

# SCIENTIFIC REPORTS



OPEN

## A signal processing method for alignment-free metagenomic binning: multi-resolution genomic binary patterns

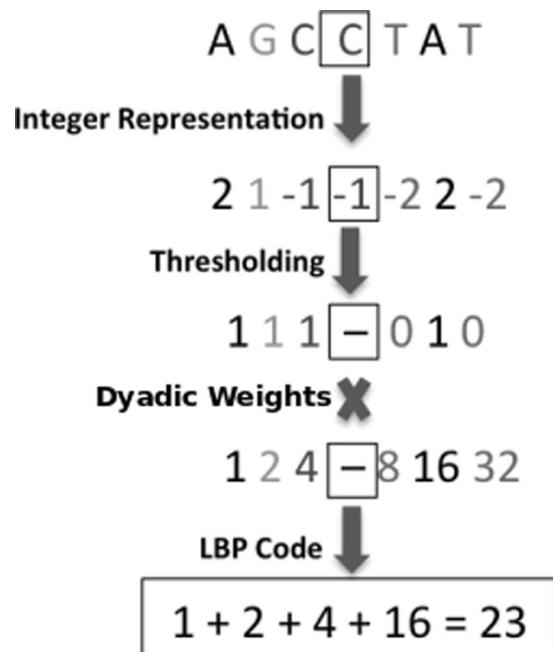
Samaneh Kouchaki<sup>1,2</sup>, Avraam Tapinos<sup>1</sup> & David L. Robertson<sup>1,3</sup> 

Algorithms in bioinformatics use textual representations of genetic information, sequences of the characters A, T, G and C represented computationally as strings or sub-strings. Signal and related image processing methods offer a rich source of alternative descriptors as they are designed to work in the presence of noisy data without the need for exact matching. Here we introduce a method, multi-resolution local binary patterns (MLBP) adapted from image processing to extract local 'texture' changes from nucleotide sequence data. We apply this feature space to the alignment-free binning of metagenomic data. The effectiveness of MLBP is demonstrated using both simulated and real human gut microbial communities. Sequence reads or contigs can be represented as vectors and their 'texture' compared efficiently using machine learning algorithms to perform dimensionality reduction to capture eigengene information and perform clustering (here using randomized singular value decomposition and BH-tSNE). The intuition behind our method is the MLBP feature vectors permit sequence comparisons without the need for explicit pairwise matching. We demonstrate this approach outperforms existing methods based on k-mer frequencies. The signal processing method, MLBP, thus offers a viable alternative feature space to textual representations of sequence data. The source code for our Multi-resolution Genomic Binary Patterns method can be found at <https://github.com/skouchaki/MrGBP>.

Algorithms in bioinformatics use textual representations of genetic information, sequences of the characters A, T, G and C represented as strings or sub-strings. For example, in genome assembly, exact substring matching of short *k*-mers of fixed length are typically used to identify related sequences/strings<sup>1,2</sup>. Although this approach works well for closely related data, it will fail predictably with divergent sequences, e.g., viruses, due to a lack of homologous regions retaining sufficient sequence identity for exact matching. While there are approaches that permit relaxed *k*-mer matching<sup>3,4</sup>, the processing methods used in signal/image processing offer an alternative feature space because they are designed to be rotation and scale invariant, and are generally less sensitive to noise by mapping data to a less detailed representation, i.e., 'texture' changes first introduced for segmenting an image in two-dimensions into several meaningful partitions<sup>8,9</sup>. It is based on assigning a code to each local window. Its implementation for one-dimensional data has been applied to other signal processing areas, specifically, speech processing<sup>10,11</sup>. Here we implement the superior multi-resolution version of LBP, called multi-resolution

We have implemented a signal processing method adapted from image comparisons (local binary patterns, LBP, Fig. 1) for the extraction of local changes in numerical representations of genetic sequence data. Preliminary results have been presented as conference papers using a linear<sup>6</sup> or non-linear<sup>7</sup> dimensionality reduction approach. LBP is a feature descriptor capturing local texture changes first introduced for segmenting an image in two-dimensions into several meaningful partitions<sup>8,9</sup>. It is based on assigning a code to each local window. Its implementation for one-dimensional data has been applied to other signal processing areas, specifically, speech processing<sup>10,11</sup>. Here we implement the superior multi-resolution version of LBP, called multi-resolution

<sup>1</sup>Evolution and Genomic Sciences, School of Biological Sciences, Faculty of Biology, Medicine and Health, The University of Manchester, Manchester, M13 9PT, UK. <sup>2</sup>Department of Engineering Science, University of Oxford, Oxford, OX3 7DQ, UK. <sup>3</sup>MRC-University of Glasgow Centre for Virus Research, Glasgow, G61 1QH, UK. Correspondence and requests for materials should be addressed to S.K. (email: [samaneh.kouchaki@eng.ox.ac.uk](mailto:samaneh.kouchaki@eng.ox.ac.uk)) or D.L.R. (email: [david.l.robertson@glasgow.ac.uk](mailto:david.l.robertson@glasgow.ac.uk))



**Figure 1.** Calculating the LBP code. A threshold of the integer numerical representation of the sequence (see Table 1) is determined by comparing the centre point (in the square) and its neighbours. The LBP code is then obtained by using dyadic weights.

Letter	Integer	EIIP	Atomic	Real
A	2	0.1260	70	-1.5
T	-2	0.1335	78	1.5
C	-1	0.1340	58	-0.5
G	1	0.0806	66	0.5

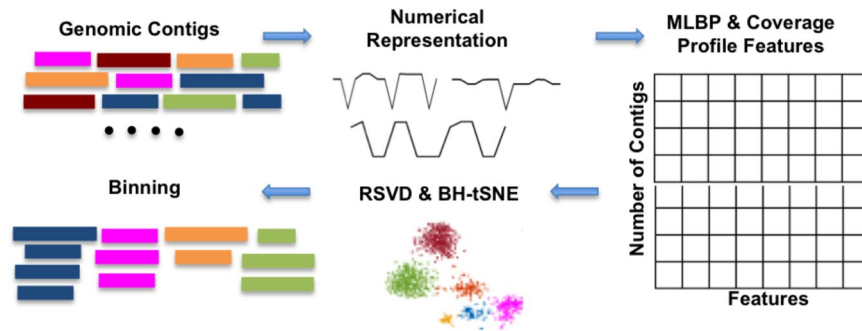
**Table 1.** The numerical representation of each letter considering Integer, EIIP, atomic and real.

LBP (MLBP), which considers texture changes at different scales<sup>12</sup> and benchmark its use in the processing of metagenomics data. We rationalise that in the same way as images, genomic sequences have ‘texture’ patterns at various scales that can be extracted using MLBP. Crucially for alignment/homology-free comparison the arbitrary location of each pattern does not affect the extracted feature vector.

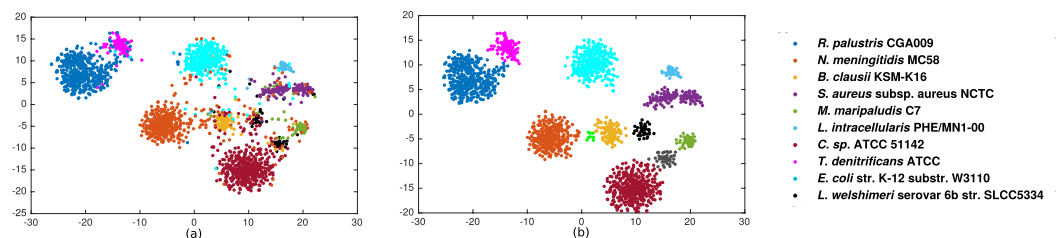
To test an application of the MBLP method and its effectiveness compared to LBP and string-based methods, we consider the problem of unsupervised grouping of genomic contigs into species-level groups (‘binning’) based on alignment-free genome composition comparisons. High-throughput/‘next-generation’ sequencing technologies have generated enormous volumes of data in metagenomic studies. In these samples, the sequence reads can be from the same or different genomes from a microbial community of viruses and bacteria, including divergent variants of the same species. Hence, reconstructing (assembling) individual genomes from this mixed data can be problematic. Moreover, sequencing errors, sequence repetition, insufficient coverage and high levels of genetic diversity can give rise to fragmented assemblies. Furthermore, comparing metagenomic data to existing reference genomes (taxonomic binning) will only identify some of the reads/contigs present. Consequently, genome composition-based techniques<sup>13,14</sup> have been introduced as an alternative way to analyse the species composition of metagenomic samples<sup>15</sup>. These methods use species-specific genomic signatures extracted by calculating the normalised frequency of  $k$ -mers of a specific size, commonly  $k = 4$ <sup>16,17</sup>. The signatures are obtained by counting the occurrences of each  $k$ -mer combination where the  $k$ -mer frequency of each sequence represents a feature vector in high-dimensional space.

A number of metagenomic binning techniques have used genomic signatures as features, for example, leveraging across-sample coverage-profiles<sup>18,19</sup>. The method emergent self organising maps (ESOM) based binning uses contour boundaries to visualise the clusters<sup>19</sup>. Unfortunately, ESOM plots are computationally very demanding. Other methods that consider coverage across multiple samples include CONCOCT<sup>18</sup> and MetaBAT<sup>20</sup>. However, they require a high number of samples to perform well, e.g., 50 or more. VizBin<sup>17</sup> is another visualisation approach that considers a single sample, but it needs manual selecting of the centroids for binning.

To perform clustering/binning we have first used singular value decomposition (SVD)<sup>21,22</sup> (specifically randomised SVD, RSVD<sup>23</sup>, for time efficiency) to reduce the dimensionality of the data, i.e., to identify the principal components of the MBLP feature vectors; termed ‘eigengene’ information<sup>24</sup>. Second, these eigengene features are passed as an input to Barnes-Hut  $t$ -distributed stochastic neighbor embedding (BH-tSNE)<sup>25</sup> for visualisation of the clusters in the data.



**Figure 2.** Schematic overview of our implementation of the MrGBP method to characterise the species relationships among metagenomic contigs.



**Figure 3.** Visualisation of the simulated metagenomic community using Integer nucleotide mapping, MLBP to extract features, RSVD for feature reduction, BH-tSNE two-dimensional representation and cluster identification using DBSCAN comparing (a) manually annotated clusters (see species names in key) to (b) the DBSCAN defined clusters.

Nucleotide mapping	Atomic	EIIP	Real	Integer
Precision	97.23	89.08	97.41	98.38
Recall	94.82	96.80	93.96	96.35
F1 score	96.01	92.78	95.65	97.35
Number of clusters	12	10	13	12

**Table 2.** Precision, recall, F1 score (%) and the number of clusters for various nucleotide mappings for a simulated low complexity metagenomic dataset.

An overview of our approach Multi-resolution Genomic Binary Patterns (MrGBP) is depicted in Fig. 2. We apply our method to both simulated and real metagenomic datasets, and demonstrate our results compare favourably to several existing binning methods. We also consider the effect of including coverage information across-samples in a hybrid approach to maximise the performance in longitudinal metagenomic samples and show improved performance. Collectively our results demonstrate the use of an image/signal processing method (MLBP) in bioinformatics, a new feature space for sequence analysis. The platform information for the reported run times is provided in the ‘Additional information’ Section.

## Results and Discussion

Calculating MLBP requires numerical data as an input (Fig. 1). Thus, genomic sequences need to be first mapped into one or several numerical representations<sup>26,27</sup>. Representation methods can be based on biochemical or biophysical properties of DNA molecules or be arbitrarily assigned numbers (Table 1). MLBP features are then extracted from these numerical representations and used to compare sequence data.

The performance of our method is tested for a low complexity simulated dataset using different numerical mappings (EIIP, atomic, real and integer nucleotide representations, Table 1) for MLBP lengths  $p \leq 6$  (Supplementary Figure 1). For example, for the integer representation our automated binning approach very closely matches the manually annotated clusters (compare panels a and b in Fig. 3). Specifically, the contigs from different species form visually separate clusters with very limited overlap with the clusters of other species.

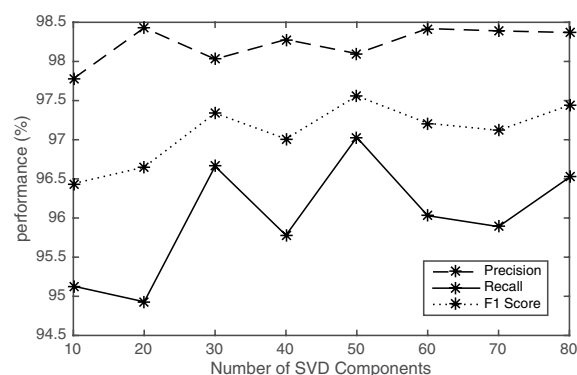
The different numerical representations provided slightly different data clusters but overall the results demonstrated similar performance (Table 2). The Integer representation was selected for subsequent analysis as it has relatively high performance and demonstrated more discrimination compared to the other representation methods. The average run time was 75.35 seconds (2184 contigs with total length 33138556 nucleotides). The run time includes loading the data, numerically representing the data, MLBP feature extraction and dimensionality reduction using BH-tSNE (Fig. 2).

Window length	2	4	6	8
Feature dimension	4	20	84	212
Precision	86.08	97.78	98.38	98.21
Recall	85.80	94.04	96.35	94.63
F1 score	85.94	95.87	97.35	96.39
Number of clusters	19	16	12	11
Run time	22.43	27.89	36.22	46.14

**Table 3.** Precision, recall, F1 score (%), the number of clusters and the run times for MLBP of various window lengths or feature dimensions ( $p \leq P$ ) and integer representation.

Number of RSVD Components	10	20	30	40	50	60	70	80
RSVD run time	0.06	0.12	0.20	0.30	0.44	0.55	0.69	0.85
BH-tSNE run time	27.76	28.83	30.04	31.20	33.05	34.65	35.63	35.80

**Table 4.** Run times of RVSD and BH-tSNE for various number of RSVD components.



**Figure 4.** Precision, recall and F1 score (%) by keeping different numbers of RSVD components. Integer representation and  $p \leq 6$  have been considered to analyse the simulated metagenomic dataset.

To check the effect of changing the window length, we considered various lengths of MLBP windows (Table 3 and Supplementary Figure 2). As the MLBP vectors are based on a histogram, the number of features is determined by the window length, which may affect final performance. Here, run time only includes the time to numerically represent the data and MLBP feature selection.

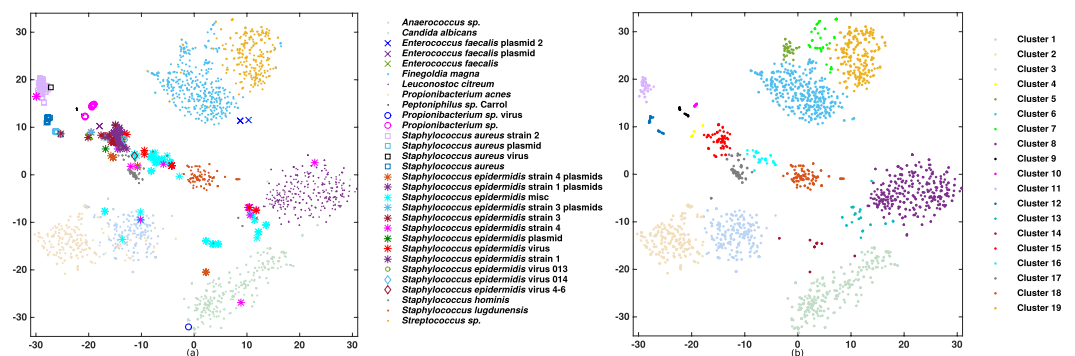
With smaller window lengths, the resulting feature vectors cannot describe the underlying structure of the metagenomic dataset, while larger feature vectors increases the time complexity (Table 3). Hence, window size should be sufficiently large to maintain the distinctness of the signal (information regarding texture changes across various contigs).

The computational complexity of our method increases as the dimensions of the feature space increase. Therefore, we considered how keeping different numbers of eigen factors can affect the performance and run time of our method (Fig. 4). We use the numerical integer representation for the nucleotide mapping and  $p \leq 6$  for feature selection. The results show that after keeping a number of eigen factors, i.e., 30, the final performance does not change significantly. However, as the number of eigen factors increases the run time of RSVD and BH-tSNE increases (Table 4). These results demonstrate that the MLBP method can analyse a metagenomic dataset in a reasonable time frame. Moreover, it is performing well considering only one sample was analysed.

**Comparison with Existing Methods for Simulated 10 and 100 Metagenomic Data.** We considered two simulated datasets with 10 and 100 genomes to compare our results to both low and complex metagenomic communities. Our results compared favourably with CONCOCT<sup>18</sup>, MetaBAT<sup>20</sup> and MaxBin2<sup>28,29</sup> (Table 5). CONCOCT bins the data by employing sequence composition and across-sample coverage. The method has been compared with a range of methods including MetaWatt<sup>30</sup>, SCIMM<sup>31</sup> and CompostBin<sup>32</sup> to show its advantage over composition based techniques. However, for available high complexity metagenomic data CONCOCT does not work as well as many samples are usually required for it to perform well. MetaBAT bins the metagenomic data using probabilistic distances of genome abundance with sequence composition. It is an efficient method for analysing complex metagenomic data. MaxBin was originally introduced for single sample data in which it bins the data based on tetra-nucleotide frequencies and it has been extended to MaxBin2 to support multiple samples. MetaBAT and MaxBin2 produce many unclassified contigs. Consequently, they have higher precision but lower recalls.

Methods	Precision	Recall	F1 score	Number of clusters
10 Genomes				
MLBP	98.38	96.35	97.36	12
CONCOCT	98.56	97.35	97.95	19
MetaBAT	90.82	94.98	92.85	9
MaxBin	93.43	96.65	95.01	10
4-mer	96.14	70.80	81.54	13
LBP	90.41	96.33	93.27	11
100 Genomes				
MLBP	91.52	83.97	87.58	116
CONCOCT	60.73	96.37	74.51	79
MetaBAT	92.34	89.62	90.96	104
MaxBin	89.83	83.96	86.80	85
4-mer	95.32	69.56	80.43	98
LBP	65.60	90.67	76.13	101

**Table 5.** Precision, recall, F1 score (%) and the number of clusters for our proposed method, CONCOCT, MetaBAT and MaxBin.



**Figure 5.** Visualisation of the infant gut metagenomic community using integer nucleotide mapping, MLBP to extract features, RSVD feature reduction, BH-tSNE two-dimensional representation and cluster identification using DBSCAN comparing (a) manually annotated clusters (see bacteria species, virus or plasmid names in key) to (b) the DBSCAN defined clusters 1 to 19.

For the 100 simulated genomes data our method performs better than CONCOCT and MetaBin methods, and quite close to MetaBAT. Lower performance is mainly because DBSCAN does not work very well for a very dense feature space (high complexity data representation). It may result in some unclustered contigs and therefore, lower performance. It still shows that the proposed pipeline can work for low and high complexity datasets. Note, alternative clustering methods could be explored.

For completeness we also ran our ‘MLBP pipeline’ (Fig. 2) replacing MLBP with (i) LBP and (ii) a k-mer text/string-based representation to compare our feature space with a commonly used 4-mer frequencies. The results indicate that the MLBP method has a more discriminative feature vector and better performance than either LBP or the string representation (Table 5).

**Real Data: Infant Human Gut.** A relatively low-complexity infant human gut dataset<sup>19</sup> was analysed to test the performance of our method with real data. A main reason for considering this dataset is to show the effectiveness of the MBLP method to bin low abundant viral community data to benchmark our texture analysis approach. The integer numerical representation was used for the nucleotide mapping,  $p \leq 8$  for feature selection and the first 60 eigen components in the dimensionality reduction stage (RSVD).

MLBP binned the data into 19 clusters with precision and recall of 88.34 and 97.22 at the species level. BH-tSNE representation of the data demonstrates the genomic contigs of the same or very similar contigs are binned together (Fig. 5). While some of the plasmids and viruses (bacteriophages) clustered with their associated host clusters, most species formed their own cluster. The bacterial species tend to form separate clusters, for example, *Anaerococcus sp.* and *C. albicans* form clusters 1 and 3 (Fig. 5). However, separating plasmid or virus from its host is less straight-forward due to their closer genome compositions. Nonetheless, our method manages to bin *S. aureus* strains, their plasmid and virus into two groups; (1) *S. aureus* strain and plasmid and (2) *S. aureus* strain 2 and virus. *Propionibacterium sp.* appears as a separate bin. *E. faecalis* and one of its plasmids forms one cluster. *S. epidermidis* has three strains, three viruses, one integrated virus (prophage) and several plasmids

Methods	Precision	Recall	F1 score	Number of clusters
MLBP	88.34	97.22	92.57	19
CONCOCT	79.5	90.62	84.69	32
MetaBAT	84.23	92.35	88.10	10
MaxBin2	82.84	93.50	87.84	10

**Table 6.** Precision, recall, F1 score (%) and the number of clusters for our proposed method, CONCOCT, MetaBAT and MaxBin2.

and the algorithm managed to bin them into five clusters where *S. epidermidis* strains 1 and 3 clustered together (including virus 13 and 14), with strain 4 forming a separate cluster (including virus 46).

Our results compare favorably with CONCOCT<sup>18</sup>, MetaBAT<sup>20</sup> and MaxBin2<sup>28,29</sup>, showing better performance on this dataset with small sample size (11 samples) in comparison with the other algorithms tested (Table 6).

To further investigate the relationship between clusters, the abundance patterns of each cluster were calculated based on the number of reads mapped to contigs at the different sampling time points (Supplementary Figure 3). Pairwise correlation coefficients were then calculated to check for any pattern among the clusters. The results suggests that there is a strong correlation between clusters of related species (Supplementary Figure 3). For example, the clusters of *Propionibacterium* and *Peptoniphilus* species have similar abundance patterns (Clusters 9–10). Similar results were also found in<sup>19</sup> where both species have proliferation in later stages and hence are well-adapted to the gut. Moreover, two clusters have been formed for *F. magna* with very similar coverage patterns (clusters 5–6). Consequently, this similarity could be analysed further to join some of the clusters. A similar pattern can be observed in the clusters of *S. aureus*, confirming the relationship between each bacteria and virus (clusters 11–12). The five clusters of *S. epidermidis* also share similar coverage patterns (clusters 13–17). A further step could be to cluster all the contigs of these five clusters separately to have a better separation of the related strains and viruses.

Finally, we checked the run time of our method. It takes about three minutes (108.67 s) to analyse this dataset (the number of contigs is 2293 and total length of them is 27594702). Although our code is relatively fast, it could be further optimised in terms of both time and memory.

## Conclusion

We have demonstrated that the image and signal processing technique, MBLP, can be adapted to numerical nucleotide sequence data comparisons and performs significantly better than LBP alone. Applied to metagenomic binning and visualisation, MLBP captures the genomic signature changes effectively, i.e., genome texture patterns, permitting alignment-free comparison and clustering of related contigs. Our results on simulated genomic fragments and contigs from infant human gut samples demonstrate that a signal processing method can capture the underlying taxonomic structure of the microbiome data and performs favourably in comparison to existing metagenomic methods. Collectively our results demonstrate the ‘signal’ in genome data can be just, if not more effectively, captured by appropriate image/signal processing algorithms as opposed to text/string-based methods. This demonstrates the potential for exploitation of an alternative feature space for alignment-free comparison of genomic sequence data either alone or combined with text/string-based representations, i.e., ‘multi-view’ representation of the data. Using other LBP/MLBP variants or features descriptors from image or signal processing will be investigated in future work.

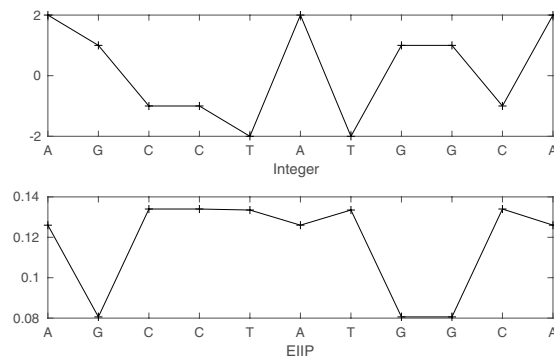
## Methods

Our methodological pipeline (Fig. 2) is comprised of several steps: (1) numerically represent the genomic contigs using a nucleotide mapping (Table 1). (2) MLBP is used to extract features from these numerical representations. If available, cross-sample coverage information (mean and standard deviation) is extracted separately using Bowtie 2<sup>33</sup> and can be considered as extra information to be added in the MLBP feature space. (3) Eigengene information is extracted using RSVD to reduce the dimensions of the feature matrix. (4) BH-tSNE is used to map RSVD features to a two-dimensional space for visualisation and data binning. (5) For quantitatively evaluating the visualisation performance, we cluster the BH-tSNE projected data using DBSCAN a density-based spatial clustering algorithm<sup>34</sup> and calculate the precision, recall and F1 score between the DBSCAN assigned labels and the original labels.

We note each step of the pipeline was based on having an appropriate analysis for the metagenomic data based on a novel feature space. Our primary purpose is to present this as a working view of the feature space, and not a novel metagenomics pipeline as such. Further optimisation in terms of implementation is of course possible and some options are already provided in the online software, e.g., changing the numerical representation.

**The Nucleotide Mapping.** The genomic reads can be represented numerically in two ways: (i) Assigning an arbitrary value to each letter A, C, G or T of the nucleotide sequence, i.e., Voss<sup>35</sup>, two or four bit binary representations<sup>36,37</sup> or (ii) defining numerical representations that correspond to certain biochemical or biophysical properties of the DNA molecules: electron ion interaction potential (EIIP)<sup>38</sup>, paired nucleotide representations<sup>39</sup> or atomic representations<sup>40</sup>.

As each numerical representation method assigns different values to each nucleotide this can lead to different results and performance when LBP/MLBP is applied. We thus compare several existing numerical representations from the literature (EIIP, atomic, real and integer nucleotide representations). Table 1 shows the value assigned to each nucleotide in each of the representations. Figure 6 shows an example of mapping a nucleotide sequence to two numerical vectors.



**Figure 6.** Two representations: integer and EIIP. Each nucleotide A, C, G or T in the sequence is assigned to a value depending on the numerical representation.

**Multi-resolution Local Binary Patterns.** LBP has found popularity not only in the field of image processing but also in signal processing<sup>41</sup>. The LBP distribution of genomic contigs was considered as the species specific genomic signatures in<sup>6</sup>. LBP examines the neighbouring points of each data point and assigns a binary code to it. By considering  $x(t)$  as the numerical representation of the  $t$ th position of a genomic segment, LBP is defined as

$$\text{LBP}(x(t)) = \sum_{i=0}^{p/2-1} \{\text{Sign}(x(t+i-p/2) - x(t))2^i + \text{Sign}(x(t+i+1) - x(t))2^{i+p/2}\}, \quad (1)$$

where  $p$  is the number of neighbouring points and Sign indicates the sign function

$$\text{Sign}(x) = \begin{cases} 0 & x < 0 \\ 1 & x \geq 0 \end{cases} \quad (2)$$

The difference between each neighbouring point and the centre point  $t$  is passed to a Sign function. Consequently, each window of length  $p+1$  is represented by a  $p$ -bit binary number. Each binary number is converted to a LBP code using a dyadic weight. An example of the LBP operator can be seen in Fig. 1 where  $p=6$ . The value of the centred point (in the square in Fig. 1) is compared with the six neighbouring points to produce the binary number and LBP code. The distribution of the LBP codes are defined using the obtained codes for each window:

$$\mathbf{h}_k = \sum_{p/2 \leq i \leq N-p/2} \delta(\text{LBP}_p(x(i), k)), \quad (3)$$

where  $\delta$  shows the Kronecker delta function,  $k=1, 2, \dots, 2^p$  is all possible values of LBP codes and  $N$  is the genomic fragment length. Considering the distribution of LBP codes makes the feature space independent of each pattern location and only dependent to frequency of each MLBP code.

MLBP is an LBP extension that combines the results of LBP distribution from various values of  $p \leq P$ . Consequently, the pattern changes of different resolution levels are considered to improve the description of the data inputs. Here, we apply MLBP to one-dimensional linear sequences to consider pattern changes of various lengths. LBP/MLBP is selected in this work due to its performance in other applications and also as it is very fast to calculate.

**Across-Samples Coverage Information.** To obtain the coverage profile for contigs across the longitudinal samples, the Illumina reads were mapped to contigs with Bowtie 2<sup>33</sup> for each time point. SAMtools<sup>42,43</sup> was then used to produce a per base depth file. As a result, our coverage feature vector for each genomic contig is the average and standard deviation of the per base depth for each contig. Coverage information provides extra information that optionally can be added to the MLBP feature space.

**Randomised Singular Value Decomposition.** A metagenomic community can be considered as a linear combination of genomic variables. The histogram of MLBP codes for each genomic fragment captures the local changes in the pattern (the “texture”) of each distinct contigs. By representing a vector of MLBP codes for each contig, low-rank matrix approximations can be used for efficient analysis of the metagenomic data. Our assumption in using SVD is that the MLBP codes of the contigs from each species have a distinct energy contribution. Therefore, the data can be represented as a linear combination of mutually independent components. RSVD a faster version of SVD was used here<sup>23</sup>.

SVD decomposition of a matrix  $\mathbf{X}$  is defined as

$$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T, \quad (4)$$

where  $\mathbf{U}$  and  $\mathbf{V}$  are the left and right singular vectors,  $\Sigma$  is singular values and  $(\cdot)^T$  denotes the transpose operator.

In metagenomic data analysis due to data complexity, SVD can be time consuming. Therefore, RSVD is used as an accurate and robust solution to estimate the dominant eigen components quickly<sup>23</sup>.

RSVD calculates the first  $i$ th eigen components of the data by using QR decomposition and mapping  $\mathbf{X}$  to a smaller matrix as

$$\begin{aligned}\Omega &= \text{randn}(N, i), \\ \mathbf{Y} &= \mathbf{X}\Omega, \mathbf{Y} = \mathbf{Q}\mathbf{R} \\ \mathbf{B} &= \mathbf{Q}^T\mathbf{X},\end{aligned}\quad (5)$$

where randn generates a random matrix of the size of its inputs and  $N$  is the number of contigs. After decomposing  $\mathbf{B}$  using SVD, the final factors are obtained using  $\mathbf{Q}$  and the eigen factors of  $\mathbf{B}$ .

**Barnes-Hut t-Distributed Stochastic Neighbor Embedding.** BH-tSNE is used in many research areas as a nonlinear technique for high dimensional data visualisation<sup>25</sup>. It works based on keeping the locality of the data in the lower dimension and was used in this paper for two-dimensional data visualisation and clustering. BH-tSNE is based on the divergence minimisation of input objects distributions and the corresponding low-dimensional data points. As a result, it can preserve the original local data structure in the final lower dimensional visualisation. Normalised Gaussian kernel has been considered as an ordinary similarity measure but it scales quadratically to the number of data points. The main objective function also has been approximated by defining the similarity function based on a number of neighbouring points<sup>25</sup>. In addition, a vantage-point tree is employed for decreasing search complexity. BH-tSNE is thus an efficient ( $O(N \log N)$ ) dimensionality reduction approach and is used in this paper for two-dimensional data visualisation and clustering.

**DBSCAN.** DBSCAN is a popular density-based clustering algorithm with the aim of discovering clusters from the approximate density distribution of corresponding data points. DBSCAN does not need the number of clusters to be specified but has two parameters that need to be determined: epsilon that indicates the closeness of the points of each cluster to each other and minPts, the minimum neighbours a point should have to be considered into a cluster. The initialisation point is a random point which has not been visited previously. The neighbourhood of this point is then retrieved and if it consists of an acceptable number of elements, a cluster is formed, otherwise the element is considered as noise. Hence, DBSCAN may result in some unclustered samples.

Usually DBSCAN parameters are not known prior to analysis and there are several ways to select their values. One way is to calculate the distance of each point to its closest nearest neighbour and use the histogram of distances to select epsilon. After selecting epsilon a histogram can be obtained of the average number of neighbours for each point using the epsilon. Some of the samples do not have enough neighbouring points and can be considered as noise. Implementation of the parameter selection is included in spark\_dbscan ([https://github.com/alitouka/spark\\_dbscan](https://github.com/alitouka/spark_dbscan)).

DBSCAN can find arbitrary shaped clusters, and is robust to outliers. However, it may not identify clusters of various densities or may fail if the data is very sparse. It is also sensitive to the selection of its parameters and the distance measure (usually Euclidean distance). The distance measure can affect any other clustering technique as well.

**Datasets.** To validate the effectiveness of our methodology we consider both simulated and real datasets. Simulated metagenomic data of Illumina sequences for 10 and 100 genomes (Supplementary Tables 1 and 3) was downloaded from [http://www.bork.embl.de/mende/simulated\\_data/](http://www.bork.embl.de/mende/simulated_data/). The data were assembled by Ray Meta<sup>44</sup> into contigs ( $k = 31$ ). Using these datasets, various aspects of our method, including MLBP window length and RSVD number of eigen components, have been analysed.

For the real data analysis, a time-series metagenomics human gut dataset comprised of 11 samples (18 runs) taken over nine days from a newborn infant<sup>19</sup> was used. The authors have assembled the data into 2329 contigs. This assembly and binning information is provided at <http://ggkbase.berkeley.edu/carroll/>. Corresponding Illumina reads can be downloaded from the NCBI, SRA052203, which consists of 18 Illumina sequencing runs (SRR492065-66 and SRR492182-97). For the real data, we mapped the reads to the contigs using Bowtie2 and coverage profiles have been obtained using SAMtools.

**Performance Evaluation.** In order to check the performance of our MrGBP method, DBSCAN<sup>34</sup> has been used to cluster the final results. The precision, recall and F1 score are calculated between the DBSCAN assigned labels and the original labels to determine the performance as a measure of a clusters "purity". Assuming there are  $m$  genomes in the dataset and it is binned to  $k$  clusters, the precision, recall and F1 score can be calculated as

$$\begin{aligned}\text{Precision} &= \frac{\sum_{i=1}^k \max_j s_{ij}}{\sum_{i=1}^k \sum_{j=1}^m s_{ij}} \\ \text{Recall} &= \frac{\sum_{j=1}^m \max_i s_{ij}}{\sum_{i=1}^k \sum_{j=1}^m s_{ij} + \sum \text{unbinned sequences}} \\ \text{F1} &= 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}\end{aligned}\quad (6)$$

where  $s_{ij}$  is the length of contigs in cluster  $i$  corresponds to genome  $j$ .



## References

- Zerbino, D. R. & Birney, E. Velvet: algorithms for de novo short read assembly using de bruijn graphs. *Genome Research* **18**, 821–829 (2008).
- Bankevich, A. *et al.* SPAdes: a new genome assembly algorithm and its applications to single-cell sequencing. *Journal of Computational Biology* **19**, 455–477 (2012).
- Healy, J. & Chambers, D. Approximate k-mer matching using fuzzy hash maps. *IEEE/ACM Transactions on Computational Biology and Bioinformatics* **11**, 258–264 (2014).
- Shajii, A., Yorukoglu, D., Yu, Y. W. & Berger, B. Fast genotyping of known SNPs through approximate k-mer matching. *Bioinformatics* **32**, i538–i544 (2016).
- Zhao, Y. Theories and applications of LBP: a survey. *International Conference on Intelligent Computing*, 112–120 (Springer, 2011).
- Kouchaki, S., Tirunagari, S., Tapinos, A. & Robertson, D. L. Local binary patterns as a feature descriptor in alignment-free visualisation of metagenomic data. *Symposium Series on Computational Intelligence (SSCI)*, 1–6 (IEEE, 2016).
- Kouchaki, S., Tirunagari, S., Tapinos, A. & Robertson, D. L. Marginalised stack denoising autoencoders for metagenomic data binning. *Conference on Computational Intelligence in Bioinformatics and Computational Biology (CIBCB)*, 1–6 (IEEE 2017).
- Pietikäinen, M. & Ojala, T. Texture analysis in industrial applications. *Image Technology*, 337–359 (Springer, 1996).
- Tirunagari, S. *et al.* Detection of face spoofing using visual dynamics. *IEEE Transactions on Information Forensics and Security* **10**, 762–777 (2015).
- Chatlani, N. & Soraghan, J. J. Local binary patterns for 1-D signal processing. *18th European Signal Processing Conference*, 95–99 (IEEE, 2010).
- Alegre, F., Amehraye, A. & Evans, N. A one-class classification approach to generalised speaker verification spoofing countermeasures using local binary patterns. *Sixth International Conference on Biometrics: Theory, Applications and Systems (BTAS)*, 1–8 (IEEE, 2013).
- Ojala, T., Pietikäinen, M. & Maenpää, T. Multiresolution gray-scale and rotation invariant texture classification with local binary patterns. *IEEE Transactions on Pattern Analysis Machine Intelligence* **24**, 971–987 (2002).
- Blaisdell, B. E. A measure of the similarity of sets of sequences not requiring sequence alignment. *Proceedings of the National Academy of Sciences USA*. **83**, 5155–5159 (1986).
- Blaisdell, B. E. Markov chain analysis finds a significant influence of neighboring bases on the occurrence of a base in eucaryotic nuclear DNA sequences both protein-coding and noncoding. *Journal of Molecular Evolution* **21**, 278–288 (1985).
- Mande, S. S., Mohammed, M. H. & Ghosh, T. S. Classification of metagenomic sequences: methods and challenges. *Briefings in Bioinformatics* **13**, 669–681 (2012).
- Lin, H.-H. & Liao, Y.-C. Accurate binning of metagenomic contigs via automated clustering sequences using information of genomic signatures and marker genes. *Scientific Reports* **6**, 24175 (2016).
- Laczny, C. C. *et al.* VizBin—an application for reference-independent visualization and human-augmented binning of metagenomic data. *Microbiome* **3**, 1 (2015).
- Alneberg, J. *et al.* Binning metagenomic contigs by coverage and composition. *Nature Methods* **11**, 1144–1146 (2014).
- Sharon, I. *et al.* Time series community genomics analysis reveals rapid shifts in bacterial species, strains, and phage during infant gut colonization. *Genome Research* **23**, 111–120 (2013).
- Kang, D. D., Froula, J., Egan, R. & Wang, Z. MetaBAT, an efficient tool for accurately reconstructing single genomes from complex microbial communities. *PeerJ* **3**, e1165 (2015).
- Golub, G. H. & Reinsch, C. Singular value decomposition and least squares solutions. *Numerische Mathematik* **14**, 403–420 (1970).
- Paukkeri, M.-S., Kivimäki, I., Tirunagari, S., Oja, E. & Honkela, T. Effect of dimensionality reduction on different distance measures in document clustering. *International Conference on Neural Information Processing*, 167–176 (Springer, 2011).
- Halko, N., Martinsson, P.-G. & Tropp, J. A. Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM Review* **53**, 217–288 (2011).
- Cleary, B. *et al.* Detection of low-abundance bacterial strains in metagenomic datasets by eigengenome partitioning. *Nature Biotechnology* **33**, 1053–1060 (2015).
- Van Der Maaten, L. Accelerating t-SNE using tree-based algorithms. *Journal of Machine Learning Research* **15**, 3221–3245 (2014).
- Lorenzo-Ginori, J. V., Rodriguez-Fuentes, A., Abalo, R. G. & Rodriguez, R. S. Digital signal processing in the analysis of genomic sequences. *Current Bioinformatics*. **4**, 28–40 (2009).
- Tapinos, A., Constantinides, B., Kell, D. B. & Robertson, D. L. Alignment by numbers: sequence assembly using compressed numerical representations. *bioRxiv* 011940 (2014).
- Wu, Y.-W., Tang, Y.-H., Tringe, S. G., Simmons, B. A. & Singer, S. W. MaxBin: an automated binning method to recover individual genomes from metagenomes using an expectation-maximization algorithm. *Microbiome* **2**, 26 (2014).
- Wu, Y.-W., Simmons, B. A. & Singer, S. W. Maxbin 2.0: an automated binning algorithm to recover genomes from multiple metagenomic datasets. *Bioinformatics* **32**, 605–607 (2015).
- Chatterji, S., Yamazaki, I., Bai, Z. & Eisen, J. A. CompostBin: A DNA composition-based algorithm for binning environmental shotgun reads. *Annual International Conference on Research in Computational Molecular Biology*, 17–28 (Springer, 2008).
- Kelley, D. R. & Salzberg, S. L. Clustering metagenomic sequences with interpolated markov models. *BMC Bioinformatics* **11**, 544 (2010).
- Kislyuk, A., Bhatnagar, S., Dushoff, J. & Weitz, J. S. Unsupervised statistical clustering of environmental shotgun sequences. *BMC Bioinformatics* **10**, 316 (2009).
- Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nature Methods* **9**, 357–359 (2012).
- Ester, M., Kriegel, H.-P., Sander, J. & Xu, X. A density-based algorithm for discovering clusters in large spatial databases with noise. *KDD* **96**, 226–231 (1996).
- Voss, R. F. Evolution of long-range fractal correlations and 1/f noise in DNA base sequences. *Physical Review Letters* **68**, 3805–3808 (1992).
- Ranawana, R. & Palade, V. A neural network based multi-classifier system for gene identification in DNA sequences. *Neural Computing and Applications* **14**, 122–131 (2005).
- Demeler, B. & Zhou, G. Neural network optimization for E. coli promoter prediction. *Nucleic Acids Research* **19**, 1593–1599 (1991).
- Nair, A. S. & Sreenadhan, S. P. A coding measure scheme employing electron-ion interaction pseudopotential (EIIP). *Bioinformation* **1**, 197–202 (2006).
- Bernaola-Galván, P., Carpena, P., Román-Roldán, R. & Oliver, J. Study of statistical correlations in DNA sequences. *Gene* **300**, 105–115 (2002).
- Holden, T. *et al.* ATCG nucleotide fluctuation of Deinococcus radiodurans radiation genes. *Optical Engineering + Applications*, 669417–669417 (International Society for Optics and Photonics, 2007).
- Ojala, T., Pietikäinen, M. & Harwood, D. A comparative study of texture measures with classification based on featured distributions. *Pattern Recognition* **29**, 51–59 (1996).
- Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).
- Li, H. A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. *Bioinformatics*. **27**, 2987–2993 (2011).
- Boisvert, S., Raymond, F., Godzaridis, E., Laviolette, F. & Corbeil, J. Ray Meta: scalable de novo metagenome assembly and profiling. *Genome Biology* **13**: R122 (2012).

## Acknowledgements

S.K. and A.T. were supported by the VIROGENESIS project. The VIROGENESIS project receives funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 634650. AT was also supported by BBSRC project grant, BB/M001121/1. We would like to thank Bede Constantinides for help with metagenomics data analysis.

## Author Contributions

S.K. designed and wrote the methods and software and performed the data analysis. A.T. provided useful comments on the nucleotide mapping and software design. S.K. and D.L.R. conceived the study. S.K. wrote the manuscript with comments from A.T. and D.L.R. All authors read and approved the final manuscript.

## Additional Information

**Supplementary information** accompanies this paper at <https://doi.org/10.1038/s41598-018-38197-9>.

**Competing Interests:** The authors declare no competing interests.

**Publisher's note:** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2019