

# Epiclomal

Probabilistic Clustering of  
Sparse DNA Methylation Data

# Motivation

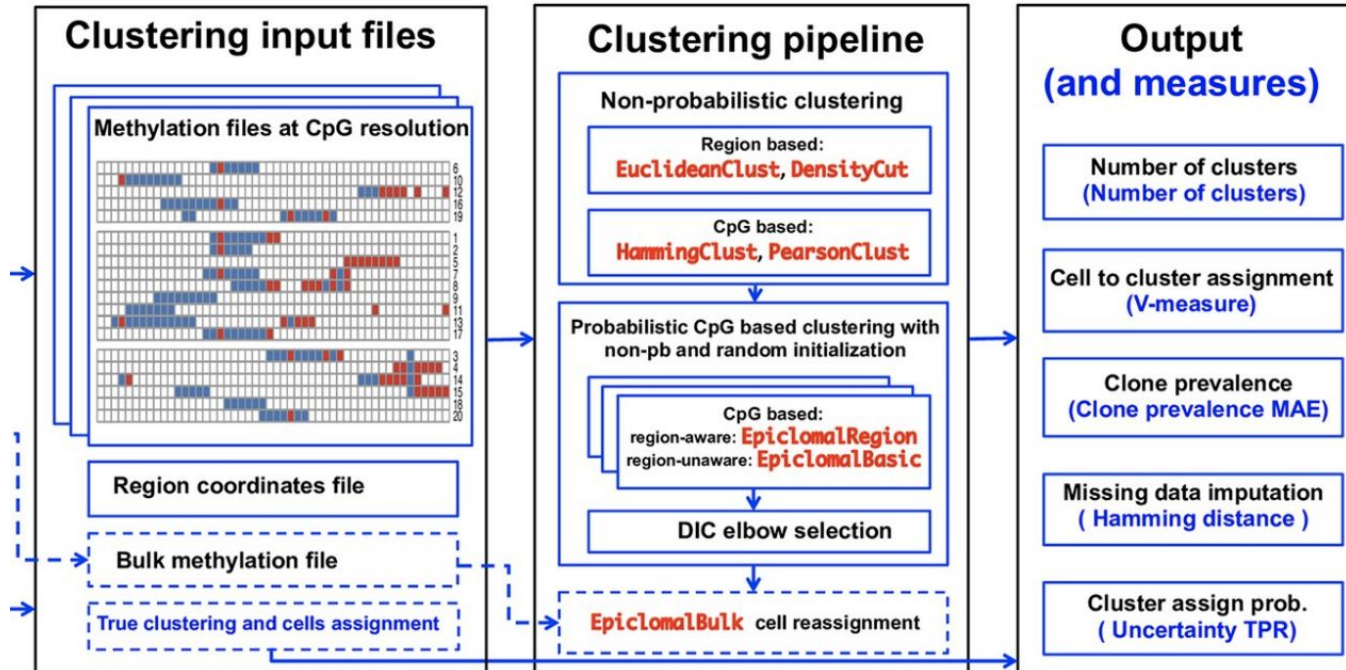
- DNA methylation at 5mC position is important for transcriptional regulation
- Single-cell whole-genome bisulfite sequencing (sc-WGBS) can assess epigenetic diversity of a cell population
- Data is sparse (lots of missing CpG sites) and subject to error
- Need for clustering according to methylation profiles
  - Identification of cancer subtypes
  - Detection of previously unknown cell types, deeper characterization of known ones
  - Imputation of missing CpG data by pooling information across clusters

# Previous Work

- A number of non-probabilistic clustering methods for CpG methylation data
  - Hou et al., Farlik et al., others
- Probabilistic approaches to imputation of missing CpG data
  - Kapourani and Sanguinetti, Angermuller et al.
- Until now: no probabilistic method with a focus on clustering based on methylation profiles
- Authors want to simultaneously cluster sc-WGBS data while inferring the missing methylation states

# Epiclomal

- Goal: cluster sparse CpG-based DNA methylation data from sc-WGBS



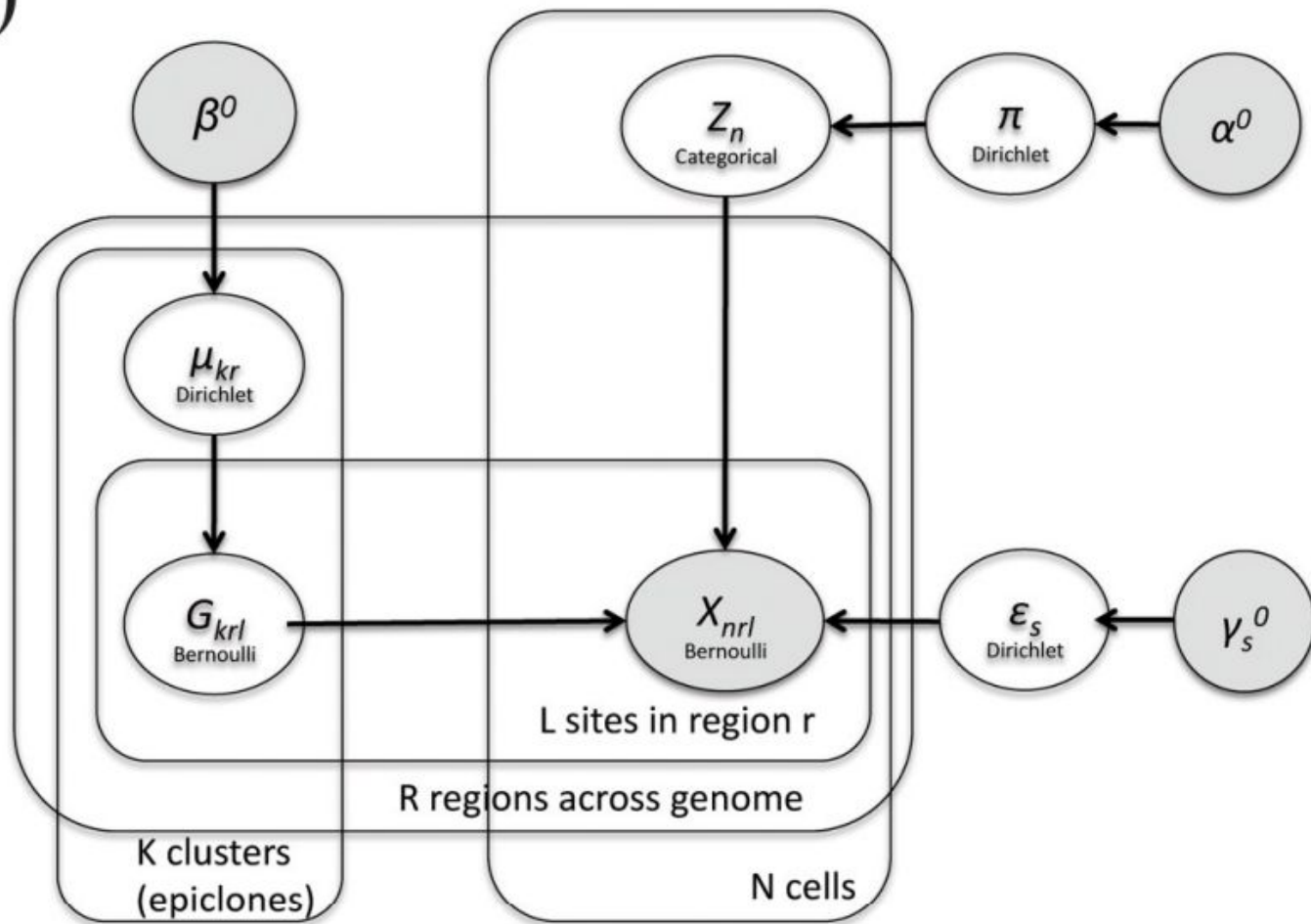
# Input

- sparse matrix of  $N$  rows (cells) and  $M$  columns (CpG sites)
- each entry is 0 (unmethylated), 1 (methylated), or missing
- partially methylated sites: 1 if methylation fraction is at least 0.5, 0 otherwise
  - median < 1.35% across datasets
- optional: select specific regions
  - EpiclomalBasic vs. EpiclomalRegion

# Methods: Non-Probabilistic Methods

- Performance comparisons
- Cluster initialization
- Region-based: EuclideanClust, DensityCut
  - Based on mean methylation for each region
  - DensityCut also infers the optimal number of clusters
- CpG-based: HammingClust, PearsonClust
  - Based on methylation of individual CpG sites
- Epiclomal considers each potential number of clusters up to a chosen  $K$ 
  - Infers optimal number of clusters (at most  $K$ )

(b)



Let  $\mathbf{X}_{nr} = (X_{nr1}, \dots, X_{nrL_r})^T$  be the vector of observed data for region  $r$  in cell  $n$ , and let  $\mathbf{X}_n = (\mathbf{X}_{nr}^T, \dots, \mathbf{X}_{nR}^T)^T$  be the vector of all observed data for cell  $n$ . Assume that  $X_{nr1}, \dots, X_{nrL_r}$  are independent for all  $n$  and  $r$ . Suppose that there are  $K \ll N$  vectors of true hidden methylation states shared across the cells. Let  $Z_n$  with values in  $\{1, \dots, K\}$  be the hidden variable indicating the true cluster (epiclonal) population of cell  $n$ . It is assumed that  $Z_1, \dots, Z_N$  are independent with  $P(Z_n = k) = \pi_k$  such that  $\sum_{k=1}^K \pi_k = 1$ . If  $Z_n = k$ , then the distribution of  $\mathbf{X}_n$  depends on the  $k$ -th vector of true hidden epigenotypes  $\mathbf{G}_k = (\mathbf{G}_{k1}^T, \dots, \mathbf{G}_{kR}^T)^T$ , where  $\mathbf{G}_{kr} = (G_{kr1}, \dots, G_{krL_r})^T$ . We assume that  $G_{kr1}, \dots, G_{krL_r}$  are independent for all  $k$  and  $r$ , with  $P(G_{krl} = s) = \mu_{krs}$  such that  $\sum_{s \in \mathcal{S}} \mu_{krs} = 1$ , that is,  $G_{krl}$  follows a categorical (Bernoulli) distribution with parameter set  $\boldsymbol{\mu}_{kr} = \{\mu_{krs} : s \in \mathcal{S}\}$ . Therefore, given the true cluster assignment and the corresponding true hidden methylation states, the observed data  $\mathbf{X}_{nr}$  are independent, with  $X_{nr1}$  following a categorical distribution with parameters depending on the hidden true state at locus  $l$  of region  $r$  for cluster population  $k$ , that is,

$$P(X_{nr1} = t | Z_n = k, G_{krl} = s) = \epsilon_{st} \text{ with } \sum_{t \in \mathcal{S}} \epsilon_{st} = 1. \quad (1)$$

We can also interpret the probability in (1) as a misclassification error, which in this context is related to sequencing error.

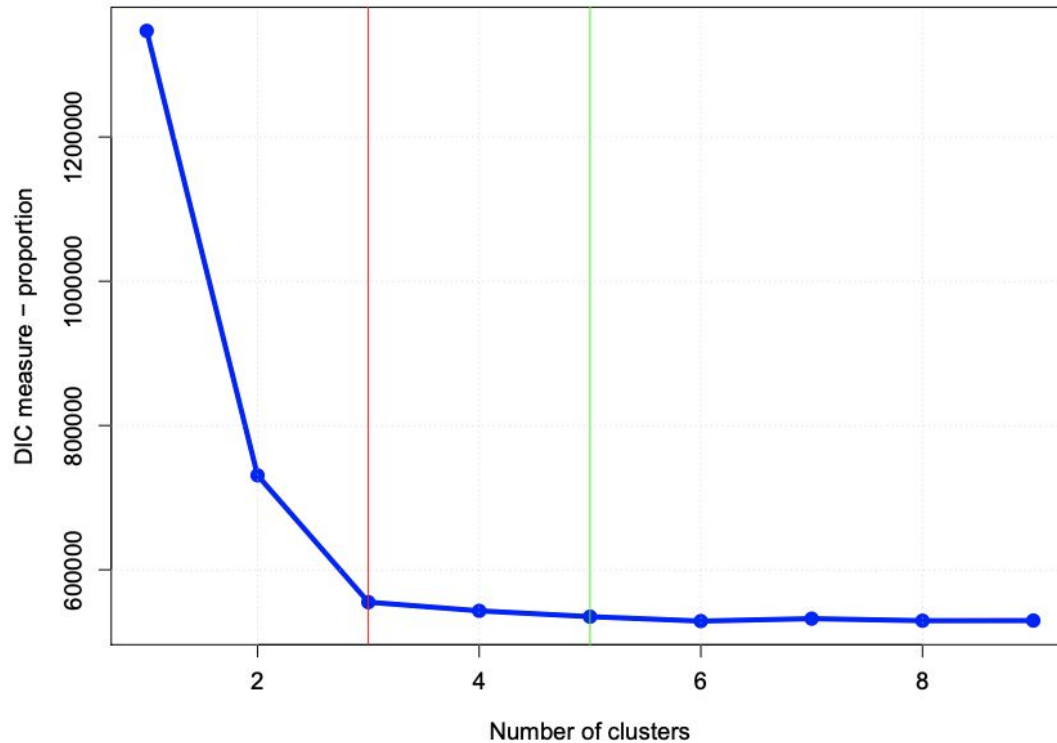


To infer  $\Theta$  and the hidden states  $\mathbf{Z} = (Z_1, \dots, Z_n)^T$  and  $\mathbf{G} = \{\mathbf{G}_1, \dots, \mathbf{G}_K\}$ , we adopt a Bayesian approach and derive a Variational Bayes (VB) algorithm [33] to approximate the posterior distribution of  $\Theta$ ,  $\mathbf{Z}$ , and  $\mathbf{G}$  given the observed data  $\mathbf{X} = \{\mathbf{X}_1, \dots, \mathbf{X}_N\}$ ,  $P(\mathbf{Z}, \mathbf{G}, \Theta|\mathbf{X})$  by finding the Variational Distribution (VD),  $q(\mathbf{Z}, \mathbf{G}, \Theta)$  with the smallest Kullback-Leibler divergence to the posterior  $P(\mathbf{Z}, \mathbf{G}, \Theta|\mathbf{X})$ , which is equivalent to maximizing the evidence lower bound (ELBO) given by

$$\text{ELBO}(q) = \mathbb{E}[\log P(\mathbf{X}, \mathbf{Z}, \mathbf{G}, \Theta)] - \mathbb{E}[\log q(\mathbf{Z}, \mathbf{G}, \Theta)]. \quad (2)$$

# Methods (cont.)

- ELBO is non-convex optimization - susceptible to local optimum
- ran VB algorithm  $T$  times (1000 for real data sets, 300 for synthetic data)
- started from different initial cluster assignments for each cell
- $K = 10$  for all runs (considers 1, 2,..., 10 clusters)
- 10 initializations each from EuclidClust, HammingClust, PearsonClust - 1 from DensityCut
- The rest ( $T - 31$ ) of the initializations are uniformly random
- Selecting the best run
  - Minimum DIC score over runs with  $c$  clusters
  - DIC elbow plot to select optimal number of clusters, best overall run



**Figure A** Example of a DIC elbow plot for EpiclomalRegion for the InHouse data set with 10 000 loci. Our DIC algorithm selects only number of clusters for which there is a decrease of at least 2% in DIC values, green vertical line. Then, the elbow value is picked, red vertical line.

# Results

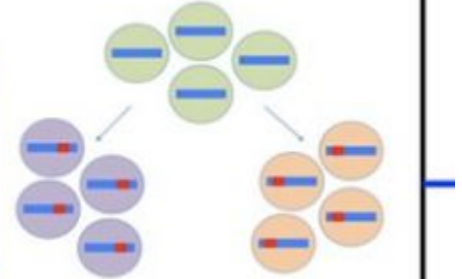
- V-measure is calculated as the harmonic mean between homogeneity (h) and completeness (c)
  - A predicted clustering has homogeneity 1 if all of the predicted clusters contain only data points which are members of a single true class.
  - A predicted clustering result has completeness 1 if it assigns all of those data points that are members of a single true class to a single predicted cluster.

# Synthetic data

## Parameters

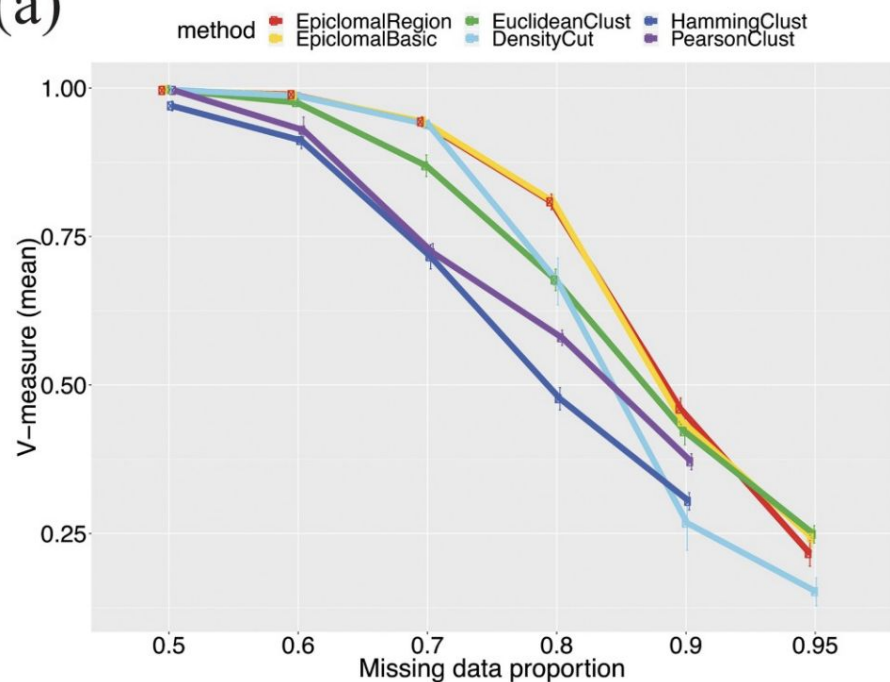
- Missing probability
- Number of cells
- Number of loci
- Number of clones
- Clone prevalence
- Number of regions
- Cell to cell variability

Random sampling  
based on a  
phylogenetic tree

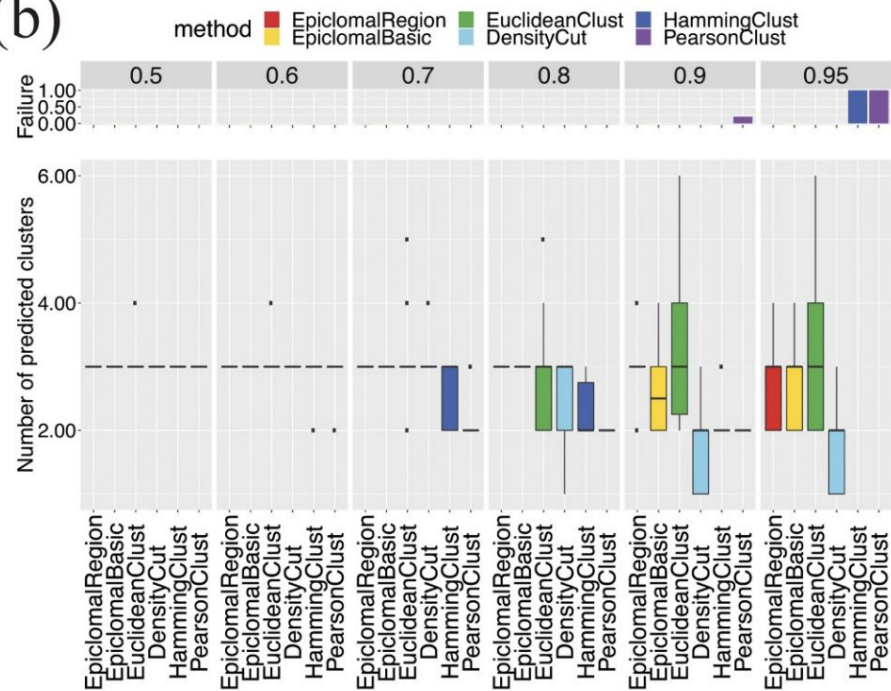


Varying parameter	Varying range
Missing proportion	0.5 to 0.95
Number of regions	25 to 200
Number of cells	12 to 2500
Cell-to-cell variability	0 to 0.3
Number of clusters (epiclones)	1 to 10
Epiclone frequencies	balanced to very unbalanced
Number of loci	5 000 to 500 000
Number of regions different between clusters	1 to 6

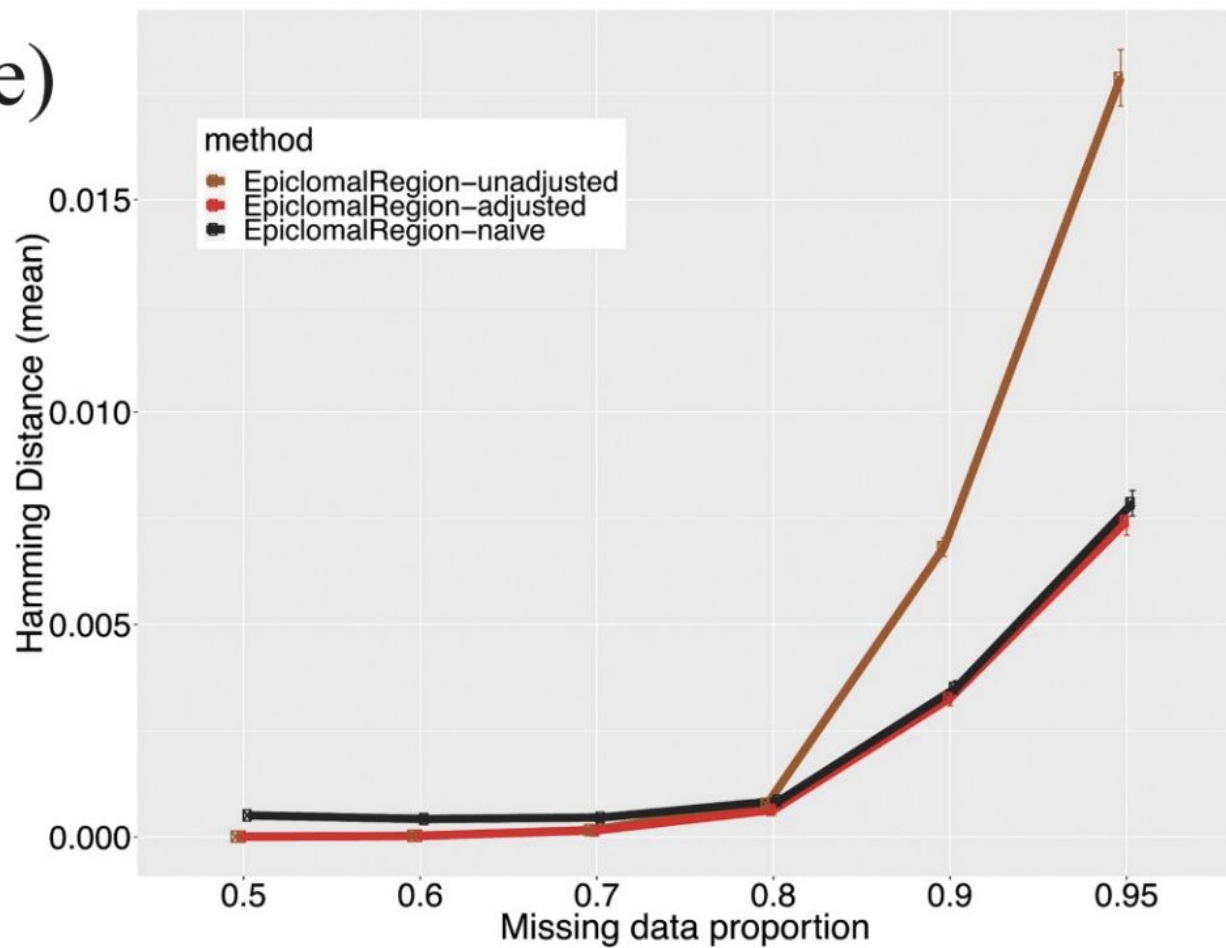
(a)

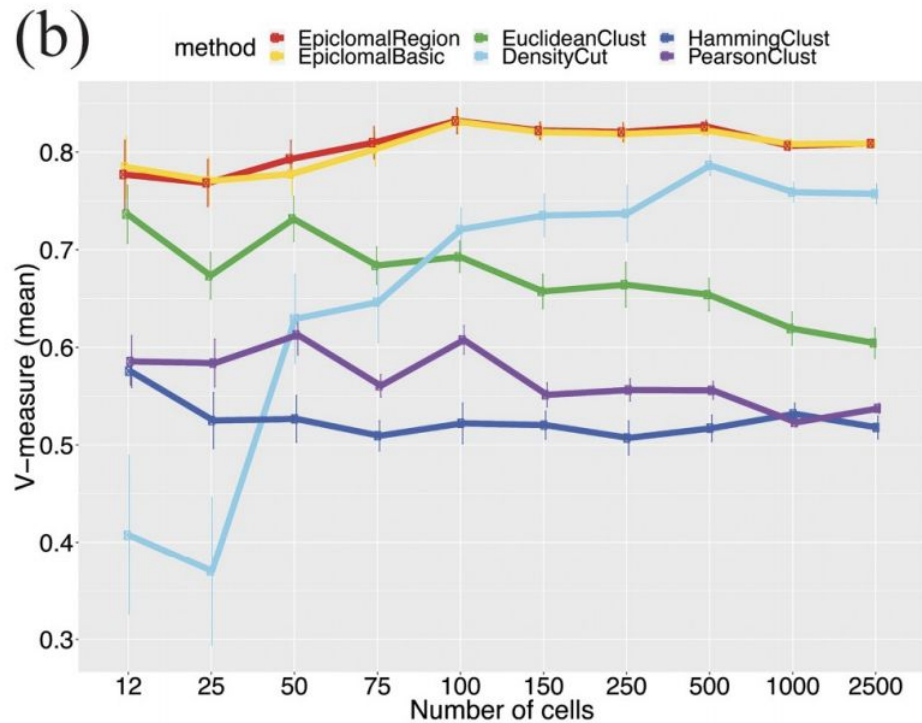
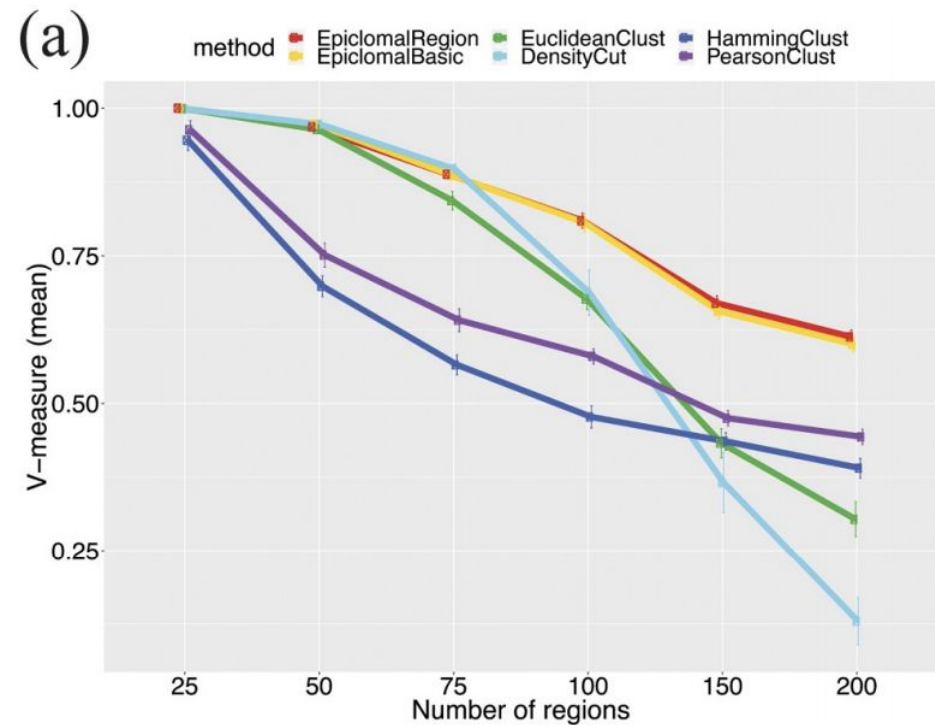


(b)

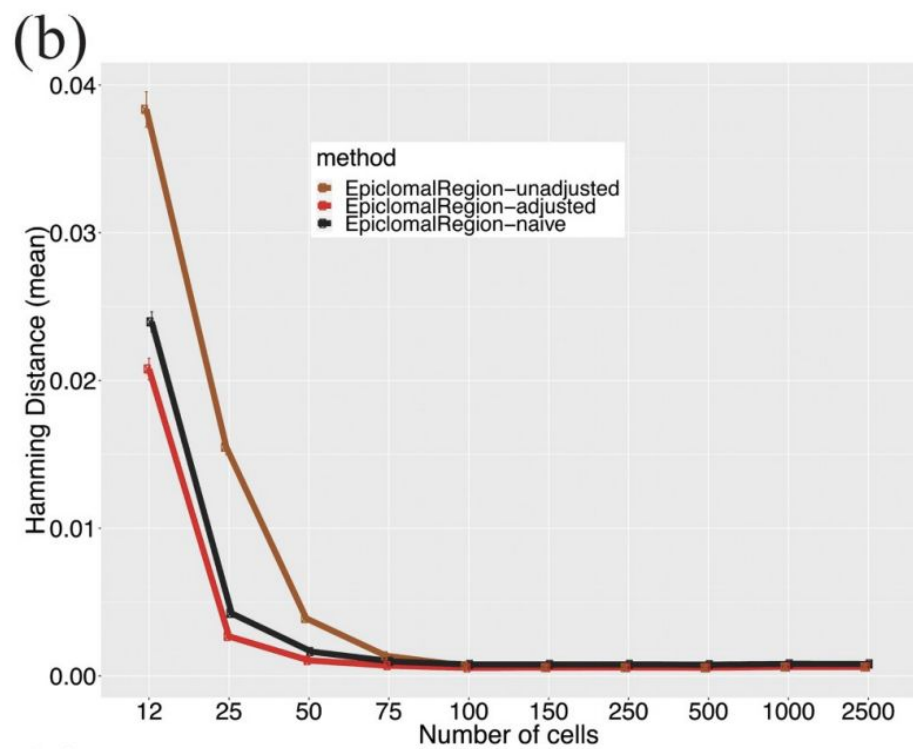
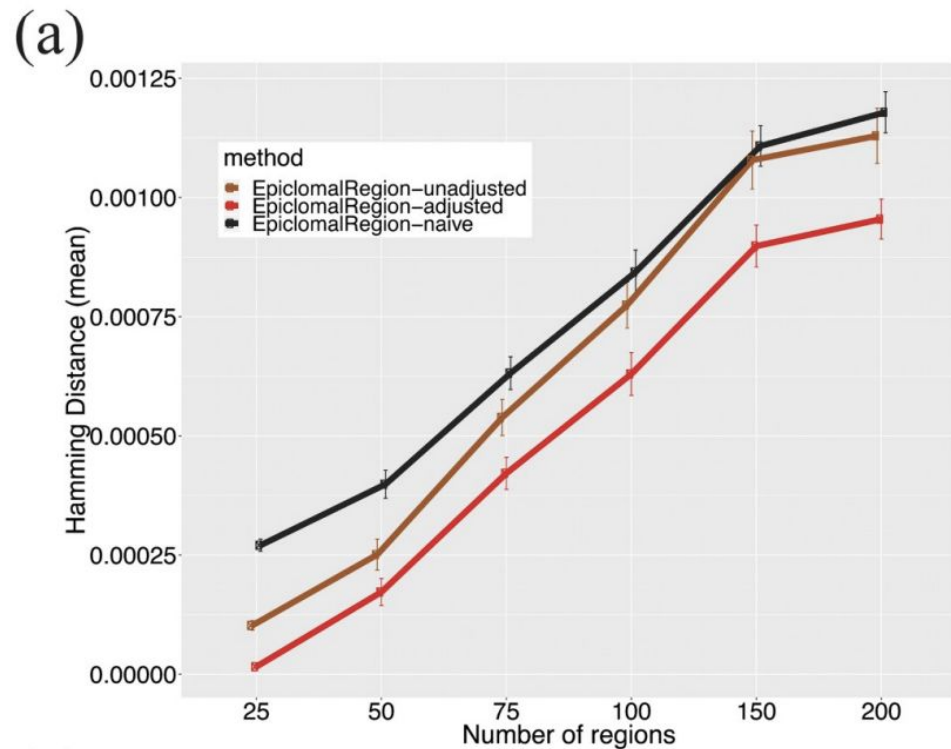


(e)

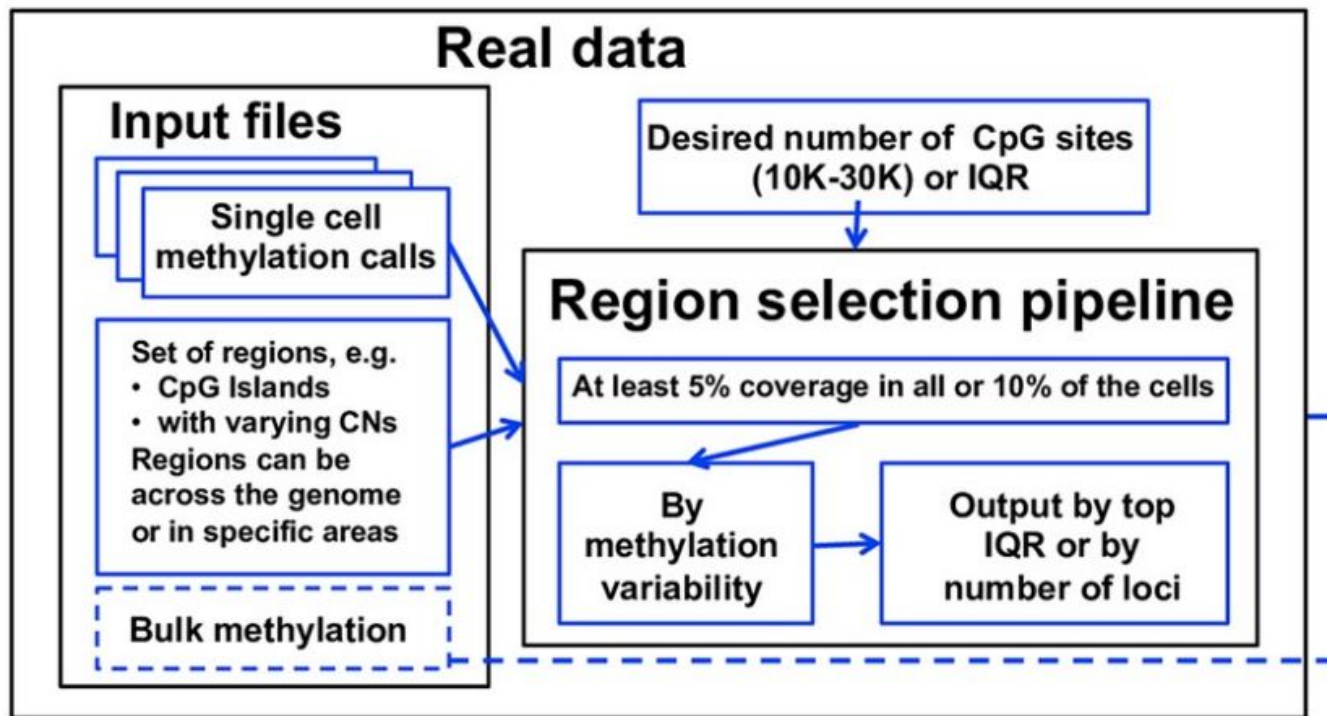




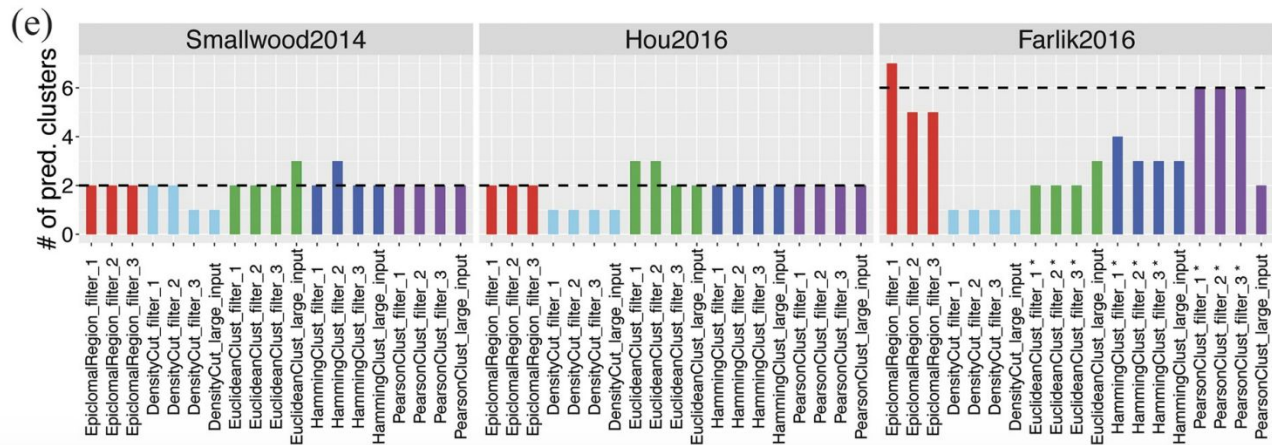
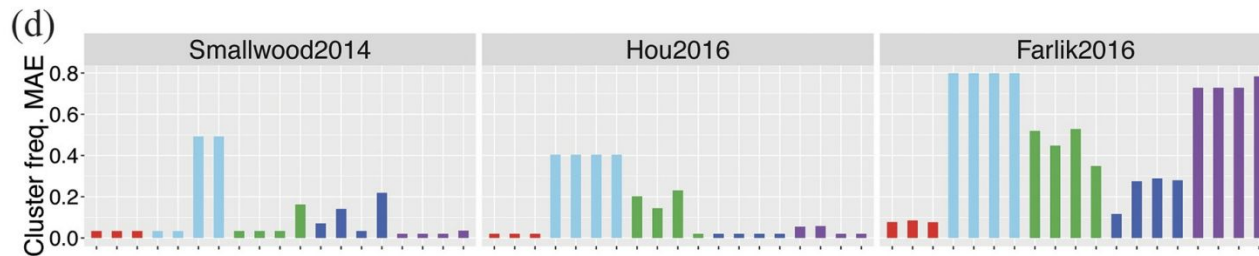
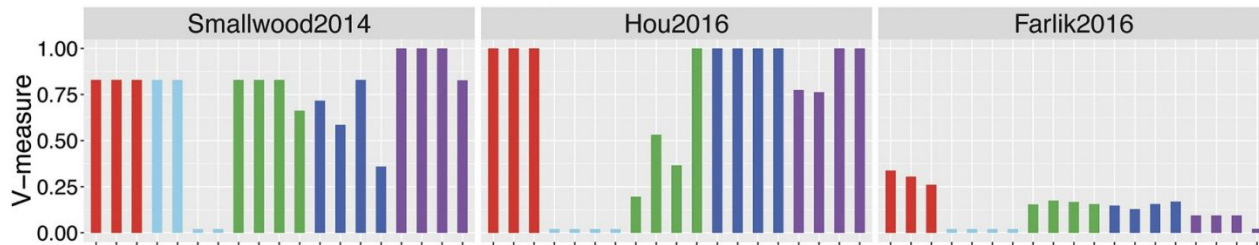




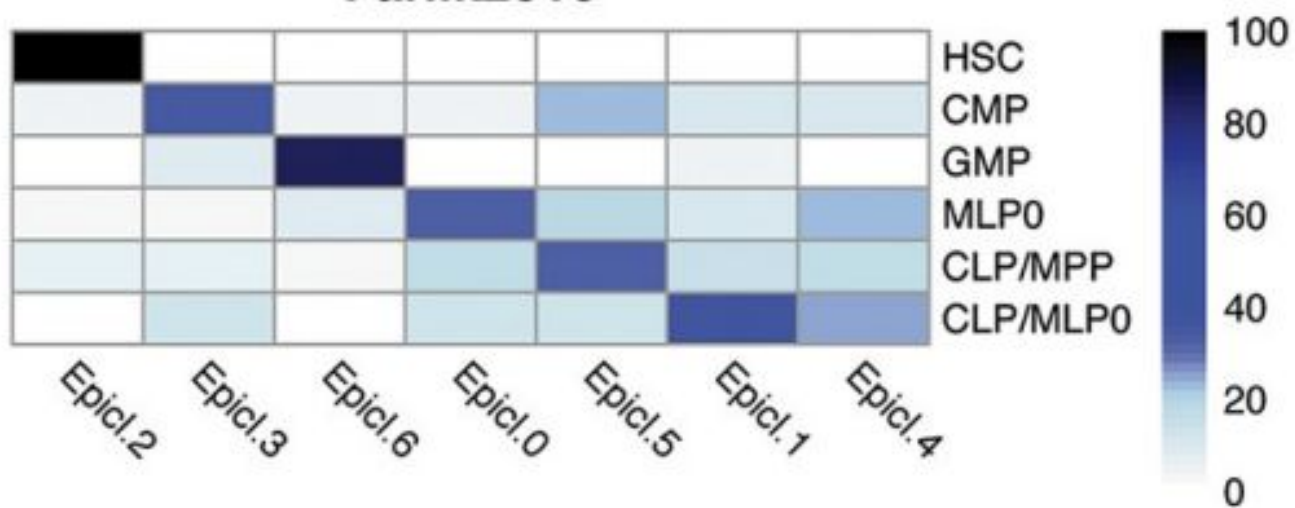
Data set	Cell type	# cells	# clusters	Regions	Miss 10K	Nloci IQR $\geq .01$	Miss IQR $\geq .01$
Smallwood2014 [16]	mouse embryonic stem cells	32	2	CGI	0.69	786 620	0.54
Hou2016 [11]	human hepatocellular carcinomas	25	2	CGI	0.87	255 136	0.90
Farlik2016 [12]	human hematopoietic cells	122	6	TFBS	0.89	512 153	0.98
InHouse	human xenografted cancer cells (3 patients)	558	NA	CGI	0.82	1 019 956	0.79



(c) method ■ EpiclomalRegion ■ EuclideanClust ■ DensityCut ■ HammingClust ■ PearsonClust



## Farlik2016

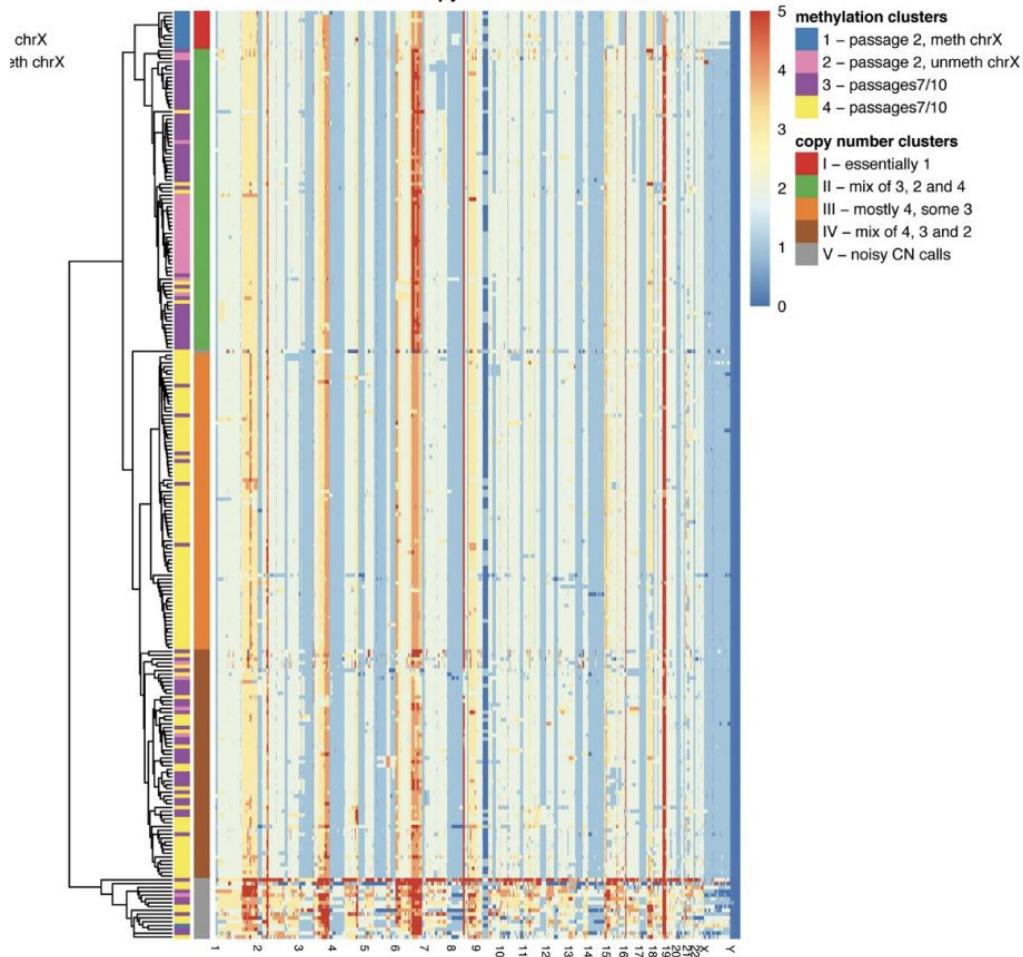


# InHouse Data

- Single-cell epigenomes generated in-house on a range of patient-derived breast tumour xenografts
- 2 patients with triple-negative breast cancer, 1 with ER+PR-Her2+ breast cancer
- For one patient: compare epiclone clusters with copy-number clones
- Some chromosomal regions may show strong copy-number influence on CpG states (X chromosome), while others may not
- Epiclones and CN clones can match or transcend each other

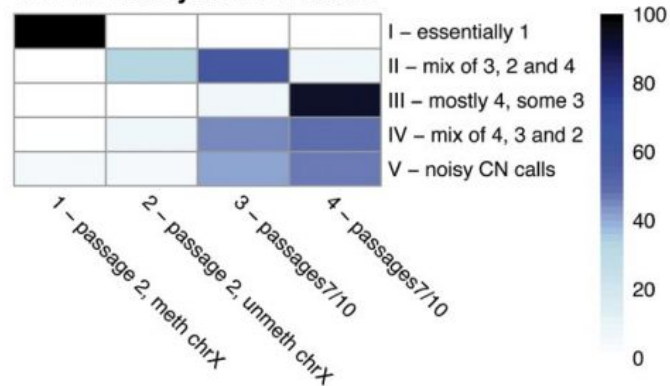
(b)

Genome-wide copy number data for SA501



(d)

CN vs. methylation clusters



(c)

