

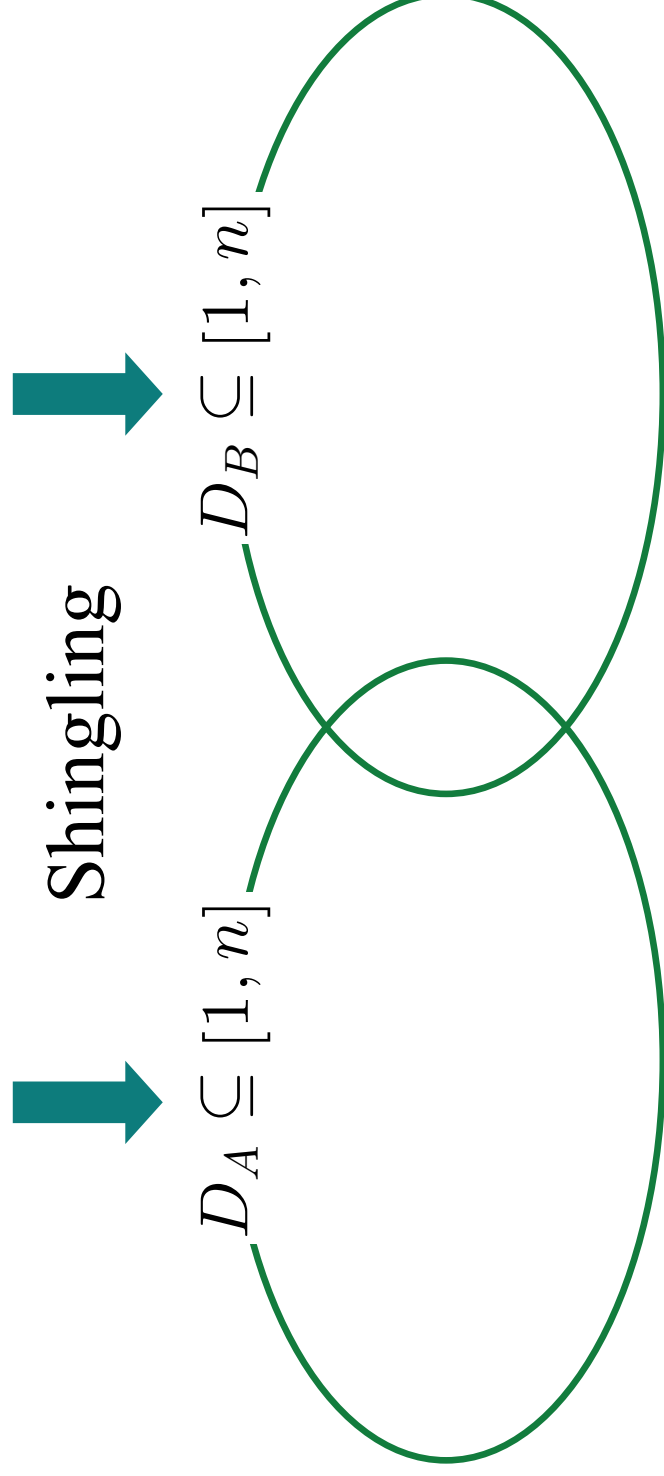
**Improved Lower Bounds for  
Families of  $\varepsilon$ -Approximate  
 $k$ -Restricted Min-Wise  
Independent Permutations**

**Toshiya Itoh    Tatsuya Nagatani**

Tokyo Institute of Technology

# Resemblance $r(A, B)$ of Documents $A$ and $B$ .

Document  $A$                       Document  $B$



$$r(A, B) = \frac{\|D_A \cap D_B\|}{\|D_A \cup D_B\|}$$

## Estimation of $r(A, B)$ .

(1) Choose  $\pi_1, \pi_2, \dots, \pi_\ell \in S_n$  independently

(2) Define sketches  $S_A$  of  $A$  and  $S_B$  of  $B$  by

$$\begin{aligned} S_A &= (\min\{\pi_1(D_A)\}, \min\{\pi_2(D_A)\}, \dots, \min\{\pi_\ell(D_A)\}) \\ &= (s_{A,1}, s_{A,2}, \dots, s_{A,\ell}) \end{aligned}$$

$$\begin{aligned} S_B &= (\min\{\pi_1(D_B)\}, \min\{\pi_2(D_B)\}, \dots, \min\{\pi_\ell(D_B)\}) \\ &= (s_{B,1}, s_{B,2}, \dots, s_{B,\ell}) \end{aligned}$$

(3) Compute  $\tilde{r}_\ell(A, B)$ , estimation of  $r(A, B)$ , by

$$\tilde{r}_\ell(A, B) = \frac{\|\{i \in [1, \ell] : s_{A,i} = s_{B,i}\}\|}{\ell}$$

$$\lim_{\ell \rightarrow \infty} \tilde{r}_\ell(A, B) = r(A, B)$$

$\mathcal{F} \subseteq S_n$  estimates  $r(A, B) \Leftrightarrow \mathcal{F} : k$ -restricted min-wise independent

## Def. 1.1

$\mathcal{F} \subseteq S_n : \varepsilon$ -Approximate  $k$ -Restricted Min-Wise Independent

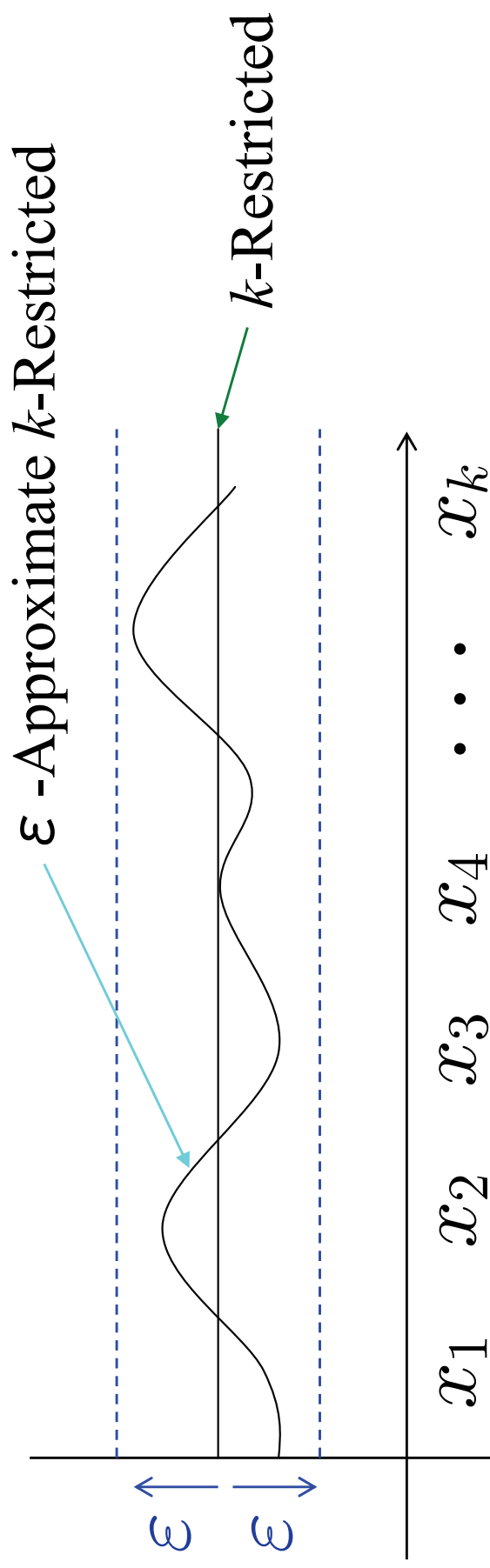
$$\forall X \subseteq [1, n] \quad \text{s. t. } \|X\| \leq k \quad \forall x \in X$$

$$\left| \Pr_{\pi \in \mathcal{F}} [\min\{\pi(X)\}] = \pi(x) \right| \leq \frac{\varepsilon}{\|X\|},$$

where  $\pi \in \mathcal{F}$  is chosen uniformly at random.

$\varepsilon = 0 \Rightarrow \mathcal{F} : k\text{-Restricted Min-Wise Independent}$

$k = n \Rightarrow \mathcal{F} : \varepsilon\text{-Approximate Min-Wise Independent}$



## Known Results (1)

	Upper Bound	Lower Bound
Min-Wise	$\text{lcm}(n, n-1, \dots, 1)$	$\text{lcm}(n, n-1, \dots, 1)$
$k$ -Restricted Min-Wise	$O(n^k e^k)$	$\Omega\left(\binom{n-1}{\lfloor (k-1)/2 \rfloor}\right)$

## Known Results (2) $\varepsilon$ -Approximate $k$ -Restricted

Upper Bounds	$2^{4k+O(k)} k^2 \log \log(n/\varepsilon) \quad (\text{constructive})$ $O\left(\frac{k^2}{\varepsilon^2} \log\left(\frac{n}{\varepsilon}\right)\right) \quad (\text{nonconstructive})$
Lower Bounds	$\Omega(k^2 (1 - \sqrt{8\varepsilon})) \quad (\text{uniform})$ $\Omega\left(\min\left\{k^2 \log\left(\frac{n}{\varepsilon}\right), \frac{\log(1/\varepsilon)(\log n - \log \log(1/\varepsilon))}{\varepsilon^{1/3}}\right\}\right) \quad (\text{biased})$ $\Omega\left(k \sqrt{\frac{1}{\varepsilon} \log\left(\frac{n}{k}\right)}\right) \quad (\text{uniform})$

$$\mathcal{F} = \{\pi_1, \pi_2, \dots, \pi_d\} \subseteq \mathcal{S}_n : \varepsilon\text{-Approx. } k\text{-Restricted Min-Wise}$$

$$s = k/3, \quad L = n/s, \quad N = L - 1$$

$$[1, n] = \{1, 2, \dots, n\} \\ = \underbrace{\{1, \dots, h, \dots, s, s+1, \dots, 2s, \dots\}}_{X_0}, \underbrace{\{s+1, \dots, 2s, \dots\}}_{X_1}, \dots, \underbrace{\{(L-1)s+1, \dots, Ls\}}_{X_N}$$

$$\forall h \in [1, s] \quad U_h = \begin{bmatrix} \pi_1 & \dots & \pi_j & \dots & \pi_d \\ \vdots & & \vdots & & \vdots \\ u_{ij}^h & \dots & \dots & & \dots \\ X_0 \cup X_1 & & \vdots & & \vdots \\ X_0 \cup X_i & & \vdots & & \vdots \\ X_0 \cup X_N & & \vdots & & \vdots \end{bmatrix}$$

$$u_{ij}^h = \begin{cases} 1/\sqrt{d} & \min\{\pi_j(X_0 \cup X_i)\} = \pi_j(h) \\ 0 & \text{otherwise} \end{cases}$$



$$X_0 = \{1, 2, 3, 4\}, X_1 = \{5, 6, 7, 8\}, X_2 = \{9, 10, 11, 12\}, \dots$$

$$U_2 = \begin{bmatrix} \pi_1 & \dots & \pi_j & \dots & \pi_d \\ u_{1j}^2 \\ u_{2j}^2 \\ \vdots \\ \vdots \end{bmatrix}$$

$$X_0 \cup X_1 = \{1, 2, 3, 4, 5, 6, 7, 8\}$$

$$X_0 \cup X_1 = \{1, 2, 3, 4, 9, 10, 11, 12\}$$

$$\vdots$$

$$\vdots$$

$$u_{1j}^2 = \begin{cases} 1/\sqrt{d} & \min\{\pi_j(\{1, 2, 3, 4, 5, 6, 7, 8\})\} = \pi_j(2) \\ 0 & \text{otherwise} \end{cases}$$

$$u_{2j}^2 = \begin{cases} 1/\sqrt{d} & \min\{\pi_j(\{1, 2, 3, 4, 9, 10, 11, 12\})\} = \pi_j(2) \\ 0 & \text{otherwise} \end{cases}$$

$$\begin{array}{c}
 \overbrace{\hspace{10em}}^N \\
 \left[ \begin{array}{cccc}
 \frac{\delta_{11}^h}{2s} & \frac{\delta_{12}^h}{3s} & \frac{\delta_{13}^h}{3s} & \frac{\delta_{1N}^h}{3s} \\
 \frac{\delta_{12}^h}{3s} & \frac{\delta_{22}^h}{2s} & \frac{\delta_{23}^h}{3s} & \dots \\
 \frac{\delta_{13}^h}{3s} & \frac{\delta_{23}^h}{3s} & \frac{\delta_{33}^h}{2s} & \dots \\
 \vdots & \vdots & \vdots & \vdots \\
 \frac{\delta_{1N}^h}{3s} & \frac{\delta_{2N}^h}{3s} & \frac{\delta_{3N}^h}{3s} & \frac{\delta_{NN}^h}{2s}
 \end{array} \right]
 \end{array}$$

$$\forall h \in [1, s] \quad V_h = (v_{ij}^h) = U_h U_h^T =$$

## Proposition 2.1

- (i)  $\forall i \in [1, N] \quad \frac{1-\varepsilon}{2s} \leq v_{ii}^h \leq \frac{1+\varepsilon}{2s}$
- (ii)  $\forall i, j \in [1, N] (i \neq j) \quad \frac{1-\varepsilon}{3s} \leq v_{ij}^h \leq \frac{1+\varepsilon}{3s}$

## Proof of Proposition 2.1

$$X_0 = \{1, 2, 3, 4\}, X_1 = \{5, 6, 7, 8\}, X_2 = \{9, 10, 11, 12\}, \dots$$

$$V_2 = (v_{ij}^2) = U_2 U_2^T$$

$$\begin{aligned} v_{11}^2 &= (u_{11}^2, u_{12}^2, \dots, u_{1d}^2) \cdot (u_{11}^2, u_{12}^2, \dots, u_{1d}^2)^T \\ &= \Pr[\min\{\underbrace{\pi(\{1, 2, 3, 4\})}_{X_0}, \underbrace{\pi(\{5, 6, 7, 8\})}_{X_1}\}] = \frac{1 \pm \varepsilon}{8} \end{aligned}$$

$$\begin{aligned} v_{12}^2 &= (u_{11}^2, u_{12}^2, \dots, u_{1d}^2) \cdot (u_{21}^2, u_{22}^2, \dots, u_{2d}^2)^T \\ &= \Pr[\min\{\underbrace{\pi(\{1, 2, 3, 4\})}_{X_0}, \underbrace{\pi(\{9, 10, 11, 12\})}_{X_2}\}] = \frac{1 \pm \varepsilon}{12} \end{aligned}$$

$$V = \begin{matrix} \overleftrightarrow{\|\mathcal{F}\|} \\ \begin{bmatrix} U_1 \\ U_2 \\ \vdots \\ U_s \end{bmatrix} \end{matrix} [U_1^T, U_2^T, \dots, U_s^T] = \begin{bmatrix} V_1 & 0 & \dots & 0 \\ 0 & V_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & V_s \end{bmatrix}$$

## Proposition 2.2

$\mathcal{F} \subseteq S_n$ :  $\varepsilon$ -Approximate  $k$ -Restricted Min-Wise

$$\begin{aligned} \|\mathcal{F}\| &\geq \text{rank}(V) \\ &= \text{rank}(V_1) + \text{rank}(V_2) + \dots + \text{rank}(V_s). \end{aligned}$$

$$A = \begin{bmatrix} a_{11} & a & a & \cdots & a \\ a & a_{22} & a & \cdots & a \\ a & a & a_{33} & \cdots & a \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a & a & a & \cdots & a_{tt} \end{bmatrix}$$

$\overbrace{\hspace{10em}}^t$ 
 $\overbrace{\hspace{10em}}^t$

### Proposition 2.3

(C1)  $\forall i, j \in [1, t] (i \neq j) \quad a_{ij} = a > 0$

(C2)  $\min \{ a_{11}, a_{22}, \dots, a_{tt} \} > a$

$\Rightarrow A : \text{nonsingular} \equiv \text{rank}(A) = t$

## Proof of Proposition 2.3

$$\begin{aligned}
 A &= \begin{bmatrix} a_{11} & a & \cdots & a \\ a & a_{12} & \cdots & a \\ \vdots & \vdots & \ddots & \vdots \\ a & a & \cdots & a_{tt} \end{bmatrix} = \begin{bmatrix} a & a & \cdots & a \\ a & a & \cdots & a \\ \vdots & \vdots & \ddots & \vdots \\ a & a & \cdots & a \end{bmatrix} + \begin{bmatrix} a_{11}-a & 0 & \cdots & 0 \\ 0 & a_{22}-a & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & a_{tt}-a \end{bmatrix} \\
 &= \underbrace{a(1, 1, \dots, 1) \cdot (1, 1, \dots, 1)^T}_{\text{positive semidefinite}} + \underbrace{\text{diag}(a_{11} - a, a_{22} - a, \dots, a_{tt} - a)}_{\text{positive definite}}
 \end{aligned}$$

$$(C2) \quad \min\{a_{11}, a_{22}, \dots, a_{tt}\} > a$$

$$\Rightarrow a_{11} - a > 0, a_{22} - a > 0, \dots, a_{tt} - a > 0$$

$$A : \text{nonsingular} \quad \equiv \quad \text{rank}(A) = t$$

# Generalized Ramsey Number

$C_m = \{c_1, c_2, \dots, c_m\}$  a set of  $m$  colors

$\chi : E \rightarrow C_m$  edge coloring of  $K_\ell = (V, E)$

$R(t_1, t_2, \dots, t_m)$

$\min \ell$  s.t.  $\exists i \subseteq [1, m] \exists K_{t_i} = (V_i, E_i) \subseteq K_\ell \forall e \in E_i \chi(e) = c_i$

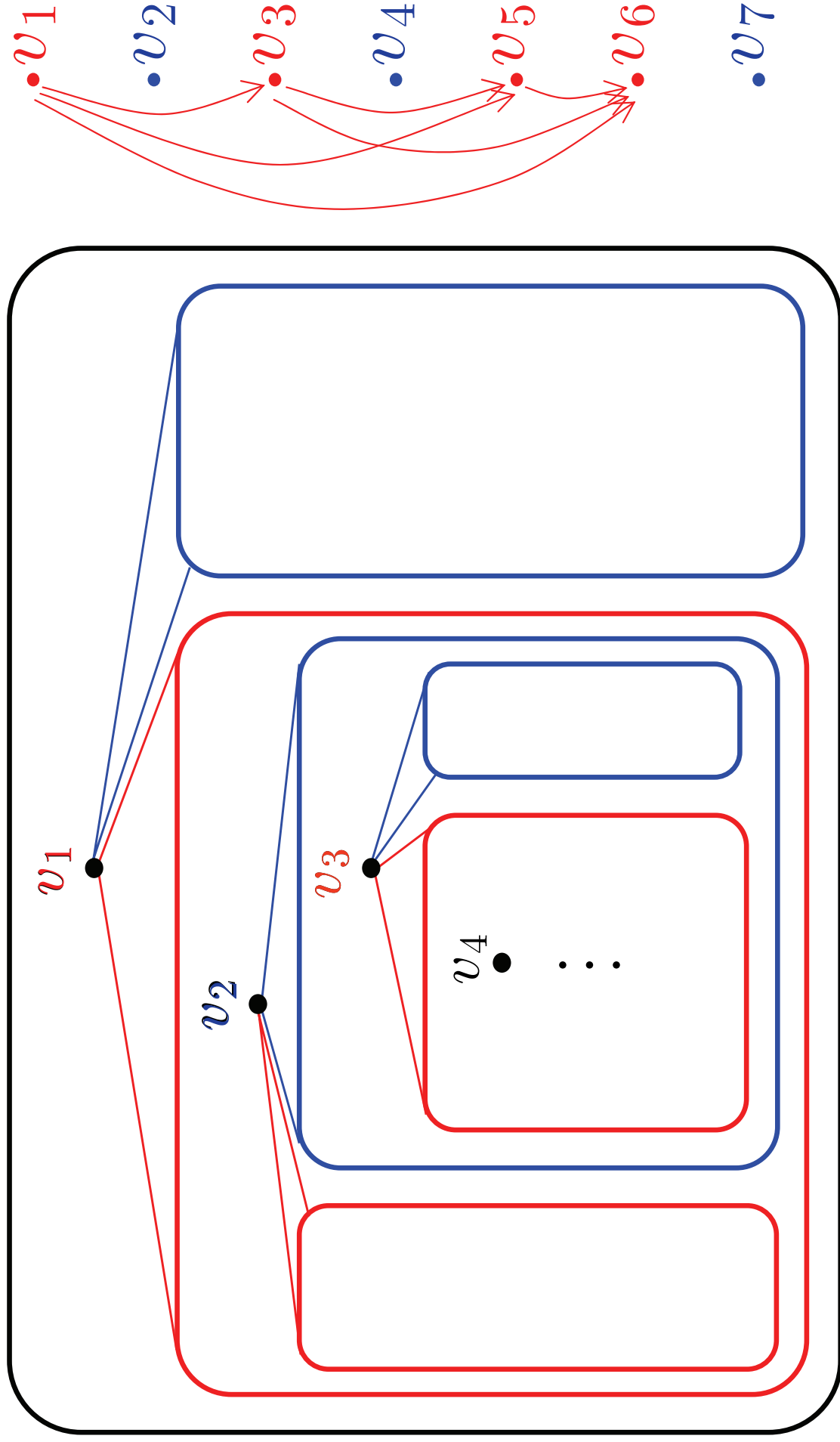
$t_1 = t_2 = \dots = t_m = t \quad R_m(t) \equiv R(t, t, \dots, t)$

**Lemma 3.1**

$\forall m \geq 2 \forall t \geq 1 \quad R_m(t) \leq m^{mt - (m-1)}$

# Proof of Lemma 3.1

$$m = 2, \quad t = 4 \quad n = 2^{2 \cdot 4 - (2-1)} = 2^7$$





Analysis for  $\text{rank}(V_h)$

$$\|\mathcal{F}\| \geq \text{rank}(V_1) + \text{rank}(V_2) + \dots + \text{rank}(V_s)$$

**Lemma 3.2**

$$\forall 0 < \varepsilon < \frac{1}{5} \quad \forall k \geq 3$$

$\mathcal{F} \subseteq S_n : \varepsilon$ -Approximate  $k$ -Restricted Min-Wise

$$\|\mathcal{F}\| < \frac{k}{2\varepsilon}m \quad (m \geq 1)$$

$$\Rightarrow \forall h \in [1, s] \quad \text{rank}(V_h) \begin{cases} = N & m = 1 \\ \geq \left\lfloor \frac{\log(3n/k)}{m \log m} \right\rfloor & m \geq 2 \end{cases}$$

## Proof of Lemma 3.2

$$\forall h \in [1, s] \quad \forall i, j \in [1, N] \quad (i \neq j)$$

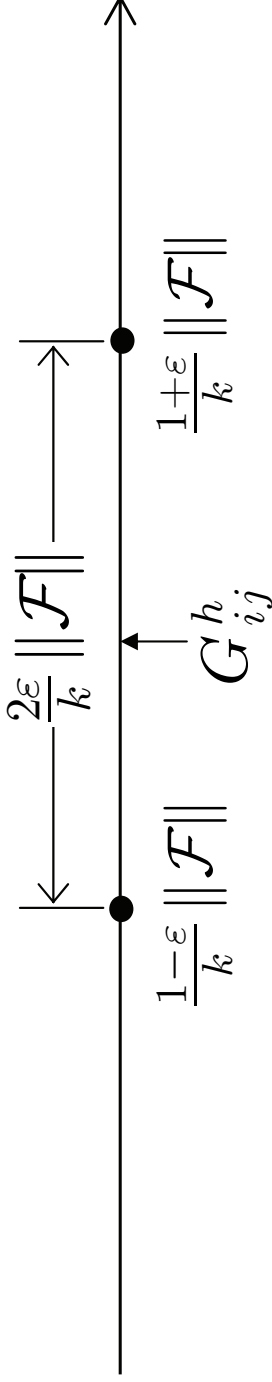
$$G_{ij}^h = \{\pi \in \mathcal{F} : \min\{\pi(X_0 \cup X_i \cup X_j)\} = \pi(h)\}$$

$$G_{ij}^h = \|G_{ij}^h\| : \text{integer}$$

$\mathcal{F} \subseteq S_n : \varepsilon$ -Approximate  $k$ -Restricted Min-Wise

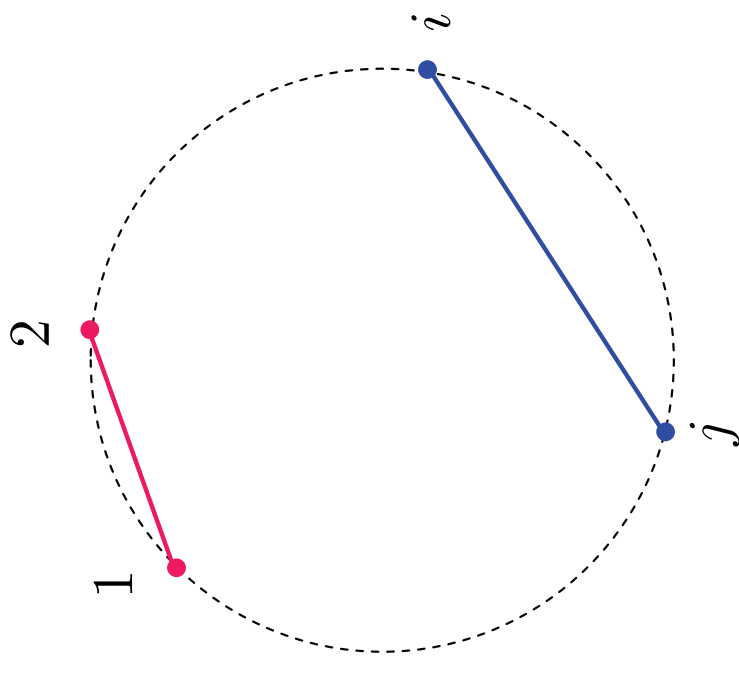
$$\Rightarrow \frac{1-\varepsilon}{3s} \leq \frac{G_{ij}^h}{\|\mathcal{F}\|} \leq \frac{1+\varepsilon}{3s}$$

$$\Rightarrow \frac{1-\varepsilon}{k} \|\mathcal{F}\| \leq G_{ij}^h \leq \frac{1+\varepsilon}{k} \|\mathcal{F}\|$$



- $\frac{2\varepsilon}{k} \|\mathcal{F}\| < m \Rightarrow G_{ij}^h$  takes at most  $m$  (integer) values
- $m$  (integer) values  $\equiv C_m = \{c_1, c_2, \dots, c_m\}$ : a set of  $m$  colors
- $V_h$  is symmetric

$$V_h = \begin{bmatrix} 1 & 2 & \dots & i & j & \dots & N \\ 1 & v_{12} & & & & & \\ 2 & v_{12} & & & & & \\ \vdots & & & & & & \\ i & & & & v_{ij} & & \\ j & & & & v_{ij} & & \\ \vdots & & & & & & \\ N & & & & & & \end{bmatrix}$$



Case 1:  $m = 1$

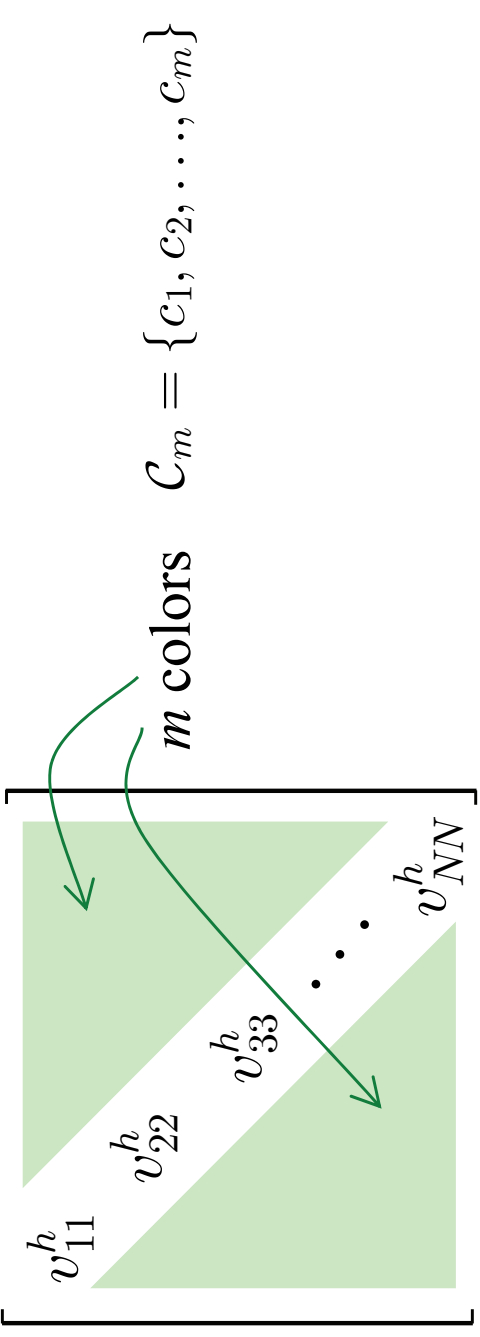
$$\begin{array}{c}
 \overbrace{\hspace{10em}}^N \\
 \left[ \begin{array}{cccc}
 v_{11}^h & v & \dots & v \\
 v & v_{22}^h & \dots & v \\
 v & v & v_{33}^h & \dots \\
 \vdots & \vdots & \vdots & \vdots \\
 v & v & \dots & v_{NN}^h
 \end{array} \right] \\
 \underbrace{\hspace{10em}}_N
 \end{array}$$

$$V_h = (v_{ij}^h) =$$

- $\forall i, j \in [1, N](i \neq j) \quad v_{ij}^h = v > 0$
- $0 < \varepsilon < \frac{1}{5} \Rightarrow \min\{v_{11}^h, v_{22}^h, \dots, v_{NN}^h\} > v$
- **Proposition 2.3**  $\Rightarrow V_h$ : nonsingular  
 $\Rightarrow \text{rank}(V_h) = N$

Case 2 :  $m \geq 2$

$$V_h = (v_{ij}^h) =$$



$$\frac{3n}{k} - 1 = N \geq m^{mt - (m-1)} \geq R_m(t)$$



**Lemma 3.1**

$$\forall \chi : E \rightarrow C_m \text{ of } K_N = (V, E)$$

$$\exists i \in [1, m] \quad \exists K_t = (V_t, E_t) \text{ s.t. } \forall e \in E_t \quad \chi(e) = c_i$$

$$0 < \varepsilon < \frac{1}{5} + \text{Proposition 2.3} \Rightarrow \text{rank}(V_h) \geq \left\lfloor \frac{\log(3n/k)}{m \log m} \right\rfloor$$

# Main Result

## Theorem 4.1

$\forall 0 < \varepsilon < \frac{1}{5} \forall k \geq 3 \quad \mathcal{F} \subseteq S_n : \varepsilon\text{-Approx. } k\text{-Restricted Min-Wise}$

$$\|\mathcal{F}\| = \Omega\left(k \sqrt{\frac{1}{\varepsilon} \log(n/k)}\right).$$

## Proof of Theorem 4.1

$$\forall m \geq 1 \quad \frac{2\varepsilon}{k} \|\mathcal{F}\| < m \quad (\Leftrightarrow) \quad \|\mathcal{F}\| < \frac{k}{2\varepsilon} m$$



## Proposition 2.2 + Lemma 3.2

$$\begin{aligned} \|\mathcal{F}\| &\geq \text{rank}(V_1) + \text{rank}(V_2) + \dots + \text{rank}(V_s) \\ &\geq s \left\lceil \frac{\log(3n/k)}{m \log m} \right\rceil = \frac{k}{3} \left\lceil \frac{\log(3n/k)}{m \log m} \right\rceil \end{aligned}$$

$$\forall m \geq 1 \left[ \left[ \|\mathcal{F}\| < \frac{k}{2\varepsilon}m \Rightarrow \|\mathcal{F}\| \geq \frac{k}{3} \left\lceil \frac{\log(3n/k)}{m \log m} \right\rceil \right] \right]$$

$\frac{k}{2\varepsilon}m$ 
 $\frac{k}{3} \left\lceil \frac{\log(3n/k)}{m \log m} \right\rceil$ 
→ contradiction



$$\forall m \geq 1 \left[ \frac{k}{2\varepsilon}m \leq \frac{k}{3} \left\lceil \frac{\log(3n/k)}{m \log m} \right\rceil \Rightarrow \|\mathcal{F}\| \geq \frac{k}{2\varepsilon}m \right]$$



$$\|\mathcal{F}\| = \Omega \left( k \sqrt{\frac{1}{\varepsilon} \log(n/k)} \right)$$

# Discussions

$$\text{(LB1)} \quad \|\mathcal{F}\| = \Omega(k^2(1 - \sqrt{8\varepsilon}))$$

$$\text{(LB2)} \quad \|\mathcal{F}\| = \Omega\left(\min\left\{k2^{k/2}\log\left(\frac{n}{\varepsilon}\right), \frac{\log(1/\varepsilon)(\log n - \log \log(1/\varepsilon))}{\varepsilon^{1/3}}\right\}\right)$$

$$\text{(Ours)} \quad \|\mathcal{F}\| = \Omega\left(k\sqrt{\frac{1}{\varepsilon}\log(n/k)}\right)$$

$$(1) \quad n \gg k \Rightarrow \text{(Ours)} > \text{(LB1)}$$

$$(2) \quad k < \frac{2}{3}\log(1/\varepsilon) \Rightarrow \text{(Ours)} > k2^{k/2}\log\left(\frac{n}{\varepsilon}\right)$$

$$(3) \quad k \geq \frac{2}{3}\log(1/\varepsilon) \Rightarrow \text{(Ours)} > \frac{\log(1/\varepsilon)(\log n - \log \log(1/\varepsilon))}{\varepsilon^{1/3}}$$

$$\Rightarrow \text{(Ours)} > \text{(LB2)}$$





*Thank you*