

# MS-DIAL: data-independent MS/MS deconvolution for comprehensive metabolome analysis

Hiroshi Tsugawa<sup>1,2</sup>, Tomas Cajka<sup>3</sup>, Tobias Kind<sup>3</sup>, Yan Ma<sup>3</sup>, Brendan Higgins<sup>4</sup>, Kazutaka Ikeda<sup>5,6</sup>, Mitsuhiro Kanazawa<sup>7</sup>, Jean VanderGheynst<sup>4</sup>, Oliver Fiehn<sup>3,8</sup> & Masanori Arita<sup>1,9</sup>

**Data-independent acquisition (DIA) in liquid chromatography (LC) coupled to tandem mass spectrometry (MS/MS) provides comprehensive untargeted acquisition of molecular data. We provide an open-source software pipeline, which we call MS-DIAL, for DIA-based identification and quantification of small molecules by mass spectral deconvolution. For a reversed-phase LC-MS/MS analysis of nine algal strains, MS-DIAL using an enriched LipidBlast library identified 1,023 lipid compounds, highlighting the chemotaxonomic relationships between the algal strains.**

Precursor- or data-independent MS/MS acquisition methods in LC-MS/MS for the untargeted analyses of biomolecules<sup>1,2</sup> have recently received considerable attention. In contrast to traditional data-dependent MS/MS acquisitions, DIA methods can obtain all fragment ions for all precursors simultaneously, thereby increasing the coverage of observable molecules and reducing the identification of false negatives. The difficulty with this approach, however, is the contamination of MS/MS spectra due to its wider isolation window (10–25 Da or more) for precursor-ion selection. Moreover, the DIA process dissociates the link between precursors and their fragment ions, compromising the molecular identification process.

In the field of proteomics, the OpenSWATH software has partly addressed these problems<sup>2</sup>. After extraction of product-ion chromatograms for the corresponding precursor range, chromatogram peaks are grouped, scored and statistically assessed on the basis of false-discovery rate (FDR) in the mProphet algorithm<sup>3</sup>. Unfortunately, this approach is not directly applicable to metabolomics. Whereas spectral similarity in shotgun proteomics is probabilistically estimated by presence or absence of peak groups,

compound annotations in metabolomics rely on overall match scores between experimental and library spectra. In addition, schemes for FDR calculation by validated decoy techniques do not exist in metabolomics<sup>4</sup>. Therefore, DIA MS/MS spectra must be purified from fragment ions of coeluting compounds and from noise ions in order for metabolomic annotations to achieve high overall library-matching scores.

A solution to this problem is mathematical deconvolution of fragment ions to extract original spectra and to reassociate the precursor-fragment links. Nikolskiy *et al.*<sup>5</sup> reported such a deconvolution approach, but their program, decoMS2, requires two different collision energies, low (usually 0 V) and high, in each precursor range to solve the mathematical equations. Interestingly, automatic mass spectral deconvolution and identification systems are routine today in gas chromatography coupled to mass spectrometry (GC-MS)<sup>6,7</sup>. DIA-type mass fragmentation schemes are the norm in hard electron-ionization GC-MS in contrast to soft electrospray-ionization LC-MS/MS. Analogous to these successful GC-MS data processing systems, we have developed the Mass Spectrometry–Data Independent AnaLysis software (MS-DIAL) that implements a new deconvolution algorithm for DIA data sets. It is a data-processing pipeline for untargeted metabolomics applicable to either data-independent or precursor-dependent MS/MS fragmentation methods.

In MS-DIAL, the raw vendor-format data or the common mzML data are first converted into the Analysis Base File (ABF) format for rapid data retrieval<sup>8</sup> (Fig. 1a). Then, precursor-ion peaks are efficiently spotted (hereafter ‘peak spotting’) by exploring two continuous data axes: retention time (Rt) and accurate mass ( $m/z$ ). Each spot represents a detected peak (Fig. 1b), and our MS<sup>2</sup>Dec algorithm is applied to each spot to deconvolute spectra in the respective precursor-ion range. The MS<sup>2</sup>Dec algorithm first extracts the product spectra for each precursor peak on all MS/MS chromatograms (the raw chromatograms are shown in regular lines in Fig. 1c) to recover the precursor-product links as a result of deconvolution, which is itself based on its established GC-MS counterpart<sup>6</sup> with substantial modifications based on accurate mass information instead of nominal masses. This enables the analyses of large-scale DIA data sets. MS<sup>2</sup>Dec uses the least-squares optimization to extract ‘model peaks’ (see Online Methods) in MS/MS chromatograms (the reconstructed model chromatograms are shown as thick lines in Fig. 1c). Finally, the pure MS/MS spectrum is determined by the peak heights of reconstructed chromatograms, thus removing the background noise and extracting the spectrum out of coeluted metabolites (Fig. 1c). Compound identification is achieved through analyses

<sup>1</sup>RIKEN Center for Sustainable Resource Science, Yokohama, Kanagawa, Japan. <sup>2</sup>Department of Biotechnology, Graduate School of Engineering, Osaka University, Suita, Osaka, Japan. <sup>3</sup>Genome Center, University of California Davis, Davis, California, USA. <sup>4</sup>Department of Biological and Agricultural Engineering, University of California Davis, Davis, California, USA. <sup>5</sup>RIKEN Center for Integrative Medical Sciences, Yokohama, Japan. <sup>6</sup>Japan Science and Technology Agency, Kawaguchi, Japan. <sup>7</sup>Reifys Inc., Minato-ku, Tokyo, Japan. <sup>8</sup>Department of Biochemistry, Faculty of Sciences, King Abdulaziz University, Jeddah, Saudi Arabia. <sup>9</sup>National Institute of Genetics, Mishima, Japan. Correspondence should be addressed to O.F. (ofiehn@ucdavis.edu) or M.A. (arita@nig.ac.jp).

RECEIVED 22 SEPTEMBER 2014; ACCEPTED 7 MARCH 2015; PUBLISHED ONLINE 4 MAY 2015; DOI:10.1038/NMETH.3393

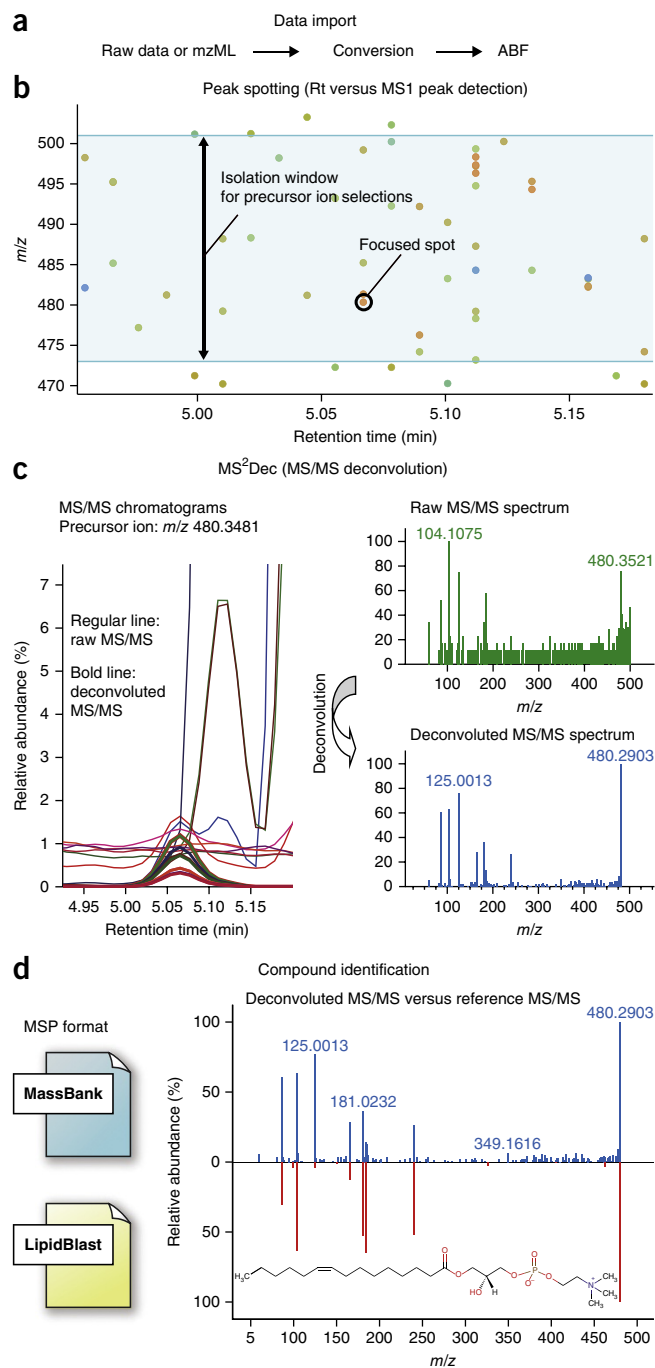
**Figure 1** | Main workflow of the MS-DIAL program. (a) MS vendor format or mzML data is converted to the ABF binary format for rapid data retrieval. (b) Peak spotting (two-dimensional peak detection, see main text) is performed to determine precursor ions for MS/MS spectra. The detected precursor ions are described as spots. The blue background shows the isolation window of precursor ions. Each focused spot is subjected to deconvolution. Rt, retention time (c) The MS<sup>2</sup>Dec deconvolution process includes chromatogram extractions (left, regular lines), model peak constructions (left, bold lines), and mass spectrum reconstructions (right). (d) The MSP format is used for matching experimental mass spectra against mass spectral libraries such as MassBank or LipidBlast. Compounds are identified by analyzing the weighted similarity score of retention time, accurate mass, isotope ratio and MS/MS spectra.

of retention time, mass accuracy, isotope ratio along with MS/MS similarity matching to libraries from publicly available databases (such as MassBank<sup>9</sup> and LipidBlast<sup>10</sup>) (Fig. 1d). MS-DIAL also implements additional functions required for untargeted metabolomics such as peak alignment, filtering and missing-value interpolation (Online Methods).

MS-DIAL is available on Windows (.NET Framework 4.0 or later; RAM: 4.0 GB or more), and the program is downloadable at the PRIME (<http://prime.psc.riken.jp/>) website and as **Supplementary Software**. It supports mzML and major MS vendor formats including those of Agilent Technologies (.D), AB Sciex (.Wiff), Thermo Fisher Scientific (.RAW), Bruker Daltonics (.D), and Waters (.RAW). The program is intended for large-scale analyses such as cohort studies; it accesses raw data sequentially and keeps only their peak information in memory. The actual processing time for an average 600 MB per assay file in our study was less than 1.2 min with an Intel Core i7-4700MQ CPU (2.4 GHz) with 8 GB RAM in Windows 8.1.

Here we used sequential window acquisition of all theoretical mass spectra (SWATH) acquisition as the DIA approach and compared the results to those from traditional data-dependent acquisition (DDA) for validation. We first showcased our MS<sup>2</sup>Dec deconvolution approach with a human plasma sample separated by hydrophilic interaction chromatography (HILIC; Fig. 2 and Online Methods). Two metabolites, metoclopramide and norcocaine, exhibited only a 1.8-s difference in their elution times (at 2.95 and 2.98 min, respectively) and fell within the same 25-Da window of the SWATH acquisition. Although the unique mass spectrum and abundance of metoclopramide could be marginally confirmed in the raw MS/MS spectrum (similarity score 0.72), the spectrum of norcocaine was thoroughly masked under the peaks from metoclopramide (similarity score 0.48) when mass spectral deconvolution was not applied. In this study, the similarity score was calculated by the dot-product scoring method (see Online Methods). MS-DIAL extracted the pure MS/MS spectrum of norcocaine (similarity score 0.80), although contamination of higher-mass peaks (for example,  $m/z$  227) was not completely suppressed. The similarity score of metoclopramide was also improved to 0.86 by deconvolution. Further examples for other metabolites are available in **Supplementary Figure 1**.

We next performed a lipidomic analysis of nine algal species using the LipidBlast library for searching<sup>10</sup>. Prior to the analysis, the library was thoroughly extended to cover major plant and algal lipids such as monogalactosyl, digalactosyl and sulfoquinovosyl diacylglycerols (MGDG, DGDG and SQDG, respectively) and diacylglycerol trimethyl homoserine (DGTS) (**Supplementary Table 1**



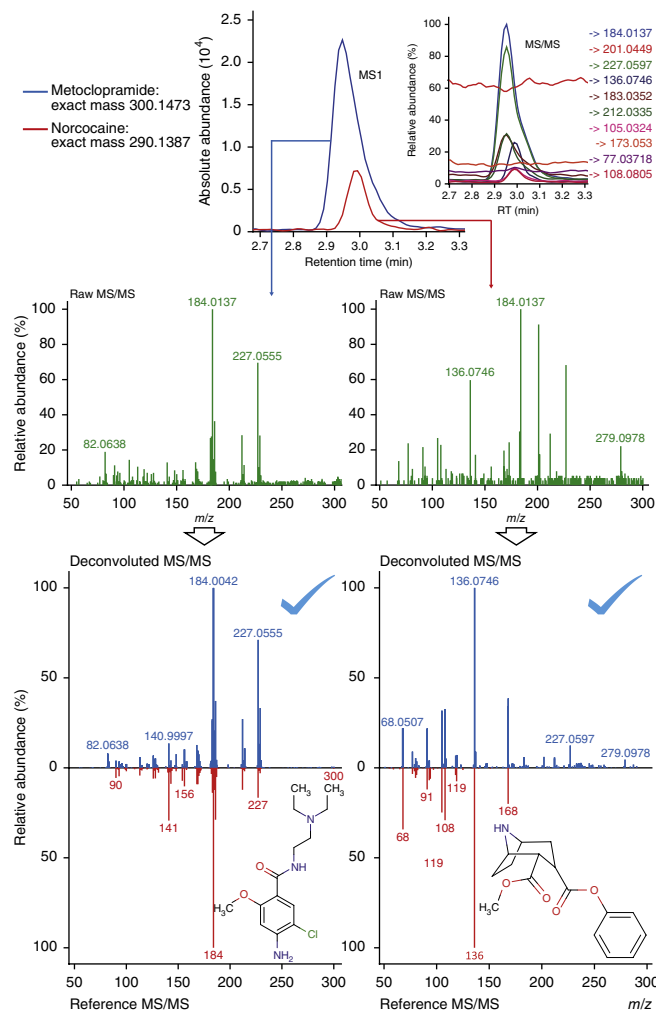
and Online Methods). Moreover, to improve identification accuracies, we predicted the retention times for all molecules in LipidBlast specifically for our chromatography method by partial least-squares regression (PLS-R)<sup>11</sup> on their PaDEL<sup>12</sup> molecular descriptors (Online Methods). Predicted retention times exhibited a standard deviation (s.d.) of 0.14 min when compared to retention times of lipid standards, which was almost equivalent to the regressed s.d. of the experimentally measured dataset (Fig. 3a and **Supplementary Data 1**).

We first tested the overall effect of using MS/MS deconvolution on spectral accuracy for lipid profiling at a 10-ms accumulation time. Indeed, spectral similarity scores were substantially improved by mass spectral deconvolution in comparison to the

**Figure 2** | A deconvolution example using SWATH acquisition with HILIC positive ion mode. Two pharmaceutical agents, metoclopramide and norcocaine, were detected in untargeted metabolomics screens and coeluted within a 1.8-s peak top difference. The MS/MS ion traces with respect to these two metabolites are also shown in the top right panel of the precursor-ion traces. The middle panels show raw MS/MS spectra of metoclopramide (left) and norcocaine (right), respectively. The spectrum of metoclopramide dominates and masks that of norcocaine, making detection of the latter highly difficult. The bottom panels show the deconvoluted MS/MS spectrum and spectra matching results of metoclopramide (left) and norcocaine (right), yielding dot-product scores of 0.80 and 0.86, respectively.

raw centroid spectra generated using a 21-Da isolation window, and approached the quality of the 1-Da isolation window spectra in targeted acquisitions (data-dependent acquisition, DDA) (Fig. 3b). Importantly, the SWATH acquisition with MS-DIAL covered a larger number of phospho- and glycolipids in both positive and negative ionization modes compared to the DDA mode (Fig. 3c and Supplementary Table 2). The only exception was SQDG lipids, whose identification scores worsened because of the low abundance of its characteristic peak ( $m/z$  225) (Supplementary Fig. 2). When we reanalyzed the same sample of *Chlamydomonas reinhardtii* under 30-ms accumulation time and a 65-Da isolation window, the number of identified lipids was notably increased not only for SQDG but also for all other lipid classes (Fig. 3c and Supplementary Table 3). Even for this wide isolation window, the deconvoluted MS/MS spectra kept >90% similarities against the targeted spectra except for SQDG (14:0/16:0) (Supplementary Figs. 2 and 3). This result implies that a wider SWATH window appears preferable for lipid profiling in negative mode, while the precursor-isolation windows and accumulation times may need further optimization. Overall, a total of 1,023 lipids were identified, of which SWATH acquisition covered >90% (Fig. 3c) and yielded 310 additional lipids that were not detected using the data-dependent MS/MS acquisition (Supplementary Data 2 and 3).

We conducted hierarchical clustering analysis (HCA) on the lipidomic profile of all 1,023 distinct lipid molecules from the nine algal species to define their overall similarities (Online Methods). The clustering result was in full concordance with the commonly accepted phylogenetic tree (Fig. 3d). The nine investigated species were found to clearly cluster, and the analysis could distinguish between the five plantae, three chromista and one protozoan species. Plantae species contained mainly 16- or 18-carbon fatty acids, whereas protozoa and chromista were comprised of very-long-chain fatty acids (>18 carbons). Among plantae, *Chlamydomonas* and *Dunaliella* (chlorophyceae) contained DGTS, whereas the two *Chlorella* species (trebouxiophyceae) and UTEX 2341 did not. Very-long-chain polyunsaturated fatty acids (PUFAs) of 20 carbons or more, such as eicosapentanoic or docosahexanoic acid, were mostly identified in *Nannochloropsis oculata* and *Euglena gracilis*. In addition, the total quantity of DGTS and phosphatidic acids (PA) were highly characteristic to these species (Supplementary Figs. 4 and 5). Note that we used the culture collection identification for the green alga strain UTEX 2341 (from The Culture Collection of Algae at The University of Texas at Austin), since its identity as either *Chlorella minutissima* (original classification) or a *Nannochloropsis* species is controversial<sup>13,14</sup>. Interestingly, our chemotaxonomy suggested that UTEX



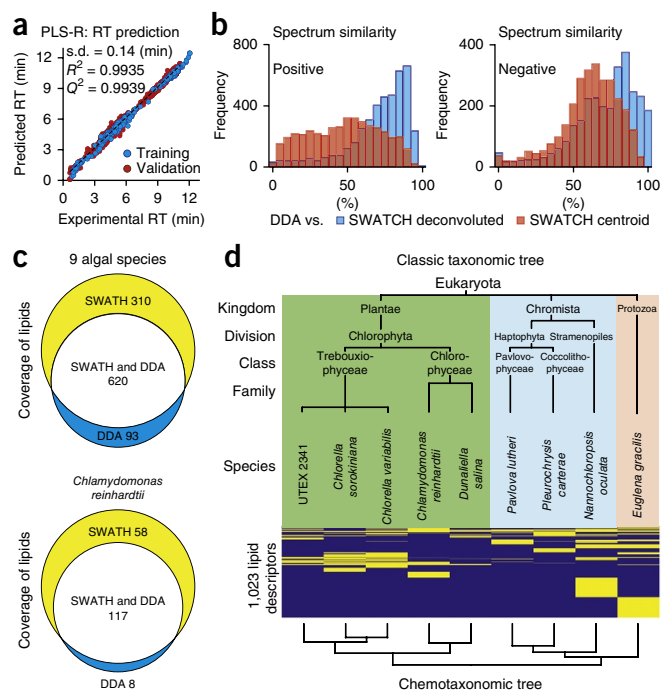
2341 is most probably *C. minutissima*. Moreover, our method led to the detection of lipids that had not been identified previously, such as 18:5 PUFA in *C. reinhardtii*, *N. oculata* and *Pleurochrysis carterae*, odd-chain lipids in all nine algal species and DGTS lipids in *E. gracilis*, *Dunaliella salina* and *N. oculata* (Supplementary Fig. 6). All major lipid classes previously reported<sup>13–16</sup> were identified in our single experiment.

In summary, MS-DIAL resolves entangled MS/MS spectra in SWATH acquisition by a two-step process: precursor-peak spotting followed by MS/MS-level deconvolution. With this software, data-independent MS/MS acquisitions can provide high efficacy and accuracy for metabolome coverage and help address a major bottleneck in metabolomics: compound identification and annotation<sup>17</sup>. It is important to note that, unlike any other software to date, MS-DIAL combines information from four sources (accurate mass, isotope ratios, retention-time prediction and MS/MS fragment matching), exceeding the two orthogonal parameters required by the Metabolomics Standards Initiative<sup>18</sup>. A more rationalized confidence score for each parameter setting and their combination will need to be explored in detail for a variety of matrices<sup>19</sup>.

Because DIA MS/MS information potentially includes all detectable ions, such data sets allow analyses a posteriori and alleviate the cost to reanalyze the same samples with different precursor selections. Although we focused on compound



**Figure 3** | System validation for lipid profiling, lipid coverage and chemotaxonomic relationship of nine algal species. **(a)** The experimental and predicted retention times (RT) of 254 (training) and 1,808 (validation) lipids were plotted along *x* and *y* axes, respectively. Prediction was performed by using PLS-R on 464 properties from the PaDEL-descriptor suite. The  $R^2$ ,  $Q^2$  and s.d. of the validation set were 0.9935, 0.9939 and 0.14 min, respectively. **(b)** Comparison of mass spectra in positive (left) and negative (right) ion modes in commonly identified lipids between SWATH (data-independent) and the traditional data-dependent (DDA) methods. Blue histogram shows the spectra similarity between the deconvoluted and DDA spectra. Red histogram shows the similarity between the centroid (non-deconvoluted) and DDA spectra. **(c)** Venn diagram showing quantitation of lipids identified using the SWATH and DDA methods. Comparative analysis of the nine algal species at 10 ms (SWATH) and 50 ms (DDA) accumulation times in both ionization modes for product-ion scanning (top panel). Lipid identification from *Chlamydomonas reinhardtii* using SWATH accumulation times of 10 ms for positive and 30 ms for negative ion mode (bottom panel). **(d)** Hierarchical clustering analysis for nine algal species and 1,023 binary variables. The top and bottom trees are from the classical taxonomies and chemotaxonomies, respectively. The yellow and blue colors between these trees show 'included' and 'not included' for each alga. UTEX 2341 is currently annotated as *Chlorella minutissima*. *C. reinhardtii* and *Dunaliella salina* are distinguished at the family level as Chlamydomonadaceae and Dunaliellaceae, respectively.



identification instead of quantitation, the MS-DIAL software also supports normalization methods required for specific needs in quantitative projects. In addition, MS-DIAL can be used with other DIA methods such as All-Ions MS/MS, MSc<sup>2</sup> and all-ion fragmentation<sup>20</sup>. The accuracy of deconvolution results, however, will depend on data acquisition scan speed and parameter settings such as scan width, overall sensitivity and data accumulation types. While SWATH data acquisition fits well with our deconvolution method, MS-DIAL also supports flexible precursor-mass windows from low to high *m/z*. Our algorithm could also benefit the proteomics community, where true mass spectral deconvolution has not been commonly used for peptide identification.

## METHODS

Methods and any associated references are available in the [online version of the paper](#).

*Note: Any Supplementary Information and Source Data files are available in the online version of the paper.*

## ACKNOWLEDGMENTS

This work was supported by the US National Science Foundation (NSF)–Japan Science and Technology Agency (JST) Strategic International Collaborative Research Program (SICORP) for Japan–United States metabolomics. We also thank the Lipid MAPS consortium for providing us the lipid SDF files; ChemAxon for a free research license for the Marvin and JChem cheminformatics tools; Z. Tietel and N. Nguyen (UC Davis) for assisting with the sample preparation of algal species; T. Bamba, Y. Izumi and T. Yamada (Osaka University) for suggestions and discussion of lipid annotation; D. Yukihiro (Kyushu University) for discussion of retention time prediction; and A. Ogiwara (Reyfigs Inc.) for development of the ABF file and for suggestions and discussion about MS-DIAL development. H.T. was also supported by Grant-in-Aid for Young Scientists (B) (Japan) 25871136. This study was also supported by the NSF (NSF MCB 113944), National Institutes of Health (NIH) (Grants P20 HL113452 and U24 DK097154), the JST-Core Research for Evolutionary Science and Technology (JST-CREST), and Database Integration Coordination Program by the National Bioscience Database Center (Japan).

## AUTHOR CONTRIBUTIONS

H.T., O.F. and M.A. designed the research. H.T. developed the MS-DIAL program. H.T. and T.C. analyzed the samples. T.C. and Y.M. contributed to the improvement of MS-DIAL program. H.T., T.K. and Y.M. performed the lipid annotations for the retention time prediction. H.T., T.K. and K.I. improved and optimized the LipidBlast library. H.T., B.H. and J.V. prepared the algal samples. M.K. developed the ABF file and the converter for this project. H.T., O.F. and M.A. thoroughly discussed this project and wrote the manuscript. T.C., T.K., B.H. and J.V. also contributed to the manuscript.

## COMPETING FINANCIAL INTERESTS

The authors declare competing financial interests: details are available in the [online version of the paper](#).

Reprints and permissions information is available online at <http://www.nature.com/reprints/index.html>.

- Zhu, X., Chen, Y. & Subramanian, R. *Anal. Chem.* **86**, 1202–1209 (2014).
- Röst, H.L. *et al. Nat. Biotechnol.* **32**, 219–223 (2014).
- Reiter, L. *et al. Nat. Methods* **8**, 430–435 (2011).
- Tsugawa, H. *et al. Anal. Chem.* **85**, 5191–5199 (2013).
- Nikolskiy, I. *et al. Anal. Chem.* **85**, 7713–7719 (2013).
- Stein, S.E. *J. Am. Soc. Mass Spectrom.* **10**, 770–781 (1999).
- Fiehn, O., Wohlgemuth, G. & Scholz, M. *Proc. Lect. Notes Bioinform* **3615**, 224–239 (2005).
- Tsugawa, H., Kanazawa, M., Ogiwara, A. & Arita, M. *Bioinformatics* **30**, 2379–2380 (2014).
- Horai, H. *et al. J. Mass Spectrom.* **45**, 703–714 (2010).
- Kind, T. *et al. Nat. Methods* **10**, 755–758 (2013).
- Wold, S., Sjostrom, M. & Eriksson, L. *Chemometr. Chemometr. Intell. Lab.* **58**, 109–130 (2001).
- Yap, C.W. *J. Comput. Chem.* **32**, 1466–1474 (2011).
- Gladu, P.K., Patterson, G.W., Wikfors, G.H. & Smith, B.C. *J. Phycol.* **31**, 774–777 (1995).
- Kind, T. *et al. J. Chromatogr. A* **1244**, 139–147 (2012).
- Haigh, W.G. *et al. Biochim. Biophys. Acta* **1299**, 183–190 (1996).
- Giroud, C., Gerber, A. & Eichenberger, W. *Plant Cell Physiol.* **29**, 587–595 (1988).
- Schymanski, E.L. & Neumann, S. *Metabolites* **3**, 412–439 (2013).
- Sumner, L.W. *et al. Metabolomics* **3**, 211–221 (2007).
- Creek, D.J. *et al. Metabolomics* **10**, 350–353 (2014).
- Egertson, J.D. *et al. Nat. Methods* **10**, 744–746 (2013).

## ONLINE METHODS

**Peak detection.** *Smoothing.* The peak detection algorithm starts with a smoothing method with respect to retention time and accurate mass. The MS-DIAL program utilizes the linearly weighted smoothing average<sup>21</sup>, which is simple and robust (**Supplementary Note** equation (1)). The software also supports several other smoothing methods: moving average<sup>21</sup>, Savitzky-Golay<sup>21</sup> and binomial filter<sup>22</sup>.

*Peak detection.* The basic concept of the peak detection algorithm consists of differential calculus and noise estimations (**Supplementary Fig. 7**). After the smoothing for retention time (against extracted ion chromatogram) or accurate mass (against mass chromatogram), peak detection is performed<sup>8</sup>. To evaluate noise, the program determines three threshold values automatically: (1) the maximum amplitude differences between two adjacent points, (2) the maxima of the first derivatives and (3) the maxima of the second derivatives in a chromatogram. The derivatives are calculated by five-point approximations (**Supplementary Note** equations (2) and (3)). From only the values below 5% of each maximum, medians of amplitude differences, first derivatives and second derivatives are computed as the threshold values for peak detection, and are hereafter called AF (amplitude filter), FF (first-order derivative filter) and SF (second-order derivative filter), respectively. When a computed median is near zero, 0.0001 is used.

The left edge of the peaks is recognized when the amplitude and the first-order derivative both exceed AF and FF in two adjacent points. In order to locate the edge more accurately, the local minimum of the adjacent 5-point window is explored by back-tracing from the detected start position. The peak top is recognized when the sign of the first-order derivative changes and the second-order derivative is less than SF. The right edge is recognized by the same criteria as the left.

**Peak spotting.** The term ‘peak spotting’ is derived from the visualization method in the MS-DIAL software and refers to peak detection based on retention time and MS1 data axes. The base peak chromatogram is formed for each mass slice of 0.1  $m/z$  with a step size of 0.05  $m/z$  (default), allowing all data points to belong to two adjacent slices (**Supplementary Fig. 8a**). Each data point of the base peak chromatogram has its scan number, retention time, base peak  $m/z$ , and base peak intensity. The peak detection algorithm as described above is applied to the base peak chromatogram and detected peak tops are shown as ‘spots’. Two spots of the same retention time and close  $m/z$  value in adjacent bins are merged by comparing their peak heights (**Supplementary Fig. 8b**). Although useful algorithms for automated noise and background reduction have been known<sup>23</sup>, we chose to exclude unwanted peaks simply by means of a user-defined exclusion mass list.

**Centroiding spectra.** When the profile mode data is analyzed in the MS-DIAL program, the spectral centroiding is performed (**Supplementary Note** equation (4)): after the same peak detection algorithm described above is performed, the ions in the user-defined region between the peak’s left and right edges are accumulated.

**MS<sup>2</sup>Dec deconvolution.** The MS<sup>2</sup>Dec procedure is applied to all spots that are detected in the peak spotting method. It consists of (1) centroiding of MS/MS peaks that are consistent to each spot, (2) extraction of the MS/MS chromatogram for each centroided spot, (3) smoothing and baseline correction, (4) model peak extraction and (5) model peak fitting for each MS/MS chromatogram by means of the least-squares method.

In contrast to the GC-MS approach<sup>6</sup>, we process high-resolution mass (instead of nominal mass) for the deconvolution. Specifically, we extract MS/MS chromatograms of the centroid MS/MS spectrum corresponding to the precursor ion detected in the spotting process, i.e., the precursor ion of the DIA MS/MS spectrum. For each spot, all MS/MS chromatograms within the following retention time (RT) range are extracted:

$$\text{RT range} \in (\text{peak top retention time} - 1.5 \times \text{peak width}, \text{peak top retention time} + 1.5 \times \text{peak width})$$

The peak width is the region between peak left and right edges of the focused peak spot.

After smoothing, each MS/MS chromatogram is subjected to the following baseline correction (**Supplementary Fig. 9**).

- Step 1. Each chromatogram is separated into a user-defined ‘segment’ value.
- Step 2. Local minimum within a user-defined ‘band width’ value is extracted and stored.
- Step 3. The median value of minimal points is computed for each segment, and points more than the median are discarded.
- Step 4. Baseline is determined by connecting the remaining minimal points.

Then, the peak detection algorithm is applied to each baseline-corrected MS/MS chromatogram. For each detected peak, two scores, the ‘ideal slope’ (**Supplementary Note** equations (5)–(7)) and ‘sharpness’ values (**Supplementary Note** equations (8)–(10)), are calculated. For these values please refer to the previous work<sup>24</sup>.

Examples of our least-squares method are shown in **Supplementary Data 4**. In our software program, an ideal slope score of >0.95 is necessary to be considered a ‘model peak’. For model peak candidates, their sharpness scores and the scan number are stored, and the second Gaussian derivative filter (**Supplementary Note** equation (11)) is fit to their array of sharpness values. Each local maximum of the Gaussian filter represents a model peak candidate. The data point region selected for each model peak candidate is described using a model peak chromatogram  $M(n)$  which has the baseline corrected chromatogram information from peak left to peak right edge. The least-squares method for the deconvolution is performed as follows:

$$C(n) = aM_1(n) + bM_2(n) + cM_3(n) + dn + e \quad (\text{eq. 1})$$

The original chromatogram  $C(n)$  is decomposed into three base vectors  $M_1(n)$ ,  $M_2(n)$  and  $M_3(n)$ . One vector,  $M_2(n)$ , corresponds to the model peak from the focused peak spot sided by two adjacent peaks on both sides,  $M_1(n)$  and  $M_3(n)$ . Note that the purpose of deconvolution is to determine  $M_k(n)$  ( $k = 1, 2, 3$ ) and coefficients  $a$ ,  $b$ ,  $c$ ,  $d$  and  $e$ . So far,  $M_k(n)$  is determined

as described above. Here, the coefficients are calculated as  $b = X^{-1}Y$  as shown below:

$$\begin{pmatrix} a \\ b \\ c \\ d \\ e \end{pmatrix} = \begin{pmatrix} \|M_1(n)\|^2 & (M_1(n), M_2(n)) & (M_1(n), M_3(n)) & (M_1(n), n) & (M_1(n), 1) \\ (M_2(n), M_1(n)) & \|M_1(n)\|^2 & (M_2(n), M_3(n)) & (M_2(n), n) & (M_2(n), 1) \\ (M_3(n), M_1(n)) & (M_3(n), M_2(n)) & \|M_3(n)\|^2 & (M_3(n), n) & (M_3(n), 1) \\ (n, M_1(n)) & (n, M_2(n)) & (n, M_3(n)) & \|n\|^2 & (1, n) \\ (1, M_1(n)) & (1, M_2(n)) & (1, M_3(n)) & (1, n) & \|1\|^2 \end{pmatrix}^{-1} \begin{pmatrix} (M_1(n), C(n)) \\ (M_2(n), C(n)) \\ (M_3(n), C(n)) \\ (n, C(n)) \\ (1, C(n)) \end{pmatrix} \quad (\text{eq. 2})$$

If the siding model peak  $M_1(n)$  or  $M_3(n)$  is not found within the extracted retention time region, equations (3), (4) and (5) will be changed to the corresponding matrix:  $X(3 \times 3)$ ,  $b(3 \times 1)$  and  $Y(3 \times 1)$  are used when both  $M_1(n)$  and  $M_3(n)$  are not found, and  $X(4 \times 4)$ ,  $b(4 \times 1)$  and  $Y(4 \times 1)$  are used when either  $M_1(n)$  or  $M_3(n)$  is not found. In the special case when  $M_2(n)$  is not found (this case would be possible when all MS/MS chromatograms are impure), an *ad hoc* model peak is inserted (instead of a *null* value) as follows. When the peak spotting algorithm is performed in RT vs. MS1 axis, model peaks of ideal slope value 1 are stored as peak candidates. Then, one model peak which has the median sharpness value within the candidates is used as the *ad hoc* model peak. Although inserting a *null* value is another option, we hypothesize that a compound was missed in the detected spot in RT vs. MS1.

**Compound identification.** The software program utilizes the NIST (US National Institute of Standards and Technology) MS format (NIST MSP ASCII) file for the reference library. Four criteria, (1) retention time, (2) accurate mass, (3) isotope ratio and (4) MS/MS spectrum information, are used for peak identification. Each score gives the standardized range from 0 to 1, meaning no similarity and a perfect match, respectively. The subscript 'act.' and 'lib.' of each equation describe the measurement value and the theoretical value, respectively.

*Accurate mass and retention time similarities.* These are calculated as follows:

Accurate mass (MS1) or RT similarities =

$$\exp\left\{-0.5 \times \left(\frac{\text{experimental value} - \text{theoretical value}}{\delta}\right)^2\right\} \quad (\text{eq. 3})$$

The background hypothesis of the equations for accurate mass and retention time similarities is that the differences between experimental and theoretical values follow the Gaussian distribution. The standard deviation  $\delta$  (user-defined) is also used as the search tolerance. If retention time information is not included in the MSP file of metabolites, the similarity value of retention time is not calculated.

*Isotope ratio.* If metabolite information in the MSP file includes the molecular formula, the theoretical isotopic distribution is

calculated from  $[M+0]$  to  $[M+5]$  by means of binomial and McLaurin expansion. An example for  $C_2H_6O$  is as follows:

$$\begin{aligned} & (^{12}C + ^{13}C)^2 (^1H + ^2H)^6 (^{16}O + ^{17}O + ^{18}O) = \\ & (^{12}C_2^1H_6^{16}O) \left(1 + \frac{^{13}C}{^{12}C}\right)^2 \left(1 + \frac{^2H}{^1H}\right)^6 \left(1 + \frac{^{17}O}{^{16}O} + \frac{^{18}O}{^{16}O}\right) \end{aligned} \quad (\text{eq. 4})$$

Here, the letters, such as  $^{12}C$ , show the natural abundance of each element. The contents except for the molecular mass  $[M+0]$  ( $^{12}C_2^1H_6^{16}O$ ) is expanded. Note that each coefficient value of expanded elements indicates the relative, i.e., isotope abundances with respect to the molecular ion ( $^{12}C_2^1H_6^{16}O$ ). Then, the relative abundances are compared between theoretical values and actual values. The intensity of  $[M+0]$  is normalized to 1. The similarity value of the isotope ratio is calculated as follows:

Isotope ratio similarity =

$$1 - \sum |r_{\text{act},i} - r_{\text{lib},i}| \quad \text{with } r_i = \frac{I_{M+i}}{I_M}, 1 \leq i \leq 5 \quad (\text{eq. 5})$$

The  $I_M$  and  $I_{M+i}$  show the intensity of the molecular ion and the isotope peak, respectively.

*Spectral similarity.* For the MS/MS spectral similarity, the MS-DIAL program utilizes the combined values of dot-product, reverse dot-product and the matched fragments ratio with the reference product ions. The weighting among the dot-product, reverse dot-product and the matched fragments ratio is 1:1:1 in the current MS-DIAL software setting. The amplitude of mass spectrum is normalized so that the highest amplitude of the product value becomes 1. The abundance  $A$  of each  $m/z$  is the integrated value within the user-defined  $MS^2$  tolerance. The dot product calculation in the MS-DIAL program are performed as follows:

$$\text{Dot product} = \frac{(\sum w A_{\text{act.}} w A_{\text{lib.}})^2}{\sum w A_{\text{act.}}^2 \sum w A_{\text{lib.}}^2} \quad \text{with } w = 1 / \left(1 + \frac{A}{\sum A - 0.5}\right) \quad (\text{eq. 6})$$

$\sum A$  is the sum of relative abundances. The reverse dot product is also calculated in the same way (**Supplementary Note** equation (12)).

The coefficient  $w$  is the weight value in order to reduce the effect of high abundance intensities. In the dot-product calculation, the half abundance of the measured spectrum is used if the corresponding mass peak does not exist in a library spectrum. Unwanted peaks derived from isotopic ions or background noise may decrease the dot-product score. In the reverse dot-product, the spectrum in the reference library is used to calculate the score. The ion abundance of the reference spectrum is halved when the pairing mass peak does not exist in the query spectrum.

**Total similarity.** The four scores are used for compound identification.

$$\text{Total score} = \frac{\text{MS/MS similarity} + \text{MS1 similarity} + \text{RT similarity} + 0.5 \times \text{isotope ratio similarity}}{3.5} \times 100 \quad (\text{eq. 7})$$

If retention time information is not available for compound identification, the total score is calculated as **Supplementary Note** equation (13). If the formula information is not available, the total score is calculated as **Supplementary Note** equation (14). If both retention time and formula information are not available, the total score is calculated as **Supplementary Note** equation (15). The compound with highest total score (above the user-defined threshold) is assigned to each focused peak. When the MS/MS spectrum is not obtained for data-dependent MS/MS acquisition, the MS/MS similarity is recognized as zero and the denominator described above is decremented by 1.

#### Peak alignment, filtering and missing value interpolation.

The algorithm of peak alignment in MS-DIAL is based on the idea of Joint Aligner implemented in MZmine<sup>25</sup>. It consists of four major steps: (1) making a reference table, (2) fitting each sample peak table to the reference peak table, (3) filtering aligned peaks and (4) interpolating missing values. The summary of MS-DIAL peak alignment algorithm is described in **Supplementary Figure 10**.

**Making a reference peak table.** As shown in **Supplementary Figure 10a**, the reference peak table consisting of retention time (RT) and  $m/z$  is created as follows:

- Step 1. A user-defined 'reference file' which is one of the aligned samples is used as the basis of the reference peak table.
- Step 2. Information of each sample peak table is inserted to the reference peak table (**Supplementary Fig. 10b**). The condition is as follows:

$$\text{If } |RT_{\text{sam.}} - RT_{\text{ref.}}| > \delta_{\text{RT}} \cup |Mass_{\text{sam.}} - Mass_{\text{ref.}}| > \delta_{\text{Mass}} \\ \text{then insert to peak table} \quad (\text{eq. 8})$$

where  $\delta_{\text{RT}}$  and  $\delta_{\text{Mass}}$  are user-defined tolerance values for RT ( $\delta_{\text{RT}}$ ) and MS1 accurate mass ( $\delta_{\text{Mass}}$ ), respectively.

- Step 3. Repeat the function for all peaks from all samples.

The reference peak table is used in order to associate each peak in each sample.

**Fitting each sample peak table to the reference peak table.** Each peak in the sample data is associated with the reference peak list using the following criterion:

$$\text{Score} = a \times \exp \left\{ -0.5 \times \left( \frac{RT_{\text{sam.}} - RT_{\text{ref.}}}{\delta_{\text{RT}}} \right)^2 \right\} \\ + b \times \exp \left\{ -0.5 \times \left( \frac{Mass_{\text{sam.}} - Mass_{\text{ref.}}}{\delta_{\text{Mass}}} \right)^2 \right\} \quad (\text{eq. 9})$$

The coefficient is user-defined RT factor ( $a$ ) and MS1 accurate mass factor ( $b$ ), respectively.  $\delta_{\text{RT}}$  and  $\delta_{\text{Mass}}$  are the same as the

above criteria to construct the reference peak table. Finally, aligned peak table including alignment ID, average RT, average  $m/z$  and intensities of all samples is generated.

**Filtering aligned peaks.** MS-DIAL provides the simple filter in order to exclude unwanted alignment ID (**Supplementary Fig. 10c**). Three-step filtering is applied for each alignment ID.

- Step 1. If all peak intensities of samples in a row are missing or undetected, the alignment information is removed.
- Step 2. If the percentage of filled peaks in an alignment ID is less than the user-defined peak count filter (default 0%), the information is removed.
- Step 3. This is optional, but if all quality control (QC) samples are not filled in an alignment ID, the information is removed.

**Interpolating missing values.** In each alignment ID, the intensity information of all samples is not always filled. As shown in **Supplementary Figure 10d**, such missing values after the above process are interpolated in MS-DIAL as follows:

- Step 1. The average retention time and average  $m/z$  of 'filled' peaks are calculated.
- Step 2. A local maximum from the following range is stored for the missing value.

$$\left( RT_{\text{average}} - \delta_{\text{RT}}, RT_{\text{average}} + \delta_{\text{RT}} \right) \cap \\ \left( Mass_{\text{average}} - \delta_{\text{Mass}}, Mass_{\text{average}} + \delta_{\text{Mass}} \right) \quad (\text{eq. 10})$$

**Databases.** The MassBank revision 173, ReSpect updated in 2012/9/25 and LipidBlast version 3 were downloaded. The spectrum data were converted to the NIST MSP format. For the hydrophilic metabolite identification, the NIST 12 MS/MS library was also used in addition to MassBank and ReSpect libraries. For the algal lipid identification, fatty acid 16:2, 16:3, 16:4 and 16:5 spectra information were added to the LipidBlast library. The position of double bonds was determined according to previous reports<sup>16</sup>. Moreover, the adduct ions and the MS/MS spectral information of formic acid were added to phosphatidylcholine (PC), lysoPC, MGDG and DGDG for the lipid identification in



negative ion mode analysis. In order to determine the ion abundances for each lipid class the heuristic model was constructed from the data sets of DDA MS/MS. The MSP format libraries (MassBank, ReSpect and LipidBlast) and the LipidBlast excel macro file are downloadable under <http://prime.psc.riken.jp/>.

**Retention time prediction for lipids.** The SDF files of all lipids in LipidBlast were constructed as follows. The SDF files of PC, lysoPC, PE, lysoPE, PG, PI, PS and PA were downloaded from LIPID MAPS<sup>26</sup>. The SDF files for the other lipid classes were created from SMILES code written in LipidBlast by ChemAxon JChem 6.3.0 molconvert (<http://www.chemaxon.com>), totaling 117,343 SDF files. They also included plasmalogen PC, PE, sphingomyelin and cholesterol ester as lipid classes, although these lipids were not the focus for algal lipid identification. The PaDEL descriptor software was used to calculate 1D and 2D molecular descriptors and PubChem fingerprints from the SDF files<sup>12</sup>. Their exact masses were also generated by ChemAxon JChem calculator. Then, redundant and uniform variables were excluded, and a total of 464 compound descriptors were used as predictor variables in the regression analysis. The in-house retention time information of 254 lipids was used for model development. Since the number of predictor variables (compound descriptors) were considerably higher than the number of data samples (the number in the training set: 254), partial least-squares regression (PLS-R) was used in order to construct the retention time prediction model<sup>11</sup>. The program of PLS-R was written in Visual Basic for Application and the source code can be downloaded at <http://prime.psc.riken.jp/>. A sevenfold cross validation was used to calculate the predictive residual sum of squares (PRESS) and  $Q^2$  value. The final model included six latent variables based on the PRESS and  $Q^2$  value and the retention time information from the training samples. In this study, retention time information of newly identified 1,808 lipids from nine algal species was used for validating that accurate precursor ion masses and MS/MS spectra were also confirmed by retention time matching.

**MS-DIAL software and data processing parameters.** MS-DIAL is available in Windows OS (.NET Framework 4.0 or later; RAM: 4.0 GB or more). Its source code was written in the C# language with the Windows Presentation Foundation (WPF) to develop the graphical user interface. The main source code such as peak detection, peak spotting and MS<sup>2</sup>Dec algorithm is downloadable at <http://prime.psc.riken.jp/>. The data processing parameter of MS-DIAL used in this study are described in **Supplementary Table 4**.

**Chemotaxonomic tree by lipid descriptors and hierarchical clustering analysis.** All lipids were annotated with subsequent manual verification of MS/MS spectral matching for compound identification (**Supplementary Table 1**). A total of 1,808 (SWATH) and 1,521 (DDA) lipids (**Supplementary Table 2**) were first integrated disregarding the acyl chain positions (*sn1*, *sn2*, *sn3*), double bond positions and stereoisomers (*E*, *Z*). For example, TAG(16:0/16:1/16:2), TAG(16:0/16:2/16:1), TAG(16:2/16:1/16:0), TAG(16:2/16:0/16:1), TAG(16:1/16:0/16:2) and TAG(16:1/16:2/16:0) were considered the same lipid. Likewise, lysoPC 16:1(7*Z*) and lysoPC 16:1(9*E*) were regarded as the same. For the remaining 1,023 lipids, presence or absence in

each of nine species was represented as a binary data matrix of size 1,023 × 9 (**Supplementary Data 2**).

Hierarchical clustering analysis was performed using the R statistical language (<http://www.R-project.org>) and the package 'amap' (<http://CRAN.R-project.org/package=amap>). The distance was calculated by 'correlation' in the package. The linkage was performed by 'average'. We cited the previous report<sup>27</sup> as the standard taxonomic tree.

**Biospecimen and algae strains.** A single human plasma sample was obtained from the Cleveland Clinic from the GeneBank study<sup>28</sup>. The cultivation procedure of *Chlamydomonas reinhardtii* followed our previous report<sup>29</sup>. The *C. reinhardtii* CC125 strain was streaked out from cryopreserved stock and cultivated in 75 mL TAP medium in 125 mL shake flasks at 25 °C under constant illumination with cool-white fluorescent bulbs at a fluence rate of 70 μmol m<sup>-2</sup> s<sup>-1</sup> and with continuous stirring (100 rpm). Four independent cultures were used for this study. The starter culture was harvested at late log phase and 1 mL cell suspensions were then shifted to 75 mL of fresh TAP medium in 125 mL shake flasks. At 0.2–0.6 OD<sub>680</sub> during the late-log phase, 1 mL cell suspensions were injected into 1 mL of –80 °C cold quenching solution composed of 70% methanol in water, centrifuged at 12,000 g for 2 min, and pellets were lyophilized and stored at –80 °C until further analysis. The same quenching procedure was used for all algae strains.

UTEX 2341 (originally classified as *Chlorella minutissima*), *Chlorella sorokiniana* (UTEX 2805), and *Chlorella variabilis* (ATCC NC64A) were plated on ATCC #5 agar and colonies were selected for inoculation into liquid cultures. All three *Chlorella* strains were cultivated simultaneously in 250 mL hybridization tubes with four independent cultures per strain. Hybridization tubes were filled with 200 mL media and maintained in a 28 °C water bath. Aeration was supplied at 125 mL per minute with 2% CO<sub>2</sub> mixed with air (v/v). Reactors were illuminated horizontally (10,000 lx) by T5 growth lamps operating on a 16:8 light/dark cycle and cultures were mixed by stir bar operating at ~150 rpm. UTEX 2341 was cultivated in N8-NH<sub>4</sub> medium<sup>30</sup>, *C. sorokiniana* in N8 medium<sup>31</sup> and *C. variabilis* in N8-NH<sub>4</sub> medium supplemented with 20 mg/L yeast extract. Culture samples (1 mL) were quenched for lipidomics analysis during the late log growth stage.

The cultures of *Euglena gracilis* (UTEX B367), *Cricosphaera carterae* (UTEX LB1014), *Nannochloropsis oculata* (UTEX LB2164), *Dunaliella salina* (UTEX LB200) and *Pavlova lutheri* (UTEX LB1293) were purchased from the UTEX culture collection of algae<sup>32</sup>. Three technical replicates for each strain were prepared from quenched samples.

**Reagent and sample preparation.** Water, isopropanol and acetonitrile were purchased from Fisher Optima. Methanol was purchased from J.T. Baker. Ammonium formate, formic acid and methyl *tert*-butyl ester (MTBE) were purchased from Sigma-Aldrich. Authentic standard compounds were purchased from Avanti Polar Lipids Inc., CDN Isotopes, Cayman Chemical and Sigma-Aldrich.

For hydrophilic interaction chromatography (HILIC)-MS/MS analysis of pharmaceutical agents present in a human plasma sample, all procedures for the metabolite extraction were kept on ice. 30 μL of human plasma was added to 1,000 μL cold mix-solvent (acetonitrile/isopropanol/water, 3:3:2, v/v/v) on ice, then vortexed for 10 s and shaken for 5 min at 4 °C using the



Orbital Mixing Chilling/Heating Plate (Torrey Pines Scientific Instruments). After 2 min centrifugation at 14,000 rcf, 300  $\mu$ L of the supernatant was transferred to a new 1.5 mL Eppendorf tube and evaporated to dryness in a Labconco Centrивap cold trap concentrator. The dried sample was resuspended with 60  $\mu$ L (80% acetonitrile in water) including 0.038  $\mu$ g/mL choline-D<sub>9</sub>, 0.050  $\mu$ g/mL TMAO-D<sub>9</sub>, 0.020  $\mu$ g/mL betaine-D<sub>9</sub>, 10.0  $\mu$ g/mL glutamine-D<sub>5</sub>, and 1.48  $\mu$ g/mL arginine-<sup>15</sup>N<sub>2</sub> and centrifuged for 5 min at 16,000 rcf. The 50  $\mu$ L aliquot was transferred to a glass amber vial (National Scientific) with a micro-insert (Supelco).

For lipid profiling, all samples for the metabolite extraction were kept on ice and performed as described previously<sup>33</sup>. 225  $\mu$ L of MeOH including 1.64  $\mu$ g/mL PE (17:0/17:0), 6.55  $\mu$ g/mL PG (17:0/17:0), 1.10  $\mu$ g/mL PC (17:0/0:0), 0.24  $\mu$ g/mL sphingosine (d17:1), 0.55  $\mu$ g/mL ceramide (d18:1/17:0), 0.44  $\mu$ g/mL SM (d18:1/17:0), 54.5  $\mu$ g/mL palmitic acid-D<sub>3</sub>, 0.44  $\mu$ g/mL PC (12:0/13:0), 22.7  $\mu$ g/mL cholesterol-D<sub>7</sub>, 0.27  $\mu$ g/mL TAG (17:0/17:1/17:0), 2.18  $\mu$ g/mL DAG (12:0/12:0/0:0), 13.1  $\mu$ g/mL DAG (18:1/2:0/0:0), 4.36  $\mu$ g/mL MAG (17:0/0:0/0:0) and 0.55  $\mu$ g/mL PE (17:1/0:0) were added to each dried algae on ice and vortexed for 10 s. Then, the MTBE including 21.8  $\mu$ g/mL cholesteryl ester (22:1) was added on ice and vortexed for 10 s. After shaking for 6 min at 4 °C in the orbital mixer, 188  $\mu$ L water was added and vortexed for 20 s. After centrifugation for 2 min at 14,000 rcf, 350  $\mu$ L of the supernatant was transferred to a new 1.5 mL Eppendorf tube and evaporated to dryness in the Labconco Centrивap cold trap concentrator. The dried sample was resuspended in 108.6  $\mu$ L MeOH:toluene 90:10 (v/v) with CUDA (12-[[[(cyclohexylamino)carbonyl]amino]-dodecanoic acid, 50 ng/mL). After vortexing for 20 s, each sample was sonicated for 5 min at room temperature. After centrifugation for 2 min at 16,000 rcf, 50  $\mu$ L of the supernatant was transferred to a glass amber vial with micro-insert. The *C. reinhardtii*, *C. sorokiniana* and *C. variabilis* samples were diluted by adding 50  $\mu$ L of MeOH:toluene 90:10 (v/v). Moreover, the *E. gracilis* sample was diluted by adding 200  $\mu$ L of MeOH:toluene 90:10 (v/v).

**Analytical conditions.** The liquid chromatography system consisted of an Agilent 1290 system (Agilent Technologies Inc.) with a pump (G4220A), a column oven (G1316C) and an autosampler (G4226A). For hydrophilic metabolite analysis, mobile phase A was 10 mM ammonium formate with 0.125% formic acid in water; mobile phase B was 95:5 acetonitrile:water (v/v) with 10 mM ammonium formate and 0.125% formic acid. An Acquity UPLC BEH Amide column (150  $\times$  2.1 mm; 1.7  $\mu$ m) coupled to a VanGuard BEH Amide pre-column (5  $\times$  2.1 mm; 1.7  $\mu$ m) (Waters; Milford, MA, USA) was used. The gradient was 0 min, 100% B; 2 min, 100% B; 7.7 min, 70% B; 9.5 min, 40% B; 10.3 min, 30% B; 12.8 min, 100% B; 16.8 min, 100% B. The column flow rate was 0.4 mL/min, autosampler temperature was 4 °C, injection volume was 2  $\mu$ L and column temperature was 45 °C. For lipid analysis, mobile phase A was 60:40 acetonitrile:water (v/v) with 10 mM ammonium formate and 0.1% formic acid; mobile phase B was 90:10 isopropanol:acetonitrile (v/v) with 10 mM ammonium formate and 0.1% formic acid.

The lipidomic LC method used an Acquity UPLC charged-surface hybrid (CSH) C18 column (100  $\times$  2.1 mm; 1.7  $\mu$ m) coupled to an Acquity CSH C18 VanGuard pre-column (5  $\times$  2.1 mm; 1.7  $\mu$ m) (Waters; Milford, MA, USA). The gradient was 0 min, 15% B; 2 min, 30% B; 2.5 min, 48% B; 11 min, 82% B,

11.5 min, 99% B; 12 min, 99% B; 12.1 min, 15% B; 15 min, 15% B. The column flow rate was 0.6 mL/min, autosampler temperature was 4 °C, injection volume was 3  $\mu$ L in positive mode and 5  $\mu$ L in negative mode, and column temperature was 65 °C.

Mass spectrometry was performed on an AB Sciex TripleTOF 5600+ system (Q-TOF) equipped with a DuoSpray ion source. All analyses were performed at the high sensitivity mode for both TOF MS and product ion scan. The mass calibration was automatically performed every 10 injections using an APCI positive/negative calibration solution via a calibration delivery system (CDS). For HILIC analysis, SWATH (sequential window acquisition of all theoretical mass spectra) acquisition with positive ion mode was used as the data independent acquisition system. The SWATH parameters were MS1 accumulation time, 50 ms; MS2 accumulation time, 30 ms; collision energy, 45 V; collision energy spread, 15 V; cycle time, 640 ms; Q1 window, 25 Da; mass range,  $m/z$  50–500. The other parameters were curtain gas, 35; ion source gas 1, 50; ion source gas 2, 50; temperature, 300 °C; ion spray voltage floating, 4.5 kV; declustering potential, 100 V; RF transmission,  $m/z$  40: 33%,  $m/z$  120: 33% and  $m/z$  390: 34%. For lipid analysis, six different methods were used; DDA (data-dependent acquisition) with positive ion mode, DDA with negative ion mode, SWATH acquisition (Q1 window, 21 Da) with positive ion mode, SWATH acquisition (Q1 window, 21 Da) with negative ion mode, SWATH acquisition (Q1 window, 65 Da) with positive ion mode and SWATH acquisition (Q1 window, 65 Da). The common parameters in both SWATH/DDA and positive/negative ion mode were collision energy, 45 V; collision energy spread, 15 V; mass range,  $m/z$  100–1,250; curtain gas, 35; ion source gas 1, 60; ion source gas 2, 60; temperature, 350 °C; declustering potential, 80 V; RF transmission,  $m/z$  80: 50%,  $m/z$  200: 50%. The ion spray voltage floating of positive/negative ion mode were +5.5/–4.5 kV, respectively. The DDA parameters in both positive and negative ion modes were MS1 accumulation time, 100 ms; MS2 accumulation time, 50 ms; cycle time, 650 ms; dependent product ion scan number, 10; intensity threshold, 500; exclusion time of precursor ion, 5 s; mass tolerance, 20 mDa; ignore peaks, within 6 Da; dynamic background subtraction, TRUE. The SWATH parameters of 21/65 Da Q1 window were MS1 accumulation time, 100/50 ms; MS2 accumulation time, 10/30 ms; cycle time, 731/640 ms; Q1 window, 21/65 Da.

21. Savitzky, A. & Golay, M.J.E. *Anal. Chem.* **36**, 1627–1639 (1964).
22. Lommen, A. *Anal. Chem.* **81**, 3079–3086 (2009).
23. Windig, W., Phalp, J.M. & Payne, A.W. *Anal. Chem.* **68**, 3602–3606 (1996).
24. Hiller, K. et al. *Anal. Chem.* **81**, 3429–3439 (2009).
25. Katajamaa, M., Miettinen, J. & Oresic, M. *Bioinformatics* **22**, 634–636 (2006).
26. Sud, M. et al. *Nucleic Acids Res.* **35**, D527–D532 (2006).
27. Cavalier-Smith, T. *Biol. Rev. Camb. Philos. Soc.* **73**, 203–266 (1998).
28. Wang, Z. et al. *Nature* **472**, 57–63 (2011).
29. Lee, Y., Park, J.-J., Barupal, D.K. & Fiehn, O. *Mol. Cell. Proteomics* **11**, 973–988 (2012).
30. Higgins, B.T. & VanderGheynst, J. *PLoS ONE* **9**, e96807 (2014).
31. Tanadul, O.U., VanderGheynst, J.S., Beckles, D.M., Powell, A.L. & Labavitch, J.M. *Biotechnol. Bioeng.* **111**, 1323–1331 (2014).
32. Brand, J.J., Andersen, R.A. & Nobles, D.R. Jr. in *Applied Phycology and Biotechnology* Second Edition, John Wiley & Sons, Ltd, Oxford, UK doi:10.1002/9781118567166.ch5 (2013).
33. Matyash, V., Liebisch, G., Kurzchalia, T.V., Shevchenko, A. & Schwudke, D. *J. Lipid Res.* **49**, 1137–1146 (2008).