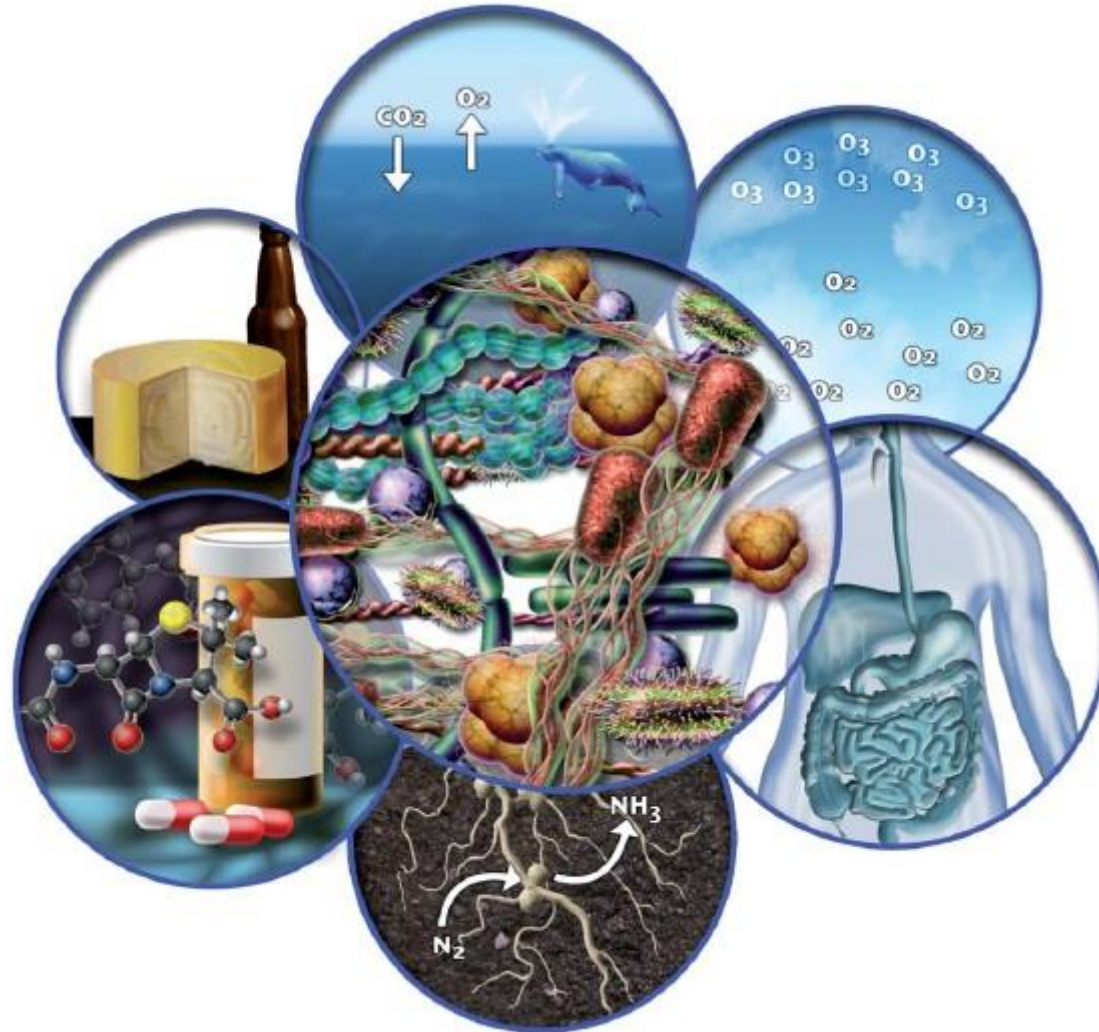


Binning metagenomic contigs using sequence
COmposition, read CoverAge, CO-alignment, and
paired-end read LinkAge

Yang Lu @ Prof. Fengzhu Sun's Lab
University of Southern California

JSM 2016

Microbes are Everywhere

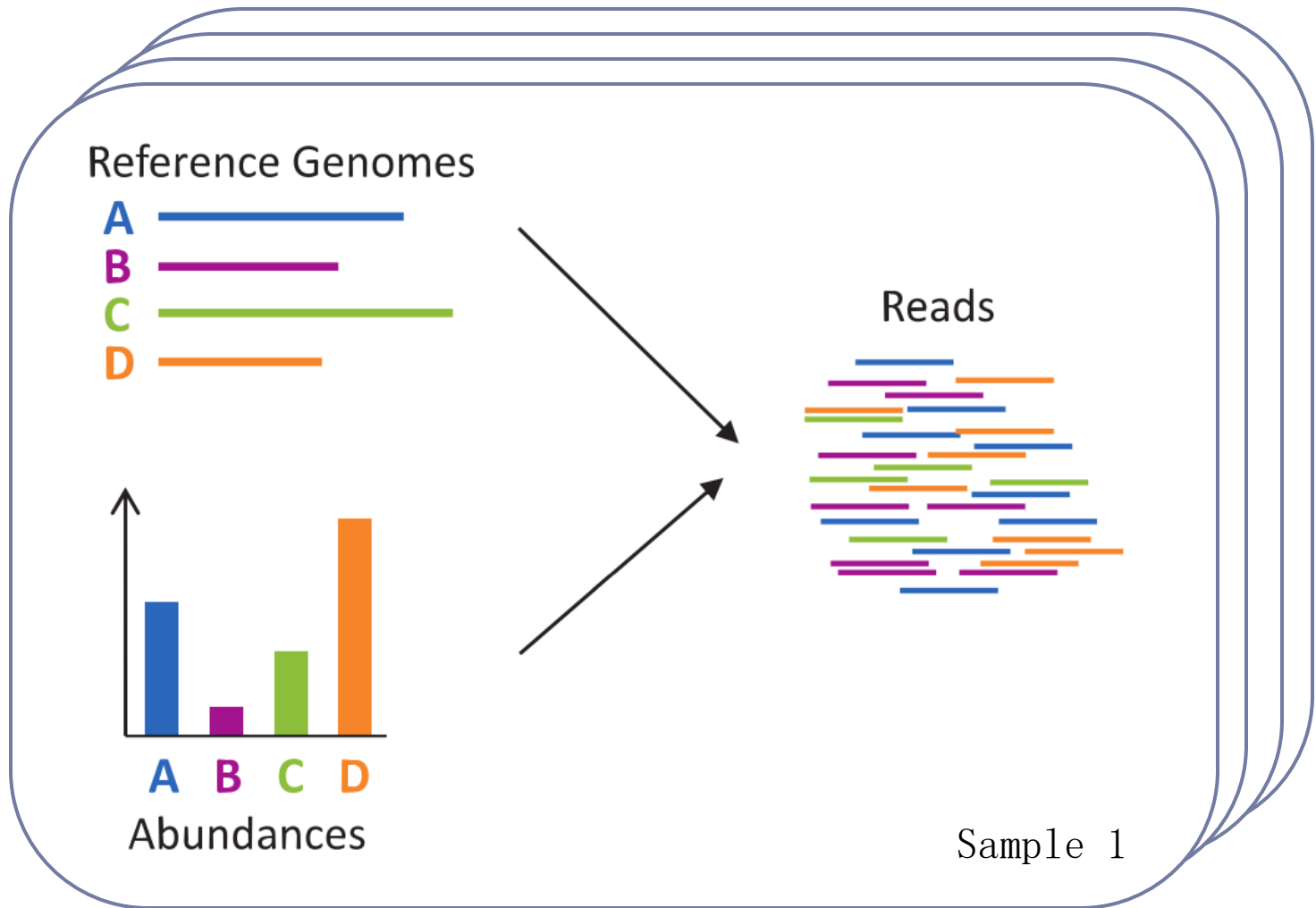


What is Metagenomics?

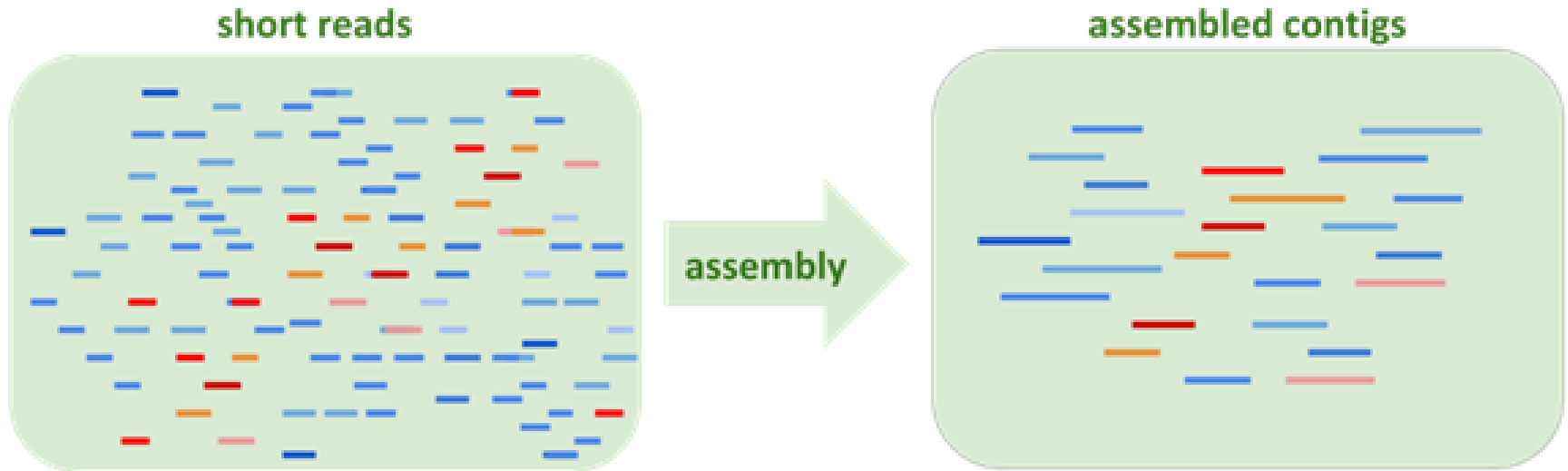
- ▶ **Metagenomics** is the study of genetic material recovered directly from environmental samples.
- ▶ Many Organisms in one sample
- ▶ Many samples from the same environment



Generative Model of Whole Genome Sequencing

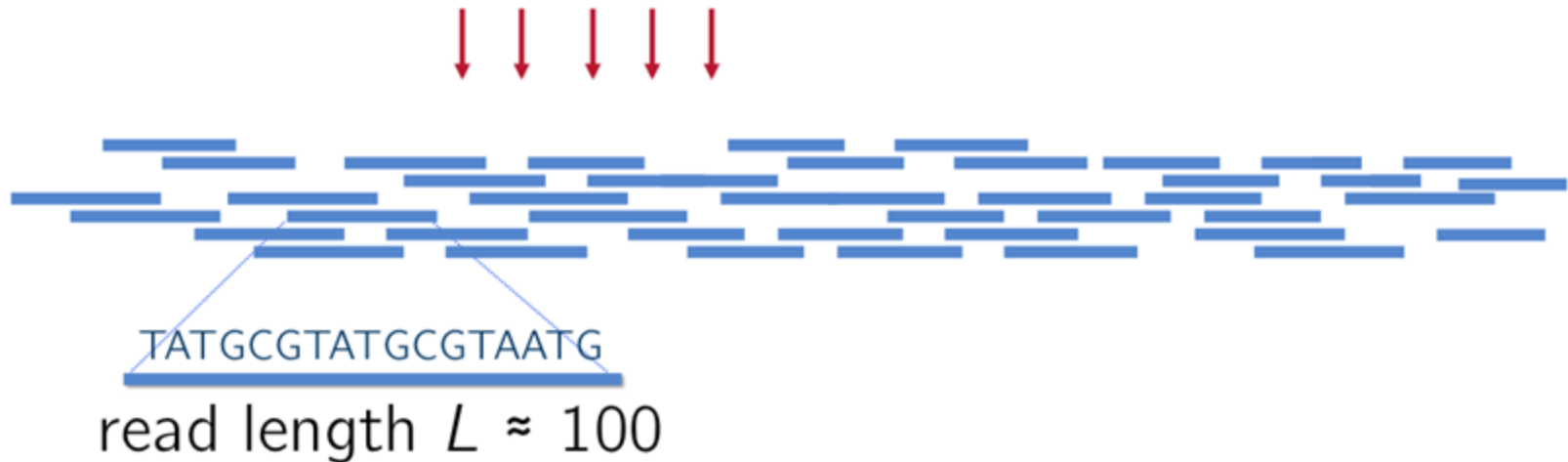


Recover by Assembly



Assembly as Huge Jigsaw Puzzle

ACGTCCTATGCGTATGCGTAATGCCACATATTGCTATGCGTAATGCGTACC

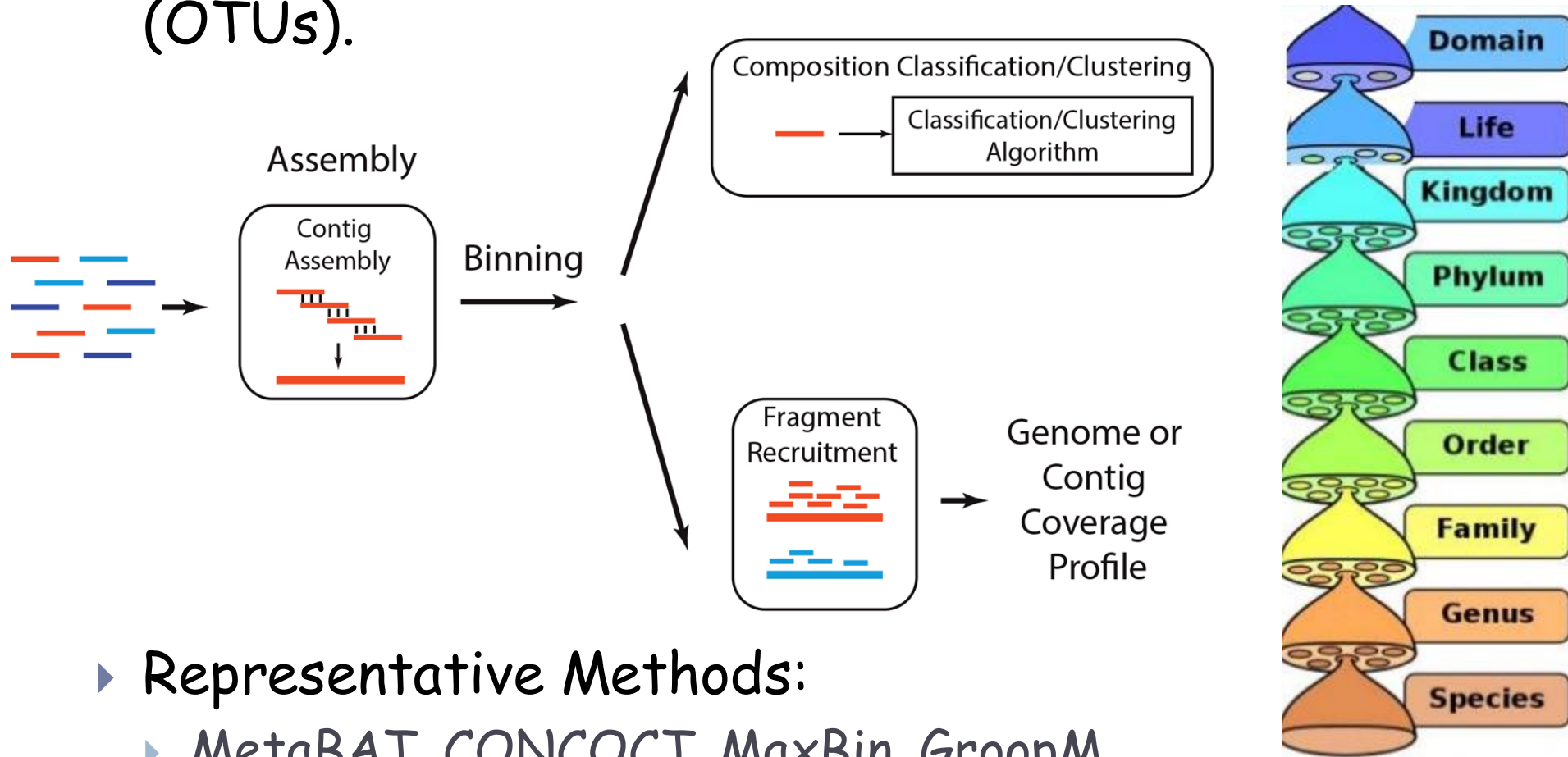


- ▶ Assembly is **Error-prone**
 - ▶ Sequencing error rate by technology limitation
 - ▶ Strain-level variation by environment complexity
 - ▶ Repetitive regions within and across genomes



Metagenomics Binning

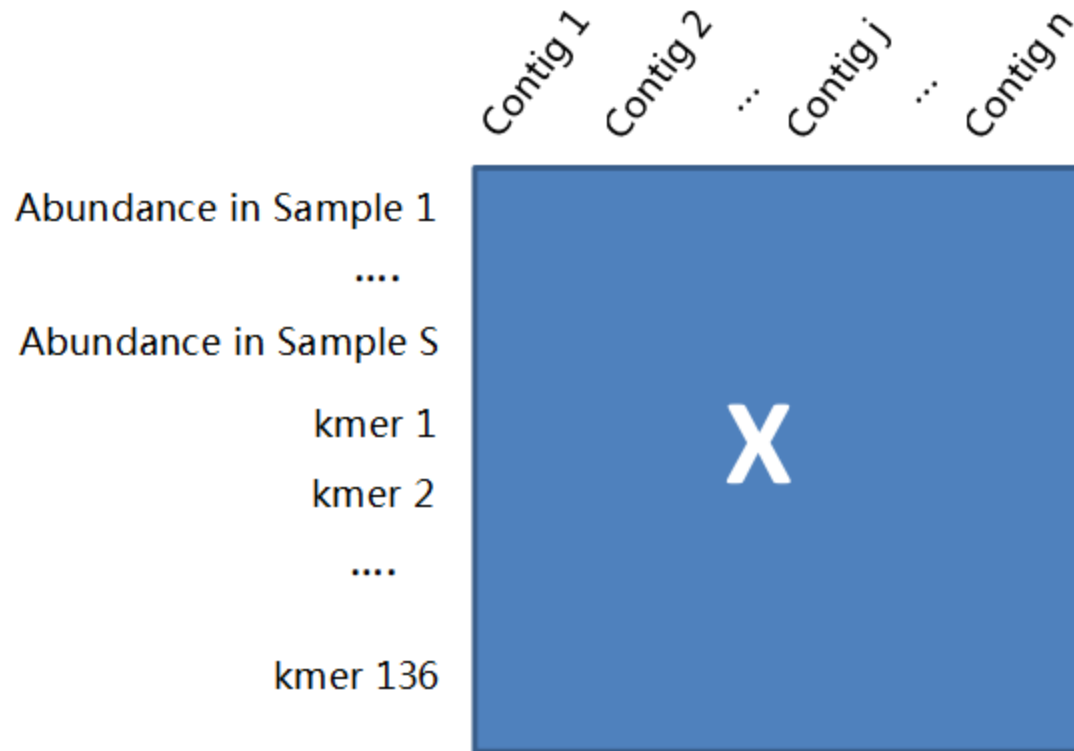
- ▶ Group contigs into Operational Taxonomic Units (OTUs).



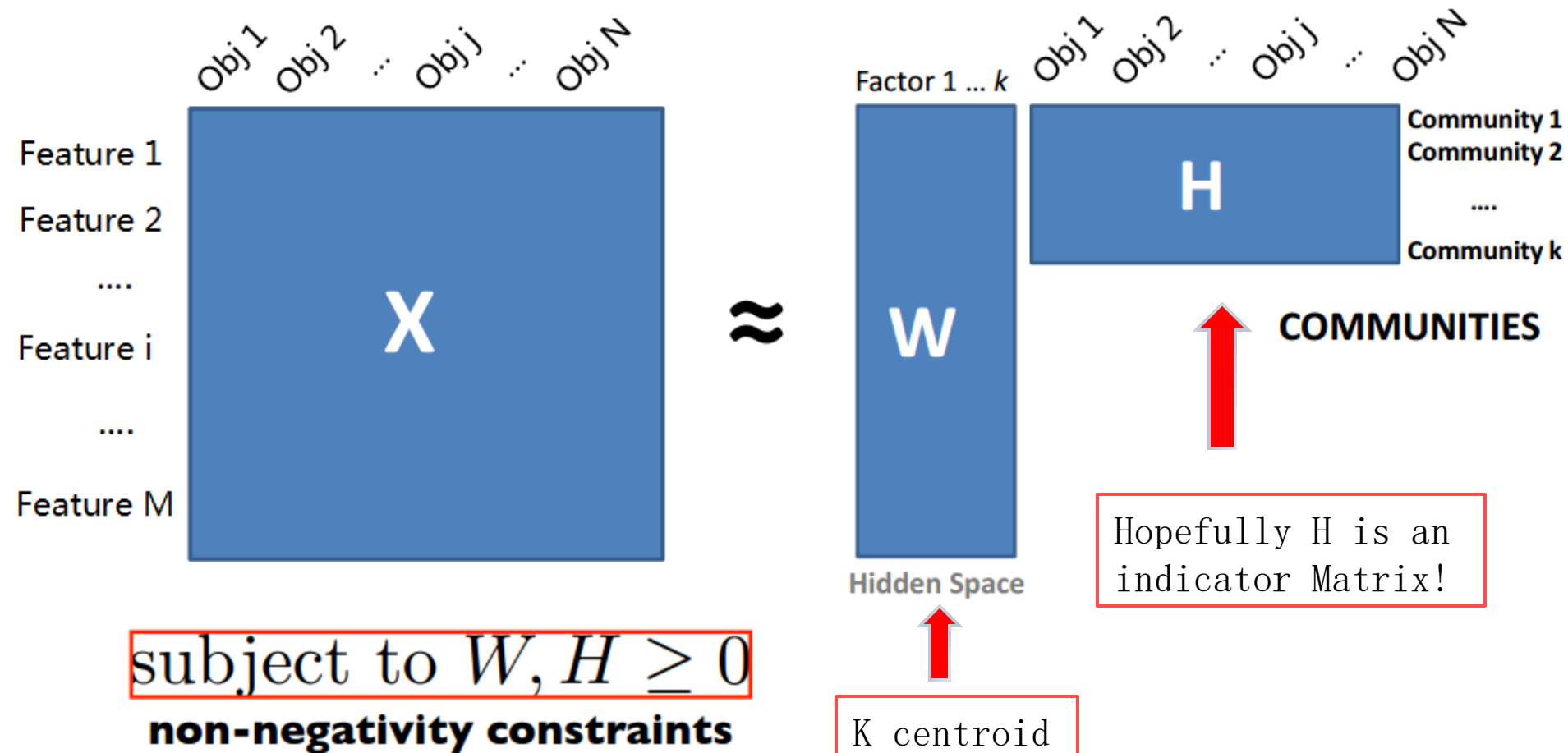
- ▶ Representative Methods:

- ▶ MetaBAT, CONCOCT, MaxBin, GroopM.

Feature-Object Matrix Representation



Illustration



Problem Formulation

$$\mathbf{x}.1 = \mathbf{h}_{11}\mathbf{w}.1 + \mathbf{h}_{21}\mathbf{w}.2 + \mathbf{h}_{31}\mathbf{w}.3 + \cdots + \mathbf{h}_{k1}\mathbf{w}.k$$

$$\mathbf{x}.2 = \mathbf{h}_{12}\mathbf{w}.1 + \mathbf{h}_{22}\mathbf{w}.2 + \mathbf{h}_{32}\mathbf{w}.3 + \cdots + \mathbf{h}_{k2}\mathbf{w}.k$$

...

$$\mathbf{x}.N = \mathbf{h}_{1N}\mathbf{w}.1 + \mathbf{h}_{2N}\mathbf{w}.2 + \mathbf{h}_{3N}\mathbf{w}.3 + \cdots + \mathbf{h}_{kN}\mathbf{w}.k$$



$$X \approx WH$$

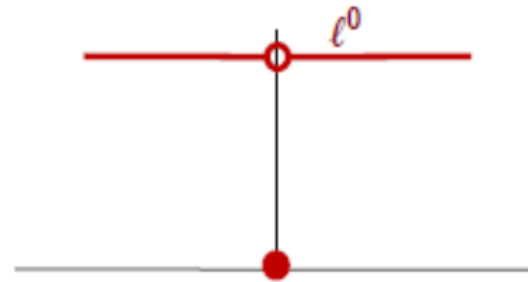
$$s.t. \quad W \geq 0, \quad \mathbb{H} \in \{0, 1\}^{K \times N}, \quad \|\mathbb{H}.j\|_0 = 1 \text{ for } j = 1, 2, \dots, N$$



Relaxation

$$\arg \min_{\substack{W \geq 0 \\ \mathbb{H} \in \{0, 1\}^{K \times N}}} \|X - W\mathbb{H}\|_F^2$$

Hard to Solve!
Need Relaxation!



$$\arg \min_{W, H \geq 0} \|X - WH\|_F^2$$

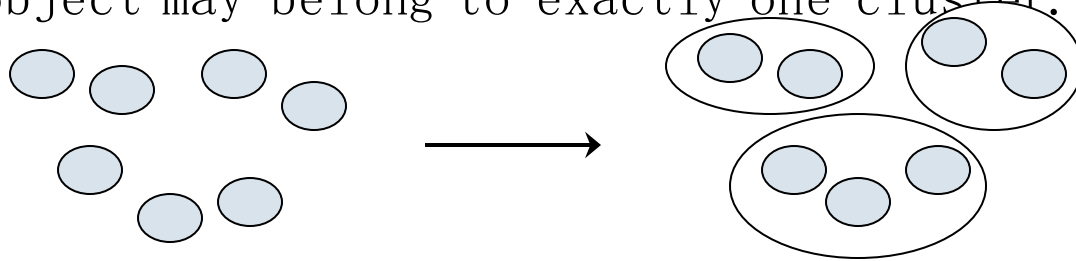


Drawback

- ▶ “hard clustering” to “soft clustering”

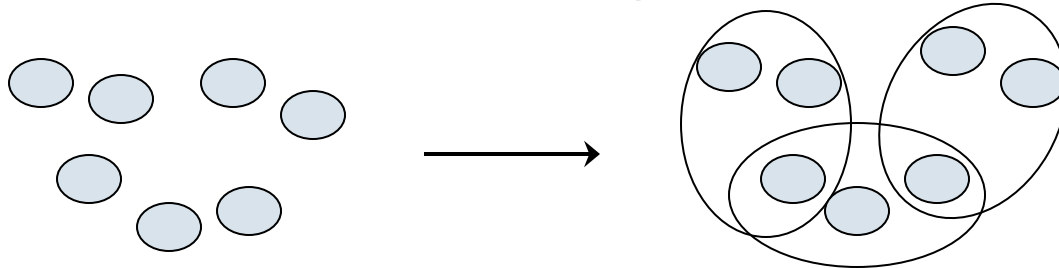
Hard Clustering

- Every object may belong to exactly one cluster.



Soft Clustering

- The membership is fuzzy – Objects may belong to several clusters with a fractional degree of membership in each.

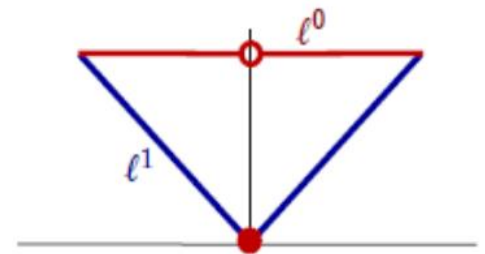


Sparsity comes to rescue

- ▶ To facilitate “hard clustering” -like behavior

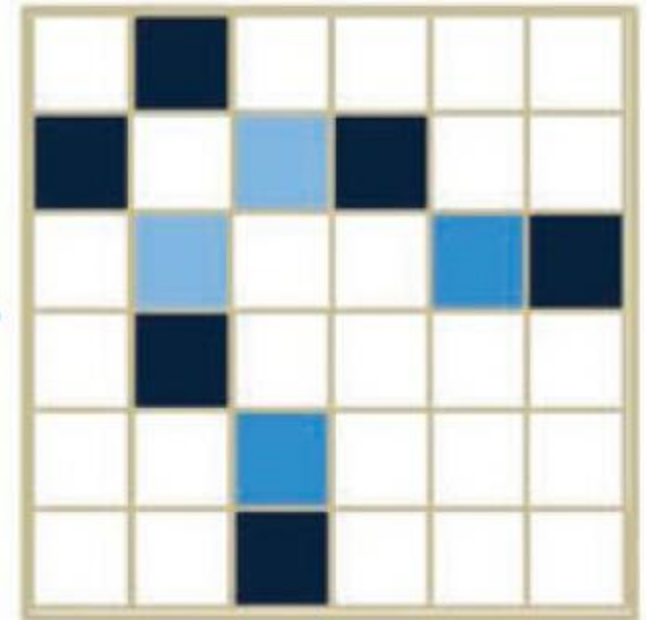
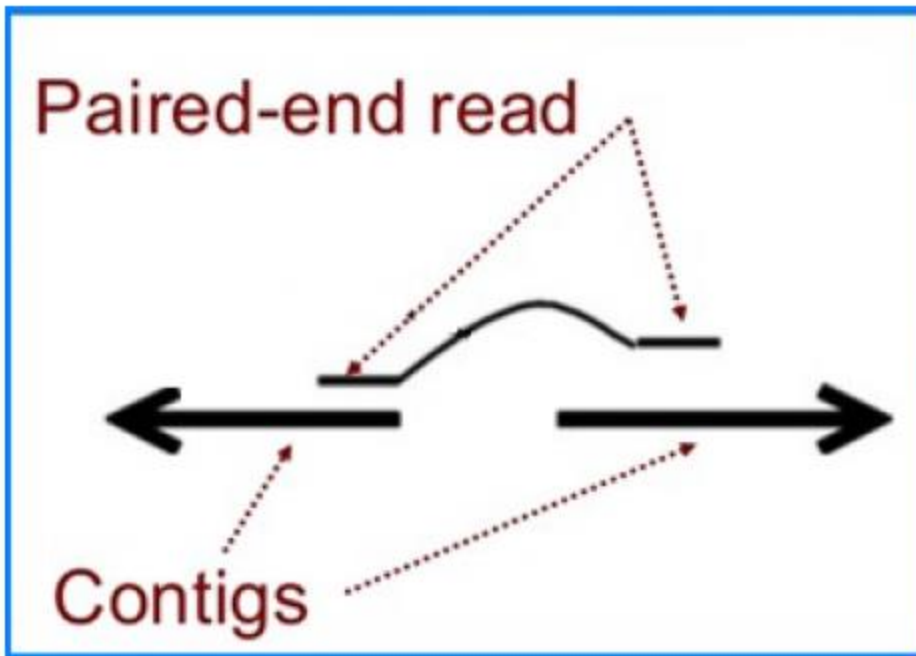
$$\arg \min_{W, H \geq 0} \|X - WH\|_F^2 + \alpha \sum_{j=1}^N \|H \cdot j\|_1^2$$

Sparse Non-negative Matrix Factorization



Incorporating Side Information

$$\arg \min_{W, H \geq 0} \|X - WH\|_F^2 + \alpha \sum_{n=1}^N \|H_{\cdot n}\|_1^2 + \beta \text{Tr}(H \mathcal{L} H^T)$$



Optimization

- ▶ By Alternating Nonnegative Least Squares

$$H \leftarrow \arg \min_{H \geq 0} \|X - WH\|_F^2 + \alpha \sum_{n=1}^N \|H_{\cdot n}\|_1^2 + \beta \text{Tr}(H\mathcal{L}H^T)$$

$$W \leftarrow \arg \min_{W \geq 0} \|X^T - H^T W^T\|_F^2$$



Block Coordinate Descent

$$\begin{aligned} & \arg \min_{H \geq 0} \|X - WH\|_F^2 + \alpha \sum_{n=1}^N \|H_{\cdot n}\|_1^2 + \beta \text{Tr}(H\mathcal{L}H^T) \\ & \approx \arg \min_{H \geq 0} \sum_{n=1}^N \left(\|X_{\cdot n} - WH_{\cdot n}\|_2^2 + \alpha \|H_{\cdot n}\|_1^2 + \beta H_{\cdot n}^T (H_{\cdot n} - 2 \sum_{n'=1}^N \mathcal{A}_{nn'} H_{\cdot n'}^{old}) \right) \\ & = \arg \min_{H \geq 0} \sum_{n=1}^N \left(\|X_{\cdot n} - WH_{\cdot n}\|_2^2 + \alpha \|H_{\cdot n}\|_1^2 + \beta \left\| H_{\cdot n} - \sum_{n'=1}^N \mathcal{A}_{nn'} H_{\cdot n'}^{old} \right\|_2^2 \right) \\ & = \arg \min_H \left\| \begin{pmatrix} X \\ 0_{1 \times N} \\ \sqrt{\beta} H^{old} \mathcal{A} \end{pmatrix} - \begin{pmatrix} W \\ \sqrt{\alpha} e_{1 \times K} \\ \sqrt{\beta} I_K \end{pmatrix} H \right\|_F^2 \end{aligned}$$



Experiments

▶ Synthetic Datasets

▶ Species Mock Community

- ▶ 101 Species, 37,628 contigs, 96 Samples,

▶ Strain Mock Community

- ▶ Mixture of *E. coli* strains, five *Bacteroides* species, five *Clostridium* genera, five other typical gut bacteria
- ▶ 9,417 contigs, 64 Samples

▶ Real Datasets

▶ Sharon

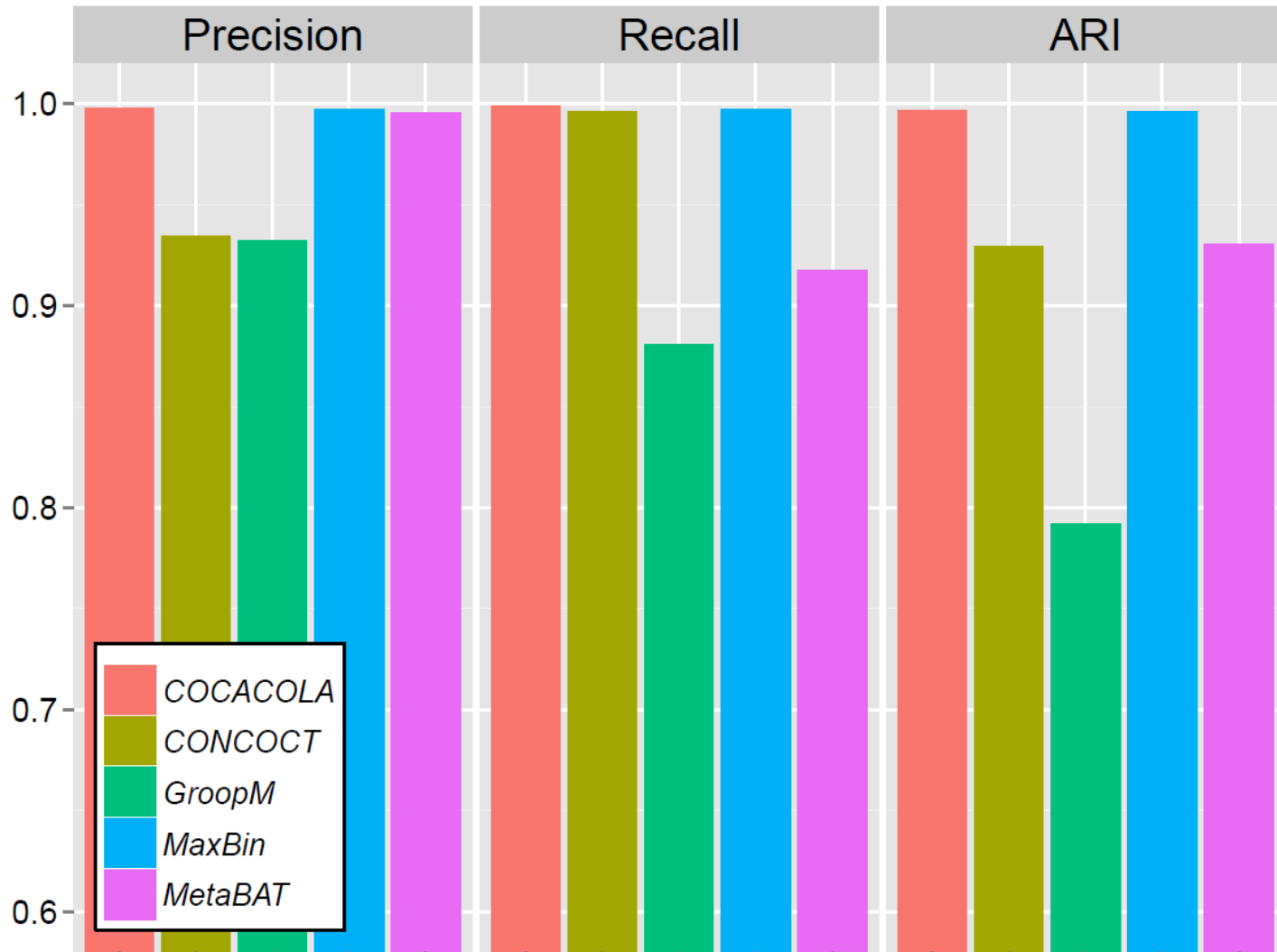
- ▶ 11 time-series samples from premature infant gut
- ▶ 2,614 out of 5,579 contigs are labelled by TAXAassign

▶ MetaHIT

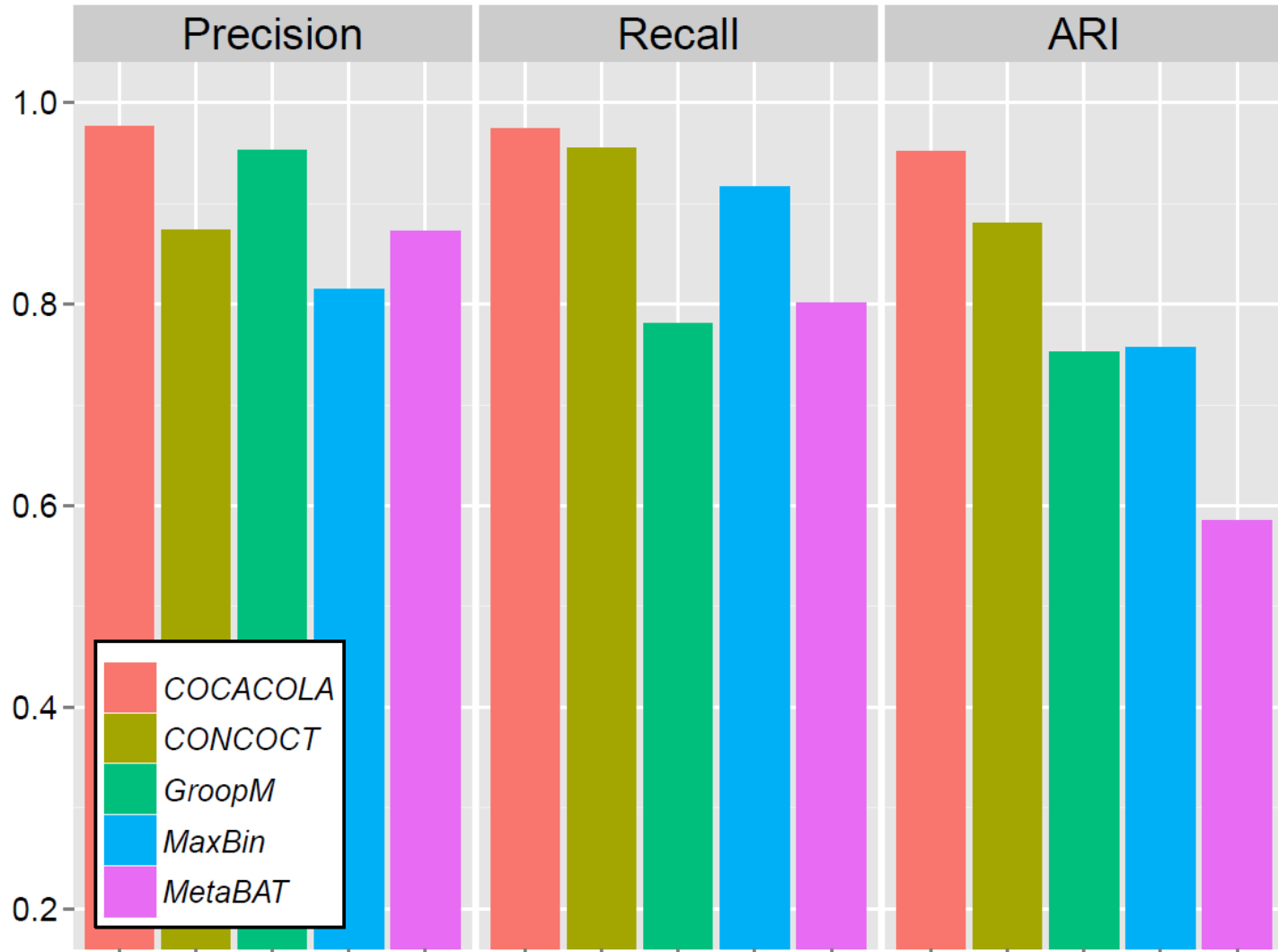
- ▶ 264 samples from MetaHIT consortium
- ▶ 17,136 out of 192,673 contigs are labelled by TAXAassign



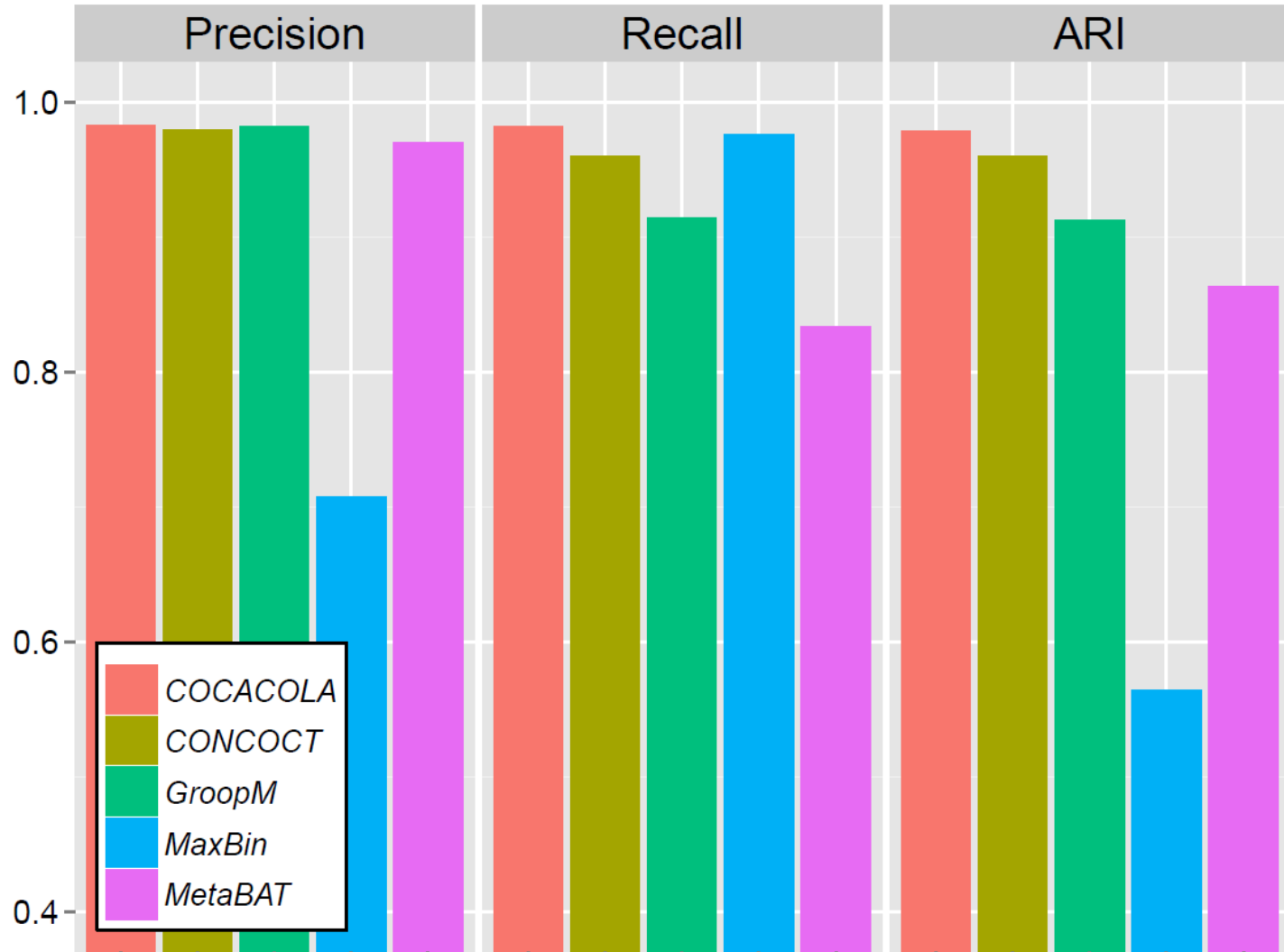
Synthetic "Species" Dataset



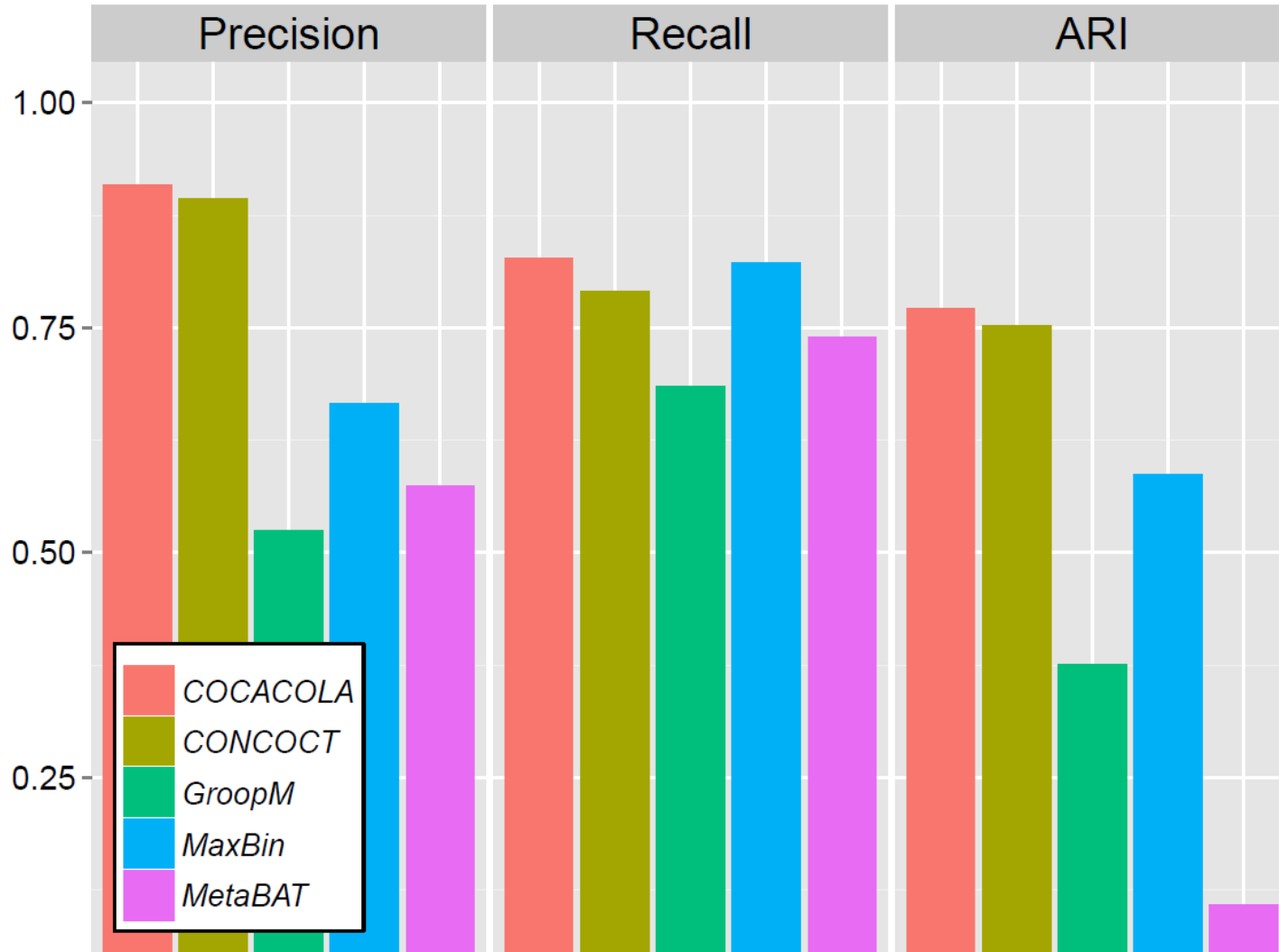
Synthetic "Strain" Dataset



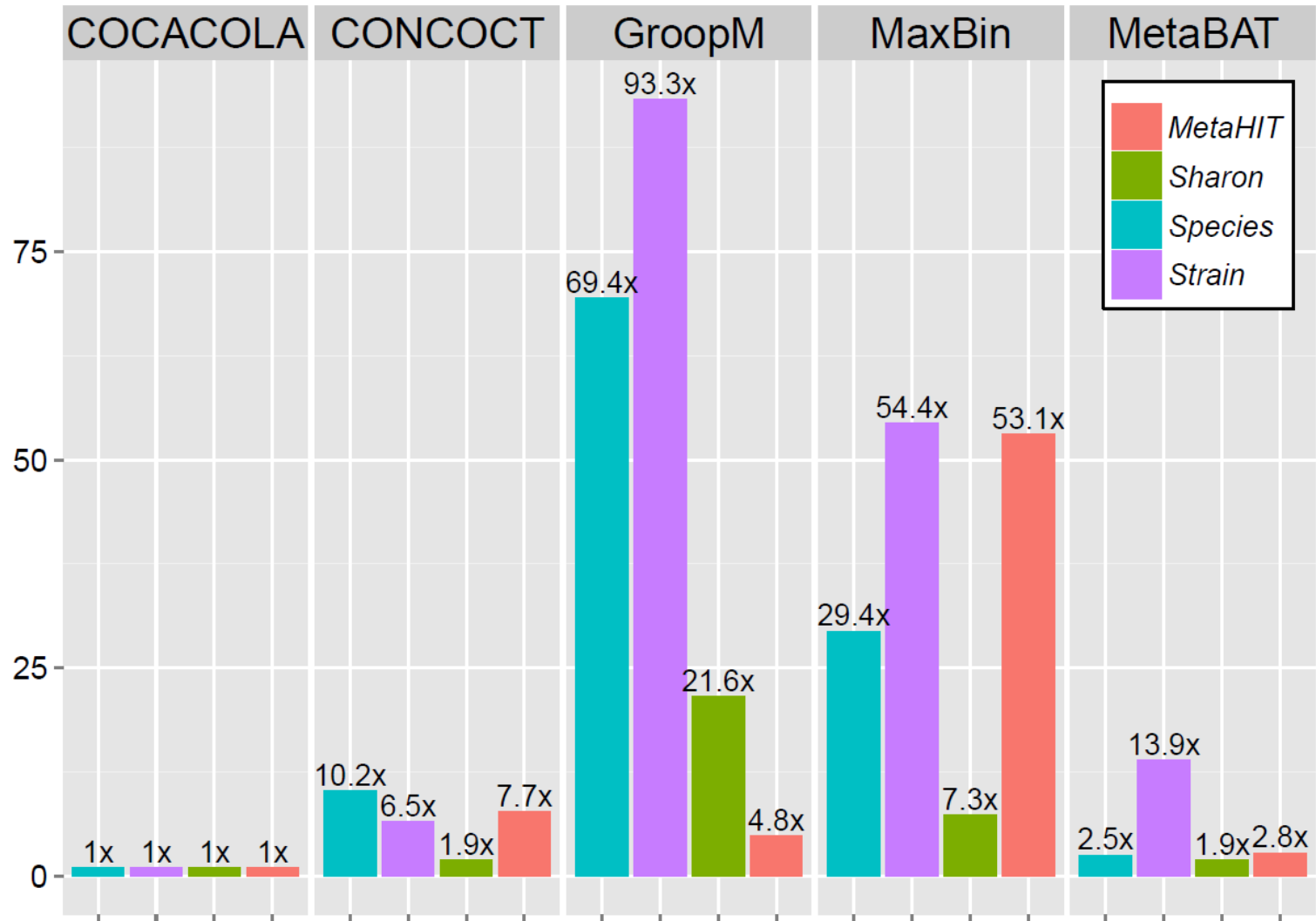
Real "Sharon" Dataset



Real "MetaHIT" Dataset



Speedup Ratio

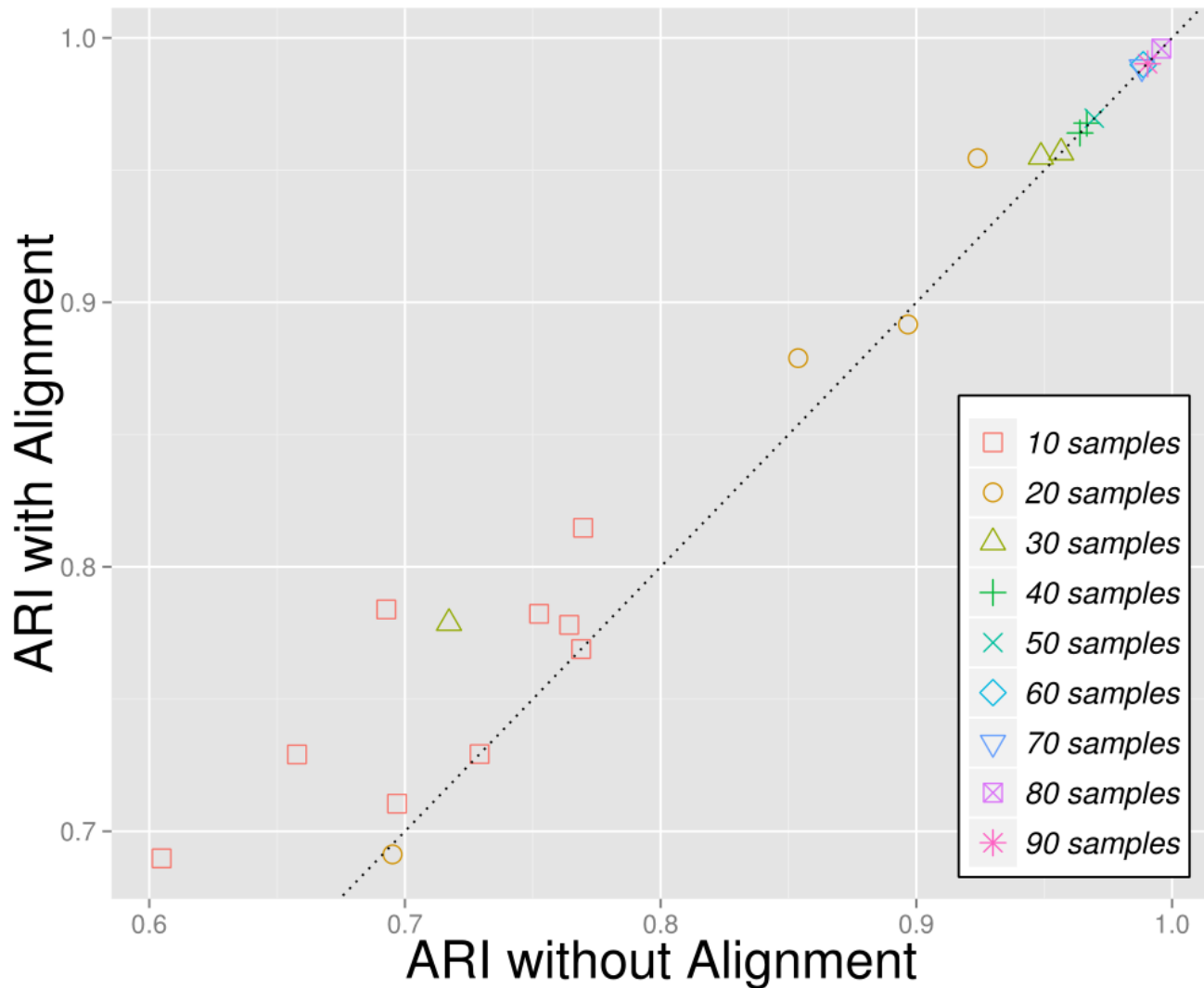


Is Side Information Useful?

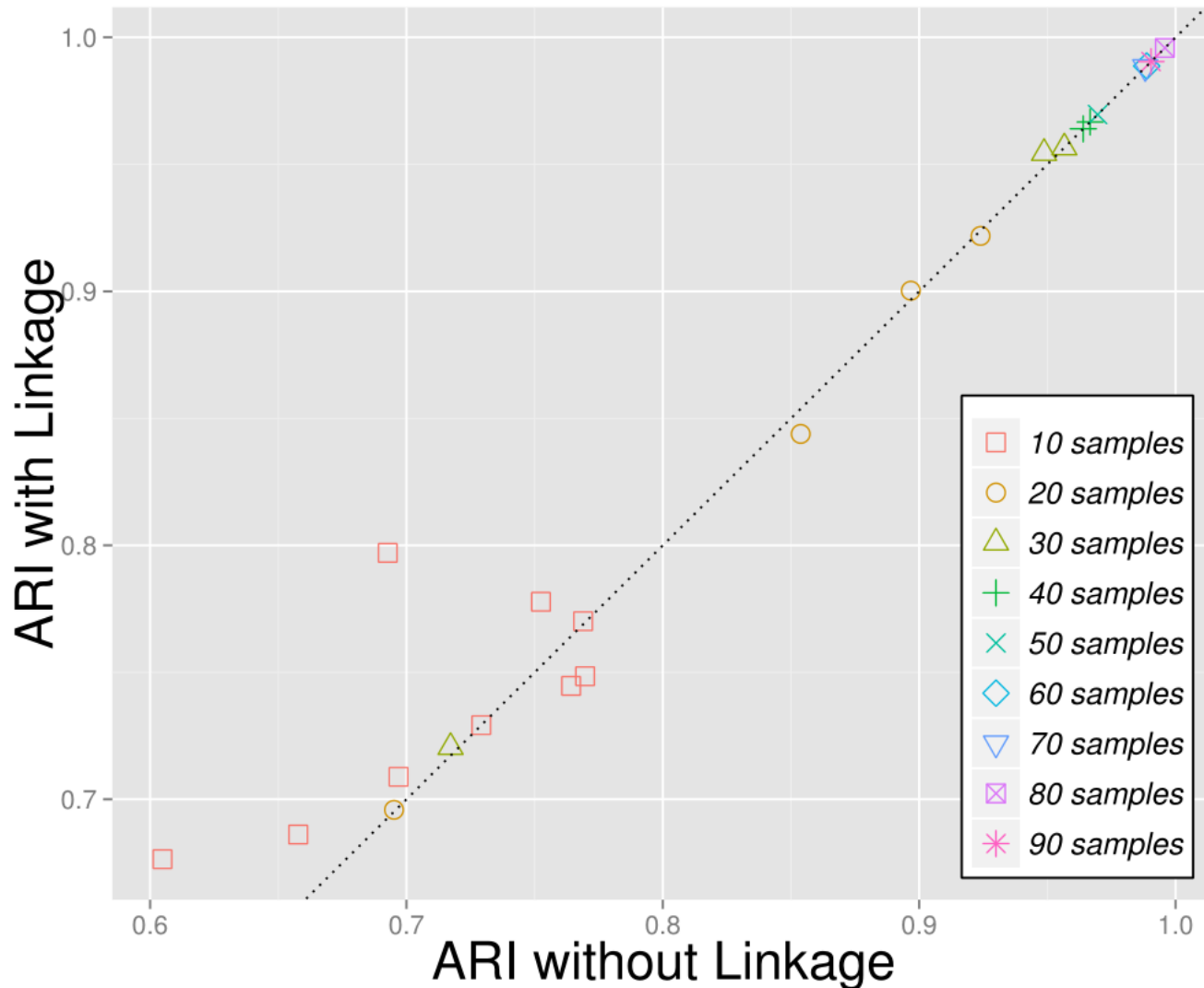
- ▶ Co-alignment to the same reference genome
- ▶ Paired-end reads linkage
 - ▶ Minimum samples support = 2
- ▶ Ensemble of both with equal weight



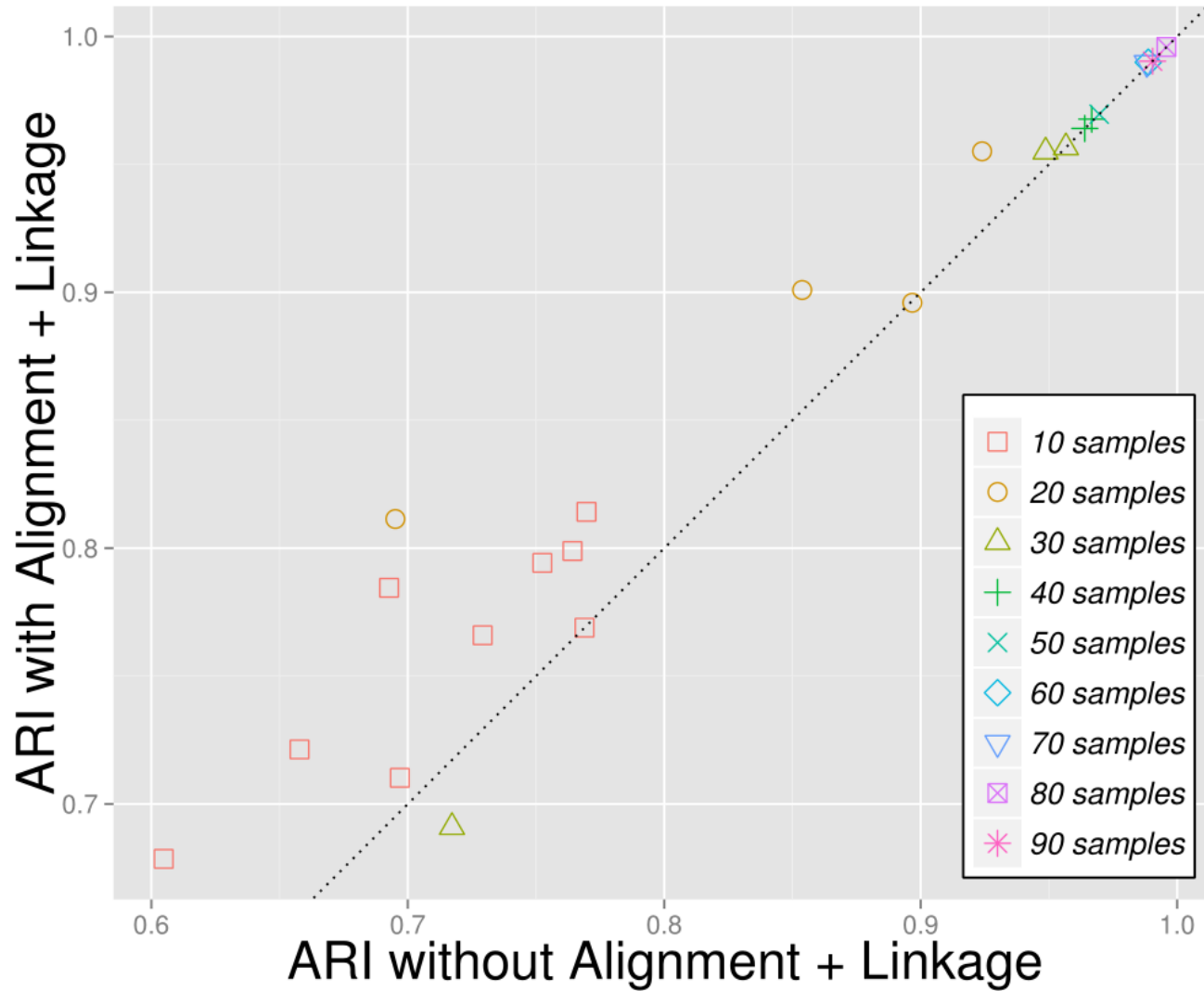
Co-Alignment as Side Information



Linkage as Side Information



Ensemble of Both



Summary so far

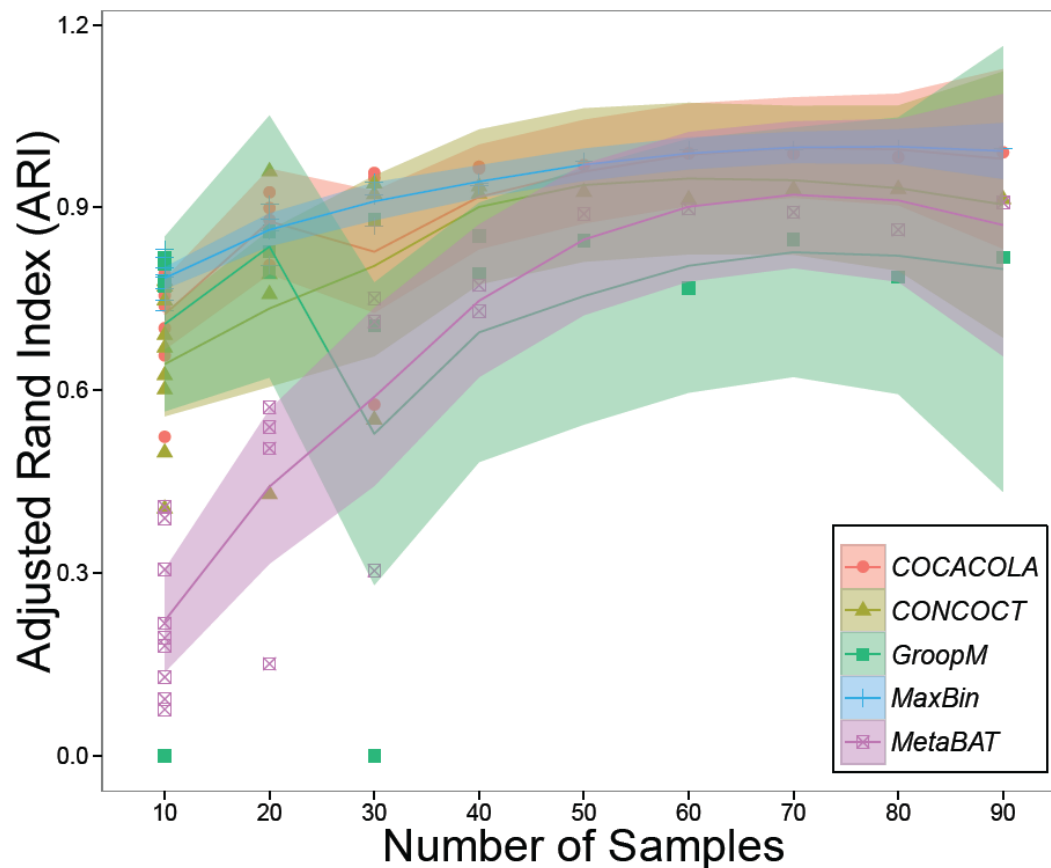
- ▶ A metagenomics contigs binning framework
 - ▶ Utilize abundance profile and sequence composition
 - ▶ Embrace additional information such as co-alignment, linkage, even customized information
 - ▶ Highly parallel and scalable

- ▶ What's next?



Limitation of Current Binning Approaches

- ▶ Observation:
 - ▶ When samples size is small, binning is unstable



Re-weight the Input Features

- ▶ **Assumption:**
 - ▶ Majority of features are neutral, i.e. with weight 1
 - ▶ Only small amount of features are either very good (weight >1) or very bad (weight <1)

- ▶ **Different from Feature Screening:**
 - ▶ Majority of features are useless (weight=0)
 - ▶ Only small amount of features are important (weight=1)

- ▶ **For each feature**
 - ▶ Tested by Multimodality dip test

Re-weighting Needs Side Information

- ▶ Let A_1 be the KNN matrix of data using heat kernel, symmetrized
- ▶ Let A_2 be the side information matrix
- ▶ Let $A = A_1 + \gamma A_2$ where $\gamma = \text{tr}(A_1' A_2) / (A_1' A_1)$ so that

$$\arg \min_{\gamma} \|A_1 - \gamma A_2\|_F^2$$

- ▶ Objective Function

$$\begin{aligned} L(W, A) &= \sum_{i,j=1}^N \|\text{diag}(W)X_{.i} - \text{diag}(W)X_{.j}\|^2 a_{ij} \\ &= \text{tr}(\text{diag}(W)X L X^T \text{diag}(W)) \end{aligned}$$



Objective Function

- ▶ Equivalent to a simple version of Mahalanobis distance Learning Formulation
- ▶ Doesn't work!
- ▶ Reformulate in terms of ΔW

$$\begin{aligned}L(W) &= \text{tr}(\text{diag}(W)XLX^T\text{diag}(W)) \\ &= \|Z\text{diag}(W)\|_F^2 \\ &= \|Z + Z\text{diag}(\Delta W)\|_F^2\end{aligned}$$



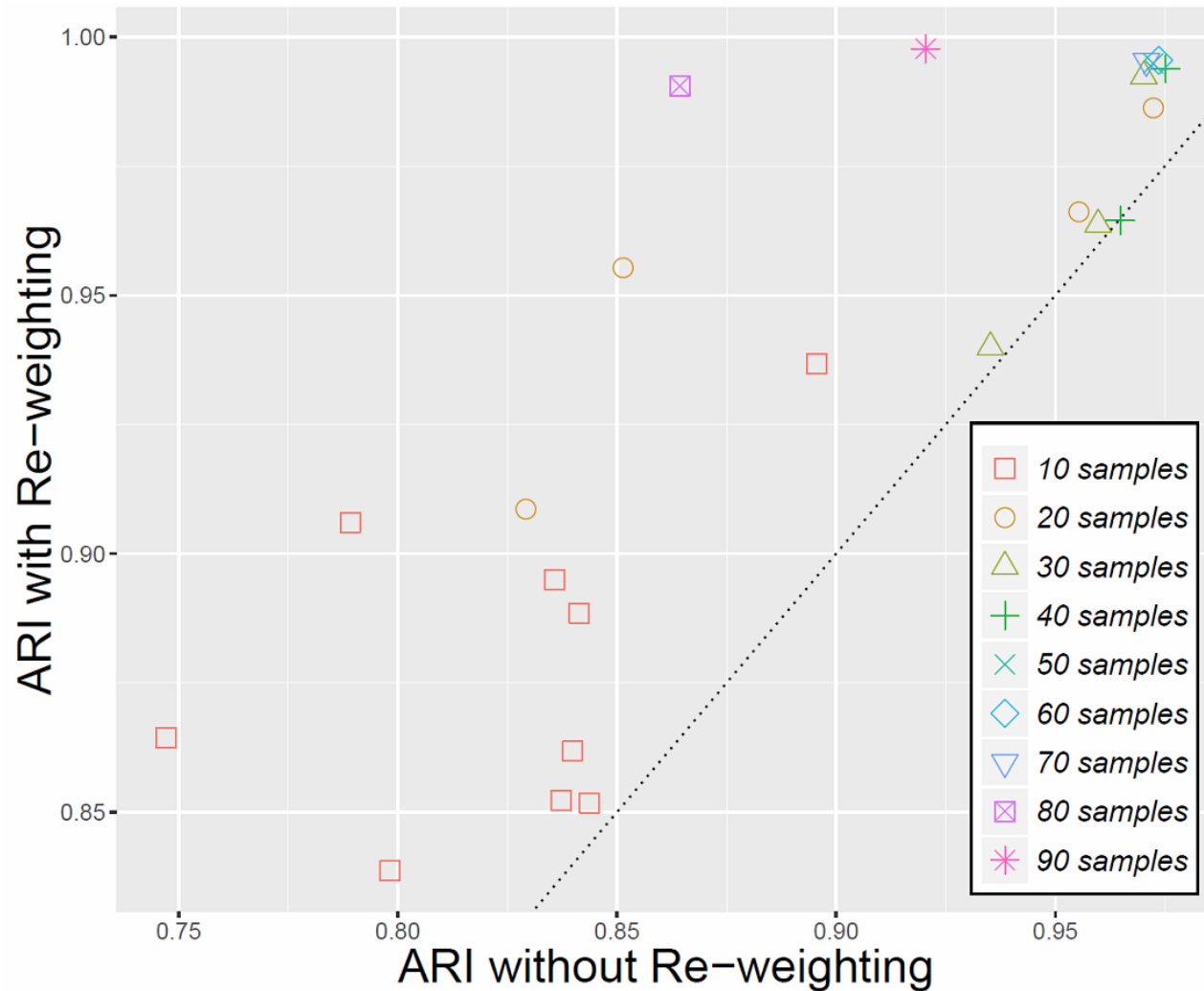
Final Objective Function

$$\arg \min_{\Delta W} \left\{ \| Z + Z \text{diag}(\Delta W) \|_F^2 + \lambda \| \Delta W \|_1 + \lambda \| \Delta W \|^2 \right\}$$

$$s.t. \quad \sum_{i=1} \Delta W_i = 0, \quad \Delta W_i \geq -1$$



Spectral Clustering after re-weighting the input



Ongoing Direction

- ▶ Encode Relative Comparison Information into regularization:
 - ▶ Contig A is closer to contig B (within the same species) than A is to C (within the same genus)
 - ▶ Incorporate the phylogenetic tree
- ▶ Feature-reweighting formulation works
 - ▶ Not only Metagenomic Binning scenario
 - ▶ Not only clustering scenario
 - ▶ Not only untransformed feature space scenario
 - ▶ More powerful combined with feature screening



Summary so far

- ▶ **A metagenomics contigs binning framework**
 - ▶ Utilize abundance profile and sequence composition
 - ▶ Embrace additional information such as co-alignment, linkage, even customized information
 - ▶ Highly parallel and scalable
- ▶ **Feature-reweighting for input data enhancement**
 - ▶ Different assumption compared to feature screening



Acknowledgement

- ▶ Research is partially supported by NSF DMS-1518001 and OCE 1136818.



Prof. Fengzhu Sun
@ Comp. Bio



Prof. Ting Chen
@ Comp. Bio



Prof. Jed Fuhrman
@Marine Bio



Prof. Jinchi Lv
@ DSO

Questions?

