

# **An Introduction to Reversible Jump MCMC for Bayesian Networks, with Application**

*Stephen Jensen, CleverSet, Inc.*

The research described in this presentation has been funded by the U.S. Environmental Protection Agency through the STAR Cooperative Agreement #CR82-9095 to Colorado State University. It has not been subjected to the Agency's review and therefore does not necessarily reflect the views of the Agency, and no official endorsement should be inferred.

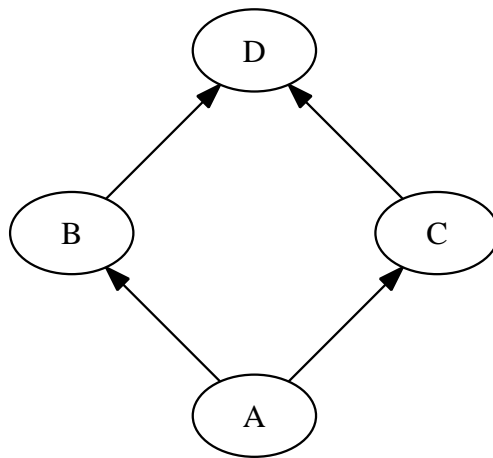
## Introduction

### Bayesian Networks

- Bayesian Networks (BNs) are a method for representing complex multivariate probability distributions.
  - Graphical Structure: Directed Acyclic Graph
    - \* Each **node** represents a variable
    - \* An **edge** represents an association between two variables
  - Parametrical Model: Multivariate Gaussian
    - \* Data are independent, multivariate Gaussian
    - \* Parameters are multivariate Gaussian
    - \* Marginal distribution of each node (variable) is a multiple regression

# Introduction

## Bayesian Networks



	C1	C2	C3	C4	C5	C6	C7
1	MULT	dis_nlcd	lrbs_bw5	lbl_bug	lnutr	lanc	phstvi
2		-0.5172	-0.6999	-0.0328	-0.8906	0.8872	0.0276
3		-0.1244	-0.8880	-1.0065	-0.3457	0.4407	0.0875
4		-1.1625	0.1419	-0.1269	-0.1849	-0.9393	-0.6673
5		-0.7843	-0.5823	-1.0254	-0.6390	1.1841	0.4350
6		1.7877	0.3384	-0.7535	1.3480	0.5916	0.2313
7		-0.8085	0.6959	1.0693	-1.4031	0.3346	0.5189
8		-0.9752	1.9179	-0.8051	-0.4138	0.3969	0.1594
9		-0.6385	0.9884	-0.1878	-1.0803	0.4192	0.3272
10		-1.2050	0.7004	-0.4008	-0.7890	-1.1716	-3.3154
11		-0.4743	1.2561	0.6133	0.3809	0.7179	-0.4636
12		-1.1888	0.3289	0.1516	-0.3833	-0.4112	-0.7153
13		-0.8671	0.4069	1.6392	-0.6531	-0.1222	-0.0443
14		-1.0799	0.3363	0.7176	-0.6891	-0.3721	-0.3554
15		-0.9426	-0.2387	-0.1882	-0.7810	0.3884	-0.0682
16		-0.2387	1.3774	-0.5428	-0.9407	0.0727	0.2792
17		-1.1721	0.1311	0.1139	-0.6531	-0.7483	-0.3678
		-0.4136	1.0537	1.3146	-0.1193	-0.3580	0.2792

- **Important point:** despite the arrows, BNs represent a probability distribution and do not imply causations, only associations.

## Introduction

### Reversible Jump Markov chain Monte Carlo

- *MCMC* is a method for simulating a probability distribution that cannot be directly simulated. MCMC is often used for model selection, especially in very large model spaces.
- *RJMCMC* is a type of MCMC that allows for dimensional changes in the probability distribution being simulated. RJMCMC can be used for model selection in cases where dimensionality may change, such as:
  - ARIMA time series models
  - Gaussian Mixtures

## RJMCMC for BNs

- RJMCMC is suited to searching for BNs because:
  - A change in the graphical structure of a BN results in a change in the number of parameters
  - The number of possible structures of BNs increases super-exponentially in the number of variables.

## RJMCMC for BNs

- How it works:
  - RJMCMC randomly “walks around” the space of possible model structures by changing one edge at a time – called structural learning.
  - At each step in its “walk”, all of the model parameters are updated – called parametrical learning.
  - At the end, you have a list of all of the model structures it visited at each step and their corresponding set of parameters.

## Example: MAHA-MAIA

- Subset of numerical data taken from the MAHA-MAIA study.
- Six variables chosen, with help of a domain expert:
  - Insect IBI - Insect Index of Biotic Integrity
  - Sediment Disturbance - A log-scale metric of excess grain-size diameter
  - Environmental disturbance - Combined percentages of disturbance (urban, agricultural, and mining disturbances)
  - pH
  - Natural logarithm of nutrients (maximum of either N or P)
  - Natural logarithm of slightly translated Acid Neutralizing Capacity



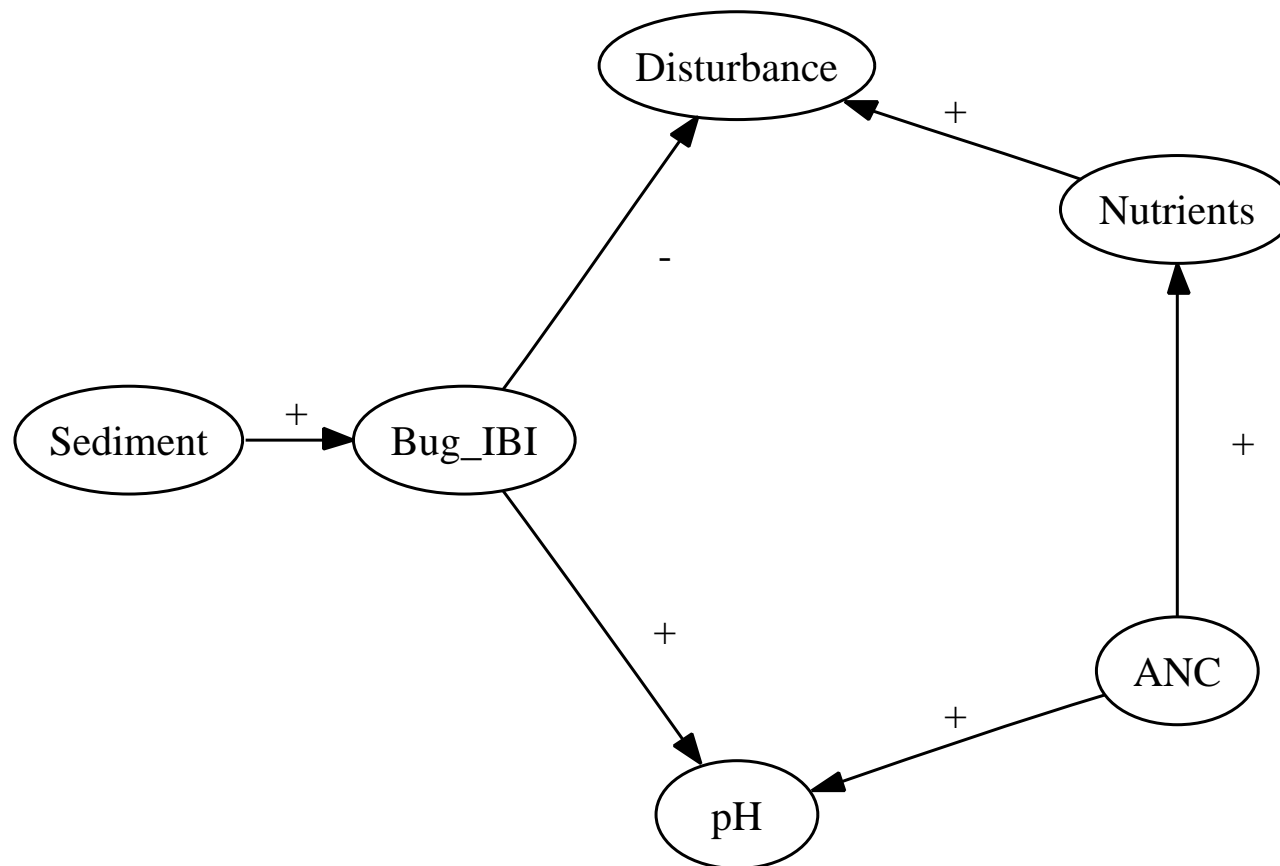
## Example: MAHA-MAIA

### Implementation notes

- To avoid numerical issues, variables were standardized.
- Three graphical structures are presented:
  - The structure with the *maximum posterior probability*,
  - The “average” model, which gives us a sense of the likelihood of an association between pairs of variables,
  - A reference model obtained from the package TETRAD IV.
- +’s and -’s, derived from the learned parameters, denote the type of quantitative association between variables.

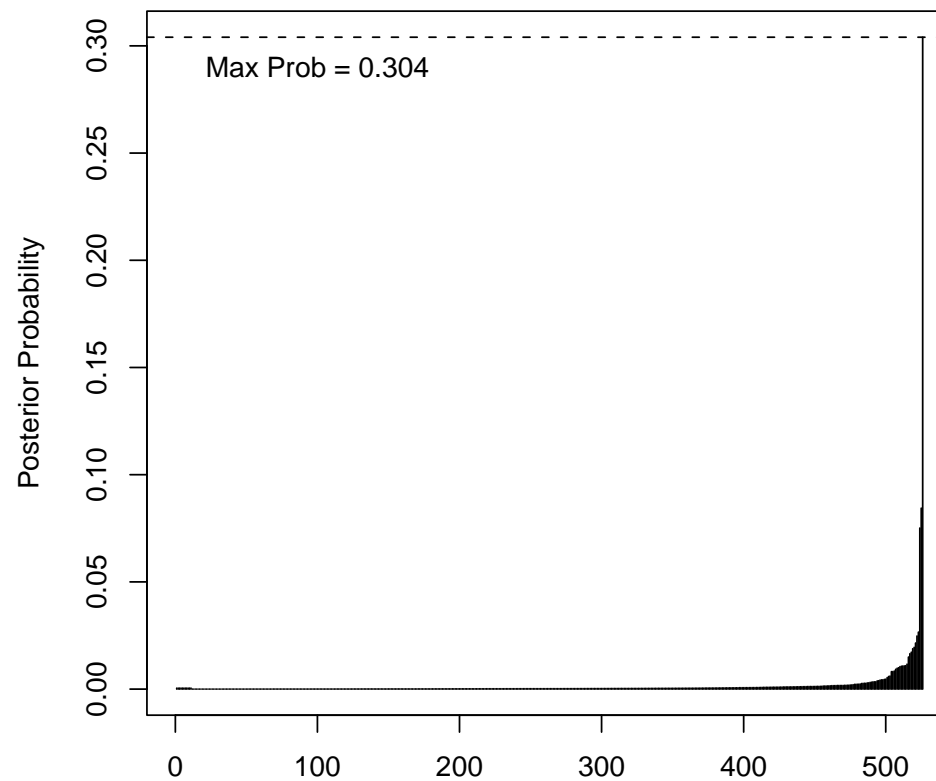
## Example: MAHA-MAIA

### Maximum posterior probability structure



## Example: MAHA-MAIA

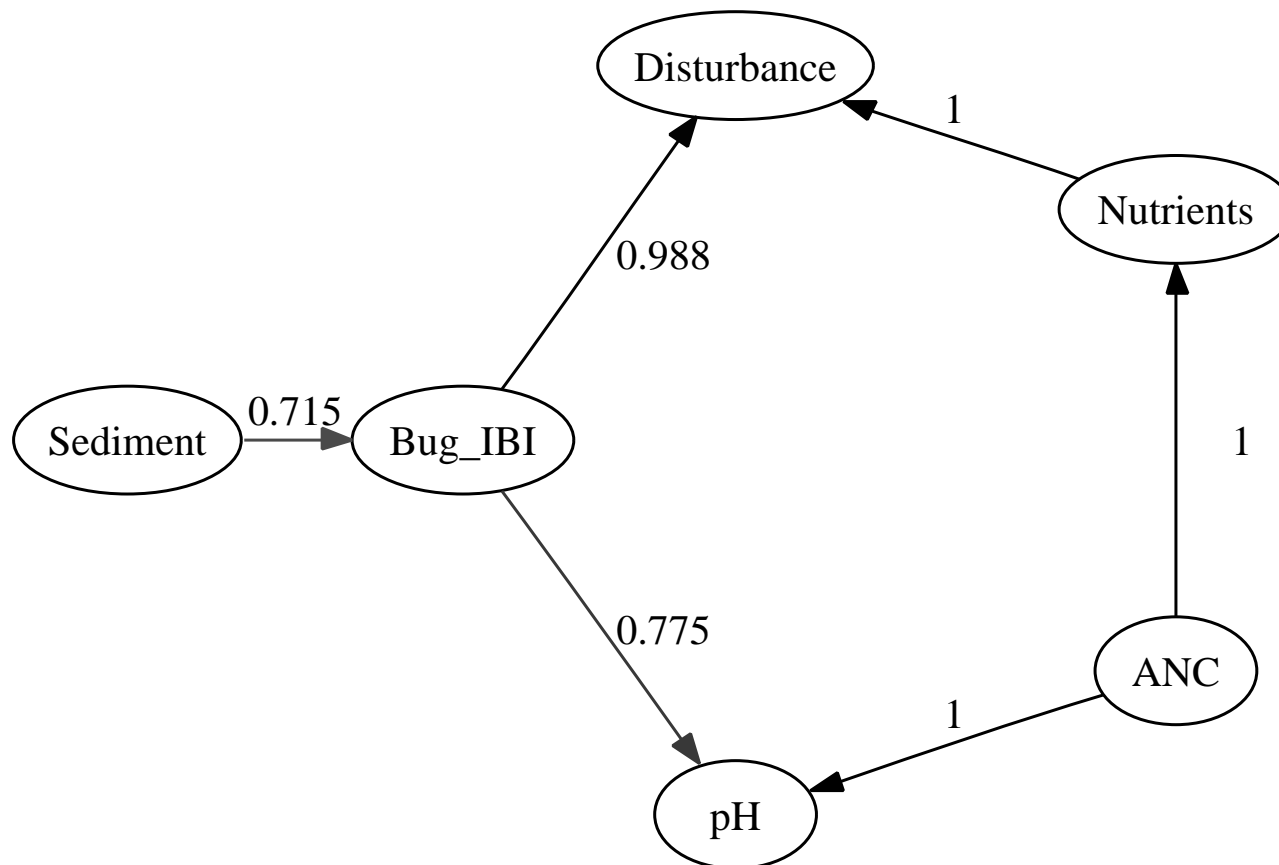
### Posterior probabilities of all visited structures



# Example: MAHA-MAIA

## “Average” structure

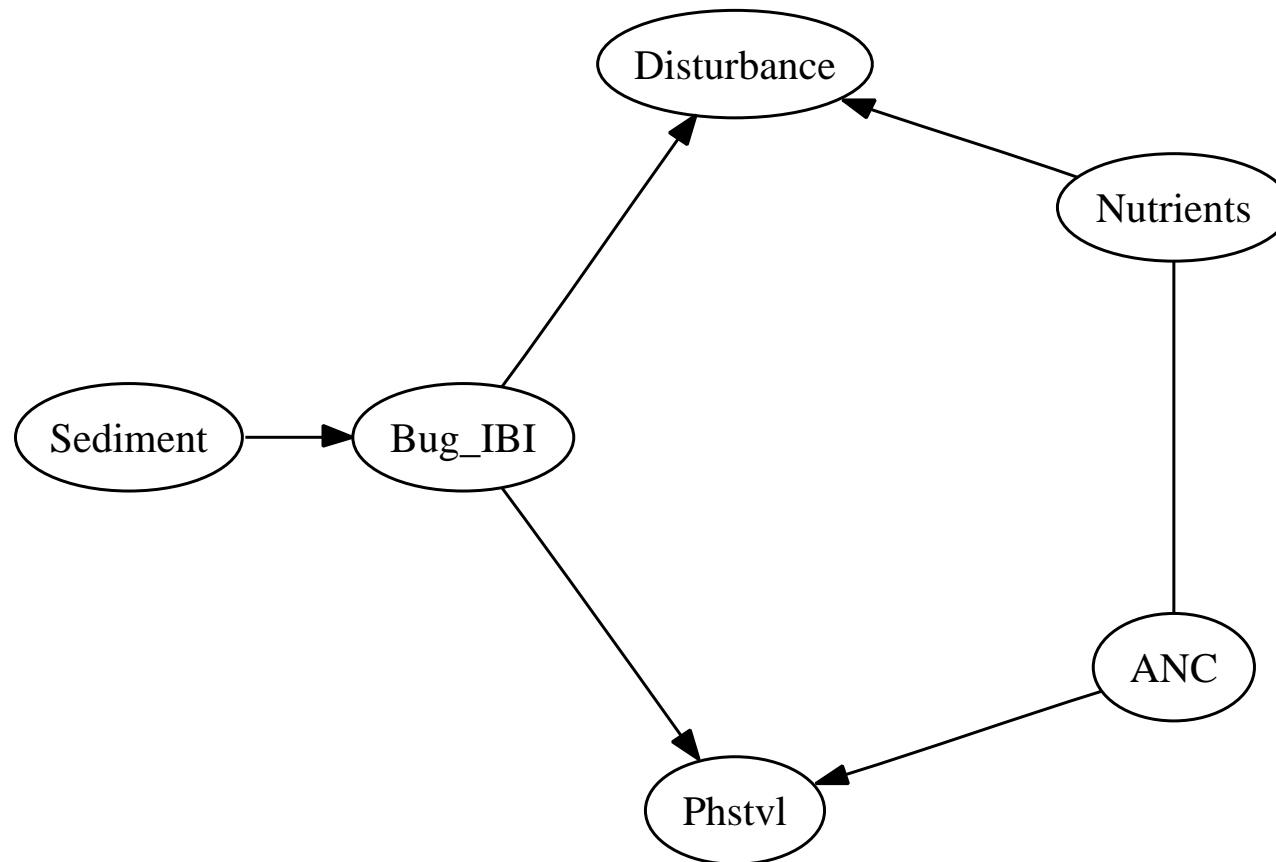
*Note: edges with posterior probability < .3 not shown*



## Example: MAHA-MAIA

### Structure discovered by TETRAD IV

*Note: undirected edge on right can be oriented either direction*



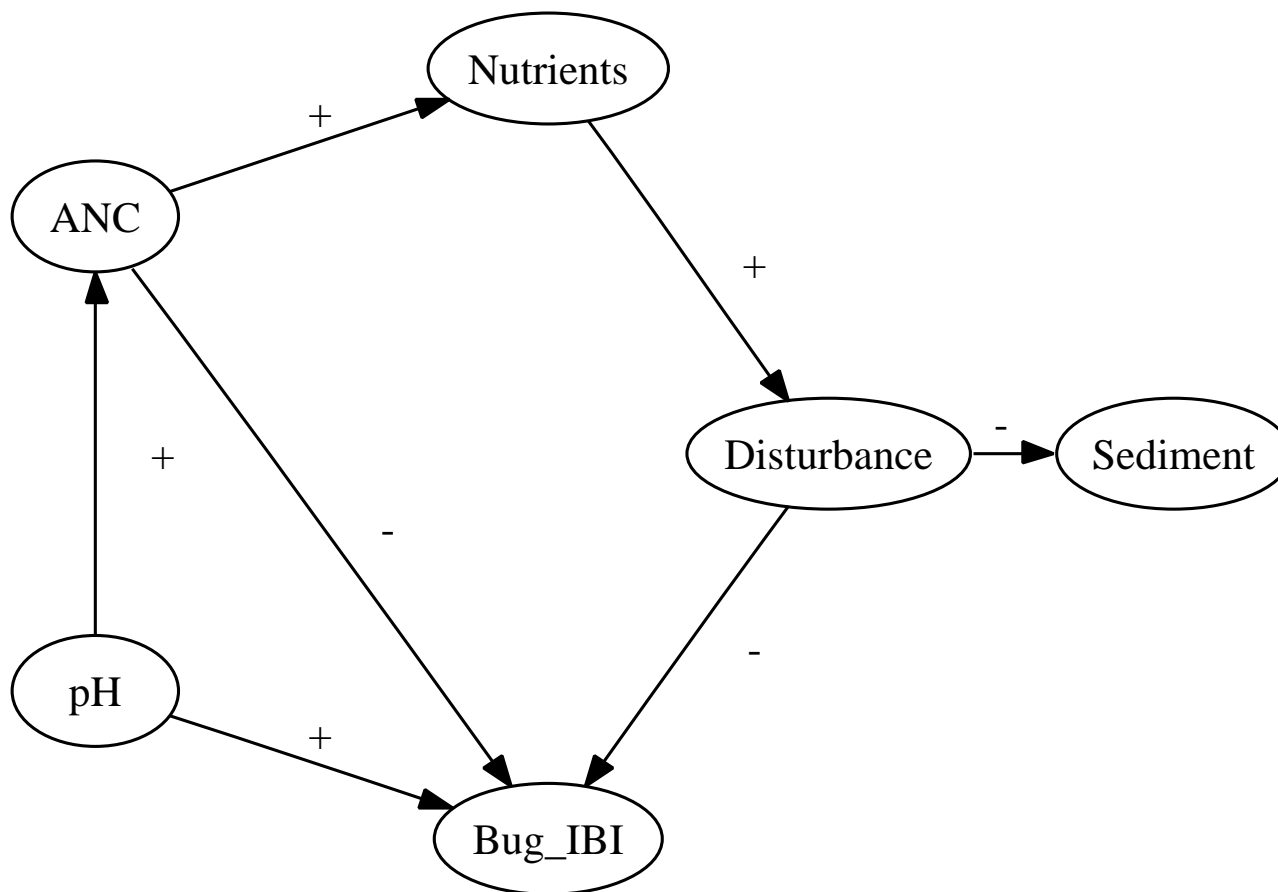
## Example: MAHA-MAIA

### Discussion

- The +'s and -'s seem reasonable, and the structures agree, but the direction of the edges from Bug\_IBI seem to preclude a causal interpretation.
- Recall that BNs encode a joint probability distribution and don't necessarily suggest a causal relation.
- Try again, but remove models from consideration that include edges emanating from Bug\_IBI.

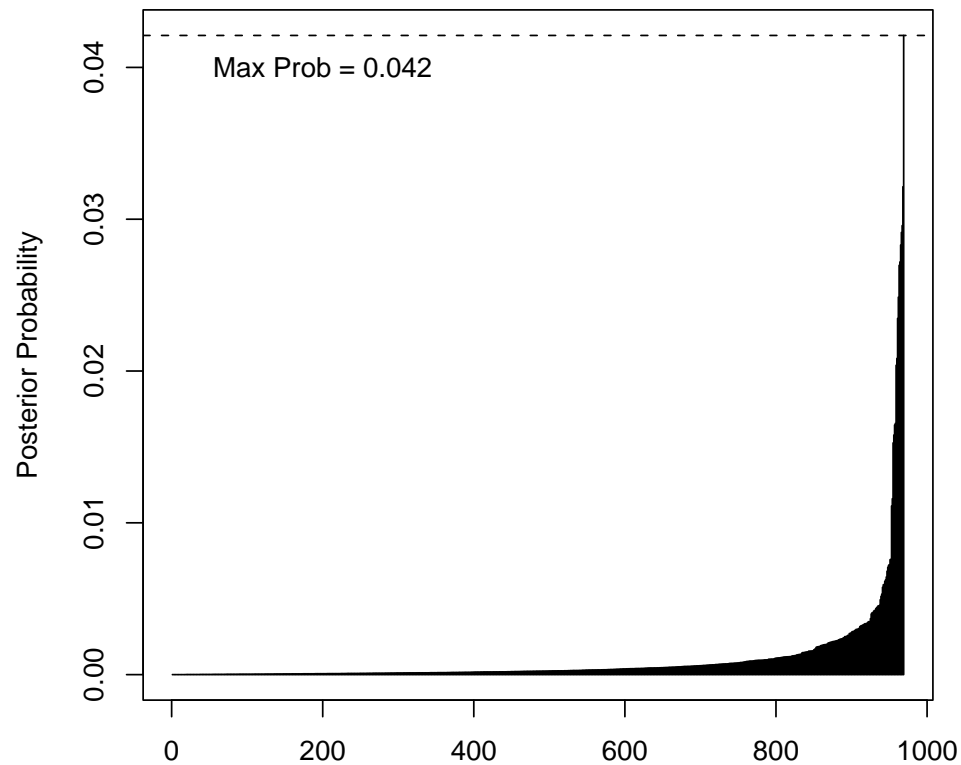
## Example: MAHA-MAIA (2)

### Maximum posterior probability structure



## Example: MAHA-MAIA (2)

### Posterior probabilities of all visited structures

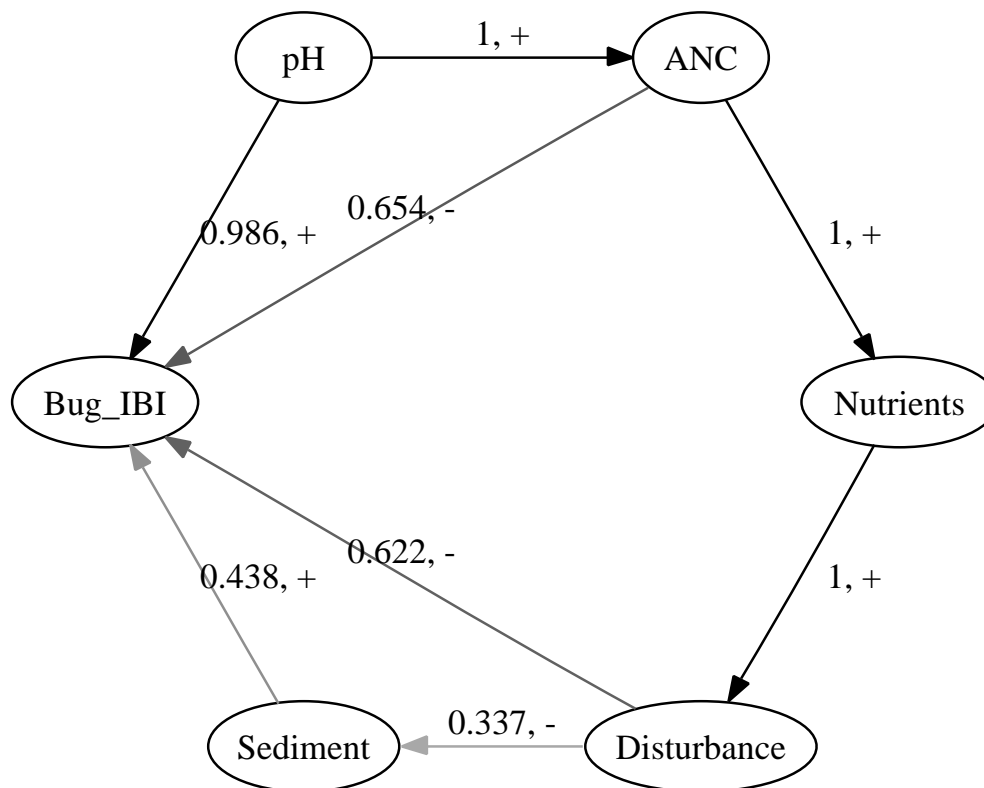




## Example: MAHA-MAIA (2)

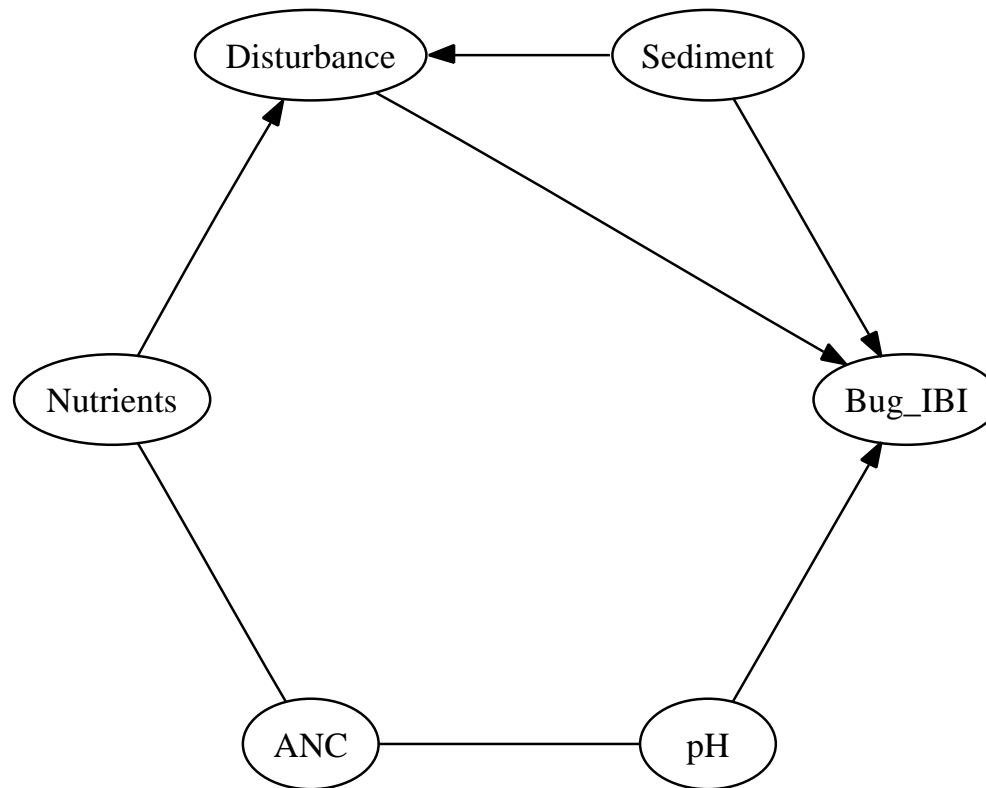
### “Average” structure

*Note: edges with posterior probability < .3 not shown*



## Example: MAHA-MAIA (2)

### Structure discovered by TETRAD IV



## Conclusion

- Finding an optimal BN that fits data is known to be a very hard problem (NP-hard, in fact), making heuristic algorithms a necessity.
- Though many of these algorithms are much faster than RJMCMC, increases in computing power and programming refinements are making its lack of speed less of an issue.
- Furthermore, other methods lack the unique strengths of RJMCMC:
  - Posterior edge probabilities give a measure of the likelihood of association
  - Combined structure discovery and parameter learning
  - Fully Bayesian solution