

Unsupervised Machine Learning for Analysis of Coexisting Lipid Phases and Domain Growth in Biological Membranes

Cesar A. López¹, Velimir V. Vesselinov², Sandrasegaram Gnanakaran^{1*}, Boian S.

Alexandrov^{1*}

¹Theoretical Division, Los Alamos National Laboratory, Los Alamos, NM 87545, USA

²Earth and Environmental Sciences Division, Los Alamos National Laboratory, Los
Alamos, NM 87545, USA

KEYWORDS: Lipid Phases, Machine Learning, Nonnegative Matrix Factorization

1
2
3 ABSTRACT: Phase separation in mixed lipid systems has been extensively studied both
4
5 experimentally and theoretically because of its biological importance. A detailed
6
7 description of such complex systems undoubtedly requires novel mathematical
8
9 frameworks that are capable to decompose and categorize the evolution of thousands if
10
11 not millions of lipids involved in the phenomenon. The interpretation and analysis of
12
13 Molecular Dynamics (MD) simulations representing temporal and spatial changes in such
14
15 systems is still a challenging task. Here, we present a new unsupervised machine
16
17 learning approach based on Nonnegative Matrix Factorization, called NMFk, that
18
19 successfully extracts physically meaningful features from neighborhood profiles derived
20
21 from coarse-grained MD simulations of ternary lipid mixture. Our results demonstrate that
22
23 leveraging NMFk can (a) determine the role of different lipid molecules in phase
24
25 separation, (b) characterize the formation of nano-domains of lipids, (c) determine the
26
27 timescales of interest and (d) extract physically meaningful features that uniquely
28
29 describe the phase separation with broad implications.
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

INTRODUCTION

Cell membranes contain mixtures of different lipid types, dynamically arranged, that play a key role in various mechanisms responsible for cell survival [1](#). In the past, membranes were thought to be homogenous systems, however, new data suggests that under different stimulus, the lipids can segregate [2-3](#) into detergent-resistant domains commonly called “rafts”. These domains are highly dynamic, varying in size and composition [4-5](#). Importantly, this lateral partitioning is responsible for activation and functioning of membrane-embedded proteins [6-7](#).

While it is possible to experimentally visualize the structure of segregated lipid domains [2, 8](#), a detailed description of such phases are inherently limited by the resolution of the experimental techniques. At this respect, molecular dynamics (MD) simulations provide a molecular understanding of membrane behavior. In fact, the presence of such lateral rearrangement has been studied extensively using coarse-grained (CG) MD of ternary lipid mixtures [9](#). These CG simulations suggest partial segregation in large lipid systems that mimic the lipid variability of the real cell membranes [10](#). MD simulations of such realistic biomolecular systems usually contain millions of particles, even when simplified models are used. When the purpose of such simulations is to gain biological or physical insights, it is challenging to identify patterns in the behavior encoded in the motion of thousands of molecules, so developing analytic tools for extracting functionally relevant features from MD generated trajectories is of great importance.

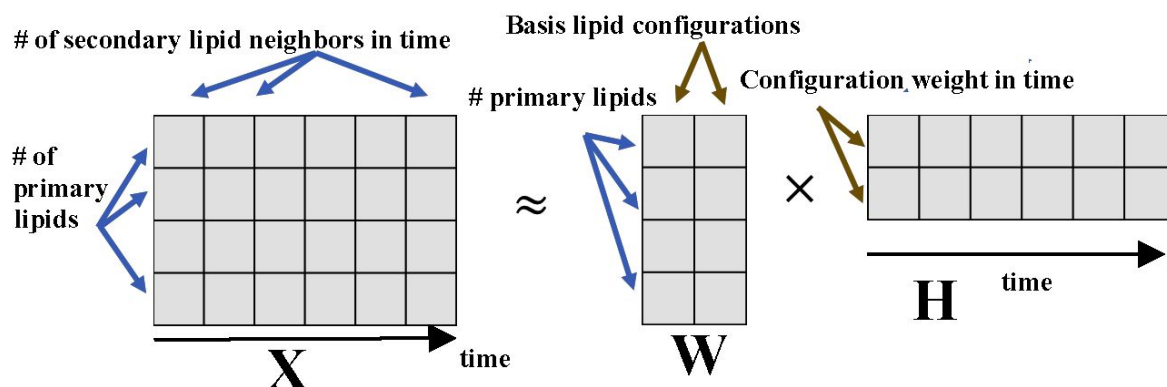
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Currently, machine learning (ML) methods have shown a lot of promise in many fields [11](#), and recently some of them have been applied for analysis and detection of classical and quantum phase transition data generated by simulations. Both supervised and unsupervised ML approaches have been used for this purpose [12-19](#), but most of these pioneer studies used small Ising-like systems for their investigations. ML has previously been coupled with MD simulations of biomolecular systems in a limited context. ML techniques were reported to predict free-energy differences when trained with MD simulation data [20](#), and unsupervised approaches such as PCA [21](#) as well as other techniques [22-23](#) have been used to reduce the dimensionality of MD generated data [24-28](#).

In general, the unsupervised ML methods learn relationships between elements in uncategorized data and classify the data without human's help but by revealing its internal structure and latent (i.e., not directly observable) features hidden in the data. The unsupervised methods include clustering [29](#), classical neural networks [30](#), and the more contemporary blind source separation (BSS) techniques [31](#). BSS is based on *factorization* which is one of the most powerful tools for extraction of latent features [32](#). BSS include principle component analysis (PCA) [33](#), singular value decomposition (SVD) [34](#), and more advanced methods, such as independent component analysis, ICA [35](#) and nonnegative matrix factorization, NMF [36](#).

A limitation shared by PCA, SVD and ICA is the difficulty to relate the extracted latent factors to physically interpretable quantities; NMF overcomes this limitation because the nonnegativity of the extracted latent factors leads to a collection of strictly additive

1
2
3 components that are sparse and parts of the data and hence are amenable to a simple
4 and meaningful interpretation without introducing prior assumptions³⁷. NMF decomposes
5 given data matrix, X_{LN} , into two matrices, W_{LK} , and H_{KM} , such that $X_{LN} \approx W_{LK} H_{KM}$. The
6 factor matrices, W_{LK} and H_{KM} , are both nonnegative and have one small dimension K that
7 represents the number of the latent features in the data (Fig. 1). A mathematically rigorous
8 formalism is given in the later sections. NMF ability to identify easy interpretable latent
9 features enables discoveries of new causal structures and unknown mechanisms hidden
10 in the data as discussed in the literature³⁸. Surprisingly, the implementation of NMF for
11 analysis of MD simulations at the interface of physics and biology has been lacking.



26
27
28
29
30
31
32
33
34
35
36
37
38
39 **Figure 1.** Illustration of a Nonnegative Matrix Factorization. The nonnegative matrix X
40 is decomposed to the product of a nonnegative matrix W , containing $K=2$ basis lipid
41 configurations, and nonnegative matrix H , containing the contributions of these two
42 configurations, in different time points.

43
44
45
46
47
48
49
50 Here, we present a new unsupervised machine learning algorithm based on NMF
51 combined with custom k-means clustering, called NMFk, capable of analyzing phase
52 separation in a system of mixed lipids directly from the pre-processed trajectories derived
53
54
55
56
57
58
59
60

1
2
3 by MD simulations. We use CG MD simulations of a physical system that comprises a 3-
4 component lipid mixture, commonly accepted to mimic the behavior of a cellular plasma
5 membrane [5](#). We show that NMFk, applied to a pre-selected data from these simulations
6
7 is able to, (a) determine the molecules that play different roles in phase separation, (b)
8 characterize the formation of lipid nano-domains, (c) reveal the timescales of interest, and
9
10 (d) extract physically meaningful features that characterize the phase separation.
11
12
13
14
15
16
17
18
19
20

21 RESULTS

22 **Generation of lipid mixture data sets using coarse-grained MD simulations**

23
24
25 For MD simulations of membranes and membrane-based biological systems, the Martini
26 coarse-grained (CG) force field [39](#), considerably reduces the computational cost of
27 calculations by nearly three orders of magnitude compared with similar MD simulations
28 using fully atomistic force fields [40](#). Particularly, the CG approach can capture relevant
29 dynamics and fluctuations of larger membrane patches, which are prohibitive with
30 atomistic simulations. Such access to larger spatial and temporal scales enables direct
31 comparisons with experimental measurements [41](#). Pioneering computational studies with
32 the Martini CG force field allowed the characterization of not only lipid segregation and
33 lipid phases, but also the relative partitioning of membrane proteins between these
34 phases [9, 42-44](#). More recently, Martini has been used in simulations of membranes with
35 lipid compositions of comparable complexity to those found in specific tissues of living
36 cells [45-46](#).
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6 Regardless of the extensive use of the Martini force field, the building-block
7
8 principle of Martini along with the 4:1 atoms to bead mapping unavoidably reduce the
9
10 accuracy due to loss of detailed description of specific molecular chemical properties.
11
12 Thus, despite the fact that many of the current Martini lipid parameters are sufficient to
13
14 guide accurate membrane simulations [39, 47-48](#), global lipid properties are compromised [49](#).
15
16 Consistent with previously published work [42, 50-51](#), the standard Martini V2.2 parameters
17
18 for DPPC, DOPC, and CHOL do not phase separate at 298 K, or even at 290 K in the
19
20 physiologically relevant L_d/L_o coexistence region of the phase diagram. Recently, we
21
22 incorporated changes in the lipid Martini force field [52](#), greatly improving the lipid
23
24 segregation in line with current experimental phase diagrams [8, 50](#).
25
26
27
28
29

30
31 We use the above two versions of the Martini force field, one lipid parameter set
32
33 that phase separates and the other that does not, to generate two data sets for the
34
35 analysis using NMFk. The first set considers the data from the coarse-grained MD
36
37 simulations of the current Martini force field (called “Standard”). As mentioned above, the
38
39 current Martini V2.2 lipid parameters for DPPC and DOPC are not able to properly
40
41 segregate the DPPC:DOPC:CHOL mixture. It serves as the prototype for non-phase
42
43 separating homogeneous lipid mixture. The second set considers the refined version of
44
45 the Martini (called “Updated”), which has been optimized to reproduce the experimental
46
47 phase separation and domain formation for this ternary system [52](#). It serves as the
48
49 prototype for phase separating lipid mixture. Using these two versions of force fields, we
50
51 carried out 20 μ s long CG MD simulations. The expectation is that we should be able to
52
53
54
55
56
57
58
59
60

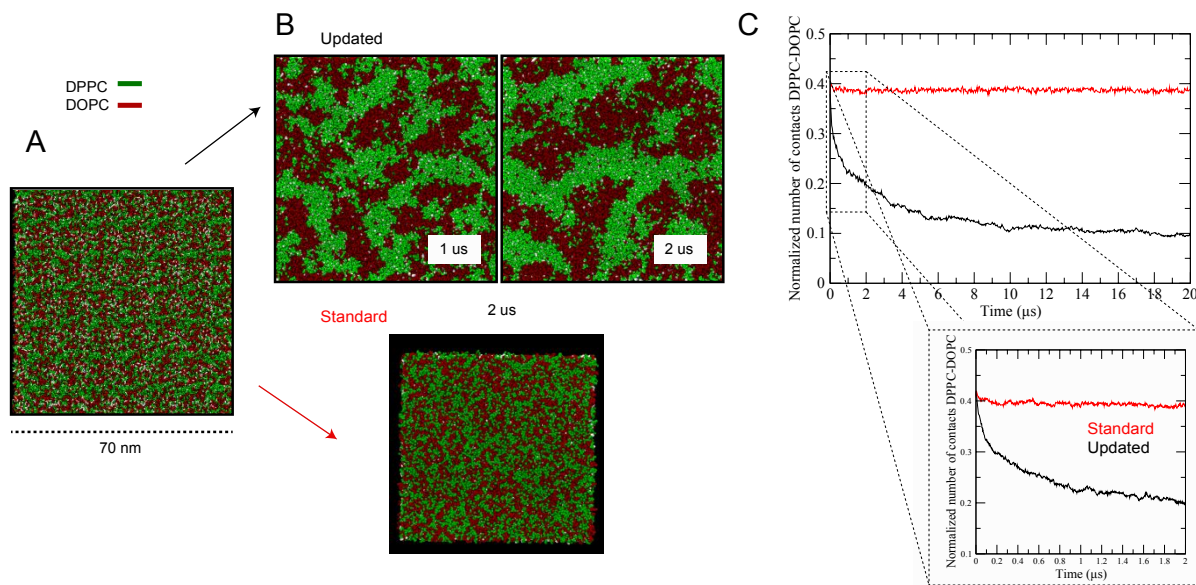
1
2
3 detect and analyze features associated with phase separation in one case but not in the
4
5 other case.
6

7 8 9 10 **Conventional Analysis of Domains in MD Simulations of Ternary lipid Mixtures**

11
12
13
14
15 The formation of lipid domains has been heavily studied both experimentally and
16
17 computationally [2-4](#), [8-10](#), [53-55](#) Computational observation of explicit lipid segregation at
18
19 nearly atomic detail dates back to almost 10 years ago [9](#) Analysis of such processes
20
21 involved direct visualization of cholesterol-rich/-poor domains, as well as physical
22
23 quantification of the area per lipid, cholesterol content, radial distribution functions and
24
25 membrane thickness mismatch [9](#). These analyses brought enough details that membrane
26
27 domains could be discriminated at the nanoscale. A more sophisticated approach can be
28
29 cited from the work of Baoukina et al. [51](#), where a Voronoi tessellation methodology was
30
31 applied in order to delineate the boundaries between ordered/disordered domains in
32
33 monolayers. This approach can be also directly combined with automated predictive tools
34
35 like Markovian based methodologies [54](#). Regardless of these published methodologies,
36
37 the analysis and prediction of such fluctuating membrane domains become inaccessible
38
39 when the complexity of lipid content increases with larger size membrane patches (e.g.
40
41 plasma membrane).
42
43
44
45
46
47
48

49 A typical process of molecular lipid phase separation is depicted in **Fig. 2**. Initially,
50
51 the lipid components are randomized, mimicking the homogeneity at a high-temperature
52
53 (**Fig. 2A**). Subsequent fast quenching of the mixture to 290 K, well below the melting
54
55
56
57
58
59
60

1
2
3 temperature of the fully saturated DPPC lipid, leads to the rapid formation of nanoscale
4 domains on a submicrosecond time scale with the Updated Martini force field (**Fig. 2B**
5
6 domains on a submicrosecond time scale with the Updated Martini force field (**Fig. 2B**
7
8 **top panel**). These nano-domains are eventually formed over the entire surface of the
9
10 membrane, but in different regions. After 0.5 microseconds, the nano-domains start to
11
12 interconnect, leading to the formation of larger cholesterol-rich regions. In agreement with
13
14 the general raft hypothesis [53](#) and previous computational studies [9](#), the “ordered” nano-
15
16 domains contain most of the saturated lipids together with cholesterol forming a Liquid
17
18 ordered (L_o) domain, whereas the “disordered” nano-domain is mainly composed of the
19
20 polyunsaturated DOPC lipid segregated in a Liquid disordered (L_d) region. Contrary to
21
22 the Updated force field, the lipid mixture based on the Standard Martini force field does
23
24 not show any tendency for phase separation or domain formation (**Fig. 2B, bottom**
25
26 **panel**), in agreement with previously published data [52](#).



31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51 **Figure 2.** Coarse-grained simulations of ternary lipid mixture. **(A)** Initial system setup for
52 the CG simulations. Lipids were initially randomly placed within the XY plane. Saturated
53 lipid, DPPC, is colored in green while unsaturated lipid, DOPC, is colored in red.
54
55
56
57
58
59
60

1
2
3 Cholesterol is colored in white. Water is not shown for clarity in depiction. **(B)** Lipid phase
4 separation within 2 us simulation. The updated CG force field shows clear phase
5 separation into L_o (green) and L_d (red). Contrary, the standard Martini lipid force field
6 does not show preferential segregation in cholesterol rich domains. **(C)** Time evolution of
7 the normalized number of contacts between saturated and unsaturated lipids, showing
8 poor separation with the standard Martini force field.
9
10
11
12
13
14
15
16
17
18
19

20 Following the conventional approaches to quantify the segregation tendency, we
21 compute the normalized total contacts between DPPC and DOPC as a function of
22 simulation time (**Fig. 2C**). Initially, these contacts are featured by larger values, meaning
23 that these two lipid types are indeed in close contact, highlighting the initial homogeneous
24 lipid mixing of the system. However, with the Updated Martini, lipid segregation leads to
25 a decrease in the total number of contacts between DPPC and DOPC (**Fig. 2C black**
26 **line**). Whereas, the contacts remain unchanged in the simulations with the standard
27 Martini (**Fig. 2C red line**). In the case of Updated Martini, contacts decay during the
28 simulations and begin to plateau with increasing simulation time. We should note here
29 that the proper convergence to a stationary state may not be achievable within the
30 simulation timescales considered here, as already published [54](#), while significant transition
31 towards segregation occurs within the first 2 us (**Fig 2C insert**). We consider this time
32 regime suitable for NMFk analysis to extract latent features associated with the phase
33 separation.
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53

54 NMFk implementation for analysis of ternary lipid mixture simulations

55
56
57
58
59
60

1
2
3
4
5
6 To perform the NMFk analysis, we first define the primary and secondary lipid types and
7
8 then calculate the neighbor matrix at a given time t , $X_L(t)$, in terms of the number of
9
10 secondary lipid type surrounding the primary lipid type. For this study, we considered
11
12 DPPC and DOPC as the primary and secondary lipid types, respectively, although, there
13
14 is no restriction of other combinations for primary and secondary lipid types. We compute
15
16 the number of closest DOPC neighbors, $X_L(t)$, around every DPPC lipid, $L = (L_1, L_2, \dots, L_M)$
17
18 , within the distance corresponding to the second peak of the DOPC-DPPC radial
19
20 distribution function (**Fig. 3 blue dashed line**). The number of neighbors within this
21
22 distance serves as the order parameter for the lipid phase separation. This neighbor data
23
24 is recorded at N consecutive simulation time points, $t = (t_1, t_2, \dots, t_N)$, corresponding to
25
26 the time evolution of the system up to first 2 us. The matrix $X_L(t)$ contains N arrangements
27
28 at N time points, presented by the columns of X , that form the matrix of the neighbors, X_L
29
30 (t); $X \in M_{MN}(I_+)$, where M is the number of the rows, and each row represents the time
31
32 evolution of the closest neighbors of one specific lipid of the primary type. The I_+ denotes
33
34 the set of nonnegative integer numbers. Next, NMFk decomposes the neighbor matrix X_L
35
36 (t) as a product of two matrices: $X_L(t) \approx W_{LK} * H_K(t)$, where the columns of W_{LK} are the
37
38 K basis lipid configurations describing the state of the lipid system and $H_K(t)$ contains
39
40 the contributions of each one of these configurations at time t (see **Fig. 1** where $K=2$).
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

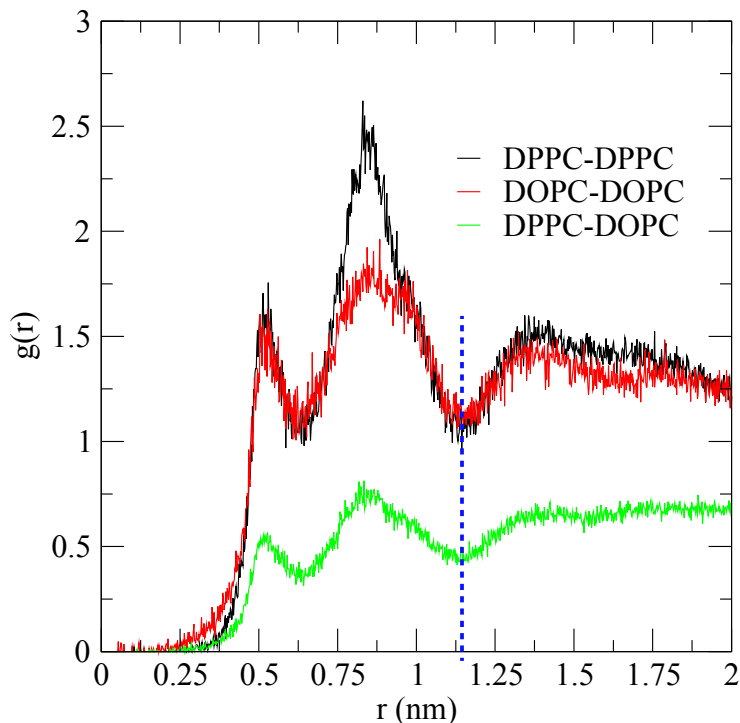


Figure 3. Lateral radial distribution function for the different lipid combinations. RDF was computed considering the center of mass of the molecules. The dashed blue line indicates the chosen cut-off distance for profiling neighbor list needed for the NMFk matrix construction.

The neighbor matrix, $X_L(t)$, contains relevant properties of the ternary lipid mixture. At the beginning of the simulations, the lipid mixture is homogeneous and there are no distinct phases, i.e., there is no specific structure in $X_L(t)$. Hence, $X_L(t)$ contains uniformly distributed integer numbers that do not have any distinct features. With the Updated Martini force field, by 2 us, the primary lipid type, DPPC, segregates into the L_o region, while the secondary lipid type, DOPC, segregates into the L_d domain. In this case, if a concrete lipid of the primary type is located deep in a L_o domain, the probability to have a lipid-neighbor of secondary type is small (close to zero). In contrast, if this primary lipid is

1
2
3 situated outside of any domain, the probability to have a neighbor of the secondary type
4
5 is much higher. Therefore, we expect the neighbor matrix to contain a structure that tracks
6
7 this phase separation and the subsequent analysis by NMFk to extract the hidden
8
9 features describing the structure and the phase separation.
10
11
12
13

14
15 The NMFk analysis, as we will see next, demonstrates exactly that: When phase
16
17 separation begins at time, $t = t_s$, for each one of the snapshots recorded at, $t > t_s$, and
18
19 for each of the M lipids of the primary type, the values of the matrix $X_L(t > t_s)$ can be
20
21 represented as a linear mixing of K basis lipid configurations presenting the probability of
22
23 the given lipid of the primary type, L_i , to have a closest neighbor-lipid of secondary type.
24
25 NMFk decomposes the matrix $X_L(t)$, to a nonnegative probability matrix, \mathbf{W} , $\mathbf{W} \in M_{MK}(\mathcal{R}_+)$
26
27), corresponding to these K basis lipid configurations, blended by the weights, presented
28
29 by the elements of a nonnegative matrix, \mathbf{H} , $\mathbf{H} \in M_{KN-s}(\mathcal{R}_+)$ that reflects how these
30
31 configurations are active and mix in time. Thus, for a given arrangement of the primary
32
33 lipids, X_L , at a time point $t > t_s$, we have,
34
35
36
37

$$38 \quad X_L(t > t_s) = \sum_{k=1}^K \mathbf{W}_k(L) \mathbf{H}_k(t > t_s) + \boldsymbol{\varepsilon}_L(t > t_s),$$

39
40 where $\boldsymbol{\varepsilon} \in M_{MN-s}(\mathcal{R}_+)$ denotes the presence of a noise or unbiased error of
41
42 decomposition. Before the time point $t = t_s$, there is no trace of a phase separation and
43
44 the pattern of the number of the closest neighbors is stochastic. Therefore, there is no
45
46 clear features that could be recognized by NMFk.
47
48
49
50

51
52
53 The reconstruction of $X_L(t > t_s)$, via the two factor matrices, $\mathbf{W}_k(L)$ and $\mathbf{H}_k(t > t_s)$,
54
55 serves as a measure for the significance and quality of the extracted latent features. For
56
57
58
59

1
2
3 the case of Updated Martini, **Table 1** presents the Pearson correlation coefficient between
4 the reconstructed lipid arrangements at each time point t (obtained by NMFk) and the
5 original arrangements at the same time point (i.e., the corresponding column of the
6 neighbor matrix $X_L(t)$). The K unique basis lipid configurations (encoded in the probability
7 matrix W), reproduced the $N - s$ lipid arrangements forming the matrix X_{MN-s} . With the
8 standard Martini, the NMFk was not able to provide a set of lipid configurations that can
9 reconstruct the simulations accurately. Indeed, and as shown in **Table 1**, the NMFk
10 analysis did not reproduce accurately the simulation data.
11
12
13
14
15
16
17
18
19
20
21
22
23

24 **NMFk derived features associated with phase separation of ternary lipid Mixture**

25
26
27

28 Here, we describe how the basis lipid configurations extracted by NMFk describe physical
29 properties of the phase separation, such as, time profile of which lipid molecules belongs
30 to which phases, the formation of nano-domain, and the spatial and temporal profiles of
31 nucleation and phase separation.
32
33
34
35
36
37
38
39

40 A typical outcome of NMFk analysis is presented on **Fig. 4**. According to our
41 Silhouette-Reconstruction criteria (see Methods Section), NMFk determines that the
42 optimal number of basis lipid configurations is 20 (**Fig. 4A**). The columns H_i of the matrix
43 H , each of which encodes the weights, that is, the participation of a basis lipid
44 configuration in time, are presented on (**Fig. 4B**). It is clear that for each H_i there is a well-
45 defined time interval containing a number of consecutive frames where the corresponding
46 basis configuration W_i is active. From **Table 1**, it can be seen that after 0.45 μ s, NMFk
47
48
49
50
51
52
53
54
55
56
57
58
59
60

reproduces the neighbor matrix $X_L(t)$ very well: the cross-correlation between the neighbor matrix and the reconstructed matrix, for each time frame, is above 0.95. The much lower cross-correlation for reconstruction of the matrix $X_L(t)$ at early times suggests a prevalence of homogeneous lipid mixture in the first 0.45 μ s. The 15 significant basis lipid configurations that are active at consecutive time intervals after the first 0.45 μ s (#6; #10; #12; #15; #17; #19; #20; #18; #16; #14; #13; #11; #7; #2; #1) represent the evolution of the nucleation and phase separation.

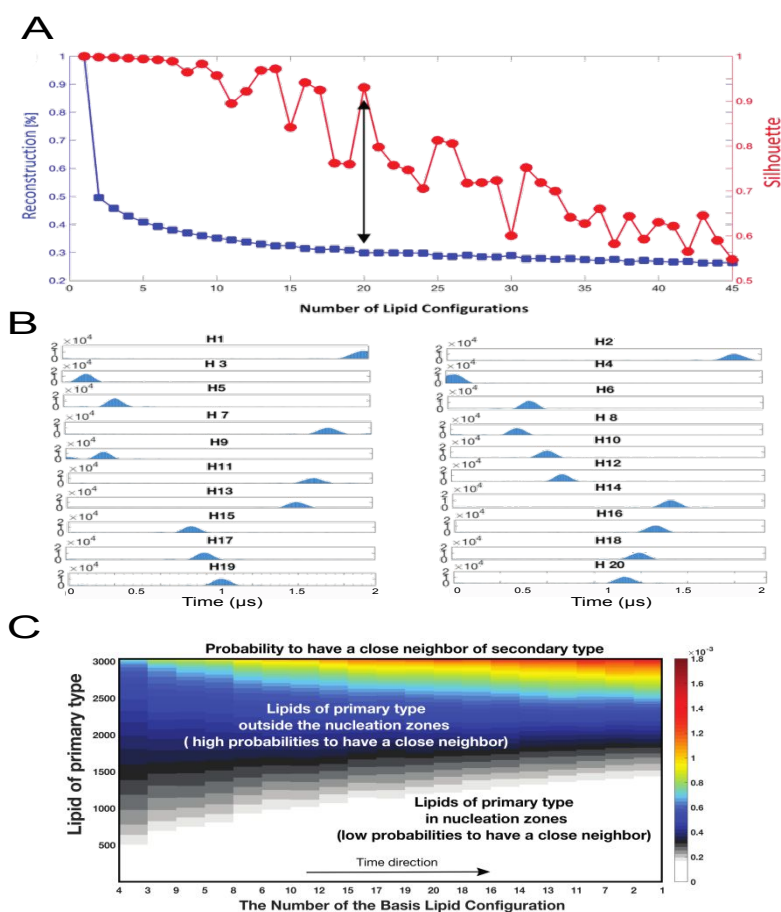


Figure 4. Outcome from NMFk analysis. **(A)** The Silhouette-Reconstruction criterium (see the Methods Section). On the x-axis is denoted the number of basis lipid configurations and on the y-axis-the average Silhouettes (right y-axis, red marking) as

1
2
3 well as the Reconstruction (left y-axis, the blue marking). The double arrow denotes
4
5 NMFk estimates of the number of basis lipid configurations. **(B)** Presentation of the
6
7 columns H_i of the matrix H that encodes the weights of the basis lipid configurations at
8
9 given time frame. **(C)** A heatmap presenting the basis lipid configurations ordered along
10
11 x-axis according to the time interval they occurred from early time to late as gathered
12
13 from the corresponding weight-columns, H_i . The color gradient from white to blue to red
14
15 represents the increase in probability for a specific primary lipid to have a secondary
16
17 lipid neighbor within each basis lipid configurations. The increase in white and gray
18
19 colors with time is indicative of the evolution of the phase separation.
20
21
22
23
24
25

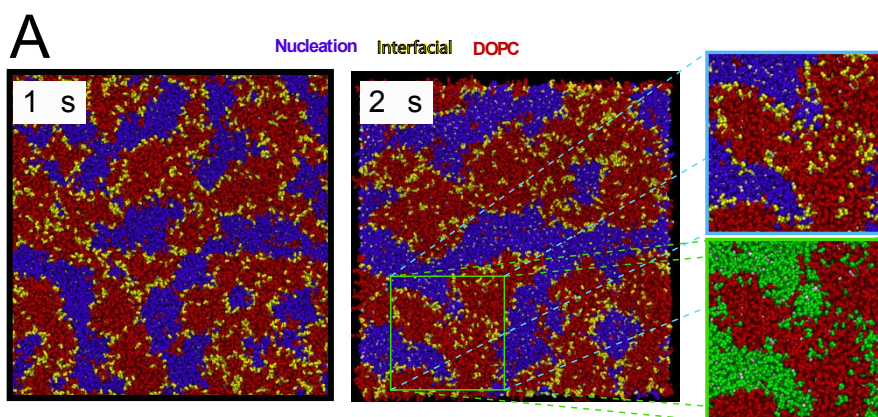
26 Each one of the 15 significant basis lipid configurations, extracted by NMFk,
27
28 contains a set of probabilities for the lipids of a primary type to have a close neighbor-lipid
29
30 of a secondary type. Each basis lipid configuration (i.e., the specific set of probabilities)
31
32 is active at a given time interval defined by the weight of this basis configuration at the
33
34 corresponding column of the matrix H . In each basis lipid configuration, we expect to have
35
36 at least two groups of probabilities: (a) the probabilities of primary lipids situated in the
37
38 nucleation domains that have relatively small number of neighbor-lipid of secondary type,
39
40 and therefore these probabilities, on average, approach zero; and (b) the probabilities for
41
42 primary lipids outside of any nucleation domain whose number of neighbor-lipid of
43
44 secondary type is much higher. There are always transition of the primary type of lipids
45
46 in and out of the nucleation domains or at the interface of the nucleation domains: at a
47
48 certain moment, a given primary lipid could be situated outside a nucleation domain
49
50 formed by primary lipids, but after a while it could reach interfacial region and eventually
51
52
53
54
55
56
57
58
59
60

1
2
3 get absorbed into the nucleation domain. Alternatively, primary lipids inside the nucleation
4 domain or at the interfacial region may venture out into the liquid disordered region
5 enriched of secondary lipids. These exchanges continuously alter the probability of a
6 given lipid of the primary type to have neighbor-lipid of secondary type, as the phase
7 separation proceeds which results in different basis lipid configurations at different time
8 points as the system goes to a phase separation. At long timescale, when the phase
9 separation has reached equilibrium, basis lipid configurations capture the exchanges of
10 the primary lipids governed by the stochasticity and diffusion.
11
12
13
14
15
16
17
18
19
20
21
22
23
24

25 To better characterize the above observations related to the extracted basis lipid
26 configurations, we applied k-means clustering on each of the 20 basis lipid configurations.
27 We combined the k-means clustering with Silhouette statistics to determine the most
28 probable number of clusters, i.e., the most probable number of groups of probabilities in
29 each basis lipid configuration. Specifically, we calculate consecutively the Silhouettes of
30 the resulting clusters, changing the number of clusters from 1 to 30 to determine the most
31 probable groups of primary (DPPC) lipids with similar probabilities to have a neighbor
32 secondary (DOPC) lipid-neighbor. **Fig. 4C** shows the basis lipid configurations ordered
33 according to the sequence of the frames corresponding to consecutive time intervals
34 determined by their respective weights. The k-means clustering procedure determined
35 that each of the extracted 20 lipid basis configurations can be separated to two clusters:
36 The first cluster contains the primary lipids with a low average probability to have a DOPC
37 lipid-neighbor and the second cluster contains the primary lipids with more than four times
38 higher average probability to have a DOPC lipid-neighbor. Further, we colored differently
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 the lipids in each of the 20 lipid basis configurations, with two clusters each, at the time
4 intervals where the respective lipid basis configuration is active. The color gradient in **Fig.**
5
6 **4C** captures these two groups of primary lipids: white to grey the primary lipids within the
7
8 nucleation domains, and blue to red-the primary lipids at the interface or outside of any
9
10 nucleation domain.
11
12
13
14
15
16
17

18 Next, we use the MD simulations trajectories to visualize and rationalize the two
19
20 primary lipid groups as extracted by the clustering of basis lipid configurations derived by
21
22 NMFk. We considered lipid basis configurations, #19 and #1 corresponding to 1us and
23
24 2us time points, respectively. All lipids contributing to those lipid basis configurations are
25
26 mapped into the trajectories at the corresponding time points in **Fig. 5**. At 2 us,
27
28 approximately 70% of primary lipids (DPPC, colored in blue) are situated in large
29
30 nucleation domains which are well packed and condensed by the high cholesterol
31
32 concentration, leading to L_d phase. These primary lipids are predominantly shielded from
33
34 the secondary (DOPC, colored red) lipids which themselves are localized to form the L_o
35
36 phase. These primary lipids correspond to the first cluster identified by NMFk and are
37
38 localized in nucleation domains.
39
40
41
42
43



1
2
3 **Figure 5.** Visual inspection of simulations trajectories for evaluation of NMFk
4 categorization primary lipids according to their localization. Primary (DPPC) lipids are
5 colored according to NMFk output as purple (in nucleation domains) and yellow (out of
6 nucleation domains) in MD configurations corresponding to two time points. Secondary
7 (DOPC) lipids are colored in red and cholesterol in white. Insets highlight a particular
8 region where lipids can be differentiated as nucleating (purple) and boundary (yellow)
9 lipids. As a comparison, the same region is represented using the conventional way of
10 addressing the distinction between saturated (green), unsaturated (red) lipids, similar to
11 Fig. 1.
12
13
14
15
16
17
18
19
20
21
22
23

24 On the other hand, at 2 us, approximately 30% of the primary lipids (DPPC, colored
25 yellow) in the same lipid basis configuration #1 are located near the interfacial regions or
26 deep into the L_d regions formed by secondary lipid types. Unlike the previous set of
27 primary lipids, these lipids are in direct contact with the secondary lipids. The **Fig. 5 insert**
28 shows how NMFk is able to distinguish these lipids from all DPPC lipids available. They
29 correspond to the second cluster identified by NMFk. Thus, each given lipid basis
30 configuration contains physical and easily interpretable features that enable us to make
31 a distinction on primary lipids depending on their location and their contribution to the lipid
32 segregation.
33
34
35
36
37
38
39
40
41
42
43
44
45
46

47 The strength of NMFk analysis is the ability to extract lipid basis configurations as
48 a function of time. These configurations enable us to determine the time dependence of
49 the latent features connected to the kinetics of phase separation without carrying out
50 tedious multiple analyses. A conventional analysis that seeks to probe the temporal profile
51
52
53
54
55
56
57
58
59
60

1
2
3 would have considered the normalized number of contacts among DPPC lipids and
4
5 between DPPC and DOPC to capture the nucleation process (**Fig. 2C**). NMFk
6
7 decomposes the nucleation process in two components, as shown in **Fig. 6**, where the
8
9 distinction is made on the temporal profiles of primary lipids from the nucleation domains
10
11 from those that are still outside the nucleation domains. In **Fig. 6A**, the purple bars
12
13 represent the total number of primary lipids in nucleation domains at consecutive time
14
15 intervals (ordered on the x-axis), while the primary lipids outside of any nucleation domain
16
17 are presented by yellow bars. At early time, rapid growth of domains is seen which is
18
19 directly correlated to the increase of nucleation of primary lipids. After this rapid growth,
20
21 a steadier behavior is observed which continues till the end of the simulation. In **Fig. 6B**
22
23 we represent the normalized number of contacts with the secondary lipids for the primary
24
25 lipids as identified from NMFk. At early times, these contacts are high due to random
26
27 encounters between lipids. This behavior begins to change around 0.6 μ s and then the
28
29 contacts between primary and secondary lipids are indicative of steady growth of the
30
31 number of primary lipids in the nucleation domains and a decrease of the total number of
32
33 lipids in the non-nucleating regions as the lipid mixture system phase separates.
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

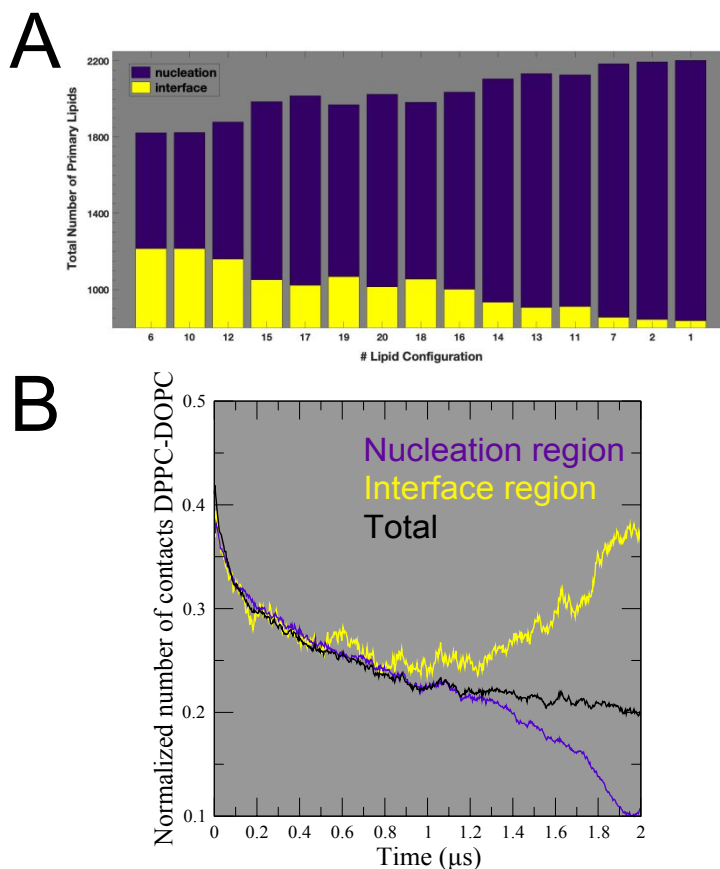


Figure 6. Time dependent variation of the total number of nucleating primary (DPPC) lipids. **(A)** The purple bars represent the total number of primary lipids in nucleation domains, as concluded from their membership in different clusters of the corresponding basis lipid configuration, while the primary lipids outside of any nucleation domain are presented by the yellow bars. The labels on x-axis (the numbers) correspond to consecutive (in time) processes extracted from 0.45 μ s to 2 μ s simulation time. **(B)** The same distinction of nucleating primary lipids as in **(A)** obtained using the normalized number of contacts of the lipids identified by NMFk. The same color scheme is used. Black line corresponds to the normalized number of contacts using the total fraction of saturated lipids (i.e., without making the distinction within the nucleating primary lipids).

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Furthermore, NMFk analysis extracts the temporal evolution of the primary lipids' *membership*: in or out of a nucleation domain. Specifically, the membership is defined by identifying primary lipids that join a nucleation domain and remain in that domain till the phase separation. We ordered the 15 significant basis lipid configurations extracted by NMFk according to the time intervals when the specific configurations are active, **Fig. 7**. We identify the primary lipids that participate in the nucleation by keeping track of the lipids in basis lipid configurations with low probabilities to have a secondary lipid-neighbor. In **Fig. 7A**, we present the *concrete* primary lipids that *remain* in the same cluster with small (purple color) or high (yellow color) probability to have a secondary lipid-neighbor at consecutive time intervals, which represents the evolution of nucleation. The inset in **Fig. 7A** demonstrates the system at much later time (~ 20 microseconds) after the initial nucleation when the phase separation has reached equilibrium and the primary lipids that are located in nucleation domains mostly preserve their membership in time. Importantly, this evolution of the primary lipids can be easily mapped at their spatial coordinates. In **Fig. 7B**, the patterns extracted by NMFk are mapped via MD trajectories to their spatial coordinates (at each time interval), to visualize the evolution of the different groups of primary lipids based on their membership: in nucleation domains (purple color) or outside those domains (yellow color). This again highlights the power of the NMFk to extract physical properties and details that can be easily visually tracked as the system undergoes localization and lipid segregation.

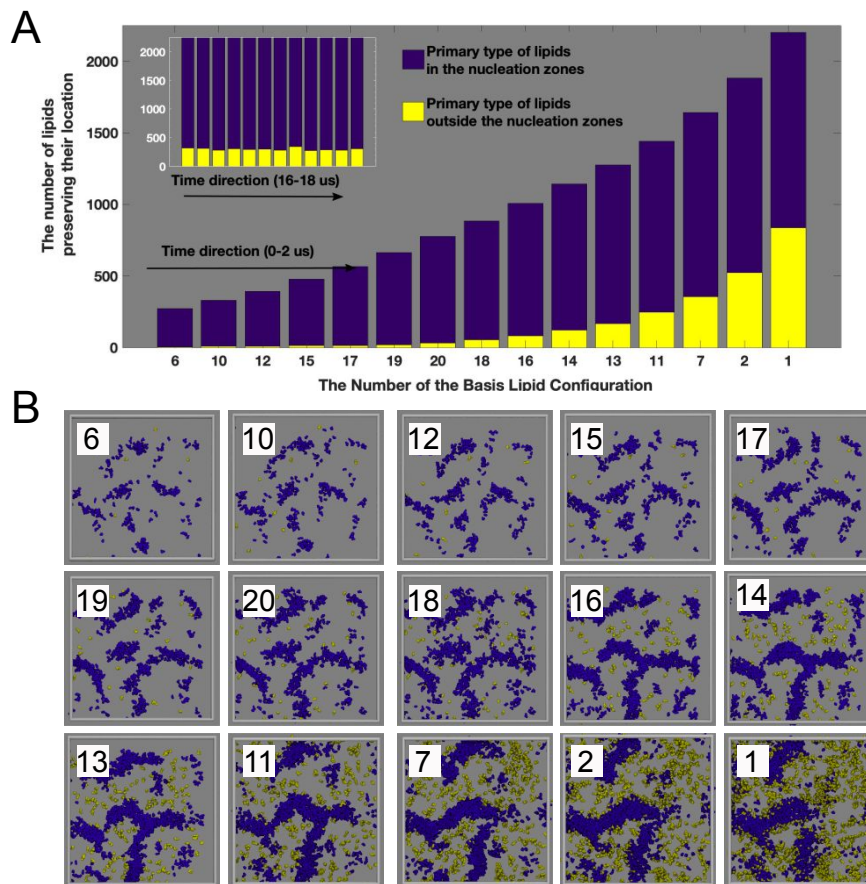


Figure 7. Temporal description of membership profile of primary lipids as captured by NMFk analysis. **(A)** The clear growth of nucleation domains with time, formed by primary lipids, is extracted from the clusters containing primary lipids with low probabilities to have lipid-neighbors of secondary type (refer **Fig. 4C**). The number of primary lipids that belong to nucleation domains are colored in purple whereas those appear at the interface (edges of the nucleation domain) or outside the nucleation domain are colored in yellow. The exponential growth of the number of the primary lipids that continue to be in the same cluster demonstrated the evolution of the steady state of the phase separation. The inset demonstrates the same process but at much later time (~ 20 microseconds) when the phase separation is in equilibrium. Here, although the minor changes still exist, primary

1
2
3 lipid membership numbers have stabilized. **(B)** Spatial visualization of the primary lipid
4 memberships extracted by NMFk analysis (between 0.45 us and 2 us). These primary
5 lipids were tracked using MD trajectories and rendered with the same color scheme as in
6 **(A)** to distinguish the primary lipids in and out of nucleation domains. Other lipids are not
7 rendered (silver background). The MD simulation box is represented as solid gray lines.

16 DISCUSSION

17
18
19
20 We introduce an unsupervised machine learning algorithm based on the nonnegative
21 matrix factorization combined with custom clustering, called NMFk, for analysis of MD
22 simulations. Specifically, we implement that algorithm here to detect and describe the
23 lateral lipid segregation in a simplistic lipid “raft” model composed of a well-characterized
24 ternary lipid mixture. Based on this study, we believe that the NMFk formalism can be
25 also implemented for extracting relevant features from a more complex biological
26 membrane.

27
28
29
30
31
32
33
34
35
36
37
38 The DOPC:DPPC:CHOL ternary lipid mixture considered here can exist as a
39 homogeneously mixed mixture or exhibit two distinct phases of L_o and L_d depending on
40 temperature. NMFk is utilized to detect and analyze this phase separation behavior of this
41 well-studied system. The distinction between the two phases, L_o and L_d , is sensitive to
42 the spatial localization of DPPC lipids and the number of DOPC neighbors. We designate
43 DPPC and DOPC as primary and secondary lipids in the NMFk formalism, respectively.
44
45 At the beginning of the simulations, there are no distinct phases and the lipid mixture is
46 homogeneous. As simulation time progresses, lipids begin to segregate. By 0.6 us, the
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 primary lipid type in this case, DPPC, begins to segregate into the L_o regions, while the
4
5 secondary lipid type, DOPC, begin to segregate into the L_d domains. Thus, the number
6
7 of secondary lipid neighbors around a given primary lipid should reflect that phase
8
9 separation.
10

11
12
13
14
15 Given that a neighborhood profile of lipids can track the DOPC:DPPC:CHOL
16
17 mixture, we first built a specific time-dependent data-matrix, $X_L(t)$ (see Methods section)
18
19 whose elements represent the number of DOPC neighbors to each one of the DPPC
20
21 lipids in the system, within a specific radius, r_{cutoff} extracted from careful analysis of the
22
23 radial distribution function. NMFk decomposed the matrix $X_L(t)$ into a product of two
24
25 matrices: (i) the matrix of the basis lipid configurations, W_{NK} , whose columns present the
26
27 configurations of DPPC lipids. NMFk determines the number of basis configuration K
28
29 based on the robustness of the decomposition. Each one of the K basis lipid
30
31 configurations in W_{NK} contains the probabilities of the set of the DPPC lipids to have
32
33 DOPC neighbors. The Silhouette-Reconstruction criterium was used to estimate the
34
35 optimal number of basis lipid configurations to be 20. K-means clustering of each these
36
37 20 basis lipid configurations is demonstrating the tendency of increasing the number of
38
39 primary lipids with a neglecting probability to have a neighbor of type DOPC when the
40
41 time advances, which correspond to an increased total number of DPPC lipids located in
42
43 nucleation domains.
44
45
46
47
48
49
50
51

52 By relating the basis lipid configurations to MD trajectories, we were able to show
53
54 details of phase separation extracted by NMFk analysis. The NMFk discriminates lipids,
55
56
57
58
59
60

1
2
3 depending on whether they belong to L_o or L_d phases or interfacial regions, as they
4 undergo phase separation. Unlike, other analyses, basis lipid configurations provide
5 details of lipids that take part in the nucleation versus those that establish line tension. It
6 tracks the complicated features of the lipid segregation process leading to L_o and L_d
7 phases. We identify lipids within the boundary of L_o phase with lipid configurations basis
8 corresponding to a signature of L_o domain where DPPC is well packed and condensed
9 by the high cholesterol concentration. Separately, another basis lipid configuration
10 captures interfacial lipids that shield the DOPC lipids from such L_o domains during phase
11 separation. Importantly, we demonstrate that the evolution of the nucleation process is
12 captured in terms of lipid membership to different basis lipid configuration active in
13 consecutive time intervals. NMFk identifies the lipids that take part in the initial nucleation
14 and remains as part of the domain towards the phase separation. Also, other lipids that
15 join such nucleation and remain till the phase separation are identified as the time
16 progresses.

39 CONCLUSIONS

40
41
42
43 The high variability and complexity of plasma membranes is still poorly understood.
44 Higher resolution spectroscopy in combination with atomic detailed computer simulations
45 are providing new insights, however, we are not yet close to fully understand or describe
46 the membrane processes regulating cellular function. A detailed description of such
47 complex systems undoubtedly requires novel mathematical frameworks that are able to

1
2
3 decompose and categorize the evolution of thousands if not millions of lipids involved in
4
5 the phenomenon.
6
7
8
9

10 Here we show the power of NMFk formalism on analyzing lipid phase separation
11 and providing a robust analysis of categorizing lipids according to their localization in the
12 membrane and elucidating the time-dependency along the nucleation process. The NMFk
13 discriminates all different types of lipids, part of L_o or L_d or interface, due to their particular
14 behavior along the trajectory and the resulting probability to have a DOPC neighbor. If
15 there is no clear pattern in the behavior of the lipids, for example, when MD simulations
16 do not produce any distinct behavior associated with phase separation, NMFk analysis
17 does not produce false features. This is the first demonstration of NMFk serving as a
18 useful tool in detecting time-dependent domain formation and lipid separation in MD
19 simulations of complex lipid mixture systems. Even though we have exhibited the
20 usefulness of NMFk in the context of a well-studied ternary lipid mixture, an extension to
21 more complicated lipid mixtures is feasible with a tensor formalism and is currently under
22 consideration.
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43

44 MATERIALS AND METHODS

45 **Membrane patch**

46
47
48
49
50 An initial configuration of a CG membrane patch was obtained using the script tools
51 provided in the Martini force field website (<http://cgmartini.nl/>). Our CG lipid system
52
53
54
55
56
57
58
59
60

1
2
3 contains DPPC:DOPC:CHOL lipids in a 37:36:27 ratio respectively, which initially were
4 randomly placed within a XY plane. This lipid ratio has been experimentally and
5 computationally observed to transitioning towards a phase separated Liquid-
6 ordered/Liquid-disordered state [54](#), [56](#). The lipids were represented using the Martini V2.2
7 force field [39](#) with a refined set of parameters, which has shown an improved phase
8 separation behavior [52](#). Similarly, simulation with the standard Martini lipid model was also
9 carried on. The total system was composed of 16366 lipids, 718830 Martini water beads
10 (175 atomistic water molecules per lipid) and 150mM NaCl to preserve an overall constant
11 ionic strength. In order to avoid spontaneous freezing of the Martini water beads (a well-
12 known artifact previously reported in the original model [39](#)), 0.1% M of the water beads
13 were replaced by anti-freeze particles.
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30

31 **Molecular dynamics protocol**

32
33 We followed a current update in parameters set-up for performing the CG simulations [57](#).
34 The equations of motion were integrated every 30fs time-step. A reaction-field
35 electrostatics algorithm was used with a Coulomb cut-off of 1.1 nm and dielectric
36 constants of 15 or 0 within or beyond this cut-off, respectively. Lennard-Jones interactions
37 were cut off at 1.1 nm, where the potential was shifted to zero. In order to accelerate the
38 lipid phase de-mixing, constant temperature was maintained at 290 K via separate
39 coupling of the solvent (water and ions) and membrane components using a velocity-
40 rescaling thermostat [58](#) with relaxation times of 1.0ps. During equilibration, the Martini
41 beads representing the phosphate groups of the lipid head regions were positionally (xyz
42 components) restraint in order to preserve the initial random positions. In this stage, the
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 solvent molecules (water and ions) were allowed to diffuse and the box pressure was
4 maintained semi-isotropically coupled at 1 bar using the Berendsen barostat [59](#) with
5 relaxation times of 12 ps and compressibilities of $3 \times 10^{-4} \text{ bar}^{-1}$. After that, production runs
6 were performed using a Parrinello-Rahman barostats [60](#). Simulations were run for 20 μs
7 using the GROMACS version 5.2.1 [61](#) and the trajectories were saved every 3 ns providing
8 the frames for the construction of the NMFk matrix (see later).
9
10
11
12
13
14
15
16
17
18

19 **Generation of the contact matrix $X_N(t)$**

20 Every frame stored within 2 μs (667 frames in total) were used for generating the
21 corresponding matrix for NMFk analysis. We rely on the implemented GROMACS tool
22 *gmx select* to output the number of DOPC lipids around every DPPC molecule within 1.1
23 nm. This cutoff-radius structurally corresponds to the second layer of neighbors, as
24 estimated by the second maximum peak of the radial distribution function $g(r)$ (**Fig. 3**).
25 Thus, each column of the data-matrix, $X_N(t)$, corresponds to a variable number of the
26 DOPC neighbors of a given DPPC lipid per frame, while the rows correspond to the
27 number of the 3038 DPPC lipids in the system. Similarly, matrix reconstruction was
28 carried out for 20 μs collection. An example of the matrix can be found as part of
29 Supporting material.
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

47 **Nonnegative Matrix Factorization and NMFk**

48 Nonnegative Matrix Factorization (NMF) is a well-known unsupervised machine learning
49 method created for parts-based representation [36](#). NMF has been successfully leveraged
50 for decomposing mixtures of various types nonnegative signals, i.e., for Blind Source
51
52
53
54
55
56
57
58
59
60

1
2
3 Separation (BSS) [38](#) problems. If the BSS problem is solved in a temporally discretized
4 framework, the goal of the NMF algorithm is to retrieve the original nonnegative signals
5 (sources), \mathbf{W} ; $\mathbf{W} \in M_{PK}(R_+)$ that produced the observational records, \mathbf{X} ; $\mathbf{X} \in M_{NP}(R_+)$,
6 detected at a set of sensors. Here R_+ denotes the set of real nonnegative numbers, N is
7 the number of the recording sensors, K is the number of unknown original signals, and P
8 is the number of discretized moments in time (time points or frames) at which the signals
9 are recorded at the sensors. Only the matrix \mathbf{X} is known initially. Thus, in a BSS problem,
10 the recorded data, \mathbf{X} , is formed by a linear mixing of K unknown original signals \mathbf{W} ,
11 blended by an unknown mixing matrix, \mathbf{H} . Since both factor matrices \mathbf{W} and \mathbf{H} are
12 unknown, and even their size K (i.e., the number of unknown original signals) is unknown
13 the problem is typically under-determined. NMF can solve such kind of problems by
14 leveraging, for example, the multiplicative update algorithm [37](#) to minimize the Frobenius
15 norm $\frac{1}{2} \|\mathbf{X} - \mathbf{W} * \mathbf{H}\|_F^2$. An additional advantage of NMF method is that it can work with
16 data in which the original signals are not independent but partially correlated [38](#).

17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39 One of the difficulties of the NMF algorithm is that it requires *prior* knowledge of K -
40 the number of the unknown original signals. Recently a new protocol called NMFk
41 addressing this limitation has been reported [62-64](#). This protocol complements classical
42 NMF with custom k-means clustering and Silhouette [65](#) statistics, which allows
43 simultaneous identification of the optimal number of the unknown basis patterns. The
44 NMFk was utilized to successfully decompose the largest available dataset of human
45 cancer [66](#) genomes, as well as for extraction of physical pressure transients [64](#) and
46 contaminants [67](#) originating from an unknown number of sources that may propagate with
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 a finite speed in nondispersive [68](#) or dispersive media [69](#) as well as for extraction of the
4 original crystal structures and phase diagram from X-ray spectra of material combinatorial
5 libraries [70](#).
6
7
8
9

10
11
12 NMFk determines the number of the unknown original signals based on the
13 robustness and reproducibility of the NMF solution. Specifically, it explores consecutively
14 the possible numbers of configurations \tilde{K} (\tilde{K} can go from 1 to N-1, where N is the total
15 number of frames), by obtaining sets of NMF minimization solutions for each \tilde{K} . Note that
16 \tilde{K} serves to index the different NMF models, and it is distinct from K , which is fixed, albeit
17 unknown number. Further, NMFk leverages a custom clustering using the cosine
18 similarity, in order to estimate the robustness of each set of NMF solutions with fixed \tilde{K}
19 but derived with different initial guesses. Comparing the quality of the derived clusters (a
20 measure how different are the extracted signals) and the accuracy of minimization among
21 the sets with various \tilde{K} , which we call a Silhouette-Reconstruction criterium, NMFk
22 determines the optimal numbers of the unknown original signals. To access the quality of
23 the clusters obtained for each set we use their average Silhouette width, S . NMFk utilizes
24 S to measure how good is a particular choice of \tilde{K} as an estimate for K . Specifically, the
25 optimal number of patterns is picked by selecting the value of \tilde{K} that leads to both: (a) an
26 acceptable reconstruction error R of the observation matrix X , where
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

$$R = \frac{\|X - W * H\|_F}{\|X\|_F},$$

47
48
49
50 and (b) a large average Silhouette width (i.e., an average Silhouette width close to one).
51
52 The combination of these two criteria is easy to understand intuitively. For solutions with
53 \tilde{K} less than the actual number of patterns ($\tilde{K} < K$) we expect the clustering to be good
54
55
56
57
58
59
60

1
2
3 (with an average Silhouette width close to 1), because several of the actual patterns could
4 be combined to produce one “super-cluster”; however, the reconstruction error will be
5 high, due to the model being too constrained (with too few degrees of freedom), and thus
6 on the under-fitting side. In the opposite limit of over-fitting, when $\tilde{K} > K$ (\tilde{K} exceeds the
7 actual number of configurations), the average reconstruction error could be quite small -
8 each solution reconstructs the observation matrix very well - but the solutions will not be
9 well-clustered (e.g., with an average Silhouette substantially less than 0.8), since there is
10 no unique way to reconstruct X with more than the actual number of configurations, and
11 no well-separated clusters will be formed.

12
13 Thus, the best estimate for the number of unknown original signals, W ,
14 corresponding to the true number of unknown original signals K , is given by the value of
15 \tilde{K} that optimizes both of these metrics simultaneously. Finally, after determining K , we
16 use the centroids of the K clusters as a final robust representation of the original signals.

17 **NMFk minimization algorithm**

18 Here we leveraged the multiplicative algorithm [37](#) based on Kullback–Leibler divergence
19 [71](#) as well as the block coordinate descent algorithm [72](#) based on Frobenius norm. We did
20 not observe any significant difference between the results obtained via these two
21 algorithms.

22 **NMFk clustering algorithm**

23 NMFk creates up to $N-1$ sets of minimizations (called NMF runs), one for each possible
24 number \tilde{K} of original patterns. In each of these runs, Q solutions (e.g., $Q = 200$) of the
25 NMF minimizations for a fixed number of patterns \tilde{K} are derived. Thus, each run results
26 in a set of solutions $U_{\tilde{K}}$:

$$U_{\tilde{K}} = \{W_{\tilde{K}}^1, H_{\tilde{K}}^1; W_{\tilde{K}}^2, H_{\tilde{K}}^2; \dots; W_{\tilde{K}}^Q, H_{\tilde{K}}^Q\},$$

where each of these “tuples” represents a distinct solution for the nominally same NMF minimization: the difference is stemming from the different (random) initial guesses. Next, NMFk performs custom clustering, assigning the \tilde{K} columns/features of each $W_{\tilde{K}}^i$ of all Q solutions to one of the \tilde{K} clusters, representing \tilde{K} basic patterns. This custom clustering is similar to k-means clustering but with an additional constraint which holds the number of elements in each of the clusters equal to the number of solutions Q . For example, with $Q = 200$ each one of the \tilde{K} identified clusters must contain exactly 200 solutions. This condition has to be enforced since each minimization (specified by a given $(W_{\tilde{K}}^i, H_{\tilde{K}}^i)$ tuple) contributes only one solution for each feature, and accordingly supplies exactly one element to each cluster. During the clustering, the similarity between patterns is measured using the cosine similarity.

Numerical codes, data and Supporting Information

The following files: COORD.gro with the coordinates; INPUT.mdp with the parameters set to start the simulation with Gromacs; PARAMETERS.top with the Martini force field-based topology for the membrane system; DPPC-DOPC-1.1nm.xvg with an example matrix generated from the MD trajectories used for NMFk calculation, are also provided as supporting information accompanied this paper.

To extract patterns of the basis lipid configurations from the MD simulations we used the NMFk method which is an extension of the original NMF [37](#) that includes a custom clustering for determination of the number of patterns [66](#). Our NMFk protocol is based on the SigProfile software created for identification of mutational signatures in human cancer

[66](#).

1
2
3 The SigProfile code is freely available at: <https://www.mathworks.com/matlabcentral>. To
4 use SigProfile, an input file should be at place. In our case, the input file is the contact
5 matrix $X_N(\mathbf{t})$ with size ($\mathbf{N} \times \mathbf{M}$), where \mathbf{N} is the number of the lipids in the MD simulations
6 and \mathbf{M} is number of the frames. A detailed description of NMFk is available elsewhere [64](#).
7
8 The input data-file, containing the contact matrix $X_N(\mathbf{t})$ as well as a script, needed to run
9 the SigProfile, are provided as supplemented information accompanied this paper.
10
11

12 A README.docx file with a list of included files, links to publicly available repositories
13 along with their brief description and instructions is also provided as supporting
14 information.
15
16

17 The simulated MD-data containing the lipids' trajectories (~ 100GB) is available
18 freely but because of its size-upon request to gnana@lanl.gov.
19
20
21
22
23

24 ACKNOWLEDGMENTS

25
26
27 This research was performed at Los Alamos National Laboratory and carried out under
28 the auspices of the National Nuclear Security Administration of the United States
29 Department of Energy. We like to acknowledge Tyler Reddy for generating synthetic data
30 to test preliminary results. The work was supported by LANL LDRD grant 20180060DR.
31
32 This work has been supported in part by the Joint Design of Advanced Computing
33 Solutions for Cancer (JDACS4C) program established by the U.S. Department of Energy
34 (DOE) and the National Cancer Institute (NCI) of the National Institutes of Health. This
35 work was performed under the auspices of the U.S. Department of Energy by Lawrence
36 Livermore National Laboratory under Contract DE-AC52-07NA27344, Los Alamos
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 National Laboratory under Contract DE-AC5206NA25396, Oak Ridge National
4
5 Laboratory under Contract DE-AC05-00OR22725, and Frederick National Laboratory for
6
7 Cancer Research under Contract HHSN261200800001E. We thank the LANL
8
9 Institutional Computing for the computing resources.
10
11
12
13
14
15

16 CONTRIBUTIONS

17
18
19

20 All authors designed the research. C. A. L. performed the MD simulations. B. S. A.
21
22 performed the NMFk analyses. All authors analyzed the data and wrote the manuscript.
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Table 1. The quality of reconstruction of X, estimated by the mean Pearson correlation between the columns of X and columns of W*H.

Coarse-grained simulations of ternary lipid mixture	Time [us]	Mean Correlation Coefficient	Standard Deviation
Standard Martini	0.45 - 2	0.54	0.040
Updated Martini	0.45 - 2	0.90	0.020
	16 - 18	0.96	0.004

REFERENCES

1. Israelachvili, J.; Wennerström, H., Role of hydration and water structure in biological and colloidal interactions. *Nature* **1996**, *379* (6562), 219.
2. Nickels, J. D.; Chatterjee, S.; Stanley, C. B.; Qian, S.; Cheng, X.; Myles, D. A. A.; Standaert, R. F.; Elkins, J. G.; Katsaras, J., The in vivo structure of biological membranes and evidence for lipid domains. *PLoS biology* **2017**, *15* (5), e2002214.
3. Kraft, M. L., Sphingolipid Organization in the Plasma Membrane and the Mechanisms That Influence It. *Frontiers in cell and developmental biology* **2016**, *4*, 154.
4. Rosetti, C. M.; Mangiarotti, A.; Wilke, N., Sizes of lipid domains: What do we know from artificial lipid membranes? What are the possible shared features with membrane rafts in cells? *Biochimica et biophysica acta* **2017**, *1859* (5), 789-802.
5. Levental, K. R.; Lorent, J. H.; Lin, X.; Skinkle, A. D.; Surma, M. A.; Stockenbojer, E. A.; Gorfe, A. A.; Levental, I., Polyunsaturated Lipids Regulate Membrane Domain Stability by

1
2
3 Tuning Membrane Order. *Biophysical journal* **2016**, *110* (8), 1800-
4 1810.

5
6 6. Dietrich, C.; Volovyk, Z. N.; Levi, M.; Thompson, N. L.;
7 Jacobson, K., Partitioning of Thy-1, GM1, and cross-linked
8 phospholipid analogs into lipid rafts reconstituted in supported
9 model membrane monolayers. *Proceedings of the National*
10 *Academy of Sciences* **2001**, *98* (19), 10642-10647.

11
12 7. Van Meer, G.; Voelker, D. R.; Feigenson, G. W., Membrane
13 lipids: where they are and how they behave. *Nature reviews*
14 *Molecular cell biology* **2008**, *9* (2), 112.

15
16 8. Veatch, S. L.; Keller, S. L., Separation of liquid phases in giant
17 vesicles of ternary mixtures of phospholipids and cholesterol.
18 *Biophysical journal* **2003**, *85* (5), 3074-83.

19
20 9. Risselada, H. J.; Marrink, S. J., The molecular face of lipid
21 rafts in model membranes. *Proceedings of the National Academy*
22 *of Sciences of the United States of America* **2008**, *105* (45), 17367-
23 72.

24
25 10. Ingólfsson, H. I.; Melo, M. N.; Van Eerden, F. J.; Arnarez, C.
26 m.; Lopez, C. A.; Wassenaar, T. A.; Periole, X.; De Vries, A. H.;
27 Tieleman, D. P.; Marrink, S. J., Lipid organization of the plasma
28 membrane. *Journal of the american chemical society* **2014**, *136*
29 (41), 14554-14559.

30
31 11. Michalski, R. S.; Carbonell, J. G.; Mitchell, T. M., *Machine*
32 *learning: An artificial intelligence approach*. Springer Science &
33 Business Media: 2013.

34
35 12. Broecker, P.; Carrasquilla, J.; Melko, R. G.; Trebst, S.,
36 Machine learning quantum phases of matter beyond the fermion
37 sign problem. *Scientific reports* **2017**, *7* (1), 8823.

38
39 13. Tanaka, A.; Tomiya, A., Detection of phase transition via
40 convolutional neural networks. *Journal of the Physical Society of*
41 *Japan* **2017**, *86* (6), 063001.

42
43 14. Zhang, Y.; Kim, E.-A., Quantum loop topography for machine
44 learning. *Physical review letters* **2017**, *118* (21), 216401.

45
46 15. Wang, L., Discovering phase transitions with unsupervised
47 learning. *Physical Review B* **2016**, *94* (19), 195105.

16. Carrasquilla, J.; Melko, R. G., Machine learning phases of matter. *Nature Physics* **2017**, *13* (5), 431.
17. Zhang, P.; Shen, H.; Zhai, H., Machine learning topological invariants with neural networks. *Physical review letters* **2018**, *120* (6), 066401.
18. Hu, W.; Singh, R. R.; Scalettar, R. T., Discovering phases, phase transitions, and crossovers through unsupervised machine learning: A critical examination. *Physical Review E* **2017**, *95* (6), 062122.
19. Wetzel, S. J., Unsupervised learning of phase transitions: From principal component analysis to variational autoencoders. *Physical Review E* **2017**, *96* (2), 022140.
20. Riniker, S., Molecular Dynamics Fingerprints (MDFP): Machine Learning from MD Data To Predict Free-Energy Differences. *Journal of chemical information and modeling* **2017**, *57* (4), 726-741.
21. Garcia, A. E., Large-amplitude nonlinear motions in proteins. *Phys Rev Lett* **1992**, *68* (17), 2696-2699.
22. Coifman R. R.; Lafon S.; Lee A. B.; Maggioni M. ; Nadler B. ; Warner F.; W., Z. S., Geometric diffusions as a tool for harmonic analysis and structure definition of data: Diffusion maps. *Proceedings of the National Academy of Sciences of the United States of America* **2005**, *102.21* 7426-7431.
23. Tenenbaum, J. B.; de Silva, V.; Langford, J. C., A global geometric framework for nonlinear dimensionality reduction. *Science* **2000**, *290* (5500), 2319-23.
24. Glazer, D. S.; Radmer, R. J.; Altman, R. B., Combining molecular dynamics and machine learning to improve protein function recognition. In *Biocomputing 2008*, World Scientific: 2008; pp 332-343.
25. Karamzadeh, R.; Karimi-Jafari, M. H.; Sharifi-Zarchi, A.; Chitsaz, H.; Salekdeh, G. H.; Moosavi-Movahedi, A. A., Machine Learning and Network Analysis of Molecular Dynamics Trajectories Reveal Two Chains of Red/Ox-specific Residue Interactions in

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Human Protein Disulfide Isomerase. *Scientific Reports* **2017**, 7 (1), 3666.

26. Gastegger, M.; Behler, J.; Marquetand, P., Machine learning molecular dynamics for the simulation of infrared spectra. *Chemical science* **2017**, 8 (10), 6924-6935.

27. Ash, J.; Fourches, D., Characterizing the chemical space of ERK2 kinase inhibitors using descriptors computed from molecular dynamics trajectories. *Journal of chemical information and modeling* **2017**, 57 (6), 1286-1299.

28. Sgourakis, N. G.; Merced-Serrano, M.; Boutsidis, C.; Drineas, P.; Du, Z.; Wang, C.; Garcia, A. E., Atomic-Level Characterization of the Ensemble of the A β (1–42) Monomer in Water Using Unbiased Molecular Dynamics Simulations and Spectral Algorithms. *Journal of Molecular Biology* **2011**, 405 (2), 570-583.

29. Hartigan, J. A.; Wong, M. A., Algorithm AS 136: A k-means clustering algorithm. *Journal of the Royal Statistical Society. Series C (Applied Statistics)* **1979**, 28 (1), 100-108.

30. Kohonen, T., The self-organizing map. *Neurocomputing* **1998**, 21 (1), 1-6.

31. Belouchrani, A.; Abed-Meraim, K.; Cardoso, J. F.; Moulines, E., A blind source separation technique using second-order statistics. *IEEE Transactions on signal processing* **1997**, 45 (2), 434-444.

32. Spearman, C., "General Intelligence," objectively determined and measured. *The American Journal of Psychology* **1904**, 15 (2), 201-292.

33. Jolliffe, I., *Principal component analysis*. Wiley Online Library: 2002.

34. Golub, G. H.; Reinsch, C., Singular value decomposition and least squares solutions. *Numerische mathematik* **1970**, 14 (5), 403-420.

35. Cichocki, A.; Yang, H. H., A new learning algorithm for blind signal separation. *Advances in neural information processing systems* **1996**, 8, 757-763.

- 1
- 2
- 3
- 4 36. Paatero, P.; Tapper, U., Positive matrix factorization: A
- 5 non-negative factor model with optimal utilization of error estimates
- 6 of data values. *Environmetrics* **1994**, 5 (2), 111-126.
- 7
- 8 37. Lee, D. D.; Seung, H. S., Learning the parts of objects by non-
- 9 negative matrix factorization. *Nature* **1999**, 401 (6755), 788-791.
- 10
- 11 38. Cichocki, A.; Zdunek, R.; Phan, A. H.; Amari, S.-i.,
- 12 *Nonnegative matrix and tensor factorizations: applications to*
- 13 *exploratory multi-way data analysis and blind source separation*.
- 14 John Wiley & Sons: 2009.
- 15
- 16 39. Marrink, S. J.; Risselada, H. J.; Yefimov, S.; Tieleman, D. P.;
- 17 de Vries, A. H., The MARTINI force field: coarse grained model for
- 18 biomolecular simulations. *The journal of physical chemistry. B*
- 19 **2007**, 111 (27), 7812-24.
- 20
- 21 40. Marrink, S. J.; Tieleman, D. P., Perspective on the Martini
- 22 model. *Chemical Society reviews* **2013**, 42 (16), 6801-22.
- 23
- 24 41. Vögele, M.; Köfinger, J.; Hummer, G., Hydrodynamics of
- 25 Diffusion in Lipid Membrane Simulations. *Physical Review Letters*
- 26 **2018**, 120 (26), 268104.
- 27
- 28 42. Domanski, J.; Marrink, S. J.; Schafer, L. V., Transmembrane
- 29 helices can induce domain formation in crowded model
- 30 membranes. *Biochimica et biophysica acta* **2012**, 1818 (4), 984-94.
- 31
- 32 43. Schafer, L. V.; de Jong, D. H.; Holt, A.; Rzepiela, A. J.; de
- 33 Vries, A. H.; Poolman, B.; Killian, J. A.; Marrink, S. J., Lipid packing
- 34 drives the segregation of transmembrane helices into disordered
- 35 lipid domains in model membranes. *Proceedings of the National*
- 36 *Academy of Sciences of the United States of America* **2011**, 108
- 37 (4), 1343-8.
- 38
- 39 44. Klingelhoefer, J. W.; Carpenter, T.; Sansom, M. S., Peptide
- 40 nanopores and lipid bilayers: interactions by coarse-grained
- 41 molecular-dynamics simulations. *Biophysical journal* **2009**, 96 (9),
- 42 3519-28.
- 43
- 44 45. Ingolfsson, H. I.; Melo, M. N.; van Eerden, F. J.; Arnarez, C.;
- 45 Lopez, C. A.; Wassenaar, T. A.; Periole, X.; de Vries, A. H.;
- 46 Tieleman, D. P.; Marrink, S. J., Lipid organization of the plasma
- 47 membrane. *J Am Chem Soc* **2014**, 136 (41), 14554-9.
- 48
- 49
- 50
- 51
- 52
- 53
- 54
- 55
- 56
- 57
- 58
- 59
- 60

- 1
2
3
4 46. Ingolfsson, H. I.; Carpenter, T. S.; Bhatia, H.; Bremer, P. T.;
5 Marrink, S. J.; Lightstone, F. C., Computational Lipidomics of the
6 Neuronal Plasma Membrane. *Biophysical journal* **2017**, *113* (10),
7 2271-2280.
8
9 47. Lopez, C. A.; Sovova, Z.; van Eerden, F. J.; de Vries, A. H.;
10 Marrink, S. J., Martini Force Field Parameters for Glycolipids.
11 *Journal of chemical theory and computation* **2013**, *9* (3), 1694-708.
12
13 48. Marrink, S. J.; de Vries, A. H.; Mark, A. E., Coarse Grained
14 Model for Semiquantitative Lipid Simulations. *The Journal of*
15 *Physical Chemistry B* **2004**, *108* (2), 750-760.
16
17 49. Khelashvili, G.; Kollmitzer, B.; Heftberger, P.; Pabst, G.;
18 Harries, D., Calculating the Bending Modulus for Multicomponent
19 Lipid Membranes in Different Thermodynamic Phases. *Journal of*
20 *chemical theory and computation* **2013**, *9* (9), 3866-3871.
21
22 50. Davis, R. S.; Sunil Kumar, P. B.; Sperotto, M. M.; Laradji, M.,
23 Predictions of phase separation in three-component lipid
24 membranes by the MARTINI force field. *The journal of physical*
25 *chemistry. B* **2013**, *117* (15), 4072-80.
26
27 51. Baoukina, S.; Mendez-Villuendas, E.; Tieleman, D. P.,
28 Molecular view of phase coexistence in lipid monolayers. *J Am*
29 *Chem Soc* **2012**, *134* (42), 17543-53.
30
31 52. Carpenter, T.; López, C. A.; Neale, C.; Montour, C.;
32 Ingólfsson, H.; Di Natale, F.; Lightstone, F.; Gnanakaran, S.,
33 Capturing Phase Behavior of Ternary Lipid Mixtures with a Refined
34 Martini Coarse-Grained Force Field. *Journal of chemical theory*
35 *and computation* **2018**, *accepted*.
36
37 53. Heberle, F. A.; Feigenson, G. W., Phase separation in lipid
38 membranes. *Cold Spring Harbor perspectives in biology* **2011**, *3*
39 (4).
40
41 54. Sodt, A. J.; Pastor, R. W.; Lyman, E., Hexagonal Substructure
42 and Hydrogen Bonding in Liquid-Ordered Phases Containing
43 Palmitoyl Sphingomyelin. *Biophysical journal* **2015**, *109* (5), 948-
44 55.
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

- 1
- 2
- 3
- 4 55. Simons, K.; Sampaio, J. L., Membrane organization and lipid
- 5 rafts. *Cold Spring Harbor perspectives in biology* **2011**, 3 (10),
- 6 a004697.
- 7
- 8 56. Garcia-Saez, A. J.; Chiantia, S.; Schwille, P., Effect of line
- 9 tension on the lateral organization of lipid membranes. *The Journal*
- 10 *of biological chemistry* **2007**, 282 (46), 33537-44.
- 11
- 12 57. de Jong, D. H.; Baoukina, S.; Ingólfsson, H. I.; Marrink, S. J.,
- 13 Martini straight: Boosting performance using a shorter cutoff and
- 14 GPUs. *Computer Physics Communications* **2016**, 199, 1-7.
- 15
- 16 58. Bussi, G.; Donadio, D.; Parrinello, M., Canonical sampling
- 17 through velocity rescaling. *The Journal of chemical physics* **2007**,
- 18 *126* (1), 014101.
- 19
- 20 59. Berendsen, H. J. C.; Postma, J. P. M.; van Gunsteren, W. F.;
- 21 DiNola, A.; Haak, J. R., Molecular dynamics with coupling to an
- 22 external bath. *The Journal of chemical physics* **1984**, 81 (8), 3684-
- 23 3690.
- 24
- 25 60. Parrinello, M.; Rahman, A., Polymorphic transitions in single
- 26 crystals: A new molecular dynamics method. *Journal of Applied*
- 27 *Physics* **1981**, 52 (12), 7182-7190.
- 28
- 29 61. Abraham, M. J.; Murtola, T.; Schulz, R.; Páll, S.; Smith, J. C.;
- 30 Hess, B.; Lindahl, E., GROMACS: High performance molecular
- 31 simulations through multi-level parallelism from laptops to
- 32 supercomputers. *SoftwareX* **2015**, 1-2, 19-25.
- 33
- 34 62. Alexandrov, L. B.; Nik-Zainal, S.; Wedge, D. C.; Campbell, P.
- 35 J.; Stratton, M. R., Deciphering signatures of mutational processes
- 36 operative in human cancer. *Cell reports* **2013**, 3 (1), 246-259.
- 37
- 38 63. Alexandrov, B. S.; Alexandrov, L. B.; Iliev, F. L.; Stanev, V. G.;
- 39 Vesselinov, V. V., Source identification by non-negative matrix
- 40 factorization combined with semi-supervised clustering. Google
- 41 Patents: 2018.
- 42
- 43 64. Alexandrov, B. S.; Vesselinov, V. V., Blind source separation
- 44 for groundwater pressure analysis based on nonnegative matrix
- 45 factorization. *Water Resources Research* **2014**, 50 (9), 7332-7347.
- 46
- 47
- 48
- 49
- 50
- 51
- 52
- 53
- 54
- 55
- 56
- 57
- 58
- 59
- 60

- 1
- 2
- 3
- 4 65. Rousseeuw, P. J., Silhouettes: a graphical aid to the
- 5 interpretation and validation of cluster analysis. *Journal of*
- 6 *computational and applied mathematics* **1987**, *20*, 53-65.
- 7
- 8 66. Alexandrov, L. B.; Nik-Zainal, S.; Wedge, D. C.; Aparicio, S.
- 9 A.; Behjati, S.; Biankin, A. V.; Bignell, G. R.; Bolli, N.; Borg, A.;
- 10 Børresen-Dale, A.-L., Signatures of mutational processes in
- 11 human cancer. *Nature* **2013**, *500* (7463), 415.
- 12
- 13 67. Vesselinov, V. V.; Alexandrov, B. S.; O'Malley, D.,
- 14 Contaminant source identification using semi-supervised machine
- 15 learning. *Journal of contaminant hydrology* **2017**.
- 16
- 17 68. Iliev, F. L.; Stanev, V. G.; Vesselinov, V. V.; Alexandrov, B. S.,
- 18 Nonnegative Matrix Factorization for identification of unknown
- 19 number of sources emitting delayed signals. *PloS one* **2018**, *13* (3),
- 20 e0193974.
- 21
- 22 69. Stanev, V. G.; Iliev, F. L.; Hansen, S.; Vesselinov, V. V.;
- 23 Alexandrov, B. S., Identification of release sources in advection-
- 24 diffusion system by machine learning combined with Green's
- 25 function inverse method. *Applied Mathematical Modelling* **2018**.
- 26
- 27 70. Stanev, V.; Vesselinov, V. V.; Kusne, A. G.; Antoszewski, G.;
- 28 Takeuchi, I.; Alexandrov, B. S., Unsupervised Phase Mapping of
- 29 X-ray Diffraction Data by Nonnegative Matrix Factorization
- 30 Integrated with Custom Clustering. *npj Computational Materials*
- 31 **2018**, 4:43.
- 32
- 33 71. Kullback, S.; Leibler, R. A., On information and sufficiency.
- 34 *The annals of mathematical statistics* **1951**, *22* (1), 79-86.
- 35
- 36 72. Xu, Y.; Yin, W., A block coordinate descent method for
- 37 regularized multiconvex optimization with applications to
- 38 nonnegative tensor factorization and completion. *SIAM Journal on*
- 39 *imaging sciences* **2013**, *6* (3), 1758-1789.
- 40
- 41
- 42
- 43
- 44
- 45
- 46
- 47
- 48
- 49
- 50
- 51
- 52
- 53
- 54
- 55
- 56
- 57
- 58
- 59
- 60