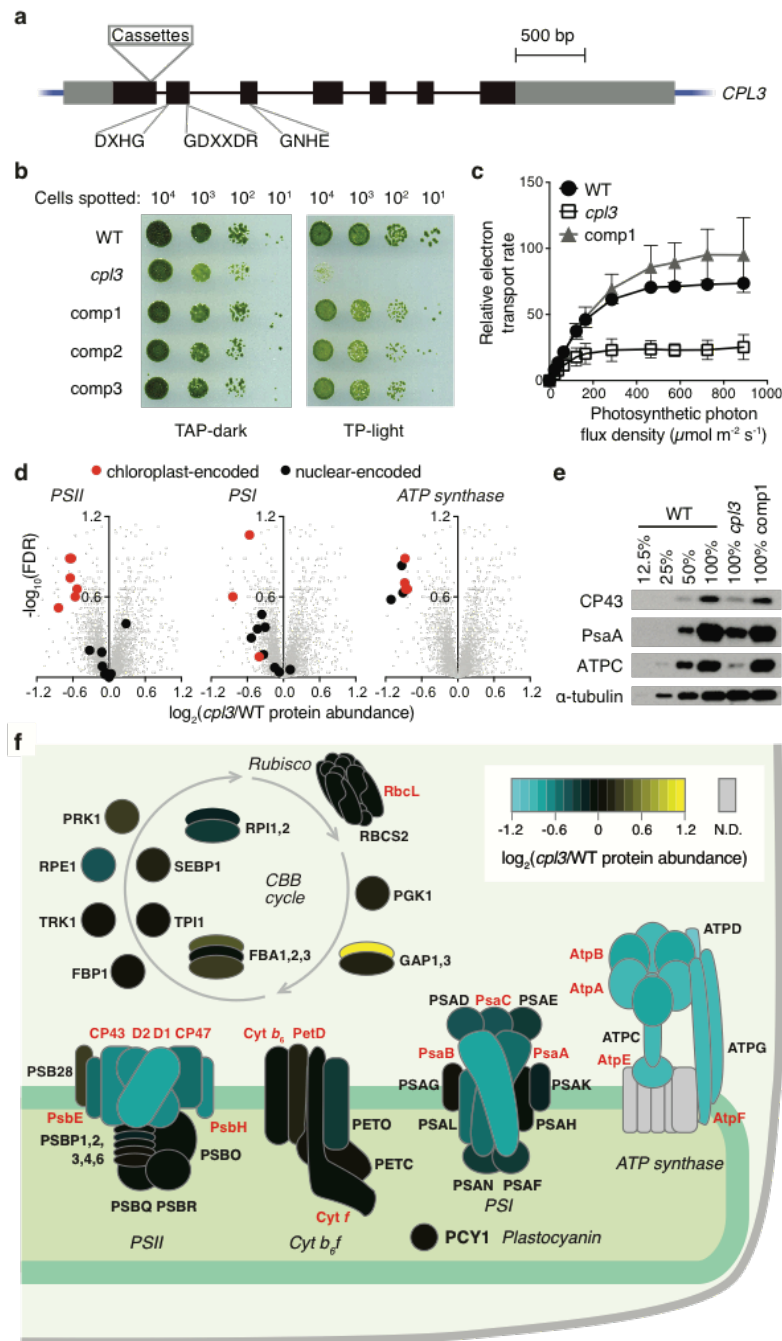
1 represent quartiles, whiskers represent 1st and 99th percentiles, and outliers are plotted as
2 crosses. Box widths are proportional to the number of genes in each bin. **d**, Insertion
3 density varies among different gene features, with the lowest density in exons. **e**, More
4 than 1,800 mutants were distributed to approximately 200 laboratories in the world
5 during the first 18 months of its availability. **f**, Distributed mutants are being used to
6 study a variety of biological processes. Only genes with some functional annotation are
7 shown.
8



1

2 **Fig. 3 | A high-throughput screen using the library identifies many genes with**
 3 **known roles in photosynthesis and reveals many novel components. a**, Unique
 4 barcodes allow screening mutants in a pool. Mutants deficient in photosynthesis can be
 5 identified because their barcodes will be less abundant after growth in photosynthetic
 6 (TP-light) relative to after growth in heterotrophic (TAP-dark) conditions. **b**, Biological
 7 replicates were highly reproducible, with a Spearman's correlation of 0.982. Each dot

1 represents one barcode. See also Supplementary Fig. 5 and Methods. **c**, The phenotype of
2 each insertion was determined by comparing its read count under TAP-dark and TP-light
3 conditions. Insertions that fell below the phenotype cutoff were considered to show a
4 defect in photosynthesis. *cpl3* alleles are highlighted in red squares or triangles. **d**,
5 Insertion phenotypes vary depending on the gene feature disrupted: exon and intron
6 insertions are most likely to show strong phenotypes, while 3'UTR insertions rarely
7 do. The plot is based on all insertions for the 43 higher-confidence genes. **e**, The TP-
8 light/TAP-dark ratio of all the alleles are shown for hit and control genes. Each column is
9 a gene; each horizontal bar is an allele, color-coded by feature. **f**, The 303 candidate
10 genes were categorized according to (1) whether or not they were previously known to
11 play a role in photosynthesis, and (2) whether the screen data yielded higher or lower
12 confidence that mutation of that gene causes a defect in photosynthetic growth. **g**, Known
13 higher-confidence genes, novel higher-confidence genes, and lower-confidence genes are
14 all enriched in predicted chloroplast-targeted proteins ($P < 0.011$). **h**, Twenty-two of the
15 higher-confidence genes and 43 of the lower-confidence genes were previously known to
16 have a role in processes related to photosynthesis. The screen additionally identified 21
17 higher-confidence and 217 lower-confidence genes that were not previously known to be
18 involved in photosynthesis.



1

2 **Fig. 4 | CPL3 is required for photosynthetic growth and accumulation of**
 3 **photosynthetic protein complexes in the thylakoid membranes. a**, The *cpl3* mutant
 4 contains cassettes inserted in the first exon of *CPL3*. The locations of conserved protein
 5 phosphatase motifs are indicated (see Supplementary Fig. 6e). Black boxes indicate

1 exons; gray boxes indicate UTRs. **b**, *cpl3* is deficient in growth under photosynthetic
2 conditions. The growth deficiency is rescued upon complementation with a wild-type
3 copy of the *CPL3* gene (comp1-3 represent three independent complemented lines). **c**,
4 *cpl3* has a lower relative photosynthetic electron transport rate than the wild-type strain
5 (WT) and comp1. Error bars indicate standard deviation ($n \geq 3$). **d**, Whole-cell
6 proteomics indicate that *cpl3* is deficient in the accumulation of PSII, PSI, and the
7 chloroplast ATP synthase. Each gray dot represents one Chlamydomonas protein. The
8 subunits of PSII, PSI and ATP synthase are highlighted as black or red symbols. See also
9 Supplementary Table 14. **e**, Western blots show that CPL3 is required for normal
10 accumulation of the PSII subunit CP43, the PSI subunit PsaA, and the chloroplast ATP
11 synthase subunit ATPC. α -tubulin was used as a loading control. **f**, A schematic summary
12 of the protein abundance of subunits in the light reactions protein complexes or enzymes
13 in the CBB cycle in *cpl3* relative to the wild type based on proteomics data. The relative
14 abundance is shown as a heatmap. Depicted subunits that were not detected by
15 proteomics are filled with gray. Nuclear-encoded proteins are labeled in black font while
16 chloroplast-encoded subunits are labeled in red font. A stack of horizontal ovals indicates
17 different isoforms for the same enzyme, such as FBA1, FBA2, and FBA3. Cyt,
18 cytochrome.

1 Tables

2 Table 1 | Higher-confidence genes from the photosynthesis screen that had a 3 previously known role in photosynthesis.

Category	Gene	Defline/ description in Phytozome ¹²	PredA lgo ^a	Alleles in two replicates			At homolog ^c	Reference and the corresponding organism(s)
				+ ^b	- ^c	FDR ^d		
Calvin- Benson- Bassham cycle	Cre03.g18 5550 (<i>SEBPI</i>)	Sedoheptulose- 1,7- biphosphatase	C	3	0	0.021	AT3G55800. 1 (<i>SBPASE</i>)	Arabidopsis ²⁹
	Cre12.g52 4500 (<i>RMT2</i>)	Rubisco small subunit N- methyltransferase	O	3	0	0.018	AT3G07670. 1	Pisum ³⁰
	Cre06.g29 8300 (<i>MRL1</i>)	Pentatricopeptide repeat protein, stabilizes rbcL mRNA	C	1 2	1 0	1.000 0.239	AT4G34830. 1 (<i>MRL1</i>)	Chlamydomonas and Arabidopsis ³¹
Carbon concentrati ng mechanism	Cre12.g49 7300 (<i>CAS1</i>)	Rhodanese-like Ca-sensing receptor	C	2 2	0 0	0.260 0.239	AT5G23060. 1 (<i>CaS</i>)	Chlamydomonas ³²
	Cre10.g45 2800 (<i>LCIB</i>)	Low-CO2- inducible protein	C	2 1	0 1	0.260 1.000	-	Chlamydomonas ³³
Chloroplast and thylakoid morphogen esis	Cre14.g61 6600	-	M	4 4	3 3	0.021 0.018	AT1G03160. 1 (<i>FZL</i>)	Arabidopsis ³⁴
	Cofactor	Cre13.g58	-	M	5	5	0.010	AT4G31390.

				2	8	1.000		
	Cre10.g42 3500 (<i>HMOX1</i>)	Heme oxygenase	C	3	0	0.021	AT1G69720. 1 (<i>HO3</i>)	Chlamydomonas ¹⁴
				3	0	0.018		
	Cre03.g18 8700 (<i>PLAP6</i>)	Plastid lipid associated protein, Fibrillin	C	3	1	0.070	AT5G09820. 2	Arabidopsis ³⁶
				3	1	0.056		
	Cre16.g65 9050	-	C	4	6	0.098	AT1G68890. 1	Chlamydomonas ³⁷
				4	6	0.075		
PSI protein synthesis and assembly	Cre12.g52 4300 (<i>CGL71</i>)	Predicted protein	C	2	0	0.260	AT1G22700. 1	Synechocystis ³⁸ , Arabidopsis ³⁹ , Chlamydomonas ⁴⁰
				2	0	0.239		
	Cre01.g04 5902	-	C	1	1	1.000	AT3G24430. 1 (<i>HCF101</i>)	Arabidopsis ^{41,42}
				2	0	0.239		
PSI RNA splicing and stabilization	Cre09.g38 9615	-	M	5	0	0.0002	AT3G17040. 1 (<i>HCF107</i>)	Chlamydomonas ⁴³ , Arabidopsis ^{42,44,f}
				5	0	0.0002		
	Cre01.g02 7150 (<i>CPLD46</i>)	DEAD/DEAH-box helicase	M	5	1	0.0004	AT1G70070. 1 (<i>EMB25</i> , <i>ISE2</i> , <i>PDE317</i>)	Arabidopsis ⁴⁵
				5	1	0.0003		
	Cre09.g39 4150 (<i>RAA1</i>)	-	M	5	1	0.0004	-	Chlamydomonas ⁴⁶
				5	1	0.0003		
	Cre12.g53 1050 (<i>RAA3</i>)	PsaA mRNA maturation factor	C	3	0	0.021	-	Chlamydomonas ⁴⁷
				3	0	0.018		
	Cre10.g44 0000 (<i>OPR120</i>)	-	C	2	0	0.260	-	Chlamydomonas ⁴⁸ , 49
				2	0	0.239		
PSII protein	Cre13.g57	Similar to	C	3	3	0.260	AT1G16720.	Arabidopsis ^{42,50,51}

				3	3	0.208		
Cre02.g07 3850 (<i>CGL54</i>)	Predicted protein	C	2	0	0.260	AT1G05385.	Arabidopsis ⁵²	
			2	0	0.239	1 (<i>LPA19</i> , <i>Psb27-H1</i>)		
Cre02.g10 5650	-	C	2	0	0.260	AT5G51545.	Arabidopsis ⁵³	
			2	0	0.239	1 (<i>LPA2</i>)		
Cre06.g27 3700 (<i>HCF136</i>)	-	C	2	0	0.260	AT5G23120.	Arabidopsis ⁴² ;	
			1	1	1.000	1 (<i>HCF136</i>)	Synechocystis ⁵⁴	
Cre10.g43 0150 (<i>LPA1</i>)	-	C	2	0	0.260	AT1G02910.	Arabidopsis ⁵⁵	
			1	1	1.000	1 (<i>LPA1</i>)		

1

2 ^aPrediction of protein localization by PredAlgo⁵⁶: C = chloroplast, M = mitochondrion, SP = secretory pathway, O =

3 other.

4 ^bThe number of exon/intron/5'UTR mutant alleles for that gene that satisfy our requirement of minimum 50 reads and
5 showed at least 10X fewer normalized reads in the TP-light sample compared to the TAP-dark sample.

6 ^cThe number exon/intron/5'UTR mutant alleles for that gene that satisfy our minimum read count requirement but did
7 not satisfy the at least 10X depletion in TP-light criterion.

8 ^dthe FDR for that gene compared to all alleles for all genes (see Methods).

9 ^eArabidopsis homolog, obtained from the "best_arabidopsis_TAIR10_hit_name" field in Phytozome¹².

10 ^fAT3G17040.1 is required for functional PSII in Arabidopsis whereas Cre09.g389615 was shown to be involved in PSI
11 accumulation in Chlamydomonas.

1 **Table 2 | Higher-confidence genes from the photosynthesis screen with no previously**
 2 **known role in photosynthesis.**

Gene	Defline/description in Phytozome	PredAlgo	Alleles in two replicates			At homolog
			+	-	FDR	
Cre01.g008550	Serine/threonine kinase-related	O	2	0	0.260	AT1G73450.1
			1	1	1.000	
Cre01.g014000	-	C	3	0	0.021	-
			3	0	0.018	
Cre01.g037800 (<i>TRX21</i>)	ATP binding protein; thioredoxin domain	O	3	3	0.260	AT2G18990.1 (<i>TXND9</i>)
			1	5	1.000	
Cre02.g073900	All-trans-10'-apo-beta-carotenal 13,14-cleaving dioxygenase	C	3	1	0.070	AT4G32810.1 (<i>ATCCD8, CCD8,</i> <i>MAX4</i>)
			3	1	0.056	
Cre02.g111550	Serine/threonine kinase-related	SP	10	8	< 10 ⁻⁶	AT4G24480.1
			6	12	0.015	
Cre03.g185200 (<i>CPL3</i>)	Metallophosphoesterase/metallo- dependent phosphatase	C	3	4	0.260	AT1G07010.1
			3	4	0.239	
Cre06.g259100	-	C	1	4	1.000	-
			3	2	0.117	
Cre06.g281800	Domain of unknown function (DUF1995)	C	3	0	0.021	-
			3	0	0.018	
Cre07.g316050 (<i>CDJ2</i>)	Chloroplast DnaJ-like protein	M	2	0	0.260	AT5G59610.1
			1	1	1.000	
Cre07.g341850 (<i>EIF2</i>)	Translation initiation factor IF-2, chloroplastic	C	2	0	0.260	AT1G17220.1 (<i>FUG1</i>)
			2	0	0.239	
Cre08.g358350 (<i>TDA1</i>)	Fast leu-rich domain-containing ^a	C	3	2	0.152	-
			3	2	0.117	
Cre09.g396250	Phosphatidate cytidyltransferase	SP	2	0	0.260	AT5G04490.1 (<i>VTE5</i>)

			1	1	1.000	
Cre10.g429650	Alpha/beta hydrolase family (Abhydrolase_5)	O	2	0	0.260	-
			1	1	1.000	
Cre10.g448950	Nocturnin	C	1	1	1.000	AT3G58560.1
			2	0	0.239	
Cre11.g467712	Structural maintenance of chromosomes smc family member; starch-binding domain	M	7	7	0.0003	AT5G05180.1
			7	7	0.0003	
Cre12.g542569	Ionotropic glutamate receptor	O	0	2	1.000	AT1G05200.1 <i>(ATGLR3.4, GLR3.4, GLUR3)</i>
			2	0	0.239	
Cre13.g566400 <i>(OPR55)</i>	Fast leu-rich domain-containing ^a	M	4	2	0.018	-
			4	2	0.015	
Cre13.g574000 <i>(GEF1)</i>	Voltage-gated chloride channel	O	1	11	1.000	AT5G26240.1 <i>(ATCLC-D, CLC-D)</i>
			4	8	0.144	
Cre13.g586750	Transportin 3 and importin	O	3	4	0.260	AT5G62600.1
			2	5	1.000	
Cre16.g658950	-	C	2	2	0.909	-
			3	1	0.056	
Cre50.g761497	Magnesium transporter mrs2 homolog, mitochondrial	M	2	0	0.260	AT5G22830.1 <i>(ATMGT10, GMN10, MGT10, MRS2-11)</i>
			2	0	0.239	

1

2 ^aThe annotation of “fast leu-rich domain-containing” cannot be confirmed by BLASTp analysis at NCBI⁵⁷.

References

1. Karpowicz, S.J., Prochnik, S.E., Grossman, A.R. & Merchant, S.S. The GreenCut2 resource, a phylogenomically derived inventory of proteins specific to the plant lineage. *J Biol Chem* **286**, 21427-39 (2011).
2. Krishnakumar, V. *et al.* Araport: the Arabidopsis information portal. *Nucleic Acids Res* **43**, D1003-9 (2015).
3. Levine, R.P. Genetic Control of Photosynthesis in *Chlamydomonas Reinhardi*. *Proc Natl Acad Sci U S A* **46**, 972-8 (1960).
4. Gutman, B.L. & Niyogi, K.K. *Chlamydomonas* and Arabidopsis. A dynamic duo. *Plant Physiol* **135**, 607-10 (2004).
5. Harris, E.H., Stern, D.B. & Witman, G.B. *The Chlamydomonas Sourcebook*, (Academic Press, 2009).
6. Rochaix, J.D. *Chlamydomonas reinhardtii* as the photosynthetic yeast. *Annu Rev Genet* **29**, 209-30 (1995).
7. Li, J.B. *et al.* Comparative genomics identifies a flagellar and basal body proteome that includes the BBS5 human disease gene. *Cell* **117**, 541-52 (2004).
8. Silflow, C.D. & Lefebvre, P.A. Assembly and motility of eukaryotic cilia and flagella. Lessons from *Chlamydomonas reinhardtii*. *Plant Physiol* **127**, 1500-7 (2001).
9. Li, X. *et al.* An Indexed, Mapped Mutant Library Enables Reverse Genetics Studies of Biological Processes in *Chlamydomonas reinhardtii*. *Plant Cell* **28**, 367-87 (2016).
10. Terashima, M., Specht, M. & Hippler, M. The chloroplast proteome: a survey from the *Chlamydomonas reinhardtii* perspective with a focus on distinctive features. *Curr Genet* **57**, 151-68 (2011).
11. Pazour, G.J., Agrin, N., Leszyk, J. & Witman, G.B. Proteomic analysis of a eukaryotic cilium. *J Cell Biol* **170**, 103-13 (2005).
12. Merchant, S.S. *et al.* The *Chlamydomonas* genome reveals the evolution of key animal and plant functions. *Science* **318**, 245-251 (2007).
13. Zones, J.M., Blaby, I.K., Merchant, S.S. & Umen, J.G. High-Resolution Profiling of a Synchronized Diurnal Transcriptome from *Chlamydomonas reinhardtii* Reveals Continuous Cell and Metabolic Differentiation. *Plant Cell* (2015).
14. Duanmu, D. *et al.* Retrograde bilin signaling enables *Chlamydomonas* greening and phototrophic survival. *Proc Natl Acad Sci U S A* **110**, 3621-6 (2013).
15. Dent, R.M. *et al.* Large-scale insertional mutagenesis of *Chlamydomonas* supports phylogenomic functional prediction of photosynthetic genes and analysis of classical acetate-requiring mutants. *Plant J* **82**, 337-351 (2015).
16. Allen, J.F., de Paula, W.B., Puthiyaveetil, S. & Nield, J. A structural phylogenetic map for chloroplast photosynthesis. *Trends Plant Sci* **16**, 645-55 (2011).
17. Giordano, M., Beardall, J. & Raven, J.A. CO₂ concentrating mechanisms in algae: mechanisms, environmental modulation, and evolution. *Annu Rev Plant Biol* **56**, 99-131 (2005).
18. Goldschmidt-Clermont, M. & Rahire, M. Sequence, evolution and differential expression of the two genes encoding variant small subunits of ribulose biphosphate carboxylase/oxygenase in *Chlamydomonas reinhardtii*. *J Mol Biol* **191**, 421-32 (1986).
19. Suzuki, Y., Arae, T., Green, P.J., Yamaguchi, J. & Chiba, Y. AtCCR4a and AtCCR4b are Involved in Determining the Poly(A) Length of Granule-bound starch synthase 1 Transcript and Modulating Sucrose and Starch Metabolism in Arabidopsis thaliana. *Plant Cell Physiol* **56**, 863-74 (2015).
20. Wang, H. *et al.* The global phosphoproteome of *Chlamydomonas reinhardtii* reveals complex organellar phosphorylation in the flagella and thylakoid membrane. *Mol Cell Proteomics* **13**, 2337-53 (2014).
21. Bassi, R., Soen, S.Y., Frank, G., Zuber, H. & Rochaix, J.D. Characterization of Chlorophyll-a/B Proteins of Photosystem-I from *Chlamydomonas-Reinhardtii*. *J Biol Chem* **267**, 25714-25721 (1992).

- 1 22. Sager, R. & Zalokar, M. Pigments and photosynthesis in a carotenoid-deficient mutant of
2 Chlamydomonas. *Nature* **182**, 98-100 (1958).
- 3 23. Baek, K. *et al.* DNA-free two-gene knockout in Chlamydomonas reinhardtii via CRISPR-Cas9
4 ribonucleoproteins. *Sci Rep* **6**, 30620 (2016).
- 5 24. Jiang, W., Brueggeman, A.J., Horken, K.M., Plucinak, T.M. & Weeks, D.P. Successful transient
6 expression of Cas9 and single guide RNA genes in Chlamydomonas reinhardtii. *Eukaryot Cell* **13**,
7 1465-9 (2014).
- 8 25. Shin, S.E. *et al.* CRISPR/Cas9-induced knockout and knock-in mutations in Chlamydomonas
9 reinhardtii. *Sci Rep* **6**, 27810 (2016).
- 10 26. Slaninová, M., Hroššová, D., Vlček, D. & Wolfgang, W. Is it possible to improve homologous
11 recombination in Chlamydomonas reinhardtii? *Biologia* **63**, 941-46 (2008).
- 12 27. Greiner, A. *et al.* Targeting of Photoreceptor Genes in Chlamydomonas reinhardtii via Zinc-finger
13 Nucleases and CRISPR/Cas9. *Plant Cell* (2017).
- 14 28. Ferenczi, A., Pyott, D.E., Xipnitou, A. & Molnar, A. Efficient targeted DNA editing and
15 replacement in Chlamydomonas reinhardtii using Cpf1 ribonucleoproteins and single-stranded
16 DNA. *Proc Natl Acad Sci U S A* (2017).
- 17 29. Liu, X.L., Yu, H.D., Guan, Y., Li, J.K. & Guo, F.Q. Carbonylation and loss-of-function analyses
18 of SBPase reveal its metabolic interface role in oxidative stress, carbon assimilation, and multiple
19 aspects of growth and development in Arabidopsis. *Mol Plant* **5**, 1082-99 (2012).
- 20 30. Klein, R.R. & Houtz, R.L. Cloning and developmental expression of pea ribulose-1,5-
21 bisphosphate carboxylase/oxygenase large subunit N-methyltransferase. *Plant Mol Biol* **27**, 249-
22 61 (1995).
- 23 31. Johnson, X. *et al.* MRL1, a conserved Pentatricopeptide repeat protein, is required for stabilization
24 of rbcL mRNA in Chlamydomonas and Arabidopsis. *Plant Cell* **22**, 234-48 (2010).
- 25 32. Wang, L. *et al.* Chloroplast-mediated regulation of CO₂-concentrating mechanism by Ca²⁺-
26 binding protein CAS in the green alga Chlamydomonas reinhardtii. *Proc Natl Acad Sci U S A* **113**,
27 12586-12591 (2016).
- 28 33. Wang, Y. & Spalding, M.H. An inorganic carbon transport system responsible for acclimation
29 specific to air levels of CO₂ in Chlamydomonas reinhardtii. *Proc Natl Acad Sci U S A* **103**, 10110-
30 5 (2006).
- 31 34. Gao, H., Sage, T.L. & Osteryoung, K.W. FZL, an FZO-like protein in plants, is a determinant of
32 thylakoid and chloroplast morphology. *Proc Natl Acad Sci U S A* **103**, 6759-64 (2006).
- 33 35. Martinis, J. *et al.* ABC1K1/PGR6 kinase: a regulatory link between photosynthetic activity and
34 chloroplast metabolism. *Plant J* **77**, 269-83 (2014).
- 35 36. Kim, E.H., Lee, Y. & Kim, H.U. Fibrillin 5 Is Essential for Plastoquinone-9 Biosynthesis by
36 Binding to Solanesyl Diphosphate Synthases in Arabidopsis. *Plant Cell* **27**, 2956-71 (2015).
- 37 37. Lefebvre-Legendre, L. *et al.* Loss of phyloquinone in Chlamydomonas affects plastoquinone pool
38 size and photosystem II synthesis. *J Biol Chem* **282**, 13250-63 (2007).
- 39 38. Wilde, A., Lunser, K., Ossenbuhl, F., Nickelsen, J. & Borner, T. Characterization of the
40 cyanobacterial ycf37: mutation decreases the photosystem I content. *Biochem J* **357**, 211-6 (2001).
- 41 39. Stockel, J., Bennewitz, S., Hein, P. & Oelmuller, R. The evolutionarily conserved tetratricopeptide
42 repeat protein pale yellow green7 is required for photosystem I accumulation in
43 Arabidopsis and copurifies with the complex. *Plant Physiol* **141**, 870-8 (2006).
- 44 40. Heinnickel, M. *et al.* Tetratricopeptide repeat protein protects photosystem I from oxidative
45 disruption during assembly. *Proc Natl Acad Sci U S A* **113**, 2774-9 (2016).
- 46 41. Lezhneva, L., Amann, K. & Meurer, J. The universally conserved HCF101 protein is involved in
47 assembly of [4Fe-4S]-cluster-containing complexes in Arabidopsis thaliana chloroplasts. *Plant J*
48 **37**, 174-85 (2004).
- 49 42. Meurer, J., Meierhoff, K. & Westhoff, P. Isolation of high-chlorophyll-fluorescence mutants of
50 Arabidopsis thaliana and their characterisation by spectroscopy, immunoblotting and northern
51 hybridisation. *Planta* **198**, 385-96 (1996).
- 52 43. Douchi, D. *et al.* A Nucleus-Encoded Chloroplast Phosphoprotein Governs Expression of the
53 Photosystem I Subunit PsaC in Chlamydomonas reinhardtii. *Plant Cell* **28**, 1182-99 (2016).
- 54 44. Felder, S. *et al.* The nucleus-encoded HCF107 gene of Arabidopsis provides a link between
55 intergenic RNA processing and the accumulation of translation-competent psbH transcripts in
56 chloroplasts. *Plant Cell* **13**, 2127-41 (2001).

- 1 45. Carlotto, N. *et al.* The chloroplastic DEVH-box RNA helicase INCREASED SIZE EXCLUSION
2 LIMIT 2 involved in plasmodesmata regulation is required for group II intron splicing. *Plant Cell*
3 *Environ* **39**, 165-73 (2016).
- 4 46. Perron, K., Goldschmidt-Clermont, M. & Rochaix, J.D. A factor related to pseudouridine
5 synthases is required for chloroplast group II intron trans-splicing in *Chlamydomonas reinhardtii*.
6 *EMBO J* **18**, 6481-90 (1999).
- 7 47. Rivier, C., Goldschmidt-Clermont, M. & Rochaix, J.D. Identification of an RNA-protein complex
8 involved in chloroplast group II intron trans-splicing in *Chlamydomonas reinhardtii*. *EMBO J* **20**,
9 1765-73 (2001).
- 10 48. Jacobs, J. *et al.* Identification of a chloroplast ribonucleoprotein complex containing trans-splicing
11 factors, intron RNA, and novel components. *Mol Cell Proteomics* **12**, 1912-25 (2013).
- 12 49. Marx, C., Wunsch, C. & Kuck, U. The Octatricopeptide Repeat Protein Raa8 Is Required for
13 Chloroplast trans Splicing. *Eukaryot Cell* **14**, 998-1005 (2015).
- 14 50. Link, S., Engelmann, K., Meierhoff, K. & Westhoff, P. The atypical short-chain dehydrogenases
15 HCF173 and HCF244 are jointly involved in translational initiation of the psbA mRNA of
16 *Arabidopsis*. *Plant Physiol* **160**, 2202-18 (2012).
- 17 51. Schult, K. *et al.* The nuclear-encoded factor HCF173 is involved in the initiation of translation of
18 the psbA mRNA in *Arabidopsis thaliana*. *Plant Cell* **19**, 1329-46 (2007).
- 19 52. Wei, L. *et al.* LPA19, a Psb27 homolog in *Arabidopsis thaliana*, facilitates D1 protein precursor
20 processing during PSII biogenesis. *J Biol Chem* **285**, 21391-8 (2010).
- 21 53. Ma, J. *et al.* LPA2 is required for efficient assembly of photosystem II in *Arabidopsis thaliana*.
22 *Plant Cell* **19**, 1980-93 (2007).
- 23 54. Komenda, J. *et al.* The cyanobacterial homologue of HCF136/YCF48 is a component of an early
24 photosystem II assembly complex and is important for both the efficient assembly and repair of
25 photosystem II in *Synechocystis* sp. PCC 6803. *J Biol Chem* **283**, 22390-9 (2008).
- 26 55. Peng, L. *et al.* LOW PSII ACCUMULATION1 is involved in efficient assembly of photosystem II
27 in *Arabidopsis thaliana*. *Plant Cell* **18**, 955-69 (2006).
- 28 56. Tardif, M. *et al.* PredAlgo: a new subcellular localization prediction tool dedicated to green algae.
29 *Mol Biol Evol* **29**, 3625-39 (2012).
- 30 57. Altschul, S.F. *et al.* Gapped BLAST and PSI-BLAST: a new generation of protein database search
31 programs. *Nucleic Acids Res* **25**, 3389-402 (1997).

32

1

2 SUPPLEMENTARY INFORMATION

3

4

5

6 **A genome-wide algal mutant library reveals a global**

7 **view of genes required for eukaryotic photosynthesis**

8

9 **Xiaobo Li, Weronika Patena, Friedrich Fauser, Robert E. Jinkerson, Shai Saroussi,**

10 **Moritz T. Meyer, Nina Ivanova, Jacob M. Robertson, Rebecca Yue, Ru Zhang,**

11 **Josep Vilarrasa-Blasi, Tyler M. Wittkopp, Silvia Ramundo, Sean R. Blum, Audrey**

12 **Goh, Matthew Laudon, Tharan Srikumar, Paul A. Lefebvre, Arthur R. Grossman,**

13 **and Martin C. Jonikas***

14

15 *e-mail: mjonikas@princeton.edu

16	Table of Contents	
17		
18	Methods.....	4
19	Supplementary Note	36
20	Accuracy of insertion mapping and number of insertions per mutant.....	36
21	Deletions, duplications, and junk fragments associated with insertions are small	37
22	Insertion sites are randomly distributed with mild cold spots and a small number of hot	
23	spots	38
24	Absence of insertions identifies over 200 genes potentially essential for growth under	
25	the propagation conditions used.....	38
26	Deleterious mutations rather than differential chromatin configuration are the major	
27	cause of insertion density variation.....	40
28	Disruption of <i>CPL3</i> is the cause of the photosynthetic deficiency in the <i>cpl3</i> mutant...40	
29	Supplementary Figures	42
30	Captions for Supplementary Table 1 to S14	55
31	References.....	62
32		

33 **Other supplementary materials for this manuscript (provided as**
34 **separate excel or text files):**

35 Supplementary Table 1. Primers and experimental design for all PCRs related to
36 library generation and mapping.

37 Supplementary Table 2. Binary codes for plate super-pooling.

38 Supplementary Table 3. Binary codes for colony super-pooling.

39 Supplementary Table 4. Read counts for each barcode in each combinatorial
40 super-pool.

41 Supplementary Table 5. List of all mapped mutants in the library.

42 Supplementary Table 6. Primers and results of PCRs used to verify the insertion
43 sites of randomly-picked mutants from the mutant library.

44 Supplementary Table 7. Statistically significant insertion hot spots and cold spots.

45 Supplementary Table 8. Statistically significant depleted functional terms.

46 Supplementary Table 9. Candidate essential genes.

47 Supplementary Table 10. Read counts of barcodes before and after pooled growth
48 in the photosynthesis screen.

49 Supplementary Table 11. Statistics of the pooled growth data for all genes.

50 Supplementary Table 12. Summary of previous characterizations of higher- and
51 lower-confidence genes' roles in photosynthesis.

52 Supplementary Table 13. Read counts of *cpl3* exon and intron alleles in the
53 pooled screens.

54 Supplementary Table 14. Proteomic characterization of the *cpl3* mutant.

55

56 **Methods**

57 **Generation of the indexed and barcoded mutant library: a conceptual overview.** A

58 three-step pipeline was developed for the generation of an indexed, barcoded library of
59 insertional mutants in *Chlamydomonas* (Fig. 1b and Supplementary Fig. 1).

60 To generate mutants, CC-4533¹ (“wild type” in text and figures) cells were
61 transformed with DNA cassettes that randomly insert into the genome, confer
62 paromomycin resistance for selection, and inactivate the genes they insert into. Each
63 cassette contained two unique 22 nucleotide barcodes, one at each end of the cassette
64 (Supplementary Fig. 1a-d). Transformants were arrayed on agar plates and each insertion
65 in a transformant would contain two barcodes. The barcode sequences as well as the
66 insertion site were initially unknown (Supplementary Fig. 1e).

67 To determine the sequences of the barcodes in each colony, combinatorial pools
68 of the individual mutants were generated, with DNA extracted, and barcodes amplified
69 and deep-sequenced. The combinatorial pooling patterns were designed so that each
70 colony was included in a different combination of pools, allowing us to determine the
71 barcode sequences associated with individual colonies based on which pools the
72 sequences were found in (Supplementary Fig. 1f and Supplementary Fig. 2a-e; Methods).
73 This procedure was similar in concept to the approach we used in our pilot study², but it
74 consumed significantly less time because we used a simple PCR amplifying only the
75 barcodes instead of a multi-step flanking sequence extraction protocol (ChlaMmeSeq¹)
76 on each combinatorial pool.

77 To determine the insertion site associated with each barcode, the library was
78 pooled into a single sample or six separate samples. The barcodes and their flanking

79 genomic DNA were PCR amplified using LEAP-Seq² (Supplementary Fig. 1g and
80 Supplementary Fig. 2f-j; Methods). The flanking sequences associated with each barcode
81 were obtained by paired-end deep sequencing^{3,4}. The final product is an indexed library
82 in which each colony has known flanking sequences that identify the genomic insertion
83 site, and barcode sequences that facilitate pooled screens in which individual mutants can
84 be tracked by deep sequencing (Fig. 3a).

85 Experimental details for this pipeline are described in paragraphs below.

86

87 **Generation of insertion cassettes.** The insertion cassette designated Cassette containing
88 Internal Barcodes 1 (CIB1) was generated in four steps: (1) generating double-stranded
89 DNAs containing random sequences (Supplementary Fig. 1a); (2) digesting the double-
90 stranded DNAs to yield cassette ends (Supplementary Fig. 1a); (3) obtaining the
91 backbone from digestion of plasmid pMJ016c that contains the sequences between the
92 two barcodes (Supplementary Fig. 1b); (4) ligating the two cassette ends with the cassette
93 backbone (Supplementary Fig. 1c).

94 Step 1: To generate each end of the cassette that contains barcodes, a long
95 oligonucleotide primer (Supplementary Fig. 1a and Supplementary Table 1) containing a
96 random sequence region of 22 nucleotides was used as a template for the extension of a
97 shorter oligonucleotide primer. Each 50- μ L reaction mixture contained 32 μ L H₂O, 10
98 μ L Phusion GC buffer, 1.5 μ L DMSO, 1 μ L 10 mM dNTP, 2.5 μ L 10 μ M long oligo, 2.5
99 μ L 10 μ M short oligo, and 0.5 μ L Phusion HS II DNA polymerase (F549L, Thermo
100 Fisher). The reaction mixtures were subjected to a single thermal cycle: 98°C for 40 sec,
101 97°C to 63°C ramp (-1°C every 10 sec), 63°C for 30 sec, 72°C for 5 min.

102 Step 2: The double-stranded product yielded from Step 1 was digested using *BsaI*
103 (R0535L, New England Biolabs). For the 5' side primer extension product, the digestion
104 yielded two bands of 87 bp (plus 4 nt of overhang) and 31 bp (plus 4 nt of overhang). For
105 the 3' side, they were 68 bp and 31 bp. The larger band from each digestion was purified
106 from a 2.5% agarose gel using D-tubes (71508-3, EMD Millipore) as previously
107 described¹ (Supplementary Fig. 1a).

108 Step 3: The synthesized plasmid pMJ016c, which contains the *HSP70-RBCS2*
109 promoter, the paromomycin resistance gene *AphVIII*, and the *PSAD* and *RPL12*
110 terminators, was digested using *BsaI*. Two bands of 2064 bp and 3363 bp were obtained.
111 The 2064 bp band (cassette backbone) was purified from a 0.8% agarose gel using the
112 QIAquick Kit (28106, Qiagen) according to the manufacturer's instructions
113 (Supplementary Fig. 1b).

114 Step 4: The two fragments and the cassette backbone were ligated using T4 DNA
115 ligase (M0202L, New England Biolabs) (Supplementary Fig. 1c). Each 30- μ L reaction
116 mixture contained 38 ng 5' cassette end, 30 ng 3' cassette end, 305 ng cassette backbone,
117 3 μ L ligase buffer, and 0.5 μ L ligase. The double-stranded product of 2,223 bp was gel
118 purified using D-tubes and used for mutant generation. The sequence of the CIB1 cassette
119 generated (Supplementary Fig. 1d) has been uploaded to the mutant ordering website:
120 <https://www.chlamylibrary.org/showCassette?cassette=CIB1>.

121

122 **Mutant generation, mutant maintenance, and medium recipes.** Chlamydomonas CC-
123 4533 strain was grown in Tris-Acetate-Phosphate (TAP) medium in a 20-L container
124 under 100 μ mol photons $m^{-2} s^{-1}$ light (measured at the periphery) to a density of 1-1.5x

125 10^6 cells/mL. Cells were collected by centrifugation at 300-1,000g for 4 min. Pellets were
126 washed once with 25 mL TAP medium supplemented with 40 mM sucrose, and then
127 resuspended in TAP supplemented with 40 mM sucrose at 2×10^8 cells/mL. 250 μ L of
128 cell suspension was then aliquoted into each electroporation cuvette (Bio-Rad) and
129 incubated at 16°C for 5-30 min. For each cuvette, 5 μ L DNA cassette CIB1 at 5 ng/ μ L
130 was added to the cell suspension and mixed by pipetting. Electroporation was performed
131 immediately as previously described¹. After electroporation, cells from each cuvette were
132 diluted into 8 mL TAP supplemented with 40 mM sucrose and shaken gently in dark for
133 6 h. After incubation, cells were plated on TAP containing 20 μ g/mL paromomycin (800
134 μ L per plate) and incubated in darkness for approximately two weeks before colony
135 picking.

136 Approximately 210,000 total mutants were picked using a Norgren CP7200
137 colony picking robot and maintained on 570 agar plates, each containing a 384-colony
138 array. We propagated this original, full library by robotically passaging the mutant arrays
139 to fresh 1.5% agar solidified TAP medium containing 20 μ g/mL paromomycin using a
140 Singer RoToR robot (Singer Instruments)². The full collection was grown in complete
141 darkness at room temperature and passaged every four weeks. In this collection, 127,847
142 of the mutants were mapped. Colonies that failed to yield barcodes or flanking sequences
143 may contain truncated insertion cassettes¹ that have lost the primer binding sites used for
144 barcode amplification or LEAP-Seq analysis. By removing the mutants that were not
145 mapped, mutants that did not survive propagation, and some of the mutants in genes with
146 20 or more insertions, we consolidated 62,389 mutants into 245 plates of 384-colony
147 arrays for long-term robotic propagation.

148 The TAP medium was prepared as previously reported⁵. The TP medium used in
149 this research was similar to TAP except that HCl instead of acetic acid was used to adjust
150 the pH to 7.5.

151

152 **Combinatorial pooling.** For combinatorial pooling and barcode determination for each
153 mutant colony, 570 plate-pools (each containing all mutants on one plate) and 384
154 colony-pools (each containing all mutants in the same colony position across all plates)
155 were generated from two separate sets of the library as previously described². Binary
156 error-correcting codes were used to design combinatorial pooling schemes, as previously
157 described². The existence of suitable binary error-correcting codes and their mathematical
158 construction methods were checked using an online database⁶. For colony super-pooling,
159 the same 384-codeword subset of the [20,10,6] code as previously employed² was used.
160 For plate super-pooling, the [21,11,6] code was generated by triple shortening of the
161 [24,14,6] code⁷. In order to ensure detection of cases of two colonies derived from a
162 single mutant, which could otherwise cause incorrect colony locations to be identified for
163 such mutants, the subset of codewords with a bit sum of 10 (708 codewords) was taken
164 from the [21,11,6] code, using the `choose_codewords_by_bit_sum` function. Both subsets
165 of codewords were checked for the possibility of such sister colony conflicts using the
166 `clonality_conflict_check` function: no conflicts were detected up to 2 errors, meaning any
167 incorrect result due to a sister colony case would have at least 2 differences compared to
168 any expected correct result. The final subset of 570 codewords for plate super-pooling
169 was chosen as previously². The final codeword lists are provided as Supplementary
170 Tables 2 and 3.

171 Generation of plate-super-pools and colony-super-pools from the plate-pools and
172 colony-pools was performed using the Biomek FX liquid handling robot (Beckman
173 Coulter) as previously described². The instruction files for the Biomek robot were
174 generated using the robotic_plate_transfer.py program.

175

176 **Barcode amplification from super-pools.** DNA was extracted from super-pool samples
177 as previously described¹ and the barcodes were amplified (Supplementary Fig. 1f) using
178 the Phusion HSII PCR system. For either 5' or 3' barcode amplifications, one primer (5'
179 R1 or 3' R1; sequences provided in Supplementary Table 1) used in the PCR was
180 common for all super-pools; the other primer (5' R2-1, 5' R2-2,...; 3' R2-1, 3' R2-2,...;)
181 contained an index sequence that allows multiplexed sequencing, i.e. combining of
182 multiple samples in one sequencing lane. Each 50 µL PCR mixture contained 125 ng
183 genomic DNA, 10 µL GC buffer, 5 µL DMSO, 1 µL dNTPs at 10 mM, 1 µL (for 5') or 2
184 µL (for 3') MgCl₂ at 50 mM, 2.5 µL of each primer at 10 µM, and 1 µL Phusion HSII
185 polymerase. The reaction mixtures were incubated at 98°C for 3 min, followed by 10
186 three-step cycles (10 sec at 98°C, 25 sec at 58°C or 63°C for 5' and 3' barcodes
187 respectively, and 15 sec at 72°C), and then 8 two-step cycles (10 sec at 98°C, and 40 sec
188 at 72°C). Similar amount of products from three to eight super-pools were combined,
189 purified using MinElute columns (28006, Qiagen), and the product bands (235 bp for 5'
190 and 209 bp for 3') were gel purified. The purified products were sequenced using the
191 Illumina HiSeq platform from a single end with a custom primer (5' Seq and 3' Seq,
192 Supplementary Table 1).

193

194 **Deconvolution of super-pool sequencing data.** The barcode sequences were extracted
195 from the Illumina sequencing data from each super-pool using the cutadapt command-
196 line program⁸, with a 13 bp expected cassette sequence, allowing 1 alignment error, and
197 taking the trimmed barcode reads between 21 and 23 bp in length. The command for 5'
198 sequences was “cutadapt -a GGCAAGCTAGAGA -e 0.1 -m 21 -M 23”, and for 3'
199 sequences “cutadapt -a TAGCGCGGGGCGT -e 0.1 -m 21 -M 23”. A barcode was found
200 in 97-99% of the sequences in each super-pool.

201 The reads for each distinct barcode sequence in each super-pool were counted
202 (Supplementary Table 4). Many of the sequenced barcodes are likely to contain PCR or
203 sequencing errors. Such barcodes were left uncorrected, because they are very unlikely to
204 appear in enough super-pools to be deconvolved and included in the final data. The
205 deconvolution based on the read count table was performed as previously described², for
206 5' and 3' data separately. A single set of optimized (N, x) parameters was chosen for each
207 dataset, with m = 0 in all cases: N = 8 and x = 0.14 for 5' plate-super-pool data, N = 8
208 and x = 0.16 for 3' plate-super-pool data, N = 6 and x = 0.12 for 5' colony-super-pool
209 data, N = 6 and x = 0.1 for 3' colony-super-pool data. Note that data for colony-super-
210 pool 14 are missing for plates 351-570, which caused imperfections in the deconvolution
211 process, but the missing data were dispensable due to the error-correction capability built
212 into the pooling scheme.

213

214 **LEAP-Seq.** To connect the flanking sequence with the corresponding barcode for each
215 insertion, we performed LEAP-Seq as reported before² except that barcodes in addition to
216 the flanking sequences were included in the amplicons (Supplementary Fig. 1g, and

217 Supplementary Fig. 2f). Genomic DNA of mutants in the library was used as the template
218 for the extension of a biotinylated primer that anneals to the insertion cassette. The
219 primer extension products were purified by binding to streptavidin-coupled magnetic
220 beads and then ligated to a single-stranded DNA adapter. The ligation products were then
221 used as templates for PCR amplification. The PCR products were gel-purified before
222 being submitted for deep sequencing.

223 We tried different combinations of primers and attempted to perform LEAP-Seq
224 either on six sub-pools (each containing mutants from one-sixth of the library) separately
225 or on the entire library in a single reaction (Supplementary Table 1). Sequencing results
226 from all the samples were used in the analyses below.

227

228 **Basic LEAP-Seq data analysis.** The LEAP-Seq samples were sequenced with Illumina
229 Hi-Seq, yielding paired-end reads. Each read pair has a proximal side, containing the
230 barcode, a part of the cassette sequence, and the immediate genomic flanking sequence;
231 and a distal side, containing the genomic sequence a variable distance away
232 (Supplementary Fig. 2f-j).

233 A newly developed method was used to separate cassette sequence from the
234 proximal reads and thus identify the barcode and genomic flanking sequence even in
235 cases where the cassette was truncated. This was done using the
236 `deepseq_strip_cassette.py` script, which uses local bowtie2 alignment⁹ to detect short
237 cassette sequence. A bowtie2 alignment was performed against the expected cassette
238 sequence (GGAGACGTGTTTCTGACGAGGGCTCGTGTGACTAGTGAGTCCAAC
239 for 5' reads and

240 ACTGACGTCGAGCCTTCTGGCAGACTAGTTGCTCCTGAGTCCAAC for 3' reads),
241 using the following bowtie2 options: "--local --all --ma 3 --mp 5,5 --np 1 --rdg 5,3 --rfg
242 4,3 --score-min C,20,0 -N0 -L5 -i C,1,0 -R5 -D30 --norc --reorder". The alignments for
243 each proximal read were filtered to only consider cases where the cassette aligns after a
244 21-23 bp barcode, at most 5 bp of expected initial cassette sequence are missing, and at
245 least 10 bp of expected cassette sequence are aligned with at most 30% errors. Out of the
246 filtered alignments, the best one was chosen in a maximally deterministic manner, in
247 order to ensure that multiple reads of the same insertion junction yield the same result.
248 The alignment with the highest alignment score is chosen (the bowtie scoring function
249 was customized to distinguish between as many cases as possible); if there were multiple
250 alignments with the same score, the one with the longer alignment was chosen.

251 The resulting cassette alignment was then removed from each proximal read, with
252 the section before the cassette being considered the barcode and the section after the
253 cassette being considered the genomic flanking region. The resulting genomic proximal
254 reads and the raw genomic distal reads were trimmed to 30 bp using the fastx_trimmer
255 command-line utility (http://hannonlab.cshl.edu/fastx_toolkit), aligned to the
256 Chlamydomonas genome (version 5.5 from Phytozome¹⁰) and the cassette, and the
257 alignments were filtered to yield a single result using deepseq_alignment_wrapper.py, as
258 previously described¹.

259 The barcode sequences and proximal and distal alignment results were merged
260 into a single dataset, with data grouped into insertion junctions based on the barcode,
261 using the add_RISCC_alignment_files_to_data function. Data relating to barcodes that
262 were not present in the combinatorial deconvolution results were discarded. The gene-

263 related information for each insertion junction was added using the
 264 `find_genes_for_mutants` and `add_gene_annotation` functions. All functions in this
 265 paragraph are methods of the `Insertional_mutant_pool_dataset` class in the
 266 `mutant_IB_RISCC_classes.py` module.

267

268 **Detecting pairs of flanking sequences that correspond to two sides of the same**
 269 **insertion (confidence levels 1 or 2).** Pairs of insertion junctions likely derived from two
 270 sides of the same insertion were identified using the
 271 `deconvolution_utilities.get_matching_sides_from_table` function, using the method
 272 previously described², with an additional distance bin of 1-10 kb. The resulting pair
 273 counts were as follows:

	0 bp	1-10 bp	11-100 bp	101 bp - 1 kb	1-10 kb	10+ kb
Inner-cassette (toward-facing)	3935	17708	7866	737	339	540
Outer-cassette (away-facing)	-	5010	188	560	58	494
Same-direction	13	17	40	158	133	1520

274

275 Additionally, there were 22,247 pairs in which the two junctions were mapped to
 276 different chromosomes.

277 The number of inner-cassette pairs is significantly larger than 50% of the number
 278 of same-direction pairs in all size ranges up to 10 kb, implying that most of the inner-

279 cassette pairs in those size ranges are derived from a single insertion with a genomic
280 deletion corresponding to the distance. This can be further confirmed by looking at the
281 indicators of the probability of correct mapping for the insertion junctions: insertions with
282 both sides mapped to the same region are almost certainly correctly mapped, and
283 therefore independent indications of their correct mapping should be higher than for other
284 insertions. As expected, the inner-cassette pairs up to 10 kb have a higher fraction of very
285 high confidence insertion pairs (with both sides having 70% or more read pairs mapping
286 to the same locus, and 500 bp or higher longest distance spanned by such read pairs): for
287 size ranges up to 10 kb, 37-41% of the pairs are very high confidence, while for 10+ kb
288 the number is only 16%.

289 The number of outer-cassette pairs is significantly larger than 50% of the number
290 of same-direction pairs in size ranges between 1 bp and 1 kb, implying that most of the
291 outer-cassette pairs in those size ranges are derived from a single insertion. There are two
292 possible physical interpretations of a single insertion yielding an outer-cassette pair of
293 insertion junctions: (1) an insertion with a genomic duplication causing the same genomic
294 DNA sequence to be present on both sides of the cassette (potentially due to single-strand
295 repair); and (2) an insertion of two cassettes flanking a “junk” fragment of genomic
296 DNA. The 1-10 bp cases must be a genomic duplication, since a 1-10 bp “junk” fragment
297 could not yield a 30 bp flanking sequences aligning to the genome. This is confirmed by
298 41% of the pairs being very high confidence. The 101 bp-1 kb cases are almost certainly
299 insertions of two cassettes flanking a “junk” fragment, based on only 3.8% of them being
300 very high confidence. The 188 11-100 bp cases, with a 27% very high confidence, are
301 likely split between the two categories; based on previous analysis¹ we used 30 bp as the

302 cutoff between cases 1 and 2 for outer-cassette pairs. The case 2 pairs, i.e. insertions of
303 two cassettes flanking a junk fragment, were used to determine the typical range of
304 lengths of junk fragments (Supplementary Fig. 3f).

305 Based on this analysis, all insertion junction pairs likely to be derived from two
306 sides of the same insertion (inner-cassette up to 10 kb and outer-cassette up to 30 bp)
307 were categorized as confidence level 1 (extremely likely to be correctly mapped) because
308 their mapping position is derived from two independent flanking sequences. They were
309 annotated in Supplementary Table 5 as confidence level 1, and the “if_both_sides”
310 column was set to “perfect” for the 0 bp distance cases, “deletion” for the remaining
311 inner-cassette cases, and “duplication” for the outer-cassette cases.

312 A similar type of analysis was performed to look for pairs of insertion junctions
313 derived from two sides of an insertion with a junk fragment. For each pair of insertion
314 junctions in one colony (except pairs of insertion junctions already identified as two sides
315 of the same insertion), we looked at the distance and relative orientation between the
316 proximal read of the first junction and each distal read from the second junction; cases
317 where the distal read was mapped to within 10 kb of the proximal read were counted as
318 matches. We repeated the process with the first and second junctions swapped. To
319 simplify the analysis, two cases were ignored: colonies with matches between more than
320 two insertions (~12% of match cases), and insertion pairs where the proximal read of one
321 insertion was a match to multiple distal reads of the other insertion with different
322 orientations (~3% of match cases). We then took the distance to the closest distal read,
323 and counted the cases by orientation and distance, as before:

	0-10 bp	11-100 bp	101 bp - 1 kb	1-10 kb
Inner-cassette (toward-facing)	11	5072	5787	289
Outer-cassette (away-facing)	28	140	152	82
Same-direction	6	185	283	195

324

325 Note that the distances are expected to be higher in this case, because if we are
326 looking at a case of two sides of one insertion with a junk fragment, the distal read will
327 be a variable distance away from the junk-genome junction which is the actual insertion
328 location. So even for insertions with no genomic deletion/duplication, the distance
329 between the proximal read on one side and the nearest distal read on the other side will
330 not be 0 bp.

331 The number of inner-cassette cases up to 1 kb is more than 10x larger than the
332 number of same-direction cases, so these insertion pairs are extremely likely to be two
333 sides of one insertion with a junk fragment (and possibly a genomic deletion). Thus, all
334 the pairs in this category were identified as confidence level 2, which are extremely likely
335 to be correctly mapped.

336 The number of inner-cassette cases with a distance of 1-10 kb and the number of
337 outer-cassette cases with a distance of 0-10 bp is also higher than the expected 50% of the
338 same-direction cases, suggesting that many of them are also two sides of the same
339 insertion, but the differences are less dramatic and thus the number of false positives
340 would be too high for us to be comfortable identifying all these pairs as confidence level
341 2.

342 The insertion position information for junk fragment sides of confidence level 2
343 insertions originally reflected the junk fragment rather than the actual genomic insertion
344 position. We corrected it to show the nearest distal read matching the non-junk side: the
345 flanking sequence and position was changed to that of that distal read; the
346 “LEAPseq_distance” field was changed to the longest distance between two distal reads
347 that mapped to the presumed real insertion position (i.e. to the same region as the
348 proximal read of the insertion junction from the other side); the remaining LEAPseq
349 fields were likewise changed to reflect the numbers of distal reads and positions mapped
350 to the presumed real insertion position. For confidence level 2 insertions, the
351 “if_both_sides” column was set to “with-junk”; for the sides with a junk fragment, the
352 “if_fixed_position” column was set to “yes_nearest_distal”, and for the sides without a
353 junk fragment it was kept as “no”.

354 The confidence level 1 and 2 insertions (counting only the non-junk side of the
355 confidence level 2 insertions) appear to be of high quality (Supplementary Fig. 2h).

356

357 **Categorizing the remaining insertions and correcting junk fragments (confidence**
358 **levels 3 and 4).** After identifying the highest-confidence insertion junctions, i.e. those
359 with two matching sides of the same insertion, we sought to separate the remaining
360 insertions (with only one side mapped) into a set with a high likelihood of having
361 correctly mapped genomic insertion positions and a set with insertion positions likely to
362 reflect junk fragments. We considered two factors to separate these two sets: (1) the
363 percentage of read pairs that map to the same locus, and (2) the longest distance spanned
364 by such a read pair (Supplementary Fig. 2, i and j). We decided to solely use the first

365 factor based on the fact that nearly all of the insertions with low distances but high
366 percentage of read pairs mapped to the same locus were ones with relatively few LEAP-
367 Seq reads, indicating that their short distances spanned are likely due to them having few
368 reads (and thus a lower chance of a long read) rather than to a junk fragment. Therefore
369 we decided to use the percentage of read pairs mapping to the same locus as the only
370 factor in distinguishing the higher and lower confidence insertion sets, because that factor
371 is independent of the number of reads. To determine what cutoff would be appropriate,
372 we took advantage of the already known confidence level 1 insertions. We calculated the
373 fraction of confidence level 1 pairs among all the colonies with exactly two insertions
374 (two insertions are required for a confidence level 1 pair) as an approximate lower bound
375 on the number of correctly mapped insertions. Over the entire dataset, this fraction is
376 65%; when calculated only on insertions with at least 50% read pairs mapping to the
377 same locus, it's 78%; for insertions with at least 60%, 70%, 80% and 90% read pairs
378 mapping to the same locus, it is 79%. Thus it is clear that using a cutoff anywhere in the
379 50-90% range significantly improves the quality of the dataset, regardless of the exact
380 position of the cutoff. This makes sense, because the 50-90% range constitutes a very
381 small fraction of all insertions. We opted to use 60% as the cutoff for confidence level 3,
382 i.e. insertions with only one mapped side but with LEAP-Seq data indicating very likely
383 correct mapping.

384 The remaining insertions, with below 60% read pairs mapping to the same locus
385 and thus with the proximal LEAP-Seq read likely to be part of a junk fragment, were
386 analyzed further to identify the most likely true insertion position. The same analysis
387 was applied to all insertions with the proximal LEAP-Seq read with no genomic

388 alignment (possibly due to a very short junk fragment resulting in the 30 bp proximal
389 read being a hybrid of the junk fragment sequence and genomic sequence from the real
390 insertion position, or simply due to PCR or sequencing errors yielding an unmappable
391 sequence), or with multiple equally good genomic alignments (which could be derived
392 from the real genomic location, but in a non-unique region of the genome, requiring the
393 use of distal reads to determine the correct insertion location), or mapped to the insertion
394 cassette (indicating a second cassette fragment inserted between the first cassette and the
395 genome, which can be treated the same way as a junk genomic DNA fragment).

396 In order to determine the best method of identifying the true insertion location
397 based on the full distal LEAP-Seq read data, we grouped the distal LEAP-Seq reads for
398 each insertion into regions no more than 3 kb in size. For each such group, we calculated
399 three measures that we thought might be the best method of identifying the real insertion
400 location: (1) the number of reads in the group, (2) the number of unique genomic
401 positions to which reads in the group were mapped, and (3) the distance spanned by the
402 reads. LEAP-Seq reads mapped to the insertion cassette, or with no unique mapping to
403 the genome, were excluded. In order to determine which method was the best, we used
404 the junk fragment sides of confidence level 2 insertions, since for those the distal reads
405 corresponding to the true genomic insertion locations had already been determined by an
406 independent method (i.e. by matching the proximal read of the other side of the
407 insertion). For each of the three methods listed above, the insertion location predicted by
408 the method was compared to the known insertion location of each confidence level 2
409 insertion with a junk fragment. The results were as follows: 90% of the known insertion
410 positions were correctly predicted by taking the region with the most total distal reads,

411 84% by taking the region with the most unique mapping positions, and 84% by taking the
412 region with the longest distance spanned by the reads. Thus, the total number of distal
413 reads was chosen as the most likely measure to yield the correct genomic insertion
414 position of insertions with a junk fragment.

415 This method was then applied to all the insertions listed in the previous paragraph,
416 yielding the most likely true location for each insertion; insertions with only a single
417 LEAP-Seq distal read in each region were excluded, because one read did not provide
418 enough data to determine the insertion position with any confidence. For some insertions,
419 the region with the most distal LEAP-Seq reads also included the proximal LEAP-Seq
420 read - in those cases, the original insertion position based on the proximal LEAP-Seq
421 read was left unchanged. It is still possible that this position reflects a relatively long junk
422 fragment rather than the true genomic insertion position, but we did not have enough data
423 to distinguish those cases from high confidence. Likewise, it is possible that the corrected
424 position with the most distal LEAP-Seq reads that do not match the proximal read reflects
425 a second long junk fragment inserted after the first junk fragment which contains the
426 proximal read (we know that insertions with multiple junk fragments can happen), but
427 given the limited length of Illumina-sequenced LEAP-Seq reads, we cannot detect those
428 cases with certainty, and have to limit ourselves to finding putative insertion positions
429 that have a reasonably high probability of being correct.

430 Additionally, it turned out that many corrected positions for insertions originally
431 mapped to the insertion cassette did not appear to be high-quality, with only a small
432 fraction of distal reads mapped to the putative real insertion position. After looking at
433 several such cases in detail, we concluded that they had not been analyzed correctly.

434 They had single LEAP-Seq reads mapped to multiple distant locations on many
435 chromosomes, compared to 100+ reads mapped to many cassette locations, with the
436 putative real insertion position identified due to two or three single LEAP-Seq reads
437 mapped close together on one chromosome. The uniformly low read numbers of genome-
438 mapped reads compared with the high read numbers of cassette-mapped reads led us to
439 conclude that the genome-mapped reads were results of PCR or sequencing errors or
440 other artifacts, rather than being derived from real LEAP-Seq products, which should
441 usually yield more than one read. Thus, those appeared to be cases where no LEAP-Seq
442 products sequenced past the additional cassette fragment - this could be expected,
443 because the full cassette is >2.2 kb in length, whereas vanishingly few LEAP-Seq reads
444 are over 1.5 kb. In contrast, junk genomic DNA fragments are mostly smaller than 500 bp
445 and all identified ones were below 1 kb, so this problem would not be expected to be
446 common in genomic junk fragment cases. Indeed a cluster of low-matching-read-percent
447 insertions was not observed in the corrected insertion positions in that category. We
448 decided to exclude this category of incorrectly mapped insertions by only including
449 corrected originally cassette-mapped insertions if >50% of the distal LEAP-Seq reads
450 mapped to the putative correct insertion location.

451 All the insertions included in the final results of this analysis were annotated as
452 confidence level 4. The final confidence level 4 insertions are of a relatively high quality
453 (Supplementary Fig. 2j). The positions, flanking sequences and LEAP-Seq data of the
454 corrected confidence level 4 insertions in Supplementary Table 5 were changed to reflect
455 the new insertion position, in the same way as for the junk fragment sides of the
456 confidence level 2 insertions above. An additional complication of the new corrected

457 insertion positions was presented by the fact that the position of the nearest distal LEAP-
458 Seq read is always at some distance from the true insertion position, depending on the
459 length of the LEAP-Seq read. We attempted to correct for this by using confidence level
460 1 insertions to determine the average distance between the proximal read (reflecting the
461 true insertion position) and the nearest distal read, separately for 5' and 3' datasets,
462 depending on the total number of LEAP-Seq reads for the insertion (binned into ranges:
463 1, 2, 3, 4-5, 6-10, 11-20, 21+ total reads). For each confidence level 4 insertion with a
464 corrected position, the position was further adjusted by the average distance for the
465 correct side and number of reads as calculated above. This distance was appended as a
466 number to the value in the "if_fixed_position" field for each insertion in Supplementary
467 Table 5.

468

469 **Insertion verification PCR.** The PCR reactions were performed in two steps to verify
470 the insertion site² (Supplementary Table 6): (1) Genomic locus amplification: genomic
471 primers that are ~1 kb away from the flanking genomic sequence reported by LEAP-Seq
472 were used to amplify the genomic locus around the flanking sequence. If wild type
473 produced the expected PCR band but the mutant did not produce it or produced a much
474 larger product, this indicated that the genomic locus reported by LEAP-Seq may be
475 disrupted by the insertional cassette and we proceeded to the second step; (2) Genome-
476 cassette junction amplification: one primer binding to the cassette (omj913,
477 GCACCAATCATGTCAAGCCT, for the 5' side and omj944,
478 GACGTTACAGCACACCCTTG, for the 3' side) and the other primer binding to
479 flanking *Chlamydomonas* genomic DNA (one of the genomic primers from the first step)

480 were used to amplify the genome-cassette junction. If the mutant produced a PCR band
481 with expected size that was confirmed by sequencing but wild type did not produced the
482 expected PCR band, we categorized this insertion as “confirmed.” In some mutants,
483 genomic primers surrounding the site of insertion did not yield any PCR products in wild
484 type or the mutant even after several trials, possibly due to incorrect reference genome
485 sequence or local PCR amplification difficulties. These cases were grouped as “failed
486 PCR” and were not further analyzed.

487 72 mutants (24 insertions each for confidence levels 1 and 2, confidence level 3
488 and confidence level 4) were chosen randomly from the library and tested. The genomic
489 DNA template was prepared from a single colony of each mutant using the DNeasy Plant
490 Mini Kit (69106, Qiagen). The PCRs were performed using the Taq PCR core kit
491 (201225, Qiagen) as described before¹. PCR products of the expected size were verified
492 by Sanger sequencing.

493

494 **Southern blotting.** Southern blotting was performed as previously described in detail².
495 Genomic DNA was digested with *StuI* enzyme (R0187L, New England Biolabs) and
496 separated on a 0.7% Tris-borate-EDTA (TBE) agarose gel. The DNA in the gel was
497 depurinated in 0.25 M HCl, denatured in a bath of 0.5 M NaOH, 1M NaCl, neutralized in
498 a bath of 1.5 M Tris-HCl, pH 7.4, 1.5 M NaCl, and finally transferred onto a Zeta-probe
499 membrane (1620159, Bio-Rad) overnight using the alkaline transfer protocol given in the
500 manual accompanying the membrane. On the next day, the membrane was gently washed
501 with saline-sodium citrate (2xSSC: 0.3 M NaCl, 0.03 M sodium citrate), dried with paper
502 towel, and UV cross-linked twice using the Stratalinker1800 (Stratagene). For probe

503 generation, the *AphVIII* gene on CIB1 was amplified using primers oMJ588
504 (GACGACGCCCTGAGAGCCCT) and oMJ589
505 (TTAAAAAATTCGTCCAGCAGGCG). The PCR product was purified and labeled
506 according to the protocol of Amersham Gene Images AlkPhos Direct Labeling and
507 Detection System (RPN3690, GE Healthcare). The membrane was hybridized at 60°C
508 overnight with 10 ng probe/mL hybridization buffer. On the next day, the membrane was
509 washed with primary and secondary wash buffers and then visualized using a CL-
510 XPosure film (34093, Thermo Fisher).

511

512 **Analyses of insertion distribution and identification of hot/cold spots.** A mappability
513 metric was defined to quantify the fraction of all possible flanking sequences from any
514 genomic region that can be uniquely mapped to that region¹. Calculation of mappability,
515 hot/cold spot analysis and simulations of random insertions were performed as described
516 previously¹, except that a 30 bp flanking sequence lengths instead of a mix of 20 bp and
517 21 bp was used (because we now use 30 bp flanking sequence data derived from LEAP-
518 Seq, rather than 20/21 bp ChlaMmeSeq sequences), and the v5.5 *Chlamydomonas*
519 genome instead of the v5.3 genome was used¹⁰. This analysis was done on the original
520 full set of mapped insertions, to avoid introducing bias from the choice of mutants into
521 the consolidated set. The hot/cold spot analysis was performed on confidence level 1
522 insertions only, to avoid introducing bias caused by junk fragments and their imperfect
523 correction. The full list of statistically significant hot/cold spots is provided in
524 Supplementary Table 7.

525

526 **Identification of underrepresented gene ontology (GO) terms.** For each GO category,
527 we calculated the total number of insertions in all genes annotated with the GO term and
528 the total mappable (mappability defined in the Supplementary Note) length of all such
529 genes, and compared them to the total number of insertions in and total mappable length
530 of the set of flagellar proteome genes¹¹. We compared these numbers using Fisher's exact
531 test, and did correction for multiple comparisons¹² to obtain the false discovery rate
532 (FDR). This analysis was done on the original full set of mapped insertions to avoid
533 introducing bias from the choice of mutants into the consolidated set. We decided to use
534 the flagellar proteome as the comparison set because flagellar genes are very unlikely to
535 be essential; we did not use intergenic insertions or the entire genome because we know
536 that the overall insertion density differs between genes and intergenic regions. The
537 statistically significant results are listed in Supplementary Table 8.

538

539 **Prediction of essential genes.** To predict essential genes in *Chlamydomonas*, we sought
540 to generate a list of genes that have fewer insertions than would be expected randomly
541 was generated. Among them, those with 0 insertion are considered candidate essential
542 genes.

543 To achieve these, for each gene, we calculated the total number of insertions in
544 that gene and the total mappable length of that gene, and compared them to the total
545 number of insertions in and total mappable length of the set of flagellar proteome genes¹¹,
546 as what we have performed on each GO category. The resulting list of genes with
547 statistically significantly fewer insertions than expected is discussed in the
548 Supplementary Note and shown in Supplementary Table 9: this includes 203 genes with

549 no insertions, and 558 genes with at least one insertion. However, only genes 5 kb or
550 longer yield an false discovery rate (FDR) of 0.05 or less when they have no insertions -
551 our overall density of insertions is not high enough to detect smaller essential genes.

552

553 **Pooled Screens.** Library plates that were replicated once every four weeks onto fresh
554 medium were switched to a 2-week replication interval to support uniform colony growth
555 before pooling. Cells were pooled from 5-days-old library plates: first, for each set of
556 eight agar plates, cells were scraped using the blunt side of a razor blade (55411-050,
557 VWR) and resuspended in 40 mL liquid TAP medium in 50-mL conical tubes. Second,
558 cells clumps were broken up by pipetting, using a P200 pipette tip attached to a 10-mL
559 serological pipette. In addition, cells were pipetted through a 100 μ m cell strainer
560 (431752, Corning). Third, these sub-pools were combined as the master pool representing
561 the full library.

562 The master pool was washed with TP, and resuspended in TP. Multiple aliquots
563 of 2×10^8 cells were pelleted by centrifugation (1,000g, 5 min, room temperature) and the
564 supernatant was removed by decanting. Some aliquots were used for inoculation of
565 pooled cultures, whereas other aliquots were frozen at -80 °C as initial pool samples for
566 later barcode extraction to enable analysis of reproducibility between technical replicates.
567 For pooled growth, 20 L TAP or TP in transparent Carboy containers (2251-0050,
568 Nalgene) were inoculated with the initial pool to a final concentration of 2×10^4 cells/mL.
569 Cultures were grown under 22°C, mixed using a conventional magnetic stir bar and
570 aerated with air filtered using a 1 μ m bacterial air venting filter (4308, Pall Laboratory).
571 The TAP culture was grown in dark. For the two replicate TP cultures, the light intensity

572 measured at the surface of the growth container was initially $100 \mu\text{mol photons m}^{-2} \text{s}^{-1}$,
573 and then increased to $500 \mu\text{mol photons m}^{-2} \text{s}^{-1}$ after the culture reached $\sim 2 \times 10^5$
574 cells/mL. When the culture reached the final cell density of 2×10^6 cells/mL after 7
575 doublings, 2×10^8 cells were pelleted by centrifugation (1,000g, 5 min, room temperature)
576 for DNA extraction and barcode sequencing.

577

578 **Barcode sequencing and data analysis for pooled screens.** Barcodes were amplified
579 and sequenced using the Illumina HiSeq platform as performed on the combinatorial
580 super-pools in library mapping (Supplementary Fig. 1f). Initial reads were trimmed using
581 cutadapt version 1.7.1⁸. Sequences were trimmed using the command "cutadapt -a <seq>
582 -e 0.1 -m 21 -M 23 input_file.gz -o output_file.fastq", where seq is
583 GGCAAGCTAGAGA for 5' data and TAGCGCGGGGCGT for 3' data. Barcodes were
584 counted by collapsing identical sequences using "fastx_collapser"
585 (http://hannonlab.cshl.edu/fastx_toolkit). The barcode read counts for each dataset were
586 normalized to a total of 100 million (Supplementary Table 10).

587 For evaluation of the quantitiveness of our barcode sequencing method,
588 barcodes obtained from two technical replicate aliquots of the same initial pool were
589 compared in read counts (Supplementary Fig. 5a). Barcodes obtained from the two TP-
590 light cultures at the end of growth were compared to assess consistency between
591 biological replicates (Fig. 3b).

592 To detect deficiency in photosynthetic growth, we compared mutant abundances
593 in TP-light with TAP-dark at the end of growth (Fig. 3c). As a quality control, different

594 barcodes in the same mutant were compared in the ratio of the TP-light read count to
595 TAP-dark read count. Highly consistent ratios were observed (Supplementary Fig. 5b).

596 For the identification of photosynthetically deficient mutants, each barcode with
597 at least 50 normalized reads in the TAP-dark dataset was classified as a hit if its ratio of
598 normalized TP-light:TAP-dark read counts was 0.1 or lower, or a non-hit otherwise. The
599 fraction of hit barcodes was 3.3% in replicate 1 and 2.9% in replicate 2. These barcodes
600 represent 2,638 and 2,369 mutants showing a growth defect in the TP-light-I and TP-
601 light-II replicates, respectively. A total of 3,109 mutants covering 2,599 genes showed a
602 growth defect in either of the TP-light sample.

603

604 **Identification and annotation of the hit genes from the screen.** To evaluate the
605 likelihood that a gene is truly required for photosynthesis, we counted the number of
606 alleles for this gene with and without a phenotype, including exon/intron/5'UTR
607 insertions. If the insertion was on the edge of one of those features, or in one of the
608 features in only one of the splice variants, it was still counted. We excluded alleles with
609 insertions in the 3' UTRs, which we observed to less frequently cause a phenotype (Fig.
610 3, d and e). In cases of multiple barcodes in the same mutant (likely two sides of one
611 insertion), the one with a higher TAP-dark read count was used for the calculation of
612 normalized TP-light:TAP-dark read counts, to avoid double-counting a single allele. For
613 each gene, a *P* value was generated using Fisher's exact test comparing the numbers of
614 alleles in that gene with and without a phenotype to the numbers of all insertions in the
615 screen with and without a phenotype (Supplementary Table 11). A false discovery rate
616 (FDR) correction was performed on the *P* values using the Benjamini-Hochberg

617 method¹², including only genes with at least 2 alleles present in the screen. Thus, genes
618 with a single allele have a *P* value but lack a FDR.

619 This process was performed for both TP-light replicates. The list of higher-
620 confidence genes was generated by taking genes with FDR of 0.27 or less in either
621 replicate - this threshold includes all genes with 2 hit alleles and 0 non-hit alleles. The
622 resulting list of hits included 37 genes in replicate 1, 34 in replicate 2, 44 total. The FDR
623 values for the higher-confidence genes in both replicates are shown in Tables 1 and 2.
624 Additionally, the list of lower-confidence genes was generated by taking genes with a *P*
625 value of 0.058 or less – this value was chosen to include genes with only one allele with a
626 phenotype and no alleles without a phenotype, but to exclude genes with one allele with
627 and one without a phenotype. The resulting list included 264 genes total (210 in replicate
628 1, 196 in replicate 2).

629 One gene in the original higher-confidence list and four genes in the original
630 lower-confidence list encode subunits of the plastidic pyruvate dehydrogenase. Mutants
631 in these genes require acetate to grow because they cannot generate acetyl-CoA from
632 pyruvate but can generate acetyl-CoA from acetate. This requirement for acetate, rather
633 than a defect in photosynthesis, likely explains why mutants in this gene showed a
634 growth defect in TP-light¹³. Removal of these genes led to a final list of 43 higher-
635 confidence genes and 260 lower-confidence genes (Fig. 3f, Tables 1 and 2, and
636 Supplementary Table 12).

637 We identified 65 (22 higher-confidence and 43 lower-confidence) out of the 303
638 hit genes as “known” genes based on genetic evidence: mutation of this gene in
639 *Chlamydomonas* or another organism caused a defect in photosynthesis. Among the

640 remaining 238 “candidate” genes (21 higher-confidence ones and 217 lower-confidence
641 ones), some genes appear to be related to photosynthesis because of their predicted
642 chloroplast localization or evolutionary conservation among photosynthetic organisms¹⁴,
643 despite lack of solid genetic evidence. For three of the candidate genes (*CGL59*, *CPL3*,
644 and *VTE5*), mutants with insertions adjacent to them were previously found to be acetate-
645 requiring or hypersensitive to oxidative stress in the chloroplast¹³.

646

647 **Analysis of candidate gene enrichment in reported transcriptional clusters related to**
648 **photosynthesis.** Two transcriptome datasets in *Chlamydomonas* were used in this
649 analysis: a diurnal regulation study¹⁵ and a dark-to-light transition study¹⁶. For the first
650 one, we chose the diurnal cluster 4 in the study that had photosynthesis-related genes
651 enriched in it¹⁵. For the second one, we chose the genes upregulated upon transition to
652 light¹⁶. In each case, the number of candidate genes included and not included in the
653 regulated gene sets was compared to the total number of *Chlamydomonas* genes included
654 and not included in the cluster, using Fisher's exact test. The resulting *P* values were
655 FDR-adjusted using the Benjamini-Hochberg method¹².

656

657 **Molecular characterization of the *cpl3* mutant.** Mutant genotyping PCRs were
658 performed as previously described². Sequences of primers represented in Supplementary
659 Fig. 6a are g1: CCGTCGTCACCTTGC-TACAAC, g2: CGTAGTTGCAAGGGGTGTTT,
660 c1: GACGTTACAGCACACCCTTG. To complement the *cpl3* mutant, the wild-type
661 *CPL3* gene was PCR amplified and cloned into the vector pRAM118 vector that contains
662 the *aph7* gene¹⁷, which confers resistance to hygromycin B. In this construct, the

663 expression of *CPL3* is under the control of the *PSAD* promoter. The construct was
664 linearized before being transformed into the *cpl3* mutant. Transformants were robotically
665 arrayed and assayed in colony sizes in the presence and absence of acetate respectively
666 (Supplementary Fig. 6, c and d). Three representative lines that showed rescued
667 photosynthetic growth were used in further phenotypic analyses (Fig. 4).

668

669 **Analyses of growth, chlorophyll, and photosynthetic electron transport.** For all
670 physiological and biochemical characterizations of *cpl3* below, we grew cells
671 heterotrophically in the dark to minimize secondary phenotypes due to defects in
672 photosynthesis.

673 For spot assays, cells were grown in TAP medium in dark to log phase to around
674 10^6 cells per mL. Cells were washed in TP and spotted onto solid TAP medium and TP
675 medium respectively. The TAP plates were incubated in dark for 12 d before being
676 imaged. The TP plates were incubated under $30 \mu\text{mol photons m}^{-2} \text{s}^{-1}$ light for 1 d, 100
677 $\mu\text{mol photons m}^{-2} \text{s}^{-1}$ light for 1 d, and then $500 \mu\text{mol photons m}^{-2} \text{s}^{-1}$ light for 4 d.

678 Chlorophyll *a* and *b* concentrations were measured as previously described¹⁸
679 using TAP-dark grown cells. We used TAP-dark-grown instead of TP-light-grown cells
680 for chlorophyll analyses, photosynthetic performance analyses, microscopy, proteomics,
681 and western blots (below) to avoid observing secondary effects due to the photosynthetic
682 defects of the *cpl3* mutants.

683 To measure photosynthetic electron transport rate, TAP-dark grown cells were
684 collected, re-suspended in fresh TAP medium, and dark acclimated for 20 min. Cells
685 were then measured in chlorophyll fluorescence under a series of increasing light

686 intensities using the “Light Curve” function on a DUAL-PAM-100 fluorometer (Walz).
687 PSII quantum yield (Φ_{PSII}) was quantified as previously described¹⁹. Relative electron
688 transport rate (rETR) was calculated according to the following equation $rETR = \Phi_{PSII} \times$
689 I. I represents the emitted irradiance.

690

691 **Microscopy.** Cells were grown under the TAP-dark condition to log phase and
692 concentrated ten-fold before microscopic analysis. Aliquots were deposited at the corner
693 of a poly-L-lysine coated microslide well (Martinsried) and spread over the bottom of the
694 well by overlaying with TAP-1% agarose at low temperature (<30°C), to minimize cell
695 motion during image acquisition. Cells were imaged at room temperature though a Leica
696 TCS SP5 laser scanning confocal microscope and an inverted 100x NA 1.46 oil
697 objective. Chlorophyll fluorescence signal was generated using 514 nm excitation, and
698 650-690 nm collection. All images were captured using identical laser and magnification
699 settings (4x zoom and single-slice through the median plane of the cell). Composite
700 images (chlorophyll fluorescence overlay with bright field) were generated with Fiji²⁰.

701

702 **Proteomics.** TAP-dark-grown cells were collected by centrifugation and flash-frozen.
703 Proteins were extracted from the frozen pellets by resuspension in lysis buffer (6M
704 guanidium Hydrochloride, 10mM tris(2-carboxyethyl)phosphine, 40mM chloroacetamide,
705 100mM Tris pH8.5, 1x MS-Safe protease inhibitor, 1x Phosphatase inhibitor cocktail II),
706 grinding with liquid nitrogen, followed by sonication. Protein lysates were then digested
707 with trypsin (Promega) into peptides. Three biological replicates were processed for each
708 strain.

709 The samples were labeled with tandem mass tags (TMTs), multiplexed and then
710 fractionated before tandem mass spectrometry analyses. Briefly, each sample was labeled
711 with the TMT labeling reagent (Thermo Fisher) according to the manufacturer's
712 instructions. The samples were then mixed in equimolar amounts and desalted using
713 C18-stage tips²¹. The dried peptide mix was then separated using strong cation exchange
714 (SCX) stage-tips²² into four fractions. Each of the four fractions were then diluted with
715 1% trifluoroacetic acid (TFA) and separated into three fractions using SDB-RPS stage
716 tips. This procedure initially resulted in a total of 12 fractions. Fractions 1-3 (the children
717 of the first SCX fraction) were pooled together yielding 10 final fractions. Each final
718 fraction was diluted and injected per run using an Easy-nLC 1200 UPLC system (Thermo
719 Fisher). Samples were loaded onto a nano capillary column packed with 1.9 μm C18-AQ
720 (Dr. Maisch) mated to metal emitter in-line with a Fusion Lumos (Thermo Fisher).
721 Samples were eluted using a split gradient of 10-20% solution B (80% ACN with 0.1%
722 FA) in 32 min and 20-40% solution B in 92 min followed column wash at 100% solution
723 B for 10 min. The mass spectrometer was operated in a data-dependent mode with the
724 60,000 resolution MS1 scan (380-1500 m/z), AGC target of $4e5$ and max injection time
725 of 50ms. Peptides above threshold $5e3$ and charges 2-7 were selected for fragmentation
726 with dynamic exclusion after 1 time for 60 s and 10 ppm tolerance. MS1 isolation
727 windows of 1.6m/z, MS2 isolation windows 2 and HCD NCE of 55% were selected.
728 MS3 fragments were detected in the Orbitrap at 50,000 resolution in the mass range of
729 120-500 with AGC $5e4$ and max injection time of 86 ms. The total duty cycle was set to
730 3.0 sec.

731 Raw files were searched with MaxQuant²³, using default settings for MS3 reporter
732 TMT 10-plex data. Files were searched against sequences of nuclear, mitochondrial, and
733 chloroplast-encoded *Chlamydomonas* proteins supplemented with common
734 contaminants^{10,24,25}. Raw files were also analyzed within the Proteome Discoverer
735 (Thermo Fisher) using the Byonic²⁶ search node (Protein Metrics). Data from Maxquant
736 and Proteome Discoverer were combined in Scaffold Q+ (Proteome Software Inc.),
737 which was used to validate MS/MS based peptide and protein identifications. Peptide
738 identifications were accepted if they could be established at greater than 80.0%
739 probability by the Scaffold Local FDR algorithm. Protein identifications were accepted if
740 they could be established at greater than 96.0% probability and contained at least 2
741 identified peptides. Scaffold Q+ un-normalized data were exported in the format of the
742 log₂ value of the reporter ion intensities, which reflect the relative abundances of the
743 same protein among different samples multiplexed. Each sample was then normalized to
744 a median of 0 (by subtracting the original median from the raw values, since the values
745 are log₂). For each gene, for each pair of samples, the normalized log₂ intensity values
746 from the three replicates of one sample were compared against those for the other sample
747 using a standard *t*-test. The resulting *P* values were adjusted for multiple testing¹²,
748 yielding a false discovery rate (FDR) for each gene in each pair of samples. We note that
749 our calculation of FDR does not take into account the spectral count of each protein
750 (provided in Supplementary Table 14), which is related to the absolute abundance of the
751 protein and impacts the accuracy of proteomic measurements. Specifically, proteins with
752 a low spectral count are likely of low abundance in cells and often exhibit a large
753 variation in the intensity value between the biological replicates.

754

755 **Western blotting.** TAP-dark grown cells were pelleted by centrifugation, resuspended in
756 an extraction buffer containing 5 mM HEPES-KOH, pH 7.5, 100 mM dithiothreitol, 100
757 mM Na₂CO₃, 2% (w/v) SDS, and 12% (w/v) sucrose, and lysed by boiling for 1 min.
758 Extracted proteins were separated on SDS-PAGE (12% precast polyacrylamide gels, Bio-
759 Rad) using tubulin as a loading and normalization control. Polypeptides were transferred
760 onto polyvinylidene difluoride membranes using a semidry blotting apparatus (Bio-Rad)
761 at 15 volts for 30 minutes. For western blot analyses, membranes were blocked for 1 h at
762 room temperature in Tris-buffered saline-0.1% (v/v) Tween containing 5% powdered
763 milk followed by a 1 h incubation of the membranes at room temperature with the
764 primary antibodies in Tris-buffered saline-0.1% (v/v) Tween containing powdered milk
765 (3% [w/v]). Primary antibodies were diluted according to the manufacturer's
766 recommendations. All antibodies were from Agrisera and the catalog numbers for the
767 antibodies against CP43, PsaA, ATPC, and α -tubulin were AS11-1787, AS06-172-100,
768 AS08-312, and AS10-680, respectively. Proteins were detected by enhanced
769 chemiluminescence (K-12045-D20, Advansta) and imaged on a medical film processor
770 (Konica) as previously described².

771

772 **Code availability.** All programs written for this work are deposited at

773 <https://github.com/Jonikas-Lab/Li-Patena-2018>.

774

775 **Supplementary Note**

776 **Accuracy of insertion mapping and number of insertions per mutant.** In

777 Chlamydomonas insertional mutants, short “junk fragments” of genomic DNA (likely
778 from lysed cells) are often inserted between the cassette and flanking genomic DNA¹.
779 The difficulty in distinguishing these junk fragments from true flanking genomic DNA
780 can lead to inaccurate mapping of the insertion to a genomic location^{1,2}. Additionally,
781 some cassettes are truncated during insertion, preventing mapping of the flanking
782 sequence on one side. We sought to help users prioritize mutants for characterization by
783 classifying insertions into categories that reflect our confidence in the mapping accuracy,
784 based on two criteria: (1) whether flanking sequences from both sides of the cassette
785 mapped to the same genomic region; and (2) whether the LEAP-Seq reads contained
786 sequences from multiple genomic regions, suggesting the presence of junk DNA
787 fragments inserted next to the cassette (Supplementary Fig. 3a and Supplementary Fig.
788 2f-j).

789 A confidence level of 1 was assigned to 19,015 insertions in which both cassette-
790 genome junctions mapped to the same genomic region and were free of junk fragments.
791 A confidence level of 2 was assigned to 5,665 insertions in which both cassette-genome
792 junctions mapped to the same genomic region, after correcting for the presence of a junk
793 fragment at one junction. A mapping confidence level of 3 was assigned to 36,600
794 insertions in which only one cassette-genome junction could be identified, with the
795 likelihood of junk DNA insertion determined to be low based on fewer than 40% of
796 LEAP-Seq reads containing sequence from multiple genomic regions. A mapping
797 confidence level of 4 was assigned to 13,643 insertions in which only one junction could

798 be identified, and that junction was likely to contain a junk fragment, or the flanking
799 sequence could not be mapped to a unique genomic location. The mapping for these
800 insertions was adjusted to reflect the most likely correct insertion site.

801 Approximately 95% of confidence level 1 and 2 insertions are mapped correctly
802 based on PCR validation of randomly chosen mutants, compared to ~73% of confidence
803 level 3 and ~58% of confidence level 4 (Supplementary Table 6; Methods).

804 Our bioinformatic analyses suggest that over 80% of the mutants harbor only one
805 mapped insertion (Supplementary Fig. 3b), consistent with Southern blot data from
806 randomly chosen mutants (Supplementary Fig. 3c).

807

808 **Deletions, duplications, and junk fragments associated with insertions are small.**

809 Random insertions in *Chlamydomonas* are sometimes also associated with deletions and
810 duplications of neighboring genomic DNA¹³. To further help users understand the quality
811 of mutants in this library, we characterized these deletions and duplications by examining
812 the sequences across both junctions of confidence level 1 insertions (Methods). Of these
813 insertions, 11% had no deletions or duplications, 74% harbored genomic deletions and
814 15% had genomic duplications. The great majority (98%) of genomic deletions were less
815 than 100 bp, but some were as large as 10 kb. While 98% of the genomic duplications
816 were shorter than 10 bp, some extended to 30bp (Supplementary Fig. 3, d and e). Both
817 the deletions and duplications likely resulted from non-homologous end joining repair
818 that occurs during cassette insertion²⁷. Additionally, examining the 651 insertions in
819 which a junk fragment separated two cassettes inserted in the same location allowed us to
820 estimate the typical junk fragment length. Most (73%) junk fragments were shorter than

821 300 bp, but some were as large as 1,000 bp (Supplementary Fig. 3f). If larger deletions,
822 duplications or junk fragments were present, they were not sufficiently frequent to allow
823 us to identify them reliably.

824

825 **Insertion sites are randomly distributed with mild cold spots and a small number of**
826 **hot spots.** While a random insertion model produced a distribution of insertion sites
827 broadly similar to the observed distribution (Fig. 1c and Supplementary Fig. 4a), we did
828 detect some cold spots and hot spots where insertion density differed significantly from
829 the random insertion model (Supplementary Fig. 4a; Supplementary Table 7; Methods).
830 Cold spots cover 26% of the genome and on average show a 48% depletion of insertions.
831 Hot spots cover 1.5% of the genome and contain 16% of insertions (Methods).

832 Hot spots fell into two distinct classes that differed in the local distribution of
833 insertions (Supplementary Fig. 4, b and c). In one class, dozens of insertions were found
834 within a region of 20-40 bp. In the other class, the insertions were distributed over a
835 much larger region of 200-1,000 bp. Our observations suggest that hot spots could be
836 caused by two distinct mechanisms; however, we did not observe a correlation between
837 specific features of the genome (e.g. sequence, exon, intron, UTR, mappability) and the
838 occurrence of either class of hot spots.

839

840 **Absence of insertions identifies over 200 genes potentially essential for growth under**
841 **the propagation conditions used.** Identification of essential genes in bacteria, fungi, and
842 mammals has revealed important molecular processes in these organisms^{3,28-30}. We
843 sought to take advantage of the very large set of mapped mutations in the library to

844 identify candidate essential *Chlamydomonas* genes based on the absence of insertions in
845 those genes (Methods). We note that our approach does not allow testing of gene
846 essentiality under all possible conditions. Therefore, it is likely that some of the candidate
847 essential genes we identify in this approach are required specifically for growth under our
848 propagation conditions, but not under all conditions. For example, mutants in respiratory
849 genes would be identified as essential if these mutants were not recovered under our
850 propagation conditions (in the dark on acetate media), although the same mutants could
851 have grown if recovery were under photosynthetic conditions.

852 Given our average density of insertions, we were able to detect a statistically
853 significant (FDR < 0.05) lack of insertions for genes with a mappable length greater than
854 5 kb. We identified 203 candidate essential genes (Supplementary Table 9). We caution
855 that this is a conservative list for two reasons: (1) if a gene has a mappable length smaller
856 than 5 kb and has no insertion, its underrepresentation is not statistically significant; (2)
857 some essential genes were not detected because there are insertions incorrectly mapped to
858 them.

859 Many of these predicted essential genes have homologs that have been shown to
860 be essential in other organisms. For example, Cre01.g029200 encodes a homolog of the
861 yeast cell cycle protease separase ESP1³¹, Cre12.g521200 encodes a homolog of yeast
862 DNA replication factor C complex subunit 1 RFC1³², and Cre09.g400553 encodes a
863 homolog of the yeast nutrient status sensing kinase Target of Rapamycin 2 TOR2³³. In
864 addition, we observed genes encoding proteins involved in acetate utilization or
865 respiration, such as acetyl-CoA synthetase/ligase³⁴ (Cre07.g353450) and components of
866 the mitochondrial F1F0 ATP synthase³⁵ (Cre15.g635850 and Cre07.g340350). As

867 discussed above, these genes may be essential under the conditions of library
868 propagation, in which acetate serves as the energy source.

869 We also observed genes on the list with nonessential homologs in other
870 organisms. One example is Cre13.g585301, which encodes monogalactosyldiacylglycerol
871 (MGDG) synthase and whose Arabidopsis homolog MGD1 is not essential³⁶. This can be
872 explained by the presence of two other isoforms of MGDG synthases in Arabidopsis but
873 not in Chlamydomonas³⁷. Comparison of our candidate Chlamydomonas essential genes
874 with those of other organisms can provide insights into evolutionary differences across
875 the tree of life.

876

877 **Deleterious mutations rather than differential chromatin configuration are the**
878 **major cause of insertion density variation.** One caveat for our above prediction of
879 essential genes is that the lack of insertions could be caused by low chromatin
880 accessibility at those loci to insertional mutagenesis. We reasoned that if chromatin
881 accessibility influenced insertion density, the 3' UTRs of these genes would also be less
882 represented; while if low insertion density primarily reflected essentiality, we would still
883 see many insertions in the 3' UTRs of these genes, because 3' UTR insertions typically
884 do not disrupt gene function (Fig. 3, d and e). For all genes in the genome, we observed
885 an insertion density of 1.1 insertions per mappable kb in exons and introns and 4.7
886 insertions per mappable kb in 3' UTRs. For the candidate essential genes, despite a lack
887 of insertions in exons and introns, the insertion density in 3' UTRs is 4.1 insertions per
888 mappable kb, similar to that of all genes. We thus conclude that low insertion density in

889 our candidate essential genes is largely caused by mutations that impair mutant fitness
890 instead of low chromatin accessibility to insertional mutagenesis.

891

892 **Disruption of *CPL3* is the cause of the photosynthetic deficiency in the *cpl3* mutant.**

893 We sought to confirm and characterize the *cpl3* insertion in detail. Our high-throughput
894 LEAP-Seq data suggested that *cpl3* contained an insertion of two back-to-back cassettes.
895 Specifically, the *cpl3* mutant contains two insertion junctions from 3' ends of two
896 cassettes in opposite orientations, within the *CPL3* gene. Junction 1 is confidence level 3
897 (no junk fragment), and junction 2 is confidence level 4 (with a junk fragment, corrected)
898 (Supplementary Fig. 6a). We successfully confirmed both junctions by PCR
899 (Supplementary Fig. 6b). Sequencing of the product from junction 2 revealed that the end
900 of the cassette has a 10-bp truncation and a 10-bp fragment of unknown origin inserted
901 between the cassette and the *CPL3* gene. The genomic flanking sequence of junction 2
902 overlaps with the flanking sequence in junction 1 by 2 bp. When we amplified across the
903 insertion site, *cpl3* yielded a product ~3 kb larger than the product from wild type
904 (Supplementary Fig. 6b). Based on these results, the most likely model for this insertion
905 is that two copies of the cassette (at least one truncated) inserted together into the *CPL3*
906 gene in opposite orientations, with a 2-bp genomic duplication at the site of insertion.

907 To confirm the involvement of *CPL3* in photosynthesis, we cloned *CPL3* genomic
908 DNA and transformed it into the *cpl3* mutant. Based on colony size, photoautotrophic
909 growth was rescued in approximately 14% of the transformants (Supplementary Fig. 6, c
910 and d), a percentage consistent with previous *Chlamydomonas* genetic studies³⁸. Three
911 rescued transformants, named comp1-3, were chosen at random for phenotypic

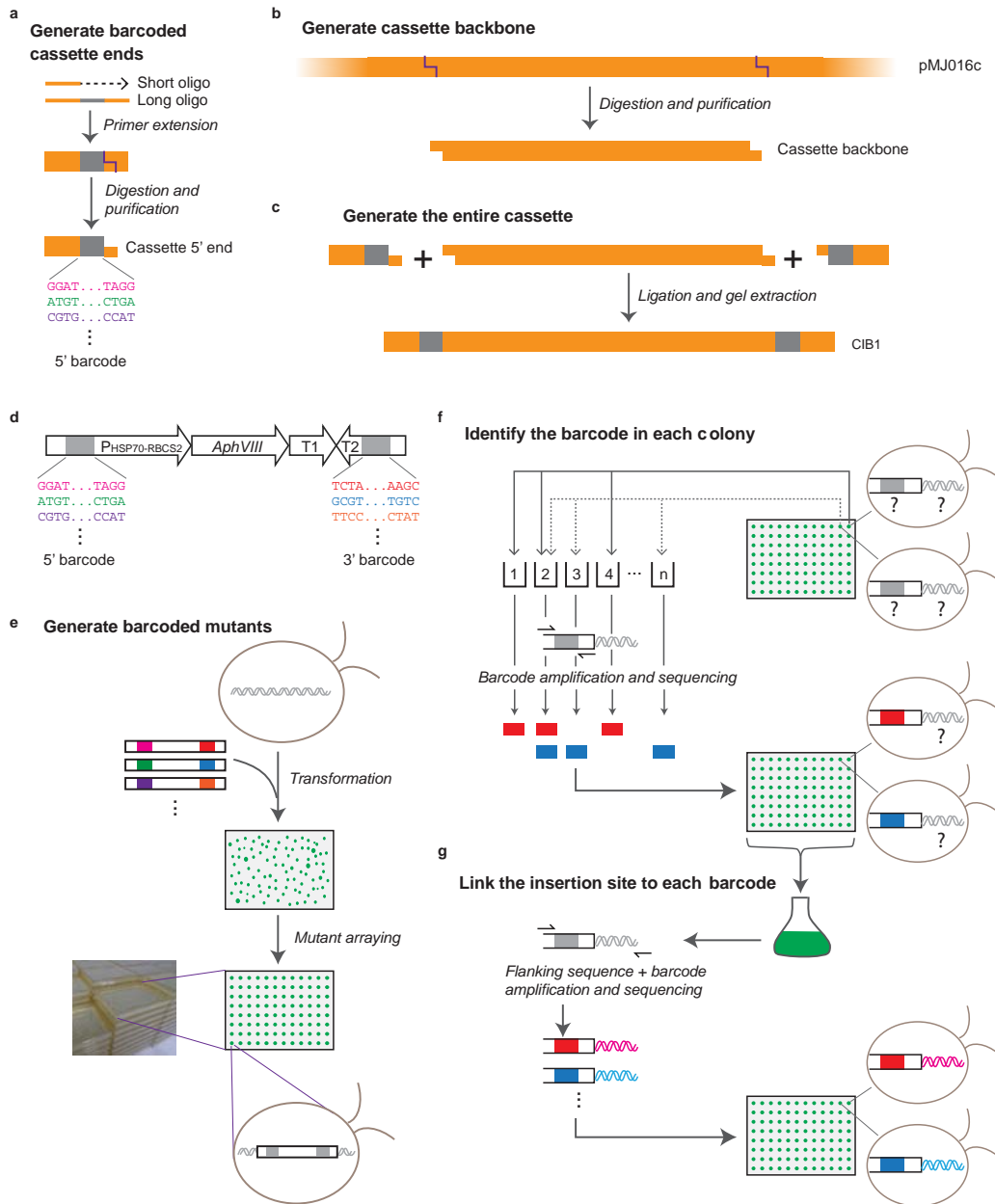
912 confirmation (Fig. 4b) and genotyping. In comp1-3, PCR across the insertion site of the
913 *cpl3* mutation with primers “g1 + g2” yielded ~1.2 kb products (expected size: 1,311 bp)
914 that indicate presence of wild-type *CPL3* sequence (from the wild-type *CPL3* in the
915 complementation construct), and weak ~4 kb bands consistent with the presence of the
916 original cassette insertion in *CPL3* (Supplementary Fig. 6b). The lower intensity of the
917 ~4 kb bands in these samples can be explained by preferential amplification of the
918 smaller template when multiple templates are present. To further confirm that comp1-3
919 still contained the original insertion in *CPL3*, we amplified the two insertion junctions in
920 the complemented lines with primers “g1 + c1” and “g2 + c1”. These genetic
921 complementation results demonstrate that the disruption of *CPL3* is the cause of the
922 growth defect of the mutant.

923

924

925 **Supplementary Figures**

926



927

928

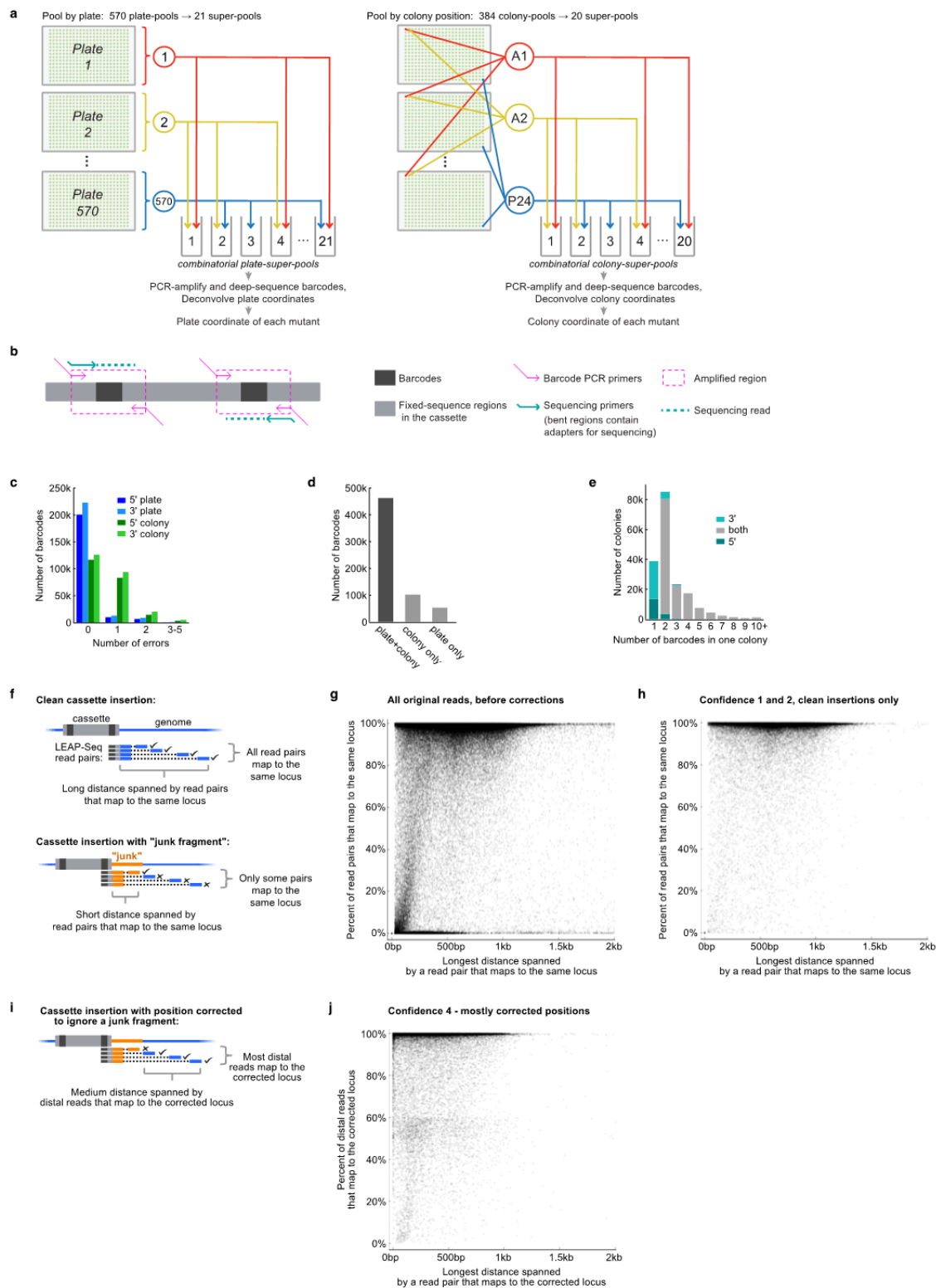
929 **Supplementary Fig. 1 | A pipeline was developed for generating barcoded cassettes**

930 **(a-d) and for generating an indexed and barcoded library of insertion mutants in**

931 **Chlamydomonas (e-g).** **a**, A long oligonucleotide primer containing a random sequence

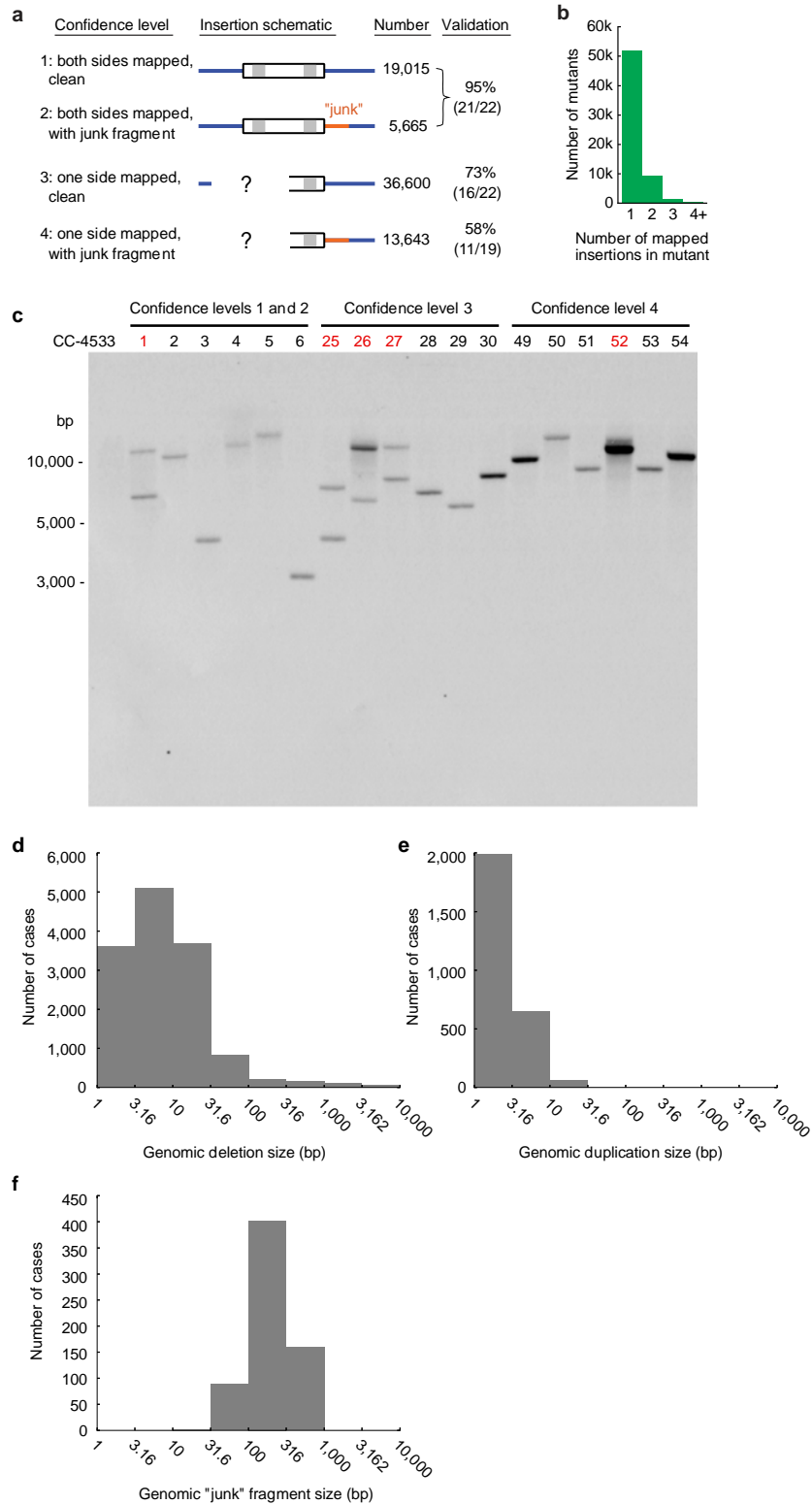
932 region (indicated in gray) was used as a template for the extension of a shorter

933 oligonucleotide primer (see Supplementary Table 1 for primer sequences). The resulting
934 double-stranded product contains a random sequence region (22 bp in length; termed
935 “barcode”). This product was restriction digested to generate a sticky end for subsequent
936 ligation. The above steps were performed to produce both the 5’ and the 3’ ends of the
937 cassette. The 5’ end of the cassette is shown as an example. **b**, The pMJ016c plasmid was
938 digested to yield the backbone of the cassette. **c**, The 5’ and 3’ ends of the cassette
939 generated above were ligated together with the cassette backbone to yield the cassette
940 CIB1. **d**, The components of the cassette CIB1 are shown. CIB1 contains the *HSP70*-
941 *RBCS2* promoter (with an intron from *RBCS2*), the *AphVIII* gene that confers resistance
942 to paromomycin, two transcriptional terminators (T1: *PSAD* terminator; T2: *RPL12*
943 terminator), and two barcodes (each 22 bp in length). **e**, Following transformation and
944 arraying of individual mutants, the sequence of the barcodes contained in each insertion
945 cassette was unique to each transformant but initially unknown for each colony. **f**,
946 Barcodes were amplified from combinatorial pools of mutants, sequenced, and traced
947 back to single colonies (Supplementary Fig. 2a-e; Methods). After this step, the barcode
948 sequence for each colony was known. For simplicity, only one side of the cassette is
949 shown. **g**, Barcodes and genomic sequences flanking the insertion cassettes were
950 amplified from a pool of the library. By pooled next-generation sequencing, the sequence
951 flanking each insertion cassette was paired with the corresponding barcode
952 (Supplementary Fig. 2f). The flanking sequences were used to determine the insertion site
953 in the genome. Because the colony location for each barcode was determined in the
954 previous step, insertion sites could then be assigned to single colonies.



956 **Supplementary Fig. 2 | Combinatorial pooling, barcode deconvolution to colony,**
957 **and determination of insertion sites. a,** To determine which plate each barcode was on,
958 each plate of mutants was pooled into one of 570 plate-pools. The plate-pools were then
959 further combinatorially pooled into 21 plate-super-pools, in such a way that each plate-
960 pool was in a unique combination of plate-super-pools. The barcodes present in each
961 plate-super-pool were determined by deep sequencing, and the barcodes were assigned to
962 plates based on the combination of plate-super-pools they were found in. A similar
963 process was applied to the colony positions of each barcode. Combining the plate and
964 colony data yielded a specific position for each barcode. **b,** The barcodes on the 5' and 3'
965 sides of the cassette were sequenced separately, each with a single-end Illumina read.
966 With the sequencing primers we used (indicated on the cassette), the reads start with the
967 barcode sequence and extend into the cassette. **c,** Most barcode colony positions were
968 identified with no errors, i.e. were found in one of the expected combinations of super-
969 pools. Some were found in a combination of super-pools that had one or more
970 differences from any expected combination, but the positions could still be identified due
971 to the redundancy built into our method. The much higher number of one-error cases in
972 the colony data compared to plate data is due to a loss of one of the colony-super-pools
973 for a significant fraction of the samples (Methods). **d,** Both a plate and a colony position
974 were identified for most barcodes. **e,** The number of barcodes mapped to an individual
975 colony varied, with 2 being the most common. For colonies with two mapped barcodes,
976 the large majority had one 5' and one 3' barcode, likely derived from two sides of one
977 cassette. **f,** LEAP-Seq reads are paired-end reads with the proximal read containing the
978 cassette barcode and immediate flanking genomic sequence, and the distal read

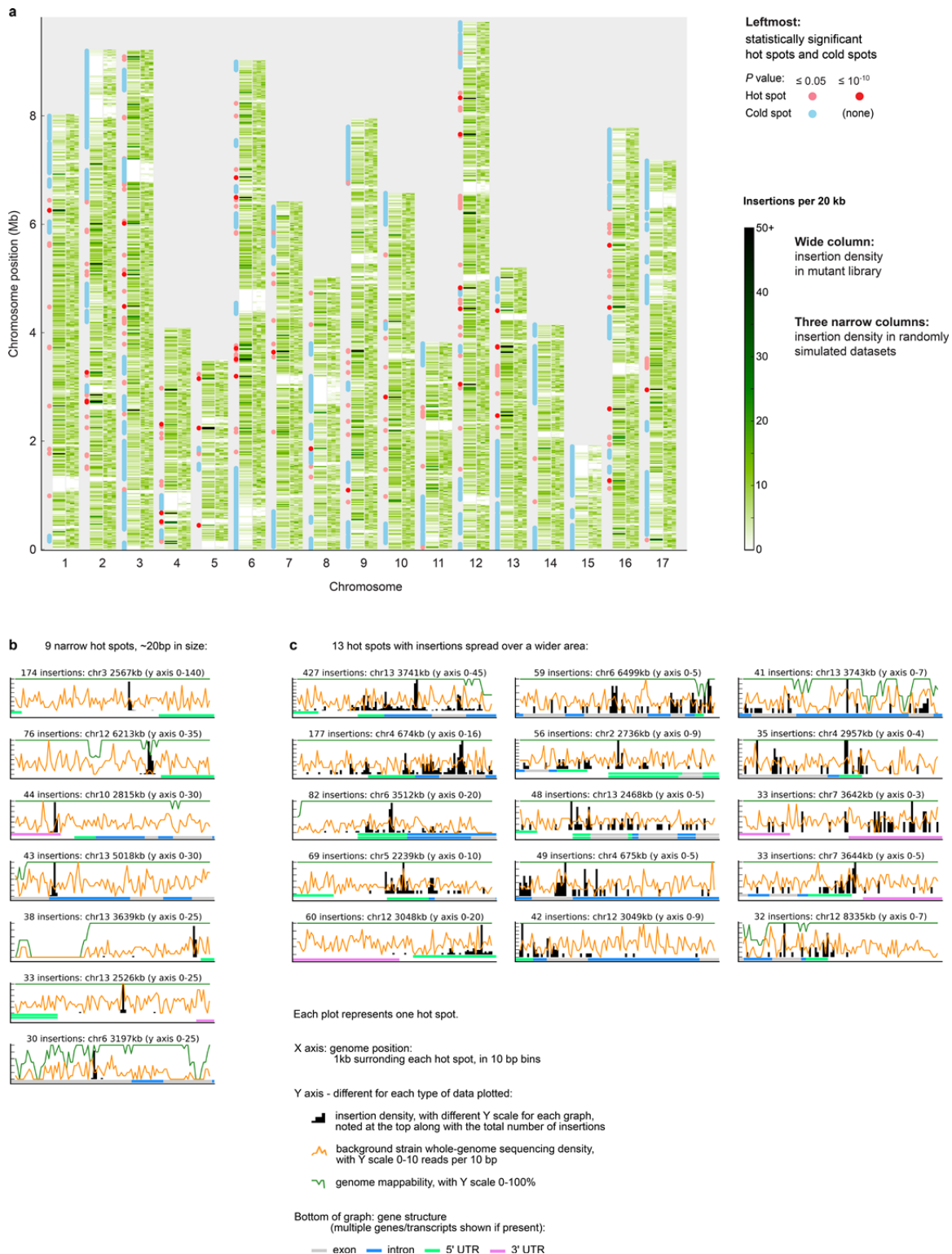
979 containing flanking genomic sequence a variable distance away from the insertion site.
980 During transformation, short fragments of genomic DNA, likely originating from lysed
981 cells, are often inserted between the cassette and the true flanking genomic DNA. We
982 refer to these short DNA fragments as “junk fragments”^{1,2}. Such junk fragments can lead
983 to incorrect insertion mapping if only the immediate flanking genomic sequence is
984 obtained. LEAP-Seq data can be used to detect presence of junk fragments at an insertion
985 junction based on two key characteristics: 1) the number of read pairs where both sides
986 aligned to the same locus and 2) the longest distance spanned by such read pairs. **g**, The
987 two key characteristics are plotted for the original full library, before any mapping
988 corrections were applied. **h**, The same two characteristics are plotted for confidence level
989 1 and 2 insertions. For confidence level 2 insertions, only the side with no junk fragment
990 is shown; for confidence level 1 insertions, one randomly chosen side is shown. **i**, LEAP-
991 Seq data can be used to correct cases of probable junk fragment insertions and determine
992 the most likely correct insertion position. The corrected data can be visualized using two
993 modified key characteristics: the number of distal reads aligned to the corrected location,
994 and the distance spanned by such reads. **j**, The modified characteristics are plotted for
995 confidence level 4 insertions.



996

997 **Supplementary Fig. 3 | Characterization of genomic disruptions in mutants in the**
 998 **library. a**, Mutants in the library were divided into four confidence levels, corresponding

999 to different mapping scenarios. The insertion sites of a number of randomly chosen
1000 mutants in each category were verified by PCR (mutants from confidence levels 1 and 2
1001 were assayed as one group; Supplementary Table 6). The numbers and percentages of
1002 confirmed insertions are shown in the last column. **b**, Most mutants have a single mapped
1003 insertion, and < 20% contain two or more mapped insertions. **c**, Eighteen randomly
1004 selected mutants from the four confidence levels were analyzed by Southern blotting
1005 using the coding sequence of *AphVIII* as the probe. Mutants are numbered and the details
1006 of their insertion sites are presented in Supplementary Table 6. The mutant number is
1007 highlighted in red when the Southern blot was interpreted to indicate at least two
1008 insertions in that mutant. The wild-type strain CC-4533 (WT) was included as a negative
1009 control. **d**, Most genomic deletions accompanying cassette insertions are smaller than 100
1010 bp, but deletions up to 10 kb are present in some mutants. Deletions larger than 10 kb
1011 may also be present, but there were not enough of them to be clearly detected based on
1012 the aggregate numbers. **e**, Most genomic duplications accompanying cassette insertion
1013 are smaller than 10 bp, but they can be up to 30 bp. Larger duplications may be present,
1014 but these are not common enough to be detected based on the aggregate numbers. **f**, The
1015 distribution of junk fragment lengths was determined using a dataset of 651 insertions of
1016 two cassettes surrounding a junk fragment, allowing us to precisely map both ends of the
1017 junk fragment using LEAP-Seq. Most junk DNA fragments are smaller than 320 bp, but
1018 we have detected some up to 1 kb in size. Larger junk fragments may be present, but are
1019 not common enough to be detected based on the aggregate numbers. Note that the x-axes
1020 for **d-f** are set to the logarithmic scale. Data presented in this figure are described in the
1021 Supplementary Note.

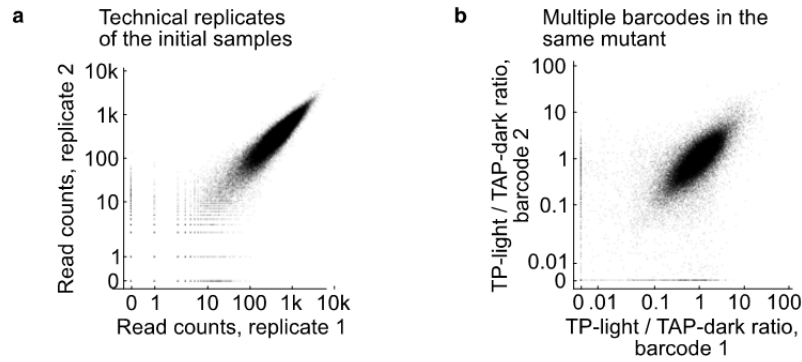


1022

1023 **Supplementary Fig. 4 | The distribution of insertions in the genome is largely**

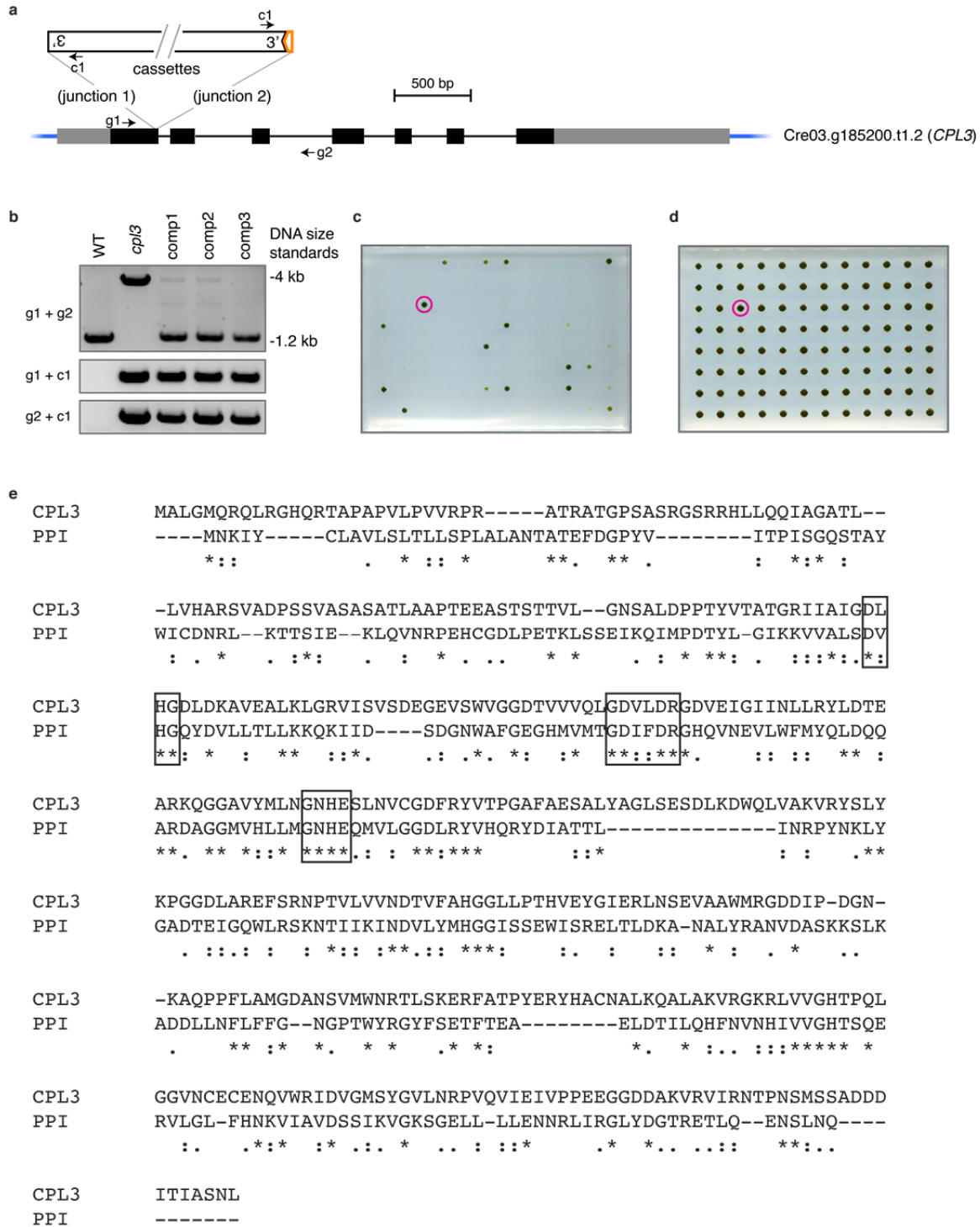
1024 **random, and the hot spots fall into two classes. a, For each chromosome, the observed**

1025 insertion density is shown as a heatmap in a wide column, followed by three narrow
1026 columns depicting three simulated datasets in which insertions were placed in randomly
1027 chosen mappable genomic locations. The simulated data provide a visual guide to the
1028 amount of variation expected from a random distribution. The large white areas present in
1029 both the observed and simulated data correspond to repetitive genomic regions in which
1030 insertions cannot be mapped uniquely. The red and blue circles/lines to the left of each
1031 chromosome show statistically significant insertion hot spots and cold spots, respectively.
1032 To ensure that we are showing true insertion density rather than artifacts caused by junk
1033 fragments or other mapping inaccuracies, the plot of insertion site distribution and
1034 identification of hot/cold spots are based on confidence level 1 insertions only. In
1035 contrast, Fig. 1c shows the distribution of insertions of all confidence levels over the
1036 genome. **b** and **c**, Each plot represents a 1-kb genomic region surrounding one hot spot,
1037 showing multiple features of that region, as listed in the legend. The plots shown are the
1038 22 1-kb regions with the highest total insertion number. The total number of insertions
1039 for each region is listed above each plot, along with the genomic position and the y-axis
1040 range. **b**, 7 of the top 22 hot spots are narrow, with 20 or more insertions in a 10-bp area,
1041 and a total width of 20-30 bp with few or no additional insertions in the surrounding 1 kb.
1042 **c**, 15 of the top 22 hot spots are wider, with multiple peaks of high insertion density
1043 spanning at least hundreds of base pairs. In either class, the insertion density peaks do not
1044 appear to reliably correlate with any of the other genomic features shown. Data presented
1045 in this figure are described in the Supplementary Note.



1046

1047 **Supplementary Fig. 5 | The barcode sequencing method is robust. a**, The barcode
1048 sequencing read counts (normalized to 100 million total reads) for each insertion were
1049 highly reproducible between technical replicates, with a Spearman's correlation of 0.978.
1050 94% of barcodes showed a normalized read count of no more than a 2-fold difference
1051 between the two replicates. **b**, The TP-light/TAP-dark ratios of multiple barcodes in the
1052 same mutant are consistent, with a Spearman's correlation of 0.744. Only 4% of insertion
1053 pairs had a greater than 5x difference between ratios. See also Fig. 3, b and c.



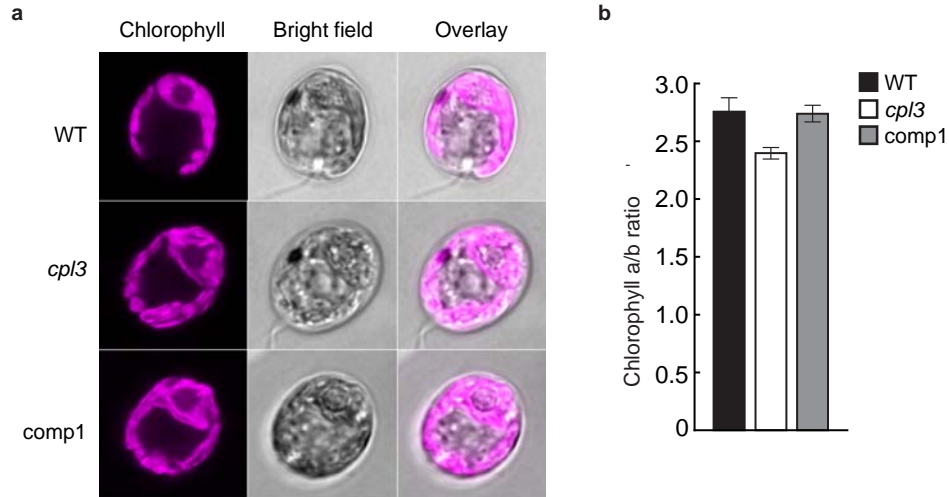
1054

1055 **Supplementary Fig. 6 | Molecular characterization of the *cpl3* mutant.** **a**, The cassette

1056 insertion site is indicated on a model of the *CPL3* gene from the *Chlamydomonas* v5.5

1057 genome. In the gene model, black boxes, gray boxes and thin lines indicate exons, UTRs,

1058 and introns respectively. Two cassettes are inserted in opposite orientations, with one of
1059 them truncated on the 3' side (indicated by a notch); the 5' ends may be intact or
1060 truncated. The orange box arrow indicates insertion of a small fragment of unknown
1061 origin. Binding sites for primers g1, g2, and c1 are indicated. **b**, PCR genotyping results
1062 of *cpl3* and complemented lines. PCR with the primer pair “g1 + g2” indicated presence
1063 of an insertion within the *CPL3* gene in the *cpl3* mutant and presence of wild-type *CPL3*
1064 sequence in the complemented lines. PCR with primer pairs “g1 + c1” and “g2 + c1”
1065 showed the presence of a cassette inserted into the *CPL3* gene in *cpl3* as well as the
1066 complemented lines. **c**, *cpl3* mutants transformed with the *CPL3* gene were arrayed and
1067 grown photosynthetically in the absence of acetate for one day under 100 $\mu\text{mol photons}$
1068 $\text{m}^{-2} \text{s}^{-1}$ light and four additional days under 500 $\mu\text{mol photons m}^{-2} \text{s}^{-1}$ light before
1069 imaging. The colony circled was a positive control strain that grows photosynthetically.
1070 Approximately 14% of transformants showed rescued photosynthetic growth, a frequency
1071 consistent with other genetic studies in *Chlamydomonas*³⁸. **d**, The same transformants
1072 were grown for five days in the presence of acetate in the medium under 50 μmol
1073 $\text{photons m}^{-2} \text{s}^{-1}$ light. All colonies grew similarly. **e**, *CPL3* contains conserved tyrosine
1074 phosphatase motifs. Sequences of *CPL3* in *Chlamydomonas* and its homolog
1075 psychrophilic phosphatase I (PPI) in *Shewanella sp.* were aligned using Clustal Omega³⁹.
1076 Asterisks (*), colon (:), and period (.) indicate conserved, strongly similar, and weakly
1077 similar amino acid residues, respectively. The motifs that are conserved among multiple
1078 protein phosphatases⁴⁰ are boxed. Data in panels **a-d** are described in the Supplementary
1079 Note. See also Fig. 4.
1080



1081

1082 **Supplementary Fig. 7 | Phenotypic characterization of the *cpl3* mutant.** **a**, *cpl3*, the
1083 wild-type strain (WT), as well as the complemented line (comp1), contain a normal cup-
1084 shaped chloroplast. Representative images of confocal chlorophyll fluorescence, bright
1085 field, and an overlay are shown for each strain. **b**, *cpl3* has a lower chlorophyll *a/b* ratio
1086 than WT and comp1 ($P < 0.03$, Student's *t*-test).

1087

1088

1089 **Captions for Supplementary Table 1 to S14 (each provided as a separate file)**

1090

1091 **Supplementary Table 1 | Primers and experimental design for all PCRs related to**

1092 **library generation and mapping**

1093 This database includes primers used for generation of the insertion cassette

1094 (Supplementary Fig. 1a), primers used for barcode amplification and sequencing (Fig. 3a,

1095 Supplementary Fig. 1f, and Supplementary Fig. 2a, b), primers and experimental design

1096 for LEAP-Seq (Supplementary Fig. 1g and Supplementary Fig. 2f-j). See also Methods.

1097

1098 **Supplementary Table 2 | Binary codes for plate super-pooling**

1099 In this database, each of the 570 rows corresponds to a plate-pool and each of the 21

1100 columns corresponds to a plate-super-pool. A 1 in row X and column Y indicates that

1101 plate-pool X was included in plate-super-pool Y, and a 0 indicates that it was not. See

1102 also Supplementary Fig. 1f, Supplementary Fig. 2a, and Methods.

1103

1104 **Supplementary Table 3 | Binary codes for colony super-pooling**

1105 Binary codes for the generation of 20 colony-super-pools from 384 colony-pools are

1106 shown in the same format as in Supplementary Table 2. See also Supplementary Fig.1f,

1107 Supplementary Fig. 2a, and Methods.

1108

1109 **Supplementary Table 4 | Read counts for each barcode in each combinatorial super-**
1110 **pool**

1111 The columns in this database are barcode, side (5' or 3'), and read counts for that barcode

1112 in each super-pool, normalized to 1 million total reads for each side and each super-pool.

1113 Only barcodes deconvolved to both a plate and a colony are included. See also
1114 Supplementary Fig. 1f, Supplementary Fig. 2a and Methods.

1115

1116 **Supplementary Table 5 | List of all mapped mutants in the library**

1117 Each line is an insertion junction (i.e. one mapped side of an insertion). Some insertions
1118 have one mapped side, some have two. Some mutants have multiple insertions. The
1119 columns are as follows (explained in detail in Methods):

1120 • mutant_ID - the ID of the mutant if it was included as part of the “consolidated”
1121 set for long-term maintenance, ‘-’ otherwise. A mutant can have multiple insertions.

1122 • side - which side of the cassette the data is derived from, 5’ or 3’.

1123 • chromosome, strand, min_position - the mapped position of the insertion,
1124 potentially corrected for a junk fragment, depending on the value of the if_fixed_position
1125 column.

1126 • gene, orientation, feature, gene_end_distances - the gene containing the
1127 insertion, the orientation with respect to the gene, the feature of the gene, and the distance
1128 from the 5’ and 3’ end of the gene. If the position is inside two overlapping genes, all
1129 fields will have two entries separated by ‘ & ’. For intergenic positions, all values are ‘-’.

1130 • intergenic_adjacent_genes, intergenic_orientations, intergenic_gene_distances –
1131 for intergenic positions, these fields note the two adjacent genes, the position of the
1132 insertion with respect to those two genes, and the distance from them. Each field will
1133 have two entries separated by ‘ & ’, unless the insertion position is on the edge of a
1134 chromosome and has no gene on one side. For insertions in genes, all values are ‘-’.

1135 • `if_both_sides` – ‘-’ if the insertion only has one mapped side, otherwise ‘perfect’,
1136 ‘deletion’ or ‘duplication’ for insertion junctions that are two sides of a confidence level
1137 1 insertion depending on whether there was a deletion/duplication in the genomic DNA,
1138 or ‘with-junk’ for insertion junctions that are two sides of a confidence level 2 insertion
1139 with a junk fragment on one side.

1140 • `confidence_level` - the confidence level for the insertion mapping, as described
1141 in Supplementary Fig. 3a.

1142 • `if_fixed_position` – ‘no’ if no junk fragment was detected in the insertion
1143 (confidence level 1, 3, the side of confidence level 2 with no junk fragment, and a small
1144 fraction of confidence level 4); ‘`yes_nearest_distal`’ if a junk fragment was detected and
1145 corrected. For the junk fragment sides of confidence level 2 insertions, the value is just
1146 ‘`yes_nearest_distal`’, indicating that the corrected insertion position for this line is the
1147 position of the nearest distal LEAP-Seq read; for most confidence level 4 insertions, the
1148 value is ‘`yes_nearest_distal_+/-X`’, indicating a further correction of X bp that was
1149 applied to the position to compensate for the average distance between the nearest distal
1150 LEAP-Seq read and the true insertion position (see Methods).

1151 • `LEAPseq_distance`, `LEAPseq_percent_confirming` - the highest distance
1152 spanned by a proximal and distal LEAP-Seq read pair mapping to the same region, and
1153 the fraction of pairs that map to the same region (see Supplementary Fig. 1g and
1154 Supplementary Fig. 2f-j).

1155 • `flanking_seq` - the flanking sequence immediately adjacent to the cassette, or, if
1156 the `if_fixed_position` column value is not ‘no’, the sequence of the distal LEAP-Seq read
1157 closest to the corrected mapping position.

1158 • barcode - the barcode sequence of the insertion.

1159 • gene_name, defline, description, etc - gene annotation from Phytozome¹⁰.

1160

1161 **Supplementary Table 6 | Primers and results of PCRs used to verify the insertion**

1162 **sites of randomly-picked mutants from the mutant library**

1163 For column definitions, see the legend of Supplementary Table 5.

1164

1165 **Supplementary Table 7 | Statistically significant insertion hot spots and cold spots**

1166 The columns give the hot spot position (chromosome, start and end base number), type
1167 (hot spot or cold spot, i.e. enriched or depleted in insertions), false discovery rate (FDR),
1168 number of insertions in the spot, and the expected number of insertions based on the
1169 length and mappability of the spot (see Methods). Only hot spots that passed the filtering
1170 are listed.

1171

1172 **Supplementary Table 8 | Statistically significant depleted functional terms**

1173 The columns give the gene ontology (GO) term, FDR, the ratio of observed vs expected
1174 insertions in all the genes annotated with the term (i.e. the effect size), the number of
1175 genes annotated with the term, the total number of insertions in those genes, the total
1176 mappable length of those genes, and the GO term definition (see Methods). Only
1177 depleted GO terms are listed - most of the enriched GO terms were due to hot spots in a
1178 single gene.

1179

1180 **Supplementary Table 9 | Candidate essential genes**

1181 The columns give the gene ID, its mappable length (excluding UTRs), the number of
1182 insertions in the gene (again excluding UTR insertions), the expected number of
1183 insertions given the mappable length if the insertion distribution was random, the FDR,
1184 and gene annotation from Phytozome¹⁰ (data described in the Supplementary Note).

1185

1186 **Supplementary Table 10 | Read counts of barcodes before and after pooled growth**
1187 **in the photosynthesis screen**

1188 The columns give the barcode, the gene in which the insertion is located (or “-” if
1189 intergenic), the gene feature, the side of cassette the data is derived from, and deep-
1190 sequencing read numbers (raw and normalized to 100 million) in the two technical
1191 replicates of the initial pool, the pool after growth in TAP-dark, and two biological
1192 replicates after growth in TP-light.

1193

1194 **Supplementary Table 11 | Statistics of the pooled growth data for all genes**

1195 The columns give:

- 1196 • the gene ID and name, the hit_category (higher-confidence candidate or lower-
1197 confidence candidate, otherwise “-”).
- 1198 • the number of alleles in the gene with and without a phenotype in replicate 1, the
1199 resulting *P* value and FDR in replicate 1 (see Methods; genes with only 1 allele have no
1200 FDR).
- 1201 • the same four numbers for replicate 2.

1202 • PredAlgo-predicted localization for the gene: C = chloroplast, M
1203 =mitochondrion, SP = secretory pathway, O = other, and “-” if no prediction could be
1204 made.

1205 • gene annotation data from Phytozome¹⁰.

1206

1207 **Supplementary Table 12 | Summary of previous characterizations of higher- and**
1208 **lower-confidence genes’ roles in photosynthesis**

1209 The columns are similar to those in Supplementary Table 11, but additionally include:

1210 • Previously reported function in photosynthesis for each gene.

1211 • The corresponding references.

1212

1213 **Supplementary Table 13 | Read counts of *cpl3* exon and intron alleles in the pooled**
1214 **screens**

1215 The numbers of read for each barcode per 100 million total reads are presented. Some of
1216 the mutants contain two barcodes and are labeled with an asterisk (*). In such cases, the
1217 barcode with a greater TAP-dark read count is used to determine whether there is a
1218 phenotype. The mutants that fall below the phenotype cutoff are labeled with an obelisk
1219 (†). Two out of the seven mutants were included in the screen but not included in the
1220 consolidated set.

1221

1222 **Supplementary Table 14 | Proteomic characterization of the *cpl3* mutant**

1223 The columns give:

1224 • ID for nuclear/chloroplast-encoded genes.

1225 • name - gene name from Phytozome bulk annotation.

1226 • annotations for proteins plotted in Fig. 4 - names of the complexes and protein

1227 subunits are shown.

1228 • spectral counts – number of spectra detected for peptides derived from each

1229 protein, related to protein abundance.

1230 • WT_repl1/2/3, *cpl3*_repl1/2/3, comp1_repl1/2/3 - log₂ intensity values for three

1231 replicates of each sample, normalized to a median of 0. WT, the wild-type parental strain

1232 CC-4533. Comp1, the complemented line.

1233 • WT_mean, *cpl3*_mean, comp1_mean - the average of the three replicate

1234 intensity values for each sample.

1235 • WT-*cpl3*_diff - *cpl3*_mean subtracted by WT_mean. The lower this value, the

1236 less abundant the protein is in *cpl3* relative to WT.

1237 • WT-*cpl3*_pval - the raw *P* value comparing the normalized replicate intensities

1238 for WT and *cpl3*, using an unpaired t-test.

1239 • WT-*cpl3*_FDR - the false discovery rate (i.e. the *P* value adjusted for multiple

1240 testing, using the Benjamini-Hochberg method¹²).

1241 • WT-comp1_diff, WT-comp1_pval, WT-comp1_FDR, *cpl3*-comp1_diff, *cpl3*-

1242 comp1_pval, *cpl3*-comp1_FDR - the same three values for the comparison between WT

1243 and comp1 samples and between *cpl3* and comp1 samples.

1244 • define, description, and all further columns - gene annotation from bulk

1245 Phytozome data.

1246

References

- 1247
1248
1249 1. Zhang, R. *et al.* High-Throughput Genotyping of Green Algal Mutants Reveals Random
1250 Distribution of Mutagenic Insertion Sites and Endonucleolytic Cleavage of Transforming DNA.
1251 *Plant Cell* **26**, 1398-1409 (2014).
1252 2. Li, X. *et al.* An Indexed, Mapped Mutant Library Enables Reverse Genetics Studies of Biological
1253 Processes in *Chlamydomonas reinhardtii*. *Plant Cell* **28**, 367-87 (2016).
1254 3. Rubin, B.E. *et al.* The essential gene set of a photosynthetic organism. *Proc Natl Acad Sci U S A*
1255 **112**, E6634-43 (2015).
1256 4. Wetmore, K.M. *et al.* Rapid quantification of mutant fitness in diverse bacteria by sequencing
1257 randomly bar-coded transposons. *MBio* **6**, e00306-15 (2015).
1258 5. Kropat, J. *et al.* A revised mineral nutrient supplement increases biomass and growth rate in
1259 *Chlamydomonas reinhardtii*. *Plant J* **66**, 770-80 (2011).
1260 6. Grassl, M. Bounds on the minimum distance of linear codes and quantum codes. Vol. 2017
1261 (<http://www.codetables.de/>, 2015).
1262 7. Simonis, J. The [23; 14; 5] Wagner code is unique. *Discrete Mathematics* **213**, 269-282 (2000).
1263 8. Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet*
1264 *J.* **17**, 10-12 (2011).
1265 9. Langmead, B. & Salzberg, S.L. Fast gapped-read alignment with Bowtie 2. *Nat Methods* **9**, 357-9
1266 (2012).
1267 10. Merchant, S.S. *et al.* The *Chlamydomonas* genome reveals the evolution of key animal and plant
1268 functions. *Science* **318**, 245-251 (2007).
1269 11. Pazour, G.J., Agrin, N., Leszyk, J. & Witman, G.B. Proteomic analysis of a eukaryotic cilium. *J*
1270 *Cell Biol* **170**, 103-13 (2005).
1271 12. Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: a practical and powerful
1272 approach to multiple testing. *J. R. Statist. Soc. B* **57**, 289-300 (1995).
1273 13. Dent, R.M. *et al.* Large-scale insertional mutagenesis of *Chlamydomonas* supports phylogenomic
1274 functional prediction of photosynthetic genes and analysis of classical acetate-requiring mutants.
1275 *Plant J* **82**, 337-351 (2015).
1276 14. Karpowicz, S.J., Prochnik, S.E., Grossman, A.R. & Merchant, S.S. The GreenCut2 resource, a
1277 phylogenomically derived inventory of proteins specific to the plant lineage. *J Biol Chem* **286**,
1278 21427-39 (2011).
1279 15. Zones, J.M., Blaby, I.K., Merchant, S.S. & Umen, J.G. High-Resolution Profiling of a
1280 Synchronized Diurnal Transcriptome from *Chlamydomonas reinhardtii* Reveals Continuous Cell
1281 and Metabolic Differentiation. *Plant Cell* (2015).
1282 16. Duanmu, D. *et al.* Retrograde bilin signaling enables *Chlamydomonas* greening and phototrophic
1283 survival. *Proc Natl Acad Sci U S A* **110**, 3621-6 (2013).
1284 17. Berthold, P., Schmitt, R. & Mages, W. An engineered *Streptomyces hygrosopicus* aph 7" gene
1285 mediates dominant resistance against hygromycin B in *Chlamydomonas reinhardtii*. *Protist* **153**,
1286 401-12 (2002).
1287 18. Porra, R.J., Thompson, W.A. & Kriedemann, P.E. Determination of accurate extinction
1288 coefficients and simultaneous equations for assaying chlorophylls a and b extracted with four
1289 different solvents: verification of the concentration of chlorophyll standards by atomic absorption
1290 spectroscopy. *Biochimica et Biophysica Acta (BBA) - Bioenergetics* **975**, 384-94 (1989).
1291 19. Saroussi, S.I., Wittkopp, T.M. & Grossman, A.R. The Type II NADPH Dehydrogenase Facilitates
1292 Cyclic Electron Flow, Energy-Dependent Quenching, and Chlororespiratory Metabolism during
1293 Acclimation of *Chlamydomonas reinhardtii* to Nitrogen Deprivation. *Plant Physiol* **170**, 1975-88
1294 (2016).
1295 20. Schindelin, J. *et al.* Fiji: an open-source platform for biological-image analysis. *Nat Methods* **9**,
1296 676-82 (2012).
1297 21. Rappsilber, J., Ishihama, Y. & Mann, M. Stop and go extraction tips for matrix-assisted laser
1298 desorption/ionization, nanoelectrospray, and LC/MS sample pretreatment in proteomics. *Anal*
1299 *Chem* **75**, 663-70 (2003).

- 1300 22. Kulak, N.A., Pichler, G., Paron, I., Nagaraj, N. & Mann, M. Minimal, encapsulated proteomic-
1301 sample processing applied to copy-number estimation in eukaryotic cells. *Nat Methods* **11**, 319-24
1302 (2014).
- 1303 23. Cox, J. & Mann, M. MaxQuant enables high peptide identification rates, individualized p.p.b.-
1304 range mass accuracies and proteome-wide protein quantification. *Nat Biotechnol* **26**, 1367-72
1305 (2008).
- 1306 24. Maul, J.E. *et al.* The *Chlamydomonas reinhardtii* plastid chromosome: islands of genes in a sea of
1307 repeats. *Plant Cell* **14**, 2659-79 (2002).
- 1308 25. Michaelis, G., Vahrenholz, C. & Pratje, E. Mitochondrial DNA of *Chlamydomonas reinhardtii*:
1309 the gene for apocytochrome b and the complete functional map of the 15.8 kb DNA. *Mol Gen
1310 Genet* **223**, 211-6 (1990).
- 1311 26. Bern, M., Kil, Y.J. & Becker, C. Byonic: advanced peptide and protein identification software.
1312 *Curr Protoc Bioinformatics* **Chapter 13**, Unit13 20 (2012).
- 1313 27. Vu, G.T. *et al.* Repair of Site-Specific DNA Double-Strand Breaks in Barley Occurs via Diverse
1314 Pathways Primarily Involving the Sister Chromatid. *Plant Cell* **26**, 2156-2167 (2014).
- 1315 28. Giaever, G. *et al.* Functional profiling of the *Saccharomyces cerevisiae* genome. *Nature* **418**, 387-
1316 91 (2002).
- 1317 29. Peters, J.M. *et al.* A Comprehensive, CRISPR-based Functional Analysis of Essential Genes in
1318 Bacteria. *Cell* **165**, 1493-506 (2016).
- 1319 30. Wang, T. *et al.* Identification and characterization of essential genes in the human genome.
1320 *Science* **350**, 1096-101 (2015).
- 1321 31. Baum, P., Yip, C., Goetsch, L. & Byers, B. A yeast gene essential for regulation of spindle pole
1322 duplication. *Mol Cell Biol* **8**, 5386-97 (1988).
- 1323 32. Cullmann, G., Fien, K., Kobayashi, R. & Stillman, B. Characterization of the five replication
1324 factor C genes of *Saccharomyces cerevisiae*. *Mol Cell Biol* **15**, 4661-71 (1995).
- 1325 33. Kunz, J. *et al.* Target of rapamycin in yeast, TOR2, is an essential phosphatidylinositol kinase
1326 homolog required for G1 progression. *Cell* **73**, 585-96 (1993).
- 1327 34. Spalding, M.H. The CO₂-Concentrating Mechanism and Carbon Assimilation. in *The
1328 Chlamydomonas Sourcebook*, Vol. 2 (eds. E.H., H., E.B., S. & G.B., W.) 257-301 (Academic
1329 Press, 2009).
- 1330 35. Devenish, R.J., Prescott, M. & Rodgers, A.J. The structure and function of mitochondrial F1F0-
1331 ATP synthases. *Int Rev Cell Mol Biol* **267**, 1-58 (2008).
- 1332 36. Jarvis, P. *et al.* Galactolipid deficiency and abnormal chloroplast development in the Arabidopsis
1333 MGD synthase 1 mutant. *Proc Natl Acad Sci U S A* **97**, 8175-9 (2000).
- 1334 37. Riekhof, W.R., Sears, B.B. & Benning, C. Annotation of genes involved in glycerolipid
1335 biosynthesis in *Chlamydomonas reinhardtii*: discovery of the betaine lipid synthase BTA1Cr.
1336 *Eukaryot Cell* **4**, 242-52 (2005).
- 1337 38. Wang, L. *et al.* Chloroplast-mediated regulation of CO₂-concentrating mechanism by Ca²⁺-
1338 binding protein CAS in the green alga *Chlamydomonas reinhardtii*. *Proc Natl Acad Sci U S A* **113**,
1339 12586-12591 (2016).
- 1340 39. Sievers, F. *et al.* Fast, scalable generation of high-quality protein multiple sequence alignments
1341 using Clustal Omega. *Mol Syst Biol* **7**, 539 (2011).
- 1342 40. Tsuruta, H., Mikami, B. & Aizono, Y. Crystal structure of cold-active protein-tyrosine
1343 phosphatase from a psychrophile, *Shewanella* sp. *J Biochem* **137**, 69-77 (2005).
- 1344