

Consensus clustering for Bayesian mixture models

Stephen Coleman^{1*}, Paul D.W. Kirk^{1,2†} and Chris Wallace^{1,2†}

Correspondence:

stephen.coleman@mrc-

su.cam.ac.uk

MRC Biostatistics Unit

University of Cambridge,

Cambridge, UK

Full list of author information is

available at the end of the article

Equal contributor

Abstract

Background: Cluster analysis is an integral part of precision medicine and systems biology, used to define groups of patients or biomolecules. ~~However, problems such as choosing the number of clusters and issues with high dimensional data arise consistently. An ensemble approach, such as consensus clustering, can overcome some of the difficulties associated with high dimensional data, frequently exploring more relevant clustering solutions than individual models. Another tool for cluster analysis, Bayesian mixture modelling, has alternative advantages, including the ability to infer the number of clusters present and extensibility. However, inference of these models is often performed using Markov-chain Monte Carlo (MCMC) methods which can suffer from problems such as poor exploration of the posterior distribution and long runtimes. This makes applying~~ Consensus clustering is an ensemble approach that is widely used in these areas, which combines the output from multiple runs of a non-deterministic clustering algorithm. Here we consider the application of consensus clustering to a broad class of heuristic clustering algorithms that can be derived from Bayesian mixture models and their extensions to 'omics data challenging. We apply consensus clustering to Bayesian mixture models to address these problems. (and extensions thereof) by adopting an early stopping criterion when performing sampling-based inference for these models. While the resulting approach is non-Bayesian, it inherits the usual benefits of consensus clustering, particularly in terms of computational scalability and providing assessments of clustering stability/robustness.

Results: ~~Consensus clustering of Bayesian mixture models successfully finds the generating structure in our simulation study and captures multiple modes in the likelihood surface. This approach also~~ In simulation studies, we show that our approach can successfully uncover the target clustering structure, while also exploring different plausible clusterings of the data. We show that, when a parallel computation environment is available, our approach offers significant reductions in runtime compared to traditional Bayesian inference when a parallel environment is available, performing sampling-based Bayesian inference for the underlying model, while retaining many of the practical benefits of the Bayesian approach, such as exploring different numbers of clusters. We propose a heuristic to decide upon ensemble size and the early stopping criterion, and then apply

3 Background

From defining a taxonomy of disease to creating molecular sets, grouping items can help us to understand and make decisions using complex biological data. For example, grouping patients based upon disease characteristics and personal omics data may allow the identification of more homogeneous subgroups, enabling stratified medicine approaches. Defining and studying molecular sets can improve our understanding of biological systems as these sets are more interpretable than their constituent members (1), and study of their interactions and perturbations may have ramifications for diagnosis and drug targets (2, 3). The act of identifying such groups is referred to as “cluster analysis”, and has been traditional been done using tools cluster analysis. Many traditional methods such as K -means clustering (4, 5) ~~or hierarchical clustering~~. However, these methods have various problems condition upon a fixed choice of K , the number of clusters. These methods are often heuristic in nature, relying on rules of thumb to decide upon a final value for K . For example, ~~in different choices of K~~ K -means clustering, its sensitivity to initialisation means multiple runs are required, with that which minimises are compared under some metric such as silhouette (6) or the within-cluster sum of squared errors (**SSE**) ~~used (7)~~. This ~~problem arises as the algorithm has no guarantees on finding the global minimum of SSE~~ as a function of K . Moreover, K -means clustering can exhibit sensitivity to initialisation, necessitating multiple runs in practice (7).

Another common problem is that traditional methods offer no measure of the ~~uncertainty in stability or robustness of~~ the final clustering, ~~a quantity of interest in many analyses~~. Returning to the stratified medicine example of clustering patients, there might be individuals ~~with almost equal probability of being allocated between several clusters which might influence decisions made~~ that do not clearly belong to any one particular cluster; however if only a point estimate is obtained, this information is not available ~~to the decision-maker~~. Ensemble methods offer a solution to this

~~problem.~~ Ensemble methods address this problem, as well as reducing sensitivity to initialisation. These approaches have had great success in supervised learning, most famously in the form of Random Forest (8) and boosting (9). In clustering, consensus clustering (10) is a popular ~~ensemble~~-method which has been implemented in ~~R~~-R (11) and to a variety of methods (12, 13) and been applied to problems such as cancer subtyping (14, 15) and identifying subclones in single cell analysis (16).

Consensus clustering uses W runs of some base ~~model or learner~~ clustering algorithm (such as K -means~~clustering~~)~~and compiles the~~. These W proposed partitions ~~are commonly compiled~~ into a *consensus matrix*, the $(i, j)^{th}$ entries of which contain the proportion of model runs for which the i^{th} and j^{th} individuals co-cluster (for this and other definitions see section 1 of the Supplementary Material), ~~although this step is not fundamental to consensus clustering and there is a large body of literature aimed at interpreting a collection of partitions (see, e.g., 17–19)~~. This consensus matrix provides an assessment of the stability of the clustering. ~~This proportion represents some measure of confidence in the co-clustering of any pair of items.~~ Furthermore, ensembles can offer reductions in computational runtime ~~– This is as the individual learners can be weaker (and thus use either less because the individual members of the ensemble are often computationally inexpensive to fit (e.g. because they are fitted using only a subset of the available data or stop before full convergence)~~ and because the learners in most ensemble methods are independent of each other and thus enable use of a parallel environment for each of the quicker model runs (20).

Traditional clustering methods ~~usually~~ condition upon a fixed choice of K , the number of clusters ~~– Choosing with the choice of K is being~~ a difficult problem ~~that haunts many analyses with researchers often relying on rules of thumb to decide upon a final model choice.~~ For example, different choices of K are compared under ~~some metric such as silhouette or SSE as a function of K .~~ (10) proposed some in

itself. In consensus clustering, Monti *et al.* (10) proposed methods for choosing K using the consensus matrix, but this means that any of the uncertainty about and Ünlü *et al.* (21) offer an approach to estimating K is not represented in the final clustering and each model given the collection of partitions, but each clustering run uses the same, fixed, number of clusters. An alternative clustering approach, model-based clustering or mixture models mixture modelling, embeds the cluster analysis within a formal, statistical framework (22). This means that models can be compared formally, and problems such as the choice of K can be addressed as a model selection problem with all the associated tools. Mixture models are also attractive, as they have great flexibility in the type of data they can be applied to due to different choice of densities. Bayesian mixture models can (23). Moreover, Bayesian mixture models can be used to try to directly infer K , treating this as another random variable that is inferred from the data. This means that the final clustering is not conditional upon a user-chosen value, but K is jointly modelled along with the clustering. Such inference can be performed through use of a Dirichlet Process (24) mixture model (24, 25), a mixture of finite mixture models (26, 27) or an over-fitted mixture model (28). These models and their extensions have a history of successful application to a diverse range of biological problems such as finding clusters of gene expression profiles (29), cell types in flow cytometry (30, 31) or scRNAseq experiments (32), and estimating protein localisation (33). Bayesian mixture models can be extended to jointly model the clustering across multiple datasets (34, 35) (section 2 of the Supplementary Material).

However, performing inference of Bayesian mixture models is a difficult task. Variational inference (VI, 36) may be used to perform approximate inference of Bayesian mixture models (37), but while VI is powerful, it can struggle with multi-modality, underestimates the variance in the posterior distribution (38) and it has been shown to have a very computationally heavy initialisation cost

~~to have good results (39). Implementation is difficult, requiring either complex derivations — (see the Appendix Supplementary Methods of 40, for an example) or black-box, approximate solutions (41).~~

Markov chain Monte Carlo (MCMC) methods are the most common tool for performing computational Bayesian inference. In Bayesian clustering ~~methods~~, they are used to ~~construct a chain of clusterings and an assessment of the convergence of this chain is made to determine if its behaviour aligns with the expected asymptotic theory.~~ draw a collection of clustering partitions from the posterior distribution. However, in practice ~~individual chains often fail to explore the full support of the posterior distribution despite the ergodicity of MCMC methods~~, chains can become stuck in local posterior modes preventing convergence (see, e.g., the Supplementary Materials of 42) ~~and can experience long runtimes/~~ or can require prohibitively long runtimes, particularly when analysing high-dimensional datasets. Some MCMC methods make efforts to overcome the problem of exploration, often at the cost of increased computational cost per iteration. ~~See, e.g., (43, 44) for examples of problems and attempted solutions for MCMC methods.~~ (43). There are MCMC methods that use parallel chains to improve the scalability or reduce the bias of the Monte Carlo estimate. However, these methods have various limitations. For instance, divide-and-conquer strategies such as Asymptotically Exact, Embarrassingly Parallel MCMC (45) use subsamples of the dataset with each chain to improve scaling with the number of items being clustered. This assumes that each subsample is representative of the population, and is less helpful in situations where we have high-dimension but only moderate sample size, such as analysis of 'omics data. Alternative approaches, such as distributed MCMC (46) and coupling (47) have to account for burn-in bias; moreover, coupling further assumes the chains meet in finite time and then stay together. In practice, a further challenge associated

with these methods is that their implementation may necessitate a substantial redevelopment of existing software.

~~We propose that applying consensus clustering to Bayesian mixture models can overcome some of the issues endemic in high dimensional Bayesian clustering. (10) suggest this application as part of their original paper, but no investigation has been attempted to our knowledge. This ensemble approach sidesteps the problems of convergence associated MCMC methods and offers computational gains through using shorter chains run in parallel. Furthermore, this approach could be directly used on any existing MCMC based implementation of Bayesian mixture models or their extensions and would avoid the re-implementation process that changing to newer MCMC methods or VI would entail.~~

Motivated by the lack of scalability of existing implementations of sampling-based Bayesian clustering (due to prohibitive computational runtimes, as well as poor exploration, as described above), here we aim to develop a general and straightforward procedure that exploits the flexibility of these methods, but extends their applicability. Specifically, we make use of existing sampling-based Bayesian clustering implementations, but only run them for a fixed (and relatively small) number of iterations, stopping before they have converged to their target stationary distribution. Doing this repeatedly, we obtain an ensemble of clustering partitions, which we use to perform consensus clustering. We propose a heuristic for deciding upon the ensemble width-size (the number of learners used, W) and the ensemble depth (the number of iterations run within each chain, D), inspired by the use of scree plots in Principal Component Analysis (PCA; 48).

We show via simulation that ensembles consisting of short chains can be sufficient to successfully recover generating structure. We also show that consensus clustering explores as many or more modes of the likelihood surface than either standard Bayesian inference or Mclust, a maximum likelihood method, all while offering

~~improvements in runtime to traditional Bayesian inference.~~ our approach can
successfully identify meaningful clustering structures.
~~We use consensus clustering of~~ We then illustrate the use of our approach to
extend the applicability of existing Bayesian clustering implementations, using as a
case study the Multiple Dataset Iteration (~~MDI~~), a (~~MDI~~; 34) model for Bayesian
integrative clustering method, to analyse multiple applied to real data. While the
simulation results serve to validate our method, it is important to also evaluate
methods on real data which may represent more challenging problems. For our real
data, we use three omics datasets relating to the cell cycle of *Saccharomyces cere-*
~~visiae to show that consensus clustering can applied to more complex MCMC-based~~
~~clustering methods and real datasets, with the aim of inferring clusters of genes~~
across datasets. As there is no ground truth available, we then validate these clusters
using knowledge external to the analysis.

150 Material and methods

151 Consensus clustering for Bayesian mixture models

We apply consensus clustering to MCMC based Bayesian clustering models using
 the method described in algorithm 1. Our application of consensus clustering has
 two main parameters at the ensemble level, the chain depth, D , and ensemble
 width, W . We infer a point clustering from the consensus matrix using the `maxpear`
 function (49) from the R package `mcclust` (50) ~~to~~ which maximises the posterior
expected adjusted Rand index between the true clustering and point estimate if the
matrix is composed of samples drawn from the posterior distribution (section 3 of
 the Supplementary Material for details). There are alternative choices of methods
to infer a point estimate which minimise different loss functions (see, e.g., 51–53).

161 *Determining the ensemble depth and width*

As our ensemble sidesteps the problem of convergence within each chain, we need an
 alternative stopping rule for growing the ensemble in chain depth, D , and number

Data: $X = (x_1, \dots, x_N)$

Input:

The number of chains to run, W

The number of iterations within each chain, D

A clustering method that uses MCMC methods to generate samples of clusterings of the data $Cluster(X, d)$

Output:

A predicted clustering, \hat{Y}

The consensus matrix \mathbf{M}

```

begin
    /* initialise an empty consensus matrix */
     $\mathbf{M} \leftarrow \mathbf{0}_{N \times N}$ ;
    for  $w = 1$  to  $W$  do
        /* set the random seed controlling initialisation and MCMC
           moves */
         $set.seed(w)$ ;
        /* initialise a random partition on  $X$  drawn from the
           prior distribution */
         $Y_{(0,w)} \leftarrow Initialise(X)$ ;
        for  $d = 1$  to  $D$  do
            /* generate a markov chain for the membership vector */
             $Y_{(d,w)} \leftarrow Cluster(X, d)$ ;
        end
        /* create a coclustering matrix from the  $D^{th}$  sample */
         $\mathbf{B}^{(w)} \leftarrow Y_{(D,w)}$ ;
         $\mathbf{M} \leftarrow \mathbf{M} + \mathbf{B}^{(w)}$ ;
    end
     $\mathbf{M} \leftarrow \frac{1}{W} \mathbf{M}$ ;
     $\hat{Y} \leftarrow$  partition  $X$  based upon  $\mathbf{M}$ ;
end

```

Algorithm 1: Consensus clustering for Bayesian mixture models.

164 of chains, W . We propose a heuristic based upon the consensus matrix to decide
 165 if a given value of D and W are sufficient. We suspect that increasing W and D
 166 might continuously improve the performance of the ensemble, but we observe in
 167 our simulations that these ~~improvements~~ changes will become smaller and smaller
 168 for greater values, ~~approaching some asymptote~~ eventually converging for each of
 169 W and D . We notice that this behaviour is analogous to PCA in that where for
 170 consensus clustering some improvement might always be expected for increasing
 171 chain depth or ensemble width, more variance will ~~always~~ be captured by increasing
 172 the number of components used in PCA. However, increasing this number beyond
 173 some threshold has diminishing returns, diagnosed in PCA by a scree plot. Following
 174 from this, we recommend, for some set of ensemble parameters, $D' = \{d_1, \dots, d_I\}$
 175 and $W' = \{w_1, \dots, w_J\}$, find the mean absolute difference of the consensus matrix
 176 for the d_i^{th} iteration from w_j chains to that for the $d_{(i-1)}^{th}$ iteration from w_j chains
 177 and plot these values as a function of chain depth, and the analogue for sequential
 178 consensus matrices for increasing ensemble width and constant depth.

179 If this heuristic is used, we believe that the consensus matrix and the resulting
 180 inference should be stable (see, e.g., 54, 55), providing a robust estimate of the
 181 clustering. In contrast, if there is still strong variation in the consensus matrix
 182 for varying chain length or number, then we believe that the inferred clustering is
 183 influenced significantly by the random initialisation and that the inferred partition
 184 is unlikely to be stable for similar datasets or reproducible for a random choice of
 185 seeds.

186 Simulation study

We use a finite mixture with independent features as the data generating model
 within the simulation study. Within this model there exist “irrelevant features” (56)
 that have global parameters rather than cluster specific parameters ~~and use the~~

~~generating model~~:- The generating model is

$$\underline{p(X, c, \theta, \pi | K)} = \underline{p(K)p(\pi|K)p(\theta|K) \prod_{i=1}^N p(c_i|\pi, K) \prod_{p=1}^P p(x_{ip}|c_i, \theta_{c_{ip}}) \phi_p p(x_{ip}|\theta_p)(1 - \phi_p)}.$$

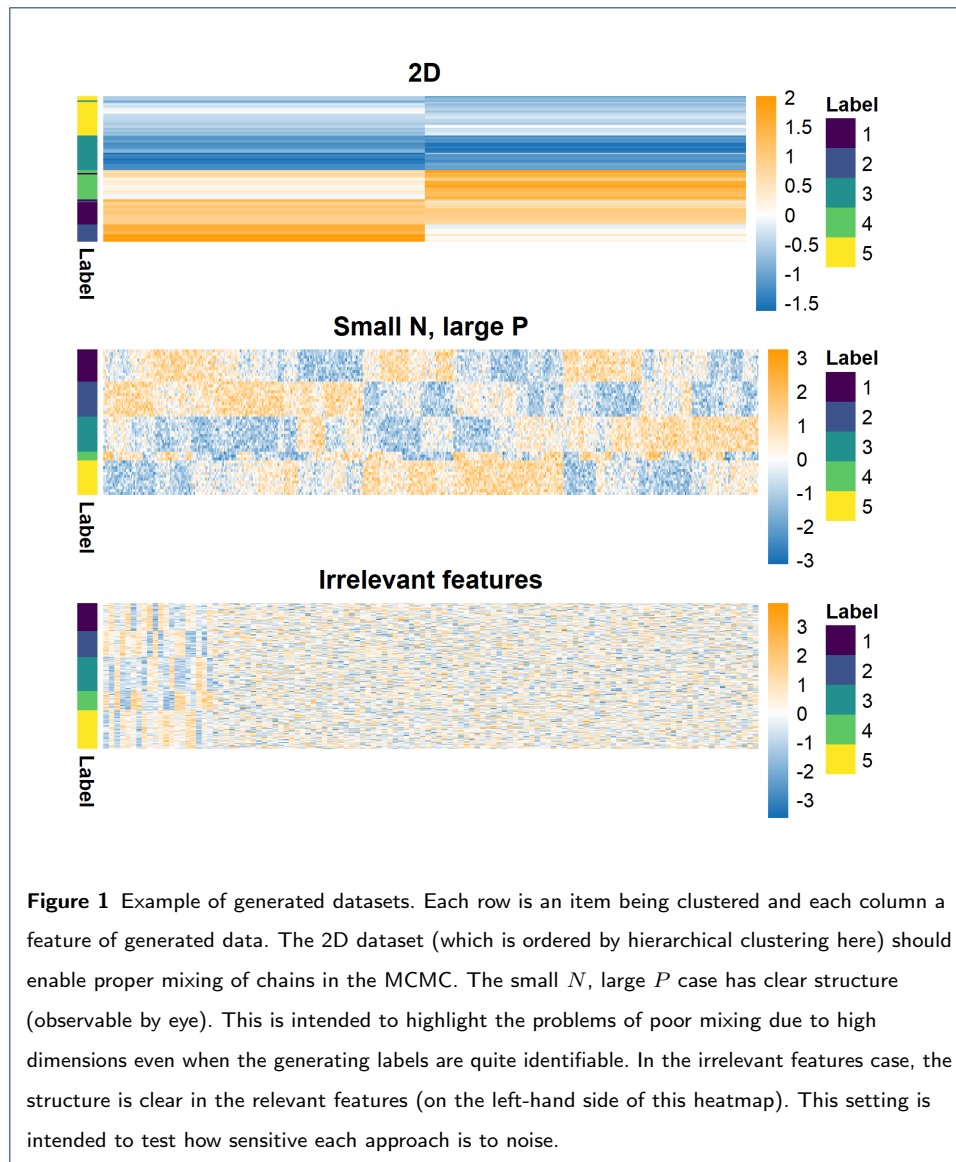
$$\underline{p(X, c, \theta, \pi | K)} = \underline{p(K)p(\pi|K)p(\theta|K) \prod_{i=1}^N p(c_i|\pi, K) \prod_{p=1}^P p(x_{ip}|c_i, \theta_{c_{ip}}) \phi_p p(x_{ip}|\theta_p)(1 - \phi_p)} \quad (1)$$

for data $X = (x_1, \dots, x_N)$, cluster label or allocation variable $c = (c_1, \dots, c_N)$,
cluster weight $\pi = (\pi_1, \dots, \pi_K)$, K clusters and the relevance variable, $\phi \in \{0, 1\}$
with $\phi_p = 1$ indicating that the p^{th} feature is relevant to the clustering. We used a
Gaussian density, so $\theta_{kp} = (\mu_{kp}, \sigma_{kp}^2)$. We defined three scenarios and simulated 100
datasets in each (~~Figure 1 and Table 1~~) figure 1 and table 1). Additional details of
the simulation process and additional scenarios are included in section 4.1 of the
Supplementary Materials.

Table 1 Parameters defining the simulation scenarios as used in generating data and labels. $\Delta\mu$ is the distance between neighbouring cluster means within a single feature. The number of relevant features (P_s) is $\sum_p \phi_p$, and $P_n = P - P_s$.

| Scenario | N | P_s | P_n | K | $\Delta\mu$ | σ^2 | π |
|---------------------|-----|-------|-------|-----|-------------|------------|---|
| 2D | 100 | 2 | 0 | 5 | 3.0 | 1 | $(\frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5})$ |
| Small N, large P | 50 | 500 | 0 | 5 | 1.0 | 1 | $(\frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5})$ |
| Irrelevant features | 200 | 20 | 100 | 5 | 1.0 | 1 | $(\frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5})$ |

In each of these scenarios we apply a variety of methods (~~listd~~ listed below) and
compare the inferred point clusterings to the generating labels using the Adjusted
Rand Index (**ARI**, 57).



- Mclust, a maximum likelihood implementation of ~~finite mixture models~~ a finite mixture of Gaussian densities (for a range of modelled clusters, K),
- 10 chains of 1 million iterations, thinning to every thousandth sample for the overfitted Bayesian mixture ~~model~~ of Gaussian densities, and
- A variety of consensus clustering ensembles defined by inputs of W chains and D iterations within each chain (see algorithm 1) with $W \in \{1, 10, 30, 50, 100\}$ and $D \in \{1, 10, 100, 1000, 10000\}$ where the base learner is an overfitted Bayesian mixture of Gaussian densities.

205 Note that none of the applied methods include a model selection step and as such
 206 there is no modelling of the relevant variables. This and the unknown value of K is
 207 what separates the models used and the generating model described in equation 1.
 208 More specifically, the likelihood of a point X_n for each method is

$$p(X_n|\mu, \Sigma, \pi) = \sum_{k=1}^K \pi_k p(X_n|\mu_k, \Sigma_k), \quad (2)$$

209 where $p(X_n|\mu_k, \Sigma_k)$ is the probability density function of the multivariate Gaussian
 210 distribution parameterised by a mean vector, μ_k , and a covariance matrix, Σ_k , and
 211 π_k is the component weight such that $\sum_{k=1}^K \pi_k = 1$. The implementation of the
 212 Bayesian mixture model restricts Σ_k to be a diagonal matrix while **Mclust** models
 213 a number of different covariance structures. Note that while we use the overfitted
 214 Bayesian mixture model, this is purely from convenience and we expect that a true
 215 Dirichlet Process mixture or a mixture of mixture models would display similar
 216 behaviour in an ensemble.

217 The ARI is a measure of similarity between two partitions, c_1, c_2 , corrected for
 218 chance, with 0 indicating c_1 is no more similar to c_2 than a random partition would
 219 be expected to be and a value of 1 showing that c_1 and c_2 perfectly align. Details of
 220 the methods in the simulation study can be found in sections 4.2, 4.3 and 4.4 of the
 221 Supplementary Material.

222 *Mclust*

223 **Mclust** (58) is a function from the R package **mclust**. It estimates Gaussian mixture
 224 models for K clusters based upon the maximum likelihood estimator of the param-
 225 eters. ~~it~~It initialises upon a hierarchical clustering of the data cut to K clusters.
 226 A range of choices of K and different covariance structures are compared and the

227 “best” selected using the Bayesian information criterion, (59) (details in section 4.2
228 of the Supplementary Material).

229 *Bayesian inference*

230 To assess within-chain convergence of our Bayesian inference we use the Geweke
231 Z -score statistic (60). Of the chains that appear to behave properly we then assess
232 across-chain convergence using \hat{R} (61) and the recent extension provided by (62).
233 If a chain has reached its stationary distribution the Geweke Z -score statistic is
234 expected to be normally distributed. Normality is tested for using a Shapiro-Wilks
235 test (63). If a chain fails this test (i.e., the associated p -value is less than 0.05), we
236 assume that it has not achieved stationarity and it is excluded from the remainder of
237 the analysis. The samples from the remaining chains are then pooled and a posterior
238 similarity matrix (**PSM**) constructed. We use the **maxpear** function to infer a point
239 clustering. For more details see section 4.3 of the Supplementary Material.

240 Analysis of the cell cycle in budding yeast

241 *Datasets*

242 The cell cycle is crucial to biological growth, repair, reproduction, and development
243 (64–66) and is highly conserved among eukaryotes (66). . This means that under-
244 standing of the cell cycle of *S. cerevisiae* can provide insight into a variety of cell
245 cycle perturbations including those that occur in human cancer (65, 67) and ageing
246 (68). We aim to create clusters of genes that are co-expressed~~in the cell cycle~~, have
247 common regulatory proteins and share a biological function. To achieve this, we use
248 three datasets that were generated using different 'omics technologies and target
249 different aspects of the molecular biology underpinning the cell cycle process.

- 250 • Microarray profiles of RNA expression from (69), comprising measurements of
251 cell-cycle-regulated gene expression at 5-minute intervals for 200 minutes (up
252 to three cell division cycles) and is referred to as the **time course** dataset.
253 The cells are synchronised at the START checkpoint in late G1-phase using

alpha factor arrest (69). We include only the genes identified by (69) as having periodic expression profiles.

- Chromatin immunoprecipitation followed by microarray hybridization (**ChIP-chip**) data from (70). This dataset discretizes p -values from tests of association between 117 DNA-binding transcriptional regulators and a set of yeast genes. Based upon a significance threshold these p -values are represented as either a 0 (no interaction) or a 1 (an interaction).
- Protein-protein interaction (**PPI**) data from BioGrid (71). This database consists of physical and genetic interactions between gene and gene products, with interactions either observed in high throughput experiments or computationally inferred. The dataset we used contained 603 proteins as columns. An entry of 1 in the $(i, j)^{th}$ cell indicates that the i^{th} gene has a protein product that is believed to interact with the j^{th} protein.

The datasets were reduced to the 551 genes with no missing data in the PPI and ChIP-chip data, as in (34).

Multiple dataset integration

We applied consensus clustering to MDI for our integrative analysis. Details of MDI are in section 2.2 of the Supplementary Material, but in short MDI jointly models the clustering in each dataset, inferring individual clusterings for each dataset. These partitions are informed by similar structure in the other datasets, with MDI learning this similarity as it models the partitions. The model does not assume global structure. This means that the similarity between datasets is not strongly assumed in our model; individual clusters or genes that align across datasets are based solely upon the evidence present in the data and not due to strong modelling assumptions. Thus, datasets that share less common information can be included without fearing that this will warp the final clusterings in some way.

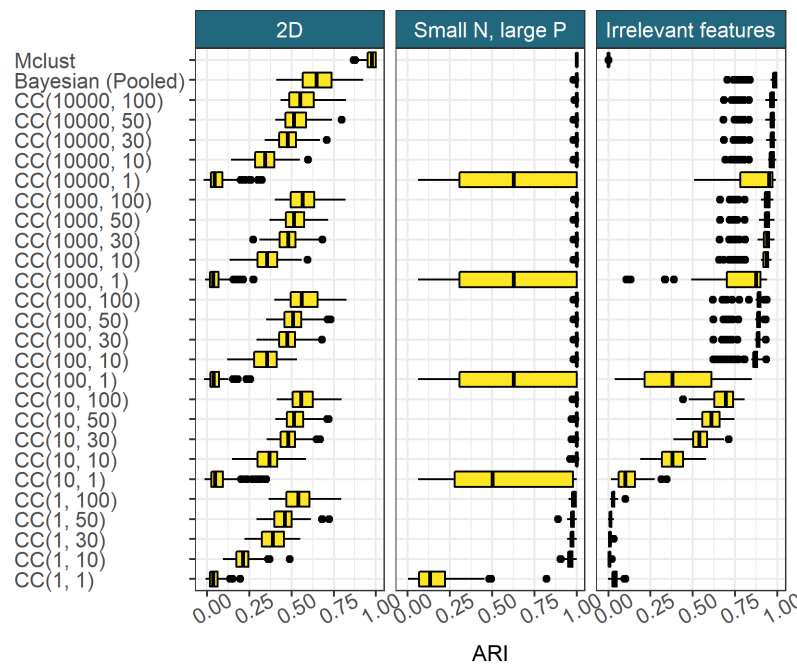


Figure 2 Model performance in the 100 simulated datasets for each scenario, defined as the ARI between the generating labels and the inferred clustering. $CC(d, w)$ denotes consensus clustering using the clustering from the d^{th} iteration from w different chains.

280 The datasets were modelled using a mixture of Gaussian processes in the time
 281 course dataset and Multinomial distributions in the ChIP-chip and PPI datasets.

282 Results

283 Simulated data

284 We use the ARI between the generating labels and the inferred clustering of each
 285 method to be our metric of predictive performance. In [Figure-figure 2](#), we see Mclust
 286 performs very well in the 2D and Small N , large P scenarios, correctly identifying the
 287 true structure. However, the irrelevant features scenario sees a collapse in performance,
 288 Mclust is blinded by the irrelevant features and identifies a clustering of $K = 1$.

289 The pooled samples from multiple long chains performs very well across all scenarios
 290 and appears to act as an upper bound on the more practical implementations of
 291 consensus clustering.

292 Consensus clustering does uncover some of the generating structure in the data,
293 even using a small number of short chains. With sufficiently large ensembles and
294 chain depth, consensus clustering is close to the pooled Bayesian samples in predictive
295 performance. It appears that for a constant chain depth increasing the ensemble
296 width used follows a pattern of diminishing returns. There are strong initial gains
297 for a greater ensemble width, but the improvement decreases for each successive
298 chain. A similar pattern emerges in increasing chain length for a constant number
299 of chains (~~Figure~~[figure 2](#)).

300 We see very little difference between the similarity matrix from the pooled samples
301 and the consensus clustering (~~Figure~~[figure 3](#)). Similar clusters emerge, and we see
302 comparable confidence in the pairwise clusterings. For the PSMs from the individual
303 chains, all entries are 0 or 1. This means only a single clustering is sampled within
304 each chain, implying very little uncertainty in the partition. However, three different
305 modes emerge across the chains showing that the chains are failing to explore the full
306 support of the posterior distribution of the clustering and are each unrepresentative
307 of the uncertainty in the final clustering. This shows that consensus clustering is
308 exploring more possible clusterings than any individual chain and, as it explores a
309 similar space to the pooled samples which might be considered more representative
310 of the posterior distribution than any one chain, it suggests it better describes the
311 true uncertainty present than any single chain. It also shows that pooling chains
312 offers robustness to multi-modality (as expected for an ensemble) and the ARI for
313 the pooled samples is an upper bound on the performance for the individual long
314 chains.

315 Figure 4 shows that chain length is directly proportional to the time taken for
316 the chain to run. This means that using an ensemble of shorter chains, as in
317 consensus clustering, can offer large reductions in the time cost of analysis when a
318 parallel environment is available compared to standard Bayesian inference. Even

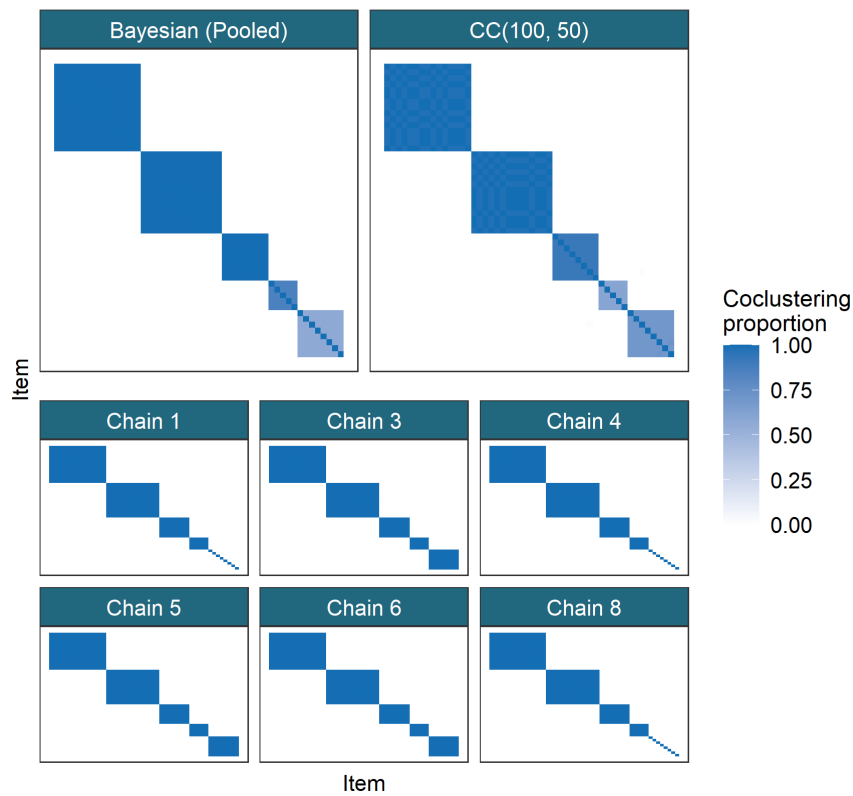


Figure 3 Comparison of similarity matrices from a dataset for the Small N , large P scenario. In each matrix, the $(i, j)^{th}$ entry is the proportion of clusterings for which the i^{th} and j^{th} items co-clustered for the method in question. In the first row the PSM of the pooled Bayesian samples is compared to the CM for CC(100, 50), with a common ordering of rows and columns in both heatmaps. In the following rows, 6 of the long chains that passed the tests of convergence are shown.

319 on a laptop of 8 cores running an ensemble of 1,000 chains of length 1,000 will
 320 require approximately half as much time as running 10 chains of length 100,000
 321 due to parallelisation, and the potential benefits are far greater when using a large
 322 computing cluster.

323 Additional results for these and other simulations are in section 4.4 of the Supple-
 324 mentary Material.

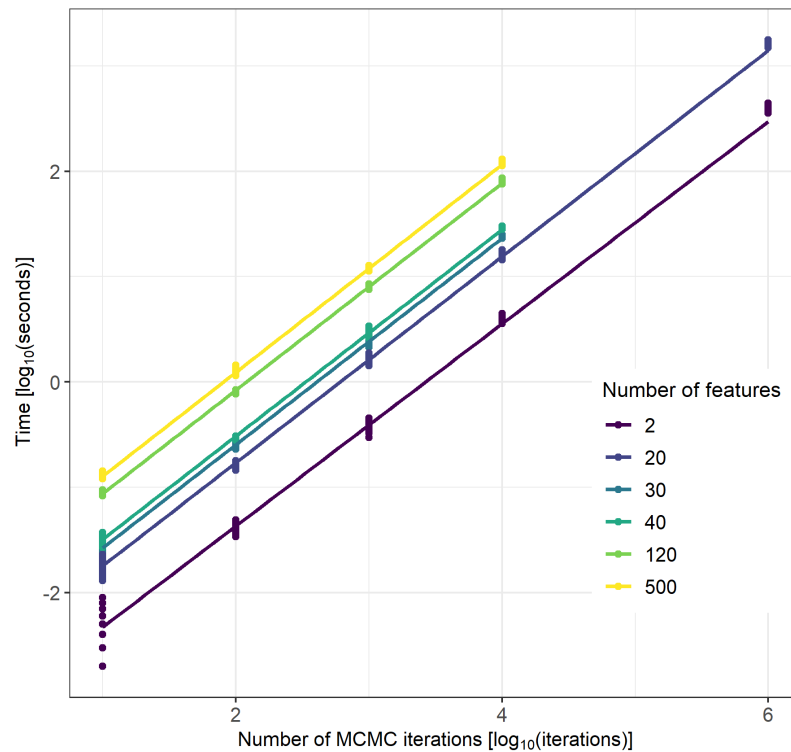
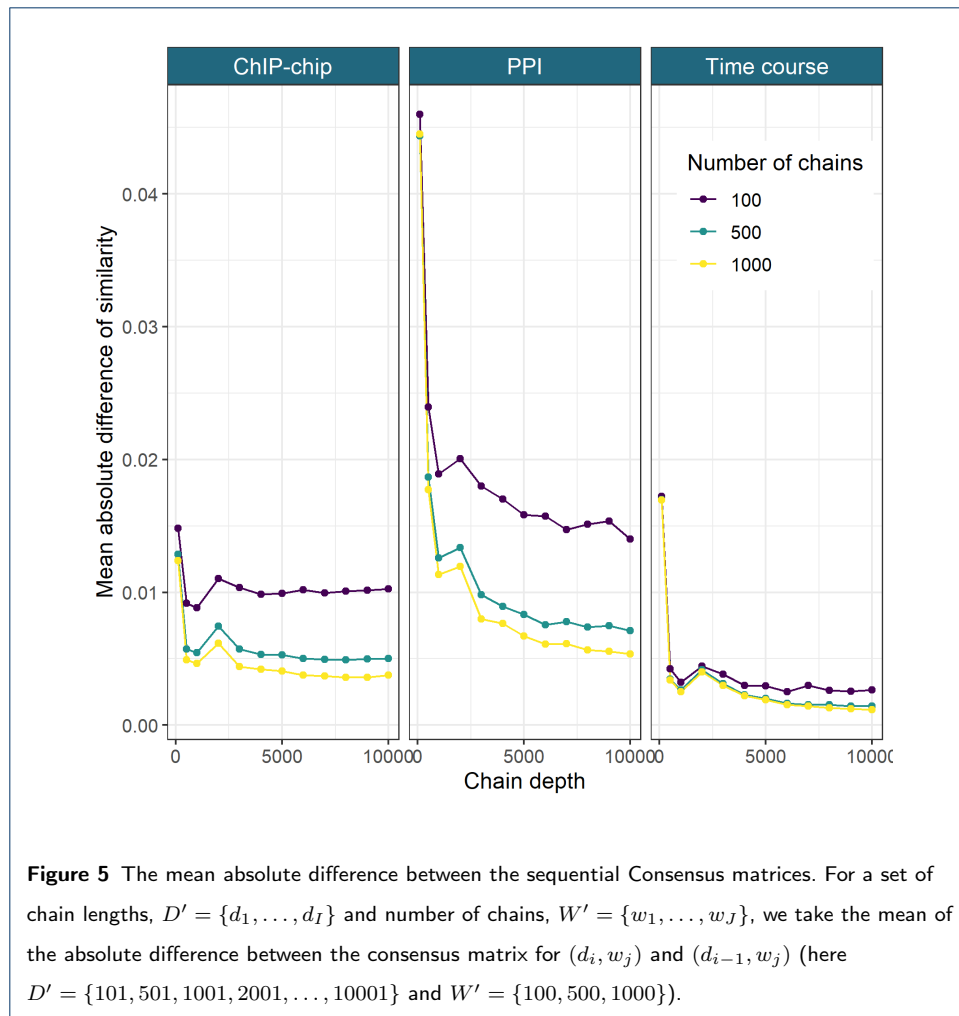


Figure 4 The time taken for different numbers of iterations of MCMC moves in $\log_{10}(\text{seconds})$. The relationship between chain length, D , and the time taken is linear (the slope is approximately 1 on the \log_{10} scale), with a change of intercept for different dimensions. The runtime of each Markov chain was recorded using the terminal command `time`, measured in milliseconds.

325 Multi-omics analysis of the cell cycle in budding yeast

326 We use the stopping rule proposed in to determine our ensemble depth and width. In
 327 [Figure-figure 5](#), we see that the change in the consensus matrices from increasing the
 328 ensemble depth and width is diminishing in keeping with results in the simulations.
 329 We see no strong improvement after $D = 6,000$ and increasing the number of learners
 330 from 500 to 1,000 has small effect. We therefore use the largest ensemble available, a
 331 depth $D = 10001$ and width $W = 1000$, believing this ensemble is stable (additional
 332 evidence in section 5.1 of the Supplementary Material).



333 We focus upon the genes that tend to have the same cluster label across multiple
 334 datasets. More formally, we analyse the clustering structure among genes for which
 335 $\hat{P}(c_{nl} = c_{nm}) > 0.5$, where c_{nl} denotes the cluster label of gene n in dataset l .

336 In our analysis it is the signal shared across the time course and ChIP-chip
 337 datasets that is strongest, with 261 genes (nearly half of the genes present) in this
 338 pairing tending to have a common label, whereas only 56 genes have a common
 339 label across all three datasets. Thus, we focus upon this pairing of datasets in
 340 the results of the analysis performed using all three datasets. We show the gene
 341 expression and regulatory proteins of these genes separated by their cluster in [Figure](#)
 342 [figure 6](#). In [Figure-figure 6](#), the clusters in the time series data have tight, unique
 343 signatures (having different periods, amplitudes, or both) and in the ChIP-chip

344 data clusters are defined by a small number of well-studied transcription factors (TFs)
 345 ~~(see Table 2 of the Supplementary Material for details of these TFs, many of which are well known to regulate ce~~
 346 (see table 2 of the Supplementary Material for details of these TFs, many of which are well known to regulate cel
 347 .

348 As an example, we briefly analyse clusters 9 and 16 in greater depth. Cluster 9 has
 349 strong association with MBP1 and some interactions with SWI6, as can be seen in
 350 ~~Figure-figure 6~~. The Mbp1-Swi6p complex, MBF, is associated with DNA replication
 351 (73). The first time point, 0 minutes, in the time course data is at the START
 352 checkpoint, or the G1/S transition. The members of cluster 9 begin highly expressed
 353 at this point before quickly dropping in expression (in the first of the 3 cell cycles).
 354 This suggests that many transcripts are produced immediately in advance of S-phase,
 355 and thus are required for the first stages of DNA synthesis. These genes' descriptions
 356 ~~(found using org.Sc.sgd.db, 74, and shown in Table 3 of the Supplementary Material)~~
 357 (found using org.Sc.sgd.db, 74, and shown in table 3 of the Supplementary Material)
 358 support this hypothesis, as many of the members are associated with DNA repli-
 359 cation, repair and/or recombination. Additionally, *TOF1*, *MRC1* and *RAD53*,
 360 members of the replication checkpoint (75, 76) emerge in the cluster as do members
 361 of the cohesin complex. Cohesin is associated with sister chromatid cohesion which
 362 is established during the S-phase of the cell cycle (77) and also contributes to
 363 transcription regulation, DNA repair, chromosome condensation, homolog pairing
 364 (78), fitting the theme of cluster 9.

365 Cluster 16 appears to be a cluster of S-phase genes, consisting of *GAS3*, *NRM1*
 366 and *PDS1* and the genes encoding the histones H1, H2A, H2B, H3 and H4. Histones
 367 are the chief protein components of chromatin (79) and are important contributors
 368 to gene regulation (80). They are known to peak in expression in S-phase (69),
 369 which matches the first peak of this cluster early in the time series. Of the other
 370 members, *NRM1* is a transcriptional co-repressor of MBF-regulated gene expression

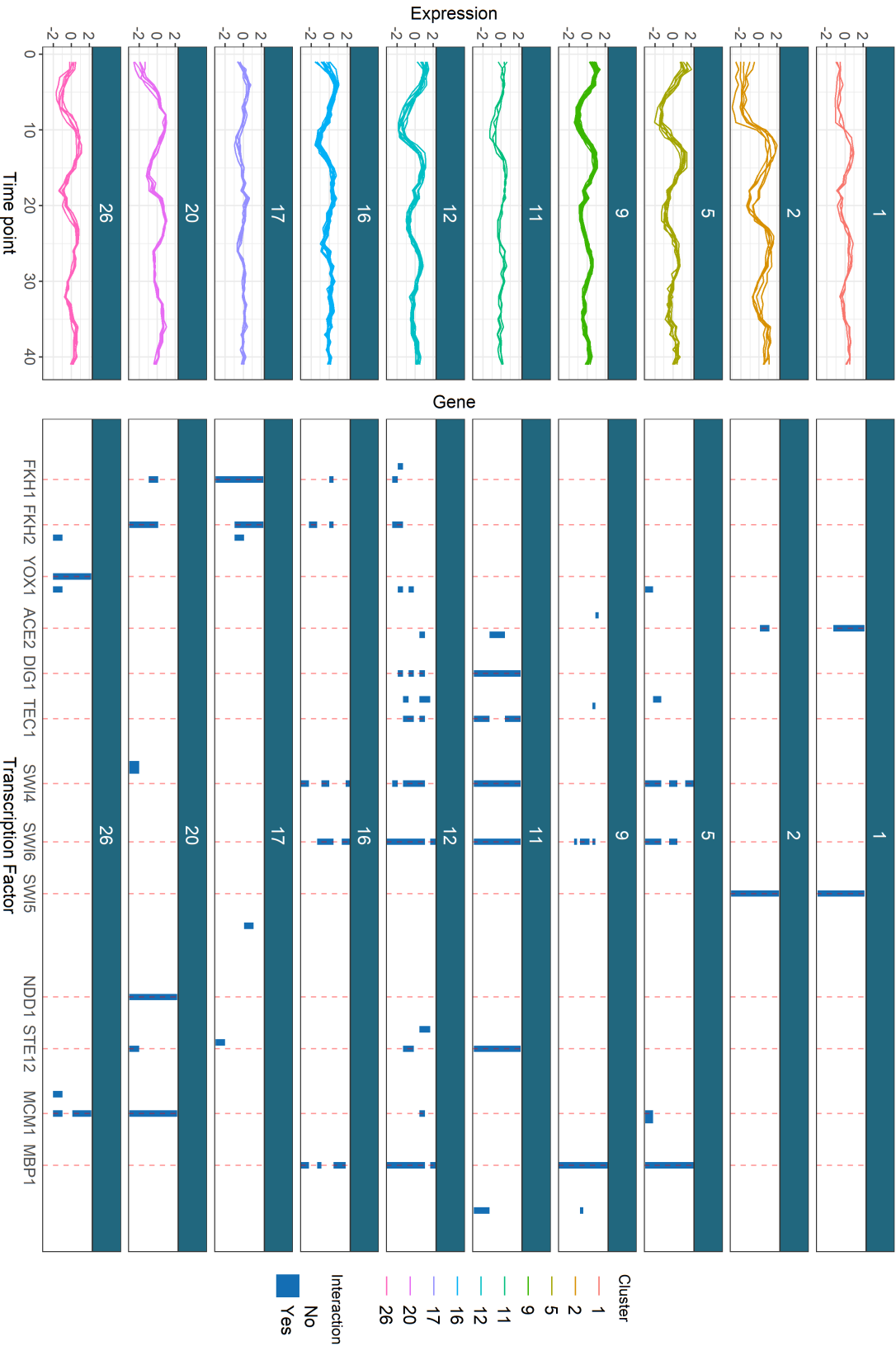


Figure 6 The gene clusters which tend to have a common label across the time course and ChIP-chip datasets, shown in these datasets. We include only the clusters with more than one member and more than half the members having some interactions in the ChIP-chip data. Red lines for the most common transcription factors are included.

371 acting at the transition from G1 to S-phase (81, 82). Pds1p binds to and inhibits
 372 the Esp1 class of sister separating proteins, preventing sister chromatids separation
 373 before M-phase (77, 83). *GAS3*, is not well studied. It interacts with *SMT3* which
 374 regulates chromatid cohesion, chromosome segregation and DNA replication (among
 375 other things). Chromatid cohesion ensures the faithful segregation of chromosomes
 376 in mitosis and in both meiotic divisions (84) and is instantiated in S-phase (77).
 377 These results, along with the very similar expression profile to the histone genes in
 378 the time course data, suggest that *GAS3* may be more directly involved in DNA
 379 replication or chromatid cohesion than is currently believed.

380 We attempt to perform a similar analysis using traditional Bayesian inference of
 381 MDI, but after 36 hours of runtime there is no consistency or convergence across
 382 chains. We use the Geweke statistic and \hat{R} to reduce to the five best behaved
 383 chains (none of which appear to be converged, see section 5.2 of the Supplementary
 384 Material for details). If we then compare the distribution of sampled values for
 385 the ϕ parameters for these long chains, the final ensemble used ($D = 10001$, W
 386 $= 1000$) and the pooled samples from the 5 long chains, then we see that the
 387 distribution of the pooled samples from the long chains (which might be believed
 388 to sampling different parts of the posterior distribution) is closer in appearance
 389 to the distributions sampled by the consensus clustering than to any single chain
 390 (figure 7). Further disagreement between chains is shown in the Gene Ontology term
 391 over-representation analysis in section 5.3 of the Supplementary Material.

392 Discussion

393 Our proposed method has demonstrated good performance on simulation studies,
 394 uncovering the generating structure ~~and approximating Bayesian inference when~~
 395 ~~the Markov chain is exploring the full support of the posterior distribution~~ in many
 396 cases and performing comparably to Mclust and long chains in many scenarios.
 397 We saw that when the chains are sufficiently deep that the ensemble approximates

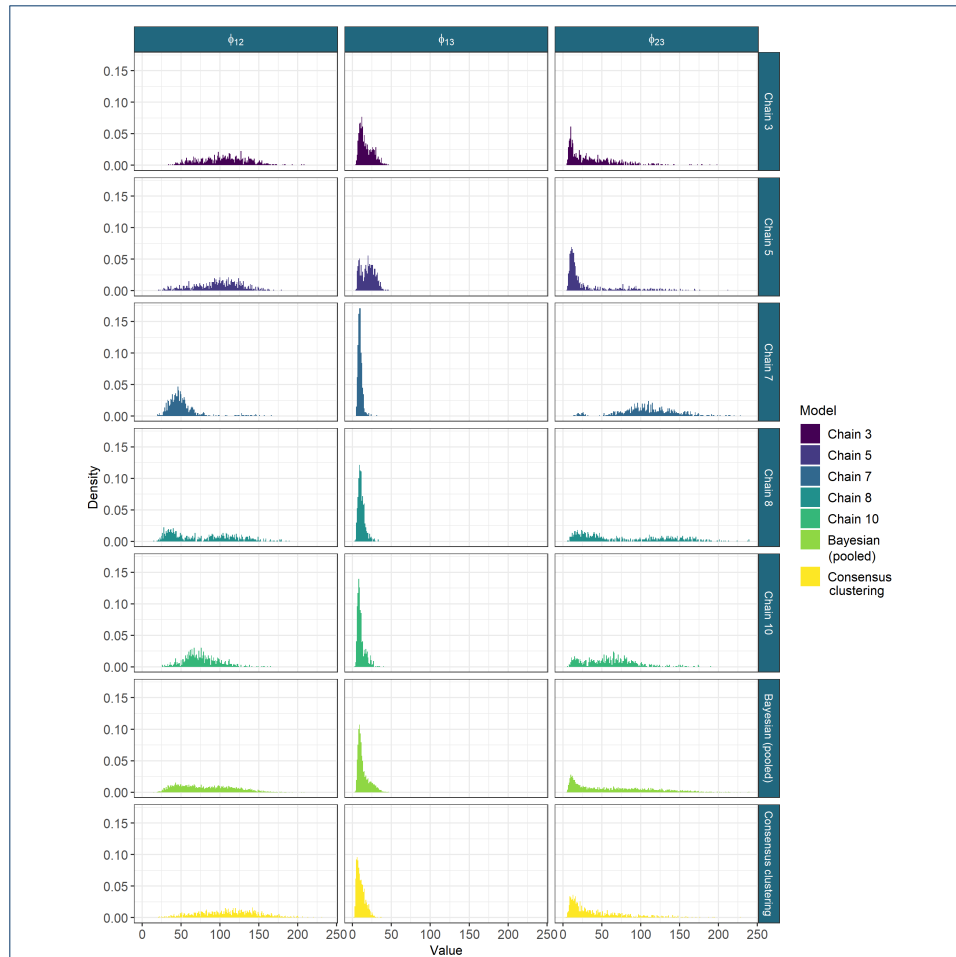


Figure 7 The sampled values for the ϕ parameters from the long chains, their pooled samples and the consensus using 1000 chains of depth 10,001. The long chains display a variety of behaviours. Across chains there is no clear consensus on the nature of the posterior distribution. The samples from any single chain are not particularly close to the behaviour of the pooled samples across all three parameters. It is the consensus clustering that most approaches this pooled behaviour.

398 Bayesian inference, as shown by the similarity between the PSMs and the CM in
 399 the 2D scenario where the individual chains do not become trapped in a single
 400 mode. However, we have shown that if a finite Markov chain fails to describe
 401 the full posterior and is itself only approximating Bayesian inference distribution,
 402 our method frequently has better ability to represent several modes in the data
 403 than individual chains and thus offers a more consistent and reproducible analysis.
 404 Furthermore, consensus clustering We also showed that the ensemble of short chains

405 is more robust to irrelevant features than Mclust. Furthermore, an ensemble of
406 short chains is significantly faster in a parallel environment than inference using
407 individual ~~chains, while retaining the ability to robustly infer the number of occupied~~
408 ~~components present~~long chains.

409 We proposed a method of assessing ensemble stability and deciding upon ensemble
410 size which we used when performing an integrative analysis of yeast cell cycle data
411 using MDI, an extension of Bayesian mixture models that jointly models multiple
412 datasets. We uncovered many genes with shared signal across several datasets and
413 explored the meaning of some of the inferred clusters ~~;~~ using data external to
414 the analysis. We found ~~sensible~~ biologically meaningful results as well as signal for
415 possibly novel biology. ~~In contrast, the traditional approach to Bayesian inference~~
416 ~~failed here. The lack of a consistent distribution across the chains made proceeding~~
417 ~~with the Bayesian analysis difficult as choosing the result of any single chain over~~
418 ~~the others would be arbitrary and thus prone to irreproducibility. The alternative of~~
419 ~~pooling the samples, which might be considered a reasonable compromise, appears~~
420 ~~to offer a very similar solution to consensus clustering, but with longer runtime and~~
421 ~~additional steps to reduce the chains to the “best-behaved” chains. We believe that~~
422 ~~the similarity between the sampled distribution of the parameters from the pooled~~
423 ~~long chains and the consensus clustering of short chains, figure 7, suggests that~~
424 ~~sufficiently deep chains within the ensemble can be used even to perform inference~~
425 ~~of continuous variables and not only the latent clustering of the data.~~ We also showed
426 that individual chains for the existing implementation of MDI do not converge in
427 a practical length of time, having run 10 chains for 36 hours with no consistent
428 behaviour across chains. This means that Bayesian inference of the MDI model is
429 not practical on this dataset with the software currently available.

430 ~~The results of our simulations and the multi-omics analysis show that consensus~~
431 ~~clustering can be successfully used in a broad context, being applicable to any~~

~~MCMC-based clustering method. It offers computational gains and improves the exploration of the clustering space, overcoming the problem of becoming trapped in specific, local extrema of the likelihood surface that emerges in high-dimensional data. This enables the application of these methods in modern 'omics datasets and, attractively, consensus clustering can be applied to existing implementations, unlike improvements to the underlying MCMC methods or alternative methods for Bayesian inference such as VI which would require re-writing software. However, consensus clustering does lose the theoretical framework of true Bayesian inference. We attempt to mitigate this with our assessment of stability in the ensemble, but this diagnosis is heuristic and subjective, and while there is empirical evidence for its success, it lacks the formal results for the tests of model convergence for Bayesian inference.~~

More generally, we have benchmarked the use of an ensemble of Bayesian mixture models, showing that this approach can infer meaningful clusterings and overcomes the problem of multi-modality in the likelihood surface even in high dimensions, thereby providing more stable clusterings than individual long chains that are prone to becoming trapped in individual modes. We also show that the ensemble can be significantly quicker to run. In our multi-omics study we have demonstrated that the method can be applied as a wrapper to more complex Bayesian clustering methods using existing implementations and that this provides meaningful results even when individual chains fail to converge. This enables greater application of complex Bayesian clustering methods without requiring re-implementation using more clever MCMC methods, a process that would involve a significant investment of human time.

We expect that researchers interested in applying some of the Bayesian integrative clustering models such as MDI and Clusternomics (35) will be enabled to do so, as consensus clustering overcomes some of the unwieldiness of existing implementations

of these complex models. More generally, we expect that our method will be useful to researchers performing cluster analysis of high-dimensional data where the runtime of MCMC methods becomes too onerous and multi-modality is more likely to be present.

Funding

This work was funded by the MRC (MC UU 00002/4, MC UU 00002/13) and supported by the NIHR Cambridge Biomedical Research Centre (BRC-1215-20014). The views expressed are those of the author(s) and not necessarily those of the NHS, the NIHR or the Department of Health and Social Care. This research was funded in whole, or in part, by the Wellcome Trust [WT107881]. For the purpose of Open Access, the author has applied a CC BY public copyright licence to any Author Accepted Manuscript version arising from this submission.

Abbreviations

ARI: Adjusted Rand Index

ChIP-chip: Chromatin immunoprecipitation followed by microarray hybridization

CM: Consensus Matrix

MCMC: Markov chain Monte Carlo

MDI: Multiple Dataset Integration

PCA: Principal Component Analysis

PPI: Protein-Protein Interaction

PSM: Posterior Similarity Matrix

SSE: Sum of Squared Errors

TF: Transcription Factor

Availability of data and materials

The code and datasets supporting the conclusions of this article are available in the github repository, <https://github.com/stcolema/ConsensusClusteringForBayesianMixtureModels>.

Competing interests

The authors declare that they have no competing interests.

Authors' contributions

SC designed the simulation study with contributions from PK and CW, performed the analyses and wrote the manuscript. PK and CW provided an equal contribution of joint supervision, directing the research and provided suggestions such as the stopping rule. All contributed to interpreting the results of the analyses. All authors revised and approved the final manuscript.

Author details

¹MRC Biostatistics Unit University of Cambridge, Cambridge, UK. ²Cambridge Institute of Therapeutic Immunology & Infectious Disease, University of Cambridge, Cambridge, UK.

References

- Hejblum BP, Skinner J, Thiébaud R. Time-course gene set analysis for longitudinal gene expression data. *PLoS computational biology*. 2015;11(6):e1004310.

- 497 2. Bai JP, Alekseyenko AV, Statnikov A, Wang IM, Wong PH. Strategic applications of gene expression: from
498 drug discovery/development to bedside. *The AAPS journal*. 2013;15(2):427–437.
- 499 3. Emmert-Streib F, Dehmer M, Haibe-Kains B. Gene regulatory networks and their applications:
500 understanding biological and medical problems in terms of networks. *Frontiers in cell and developmental*
501 *biology*. 2014;2:38.
- 502 4. Lloyd S. Least squares quantization in PCM. *IEEE transactions on information theory*. 1982;28(2):129–137.
- 503 5. Forgy EW. Cluster analysis of multivariate data: efficiency versus interpretability of classifications.
504 *biometrics*. 1965;21:768–769.
- 505 6. Rousseeuw PJ. Silhouettes: A Graphical Aid to the Interpretation and Validation of Cluster Analysis.
506 *Journal of Computational and Applied Mathematics*. 1987 Nov;20:53–65.
- 507 7. Arthur D, Vassilvitskii S. k-means++: The advantages of careful seeding. Stanford; 2006.
- 508 8. Breiman L. Random forests. *Machine learning*. 2001;45(1):5–32.
- 509 9. Friedman JH. Stochastic gradient boosting. *Computational statistics & data analysis*. 2002;38(4):367–378.
- 510 10. Monti S, Tamayo P, Mesirov J, Golub T. Consensus clustering: a resampling-based method for class
511 discovery and visualization of gene expression microarray data. *Machine learning*. 2003;52(1-2):91–118.
- 512 11. Wilkerson, D M, Hayes, Neil D. ConsensusClusterPlus: a class discovery tool with confidence assessments
513 and item tracking. *Bioinformatics*. 2010;26(12):1572–1573.
- 514 12. John CR, Watson D, Russ D, Goldmann K, Ehrenstein M, Pitzalis C, et al. M3C: Monte Carlo
515 reference-based consensus clustering. *Scientific reports*. 2020;10(1):1–14.
- 516 13. Gu Z, Schlesner M, Hübschmann D. cola: an R/Bioconductor package for consensus partitioning through
517 a general framework. *Nucleic Acids Research*. 2020 12;Gkaa1146. Available from:
518 <https://doi.org/10.1093/nar/gkaa1146>.
- 519 14. Lehmann BD, Bauer JA, Chen X, Sanders ME, Chakravarthy AB, Shyr Y, et al. Identification of human
520 triple-negative breast cancer subtypes and preclinical models for selection of targeted therapies. *The*
521 *Journal of clinical investigation*. 2011;121(7):2750–2767.
- 522 15. Verhaak RG, Hoadley KA, Purdom E, Wang V, Qi Y, Wilkerson MD, et al. Integrated genomic analysis
523 identifies clinically relevant subtypes of glioblastoma characterized by abnormalities in PDGFRA, IDH1,
524 EGFR, and NF1. *Cancer cell*. 2010;17(1):98–110.
- 525 16. Kiselev VY, Kirschner K, Schaub MT, Andrews T, Yiu A, Chandra T, et al. SC3: consensus clustering of
526 single-cell RNA-seq data. *Nature methods*. 2017;14(5):483–486.
- 527 17. Li T, Ding C. Weighted Consensus Clustering. In: *Proceedings of the 2008 SIAM International Conference*
528 *on Data Mining*. Society for Industrial and Applied Mathematics; 2008. p. 798–809.
- 529 18. Carpineto C, Romano G. Consensus Clustering Based on a New Probabilistic Rand Index with Application
530 to Subtopic Retrieval. *IEEE Transactions on Pattern Analysis and Machine Intelligence*. 2012
531 Dec;34(12):2315–2326.
- 532 19. Strehl A, Ghosh J. Cluster Ensembles – A Knowledge Reuse Framework for Combining Multiple Partitions.
533 *Journal of Machine Learning Research*. 2002;3:583–617.
- 534 20. Ghaemi R, Sulaiman MN, Ibrahim H, Mustapha N, et al. A survey: clustering ensembles techniques. *World*
535 *Academy of Science, Engineering and Technology*. 2009;50:636–645.
- 536 21. Ünlü R, Xanthopoulos P. Estimating the Number of Clusters in a Dataset via Consensus Clustering.
537 *Expert Systems with Applications*. 2019 Jul;125:33–39.
- 538 22. Fraley C, Raftery AE. Model-based clustering, discriminant analysis, and density estimation. *Journal of the*
539 *American statistical Association*. 2002;97(458):611–631.

- 540 23. Fraley C. How Many Clusters? Which Clustering Method? Answers Via Model-Based Cluster Analysis. *The*
 541 *Computer Journal*. 1998 Aug;41(8):578–588.
- 542 24. Ferguson TS. A Bayesian analysis of some nonparametric problems. *The annals of statistics*. 1973;p.
 543 209–230.
- 544 25. Antoniak CE. Mixtures of Dirichlet Processes with Applications to Bayesian Nonparametric Problems. *The*
 545 *Annals of Statistics*. 1974 Nov;2(6):1152–1174.
- 546 26. Richardson S, Green PJ. On Bayesian analysis of mixtures with an unknown number of components.
 547 *Journal of the Royal Statistical Society: series B*. 1997;59(4):731–792.
- 548 27. Miller JW, Harrison MT. Mixture models with a prior on the number of components. *Journal of the*
 549 *American Statistical Association*. 2018;113(521):340–356.
- 550 28. Rousseau J, Mengersen K. Asymptotic behaviour of the posterior distribution in overfitted mixture models.
 551 *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*. 2011;73(5):689–710.
- 552 29. Medvedovic M, Sivaganesan S. Bayesian infinite mixture model based clustering of gene expression profiles.
 553 *Bioinformatics*. 2002;18(9):1194–1206.
- 554 30. Chan C, Feng F, Ottinger J, Foster D, West M, Kepler TB. Statistical mixture modeling for cell subtype
 555 identification in flow cytometry. *Cytometry Part A: The Journal of the International Society for Analytical*
 556 *Cytology*. 2008;73(8):693–701.
- 557 31. Hejblum BP, Alkhassim C, Gottardo R, Caron F, Thiébaud R, et al. Sequential Dirichlet process mixtures
 558 of multivariate skew *t*-distributions for model-based clustering of flow cytometry data. *The Annals of*
 559 *Applied Statistics*. 2019;13(1):638–660.
- 560 32. Prabhakaran S, Azizi E, Carr A, Pe'er D. Dirichlet process mixture model for correcting technical variation
 561 in single-cell gene expression data. In: *International Conference on Machine Learning*; 2016. p. 1070–1079.
- 562 33. Crook OM, Mulvey CM, Kirk PD, Lilley KS, Gatto L. A Bayesian mixture modelling approach for spatial
 563 proteomics. *PLoS computational biology*. 2018;14(11):e1006516.
- 564 34. Kirk P, Griffin JE, Savage RS, Ghahramani Z, Wild DL. Bayesian correlated clustering to integrate
 565 multiple datasets. *Bioinformatics*. 2012;28(24):3290–3297.
- 566 35. Gabasova E, Reid J, Wernisch L. Clusternomics: Integrative context-dependent clustering for
 567 heterogeneous datasets. *PLoS computational biology*. 2017;13(10):e1005781.
- 568 36. Blei DM, Jordan MI, et al. Variational inference for Dirichlet process mixtures. *Bayesian analysis*.
 569 2006;1(1):121–143.
- 570 37. Martin GM, Frazier DT, Robert CP. Computing Bayes: Bayesian Computation from 1763 to the 21st
 571 Century. *arXiv preprint arXiv:200406425*. 2020;.
- 572 38. Turner RE, Sahani M. Two problems with variational expectation maximisation for time-series models. In:
 573 Barber D, Cemgil AT, Chiappa S, editors. *Bayesian time series models*. 1st ed. Cambridge University Press;
 574 2011. .
- 575 39. Wang L, Dunson DB. Fast Bayesian inference in Dirichlet process mixture models. *Journal of*
 576 *Computational and Graphical Statistics*. 2011;20(1):196–216.
- 577 40. Argelaguet R, Velten B, Arnol D, Dietrich S, Zenz T, Marioni JC, et al. Multi-Omics Factor Analysis—a
 578 framework for unsupervised integration of multi-omics data sets. *Molecular systems biology*.
 579 2018;14(6):e8124.
- 580 41. Kucukelbir A, Tran D, Ranganath R, Gelman A, Blei DM. Automatic differentiation variational inference.
 581 *The Journal of Machine Learning Research*. 2017;18(1):430–474.

- 582 42. Strauss ME, Kirk PD, Reid JE, Wernisch L. GPseudoClust: deconvolution of shared pseudo-profiles at
583 single-cell resolution. *Bioinformatics*. 2020;36(5):1484–1491.
- 584 43. Robert CP, Elvira V, Tawn N, Wu C. Accelerating MCMC algorithms. *Wiley Interdisciplinary Reviews:*
585 *Computational Statistics*. 2018;10(5):e1435.
- 586 44. Yao Y, Vehtari A, Gelman A. Stacking for non-mixing Bayesian computations: the curse and blessing of
587 multimodal posteriors. *arXiv preprint arXiv:200612335*. 2020;.
- 588 45. Neiswanger W, Wang C, Xing E. Asymptotically Exact, Embarrassingly Parallel MCMC. *arXiv:13114780*
589 *[cs, stat]*. 2014 Mar;.
- 590 46. Murray L. Distributed Markov Chain Monte Carlo. In: *Proceedings of Neural Information Processing*
591 *Systems workshop on learning on cores, clusters and clouds*. vol. 11; 2010. .
- 592 47. Jacob PE, O'Leary J, Atchadé YF. Unbiased Markov Chain Monte Carlo Methods with Couplings. *Journal*
593 *of the Royal Statistical Society: Series B (Statistical Methodology)*. 2020;82(3):543–600.
- 594 48. Wold S, Esbensen K, Geladi P. Principal component analysis. *Chemometrics and intelligent laboratory*
595 *systems*. 1987;2(1-3):37–52.
- 596 49. Fritsch A, Ickstadt K, et al. Improved criteria for clustering based on the posterior similarity matrix.
597 *Bayesian analysis*. 2009;4(2):367–391.
- 598 50. Fritsch A. mcclust: process an MCMC sample of clusterings; 2012. R package version 1.0. Available from:
599 <https://CRAN.R-project.org/package=mcclust>.
- 600 51. Wade S, Ghahramani Z. Bayesian Cluster Analysis: Point Estimation and Credible Balls (with Discussion).
601 *Bayesian Analysis*. 2018 Jun;13(2):559–626.
- 602 52. Lourenço A, Rota Bulò S, Rebagliati N, Fred ALN, Figueiredo MAT, Pelillo M. Probabilistic Consensus
603 Clustering Using Evidence Accumulation. *Machine Learning*. 2015 Jan;98(1):331–357.
- 604 53. Dahl DB, Johnson DJ, Mueller P. Search Algorithms and Loss Functions for Bayesian Clustering.
605 *arXiv:210504451 [stat]*. 2021 May;.
- 606 54. Von Luxburg U, Ben-David S. Towards a statistical theory of clustering. In: *Pascal workshop on statistics*
607 *and optimization of clustering*. Citeseer; 2005. p. 20–26.
- 608 55. Meinshausen N, Bühlmann P. Stability selection. *Journal of the Royal Statistical Society: Series B*
609 *(Statistical Methodology)*. 2010;72(4):417–473.
- 610 56. Law MH, Jain AK, Figueiredo M. Feature selection in mixture-based clustering. In: *Advances in neural*
611 *information processing systems*; 2003. p. 641–648.
- 612 57. Hubert L, Arabie P. Comparing partitions. *Journal of classification*. 1985;2(1):193–218.
- 613 58. Scrucca L, Fop M, Murphy BT, Raftery AE. mclust 5: clustering, classification and density estimation
614 using Gaussian finite mixture models. *The R Journal*. 2016;8(1):289–317. Available from:
615 <https://doi.org/10.32614/RJ-2016-021>.
- 616 59. Schwarz G, et al. Estimating the dimension of a model. *The annals of statistics*. 1978;6(2):461–464.
- 617 60. Geweke J, et al. Evaluating the accuracy of sampling-based approaches to the calculation of posterior
618 moments. vol. 196. Federal Reserve Bank of Minneapolis, Research Department Minneapolis, MN; 1991.
- 619 61. Gelman A, Rubin DB, et al. Inference from iterative simulation using multiple sequences. *Statistical*
620 *science*. 1992;7(4):457–472.
- 621 62. Vats D, Knudson C. Revisiting the Gelman-Rubin diagnostic. *arXiv preprint arXiv:181209384*. 2018;.
- 622 63. Shapiro SS, Wilk MB. An analysis of variance test for normality (complete samples). *Biometrika*.
623 1965;52(3/4):591–611.

- 624 64. Tyson JJ, Chen KC, Novák B. Cell Cycle, Budding Yeast. In: Dubitzky W, Wolkenhauer O, Cho KH,
 625 Yokota H, editors. Encyclopedia of Systems Biology. New York, NY: Springer New York; 2013. p. 337–341.
- 626 65. Chen KC, Calzone L, Csikasz-Nagy A, Cross FR, Novak B, Tyson JJ. Integrative analysis of cell cycle
 627 control in budding yeast. *Molecular biology of the cell*. 2004;15(8):3841–3862.
- 628 66. Alberts B, Johnson A, Lewis J, Raff M, Roberts K, Walter P. The cell cycle and programmed cell death.
 629 *Molecular biology of the cell*. 2002;4:983–1027.
- 630 67. Ingalls B, Duncker B, Kim D, McConkey B. Systems level modeling of the cell cycle using budding yeast.
 631 *Cancer informatics*. 2007;3:117693510700300020.
- 632 68. Jiménez J, Bru S, Ribeiro M, Clotet J. Live fast, die soon: cell cycle progression and lifespan in yeast cells.
 633 *Microbial Cell*. 2015;2(3):62.
- 634 69. Granovskaia MV, Jensen LJ, Ritchie ME, Toedling J, Ning Y, Bork P, et al. High-resolution transcription
 635 atlas of the mitotic cell cycle in budding yeast. *Genome biology*. 2010;11(3):1–11.
- 636 70. Harbison CT, Gordon DB, Lee TI, Rinaldi NJ, Macisaac KD, Danford TW, et al. Transcriptional
 637 regulatory code of a eukaryotic genome. *Nature*. 2004;431(7004):99–104.
- 638 71. Stark C, Breitkreutz BJ, Regul T, Boucher L, Breitkreutz A, Tyers M. BioGRID: a general repository for
 639 interaction datasets. *Nucleic acids research*. 2006;34(suppl_1):D535–D539.
- 640 72. Simon I, Barnett J, Hannett N, Harbison CT, Rinaldi NJ, Volkert TL, et al. Serial regulation of
 641 transcriptional regulators in the yeast cell cycle. *Cell*. 2001;106(6):697–708.
- 642 73. Iyer VR, Horak CE, Scafe CS, Botstein D, Snyder M, Brown PO. Genomic binding sites of the yeast
 643 cell-cycle transcription factors SBF and MBF. *Nature*. 2001;409(6819):533–538.
- 644 74. Carlson M, Falcon S, Pages H, Li N. Org. sc. sgdb: Genome wide annotation for yeast. R package
 645 version. 2014;2(1).
- 646 75. Bando M, Katou Y, Komata M, Tanaka H, Itoh T, Sutani T, et al. Csm3, Tof1, and Mrc1 form a
 647 heterotrimeric mediator complex that associates with DNA replication forks. *Journal of Biological*
 648 *Chemistry*. 2009;284(49):34355–34365.
- 649 76. Lao JP, Ulrich KM, Johnson JR, Newton BW, Vashisht AA, Wohlschlegel JA, et al. The yeast DNA
 650 damage checkpoint kinase Rad53 targets the exoribonuclease, Xrn1. *G3: Genes, Genomes, Genetics*.
 651 2018;8(12):3931–3944.
- 652 77. Tóth A, Ciosk R, Uhlmann F, Galova M, Schleiffer A, Nasmyth K. Yeast cohesin complex requires a
 653 conserved protein, Eco1p (Ctf7), to establish cohesion between sister chromatids during DNA replication.
 654 *Genes & development*. 1999;13(3):320–333.
- 655 78. Mehta GD, Kumar R, Srivastava S, Ghosh SK. Cohesin: functions beyond sister chromatid cohesion.
 656 *FEBS letters*. 2013;587(15):2299–2312.
- 657 79. Fischle W, Wang Y, Allis CD. Histone and chromatin cross-talk. *Current opinion in cell biology*.
 658 2003;15(2):172–183.
- 659 80. Bannister AJ, Kouzarides T. Regulation of chromatin by histone modifications. *Cell research*.
 660 2011;21(3):381–395.
- 661 81. de Bruin RA, Kalashnikova TI, Chahwan C, McDonald WH, Wohlschlegel J, Yates III J, et al. Constraining
 662 G1-specific transcription to late G1 phase: the MBF-associated corepressor Nrm1 acts via negative
 663 feedback. *Molecular cell*. 2006;23(4):483–496.
- 664 82. Aligianni S, Lackner DH, Klier S, Rustici G, Wilhelm BT, Marguerat S, et al. The fission yeast
 665 homeodomain protein Yox1p binds to MBF and confines MBF-dependent cell-cycle transcription to G1-S
 666 via negative feedback. *PLoS Genet*. 2009;5(8):e1000626.

83. Ciosk R, Zachariae W, Michaelis C, Shevchenko A, Mann M, Nasmyth K. An ESP1/PDS1 complex regulates loss of sister chromatid cohesion at the metaphase to anaphase transition in yeast. *Cell*. 1998;93(6):1067–1076.
84. Cooper KF, Mallory MJ, Guacci V, Lowe K, Strich R. Pds1p is required for meiotic recombination and prophase I progression in *Saccharomyces cerevisiae*. *Genetics*. 2009;181(1):65–79.

Additional Files

Additional file 1 — Supplementary materials

Additional relevant theory, background and results. This includes some more formal definitions, details of Bayesian mixture models and MDI, the general consensus clustering algorithm, additional simulations and the generating algorithm used, steps in assessing Bayesian model convergence in both the simulated datasets and yeast analysis, a table of the transcription factors that define the clustering in the ChIP-chip dataset, a table of the gene descriptions for some of the clusters that emerge across the time course and ChIP-chip datasets and Gene Ontology term over-representation analysis of the clusterings from the yeast datasets. (PDF, 10MB)