1  2 # Lipidome visualisation, comparison, and analysis in a vector space

3

4  Short title: Lipidome Analysis

5

6  Timur Olzhabaev[1,2], Lukas Müller[1,2], Daniel Krause[2], Dominik Schwudke[2,3,4*], Andrew

7  Ernest Torda[1*]

8

9  [1] Centre for Bioinformatics, University of Hamburg, Hamburg, Germany

10  [2] Bioanalytical Chemistry, Research Center Borstel Leibniz Lung Center, Borstel,
11  23845, Germany

12  [3] German Center for Infection Research, Thematic Translational Unit Tuberculosis,
13  23845 Borstel, Germany

14  [4] German Center for Lung Research (DZL), Airway Research Center North (ARCN),
15  23845 Borstel, Germany

16

17  * Corresponding authors

18  E-mail: andrew.torda@uni-hamburg.de

19  For submission to PLoS Computational Biology

## Abstract

A shallow neural network was used to embed lipid structures in a 2- or 3-dimensional space with the goal that structurally similar species have similar vectors. Tests on complete lipid databanks show that the method automatically produces distributions which follow conventional lipid classifications. The embedding is accompanied by the web-based software, Lipidome Projector. This displays user lipidomes as 2D or 3D scatterplots for quick exploratory analysis, quantitative comparison and interpretation at a structural level.

## Author summary

Lipids are not just the basis of membranes. They carry signals and metabolic energy. This means that the presence, absence, and quantity of lipids reflects a cell's biochemical state - starving, nourished, sick or healthy. Lipidomics (measuring all lipids in a biological specimen) provides lists of the chemical species and their quantities.

We have used a shallow neural network from natural language modelling to embed lipids in a continuous vector space. Firstly, this means that similar molecules have similar positions in this space. Conventional lipid categories cluster automatically. Secondly, the accompanying web-based software, Lipidome Projector imports a lipidome and displays it as a set of points. Reading several lipidomes at once allows quantitative and structural comparisons. Combined with the ability to show structure and abundance diagrams, the software allows exploratory analysis and interpretation of lipidomics datasets.

2

# Introduction

42

43 Lipids remind one of membranes or fats, but they also carry energy and signals, so

44 one may assume that the set of lipids in a sample reflects the health and metabolic

45 state of a tissue or organism. Mass spectrometry provides lipidome information, but a

46 list of $10^2$-$10^4$ lipids and quantities is not easily interpretable. For exploratory analysis,

47 one would like a method that highlights chemical trends and shows how samples

48 differ with respect to lipid structures and quantities. Given a set of mass spectrometry

49 peaks that have been assigned to lipids, the idea is to display lipidomes as

50 scatterplots in a 2- or 3-dimensional space. This requires two steps. First, there must

51 be a continuous vector space such that each lipid gets distinct coordinates. Second,

52 one needs software to display and compare plots interactively. The software should

53 make it easy to relate points back to their names and chemical structures.

54 The aims here are different to those of other lipidomics software packages. If one

55 wants to treat a lipidome similarly to gene expression data, one can look for changed

56 levels of lipids or focus on molecules whose abundances are correlated [1–3]. If one

57 wants to see a lipidome in terms of networks, there is network construction and

58 display software [4]. Our focus is different. Lipidome Projector lets one quickly

59 highlight and interactively explore differences between groups of samples, with the

60 simultaneous display of abundances and structures.

61 The first challenge is finding vectors for molecules for the two- and three-dimensional

62 plots. Previous attempts applied ideas from string comparisons [5], but this was not

63 without problems. Whatever notation one uses, a small change to a molecule can

64 lead to a large change in a string representation such as SMILES [6], so the similarity

65 metrics are fundamentally unstable. Kopczynski et al approached the problem with

66  elegant distance metrics, but this required some preconceptions about lipid structures

67  and used expensive graph similarity methods [7].

68  We come to the problem with slightly different ideas and some specific goals. The

69  method should be objective, unsupervised and require minimal chemical

70  preconceptions. Coordinates should be quite different for unrelated molecules, but

71  systematic changes such as extending the length of an aliphatic chain should give a

72  series of points near each other. Adding a phosphate or alcohol group to two different

73  molecules should change both coordinates in a similar manner. Our method for lipids

74  is a modified version of Mol2Vec [8], a technique from the small-molecule literature

75  which is, in turn, based on Word2Vec [9] a word embedding method from natural

76  language processing. To embed words, one first defines a vocabulary and gives

77  each word a unique token. In a text corpus, similar tokens appear in similar contexts

78  with reasonable probability, such that a token / context prediction task can be used to

79  train semantic vector representations. To apply the idea in chemistry, one constructs

80  a vocabulary of chemical fragments and trains a shallow network on a large set of

81  molecules to recognise surrounding contexts. Input fragments are represented by

82  integer identifiers derived from computed sparse connectivity fingerprints [10].

83  Fragment vectors come from hidden layer weights of the trained network and are

84  summed to produce vector representations of entire molecules.

85  Calculating the vector space model is performed once on a large set of lipid

86  structures and takes several hours. User lipidome data is simply matched to

87  precomputed vectors. Lipidome Projector, the browser-based application for

88  visualization and analysis, allows one to interactively explore lipidomes in the vector

89  space and additionally displays lipid abundance charts and molecular structures.

90   To judge our methods, we consider the distributions of lipids in the computed vector

91   space and apply Lipidome Projector visualizations on three published lipidome

92   datasets.

# Materials and Methods

## Lipid Vector Space

95   For training, the Lipid Maps Structure Database (LMSD) [11] and SwissLipids [12]

96   (both accessed Jan 2023) were combined. SwissLipids entries were filtered to obtain

97   lipids with valid SMILES at isomeric subspecies level. The combination of databases

98   resulted in over 620 000 unique structures. RDKit [13] was used to convert all

99   database entries to a consistent charge state and RDKit's implementation of

100  extended connectivity fingerprints [10] was used to assign a unique identifier to each

101  substructure of a specified radius around each atom. Substructure identifiers were

102  ordered according to the position of the substructure's central atom within the

103  molecule's canonical SMILES string.

104  A few small modifications to Mol2Vec were necessary. First, chirality was explicitly

105  considered. Secondly, a parameter had to be adapted to capture differences in long

106  alkyl chains. Mol2Vec descriptors for small molecules are usually built from

107  fragments using atoms (radius 0) and their immediate neighbours (radius 1). For the

108  much larger lipid structures, radii of size 0, 1, 2, 3, 4, 5, 10, 15, 20, 25, 30, 35, 40, 45

109  and 50 were used, resulting in just under three million unique fragments for the

110  combination of databases. For each lipid, the set of fragments for each radius was

111  used as a separate training sentence.

112  Gensim [14] was used to train the Word2Vec model with training parameters listed in

113  Table S1. The network generated 100-dimensional substructure vectors, which were

114  summed for each molecule. For visualization, the Barnes-Hut [15] version of t-

115 distributed stochastic neighbour embedding [16] as implemented in OpenTSNE [17]

116 was used to reduce the 100-dimensional vector space to generate 2- and 3-

117 dimensional vector sets (parameters listed in Table S2). The embedding process is

118 summarised in Fig 1 A.

119 **Fig 1. Vector Space Generation and Matching.** (A) A lipid structure is decomposed

120 into its substructures of different sizes represented by Morgan sparse fingerprint

121 integers, which constitute the training data for Word2Vec. A molecule's vector is the

122 sum of its substructure vectors and is projected to 2D or 3D with stochastic neighbour

123 embedding. (B) The user provides a list of lipid species names and component

124 constraints. Lipid names are parsed and matched to appropriate isomer names from

125 the pre-parsed database. The component constraints are applied to filter the

126 matches. Vectors of the remaining isomers are averaged for each lipid. Not illustrated

127 is an additional step, in which database matching is attempted on the original names

128 of unparsed lipid species.

## Lipidome Processing

130 As part of building the system, entries from the lipid databases are stored along with

131 their corresponding vectors and higher-level abbreviations for each isomer following

132 previously defined levels [18]. When a user lipidome is imported, entries are matched

133 against pre-calculated vectors (Fig 1 B). Goslin [19] is used to parse both databases

134 and user data. It accepts common nomenclature, but should it fail, the process will

135 look for a match based on user-provided names. This means that Lipidome Projector

136 covers at least all entries from the union of SwissLipids and the LMSD that were

137 successfully parsed by Goslin (S1 Dataset gives a list of translated class names).

138 Mass spectrometry often does not identify a lipid at the complete structure level [18]

139 so additional steps are necessary to deal with this ambiguity. The software finds the

6

140    set of isomers that match the higher-level abbreviation, but not all members of this

141    set will be plausible for the organism under consideration. To filter the list of possible

142    lipids, Lipidome Projector expects a constraints list with allowed fatty acyls and long-

143    chain bases. The remaining isomer vectors are averaged to produce a single

144    representative vector.

## Visualization and Analysis Software

146    Plots are generated using Plotly.py [20]. Marker sizes are derived from respective

147    lipid abundances, to which either linear or min-max scaling is applied. Dash [20] is

148    used to build the web-application front end. The rest of the application was built in

149    Python [21] with pandas [22] used for data-table storage and manipulation. Parsing

150    and matching are performed server-side. The original lipidome dataset together with

151    the newly derived lipid names and computed vectors is stored inside the user's

152    browser session and sent to the server for temporary processing operations such as

153    averaging of samples or plot updates. Lipidome datasets and constraints are read in

154    a simple table format.

## Datasets

156    Publicly available lipidome datasets from drosophila [23], yeast [24] and mouse [25]

157    were used for development and analysed as user cases. Python scripts for the

158    extraction of the original data and formatting into formats appropriate for Lipidome

159    Projector, as well as manually constructed respective FA and LCB constraint files are

160    given in S2 dataset.

## Results

## Lipid Vector Space

163    We first consider the projection of lipids into a vector space by looking at the

164    distributions of points for entries from the combined databases with a valid structure

7

165   and class. Are the vectors consistent with chemical intuition and database

166   classification? Fig 2 A shows the entire lipid set in two dimensions (see S2 Fig for 3D

167   version). With some exceptions, lipids within a category are grouped together in the

168   vector space despite the underlying structural diversity. For the largest categories,

169   glycerolipids (GL), glycerophospholipids (GP) and sphingolipids (SP) a clear

170   separation can be observed with some overlap and outliers at some borders. To look

171   in more detail, one can focus on the class level with the example of selected

172   glycerophospholipid classes. Fig 2 B marks three clusters, which largely correspond

173   to diacyl, mono-alkyl and plasmalogen glycerophospholipids respectively. This

174   suggests that the embedding has mostly captured the chemical connectivity at the

175   glycerol. Within each large cluster, phosphatidylinositols (PI) and

176   phosphatidylcholines (PC) form their own subgroups with some local exceptions. For

177   the other classes there are numerous smaller, intertwined clusters spread across the

178   vector space. Also marked are a few unusual molecules with uncommon fatty acyl

179   double bond structures such as (5E, 9E) or chains which are heavily methylated or

180   even contain ladderane, a structural moiety seen in bacteria. These are positioned

181   outside the main group as one might expect since the database is dominated by the

182   biochemistry of mammals. The remaining plots in Fig 2 show how the lipid vectors

183   capture chemical functional groups and their structural context. In Fig 2 C there is a

184   general trend of more double bonds from left to right. Focusing on a local region

185   shows that clustering is determined by lipid class (Fig 2 D) and fatty acyl double bond

186   location and number (Fig 2 E). Additionally, one can see a systematic change in

187   mass as one moves along clusters (Fig 2 F). These patterns suggest that the

188   embedding captures gradual structural changes. This was further assessed using a

189   contrived example borrowed from the literature [5]. Three sets of manually generated

190   structures were added to the training data. The first two consist of series of

8

191   phosphatidylinositols with a successively longer fatty acyl chain. The sets are the

192   same, except for the presence / absence of a double bond in the lengthening chain.

193   Fig 3 B shows that growing an aliphatic chain gives progressively changing vector

194   positions, while the presence of the double bond leads to a large, but consistent

195   displacement. The third set consists of a series of ceramides, each of which is

196   hydroxylated at a different position within its fatty acyl chain (Fig 3 A). The steps of

197   the hydroxylated position translate into an almost linear series of vectors with the

198   exception of an outlier near the acyl bond.

199   **Fig 2. Vector Space (2D).** (A) Entire vector space. Marker colour represents lipid

200   category: Fatty acids (FA), glycerolipids (GL), glycerophospholipids (GP),

201   sphingolipids (SP), sterol lipids (ST), prenol lipids (PR), saccharolipids (SL) and

202   polyketides (PK). (B) Region of the vector space focused on selected

203   glycerophospholipids: Glycerophosphates (PA), glycerophosphocholines (PC),

204   glycerophosphoethanolamines (PE), glycerophosphoglycerols (PG),

205   glycerophosphoinositols (PI) and glycerophosphoserines (PS). Marker colour: Lipid

206   class. (C) Same region as in B, marker colour represents the number of fatty acyl

207   double bonds. (D) Zoomed-in region of selected glycerophospholipids, marker colour

208   represents lipid class. (E) Same region as in D, marker colour represents the double

209   bond profile of the 2-sn fatty acyl. (F) Same region as in D, marker colour represents

210   molecule mass.

211   **Fig 3. Impact of Stepwise Structural Changes.** (A) Local vector space region of

212   manually added ceramide structures. Marker annotations denote the fatty acyl

213   hydroxylation position. (B) Local vector space region of manually added

214   phosphatidylinositol structures. Marker annotations denote the length of the 2-sn fatty

215   acyl.

9

216   Another aspect of the quality of the vector space is its coverage of lipid classes, fatty

217   acyls, and long-chain bases, which in our case, is completely dependent on the

218   underlying databases and the parser. When lipidomes are imported, entries are

219   discarded if they cannot be matched or if they are rejected by the constraint-based

220   filtering. For the three example literature datasets used here, we implemented

221   plausible FA / LCB constraints and performed the matching to the database.

222   Reasonable manual preprocessing steps, such as re-formatting the data, removing

223   duplicate entries, and adjusting unusual nomenclature were performed beforehand,

224   and are available as Python scripts in S2 dataset. The processing statistics are listed

225   in Table 1.

226   **Table 1.** Matching statistics for development datasets.

| Dataset | Num. lipids | Successfully matched | Parsed - not matched | Not parsed - not matched | Filtered |
|---|---|---|---|---|---|
| *Drosophila* | 359 | 324 (90.3%) | 9 (2.5%) | 4 (1.1%) | 22 (6.1%) |
| Yeast | 249 | 235 (94.4%) | 14 (5.6%) | 0 | 0 |
| LAMP3 | 209 | 199 (95.2%) | 3 (1.4%) | 0 | 7 (3.3%) |

227

# Visualization

229   One has to look at complete databases to judge the vector space and embedding of

230   lipids. A user, however, would be interested in what one sees in their lipidome. We

231   take three examples from the literature and look at the scatterplots in the light of the

232   biochemistry noted by the original authors.

233   The first dataset consists of lipidomes of different *Drosophila melanogaster* larval

234   tissue types (brain, fat body, gut, lipoprotein, salivary gland, wing disc) fed with

235   different diets (plant food or yeast food) [23]. For our quick analysis, we averaged the

236    lipidome samples by tissue type. Carvalho et al noted that hexosyl ceramides

237    (HexCer) and ether glycerophospholipids (O-) were only detected in gut and brain

238    tissues respectively. Fig 4 A shows how this kind of feature can be easily observed

239    and highlighted. Fig 4 B displays a comparison of fat body and lipoprotein tissue

240    types focused on a glycerolipid region and highlights the expected large amounts of

241    triacylglycerol (TG) species in the fat body and conversely an overabundance of

242    diacylglycerols (DG) in the lipoprotein tissue, both noted in the original publication.

243    **Fig 4. Lipidome Dataset Projections.** (A) *Drosophila* dataset averaged over tissue

244    type. HexCer and ether-linked GPs are only present in gut and brain tissues

245    respectively. Min-max scaling of abundances was used to calculate marker area. (B)

246    *Drosophila* dataset zoomed in to a glycerolipid region of the vector space showing

247    selected tissue samples (same marker scaling as in A). (C) Yeast lipidomes –

248    comparison between the means of the wildtype and the Elo2 and Elo3 strains with

249    min-max marker scaling. (D) Yeast dataset zoomed in on a region of partially

250    annotated sphingolipids (same marker scaling as in C). Elo2 and Elo3 strains contain

251    species with shorter fatty acyls. (E) Mouse lung lipidome dataset lipids coloured by

252    the $\log_2$ abundance fold change between the wildtype and LAMP3-KO asthma

253    conditions. Certain lipids with relatively high change values are annotated. (F) PG

254    region comparison between wildtype and LAMP3-KO asthma conditions. Linear

255    scaling applied to marker sizes.

256    The second example is focussed on a yeast study comparing the wildtype strain

257    (BY4741) and mutants that were defective in fatty acyl elongation (Elo1, Elo2, Elo3)

258    [24]. Two different growth temperatures (24°C and 37°C) were considered. The study

259    showed that the Elo2 and Elo3 strains produce sphingolipids with shorter fatty acyl

260    chains. We averaged the samples by strain, filtered Elo1, and projected the full

11

261    results onto our vector space (Fig 4 C). Fig 4 D displays sphingolipid abundances

262    from the wildtype strain compared to average abundances from the Elo2 and Elo3

263    group, clearly showing that species with shorter fatty acyls occurring in the Elo strains

264    with a higher prevalence.

265    The third dataset is taken from a study of LAMP3-deficient mice, evaluating the role

266    of this protein in the lung [25]. The two different conditions genotype (wildtype /

267    LAMP3-KO) and challenge (none / allergen induced asthma) resulted in four groups

268    of mice. Fig 4 E and F show that if we average the samples by genotype and

269    challenge and compare the wildtype to the LAMP3-KO genotypes in the asthma

270    group, there is a large reduction in phosphatidylglycerols in the LAMP3-KO group, as

271    noted by the authors. Fig 4 E also shows the increased abundance of diacylglycerols

272    and decreased amounts of certain sphingolipids and phosphatidylinositols in the

273    wildtype group.

# Discussion

275    There are two aspects to this work. Firstly, there is the fundamental embedding of

276    molecules in a low-dimensional space. Secondly, there are practical issues and the

277    software implementation.

278    From the point of view of the vector space, there are some surprising observations.

279    The lipid coordinates agree with chemical intuition, although the training was

280    completely unsupervised. The lipids compositions from myriads of substructure

281    vectors on their own produce a systematically organized vector space, which is

282    improved by substructure vector training. Not only were classic lipid categories

283    separated, but unusual structures are given coordinates on the edges of the common

284    lipid classes (Fig 2 B). The local and global structure of the embedding is interesting.

285    Globally, the space reflects broad classes, but locally, it is remarkable that moving a

286    hydroxylation along a chain gives a set of points near each other and almost lying on

287    a smooth curve. There is reason to say this is unexpected. Consider the space as

288    first calculated in 100 dimensions. Maybe there are directions corresponding to

289    phosphorylation, chain extension, moving bonds and other chemical properties.

290    When we project the space to two or three dimensions, one will inevitably lose

291    information. The local structure is a tribute to stochastic nearest neighbour-

292    embedding rather than any invention on our part.

293    There are also differences compared to other vector spaces for lipids. Marella et al

294    calculated the differences between molecules using the differences between string

295    representations of the molecules [5]. This suffers from the instability of string

296    representations. Kopczynski et al avoided this problem by using graph-based

297    similarity [7]. There is a less obvious difference in the methods. Kopczynski et al

298    calculated distances between lipids and used principal coordinate analysis to get low

299    dimensional coordinates from the distance matrix. This is deterministic, but

300    discarding everything after the few most important eigenvectors is a brutal truncation.

301    Our method also requires dimensional reduction, but our experiments with principal

302    component analysis suggested that too much local structure was lost. We would

303    concede that stochastic neighbour embedding is not deterministic, the cost function

304    details are ad hoc and it does not have the geometric rigour of principal component

305    analysis. It does, however, seem to preserve relationships between neighbouring

306    molecules.

307    Kopczynski et al's approach does admit one feature that we lack. We construct a

308    space based on all known lipids and then show all lipidomes in this context. In

309    contrast, Kopczysnki et al build a new space for each set of lipidomes. This allows

310    them to construct a very natural measure for the similarity of lipidomes and lends

311    itself to clustering of datasets.

312    Continuing in this self-critical vein, the non-determinism of our approach might be

313    considered a disadvantage. Repeating the training and dimensional reduction always

314    gives slightly different results. With more training time or different parameters, one

315    might get even better results. Having experimented in this direction, we suspect that

316    this is not a useful pursuit. It would be more profitable to consider completely different

317    strategies. Graph convolutional networks would be a natural fit to molecular

318    structures [26] and one could experiment with novel dimensionality reduction

319    methods such as UMAP [27].

320    Besides the embedding, other issues should be addressed. We are not the first

321    group to lament the inherent inconsistency of lipid nomenclature [18]. Synonyms

322    such as SM(d18:1/14:0) and SM 18:1;2/14:0 are tedious but can be handled

323    mechanically by packages such as Goslin. A more fundamental problem are lipid

324    notation ambiguities which cannot be solved by any parser.

325    In this study we encountered ambiguities in the position, number and precise location

326    of double bonds and hydroxylations of sphingolipids. Some line notations would allow

327    one to denote some ambiguities [28], but lipidome data is typically not stored in such

328    formats. Another problem is that a user lipidome may contain species that are not in

329    the training set (SwissLipids + LMSD). This problem will be alleviated when we

330    implement an on-the-fly method to generate structures and respective vectors from

331    nomenclature only.

332    The second half of this work is the software. With the vector space precomputed, it is

333    not too demanding to run on an ordinary laptop. The web application stores lipidome

334  data on the client side and sends it to the server for processing operations. This does

335  require a fair amount of client-server communication, but we are currently moving

336  more processing tasks to the client's browser. Software is also a matter of taste. The

337  current release displays properties such as relative abundances using very compact

338  methods, but these might at first seem foreign to a user.

339  There are clear directions for the future. There will be improvements to the underlying

340  vector space as we experiment with the embedding model and as the databases are

341  updated. The software will change as a result of user experience, and it will

342  automatically benefit from the evolution of the parsing package [19]. Finally, we plan

343  proper integration with biochemical pathway software. As it stands, the vector space

344  is conceptually useful, and the software fills a practical niche.

## Acknowledgments

## Availability

350  Lipidome Projector is available for download

351  (https://www.github.com/olzhabaev/lipidome_projector) and released under the MIT

352  license. It is a web-application that can be run locally or deployed to a server. The

353  repository has pre-computed vectors for and pre-parsed versions of the Lipid Maps

354  and SwissLipids databases. The software distribution also includes modules for the

355  pre-processing of the databases and a complete recalculation of the vector space.

356  An instance of Lipidome Projector is available at: https://lipidomeprojector.zbh.uni-

357  hamburg.de/

15

# References

358

359  1.  Mohamed A, Hill MM. LipidSuite: interactive web server for lipidomics

360      differential and enrichment analysis. Nucleic Acids Res. 2021;49: W346–W351.

361      doi:10.1093/nar/gkab327

362  2.  Mohamed A, Molendijk J, Hill MM. lipidr: A Software Tool for Data Mining and

363      Analysis of Lipidomics Datasets. J Proteome Res. 2020;19: 2890–2897.

364      doi:10.1021/acs.jproteome.0c00082

365  3.  Kyle JE, Aimo L, Bridge AJ, Clair G, Fedorova M, Helms JB, et al. Interpreting

366      the lipidome: bioinformatic approaches to embrace the complexity.

367      Metabolomics. 2021;17: 55. doi:10.1007/s11306-021-01802-6

368  4.  Köhler N, Rose TD, Falk L, Pauling JK. Investigating Global Lipidome

369      Alterations with the Lipid Network Explorer. Metabolites. 2021;11: 488.

370      doi:10.3390/metabo11080488

371  5.  Marella C, Torda AE, Schwudke D. The LUX Score: A Metric for Lipidome

372      Homology. PLoS Comput Biol. 2015;11: e1004511.

373      doi:10.1371/journal.pcbi.1004511

374  6.  Weininger D. SMILES, a chemical language and information system. 1.

375      Introduction to methodology and encoding rules. J Chem Inf Comput Sci.

376      1988;28: 31–36. doi:10.1021/ci00057a005

377  7.  Kopczynski D, Hoffmann N, Troppmair N, Coman C, Ekroos K, Kreutz MR, et

378      al. LipidSpace: Simple Exploration, Reanalysis, and Quality Control of Large-

379      Scale Lipidomics Studies. Anal Chem. 2023;95: 15236–15244.

380      doi:10.1021/acs.analchem.3c02449

381    8.    Jaeger S, Fulle S, Turk S. Mol2vec: Unsupervised Machine Learning Approach

382          with Chemical Intuition. J Chem Inf Model. 2018;58: 27–35.

383          doi:10.1021/acs.jcim.7b00616

384    9.    Mikolov T, Chen K, Corrado G, Dean J. Efficient Estimation of Word

385          Representations in Vector Space. 2013.

386    10.   Rogers D, Hahn M. Extended-Connectivity Fingerprints. J Chem Inf Model.

387          2010;50: 742–754. doi:10.1021/ci100050t

388    11.   Sud M, Fahy E, Cotter D, Brown A, Dennis EA, Glass CK, et al. LMSD: LIPID

389          MAPS structure database. Nucleic Acids Res. 2007;35: D527–D532.

390          doi:10.1093/nar/gkl838

391    12.   Aimo L, Liechti R, Hyka-Nouspikel N, Niknejad A, Gleizes A, Götz L, et al. The

392          SwissLipids knowledgebase for lipid biology. Bioinformatics. 2015;31: 2860–

393          2866. doi:10.1093/bioinformatics/btv285

394    13.   RDKit: Open-source cheminformatics. https://www.rdkit.org/.

395    14.   Řehuřek R, Sojka P. Software Framework for Topic Modelling with Large

396          Corpora. Proceedings of the LREC 2010 Workshop on New Challenges for

397          NLP Frameworks. Valletta, Malta: ELRA; 2010. pp. 45–50.

398    15.   Van Der Maaten L. Accelerating t-SNE using tree-based algorithms. The

399          journal of machine learning research. 2014;15: 3221–3245.

400    16.   der Maaten L, Hinton G. Visualizing data using t-SNE. Journal of machine

401          learning research. 2008;9.

402  17.  Poličar PG, Stražar M, Zupan B. OpenTSNE: A modular Python library for t-

403       SNE dimensionality reduction and embedding. bioRxiv. 2019.

404       doi:10.1101/731877

405  18.  Liebisch G, Fahy E, Aoki J, Dennis EA, Durand T, Ejsing CS, et al. Update on

406       LIPID MAPS classification, nomenclature, and shorthand notation for MS-

407       derived lipid structures. J Lipid Res. 2020;61: 1539–1555.

408       doi:10.1194/jlr.S120001025

409  19.  Kopczynski D, Hoffmann N, Peng B, Ahrends R. Goslin: A Grammar of

410       Succinct Lipid Nomenclature. Anal Chem. 2020;92: 10957–10960.

411       doi:10.1021/acs.analchem.0c01690

412  20.  Inc. PT. Collaborative data science. Montreal, QC: Plotly Technologies Inc.;

413       2015. Available: https://plot.ly

414  21.  Van Rossum G, Drake Jr FL. Python tutorial. Centrum voor Wiskunde en

415       Informatica Amsterdam, The Netherlands; 1995.

416  22.  McKinney W. Data Structures for Statistical Computing in Python. 2010. pp.

417       56–61. doi:10.25080/Majora-92bf1922-00a

418  23.  Carvalho M, Sampaio JL, Palm W, Brankatschk M, Eaton S, Shevchenko A.

419       Effects of diet and development on the Drosophila lipidome. Mol Syst Biol.

420       2012;8. doi:10.1038/msb.2012.29

421  24.  Ejsing CS, Sampaio JL, Surendranath V, Duchoslav E, Ekroos K, Klemm RW,

422       et al. Global analysis of the yeast lipidome by quantitative shotgun mass

423       spectrometry. Proc Natl Acad Sci U S A. 2009;106.

424       doi:10.1073/pnas.0811700106

425   25.   Lunding LP, Krause D, Stichtenoth G, Stamme C, Lauterbach N, Hegermann J,

426         et al. LAMP3 deficiency affects surfactant homeostasis in mice. PLoS Genet.

427         2021;17. doi:10.1371/journal.pgen.1009619

428   26.   Hamilton WL, Ying R, Leskovec J. Representation Learning on Graphs:

429         Methods and Applications. 2017.

430   27.   Sainburg T, McInnes L, Gentner TQ. Parametric UMAP Embeddings for

431         Representation and Semisupervised Learning. Neural Comput. 2021; 1–27.

432         doi:10.1162/neco_a_01434

433   28.   Homer RW, Swanson J, Jilek RJ, Hurst T, Clark RD. SYBYL Line Notation

434         (SLN): A Single Notation To Represent Chemical Structures, Queries,

435         Reactions, and Virtual Libraries. J Chem Inf Model. 2008;48: 2294–2307.

436         doi:10.1021/ci7004687

437

# Supporting information

## Figures

**S1 Fig. Lipidome Projector Interface, *Drosophila* dataset.** Top left: Lipidome

dataset scatter plot; Top right: Settings, data operations and abundance charts.

Bottom: Abundance and feature tables.


**S2 Fig. Vector Space (3D)**. (A) Projection of the entire vector space. Marker colour

represents lipid category: Fatty acids (FA), glycerolipids (GL), glycerophospholipids

(GP), sphingolipids (SP), sterol lipids (ST), prenol lipids (PR), saccharolipids (SL) and

polyketides (PK). (B) Region of the vector space focused on a set of selected

glycerophospholipids: Glycerophosphates (PA), glycerophosphocholines (PC),

glycerophosphoethanolamines (PE), glycerophosphoglycerols (PG),

glycerophosphoinositols (PI) and glycerophosphoserines (PS). Marker colour: Lipid

class. (C) Same region as in B. Marker colour: Number of fatty acyl double bonds.

(D) Zoomed in region of selected glycerophospholipids. Marker colour: Lipid class.

(E) Same region as in D. Marker colour: Double bond profile of the 2-sn fatty acyl. (F)

Same region as in D. Marker colour: Molecule mass.

## Tables

**S1 Table. Word2Vec Embedding Parameters.**

**S2 Table. Stochastic Neighbour Embedding Parameters.**

## Data

**S1 Dataset.** List of classes present in LMSD and SwissLipids recognised by the Goslin parser in translated representation.

**S2 Dataset.** Python scripts with instructions for the extraction and transformation of original datasets; Transformed datasets; Dataset FA / LCB constraints.

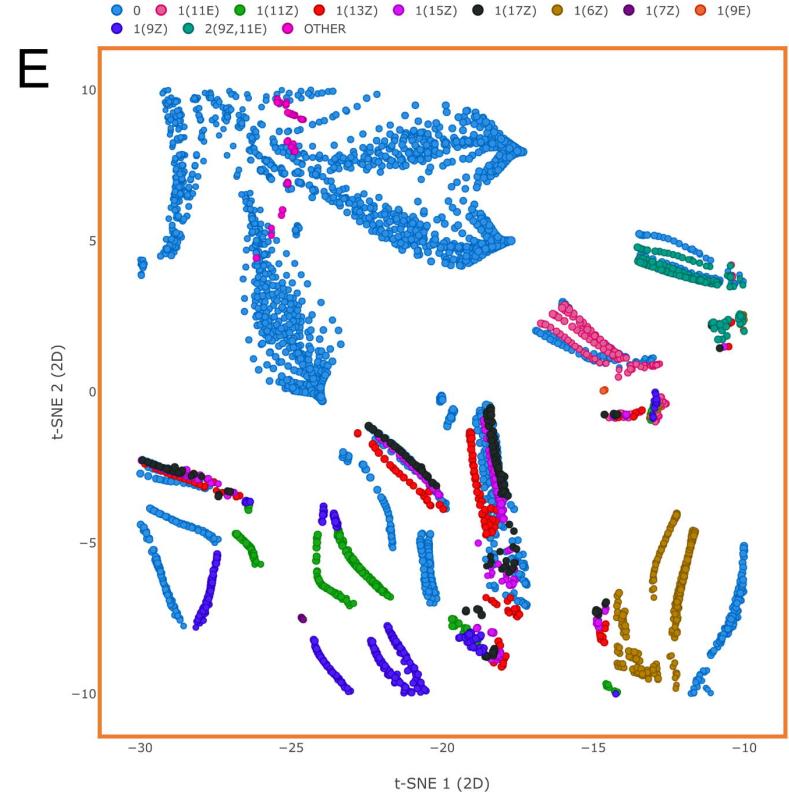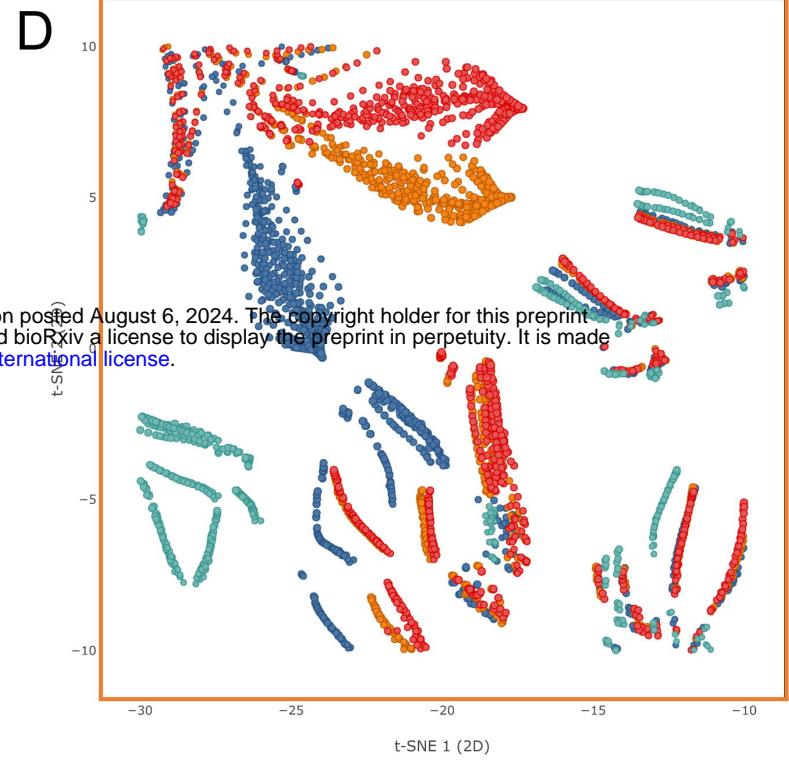**S3 Dataset.** Partially interactive HTMLs of vector space and dataset projection scatter plots.

**A**



**B**