

1 **Green genes: bioinformatics and systems biology innovations drive algal**  
2 **biotechnology**

3

4 **Authors' names:**

5 Maarten J.M.F. Reijnders<sup>1</sup>, Ruben van Heck<sup>1</sup>, Carolyn M.C. Lam<sup>1,2</sup>, Mark A. Scaife<sup>3</sup>, Vitor A.P.  
6 Martins dos Santos<sup>1,2</sup>, Alison G. Smith<sup>3</sup>, Peter J. Schaap<sup>1</sup>

7 **Authors' addresses:**

8 <sup>1</sup>Laboratory of Systems and Synthetic Biology, Wageningen University, Dreijenplein 10, Building  
9 number 316, 6703 HB Wageningen, The Netherlands.

10 <sup>2</sup>LifeGlimmer GmbH, Markelstr. 38, D-12163 Berlin, Germany.

11 <sup>3</sup>Department of Plant Sciences, University of Cambridge, Downing Street, Cambridge CB2 3EA,  
12 UK.

13 **Corresponding author:** Schaap, Peter J. ([peter.schaap@wur.nl](mailto:peter.schaap@wur.nl))

14

15

## Abstract (120 words)

Many species of microalgae possess the capacity to produce hydrocarbons, polysaccharides, and other valuable products in significant amounts. However, large-scale production of algal products is not yet competitive against non-renewable alternatives from fossil fuel. Metabolic engineering approaches will help improve productivity, but the exact metabolic pathways and identity of the majority of the genes involved remain unknown. Recent advances in bioinformatics and systems biology modeling coupled with increasing numbers of algal genome sequencing projects are providing the means to address this. A multi-disciplinary integration of methods will provide synergy for a systems-level understanding of microalgae, and thereby speed up the improvement of industrially valuable strains. In this review we highlight recent advances, challenges, their application in microalgae research, and potentials.

## Introduction

Microalgae are simple photosynthetic eukaryotes that are among the most diverse of all organisms. Microalgae inhabit all aquatic ecosystems, from oceans, lakes and rivers to even snow and glaciers, as well as terrestrial systems including rocks and other hard surfaces. Microalgae exhibit significant variation in physiology and metabolism, which is reflected in the genetic diversity that exists both between species, and within a single genome, where a seeming mosaic of genetic origins can be observed [1]. Mining the genomes of these organisms is a great opportunity to identify novel pathways of biotechnological importance. In particular, microalgae are of considerable interest for the synthesis of a range of industrially useful products, such as hydrocarbons and polysaccharides [2, 3], due to rapid growth rates, amenability for large scale fermentation, and the potential for sustainable process development [4].

1 Algae as a source of biofuel molecules, such as triacylglycerides (TAGs), the precursor for  
2 biodiesel [5], have been a focus in recent years, with potential yields an order of magnitude  
3 greater than competing agricultural processes [6]. Currently, in order for microalgae to  
4 synthesize TAG it is necessary to expose them to stress conditions such as nutrient limitation,  
5 which reduces growth and increases energy dissipation. This trade-off between biosynthesis of  
6 TAG and cell growth is therefore a severely limiting factor [7]. If a better understanding of the  
7 metabolic and regulatory networks were available, they could be rewired for increased TAG  
8 synthesis, with fewer drawbacks than for existing algal cells.

9 The production of other interesting algal products will also benefit from a better understanding  
10 of microalgae on a systems level. For example, polysaccharides such as starch and cell wall  
11 materials can be used for biotechnological applications [8]. These carbohydrates can be  
12 degraded to fermentable sugars for bio-ethanol production [9], or serve as chemical building  
13 blocks for renewable materials, but the composition and proportions of the different sugar  
14 components require optimization. Similarly, various valuable secondary metabolites produced  
15 by microalgae are of interest in the food, nutrition, and cosmetics industries [2], but often they  
16 are produced in trace amounts, or only under conditions that are not amenable to industrial  
17 cultivation.

18 Over 30 microalgal genomes have been sequenced and numerous transcriptomics, proteomics  
19 and other systems biology studies performed. Yet our understanding of metabolic pathways  
20 within these microalgae remains limited [10]. A significant knowledge gap needs to be filled  
21 between omics data, the annotation thereof and our systems level understanding. This will  
22 allow the conversion of these resources into useable genome scale models, and provide the  
23 basis for effective metabolic engineering, synthetic biology and biotechnology. Here we  
24 consider novel approaches to improve annotation of algal omics data, the establishment of

genome-scale metabolic models, and ways to integrate these for effective exploitation of microalgae for biotechnology.

#### Annotation challenges for microalgae

The nuclear genome of the model green alga *Chlamydomonas reinhardtii*, sequenced in 2007 [1], is approximately 120 Mb and encodes ~15,000 genes. Although *C. reinhardtii* is commonly used as a reference for the annotation of other microalgae, only a subset of ~50 proteins have experimentally validated functions according to the UniProt database (<http://www.uniprot.org>), compared to ~6800 proteins for the model plant *Arabidopsis thaliana*. Consequently, most *C. reinhardtii* genes have been computationally annotated by inferred homology with other organisms, including *A. thaliana*, other plants and microbes [1], using BLAST or family-wise alignment methods such as HMMER and InterProScan (Table 2). Blast-based methods often use the principle of one-to-one recognition meaning that annotation of a query gene is based on the annotation of a single known gene. This limits the success rate for recognition and correct annotation of distantly related *C. reinhardtii* genes, but becomes even more problematic when an *in silico* derived annotation of *C. reinhardtii* is subsequently used for the annotation of other algal species. This is because two algal species can be more diverse than any two plant species, for example. Therefore, these methods, which are highly suitable for high-throughput analysis due to their simplicity, are not suitable for accurate in-depth annotation of algal genomes. In the Critical Assessment of protein Function Annotation (CAFA) experiment [11], the accuracy of more advanced function annotation algorithms was assessed. The CAFA concluded that 33 of 54 tested function annotation algorithms outperformed the standard BLAST-based method (Table 2). The substantial improvement can be explained by the fact that these 2<sup>nd</sup> generation methods do not apply the one-to-one

1 recognition principal, but rather use a one-to-many recognition strategy to increase their  
2 success rate and include context aware principles for annotation. An example is Argot2 (Box 1)  
3 [12], which applies the one-to-many recognition strategy by calculating the statistical  
4 significance of all candidate homologous genes found by BLAST [13] and HMMER [14],  
5 combined with an assessment of semantic similarities of associated GO terms. In a context  
6 aware multi-level approach, annotation is not merely based on sequence similarity, but also  
7 other factors such as protein-protein interactions [15], transcript expression patterns [15],  
8 phylogenetic trees [16], compartmentalization [17], and literature [18] are taken into account.  
9 FFPred2 from UCL-Jones [17] is the prime example of such a homology-independent function  
10 annotation algorithm.

11

12 Advanced multi-level annotation methods effectively increase the recall of function prediction  
13 while maintaining a reasonable precision. The challenge in genome annotation of microalgae  
14 lies in the small number of experimentally validated algal genes and the lack of algae-specific  
15 contextual data such as protein interaction data. This results in a relatively low number of  
16 genes predicted to have a specific biological function. To overcome this, multiple annotation  
17 methods and data sources should be combined. The combined result increases the number of  
18 annotated genes, whilst a consensus prediction among the different methods improves the  
19 accuracy of the annotation [19]. Due to their simplicity and speed, 1<sup>st</sup>-generation methods can  
20 be used for initial high-throughput analysis of a large set of genes. 2<sup>nd</sup>-generation methods can  
21 then be used for a refined analysis of these genes. However, in order to utilize these advanced  
22 methods fully, a significant amount of experimentally determined contextual data is required.  
23 Whilst increasing amounts of gene expression data are being generated, there is still little  
24 structural and protein interaction data available for algae. In the absence of such experimental  
25 facts, it is still possible to generate this contextual information by *in silico* prediction methods

[20, 21]. Although studies have shown that this is a feasible option [22], caution is necessary, since there is a high risk of error propagation.

Apart from functional annotation it is also important to establish the cellular location of a protein. For this there are several tools available, including Argot2 (Box 1) [12], TargetP [23], SignalP [24], PSORTb [25] and PredAlgo [26]. The latter of which is a tailor-made multi-subcellular localization prediction tool dedicated to three compartments of green algae: the mitochondrion, the chloroplast, and the secretory pathway. However, due to the limited number of algal proteins with a known cellular localization, the algorithm is trained with a relatively small *C. reinhardtii* data set. This raises questions regarding reliability for other algal species that are more distantly related. Therefore it is advisable to use PredAlgo in combination with non-algal specific tools in a similar way as for functional annotation.

To support large-scale annotation of algal sequence data, up to date databases and readily available supporting tools are required. Online databases give the means to share data easily so that the scientific community can profit as a whole. Supporting tools can assist in annotating genes, pathways, and performing statistical analysis. While genomic data for various algae are available in NCBI and UniProt, the amount of public data is lagging behind in comparison to plant and bacterial species. Additionally, tools and databases that do more than storing the available sequencing data are needed. There are a small number of tools available, although these are often limited to *C. reinhardtii*. One such tool is ChlamyCyc [27] a *C. reinhardtii* specific pathway/genome database of the MetaCyc [28] facility for metabolic pathway analysis. Additionally, the Augustus tool, which is commonly used for prediction of eukaryotic genes [29], has a tailor-made section for *C. reinhardtii*. Finally, the Algal Functional Annotation Tool [30] incorporates annotation data for a few microalgal species from several pathway databases,

1 ontologies, and protein families. Broadening the scope of these, annotation tools for a range of  
2 microalgae would allow comparative analysis, which is useful for easy mapping of various  
3 differences between microalgae. In this context, a useful tool which has been applied to plant  
4 research is Phytozome (<http://www.phytozome.net>) [31], a comparative hub for analysis of  
5 plant genomes and gene families. It acts as a reference for the key data of many plant species,  
6 and provides click-to-go features such as BLAST and summarizing key data. Phytozome has  
7 grown to be a major asset to the plant science community. Although it contains data from a few  
8 green algae, an expanded web-portal focused on algal systems-bioinformatics research could  
9 be of immense benefit to the field, particularly for those studying the more industrially-relevant  
10 diatoms and heterokont species (Table 1). Such a web-portal would provide the means for new  
11 and existing tools specifically useful for algal species to facilitate exposure to a broad audience.  
12 Additionally, it could act as a hosting platform for small but useful tools like a refined algal  
13 literature research algorithm, and tools that suggest genes to fill gaps in metabolic or  
14 regulatory pathways for microalgae. Adopting an algal web-portal would provide a good  
15 overview of all available data and tools, and help in reducing the redundancy that is often seen  
16 in biology and bioinformatics.

17

1 **Table 1: A list of selected industrially useful microalgae.**

Species	Genome size <sup>‡</sup> (Mb)	Proteins in UniProt <sup>‡</sup>	Characteristics <sup>‡</sup>	Ref <sup>‡</sup>
<i>Chlamydomonas reinhardtii</i>	120	15,144	<ul style="list-style-type: none"> <li>• Model system for unicellular green algae</li> <li>• Rapid growth</li> </ul>	
<i>Monoraphidium neglectum</i>	68	16,761	<ul style="list-style-type: none"> <li>• Biofuel production candidate</li> </ul>	[32]
<i>Nannochloropsis</i> sp.	44	16,226	<ul style="list-style-type: none"> <li>• Produces high amounts of omega-3-long-chain polyunsaturated fatty acids</li> </ul>	
<i>Phaeodactylum tricornutum</i>	27	10,673	<ul style="list-style-type: none"> <li>• Production of antibacterial fatty acids</li> </ul>	[33]
<i>Chlorella variabilis</i>	46	9,831	<ul style="list-style-type: none"> <li>• Contains several essential nutrients</li> <li>• Rich source of lutein</li> </ul>	
<i>Ostreococcus tauri</i>	12.6	9,050	<ul style="list-style-type: none"> <li>• Smallest microalgal genome</li> </ul>	
<i>Chlorella vulgaris</i>	n.a.	292	<ul style="list-style-type: none"> <li>• High lipid content under nitrogen limitation</li> </ul>	[34]
<i>Dunaliella salina</i>	n.a.	238	<ul style="list-style-type: none"> <li>• High concentration of beta-carotene</li> </ul>	[35]
<i>Chlorella protothecoides</i>	n.a.	96	<ul style="list-style-type: none"> <li>• High lipid content in heterotrophic growth [3]</li> <li>• Highest published biomass yield [36]</li> </ul>	[3], [36]
<i>Haematococcus pluvialis</i>	n.a.	60	<ul style="list-style-type: none"> <li>• Antioxidant astaxanthin production</li> </ul>	[37]
<i>Botryococcus braunii</i>	~166 – 211 [38]	30	<ul style="list-style-type: none"> <li>• High levels of liquid hydrocarbons and exopolysaccharides [39]</li> </ul>	[38], [39]
<i>Neochloris oleoabundans</i>	n.a.	0	<ul style="list-style-type: none"> <li>• High lipid content</li> </ul>	[7]

2 <sup>‡</sup>: Genome size, estimated protein numbers and characteristics are according to NCBI and  
3 UniProt, unless otherwise specified.  
4 n.a.: Not available.

5



1 **Table 2: Comparison of the features of commonly used functional annotation tools.**

Methods	Success rate*	Computational speed	Availability	Additional notes	Ref
Standard BLAST	Limited	Fast	Online/ offline	<ul style="list-style-type: none"> <li>• Dependent on global sequence similarity for success</li> <li>• Suitable for high throughput analysis</li> </ul>	[13]
HMMER	Moderate	Fast	Online/ offline	<ul style="list-style-type: none"> <li>• Family-wise alignment method</li> <li>• Suitable for high throughput analysis</li> </ul>	[14]
InterProScan	Moderate	Slow	Online/ offline	<ul style="list-style-type: none"> <li>• Family-wise alignment method</li> <li>• Uses pre-computed protein domains</li> </ul>	[40]
FFPred2	High	Slow	Limited online/ offline	<ul style="list-style-type: none"> <li>• Algorithms currently trained on non-algal datasets</li> <li>• Not suitable for high throughput analysis</li> </ul>	[17, 20]
Argot2	High	Moderate	Limited online	<ul style="list-style-type: none"> <li>• Initial selection is dependent on BLAST and HMMER output</li> <li>• Additionally predicts compartmentalization</li> <li>• User-friendly interface</li> </ul>	[12]

\* For distantly related sequences

## 1    **Box 1: Argot2**

2    One of the top performers in the CAFA experiment is Argot2 [12]. It stands out when it comes  
3    to simplicity, as well as its incorporation of BLAST and HMMER. This method combines an easy  
4    interface with multi-layer analysis, making it a perfect starting point for biologists who want to  
5    annotate their data.

6    Argot2 requires a nucleotide or protein sequence as input. It queries the UniProt and Pfam  
7    databases using BLAST and HMMER respectively, providing an initial high-throughput  
8    sequence analysis. A weighting scheme and clustering algorithm are then applied to the results  
9    to select the most accurate GO terms for each query sequence. The user can choose to perform  
10   this entire process online at the Argot2 webserver, limited to a hundred sequences per query.  
11   Alternatively, if the BLAST and HMMER steps are performed by the user locally and provided to  
12   the webserver, over a thousand sequences can be submitted per query. After the analysis is  
13   completed, which can take several hours depending on the amount of input data, the user is  
14   provided with the prediction results as well as the intermediate BLAST and HMMER files. These  
15   predictions include molecular function, biological processes and cellular component GO terms  
16   for each query. Predicted GO terms are ranked by a score based on statistical significance and  
17   specificity. Optionally, the user can choose to compute protein clusters based on functional  
18   similarity.

19

## Box 2: Flux analysis in microalgae

Flux balance analysis (FBA) [41] is the most commonly applied method to simulate metabolism in genome-scale metabolic models. It identifies a theoretically optimal use of metabolic capabilities for a selected metabolic objective, in a specific environment. As some microalgae can grow autotrophically in chemically defined medium, the boundary conditions for consumption of all medium components are well-specified in those cases. This is advantageous for *in silico* metabolic flux analysis using metabolic models *e.g.* how a microalga can achieve maximal growth under defined illumination. In addition, disabling metabolic capabilities associated with a gene allows simulation of mutant strains. FBA can thus assess the potential of different strains and different environmental conditions. In order to run FBA, all reactions are organized in a stoichiometric matrix  $S$ . Each column in  $S$  represents a different reaction, and each row a different metabolite. A nonzero value at position  $[i,j]$  thus indicates the stoichiometric coefficient of metabolite  $i$  in reaction  $j$ . FBA then employs two different constraints: (i) Metabolism is assumed to be in steady-state; production/degradation of intermediate compounds is not possible. (ii) Thermodynamics (reversibility) and substrate availability both dictate lower and upper flux bounds for individual reactions. Finally, one or more reactions are selected to represent the metabolic objective of, for example, algal biomass production. Together, the  $S$  matrix, the constraints, and the objective function form a linear programming problem:

$$\max(\mathbf{x} * \mathbf{c})$$

$$s.t. \quad S * \mathbf{x} = \mathbf{0}$$

$$\mathbf{x} \geq \mathbf{lb}$$

$$\mathbf{x} \leq \mathbf{ub}$$

1 With:  $\mathbf{x}$  — flux vector,  $\mathbf{c}$  — objective vector,  $\mathbf{0}$  — a null vector ensuring steady state, and  $\mathbf{lb}/\mathbf{ub}$  —  
2 lower/upper bounds for each reaction. The vector  $\mathbf{x}$  represents a flux distribution with the  
3 theoretically maximal value for the metabolic objective. However, due to the presence of  
4 alternative/cyclic pathways, there are often alternative flux distributions with equally high  
5 values for the objective function. Flux variability analysis [34] explores for each reaction to  
6 what extent the flux can vary while permitting only a small reduction in the obtained objective  
7 value. In addition, experimental data can be used to provide additional constraints. For  
8 example,  $^{13}\text{C}$ -labelling experiments provides experimentally measured fluxes as inputs for the  
9 model simulations [42, 43]. Several FBA-based methods also facilitate the integration of  
10 transcriptomic and proteomic data with metabolic models to constrain reactions based on the  
11 gene expression levels [44]. Thereby, flux distributions are identified which are most consistent  
12 with the expression data [45]. Metabolic models are thus integrated with, and their predictions  
13 compared to, experimental data yielding new insights in metabolic functioning.

14

15

## Understanding algal metabolism on a systems level

The sheer number of genes for metabolic enzymes, combined with the complexity of cellular metabolism, means that it is not straightforward to establish metabolic capability, even for well-annotated species. This limitation has led to the development of metabolic models, which represent a snapshot of metabolism of an organism in a network format. Once an annotated algal genome or transcriptome is available, a corresponding genome-scale metabolic model (GSMM) can be reconstructed and the topology of the metabolic network of the algal species can be analyzed. An initial draft model can be generated directly from the genome annotation and is then adjusted and expanded based on experimental data, literature and gap-filling procedures. The final model then includes all reactions the alga is known to perform and the associated genes and constraints, e.g. reaction directionalities and rate limits. Due to their comprehensive representation of metabolism, metabolic models form the basis for a large and diverse set of mathematical methods predicting metabolic behavior. These methods include the widely employed Flux Balance Analysis (FBA) [41] and Flux Variability Analysis (FVA) [46], but also methods integrating fluxomic, transcriptomic, or proteomic data (see Box 2) [41]. For an extensive overview of mathematical methods using metabolic models we refer to Zomorodi *et al.* [47]. Here, we focus on recent developments in the modeling of microalgae specifically.

Metabolic models of microalgae reflect the modeling counterpart of their current annotation; therefore, inconsistencies between model predictions and experimental findings indicate missing and/or poor annotations. For example, experimentally identified metabolites were compared to metabolites that could be produced in metabolic reconstructions of *C. reinhardtii* [48, 49] (Figure 2). Metabolites found experimentally but not in the models initiated pathway elucidation and identification of the corresponding genes, and thereby led to an improved genome annotation [48]. This procedure was automated by Christian *et al.*, who designed a

1 gap-filling method to identify reactions allowing production in a model of experimentally-  
2 detected metabolites [49]. These updated reactions and annotations [48, 49] were  
3 subsequently stored in ChlamyCyc [27], allowing continuous expansion of the database.  
4 Concurrently, a separate *C. reinhardtii* metabolic model, iAM303, was created, in which the  
5 included open reading frames were experimentally validated. This led both to improved  
6 structural genome annotation, and to additional support for the reactions included in the  
7 model [50]. This model was greatly expanded in iRC1080 in 2011 and additional ORFs were  
8 validated [51]. The predictive power of the latter model was tested for 30 environmental  
9 conditions and 14 gene knockouts. In addition, iRC1080 predicted essential genes (lethal  
10 phenotype upon knockout) under varying experimental conditions, although these predictions  
11 remain to be verified [51]. Recently GSMMs for *Ostreococcus tauri* and *Ostreococcus lucimarinus*  
12 have been constructed [52] (Figure 2), demonstrating expansion in the field. The initial models,  
13 based on the available gene annotations, revealed that these could not account for the  
14 production of many biomass constituents [52]. The gap-filling method designed in [49] was  
15 subsequently employed to find suitable reactions for the production of these metabolites [52].

16 It is well recognized that the exact choice of growth conditions is highly important to attain  
17 desired metabolic activities. Metabolic models can explore how different growth conditions  
18 affect metabolism and identify theoretically optimal conditions for a given metabolic objective.  
19 For example, multiple metabolic models of *C. reinhardtii* were used to simulate metabolism  
20 under autotrophic, heterotrophic and mixotrophic conditions to verify model predictions [53],  
21 to investigate how metabolite production is influenced [53, 54], and to contrast mutant strains  
22 [51]. *C. reinhardtii* metabolic models were also used to determine how quantity of light [51, 55,  
23 56] and its spectral composition [51] affect metabolism. Of particular interest is the possibility  
24 to predict an optimal light spectrum for a given metabolic goal [51]. In contrast to these

1 successful models of *C. reinhardtii*, the metabolism of other algae is only poorly understood. For  
2 example, some industrially relevant algae can currently not be grown efficiently without  
3 bacterial presence [57]. Potentially, these algae and associated bacteria can be modeled  
4 simultaneously in order to deduce their relationship, as has been done for microbial  
5 communities [58, 59].

6 The most comprehensive algal metabolic models to date are iRC1080 [51] and AlgaGEM [53],  
7 which are GSMMs and account for various cellular compartments. However, they vary in degree  
8 of compartmentalization (Figure 2). In iRC1080, half (865/1730) of the non-transport  
9 reactions occur in cellular compartments other than the cytosol. In contrast, this is only about  
10 12% (201/1617) for AlgaGEM. This reflects that independently generated GSMMs for the same  
11 organism can differ significantly in their representation of metabolism, as different sources of  
12 information are included. By combining the information from all currently available *C.*  
13 *reinhardtii* metabolic models, as well as from improved annotation methods, a single and more  
14 comprehensive GSMM may be obtained. This consensus *C. reinhardtii* GSMM would be an  
15 important starting point for the generation of GSMMs for other interesting microalgae, with the  
16 proviso mentioned earlier that it might not be applicable to distantly related microalgae.  
17 Alternatively, *ab initio* models can be made using genome data for the alga in question, but  
18 employing the strategies and tools developed for *C. reinhardtii*, as has been done for  
19 *Ostreococcus* [52]. Ultimately, GSMMs of various microalgae will be valuable for designing  
20 strategies that increase the production of compounds of interest [47, 60]. This combined with  
21 the design of novel synthetic pathways, such as the species-independent prediction  
22 demonstrated for novel isobutanol, 3-hydroxypropionate, and butyryl-CoA biosynthesis [61],  
23 will pave the way for model-driven engineering of algal species.

**Figure 1: Overview of metabolic models of microalgae.** Green boxes represent *C. reinhardtii*

GSMs, the red box represents an *Arabidopsis thaliana* GSM associated with a *Chlamydomonas reinhardtii* GSM, and the blue box represents two *Ostreococcus* GSMs. A connection between two GSMs indicates that the former was used in the reconstruction of the latter. The boxes are annotated with the model names if available and otherwise with the author name(s): Christian *et al.* [49], Boyle & Morgan [54], AraGEM [62], iAM303 [50], Cogne *et al.* [56], Kliphuis *et al.* [55], AlgaGEM [53], iRC1080 [51], Krumholz *et al.* [52]. The numbers in each box indicate the total number of reactions (R), total number of genes (G), unique decompartmentalized metabolites (M), and biological cellular compartments (C) found in available model files. The pie charts depict the distribution of biochemical reactions among different compartments as well as compartment-spanning transport reactions (reaction categories are shown in the legend). The category 'others' refers to the following compartments: flagellum, Golgi apparatus, thylakoid lumen, nucleus, and eyespot.

\*: Additional information obtained from authors.

\*\*: Gene information not available from model files.



## Integrating bioinformatics and modeling for algal biotechnology

Improvements in algal annotations will need to interact closely with systems modeling of the metabolic and regulatory networks in order to refine our understanding of the capabilities of a specific alga, and to provide a basis for applications in biotechnology. Figure 2 shows the connection between the various stages in bioinformatics and systems biology modeling. New algal genomic, transcriptomic, and proteomic data are collected (step 1), allowing identification of genes and proteins (step 2). After 1st-generation high-throughput functional annotation (step 3), a refinement step using 2<sup>nd</sup>-generation function annotation algorithms (step 4) is applied. The bioinformatics annotation itself is an iterative process for genes and proteins until they are deemed sufficient (step 5). These annotations (step 6), as well as data available from public databases and the literature (step 7), are then used by systems biology modeling to reverse-engineer a GSMM (step 8) to study metabolic interactions in different circumstances in detail. After attaining a GSMM, experimental validation of the metabolic model (step 9) should be performed to validate model predictions or pinpoint inaccuracies and knowledge gaps. Depending on these results, additional omics data or refinement of annotation is required. Due to the low number of experimentally validated algal proteins, the feedback loop from algal modeling back to genes/proteins function prediction plays a significant role in strengthening the knowledge foundation, which will ultimately underpin efficient engineering of algal genomes for industrial product synthesis. Once an algal GSMM is constructed, it should be made available in a common public database and literature.

The GSMMs provide a basis for both computational and lab-driven experiments, assisting in the discovery of biotechnology-driven solutions for genetic bottlenecks in algae. For example, to enable microalgae to become a viable industrial biosynthesis platform, their photosynthetic

1 efficiency, product yield, and growth rates under conditions for product synthesis, will need to  
2 be addressed. Photosynthetic efficiency, with an estimated maximum of 8-9% in wild-type  
3 algae [63, 64], sets a limit to both product synthesis and growth rate. Because of efficient light-  
4 harvesting antenna, algal cells can absorb much more light than they are able to use for  
5 photosynthesis [56], with the excess lost as heat or fluorescence. In dense algal cultures, such  
6 as might be found in industrial cultivation systems, this reduces light penetration, placing a  
7 limit on the depth of the culture, and increasing the surface area to volume ratio required for  
8 maximum productivity. Truncated light-harvesting chlorophyll antenna size (*tla*) mutation with  
9 reduced antenna size in *C. reinhardtii* has been shown to improve solar energy conversion  
10 efficiency and photosynthetic productivity in mass culture and bright light [65]. Another study  
11 has modeled different pathways for the process of carbon fixation [66], as a means to overcome  
12 the low oxygenase activity of Rubisco [67]. Bar-Even *et al.* [66] computationally identified  
13 alternative carbon fixation pathways by using approximately 5,000 known metabolic enzymes,  
14 hoping to find carbon fixation pathways with superior kinetics, energy efficiency, and topology.  
15 Some of their proposed pathways were estimated to be up to two to three times more efficient  
16 than the conventional Calvin-Benson cycle. Using an algal GSM to study these pathways would  
17 help to understand how these predictions may affect biomass and product synthesis in  
18 microalgae.

19 As explained earlier, nitrogen limitation is a necessary stimulus for TAG accumulation by  
20 microalgae [7]. This also triggers a reduction in photosynthetic membrane lipids and cessation  
21 of cell growth. The link between TAG accumulation and macronutrient stress has been  
22 investigated using a systems approach, leading to the identification of a putative N-triggered  
23 transcription factor in *C. reinhardtii*, the overexpression of which resulted in about 50%  
24 increase in lipid production under specific experimental conditions [68]. In another approach,  
25 in the diatom, *Thalassiosira pseudonana*, TAG production was increased via an RNAi knockdown

1 strategy targeting not the biosynthesis of lipids, or the production of competing energy sinks,  
2 but instead targeting lipases, involved in glycerolipid catabolism [69]. The integration of  
3 knowledge gained from GSMMs and similar metabolic engineering offers scope for improved  
4 efficiency, based on rational design. For example, farnesyl pyrophosphate is a precursor of  
5 terpenoids, steroids, and carotenoids, and the metabolite itself is also a product of interest in  
6 algae. Bacterial promoters responsive to the toxic accumulation of farnesyl pyrophosphate have  
7 been identified and used to regulate the expression of the precursor biosynthesis operon. This  
8 increased the yield of amorphadiene two fold over chemically inducible and constitutive gene  
9 expression [70]. Such an approach in microalgae would be foreseeable in the future when  
10 promoters in various algal species are better understood, through model-driven design that  
11 incorporates systems data.

12 In contrast to the use of bacteria and yeasts for industrial production, algal biotechnology is in  
13 its infancy. Alongside genome sequence information, a key requirement is the ability to carry  
14 out genetic transformation, and while this is only routine for *C. reinhardtii*, and the diatoms *P.*  
15 *tricornutum* and *Thalassiosira pseudonana*, in the last few years there has been a rapid increase  
16 in published methods for transformation of several species of industrial interest including  
17 *Nannochloropsis* sp. [71]. Moreover, the ability to engineer the chloroplast genome offers  
18 considerable opportunities for metabolic engineering, given the focus of this organelle on  
19 biosynthesis [72]. But for predictive metabolic engineering there is an urgent need to expand  
20 the toolbox, particularly for the regulation of transgene expression. In this context, there are a  
21 number of well-established systems for inducible gene expression in *C. reinhardtii*, most  
22 notably promoters that are regulated in response to nitrate (*NIT1* or *NIA1*) [73] or , copper  
23 (*CYC6*) [74]. More recently, vitamin responsive *cis*-elements have been identified, namely a  
24 cobalamin (vitamin B<sub>12</sub>) responsive promoter [75], as well as a thiamine (vitamin B<sub>1</sub>)

responsive riboswitch [76], have been demonstrated as useful regulatory tools. Vitamins present advantages of being benign, cheap and effective at low concentrations. However, the majority of these have been discovered by coincidence rather than design, and a more rational approach will come from use of transcriptomic data to provide promoters responsive to particular regulators, for example in response to CO<sub>2</sub> levels. [77]. Further facilitation of transgene expression comes from the use of 2A peptides [78], which cause self-cleavage to release individual domains from a fusion protein. They thus provide the capacity for operon-like transgene expression within the nucleus. Marker recycling methods for chloroplast engineering have also been developed for *C. reinhardtii* [72, 79]. In spite of these developments, progress remains parallel in nature, and heavily focused upon the development of *C. reinhardtii*.

For microalgae to develop as a biotechnology platform, rational design to address their current shortcomings must be achieved through the development of fit-for-purpose metabolic engineering or synthetic biology resources. The relative immaturity of the field combined with the enticing potential of integrating predictive design of microalgae with the bioinformatics and systems biology modeling framework (Figure 2) offers new perspectives for future improvements in algal biotechnology. The current prominent challenges in algal bioinformatics and genome-scale modeling are the foundation for overcoming the knowledge barrier to enable predictive modifications of various algal genomes in the future.

**Figure 2: A multidisciplinary workflow for integrative and systematic understanding of algae.**

Black arrow: *in silico* data or predictions; white arrow: experimental (wet-lab) data.

## 1    **Concluding remarks**

2    The significant gap of unknown and non-validated gene and protein functions in algae remains  
3    one of the top challenges faced by scientists wanting to tap further into the potential of these  
4    organisms for sustainable biosynthesis. Predictive design of metabolic engineering strategies  
5    for microalgae still has a long journey ahead. An improved understanding of the metabolism,  
6    regulation, and growth of algae, together with their interactions with co-existing bacteria, is a  
7    crucial first step. Extending bioinformatics approaches for function prediction through  
8    incorporation of new methodology, integrated and flexible databases, and combination with  
9    metabolic modeling and model-driven design of experiments at the systems biology level, will  
10    underpin this process, and enable the future era of algal industrial biotechnology.

11

## 12    **Acknowledgement**

13    We acknowledge support from the EU FP7 project SPLASH (Sustainable PoLymers from Algae Sugars  
14    and Hydrocarbons), grant agreement number 311956.

## 1   References

- 2   1       Merchant, S.S., *et al.* (2007) The *Chlamydomonas* genome reveals the evolution of key animal and  
3       plant functions. *Science* 318, 245-250
- 4   2       Borowitzka, M.A. (2013) High-value products from microalgae — their development and  
5       commercialisation. *J. Appl. Phycol.* 25, 743-756
- 6   3       Scott, S.A., *et al.* (2010) Biodiesel from algae: challenges and prospects. *Curr. Opin. Biotechnol.* 21,  
7       277-286
- 8   4       Wijffels, R.H., *et al.* (2010) Microalgae for the production of bulk chemicals and biofuels. *Biofuels*,  
9       *Bioprod. Biorefin.* 4, 287-295
- 10  5       Merchant, S.S., *et al.* (2012) TAG, you're it! *Chlamydomonas* as a reference organism for  
11       understanding algal triacylglycerol accumulation. *Curr. Opin. Biotechnol.* 23, 352-363
- 12  6       Mata, T.M., *et al.* (2010) Microalgae for biodiesel production and other applications: a review.  
13       *Renewable Sustainable Energy Rev.* 14, 217-232
- 14  7       Klok, A.J., *et al.* (2013) Simultaneous growth and neutral lipid accumulation in microalgae.  
15       *Bioresour. Technol.* 134, 233-243
- 16  8       Busi, M.V., *et al.* (2013) Starch metabolism in green algae. *Starch - Stärke* 66, 28-40
- 17  9       Ho, S.H., *et al.* (2013) Bioethanol production using carbohydrate-rich microalgae biomass as  
18       feedstock. *Bioresour. Technol.* 135, 191-198
- 19  10       Hildebrand, M., *et al.* (2013) Metabolic and cellular organization in evolutionarily diverse  
20       microalgae as related to biofuels production. *Curr. Opin. Chem. Biol.* 17, 506 - 514
- 21  11       Radivojac, P., *et al.* (2013) A large-scale evaluation of computational protein function prediction.  
22       *Nat. Methods* 10, 221-227
- 23  12       Falda, M., *et al.* (2012) Argot2: a large scale function prediction tool relying on semantic  
24       similarity of weighted Gene Ontology terms. *BMC Bioinf.* 13, S14
- 25  13       Altschul, S.F., *et al.* (1997) Gapped BLAST and PSI-BLAST: a new generation of protein database  
26       search programs. *Nucleic Acids Res.* 25, 3389-3402

1 14 Finn, R.D., *et al.* (2011) HMMER web server: interactive sequence similarity searching. *Nucleic*  
2 *Acids Res.* 39, W29-W37

3 15 Kourmpetis, Y.A., *et al.* (2010) Bayesian Markov Random Field analysis for protein function  
4 prediction based on network data. *PLoS One* 5, e9293

5 16 Engelhardt, B.E., *et al.* (2011) Genome-scale phylogenetic function annotation of large and  
6 diverse protein families. *Genome Res* 21, 1969-1980

7 17 Cozzetto, D., *et al.* (2013) Protein function prediction by massive integration of evolutionary  
8 analyses and multiple data sources. *BMC Bioinf.* 14, S1

9 18 Wong, A. and Shatkay, H. (2013) Protein function prediction using text-based features extracted  
10 from the biomedical literature: The CAFA Challenge. *BMC Bioinf.* 14, S14

11 19 Rentzsch, R. and Orengo, C.A. (2009) Protein function prediction — the power of multiplicity.  
12 *Trends Biotechnol.* 27, 210-219

13 20 Buchan, D.W., *et al.* (2010) Protein annotation and modelling servers at University College  
14 London. *Nucleic Acids Res.* 38, W563-568

15 21 Rodgers-Melnick, E., *et al.* (2013) Predicting whole genome protein interaction networks from  
16 primary sequence data in model and non-model organisms using ENTS. *BMC Genomics* 14, 608

17 22 Franceschini, A., *et al.* (2013) STRING v9.1: protein-protein interaction networks, with increased  
18 coverage and integration. *Nucleic Acids Res.* 41, D808-815

19 23 Emanuelsson, O., *et al.* (2000) Predicting subcellular localization of proteins based on their N-  
20 terminal amino acid sequence. *J. Mol. Biol.* 300, 1005-1016

21 24 Petersen, T.N., *et al.* (2011) SignalP 4.0: discriminating signal peptides from transmembrane  
22 regions. *Nat. Methods* 8, 785-786

23 25 Yu, N.Y., *et al.* (2010) PSORTb 3.0: improved protein subcellular localization prediction with  
24 refined localization subcategories and predictive capabilities for all prokaryotes. *Bioinformatics*  
25 26, 1608-1615

26 26 Tardif, M., *et al.* (2012) PredAlgo: A new subcellular localization prediction tool dedicated to  
27 green algae. *Mol. Biol. Evol.* 29, 3625-3639

1 27 May, P., *et al.* (2009) ChlamyCyc: an integrative systems biology database and web-portal for  
2 *Chlamydomonas reinhardtii*. *BMC Genomics* 10, 209

3 28 Caspi, R., *et al.* (2014) The MetaCyc database of metabolic pathways and enzymes and the BioCyc  
4 collection of Pathway/Genome Databases. *Nucleic Acids Res.* 42, D459-D471

5 29 Keller, O., *et al.* (2011) A novel hybrid gene prediction method employing protein multiple  
6 sequence alignments. *Bioinformatics* 27, 757-763

7 30 Lopez, D., *et al.* (2011) Algal Functional Annotation Tool: a web-based analysis suite to  
8 functionally interpret large gene lists using integrated annotation and expression data. *BMC*  
9 *Bioinf.* 12, 282

10 31 Goodstein, D.M., *et al.* (2012) Phytozome: a comparative platform for green plant genomics.  
11 *Nucleic Acids Res.* 40, D1178-D1186

12 32 Bogen, C., *et al.* (2013) Reconstruction of the lipid metabolism for the microalga *Monoraphidium*  
13 *neglectum* from its genome sequence reveals characteristics suitable for biofuel production. *BMC*  
14 *Genomics* 14, 926

15 33 Desbois, A., *et al.* (2008) Isolation and structural characterisation of two antibacterial free fatty  
16 acids from the marine diatom, *Phaeodactylum tricornutum*. *Appl. Microbiol. Biotechnol.* 81, 755-  
17 764

18 34 Feng, Y., *et al.* (2011) Lipid production of *Chlorella vulgaris* cultured in artificial wastewater  
19 medium. *Bioresour. Technol.* 102, 101-105

20 35 Prieto, A., *et al.* (2011) Assessment of carotenoid production by *Dunaliella salina* in different  
21 culture systems and operation regimes. *J. Biotechnol.* 151, 180-185

22 36 Doucha, J. and Lívanský, K. (2012) Production of high-density *Chlorella* culture grown in  
23 fermenters. *J. Appl. Phycol.* 24, 35-43

24 37 Ambati, R.R., *et al.* (2014) Astaxanthin: sources, extraction, stability, biological activities and its  
25 commercial applications — a review. *Mar. Drugs.* 12, 128-152

26 38 Weiss, T.L., *et al.* (2011) Genome size and phylogenetic analysis of the A and L races of  
27 *Botryococcus braunii*. *J. Appl. Phycol.* 23, 833-839



1 39 Weiss, T.L., *et al.* (2012) Colony organization in the green alga *Botryococcus braunii* (race B) is  
2 specified by a complex extracellular matrix. *Eukaryot. Cell* 11, 1424-1440

3 40 Jones, P., *et al.* (2014) InterProScan 5: genome-scale protein function classification.  
4 *Bioinformatics* 30, 1236-1240

5 41 Orth, J.D., *et al.* (2010) What is flux balance analysis? *Nat. Biotechnol.* 28, 245-248

6 42 Wiechert, W. (2001) <sup>13</sup>C metabolic flux analysis. *Metab. Eng.* 3, 195-206

7 43 Weitzel, M., *et al.* (2013) 13CFLUX2 — high-performance software suite for <sup>13</sup>C-metabolic flux  
8 analysis. *Bioinformatics* 29, 143-145

9 44 Jensen, P., *et al.* (2011) TIGER: toolbox for integrating genome-scale metabolic models,  
10 expression data, and transcriptional regulatory networks. *BMC Syst. Biol.* 5, 147

11 45 Blazier, A.S. and Papin, J.A. (2012) Integration of expression data in genome-scale metabolic  
12 network reconstructions. *Front. Physiol.* 3, 299

13 46 Mahadevan, R. and Schilling, C.H. (2003) The effects of alternate optimal solutions in constraint-  
14 based genome-scale metabolic models. *Metab. Eng.* 5, 264-276

15 47 Zomorodi, A.R., *et al.* (2012) Mathematical optimization applications in metabolic networks.  
16 *Metab. Eng.* 14, 672-686

17 48 May, P., *et al.* (2008) Metabolomics- and proteomics-assisted genome annotation and analysis of  
18 the draft metabolic network of *Chlamydomonas reinhardtii*. *Genetics* 179, 157-166

19 49 Christian, N., *et al.* (2009) An integrative approach towards completing genome-scale metabolic  
20 networks. *Mol. Biosyst.* 5, 1889-1903

21 50 Manichaikul, A., *et al.* (2009) Metabolic network analysis integrated with transcript verification  
22 for sequenced genomes. *Nat. Methods* 6, 589-592

23 51 Chang, R.L., *et al.* (2011) Metabolic network reconstruction of *Chlamydomonas* offers insight into  
24 light-driven algal metabolism. *Mol. Syst. Biol.* 7, 518

25 52 Krumholz, E.W., *et al.* (2012) Genome-wide metabolic network reconstruction of the picoalga  
26 *Ostreococcus*. *J. Exp. Bot.* 63, 2353-2362

1 53 Gomes de Oliveira Dal'Molin, C., *et al.* (2011) AlgaGEM — a genome-scale metabolic  
2 reconstruction of algae based on the *Chlamydomonas reinhardtii* genome. *BMC Genomics*  
3 12(suppl. 4), S5

4 54 Boyle, N. and Morgan, J. (2009) Flux balance analysis of primary metabolism in *Chlamydomonas*  
5 *reinhardtii*. *BMC Syst. Biol.* 3, 4

6 55 Kliphuis, A.J., *et al.* (2012) Metabolic modeling of *Chlamydomonas reinhardtii*: energy  
7 requirements for photoautotrophic growth and maintenance. *J. Appl. Phycol.* 24, 253-266

8 56 Cogne, G., *et al.* (2011) A model-based method for investigating bioenergetic processes in  
9 autotrophically growing eukaryotic microalgae: application to the green algae *Chlamydomonas*  
10 *reinhardtii*. *Biotechnol. Progr.* 27, 631-640

11 57 Rivas, M.O., *et al.* (2010) Interactions of *Botryococcus braunii* cultures with bacterial biofilms.  
12 *Microb. Ecol.* 60, 628-635

13 58 Zomorodi, A.R. and Maranas, C.D. (2012) OptCom: a multi-level optimization framework for the  
14 metabolic modeling and analysis of microbial communities. *PLoS Comput. Biol.* 8, e1002363

15 59 Zomorodi, A.R., *et al.* (2014) d-OptCom: dynamic multi-level and multi-objective metabolic  
16 modeling of microbial communities. *ACS Synth. Biol.* 3, 247-257

17 60 Tomar, N. and De, R.K. (2013) Comparing methods for metabolic network analysis and an  
18 application to metabolic engineering. *Gene* 521, 1-14

19 61 Cho, A., *et al.* (2010) Prediction of novel synthetic pathways for the production of desired  
20 chemicals. *BMC Syst. Biol.* 4, 35

21 62 Gomes de Oliveira Dal'Molin, C., *et al.* (2010) AraGEM, a genome-scale reconstruction of the  
22 primary metabolic network in *Arabidopsis*. *Plant Physiol.* 152, 579-589

23 63 Chisti, Y. (2013) Constraints to commercialization of algal fuels. *J. Biotechnol.* 167, 201-214

24 64 Lehr, F. and Posten, C. (2009) Closed photo-bioreactors as tools for biofuel production. *Curr.*  
25 *Opin. Biotechnol.* 20, 280-285

1 65 Kirst, H., *et al.* (2012) Truncated photosystem chlorophyll antenna size in the green microalga  
2 *Chlamydomonas reinhardtii* upon deletion of the TLA3-CpSRP43 gene. *Plant Physiol.* 160, 2251-  
3 2260

4 66 Bar-Even, A., *et al.* (2010) Design and analysis of synthetic carbon fixation pathways. *Proc. Natl.*  
5 *Acad. Sci. USA* 107, 8889-8894

6 67 Whitney, S.M., *et al.* (2011) Advancing our understanding and capacity to engineer nature's CO<sub>2</sub>-  
7 sequestering enzyme, Rubisco. *Plant Physiol.* 155, 27-35

8 68 Yohn, C., *et al.* (2011) Sapphire Energy, Inc. Stress-induced lipid trigger. Application no.:  
9 11740278.4. Published source: European Patent Register.

10 69 Trentacoste, E.M., *et al.* (2013) Metabolic engineering of lipid catabolism increases microalgal  
11 lipid accumulation without compromising growth. *Proc. Natl. Acad. Sci. USA* 110, 19748-19753

12 70 Dahl, R.H., *et al.* (2013) Engineering dynamic pathway regulation using stress-response  
13 promoters. *Nat. Biotechnol.* 31, 1039-1046

14 71 Kilian, O., *et al.* (2011) High-efficiency homologous recombination in the oil-producing alga  
15 *Nannochloropsis* sp. *Proc. Natl. Acad. Sci. USA* 108, 21265-21269

16 72 Purton, S., *et al.* (2013) Genetic engineering of algal chloroplasts: progress and prospects. *Russ. J.*  
17 *Plant Physiol.* 60, 491-499

18 73 Ohresser, M., *et al.* (1997) Expression of the arylsulphatase reporter gene under the control of  
19 the *nit1* promoter in *Chlamydomonas reinhardtii*. *Curr. Genet.* 31, 264-271

20 74 Quinn, J.M., *et al.* (2003) Copper response element and Crr1-dependent Ni<sup>2+</sup>-responsive  
21 promoter for induced, reversible gene expression in *Chlamydomonas reinhardtii*. *Eukaryot. Cell* 2,  
22 995-1002

23 75 Helliwell, K.E., *et al.* (2014) Unraveling vitamin B<sub>12</sub>-responsive gene regulation in algae. *Plant*  
24 *Physiol.* 165, 388-397

25 76 Ramundo, S., *et al.* (2013) Repression of essential chloroplast genes reveals new signaling  
26 pathways and regulatory feedback loops in *Chlamydomonas*. *Plant Cell* 25, 167-186

1    77    Fang, W., *et al.* (2012) Transcriptome-wide changes in *Chlamydomonas reinhardtii* gene  
2           expression regulated by carbon dioxide and the CO<sub>2</sub>-concentrating mechanism regulator  
3           CIA5/CCM1. *Plant Cell* 24, 1876-1893

4    78    Rasala, B.A., *et al.* (2012) Robust expression and secretion of xylanase 1 in *Chlamydomonas*  
5           *reinhardtii* by fusion to a selection gene and processing with the FMDV 2A peptide. *PLoS One* 7,  
6           e43349

7    79    Day, A. and Goldschmidt-Clermont, M. (2011) The chloroplast transformation toolbox: selectable  
8           markers and marker removal. *Plant Biotechnol. J.* 9, 540-553

9

10

1    **Highlights**

- 2        •    Microalgae are potential hosts for industrial biosynthesis of valuable compounds.
- 3        •    Genome sequences of many microalgae are available, but annotation lags behind.
- 4        •    We propose an integrative approach to improve algal proteins annotation.
- 5        •    Systems biology modeling of microalgae is crucial to facilitate algal engineering.

6

Figure 1

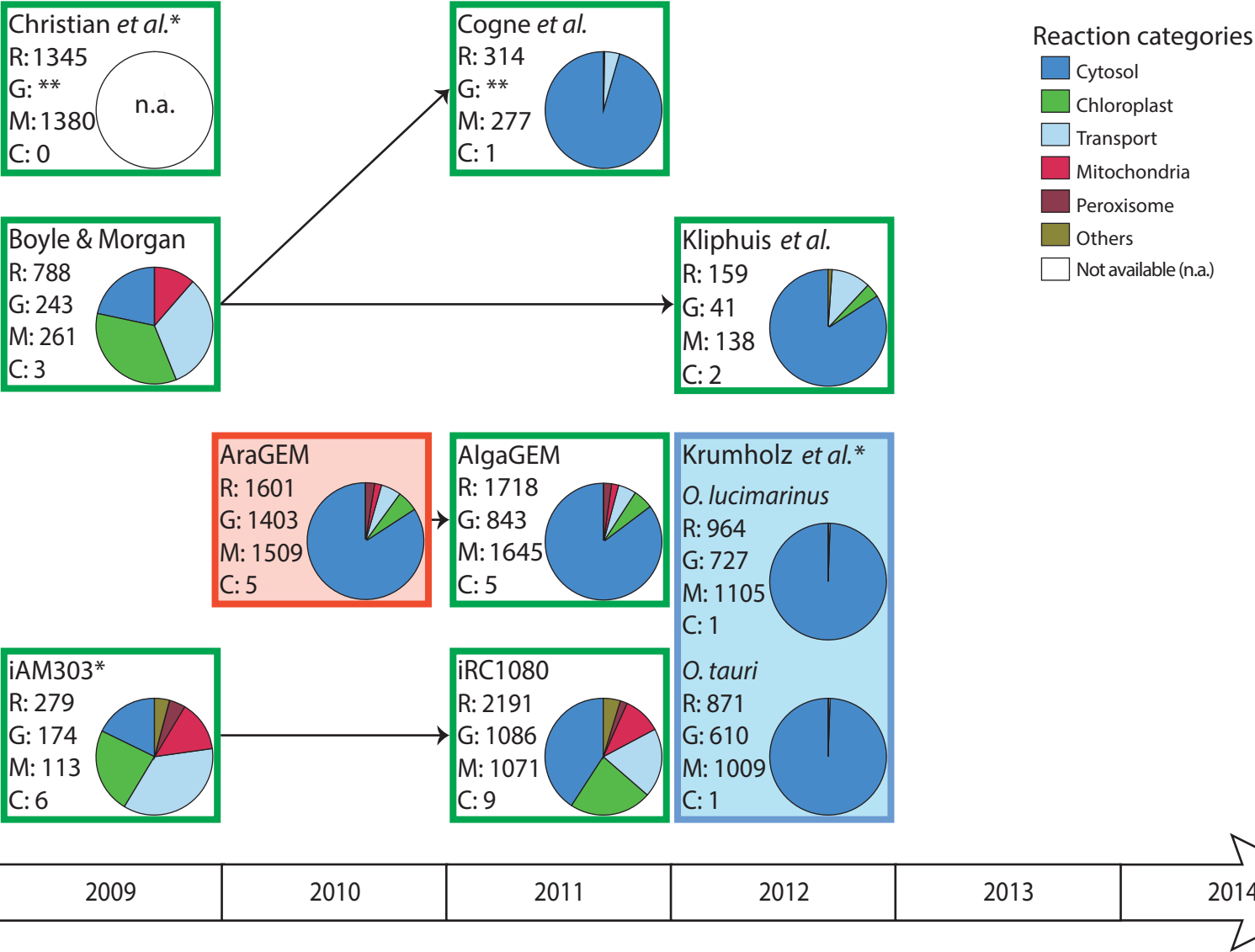
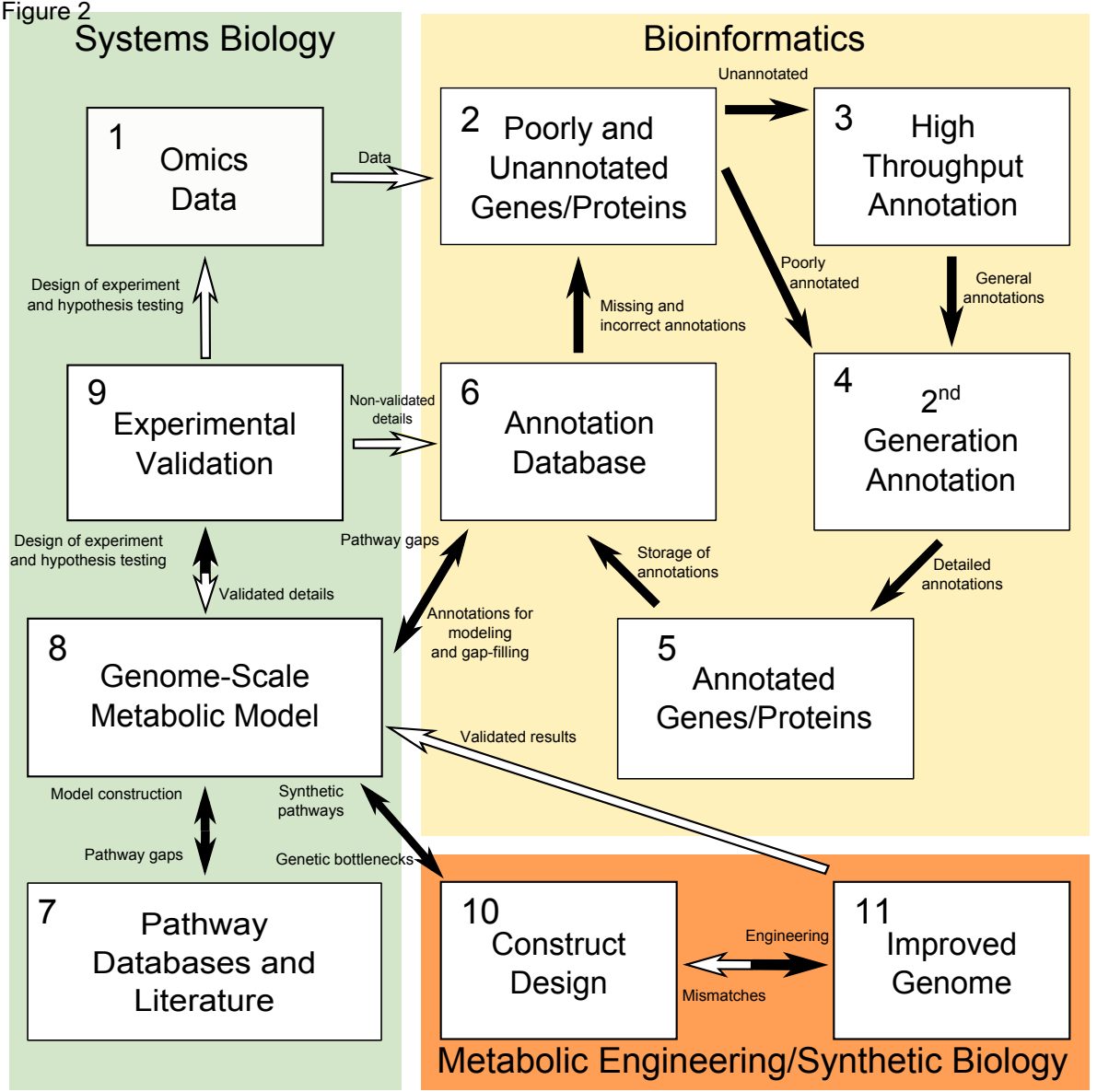


Figure 2



Original Figure File\_Fig.1

[Click here to download Original Figure File: Figure\\_1\\_Adobe Illustrator.ai](#)



Original Figure File

[Click here to download Original Figure File: Figure\\_2\\_Inkscape.svg](#)