



HAL
open science

BitQT: a graph-based approach to the quality threshold clustering of molecular dynamics

Roy González-Alemán, Daniel Platero-Rochart, David Hernández-Castillo,
Erix Hernández-Rodríguez, Julio Caballero, Fabrice Leclerc, Luis
Montero-Cabrera

► To cite this version:

Roy González-Alemán, Daniel Platero-Rochart, David Hernández-Castillo, Erix Hernández-Rodríguez, Julio Caballero, et al.. BitQT: a graph-based approach to the quality threshold clustering of molecular dynamics. *Bioinformatics*, 2022, 38 (1), pp.73-79. 10.1093/bioinformatics/btab595 . hal-03767725

HAL Id: hal-03767725

<https://hal.science/hal-03767725>

Submitted on 23 Nov 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Structural Bioinformatics

BitQT: A Graph-Based Approach to the Quality Threshold Clustering of Molecular Dynamics

Roy González-Alemán^{1, 2,*}, Daniel Platero-Rochart¹, David Hernández-Castillo³, Erix W. Hernández-Rodríguez⁵, Julio Caballero⁴, Fabrice Leclerc^{2,*} and Luis Montero-Cabrera¹

¹Laboratorio de Química Computacional y Teórica (LQCT), Facultad de Química, Universidad de La Habana, La Habana, 10400, Cuba

²Institute for Integrative Biology of the Cell (I2BC), CEA, CNRS, Université Paris Saclay, Gif-sur-Yvette, F-91198, France

³Institute of Theoretical Chemistry, University of Vienna, Währinger Str. 17, 1090 Vienna, Austria

⁴Departamento de Bioinformática, Facultad de Ingeniería, Centro de Bioinformática, Simulación y Modelado (CBSM), Universidad de Talca, Talca, Chile, and

⁵Laboratorio de Bioinformática y Química Computacional, Escuela de Química y Farmacia, Facultad de Medicina, Universidad Católica del Maule, 3460000 Talca, Chile.

*To whom correspondence should be addressed.

Associate Editor: XXXXXXX

Received on XXXXX; revised on XXXXX; accepted on XXXXX

Abstract

Motivation: Classical Molecular Dynamics is a standard computational approach to model time-dependent processes at the atomic level. The inherent sparsity of increasingly huge generated trajectories demands clustering algorithms to reduce other post-simulation analysis complexity. The quality threshold (QT) variant is an appealing one from the vast number of available clustering methods. It guarantees that all members of a particular cluster will maintain a collective similarity established by a user-defined threshold. Unfortunately, its high computational cost for processing big data limits its application in the molecular simulation field.

Results: In the present work, we propose a methodological parallel between QT clustering and another well-known algorithm in the field of Graph Theory, the Maximum Clique Problem. Molecular trajectories are represented as graphs whose nodes designate conformations, while unweighted edges indicate mutual similarity between nodes. The use of a binary-encoded RMSD matrix coupled to the exploitation of bitwise operations to extract clusters significantly contributes to reaching a very affordable algorithm compared to the few implementations of QT for Molecular Dynamics available in the literature. Our alternative provides results in good agreement with the exact one while strictly preserving the collective similarity of clusters. [The source code and documentation of BitQT are free and publicly available on GitHub \(https://github.com/LQCT/BitQT.git\)](https://github.com/LQCT/BitQT.git) and [ReadTheDocs \(https://bitqt.readthedocs.io/en/latest/\)](https://bitqt.readthedocs.io/en/latest/) respectively.

Contact: roy_gonzalez@fq.uh.cu, fabrice.leclerc@i2bc.paris-saclay.fr

Supplementary information: Supplementary data are available at *Bioinformatics* online.

1 Introduction

Molecular dynamics (MD) is a powerful tool to gain insight into the conformational behavior of nanoscopic systems. Nowadays,

methodologies like coarse-grained MD, accelerated MD, and replica-exchange MD are common ways to reach a representative sampling of dynamically meaningful states. As the computational power grows, the size of trajectories generated by these techniques represents a massive amount of information that is potentially difficult to analyze. Geometrical clustering is a classical way to simplify those trajectories by grouping

1 similar conformations into sets known as clusters. In such a way, 62
2 conformations inside a cluster are more similar between them than those 63
3 from other clusters. 64

4 Many clustering algorithms exist for analyzing MD (Peng *et al.*, 65
5 2018), having benefits and shortcomings that make them suitable for 66
6 particular applications and inappropriate for others (Röttger, 2016). Due 67
7 to the inherent subjectivity associated with classification (the same 68
8 set of elements can be grouped according to many different criteria), 69
9 some authors consider clustering as an art (von Luxburg *et al.*, 2012). 70
10 However, in those particular cases where strongly geometrically correlated 71
11 conformations are needed to be returned as clusters, the Quality Threshold 72
12 (QT) algorithm (Heyer *et al.*, 1999) stands out as an ideal option. 73

13 QT appeared in the context of clustering gene expression patterns. 74
14 Since then, it has been applied to many areas other than microbiology 75
15 (Tang *et al.*, 2010; Yaakob *et al.*, 2010; Olson *et al.*, 2011; Dutta and 76
16 Overbye, 2011; Yaakob and Jain, 2012), including the MD field (Procacci 77
17 *et al.*, 1997; Danalis *et al.*, 2012). Two remarkable features of this algorithm 78
18 are the guarantee that no pair of frames having a similarity value greater 79
19 than a user-specified cutoff will coalesce into the same cluster and that the 80
20 number of clusters to retrieve must not be known *a priori*. However, QT 81
21 has an expensive computational cost (Danalis *et al.*, 2012) that currently 82
22 limits its applicability. 83

23 Several popular software have inaccurately qualified their clustering 84
24 implementations as QT or QT-like variants in the past (González-Alemán 85
25 *et al.*, 2020b). These pseudo-QT alternatives correspond to another simple 86
26 and largely disseminated algorithm (Daura *et al.*, 1999) that has been 87
27 recently optimized for the efficient treatment of long molecular trajectories 88
28 (González-Alemán *et al.*, 2020a). 89

29 After careful inspection of current literature, we have found only 90
30 two valid and ready-to-use attempts to implement QT to analyze MD 91
31 trajectories. The first one corresponds to the *qtcluster* command of 92
32 the ORAC suite (Procacci *et al.*, 1997) while the second one is an 93
33 implementation previously published by authors of this study (González- 94
34 Alemán *et al.*, 2020b), referred to as *QTPy* from now on. 95

35 While *QTPy* can be stated as an exact version of the QT proposed 96
36 by Heyer in 1999, it should be emphasized that *qtcluster* only partially 97
37 complies with the original algorithm. Perhaps the most essential feature 98
38 that makes *qtcluster* a fast QT implementation lies in the fact that it 99
39 is not an exact QT attempt, only preserving one condition from the 100
40 exact algorithm; the one assuring the collective similarity of retrieved 101
41 clusters. It is also worth noting that *qtcluster* uses the maximum difference 102
42 between corresponding pairs of atoms as the similarity measure while 103
43 *QTPy* employs the more customary optimal RMSD. Both of them are 104
44 marked by a run time and RAM consumption that impedes the processing 105
45 of relatively long trajectories. 106

46 Here we propose a heuristic variation of QT that can output equivalent 107
47 results to the exact algorithm at a much less computational cost. It has 108
48 been devised using a parallel with the Maximum Clique Problem (MCP). 109
49 A clique is a fully connected sub-graph, i.e. all pairs of nodes in it are 110
50 connected by an edge, so the MCP is concerned with searching for the 111
51 biggest clique in a graph. In our workflow, molecular trajectories are 112
52 represented as graphs in which each frame is depicted as a node. The 113
53 similarity between frames is encoded as binary (unweighted) edges, and 114
54 clusters are found following a heuristic search of big cliques. 115

55 The construction of a binary-encoded similarity matrix, instead of the 116
56 classical half/single-precision float matrix, leads to considerable RAM 117
57 savings regarding the existing QT implementations. This binary matrix 118
58 also allows implementing the fundamental clustering steps as bitwise 119
59 operations faster than the corresponding set operations when dealing with 120
60 considerable amounts of data. Our proposal, *BitQT*, is free and publicly 121
61 available at GitHub (<https://github.com/LQCT/BitQT.git>). 122

2 Computational Details

BitQT heuristic has been coded in Python 3 programming language and makes heavy use of two third-party libraries: version 1.9.4 of MDTraj (McGibbon *et al.*, 2015) for the fast RMSD calculations and version 1.6.1 of bitarray for all the binary-related operations (<https://github.com/ilanschnell/bitarray>). The two QT implementations used in this work for comparisons against *BitQT* correspond to the *QTPy* code, previously published by authors of this study and available at GitHub (<https://github.com/rglez/QT>) and the *qtcluster* command distributed in version 6.0.1 of the ORAC package (Procacci *et al.*, 1997).

We selected MD trajectories of different sizes and compositions to benchmark the performance of these algorithms. They are referred generically by their size as follows: 6K- a 6001 frames REMD simulation of the Tau peptide (Shea and Levine, 2016), 30K- a 30605 frames MD of villin headpiece based on PDB 2RJY (Melvin *et al.*, 2016), 50K- a 50500 frames MD of serotype 18C of *Streptococcus Pneumoniae*, 100K- a 100500 frames MD of Cyclophilin A based on PDB 2N0T, and 250K- a 250000 frames MD of four chains of the Tau peptide that corresponds to the MD simulation of an extended Tau peptide (PHF8) during 1 μ s (Álvarez-Ginarte *et al.*, unpublished work). Not referenced trajectories were obtained by the authors of this work. The details of the MD are available in the Supporting Information (S2: Details of the Molecular Dynamics Simulations). All trajectory and topology files used in this work can be found online at the following addresses: 6K, 50K, 100K, and 250K at <https://doi.org/10.6084/m9.figshare.c.5403930.v1>, and 30K at <https://doi.org/10.6084/m9.figshare.3983526.v1>.

QTPy and *BitQT* used the same quality threshold value k for each trajectory; 4 Å for trajectories 6K and 30K, 3 Å for trajectory 250K and 2 Å for trajectories 50K and 100K. These values were set after a trial/error procedure aided by visual inspection of the generated clusters uniformity. However, as *qtcluster* does not use the RMSD metric (Steipe, 2002), we adjusted the k values for each trajectory ran with this software. We multiplied the corresponding k by 2.4, in analogy with a previously published report of *qtcluster*'s authors (see S.I of Guardiani *et al.* (2012)).

All calculations were performed on an AMD Ryzen5 Hexa-core Workstation with a processor speed of 3.6 GHz and 64 GB RAM under a 64-bit Xubuntu 18.04 operating system. Run times and RAM peaks were recorded with the `/usr/bin/time` Linux command.

3 Approaching QT Clustering from a Maximum Clique Problem Perspective

If we define the *diameter* of a cluster C as the maximum distance between any pair of its elements (equation 1), the exact QT algorithm applied to an MD trajectory can be described as follows: After the user sets a similarity threshold k , one arbitrary frame is selected and marked as a candidate cluster C_1 . The remaining frames are iteratively added to C_1 if and only if two conditions hold: *Condition 1*- the entering frame increases the diameter of C_1 by the minimum amount, and *Condition 2*- the diameter of C_1 does not exceed the threshold k . A second candidate cluster is formed by starting with another frame and repeating the procedure. Note that all frames are made available to the second candidate cluster (frames from the first candidate cluster are not discarded from consideration). This process continues for all frames n in the trajectory until C_n candidate clusters have been formed. The one with more frames is set as a cluster, its elements removed from further consideration, and the entire process repeated until no more clusters can be discovered.

$$diam(C) = \max(d_{ij}) \mid \forall (i, j) \in C \quad (1)$$

1 The crucial aspect of the above-described algorithm lies in its ability 61
 2 to guarantee that all pairwise similarities inside a cluster will remain under 62
 3 the threshold k . This aspect is assured entirely by *Condition 2*, whose 63
 4 relevance has been previously discussed (González-Alemán *et al.*, 2020b). 64
 5 It is worth noting that *Condition 1* merely limits the size of retrieved clusters 65
 6 but has no impact in maintaining their collective similarity. 66

7 Concepts and tools from graph theory have been widely used to 67
 8 represent numerous situations in which several objects are mutually 68
 9 related. Before showing how QT can be approached from a graph- 69
 10 theoretical perspective, we will briefly define some basic underlying 70
 11 concepts. 71

12 A graph $G = (V, E)$ is a pair of a set of vertices (nodes) V and 72
 13 a set of edges E . Each edge is a two-element subset of V and denotes 73
 14 the adjacency between the nodes it connects. Two connected nodes are 74
 15 called *neighbors*, and the number of neighbors of a given node constitutes 75
 16 its *degree*. Connectivity of simple graphs can be represented using its 76
 17 *adjacency matrix*, a square symmetric matrix M in which $M_{ij} = 1$ 77
 18 if nodes i and j are connected and $M_{ij} = 0$ otherwise. If there is no
 19 directionality in the definition of the edges and there is no data associated
 20 to them, it is said that the graph is *undirected* and *unweighted*.

21 A *clique* is a subgraph in which vertices are all pairwise adjacent. If a 78
 22 clique is not contained in any other clique, it is said to be *maximal*, while 79
 23 the term *maximum clique* denotes the maximal clique with a maximum 80
 24 number of nodes (maximum *cardinality*). The maximum clique problem 81
 25 (MCP) solves the challenge of finding the maximum clique inside a given 82
 26 graph. 83

27 A central idea of MCP algorithms is the notion of *vertex coloring*. 84
 28 A *proper vertex coloring* refers to assigning a particular color (or any 85
 29 other unique label) to each vertex of a graph so that adjacent vertices do 86
 30 not share the same color. The *vertex coloring problem* consists of finding 87
 31 a proper coloring that uses the fewest number of colors, known as the 88
 32 graph's chromatic number (χ). It is common to employ coloring techniques 89
 33 because χ is an upper bound to the maximum clique's size of a graph. This 90
 34 property is exploited to discard impossible solutions and guide the search 91
 35 of cliques (San Segundo and Tapia, 2014). As exact coloring itself is an 92
 36 NP-hard problem, heuristics are usually applied. 93

37 To make a parallel between QT and MCP, it is possible to represent 94
 38 each frame of an MD trajectory as a node of an undirected graph in which 95
 39 edges depict RMSD similarity between nodes. Only edges with an RMSD 96
 40 less or equal to the threshold k are allowed, so there would be no weights
 41 associated with them. In that context, QT can be declared as an iterative 97
 42 search of cliques. However, QT cliques are not necessarily maximum due
 43 to *Condition 1* of the algorithm, which ensures that they should have a 98
 44 minimum weight instead of a maximum cardinality. *Condition 1* requires 99
 45 the diameter of the clusters to be minimum. Still, it is *Condition 2* that 100
 46 assures the respect of a quality threshold in the pairwise similarity of 101
 47 retrieved clusters. 102

48 Conveniently, a redefinition of the QT algorithm can be made to 103
 49 search for maximum-sized clusters instead of minimum-weighted without 104
 50 compromising the pairwise similarity assured by the second condition. 105
 51 In most clustering applications, maximizing the size of the clusters is 106
 52 a desirable feature. Relaxation of *Condition 1* in this way automatically 107
 53 converts QT in an MCP problem, accessible by the graph theory tools. This 108
 54 approach profoundly impacts how molecular similarity can be encoded and 109
 55 the efficiency of algorithms used to solve the problem, as discussed in the 110
 56 following sections. 111

57 3.1 Binary encoding of RMSD pairwise similarity 112

58 As the ultimate goal of our clustering proposal is to partition all MD 115
 59 trajectory frames, all the pairwise similarities should be analyzed. This 116
 60 information can be saved in RAM as a matrix to accelerate the algorithm's 117

run time. As $RMSD_{ij} = RMSD_{ji}$, the similarity matrix is symmetric. Although the valuable information is contained in one of the triangles, many current MD clustering software preserve the whole matrix to avoid the performance penalty of working with "triangular" data structures.

The amount of RAM needed for the storage of the matrix expressed in GB can be calculated using the equation 2, where N is the total number of frames in the trajectory and m is the size of the numeric type used to express the RMSD values (in bytes). Being the RMSD a float number ranging from 0.0 to infinite, the common choice is to use float numeric types to represent inter-frame similarity. Some clustering alternatives like TTClust (Tubiana *et al.*, 2018) use the costly choice of double-precision float ($m=8$). Other options like GROMACS (Abraham *et al.*, 2015), and WORDOM (Seeber *et al.*, 2007) packages use single-precision floats ($m=4$), saving half of RAM just by adjusting the precision used to express RMSD. It is worth noting that the minimum size of standard available floats is a half-precision value ($m=2$), which is enough for most MD clustering applications and the one used in QTPy.

$$V_{RAM} = \frac{m * N^2}{2^{30}} \quad (2)$$

Here we followed a different approach to diminish the value of m . If we conceive the QT algorithm as an MCP problem, after considering the relaxation of *Condition 1* our search will be focused on finding cliques of maximum cardinality, and no helpful information is extracted from the weight of the edges other than its absence or existence. This information can therefore be encoded as a binary matrix M where $M_{ij} = 1$ if nodes i and j are similar ($RMSD_{ij} \leq k$) or 0 otherwise. Note that M contains the same information that the adjacency matrix of the graph except for the diagonal, which in this case will always be one instead of zero ($RMSD_{ii} \equiv 0.0$). For the sake of simplicity, we will refer to M as the adjacency matrix of the trajectory graph.

By using the binary adjacency matrix, we reduce the RAM consumption of this object in 16, 32, or 64 times ($m = \frac{1}{8}$) compared to other software that deals with half, single or double-precision float values to represent the pairwise RMSD distance. Besides the RAM saving, expressing similarity as a binary matrix offers the possibility to perform the search of cliques using binary operators (AND and XOR), contributing to the speedup of the heuristic clique search algorithm we propose in the following section.

3.2 QT as a heuristic search of big cliques

Since MCP is an NP-hard problem, no efficient exact polynomial-time algorithms are expected to be found. Nevertheless, exact proposals exist to treat the MCP relatively fast for real problems of limited size. San Segundo and co-workers' efforts are of particular relevance for us as they also use the binary adjacency matrix of graphs and bitwise operations to develop their algorithms (San Segundo *et al.*, 2010; San Segundo *et al.*, 2013; San Segundo and Artieda, 2015; San Segundo *et al.*, 2016, 2017b,a). However, they are mainly focused on exact solutions rather than approachable heuristics.

A heuristic method tries to find a satisfactory solution to a complex problem using logical assumptions. While heuristics for MCP reduce the time of finding cliques, there is no guarantee that found cliques would be maximum. Nevertheless, heuristics are widely used in applications where a marginal error is not of great importance. In our case, we want to keep the common similarity of clusters, but their size is not of a big concern. After all, the original QT does not provide either maximum cliques. We are interested in a cheaper way to keep pairwise similarity, and for that purpose, a heuristic approach may suffice. Next we describe the workflow of the BitQT clustering algorithm, which is built upon a not previously published heuristic for searching big cliques (see the pseudocode in the

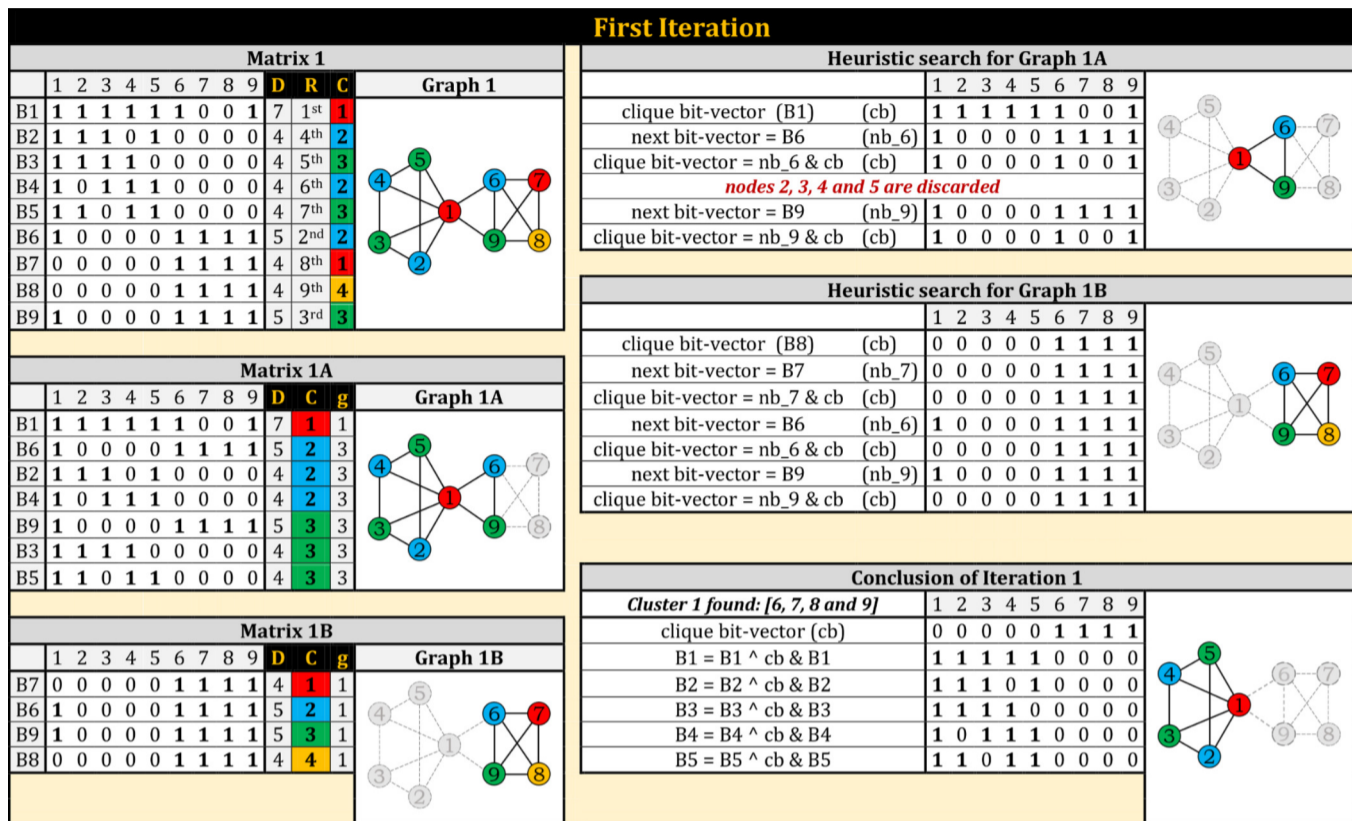


Fig. 1. First iteration of the binary heuristic for searching cliques implemented in BitQT.

1 Supplementary Information "S1: Pseudocode of BitQT algorithm"). A 29
2 formal review of the many MCP heuristics available is out of this paper's 30
3 scope and can be found elsewhere (Wu and Hao, 2015). 31
4 We start with the calculation of the binary similarity matrix that will 32
5 be stored in RAM. The float vector containing the one-versus-all RMSD 33
6 similarity of each frame is transformed into a bit-vector B_i (B1 to B9 in 34
7 Matrix 1, Figure 1) in which $B_{ij} = 1$ if $RMSD_{ij} \leq k$, zero otherwise. 35
8 Each vertex's degree is calculated as the total number of switched-on 36
9 positions in the B_i vector (D column in Matrix 1, Figure 1). Note that 37
10 B_i vectors always have 1 at the i^{th} position ($RMSD_{ii} = 0 \leq k$), so 38
11 D column actually contain $degree + 1$ of each vertex in the trajectory 39
12 graph. Then, the subsequent steps are followed. 40
13 **1- Vertex coloring:** Each vertex of the input graph (Graph 1, Figure 41
14 1) is ranked (column R, Matrix 1, Figure 1) in descending order of their 42
15 corresponding degrees (column D, Matrix 1, Figure 1). Following the rank 43
16 order, each vertex takes a color label that it shares with all other vertices 44
17 that are neither colored nor neighbors (column C, Matrix 1, Figure 1). 45
18 **2- Clique search from the maximum degree node:** After all vertices 46
19 are colored, the search of a clique starts considering only neighbors of 47
20 the maximum degree node of the graph (Graph 1A, Figure 1), which is 48
21 called the *seed* of the clique (node 1 in Matrix 1A, Graph 1A, Figure 1). 49
22 Neighbors of the seed are strictly ordered for further processing following 50
23 three criteria (DCg ordering); descending order of their degrees, ascending 51
24 order of their color class, and ascending order of the degeneracy of the color 52
25 class (columns D, C, and g, respectively, Matrix 1A, Figure 1). Note that 53
26 for our purposes, degeneracy is perceived as the number of nodes of the 54
27 color class in the context of the neighbors of a seed node, not in the entire 55
28 graph (in which case using it for order would be meaningless). 56

Following this ordering, the first node is selected to start a clique, and subsequent nodes will be added to that clique if they have a still-not-explored color and if they are adjacent to previously explored nodes (clique propagation).

BitQT performs this search using bitwise operations. The bit-vector B_i corresponding to the maximum degree node is set as the clique bit-vector (B1 in Heuristic search of Graph 1A, Figure 1). Following the DCg ordering, an AND operation is performed between the clique bit-vector and the next node bit-vector if it has a new color (B6 in Heuristic search of Graph 1A, Figure 1). Indices corresponding to bits that become zero by this operation are discarded from further consideration (B2, B3, B4, and B5) as they are not adjacent to processed nodes (B1 and B6). The resulting bit-vector becomes the new clique bit-vector used for the AND operation with the next candidate following the DCg ordering (B9). The bit-vector resulting from the iterative AND operations contains the members of the first clique.

3- Clique search from promising nodes: Once the clique retrieved by using the maximum degree node as the seed is found in the previous step, the same exploration strategy is conducted for every promising node in the original graph (Graph 1). A promising node (B8 in Graph 1, Figure 1) is defined as a node with a color not present in the first clique and whose degree is higher than the number of nodes in the first clique. Using such nodes as seeds for propagation might lead to the formation of a bigger clique (Heuristic search of Graph 1B, Scheme 1).

4- Conclusion and updating: When the maximum degree node and all promising nodes have been used as seeds, the maximum clique found is picked as a cluster and their members removed from the input graph (the corresponding B_i vectors removed from the binary matrix). An updating of the remaining bit-vector is necessary to set as zero all entries 57

Table 1. Run time and RAM consumption of analyzed QT implementations on different trajectories.¹

Traj. Name	# atoms (selection)	BitQT		qtcluster		QTPy	
		Run time h:mm:ss	RAM peak GB	Run time h:mm:ss	RAM peak GB	Run time h:mm:ss	RAM peak GB
6K	217 (all)	0:00:08	0.101	0:08:21	0.529	0:04:36	0.181
30K	64 (CA)	0:02:15	0.470	0:18:55	0.270	3:41:11	2.710
50K	78 (no H)	0:12:34	0.435	1:14:08	1.526	181:51:57	7.101
100K	660 (backbone)	1:15:37	4.355	0:00:49	>81.014	>200:00:00	18.626
250K	160 (backbone)	6:36:04	8.128	130:18:06	17.476	0:00:03	>117.000

¹ Bold entries denote either a time crash (job taking more than 200 h) or a memory crash (job taking more than 64GB)

corresponding to nodes that formed the cluster, which will not be available for subsequent iterations. This updating is bitwise encoded as a consecutive AND/XOR operation between remaining bit-vectors and the clique bit-vector (Conclusion of iteration 1, Figure 1). The same steps are repeated from Step 2 until no more cliques can be found.

During the execution of *BitQT*, some scenarios leading to ties may arise, for instance, selecting the node of the highest degree as seed (in “2-Clique search from the maximum degree node” and “3-Clique search from promising nodes”), or selecting the maximum clique (in “4-Conclusion and updating”). *BitQT* solves these cases by choosing the element with the lowest index among the available options as the “winner” of the tie. These ties can also appear in the original QT algorithm (see Section 3 when selecting the candidate cluster with most neighbors as a cluster). *QTPy* also picks as “winner” of the tie the element with the lowest index from the available options. Choosing one or another “winner” does impact the outcome of algorithms in terms of cluster composition. However, the choice of a “winner” in a tied scenario will never invalidate the discussed guarantees of *BitQT* or *QTPy*.

4 BitQT benchmark

4.1 Performance

In this section, we compare the run time and memory usage of *BitQT*, *QTPy* and *qtcluster*, which are the only QT implementations for MD we have found in the literature. These parameters are shown in Table 1 for the clustering of the six different MD trajectories that we described in Section 2 (6K, 30K, 50K, 100K, and 250K). Given that these software are programmed following distinct algorithms and also using different programming languages (Fortran 90 for *qtcluster* and Python 3 for *BitQT* and *QTPy*), we are only able to provide general insights into the disparate performances observed in Table 1.

From the three options, *QTPy* is the only one that always creates a square float matrix for saving the RMSD distances, so its RAM peak is expected to be the highest. The only exception is 6K, where the pairwise matrix uses only about 69 MB of RAM, so other data structures (or merely the molecular trajectory) will be responsible for the peak. RAM usage of *BitQT* also grows quadratically with the number of frames in the trajectory. However, as it uses bits instead of half-precision floats, there is a 16X memory saving in this object’s construction compared to *QTPy*.

The memory usage of *qtcluster* may be confusing at first sight, as it can process a 250K trajectory but produced a memory crash when dealing with a simulation of 100K frames. This behavior is a direct consequence of the similarity metric, the maximum difference between corresponding pairs of atoms. As expressed in equation 3, under this metric, the similarity of two frames S_m and S_n is assessed by the absolute maximum value of the difference between their inter-atomic distances.

$$d_{S_m, S_n} = \max_{i,j} |d_{ij}(S_m) - d_{ij}(S_n)| \quad (3)$$

This means that it is necessary to hold the square matrix of the selected inter-atomic distances for each conformation in RAM. In practice, *qtcluster* allocates the values of only one triangle of that matrix for every conformation.

The RAM used by the *qtcluster* similarity matrix (in GB) is expressed by equation 4, in which N is the total number of frames in the trajectory, m is the size of the numeric type used to express the similarity values (in bytes), and $natoms$ is the number of selected atoms. It is clear why *qtcluster* crashed at 100K but could process 250K; the 100K trajectory contained 660 atoms and 250K only 160. Substituting in equation 4 and taking $m = 4$ we obtain 81 GB for 100K and about 12 GB for 250K. Inconveniently, *qtcluster* can analyze big trajectories only when the number of selected atoms is relatively small.

$$V_{RAM_{qtcluster}} = \frac{m * N * \frac{natoms * (natoms - 1)}{2}}{2^{30}} \quad (4)$$

In a nutshell, while the three algorithms have quadratic memory complexity, the costs of *BitQT* and *QTPy* are governed by the trajectory size. In contrast, *qtcluster* is dominated by the size of the atomic selection.

Run time reported in Table 1 exhibits a general trend; *QTPy* is the slowest choice, followed by *qtcluster*, which is greatly outperformed by *BitQT*. It is worth noting that *QTPy* is the only one that implements the exact version of QT (Heyer *et al.*, 1999). As we have commented before, the exact QT has a very high computational cost evinced in the *QTPy* run times. The RMSD computation step can be safely discarded as the main contributor to the slow time performance of *QTPy* because it employs the same library that *BitQT* for this purpose (MDTraj). Given its slowness, *QTPy* applications are limited to the processing of small trajectories or as a reference for the development of future QT algorithms applied to the MD field.

qtcluster was designed as a high-speed alternative for the QT partitioning of MD. The similarity metric employed by this script (equation 3) is cheaper than the more customary RMSD and avoids any alignment. Somewhat similar to *BitQT*, *qtcluster* only preserves the original condition assuring the collective similarity of retrieved clusters. For big trajectories, however, *qtcluster* is not a fast option.

Comparatively, *BitQT* has the best run time performance allowing it to handle relatively long MD trajectories. The accelerated computing of optimal RMSD distances through the MDTraj engine joined to the developed binary-based heuristic for searching cliques are the cornerstones of its cheap cost.

4.2 Preservation of the Quality Threshold

As we discussed earlier in section 3, there are two fundamental restrictions in the QT original algorithm: *Condition 1*, which requires the diameter of clusters to be of minimum size, and *Condition 2*, which ensures the respect of a quality threshold in the values of intra-cluster similarity. *BitQT* conveniently relaxed the former, but it carefully does preserve the latter one. The previous claim implies that all clusters returned by *BitQT* must have a diameter less equal than the quality threshold k .

Figure 2 shows the distribution of all clusters’ diameter for every analyzed trajectory. As it is appreciated, pairwise distances between frames of the same cluster never surpass the predefined quality threshold k (4 Å for 6K and 30K, 3 Å for 250K, and 2 Å for 50K and 100K).

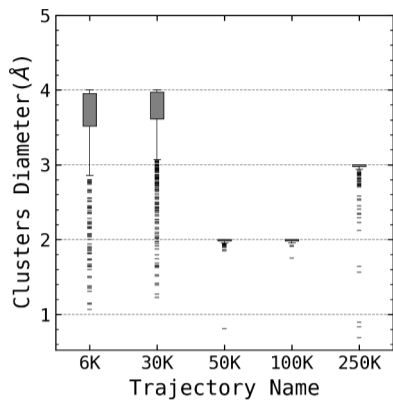


Fig. 2. Distributions of cluster diameters returned by BitQT for each analyzed trajectory.

Figure 2 also demonstrates that *BitQT* clusters are cliques in the MD trajectory graph. As we discussed in Section 3, an edge between two nodes i and j of the trajectory graph is set if and only if $d_{ij} \leq k$. If all pairwise distances between frames in every cluster are under k , then the corresponding nodes of the trajectory graph are pairwise connected, implying that clusters are indeed cliques.

4.3 Equivalence between *BitQT* and *QTPy*

If we consider an MD trajectory T as a set of N elements (frames) $T = \{t_1, t_2, \dots, t_N\}$, the outcome of applying a given clustering algorithm on T is a partition P of the N objects into C clusters, $P = \{p_1, p_2, \dots, p_C\}$, such that the union of all the subsets in P is equal to T and the intersection of any two subsets in P is empty. *QTPy* and *BitQT* produced such partitions ($Q = \{q_1, q_2, \dots, q_C\}$ and $B = \{b_1, b_2, \dots, b_C\}$ respectively) for the 6K, 30K, and 50K trajectories.

Considering $\binom{N}{2} = N(N-1)/2$ as the total number of element pairs (t_i, t_j) in T , there exist four classifications of pairs when comparing Q and B outcomes; a-) elements in a pair are placed in the same group in Q and in the same group in B (true positives), b-) elements in a pair are placed in the same group in Q and in different groups in B (false negatives), c-) elements in a pair are placed in the same group in B and in different groups in Q (false positives), and d-) elements in a pair are placed in different groups in Q and B (true negatives). It is possible to assess the equivalence between Q and B based on the number of pairs of elements lying in any of these four categories.

The Rand Index (Rand, 1971) (Equation 5) expresses the fraction of pairs of elements on which two clusterings coincide (from 0 for unrelated to 1 in a perfect match). However, RI approaches its upper limit as the number of clusters increases because d tends to grow even for poorly related partitions, giving a high score. An Adjusted RI (Hubert and Arabie, 1985; Steinley, 2004) corrected against "agreements-by-chance" (ARI) has been extensively used (Equation 6) to measure the correspondence between partitions created by clustering algorithms. ARI values extend from -1 (poorly related partitions) to 1 (highly similar partitions).

$$RI = \frac{a + d}{a + b + c + d} \quad (5)$$

$$ARI = \frac{\binom{N}{2}(a + d) - [(a + b)(a + c) + (c + d)(b + d)]}{\binom{N}{2}^2 - [(a + b)(a + c) + (c + d)(b + d)]} \quad (6)$$

An ARI analysis between partitions obtained with *QTPy* (Q) and *BitQT* (B) for trajectories 6K, 30K, and 50K is shown at Figure 3. Note that

instead of reporting just the global ARI between Q and B , we explicitly compared the ARI between both partitions at the top- X clusters (Q_X and B_X), taking X from 1 (the first cluster) to C (the total number of clusters). Consequently, the global ARI between Q and B corresponds to the last point of each curve. The remaining points indicate the correspondence between the first X clusters of Q and B .

For trajectories 6K and 30K, the global ARI is 0.87, indicating a good agreement between clusters produced by *QTPy* and *BitQT*. An even higher index is reported for the first X clusters with sizes bigger than 1% of the trajectory size ($ARI_{1\%}$). These most populated clusters are often considered the most relevant of the trajectory as they groups the representative conformational states explored in an MD simulation. $ARI_{1\%}$ (represented by a bold point in Figure 3A-C) is 0.96, 0.88 and 0.88 for trajectories 6K, 30K, and 50K, respectively. This is indicative of a very good agreement between the most popular clusters obtained by *QTPy* and *BitQT*.

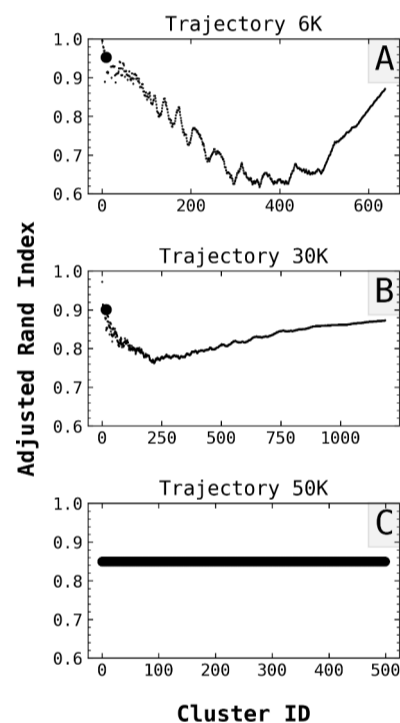


Fig. 3. To do

Observed ARI fluctuations at different top- X are expected because both algorithms pick their seeds to form clusters differently. It is possible that at a given value of X , clusters formed by *QTPy* were still not recovered by *BitQT* or vice versa. However, fluctuations are more pronounced for the less populated clusters.

5 Conclusions

The QT algorithm is an appealing option for partitioning MD trajectories as it assures a collective similarity of frames in recovered clusters. However, its inherent complexity currently limits its application. In the present work, we have relaxed a condition in the original formulation of QT

Instead of looking for minimum-sized clusters where all pairwise similarity values were under a threshold, we reformulated the problem to maximize the size of those clusters. This trivial change allowed us to

- 1 approach QT from an MCP perspective. The use of a similarity binary 57
 2 matrix (rather than a float-encoded one) **greatly** diminished the RAM 58
 3 resources. It made it possible to implement most clustering steps as fast 59
 4 bitwise operations. 60
- 5 Rather than an exact implementation of the MCP, we developed our 61
 6 modified version of QT called *BitQT* using an MCP heuristic whose 62
 7 out-coming clusters are in **good agreement** with those obtained by the 63
 8 original version *QTPy*. *BitQT* strictly guaranteed the preservation of the 64
 9 user-defined quality threshold in all reported clusters. 65
- ## 10 Acknowledgements
- 11 D.H.C. thanks Joan-Emma Shea and Zach Levine for providing the 6K 70
 12 trajectory used in this work. 71
- ## 13 Funding
- 14 This work was supported by the Eiffel Scholarship Program of Excellence 74
 15 of Campus France [P744468L to R.G.A.]; the Project Hubert Curien- 75
 16 Carlos J. Finlay [41814TM to R.G.A, F.L, and L.M.C]; and the 76
 17 Fondo Nacional de Desarrollo Científico y Tecnológico [CONICYT 77
 18 FONDECYT/INACH/POSTDOCTORADO/No. 3170107 to E.W.H.R.]. 78
- ## 19 References
- 20 Abraham, M. J. *et al.* (2015). Gromacs: High performance 82
 21 molecular simulations through multi-level parallelism from laptops to 83
 22 supercomputers. 84
- 23 Danalis, A. *et al.* (2012). Efficient quality threshold clustering for parallel 85
 24 architectures. *Proceedings of the 2012 IEEE 26th International Parallel 86
 25 and Distributed Processing Symposium, IPDPS 2012*, pages 1068–1079. 87
- 26 Daura, X. *et al.* (1999). Peptide Folding: When Simulation Meets 88
 27 Experiment. *Angewandte Chemie International Edition*, **38**(1/2), 89
 28 236–240. 90
- 29 Dutta, S. and Overbye, T. (2011). A clustering based wind farm collector 91
 30 system cable layout design. In *2011 IEEE Power and Energy Conference 92
 31 at Illinois*, pages 1–6. IEEE. 93
- 32 González-Alemán, R. *et al.* (2020a). BitClust: Fast Geometrical Clustering 94
 33 of Long Molecular Dynamics Simulations. *Journal of Chemical 95
 34 Information and Modeling*, **60**(2), 444–448. 96
- 35 González-Alemán, R. *et al.* (2020b). Quality Threshold Clustering 97
 36 of Molecular Dynamics: A Word of Caution. *Journal of Chemical 98
 37 Information and Modeling*, **60**(2), 467–472. 99
- 38 Guardiani, C. *et al.* (2012). Conformational Landscape of N-Glycosylated 100
 39 Peptides Detecting Autoantibodies in Multiple Sclerosis, Revealed by 101
 40 Hamiltonian Replica Exchange. *J. Phys. Chem. B*, **116**(18), 5458–5467. 102
- 41 Heyer, L. J. *et al.* (1999). Exploring expression data identification and 103
 42 analysis of coexpressed genes. *Genome Research*, **9**(11), 1106–1115. 104
- 43 Hubert, L. and Arabie, P. (1985). Comparing partitions. *J. Classif.*, **2**(1), 105
 44 193–218. 106
- 45 McGibbon, R. T. *et al.* (2015). MDTraj: A Modern Open Library for 107
 46 the Analysis of Molecular Dynamics Trajectories. *Biophysical Journal*, 108
 47 **109**(8), 1528–1532. 109
- 48 Melvin, R. L. *et al.* (2016). Uncovering Large-Scale Conformational 110
 49 Change in Molecular Dynamics without Prior Knowledge. *Journal of 111
 50 Chemical Theory and Computation*, **12**(12), 6130–6146. 112
- 51 Olson, M. T. *et al.* (2011). Production of reliable MALDI spectra with 113
 52 quality threshold clustering of replicates. *Journal of the American 114
 53 Society for Mass Spectrometry*, **22**(6), 969–975. 115
- 54 Peng, J.-h. *et al.* (2018). Clustering algorithms to analyze molecular 116
 55 dynamics simulation trajectories for complex chemical and biological 117
 56 systems. *Chinese J. Chem. Phys.*, **31**(4), 404–420. 118
- Procacci, P. *et al.* (1997). ORAC: A molecular dynamics program 119
 to simulate complex molecular systems with realistic electrostatic 120
 interactions. *Journal of Computational Chemistry*, **18**(15), 1848–1862. 121
- Rand, W. M. (1971). Objective criteria for the evaluation of clustering 122
 methods. *J. Am. Stat. Assoc.*, **66**(336), 846–850. 123
- Röttger, R. (2016). Clustering of Biological Datasets in the Era of Big 124
 Data. *Journal of integrative bioinformatics*, **13**(1), 300. 125
- San Segundo, P. and Artieda, J. (2015). A novel clique formulation for the 126
 visual feature matching problem. *Applied Intelligence*, **43**(2), 325–342. 127
- San Segundo, P. and Tapia, C. (2014). Relaxed approximate coloring in 128
 exact maximum clique search. *Computers and Operations Research*, **44**, 129
 185–192. 130
- San Segundo, P. *et al.* (2013). An improved bit parallel exact maximum 131
 clique algorithm. *Optimization Letters*, **7**(3), 467–479. 132
- San Segundo, P. *et al.* (2016). A new exact maximum clique algorithm for 133
 large and massive sparse graphs. *Computers and Operations Research*, 134
66, 81–94. 135
- San Segundo, P. *et al.* (2017a). A parallel maximum clique algorithm for 136
 large and massive sparse graphs. *Optimization Letters*, **11**(2), 343–358. 137
- San Segundo, P. *et al.* (2017b). An enhanced bitstring encoding for exact 138
 maximum clique search in sparse graphs. *Optimization Methods and 139
 Software*, **32**(2), 312–335. 140
- San Segundo, P. *et al.* (2010). Fast exact feature based data correspondence 141
 search with an efficient bit-parallel MCP solver. *Applied Intelligence*, 142
32(3), 311–329. 143
- Seeber, M. *et al.* (2007). Wordom: A program for efficient analysis of 144
 molecular dynamics simulations. *Bioinformatics*, **23**(19), 2625–2627. 145
- Shea, J.-E. and Levine, Z. A. (2016). Studying the early stages of protein 146
 aggregation using replica exchange molecular dynamics simulations. In 147
Protein Amyloid Aggregation, pages 225–250. Springer. 148
- Steinley, D. (2004). Properties of the Hubert-Arabie adjusted Rand index. 149
Psychol. Methods, **9**(3), 386–396. 150
- Steipe, B. (2002). A revised proof of the metric properties of optimally 151
 superimposed vector sets. *Acta Crystallographica Section A*, **58**, 506. 152
- Tang, Z. *et al.* (2010). A new method for alignment of lc-maldi-tof data. In 153
*2010 IEEE International Conference on Bioinformatics and Biomedicine 154
 (BIBM)*, pages 346–351. IEEE. 155
- Tubiana, T. *et al.* (2018). TTClust: A Versatile Molecular Simulation 156
 Trajectory Clustering Program with Graphical Summaries. *Journal of 157
 Chemical Information and Modeling*, **58**(11), 2178–2182. 158
- von Luxburg, U. *et al.* (2012). Clustering: Science or Art? *JMLR: 159
 Workshop and Conference Proceedings*, **27**, 6579. 160
- Wu, Q. and Hao, J.-K. (2015). A review on algorithms for maximum 161
 clique problems. *European Journal of Operational Research*, **242**(3), 162
 693–709. 163
- Yaakob, S. N. and Jain, L. (2012). An insect classification analysis based on 164
 shape features using quality threshold ARTMAP and moment invariant. 165
Applied Intelligence, **37**(1), 12–30. 166
- Yaakob, S. N. *et al.* (2010). A novel Euclidean quality threshold ARTMAP 167
 network and its application to pattern classification. *Neural Computing 168
 and Applications*, **19**(2), 227–236. 169