

### Supplementary Method:

#### Mathematical proof of the empirical null distribution of $\widehat{MR}_{M \times M}(i, j)$

##### Lemma 1:

For the  $M \times M$  randomly assigned matrix  $C$ ,  $\text{Rank}(i \rightarrow j) \perp \text{Rank}(j \rightarrow i) | C(i, j)$ .

Proof: With given  $C(i, j)$ ,  $\text{Rank}(i \rightarrow j)$  only depends on  $C(i, \cdot), \cdot \neq j$  while  $\text{Rank}(j \rightarrow i)$  only depends on  $C(\cdot, j), \cdot \neq i$ . Since  $C(i, \cdot) \perp C(\cdot, j)$  for  $i, \cdot \neq i, j$  and  $\cdot, j! = i, j$ , we have  $\text{Rank}(i \rightarrow j) \perp \text{Rank}(j \rightarrow i) | C(i, j)$   $\square$

##### Proposition 1:

For the  $M \times M$  matrix randomly assigned from the distribution with pdf  $f$  and cdf  $F$ , the joint distribution of  $\text{Rank}(i \rightarrow j)$  and  $\text{Rank}(j \rightarrow i)$  is:

$$P(\text{Rank}(i \rightarrow j) = x, \text{Rank}(j \rightarrow i) = y) = \int \binom{M-1}{x-1} ((1-F(k))^{x-1} F(k)^{M-x}) \cdot \binom{M-1}{y-1} ((1-F(k))^{y-1} F(k)^{M-y}) \cdot dF(k)$$

Proof:

$$\begin{aligned} & \text{Prob}(\text{Rank}(j \rightarrow i) = x | C(i, j) = k) = \\ & \text{Prob}(C_{i(j_m)} > k, C_{i(j_n)} < k, m = 1 \dots x-1, n = x+1 \dots N) = \\ & C_{N-1}^{x-1} (1-F(k))^{x-1} F(k)^{N-x} \end{aligned}$$

, in which  $j_m = 1 \dots j-1$  and  $j_n = j+1 \dots N$   
, and  $C_{i(j)}$  is the  $j$ th value in  $C(i, \cdot)$  in decreasing order

By Lemma 1,

$$\begin{aligned} & P(\text{Rank}(i \rightarrow j), \text{Rank}(j \rightarrow i) | C(i, j)) = \\ & P(\text{Rank}(i \rightarrow j) | C(i, j)) * P(\text{Rank}(j \rightarrow i) | C(i, j)) \end{aligned}$$

$$\begin{aligned} & P(\text{Rank}(i \rightarrow j) = x, \text{Rank}(j \rightarrow i) = y | C(i, j) = k) = \\ & (C_{N-1}^{x-1} (1-F(k))^{x-1} F(k)^{N-x}) \cdot (C_{N-1}^{y-1} (1-F(k))^{y-1} F(k)^{N-y}) \end{aligned}$$

Hence

$$P(\text{Rank}(i \rightarrow j) = x, \text{Rank}(j \rightarrow i) = y) = \int \binom{M-1}{x-1} ((1-F(k))^{x-1} F(k)^{M-x}) \cdot \binom{M-1}{y-1} ((1-F(k))^{y-1} F(k)^{M-y}) \cdot dF(k)$$

$\square$

##### Proposition 2:

The empirical null distribution of  $\widehat{MR}_{M \times M}(i, j)$  does not depend on the empirical distribution of  $C(i, j)$

Proof: Since  $C(i, j)$  follows the pdf  $f$ , by Proposition 1:

$$\begin{aligned} & P(\text{Rank}(i \rightarrow j) = x, \text{Rank}(j \rightarrow i) = y | C(i, j) = k) = \\ & (C_{N-1}^{x-1} (1-F(k))^{x-1} F(k)^{N-x}) \cdot (C_{N-1}^{y-1} (1-F(k))^{y-1} F(k)^{N-y}) \end{aligned}$$

Hence the  $MR(i, j)^2 = \text{Rank}(i \rightarrow j) \cdot \text{Rank}(j \rightarrow i)$  with given  $C(i, j) = k$  follows a multiplication of two identical binomial distributions with  $n = N-1$  and  $P = 1-F(S_{ij}) \sim U(0, 1)$ , which does not depend on the empirical distribution of  $C(i, j)$ .  $\square$

To generate an empirical null distribution of  $\widehat{MR}_{M \times M}(i, j)$ , we can randomly draw  $\tilde{p}$  from

$U(0,1)$  and generate  $Rank(i \rightarrow j)$  and  $Rank(j \rightarrow i)$  from two random numbers  $X_{i,j,1}$  and  $X_{i,j,2}$  from the binomial distribution  $Bin(\tilde{p}, n - 1)$  and  $\widehat{MR}(i, j)' = \sqrt{(X_{i,j,1} + 1) \cdot (X_{i,j,2} + 1)}$ . Hence the  $P(Rank(i \rightarrow j) = x, Rank(j \rightarrow i) = y)$ .

**Theorem 1:**

Empirical null distribution of  $MR(i, j)$  can be simulated by

$$MR_{i,j} = \sqrt{(X_{i,j,1} + 1) \cdot (X_{i,j,2} + 1)}$$

$$X_{i,j,1} \sim Binom(n - 2, p_{ij}), X_{i,j,2} \sim Binom(n - 2, p_{ij})$$

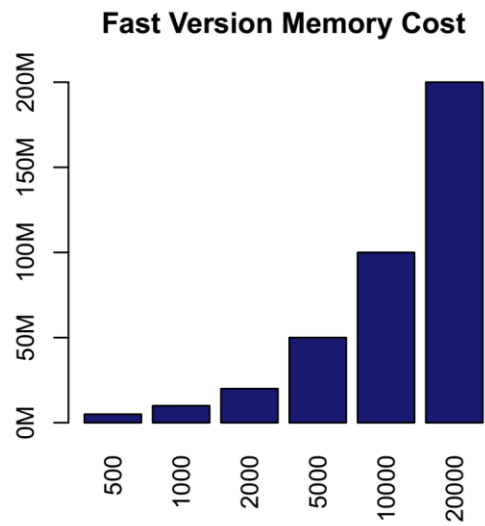
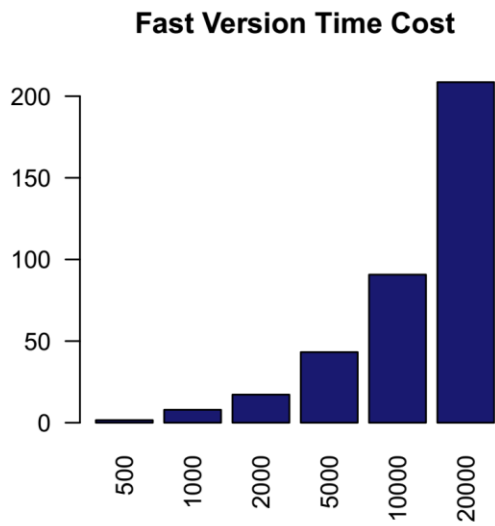
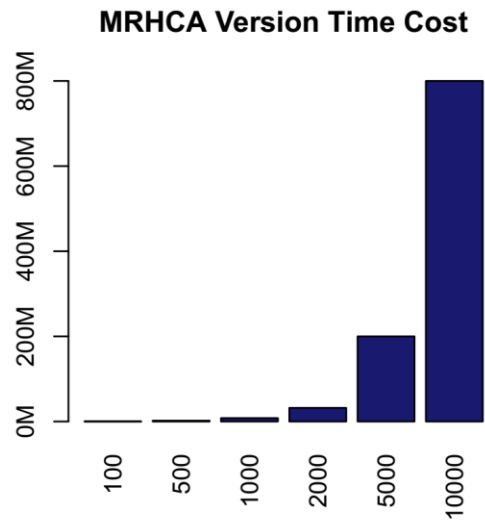
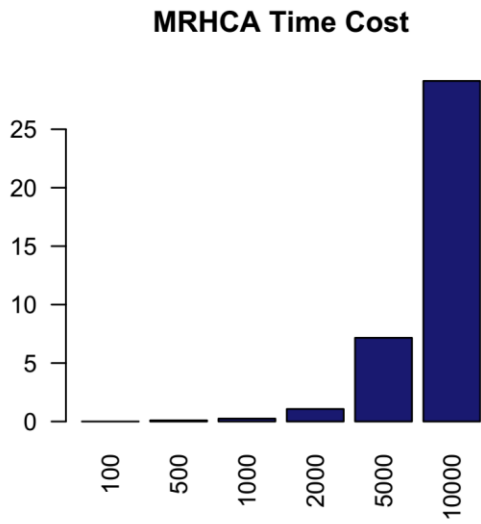
$$p_{ij} \sim U(0,1)$$

Proof:

By Proposition 1 and 2, and the discussion above, the random numbers simulated as described in the theorem form empirical null distribution of  $\widehat{MR}_{M \times M}(i, j)$ .  $\square$

**Memory and Computational cost analysis:**

We have analyzed the memory and computational cost of MRHCA and the fast version over randomized gene expression profile of 100, 500, 1000, 2000, 5000, 10000 and 20000 genes with 300 samples. The analyses were conducted by R x64 3.4.0 on a desktop with Intel Core i7-7700K CPU and 16.0GB RAM. The memory and algorithm running time are shown in the figure below. The loops in this fast version of the Rank\_Index generation was not accelerated by a C library. Hence the time cost of the fast version is more than the MRHCA algorithm. See more details on <https://github.com/zy26/mrct>.



The Figures show the time (s) and memory (M) costs (y-axis) of the MRHCA and the fast version on simulated data with number of genes shown on x-axis.