

Info-Theo 2 Fall 2024

CODED CACHING

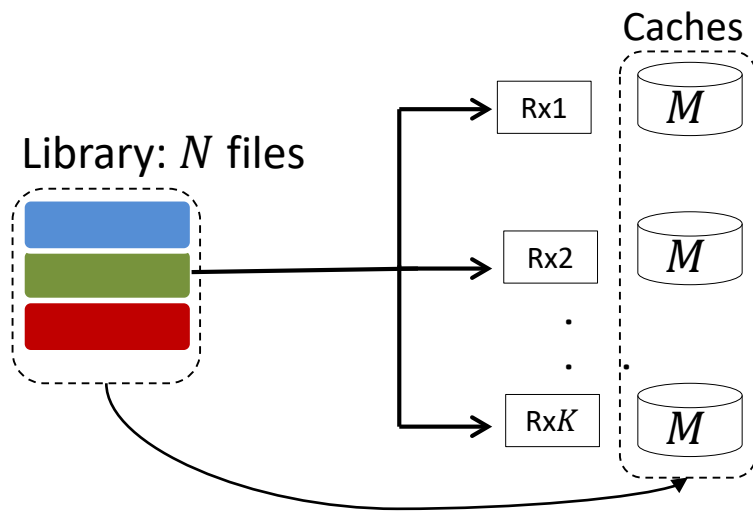
A powerful Information Theoretic Approach

PETROS ELIA

Two Problems Sharing Some Fundamentals

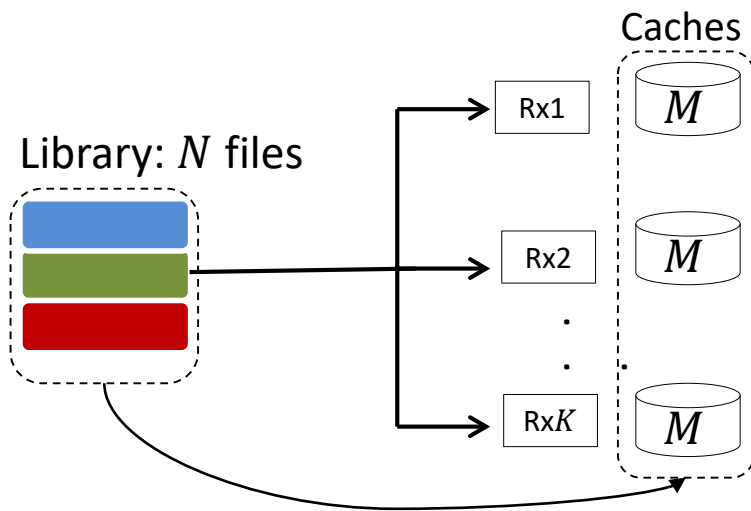
Two Problems Sharing Some Fundamentals

Cache-aided Communications

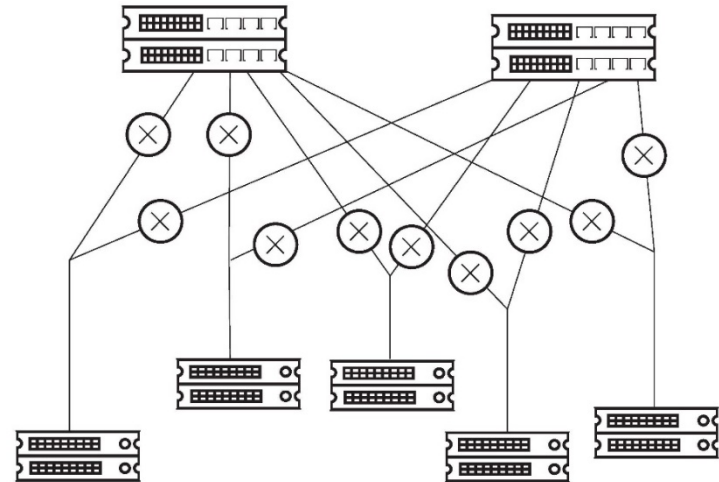


Two Problems Sharing Some Fundamentals

Cache-aided Communications

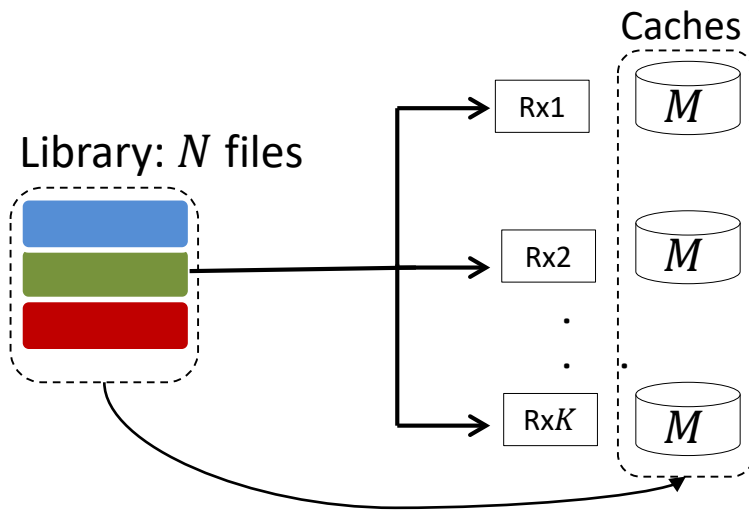


Communications in Distributed Computing

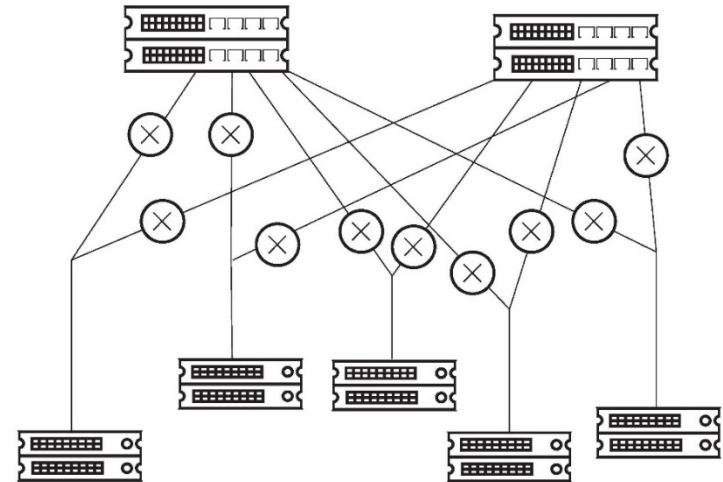


Two Problems Sharing Some Fundamentals

Cache-aided Communications



Communications in Distributed Computing

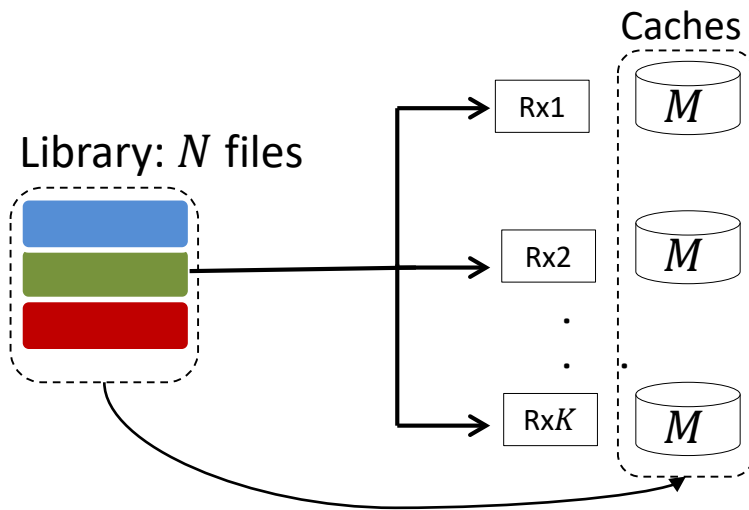


Common ingredients:

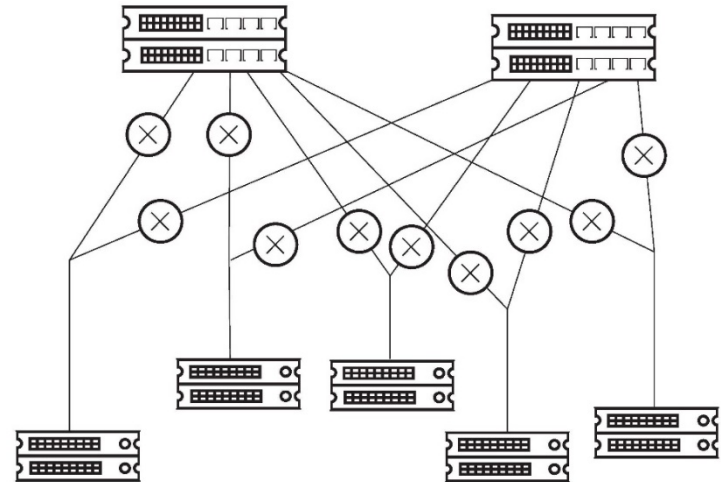
Communication cost and redundancy of stored/computed data

Two Problems Sharing Some Fundamentals

Cache-aided Communications



Communications in Distributed Computing



Common ingredients:

Communication cost and redundancy of stored/computed data

Common tool: Clique-based coding

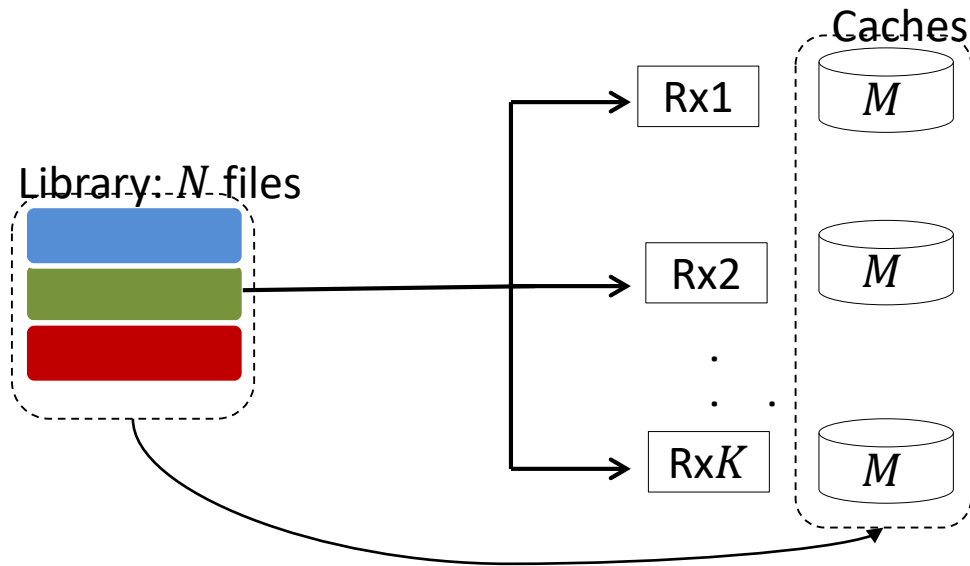
Common Bottlenecks/Challenges

Outline

- Basic elements of coded caching (first communications)
 - Basic properties
 - Main gains & Important variants
- Ramifications of coded caching in distributed computing
 - Coded map reduce
 - Tradeoff between computing and communicating
 - Main bottlenecks
- Main challenges
 - Subpacketization: The bottleneck and new algorithmic solutions
 - Non uniformity: The bottleneck and new solutions
- Exploiting multiple dimensions to alleviate bottlenecks
 - Multiplicative gains by reducing subpacketization
 - Reducing effect of non uniformity
 - New coding structures
- Open problems and closing remarks

Simple Caching

First: Cache-Aided Communications

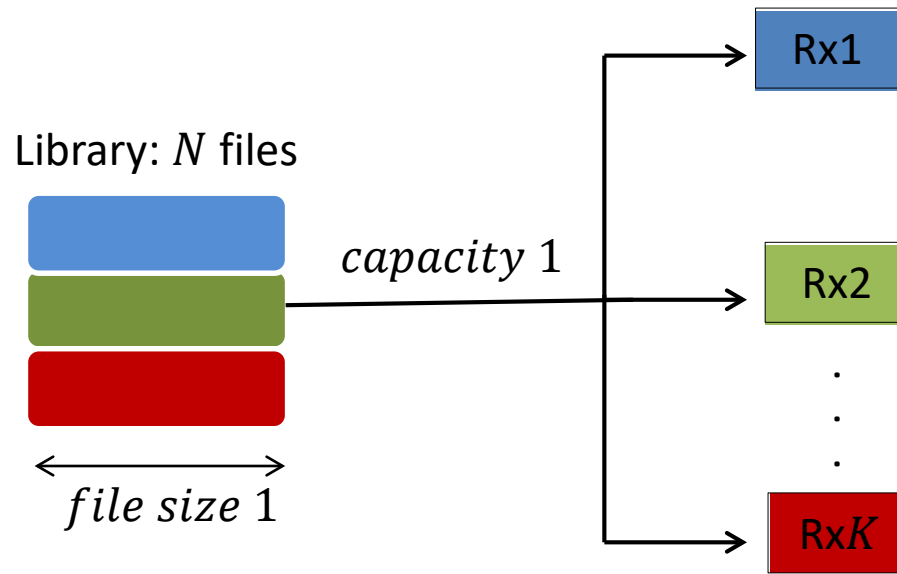


$$\gamma \stackrel{\text{def}}{=} \frac{M}{N} \stackrel{\text{def}}{=} \frac{\text{individual cache size}}{\text{library size}}$$

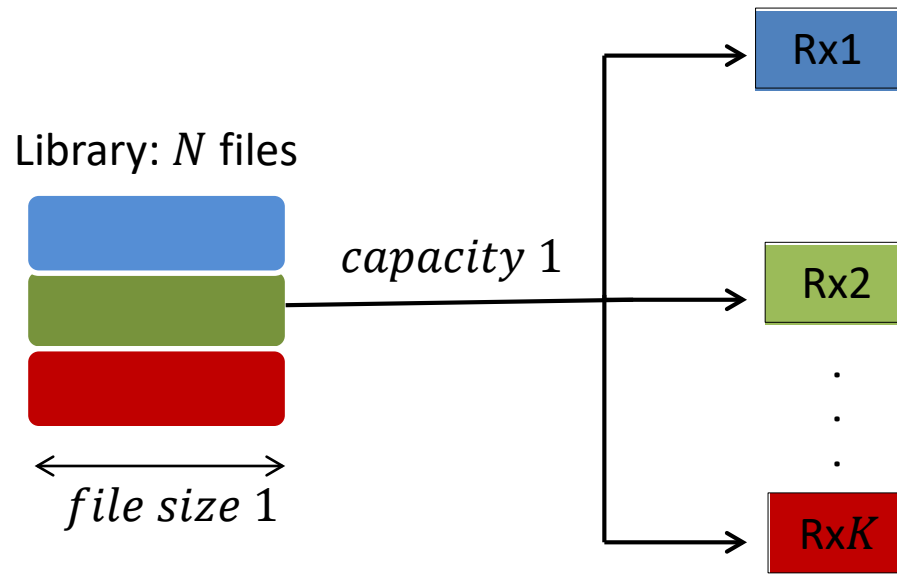
$T(\gamma)$: duration of delivery phase

OBJECTIVE: reduce $T(\gamma)$

Single-stream Channel: No Caching ($M = 0$)

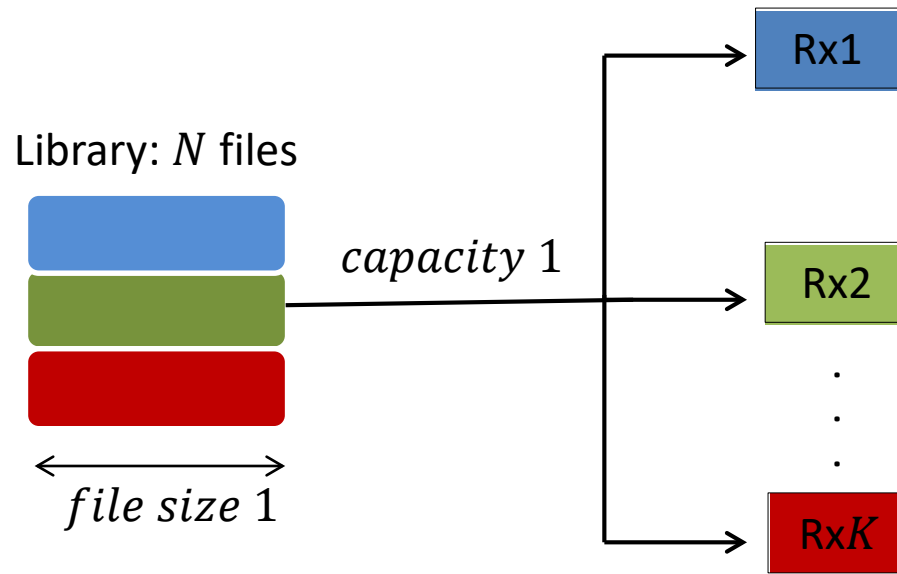


Single-stream Channel: No Caching ($M = 0$)



- Transmission sequence:

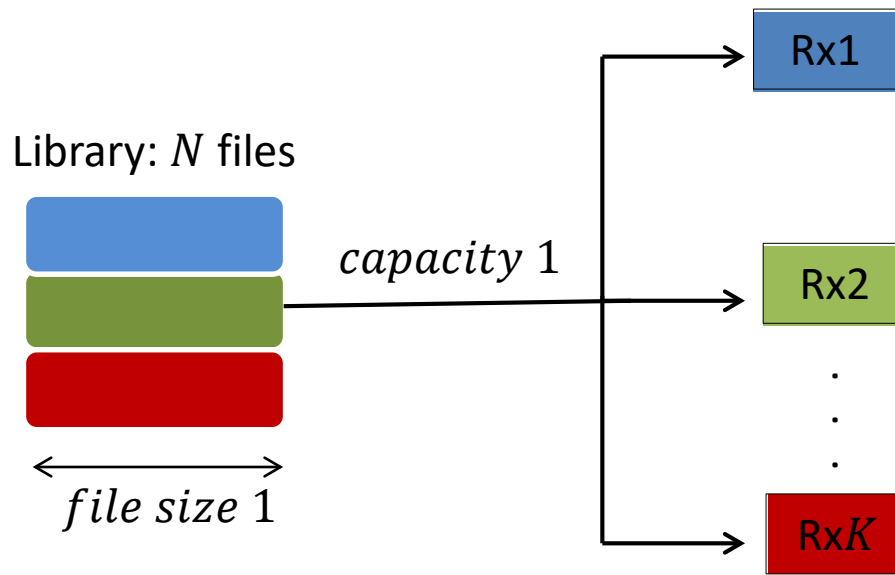
Single-stream Channel: No Caching ($M = 0$)



- Transmission sequence:

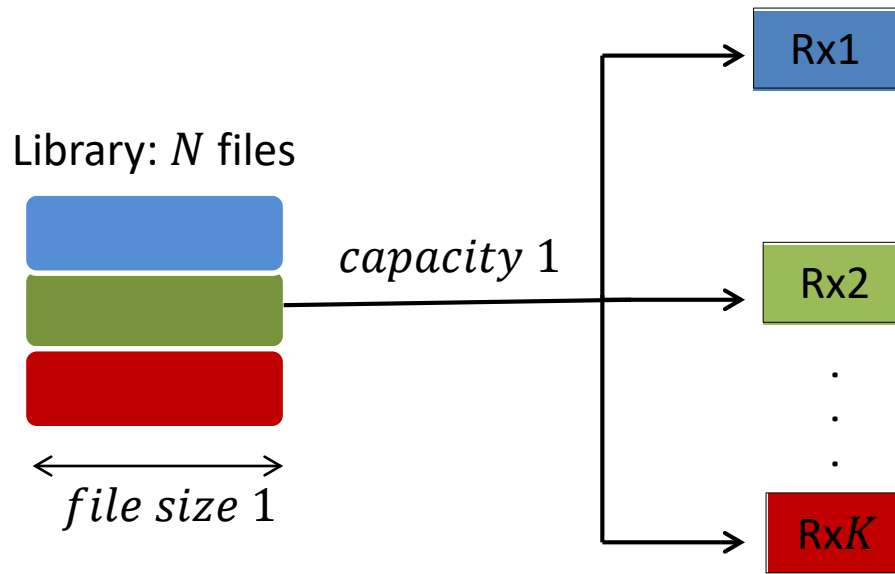


Single-stream Channel: No Caching ($M = 0$)



- Transmission sequence:  

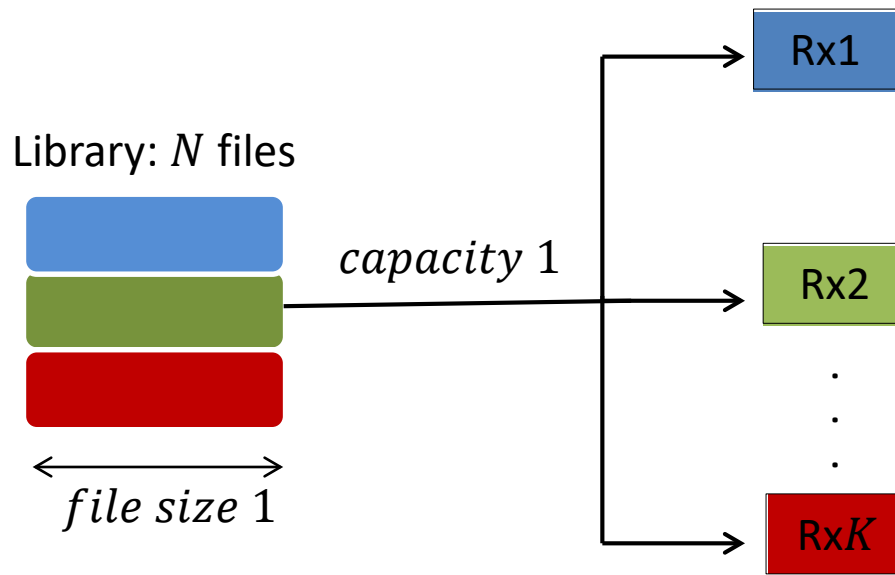
Single-stream Channel: No Caching ($M = 0$)



- Transmission sequence:



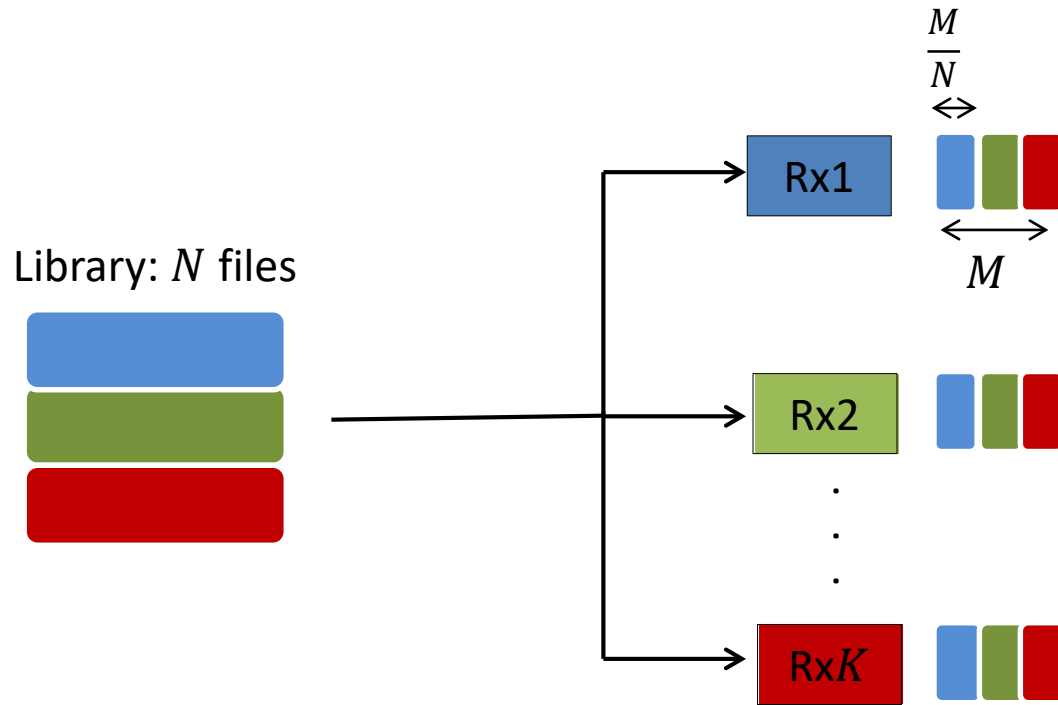
Single-stream Channel: No Caching ($M = 0$)



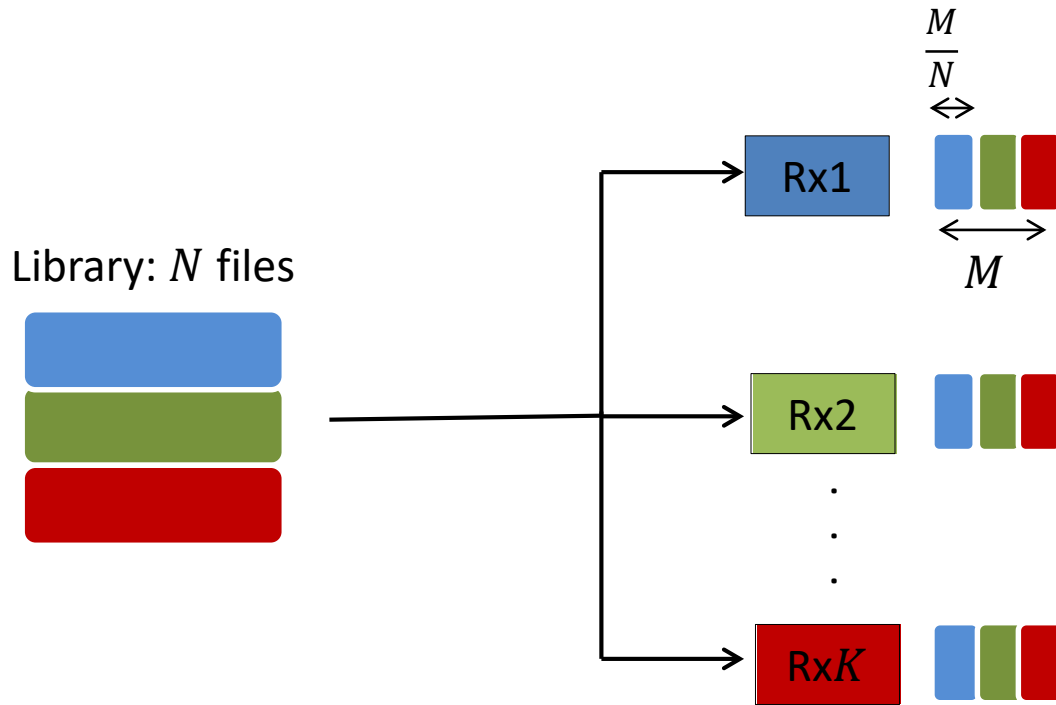
- Transmission sequence:   

$$T = K$$

Traditional Caching (Worst-Case Consideration)

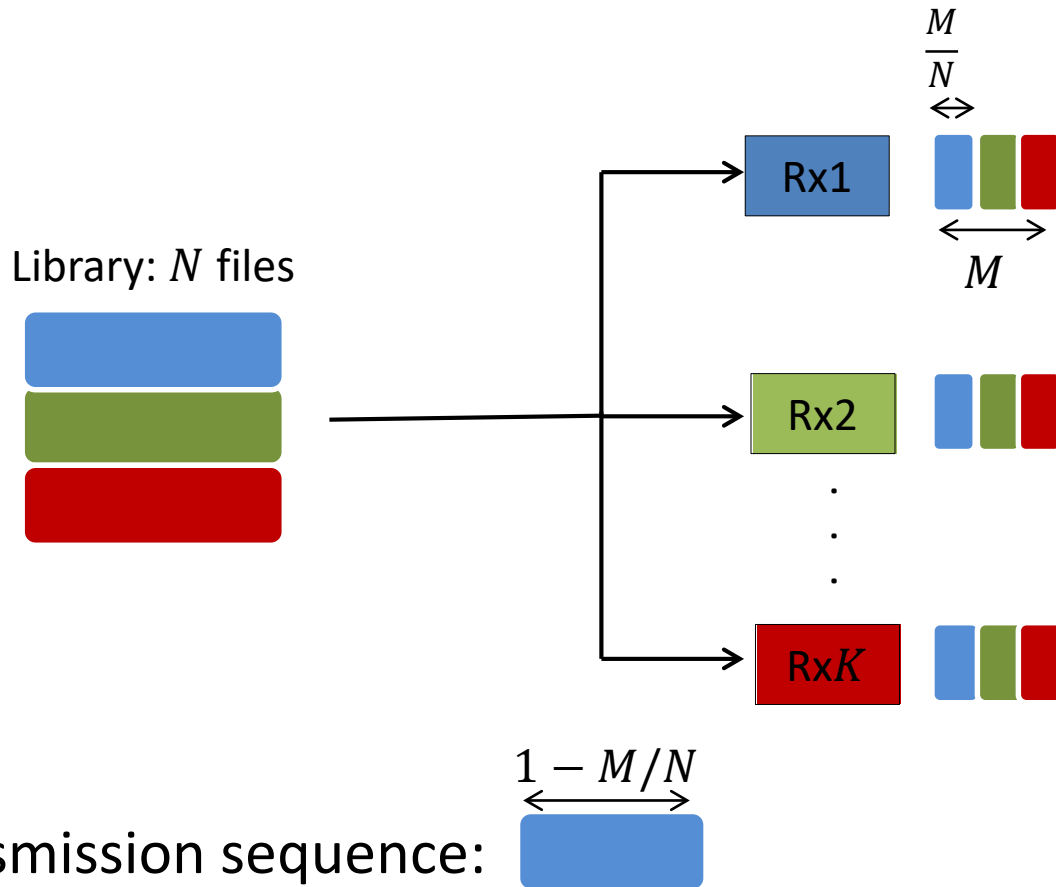


Traditional Caching (Worst-Case Consideration)

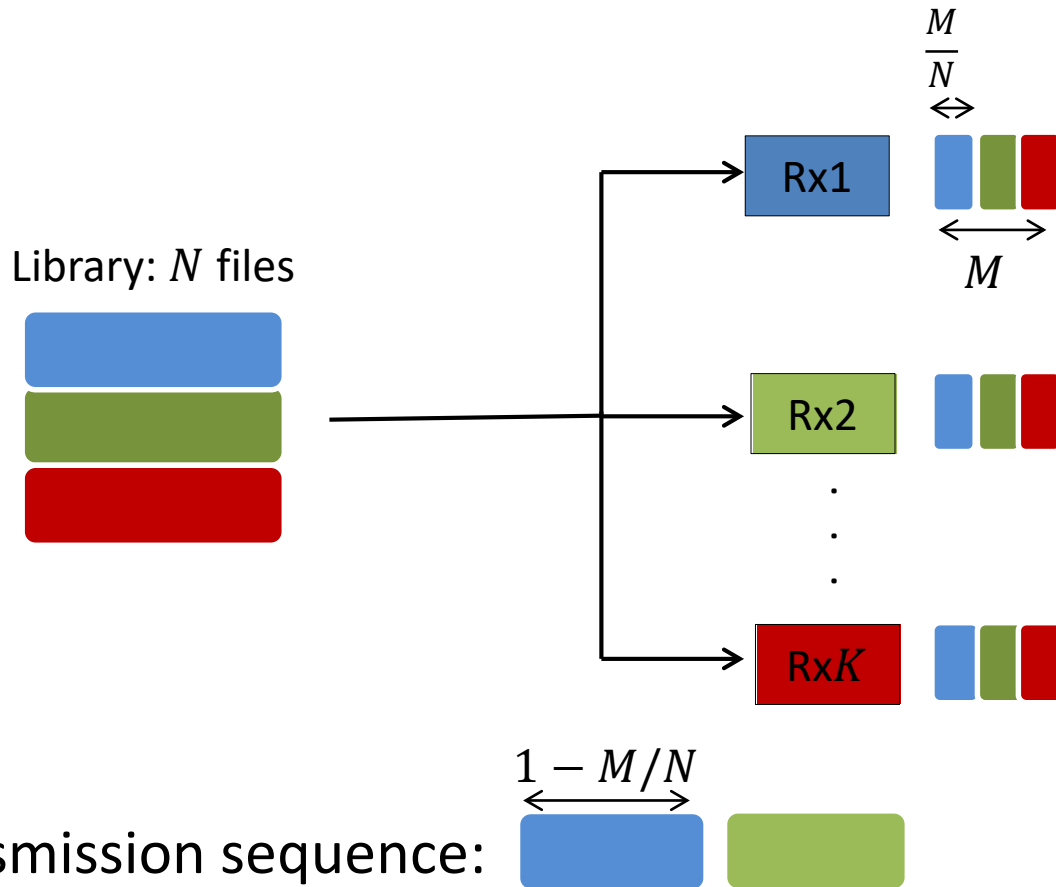


- Transmission sequence:

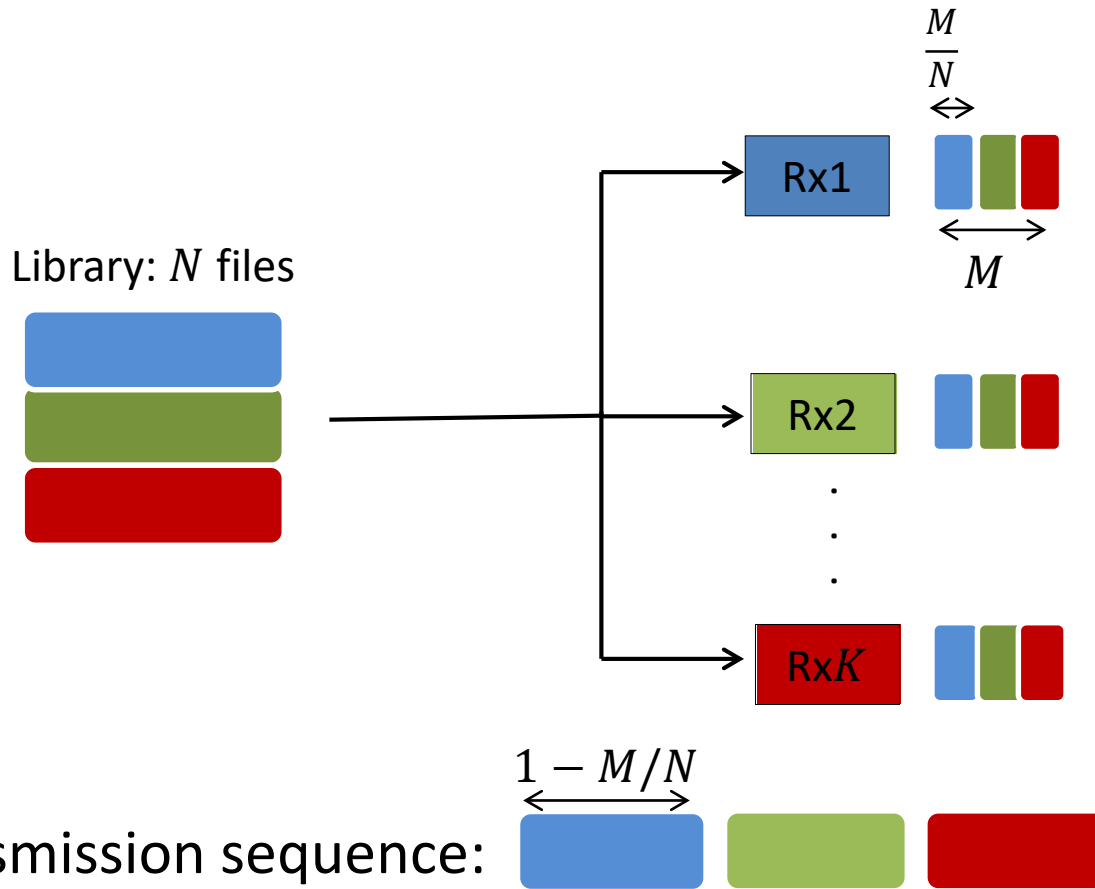
Traditional Caching (Worst-Case Consideration)



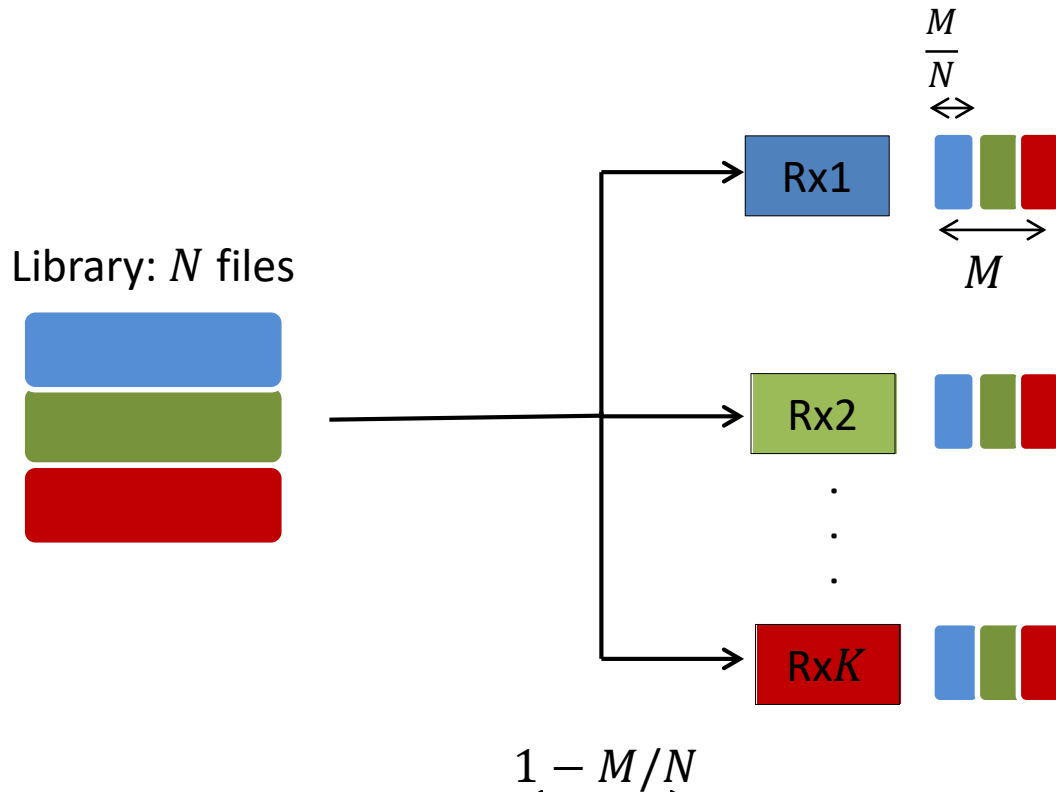
Traditional Caching (Worst-Case Consideration)




Traditional Caching (Worst-Case Consideration)



Traditional Caching (Worst-Case Consideration)



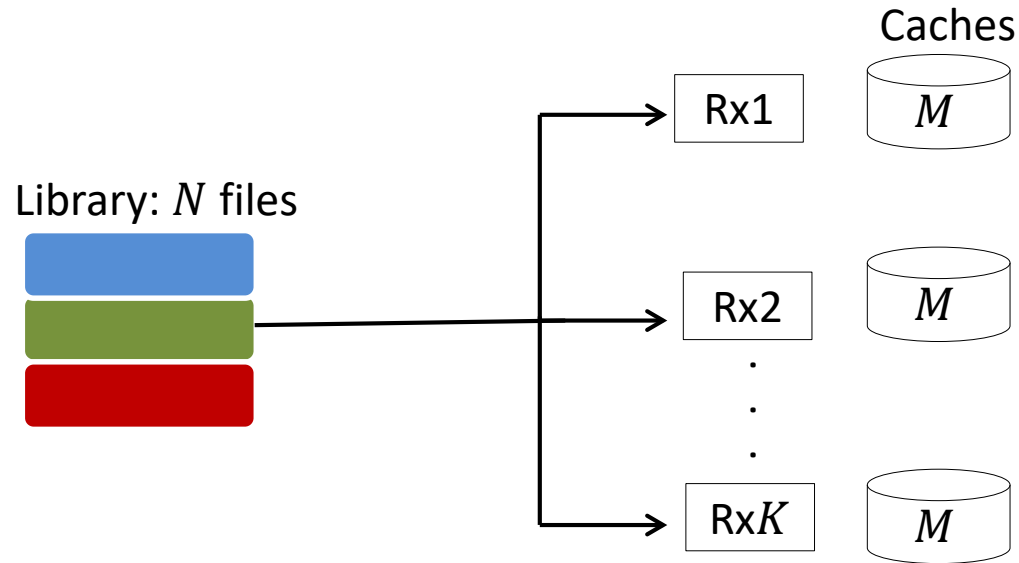
- Transmission sequence: 
- Local cache gain: $(1 - M/N)$ for each user
- The rate:

$$T = K(1 - M/N) = K(1 - \gamma),$$

$$\gamma \stackrel{\text{def}}{=} \frac{M}{N}$$

The power of advanced caching

Coded Caching



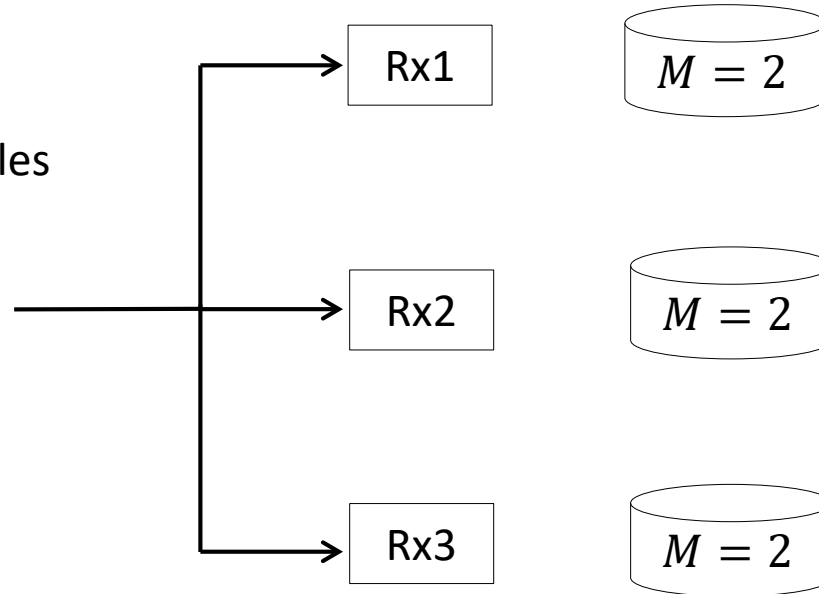
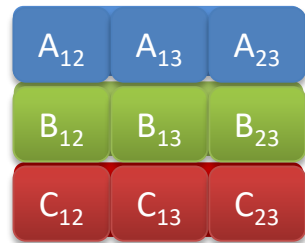
Key breakthrough: USE CACHES TO CANCEL INTERFERENCE

- Cache so that one transmission is useful to many
 - Even if requested files are different
- Large decreases in delay

Result: Maddah-Ali, Niesen (2013)

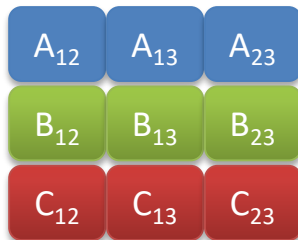
Example: $N = K = 3, M = 2$ ($\gamma = \frac{2}{3}$)

Library: $N = 3$ files

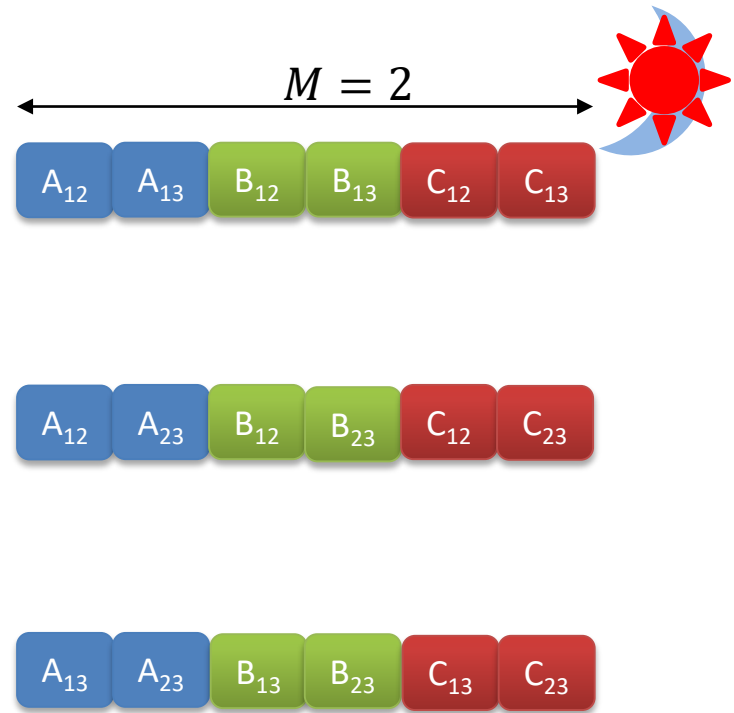
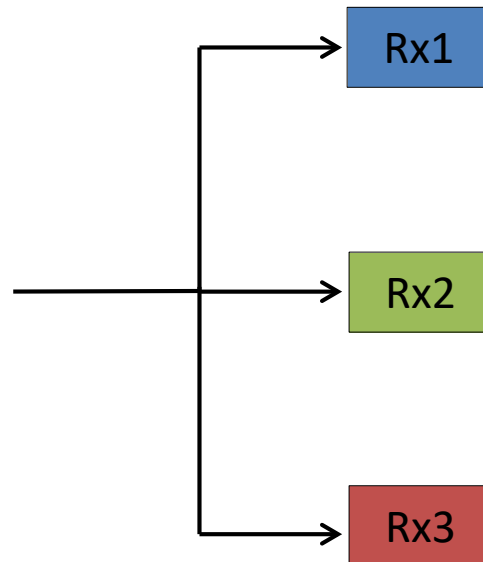


Example: $N = K = 3, M = 2$ $(\gamma = \frac{M}{N} = \frac{2}{3})$

Library: N files

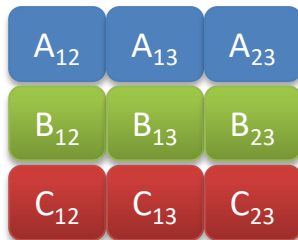


\longleftrightarrow
 $\frac{1}{3}$

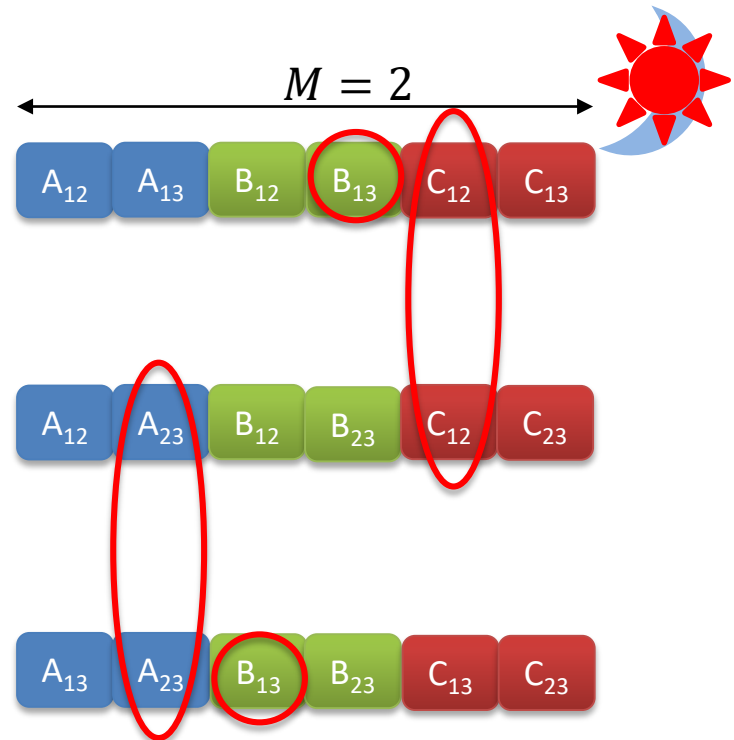
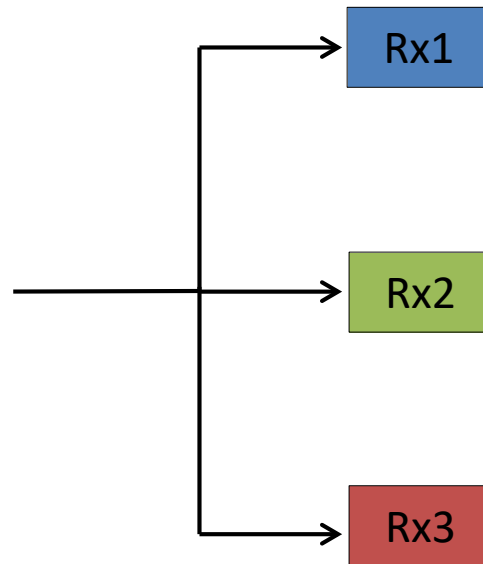


Example: $N = K = 3, M = 2$ $(\gamma = \frac{M}{N} = \frac{2}{3})$

Library: N files

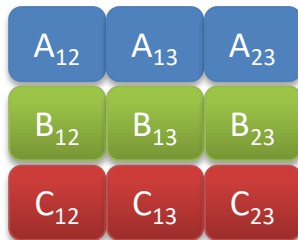


\longleftrightarrow
 $\frac{1}{3}$

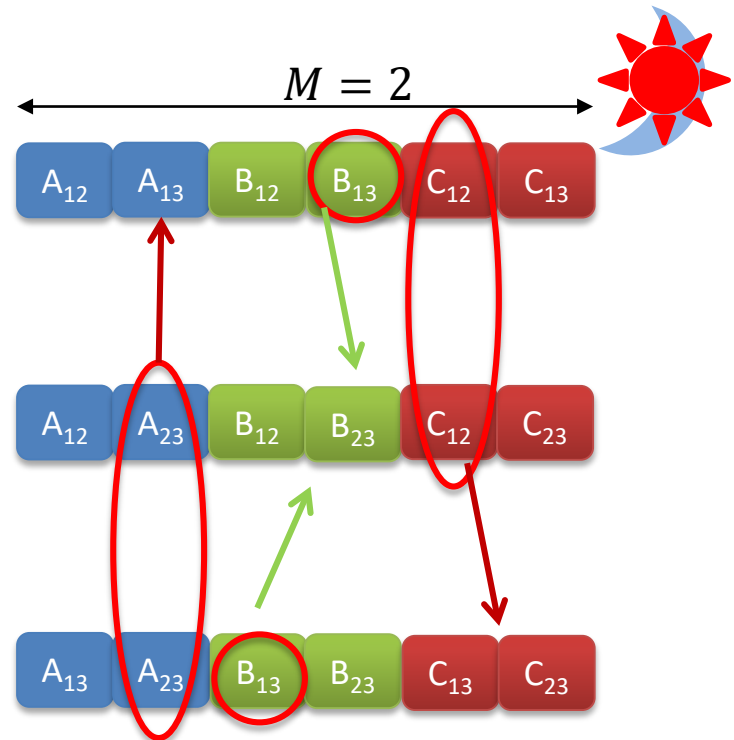
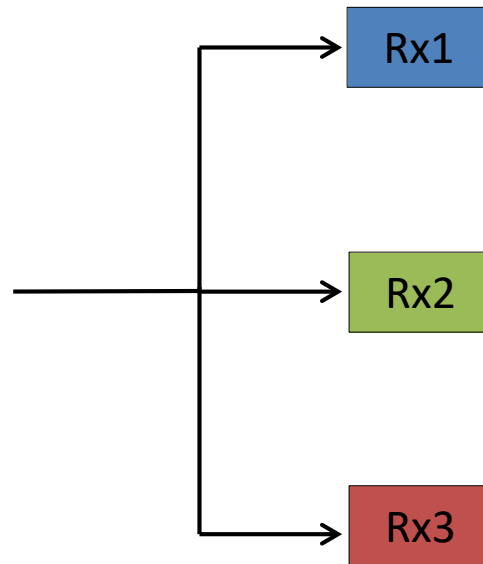


Example: $N = K = 3, M = 2$ $(\gamma = \frac{M}{N} = \frac{2}{3})$

Library: N files

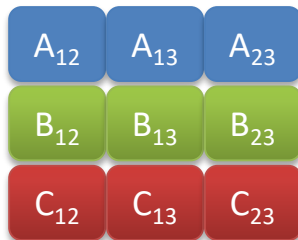


\longleftrightarrow
 $\frac{1}{3}$

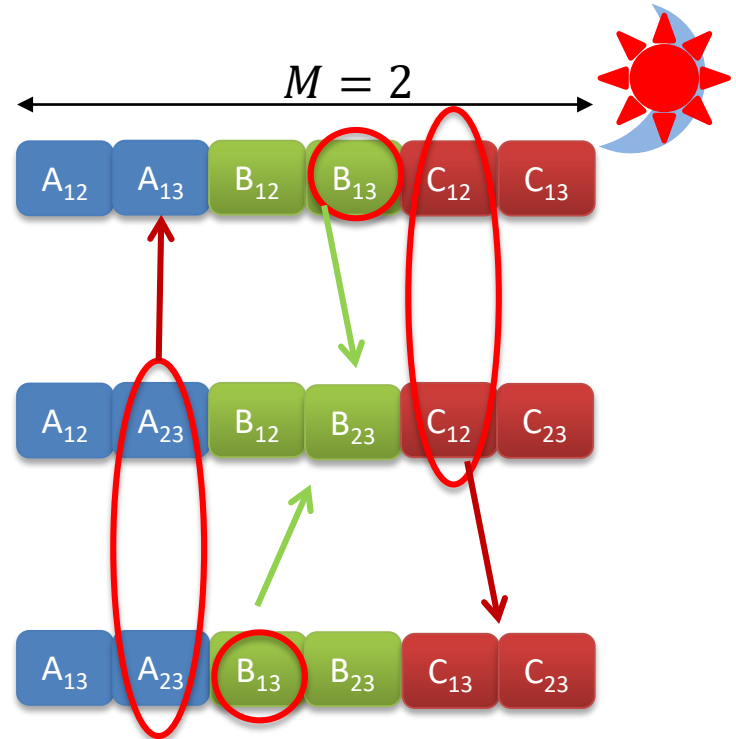
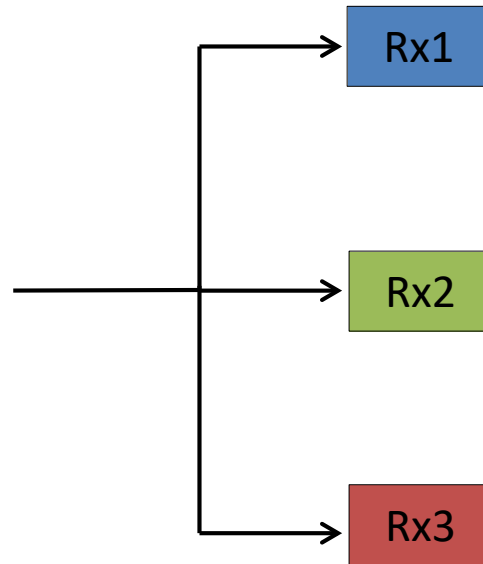


Example: $N = K = 3, M = 2$ $(\gamma = \frac{M}{N} = \frac{2}{3})$

Library: N files

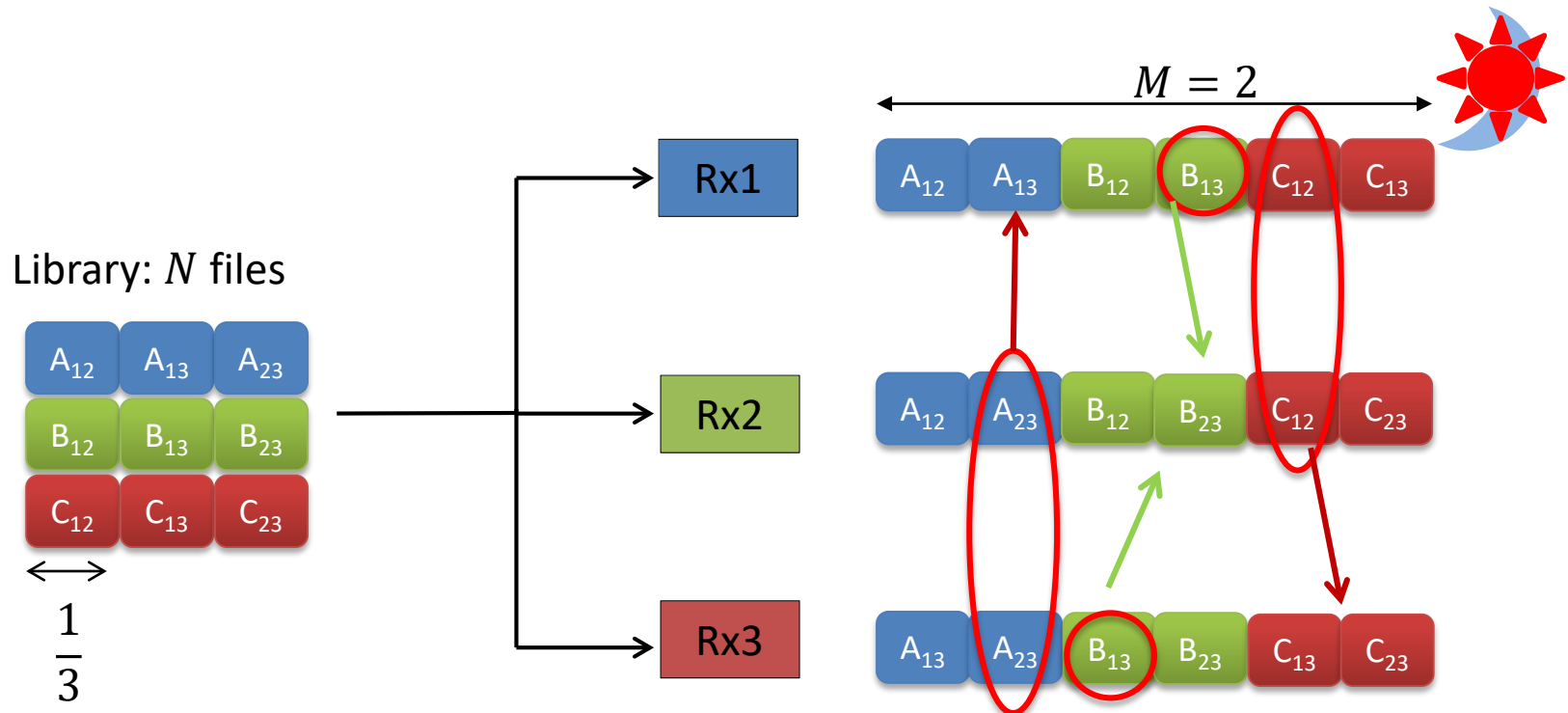


\longleftrightarrow
 $\frac{1}{3}$



- Transmit : $A_{23} \oplus B_{13} \oplus C_{12}$ (a common message for all)

Example: $N = K = 3, M = 2$ $\left(\gamma = \frac{M}{N} = \frac{2}{3}\right)$



- Transmit : $A_{23} \oplus B_{13} \oplus C_{12}$ (a common message for all)

$$Gain = 3 = \frac{KM}{N} \text{ users at a time}$$

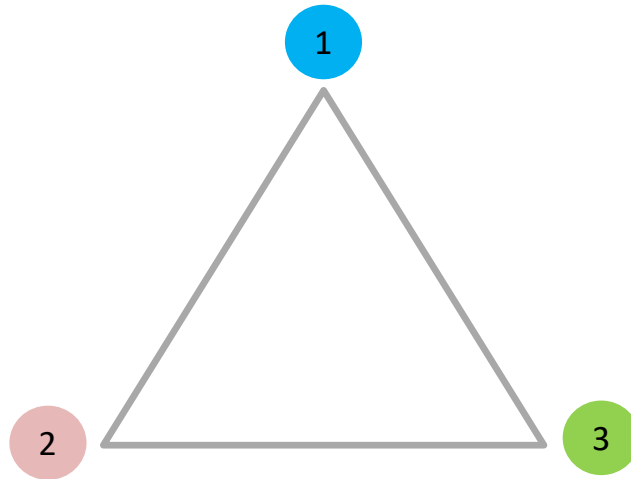
Coded Caching: Intuition - Clique

- *Deliver to $K\gamma + 1$ users at a time*
- *Via XORs with $K\gamma + 1$ subfiles.*
 - ***Each user (out of the $K\gamma + 1$ now served) knows all summands except its own***

$$T = \frac{K(1 - \gamma)}{1 + K\gamma} \approx \frac{1 - \gamma}{\gamma}$$

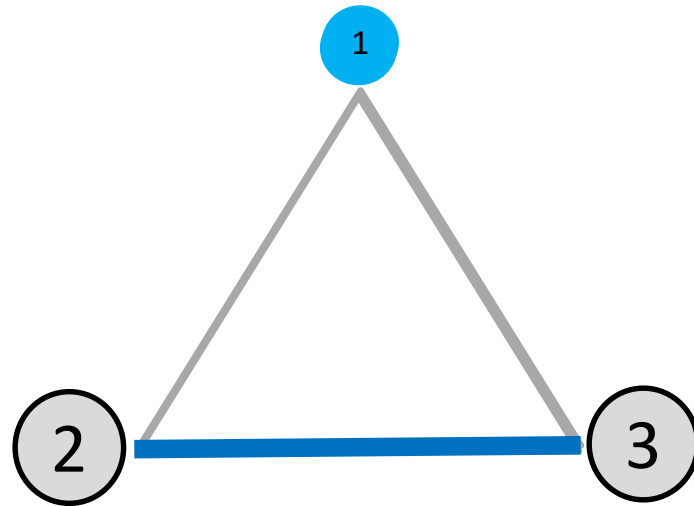
Coded Caching: Intuition - Clique

$$K = 3, K\gamma = 2$$



Coded Caching: Intuition - Clique

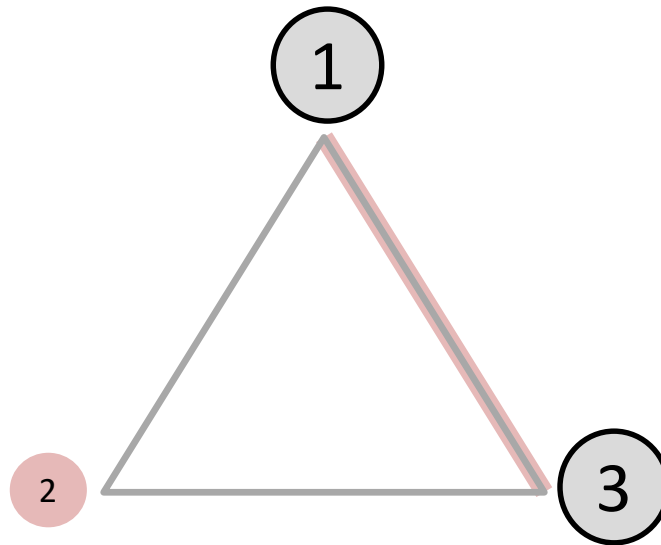
$$K = 3, K\gamma = 2$$



A_{23}

Coded Caching: Intuition - Clique

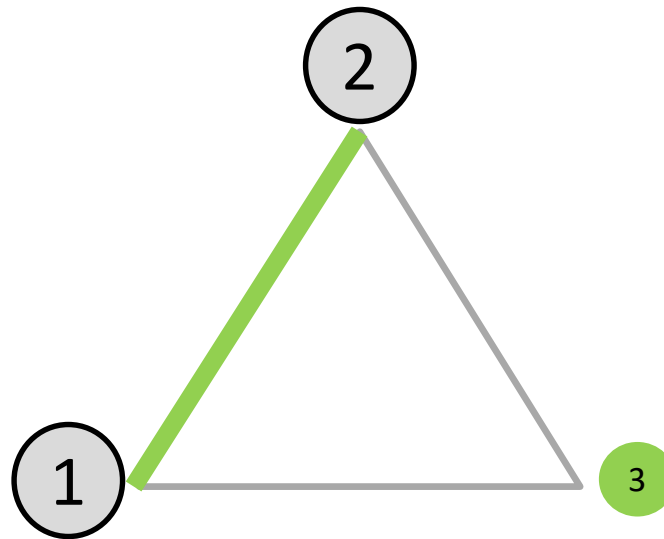
$$K = 3, K\gamma = 2$$



$$\boxed{A_{23}} \oplus \boxed{B_{13}}$$

Coded Caching: Intuition - Clique

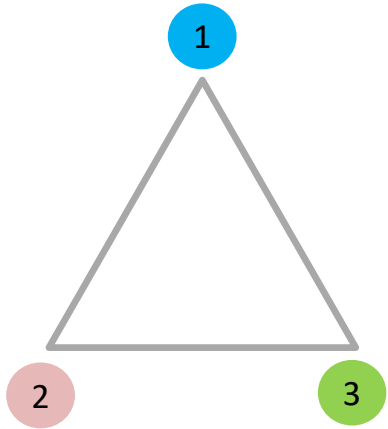
$$K = 3, K\gamma = 2$$



$$A_{23} \oplus B_{13} \oplus C_{12}$$

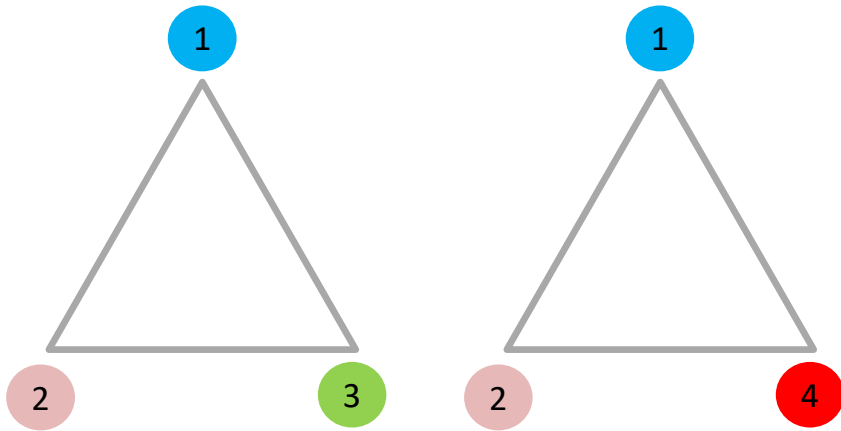
Coded Caching: Intuition - Cliques

$$K = 4, K\gamma = 2$$



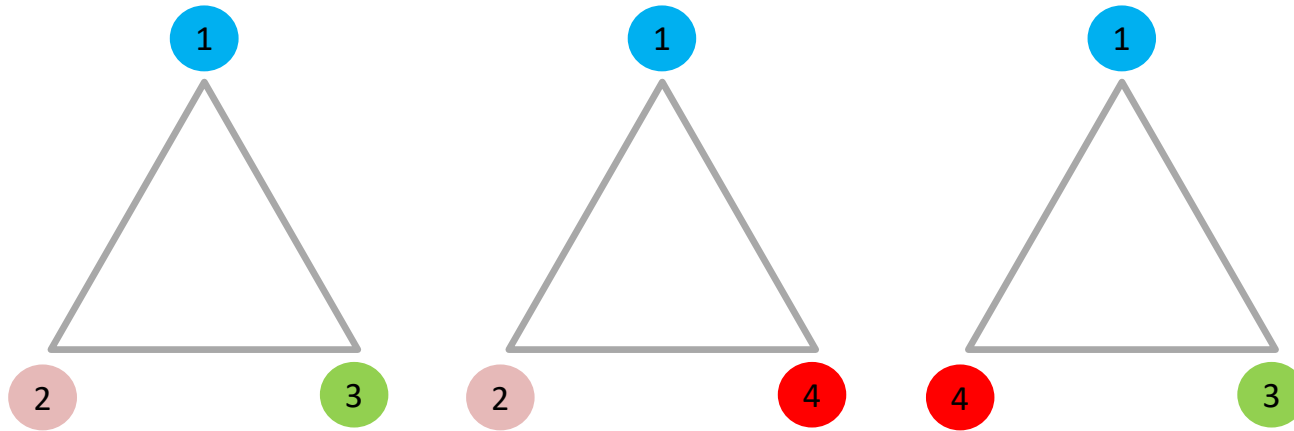
Coded Caching: Intuition - Cliques

$$K = 4, K\gamma = 2$$



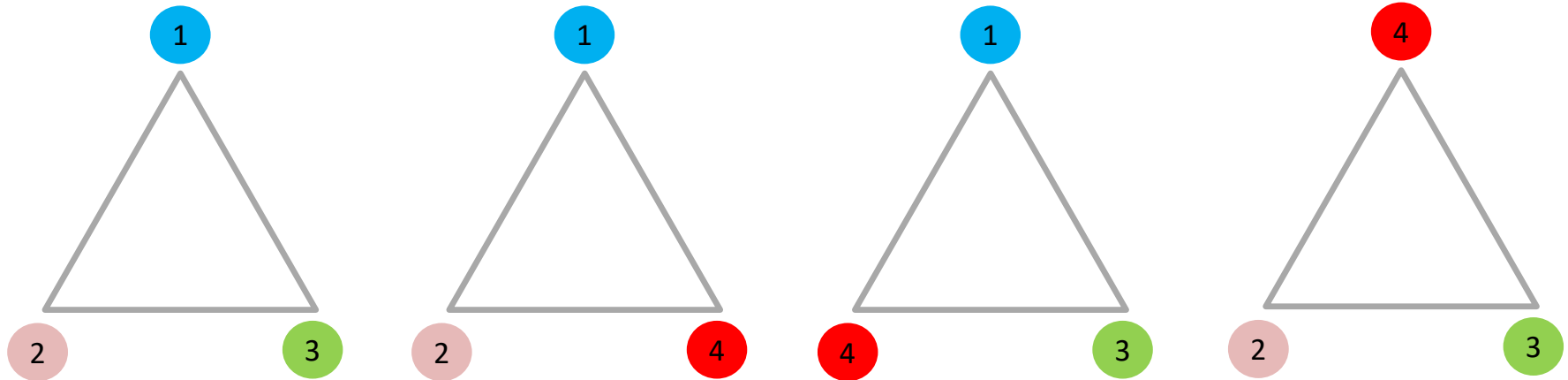
Coded Caching: Intuition - Cliques

$$K = 4, K\gamma = 2$$



Coded Caching: Intuition - Cliques

$$K = 4, K\gamma = 2$$

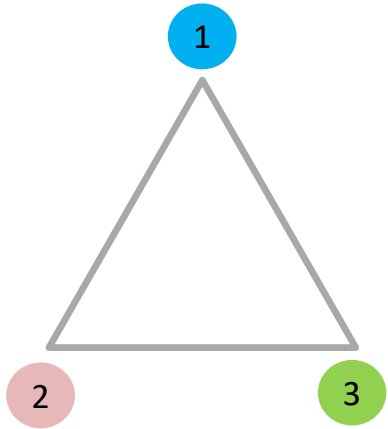


Coded Caching: Intuition - Cliques

$$K = 4, K\gamma = 2$$

Coded Caching: Intuition - Cliques

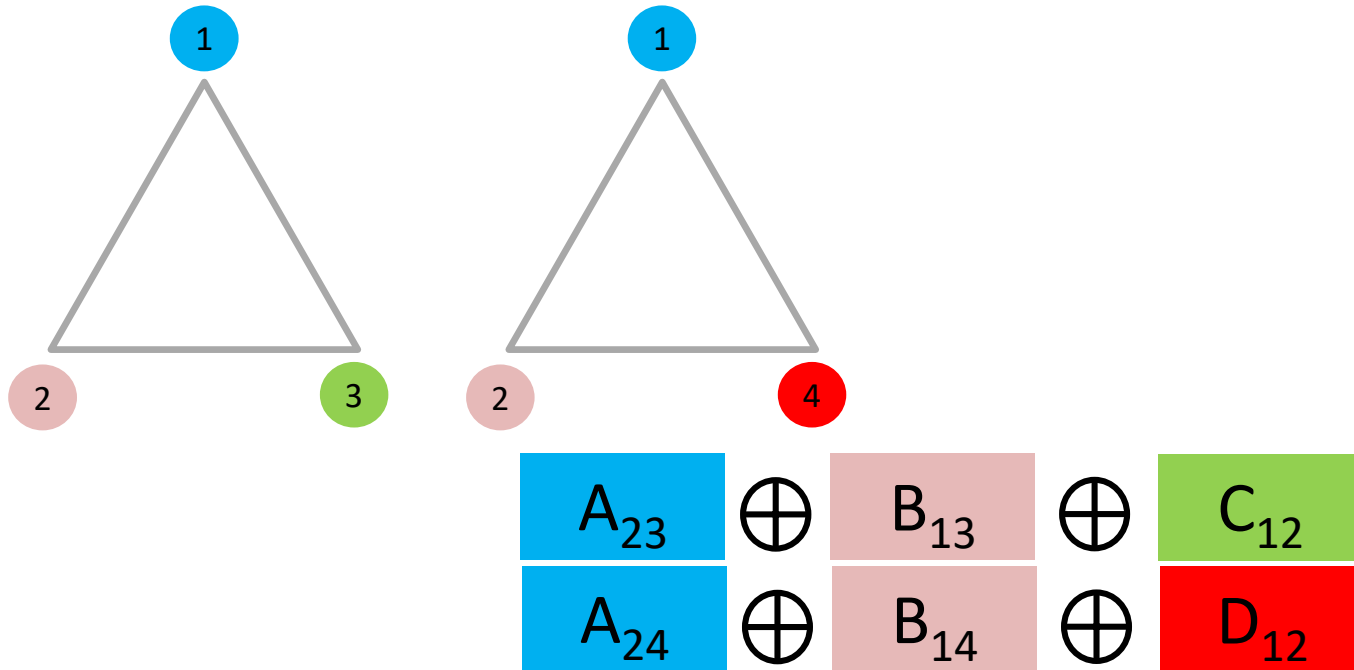
$$K = 4, K\gamma = 2$$



$$A_{23} \oplus B_{13} \oplus C_{12}$$

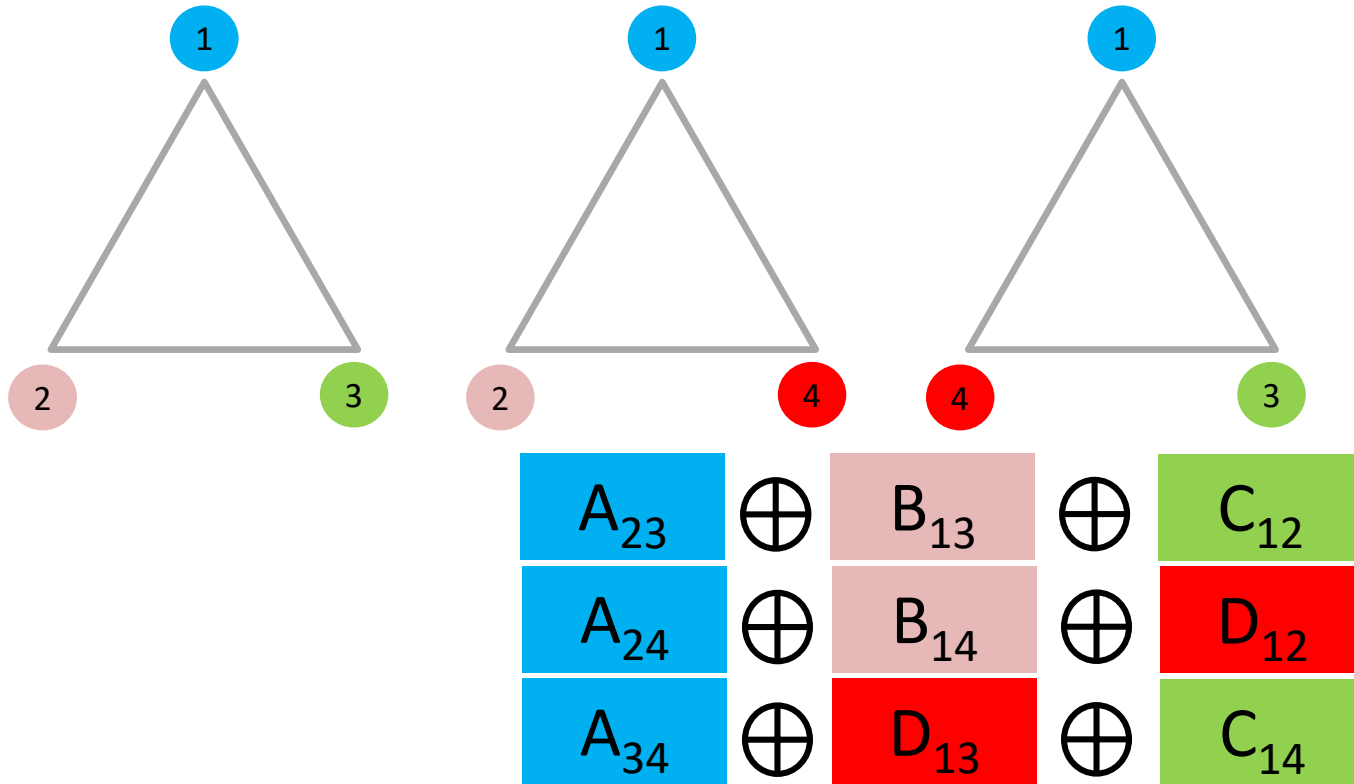
Coded Caching: Intuition - Cliques

$$K = 4, K\gamma = 2$$



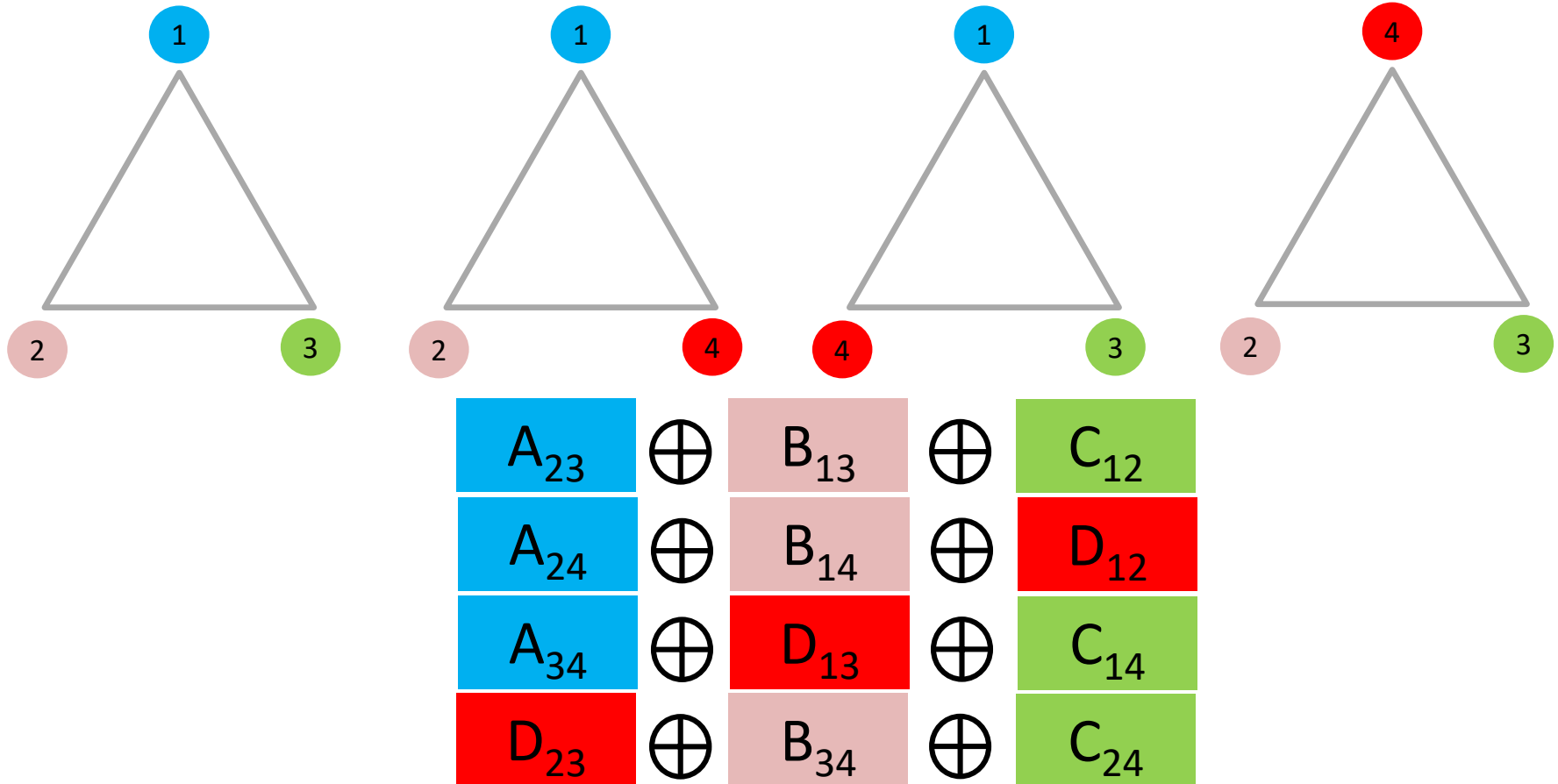
Coded Caching: Intuition - Cliques

$$K = 4, K\gamma = 2$$



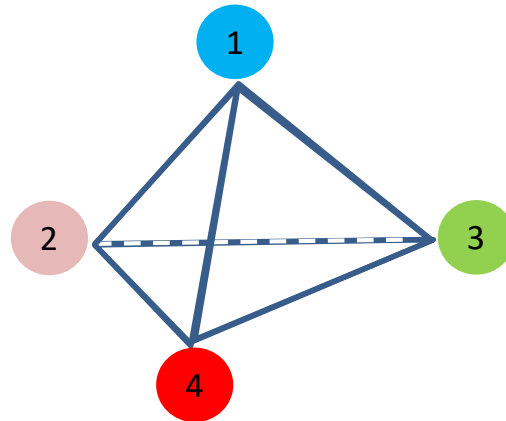
Coded Caching: Intuition - Cliques

$$K = 4, K\gamma = 2$$



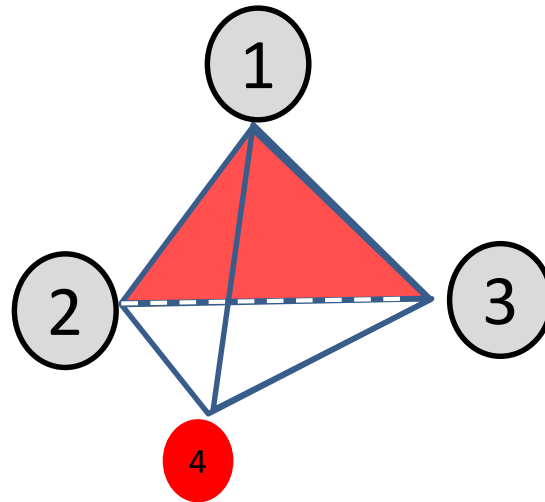
Coded Caching: Intuition - Clique

$$K = 4, K\gamma = 3$$



Coded Caching: Intuition - Clique

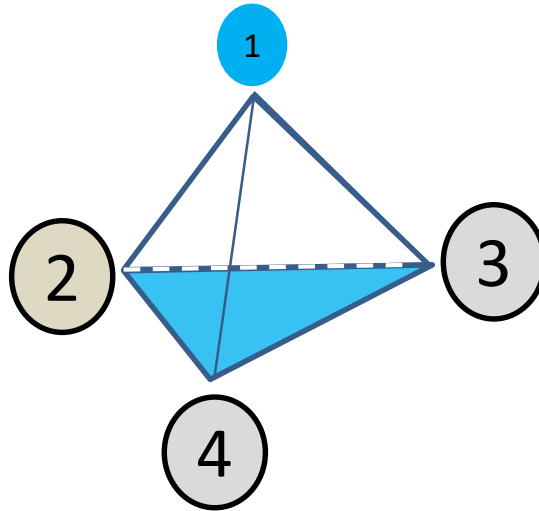
$$K = 4, K\gamma = 3$$



D_{123}

Coded Caching: Intuition - Clique

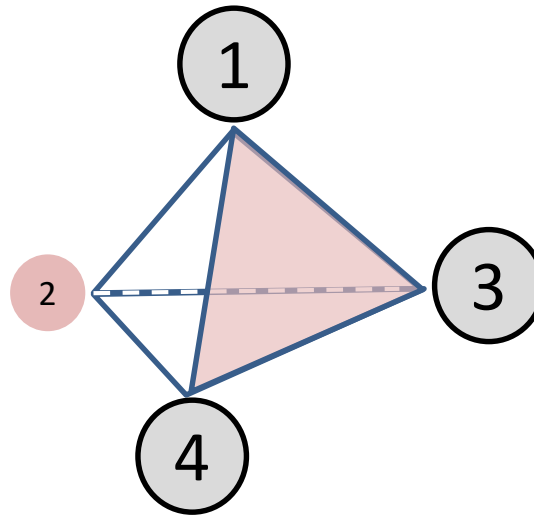
$$K = 4, K\gamma = 3$$



$$\boxed{D_{123}} \oplus \boxed{A_{234}}$$

Coded Caching: Intuition - Clique

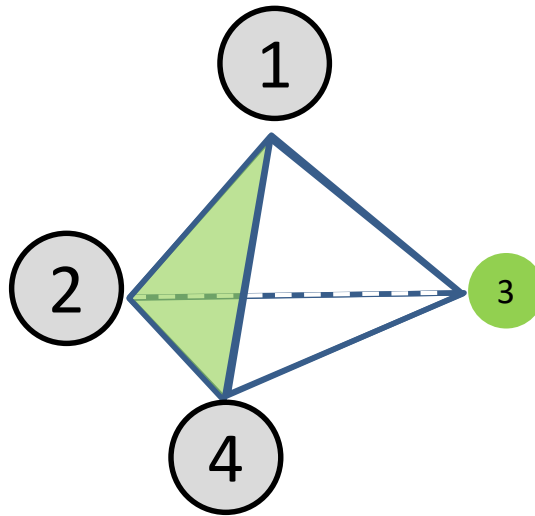
$$K = 4, K\gamma = 3$$



$$\boxed{D_{123}} \oplus \boxed{A_{234}} \oplus \boxed{B_{134}}$$

Coded Caching: Intuition - Clique

$$K = 4, K\gamma = 3$$



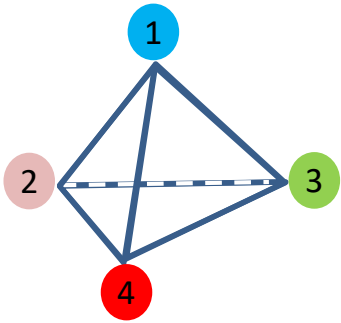
$$\boxed{D_{123}} \oplus \boxed{A_{234}} \oplus \boxed{B_{134}} \oplus \boxed{C_{124}}$$

Coded Caching: Intuition - Cliques

$$K = 5, K\gamma = 3$$

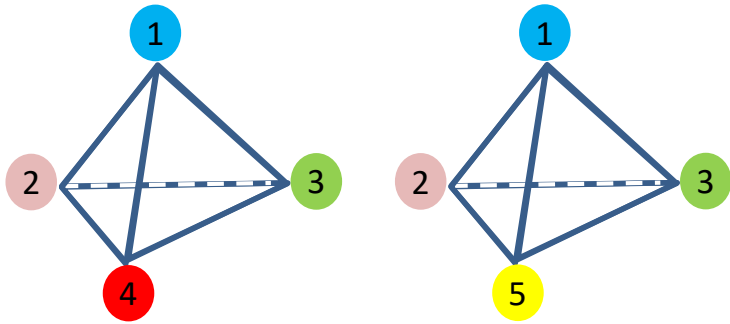
Coded Caching: Intuition - Cliques

$$K = 5, K\gamma = 3$$



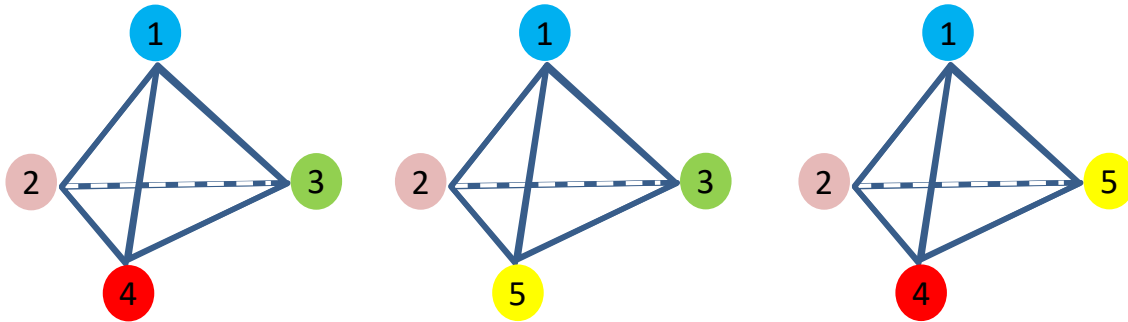
Coded Caching: Intuition - Cliques

$$K = 5, K\gamma = 3$$



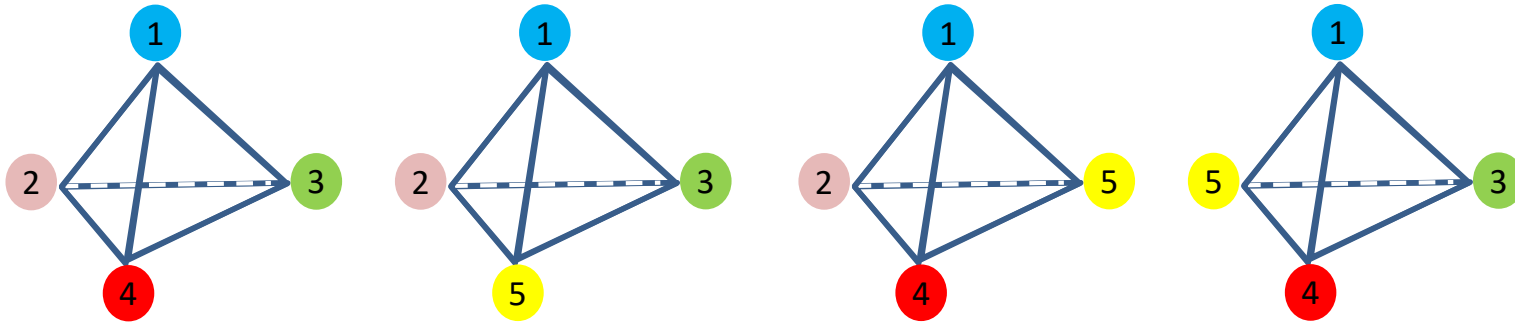
Coded Caching: Intuition - Cliques

$$K = 5, K\gamma = 3$$



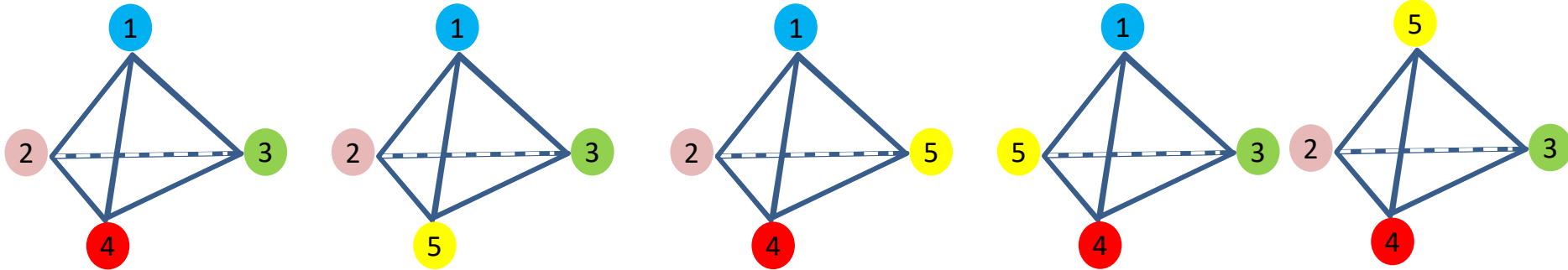
Coded Caching: Intuition - Cliques

$$K = 5, K\gamma = 3$$



Coded Caching: Intuition - Cliques

$$K = 5, K\gamma = 3$$

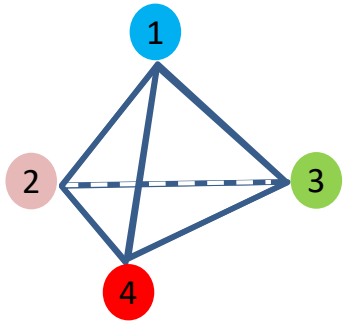


Coded Caching: Intuition - Cliques

$$K = 5, K\gamma = 3$$

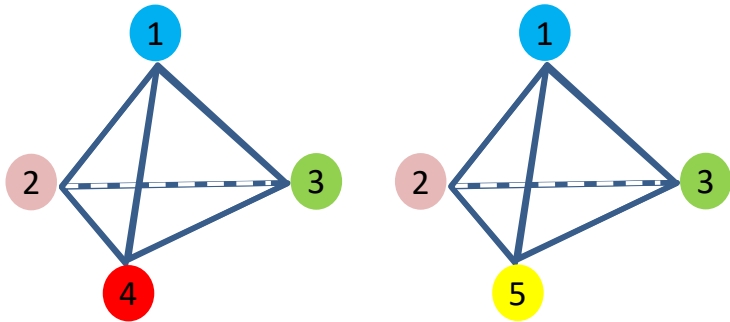
Coded Caching: Intuition - Cliques

$$K = 5, K\gamma = 3$$



$$D_{123} \oplus A_{234} \oplus B_{134} \oplus C_{124}$$

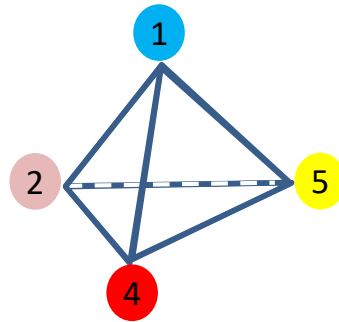
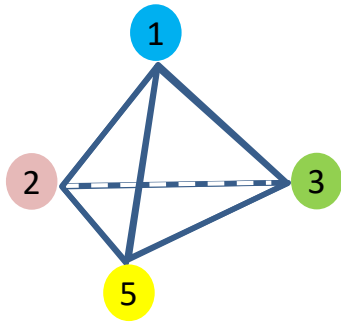
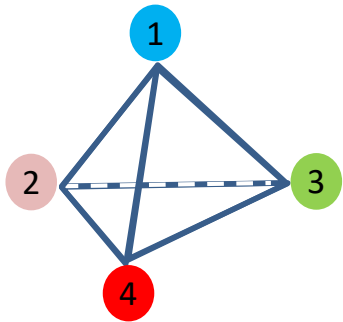
Coded Caching: Intuition - Cliques



$$\begin{array}{ccccccc} D_{123} & \oplus & A_{234} & \oplus & B_{134} & \oplus & C_{124} \\ E_{123} & \oplus & A_{235} & \oplus & B_{135} & \oplus & C_{125} \end{array}$$

Coded Caching: Intuition - Cliques

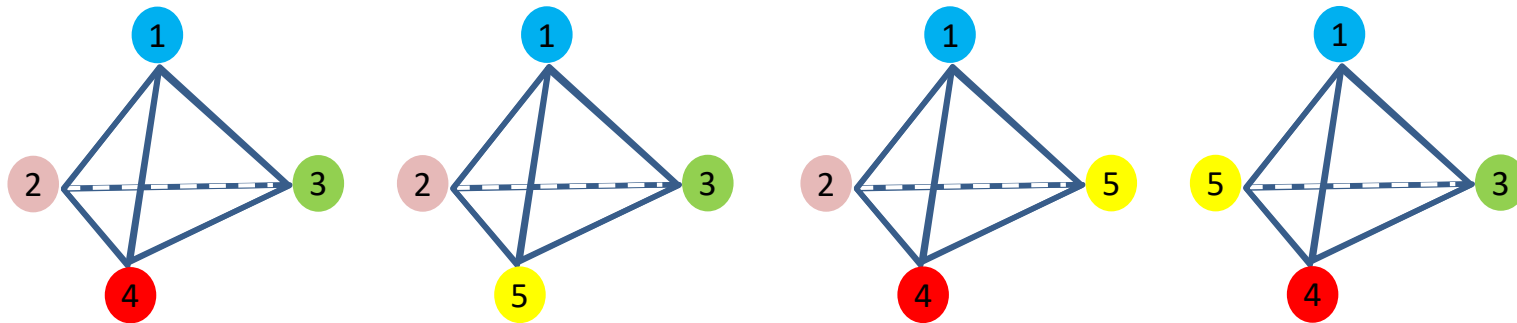
$$K = 5, K\gamma = 3$$



D_{123}	\oplus	A_{234}	\oplus	B_{134}	\oplus	C_{124}
E_{123}	\oplus	A_{235}	\oplus	B_{135}	\oplus	C_{125}
E_{124}	\oplus	A_{245}	\oplus	D_{125}	\oplus	B_{145}

Coded Caching: Intuition - Cliques

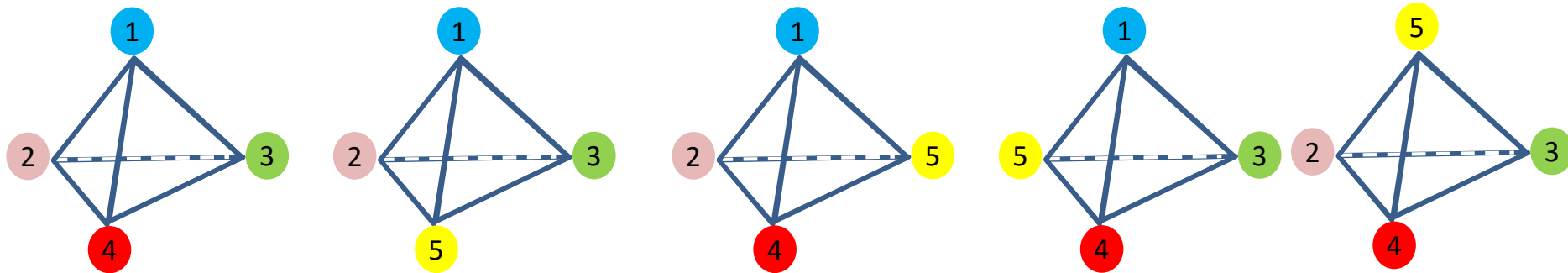
$$K = 5, K\gamma = 3$$



D_{123}	\oplus	A_{234}	\oplus	B_{134}	\oplus	C_{124}
E_{123}	\oplus	A_{235}	\oplus	B_{135}	\oplus	C_{125}
E_{124}	\oplus	A_{245}	\oplus	D_{125}	\oplus	B_{145}
E_{134}	\oplus	A_{345}	\oplus	D_{135}	\oplus	C_{145}

Coded Caching: Intuition - Cliques

$$K = 5, K\gamma = 3$$



D_{123}	\oplus	A_{234}	\oplus	B_{134}	\oplus	C_{124}
E_{123}	\oplus	A_{235}	\oplus	B_{135}	\oplus	C_{125}
E_{124}	\oplus	A_{245}	\oplus	D_{125}	\oplus	B_{145}
E_{134}	\oplus	A_{345}	\oplus	D_{135}	\oplus	C_{145}
E_{234}	\oplus	B_{345}	\oplus	D_{235}	\oplus	C_{245}

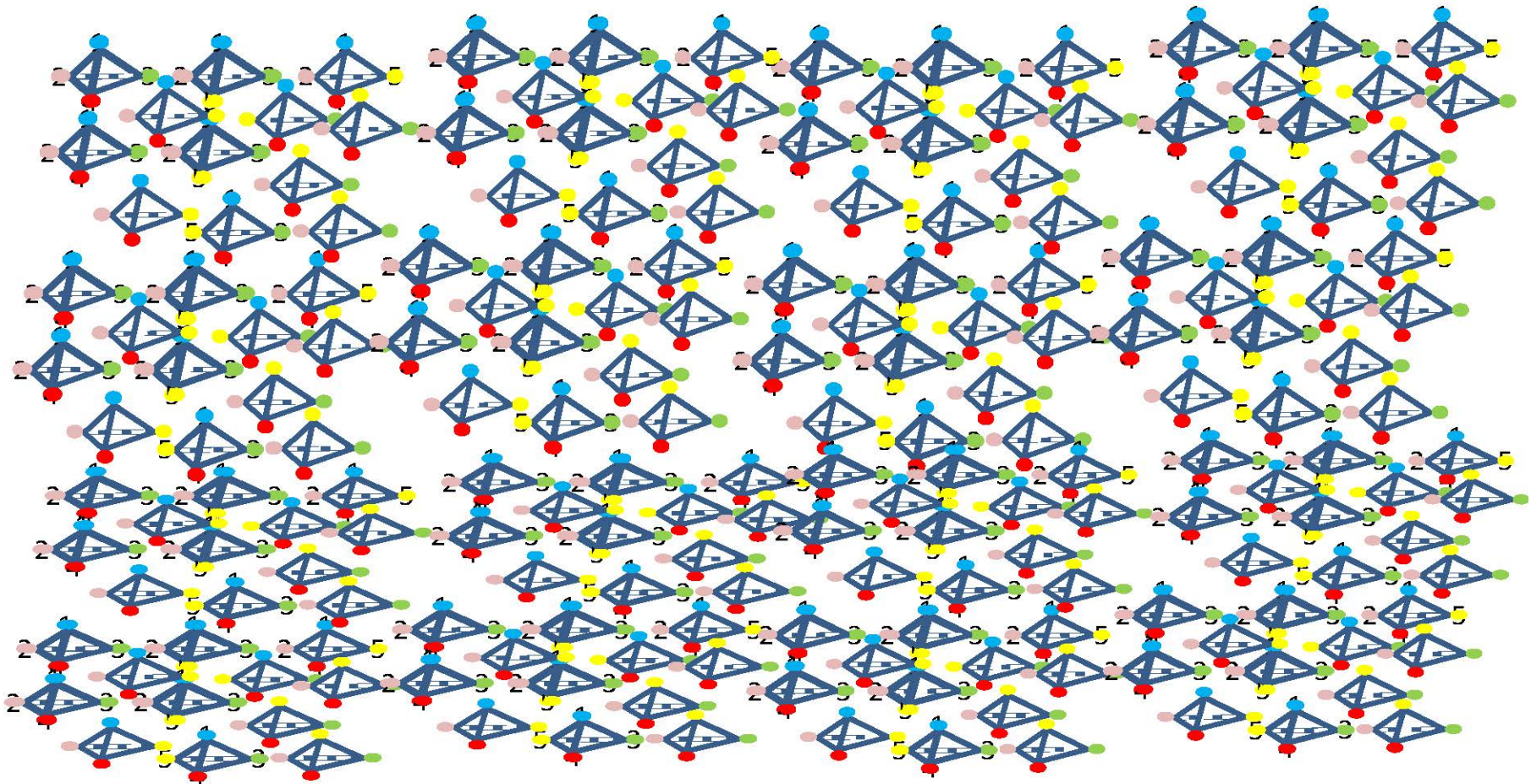
Coded Caching: Intuition – Problematically Many Cliques

$$K = 100, \quad K\gamma = 9$$

Coded Caching: Intuition – Problematically Many Cliques

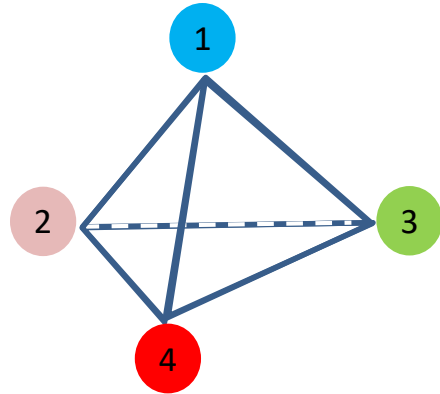
$$K = 100, \quad K\gamma = 9$$

9 users at a time, $2 \cdot 10^{13}$ cliques



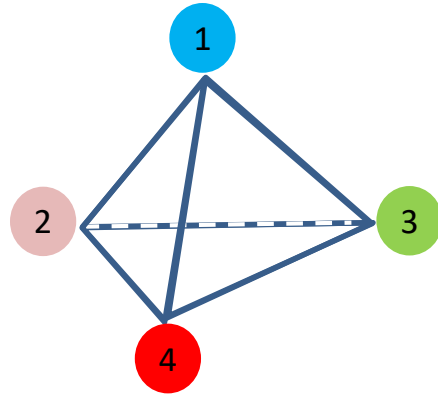
Extension of Clique-based Idea

Extension of Clique-based Idea

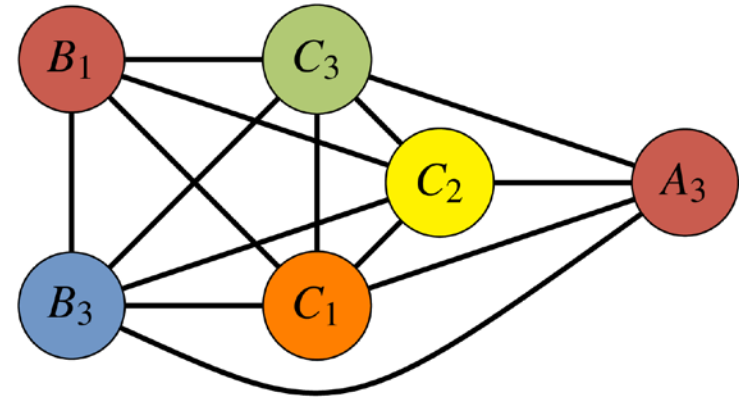


$$D_{123} \oplus A_{234} \oplus B_{134} \oplus C_{124}$$

Extension of Clique-based Idea

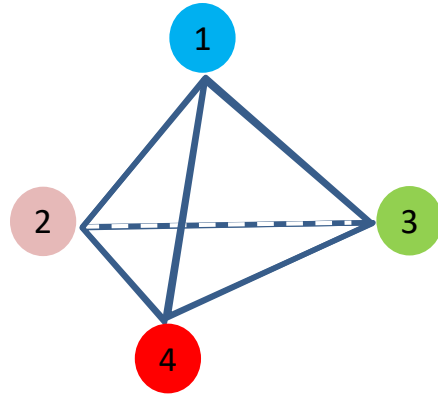


$$D_{123} \oplus A_{234} \oplus B_{134} \oplus C_{124}$$

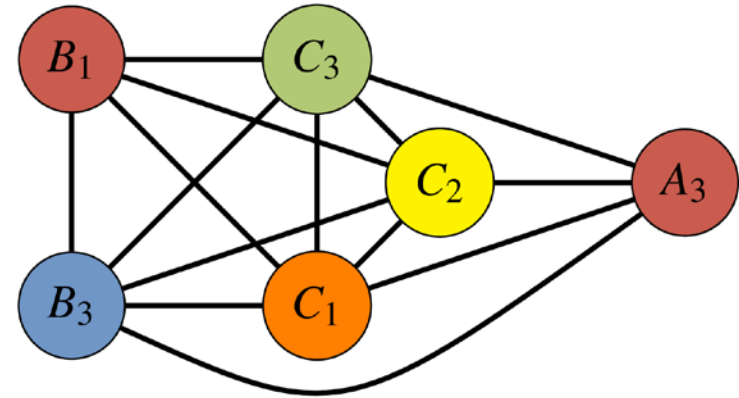


$$A_3 \oplus B_1 \oplus B_3 \oplus C_1 \oplus C_2 \oplus C_3$$

Extension of Clique-based Idea



$$D_{123} \oplus A_{234} \oplus B_{134} \oplus C_{124}$$



$$A_3 \oplus B_1 \oplus B_3 \oplus C_1 \oplus C_2 \oplus C_3$$

Variety of settings:

- Non-uniform demands
- Decentralized placement

Example: Non-uniform File Popularity (random placement)

Example: Non-uniform File Popularity (random placement)

- 3 files $\{A, B, C\}$. A very popular, File B popular, File C not popular

Example: Non-uniform File Popularity (random placement)

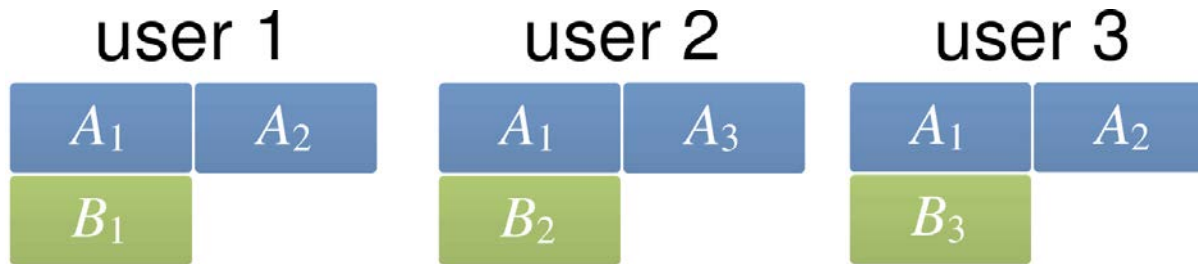
- 3 files $\{A, B, C\}$. A very popular, File B popular, File C not popular
- Split each file into 3 parts each. E.g. $A = \{A_1, A_2, A_3\}$

Example: Non-uniform File Popularity (random placement)

- 3 files $\{A, B, C\}$. A very popular, File B popular, File C not popular
- Split each file into 3 parts each. E.g. $A = \{A_1, A_2, A_3\}$
- Cache distribution $\mathbf{p} = \{A = \frac{2}{3}, B = \frac{1}{3}, C = 0\}$

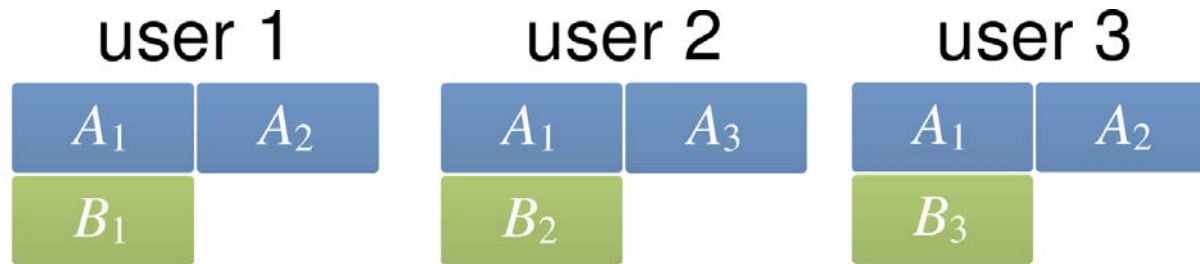
Example: Non-uniform File Popularity (random placement)

Cache realization \mathcal{C}



Example: Non-uniform File Popularity (random placement)

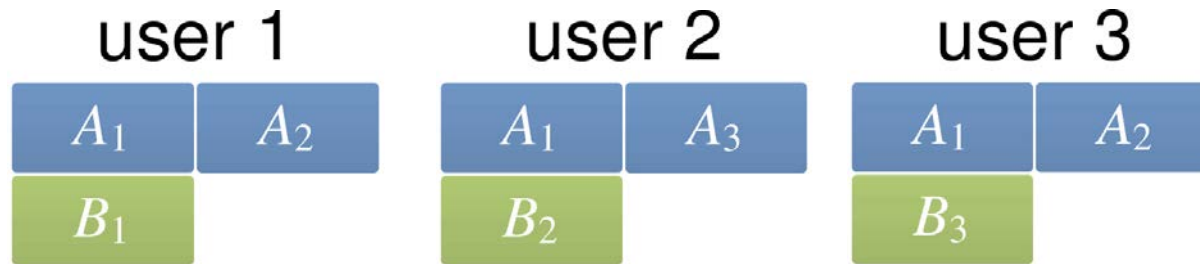
Cache realization \mathcal{C}



Request: $\text{user1} \rightarrow A$, $\text{user2} \rightarrow B$, $\text{user3} \rightarrow C$

Example: Non-uniform File Popularity (random placement)

Cache realization \mathcal{C}



Request: $\text{user1} \rightarrow A$, $\text{user2} \rightarrow B$, $\text{user3} \rightarrow C$

Needed subfiles : $\mathcal{Q} = \{A_3, B_1, B_3, C_1, C_2, C_3\}$

Non-uniform Example: Conflict Graph

Vertex for each requested subpart:

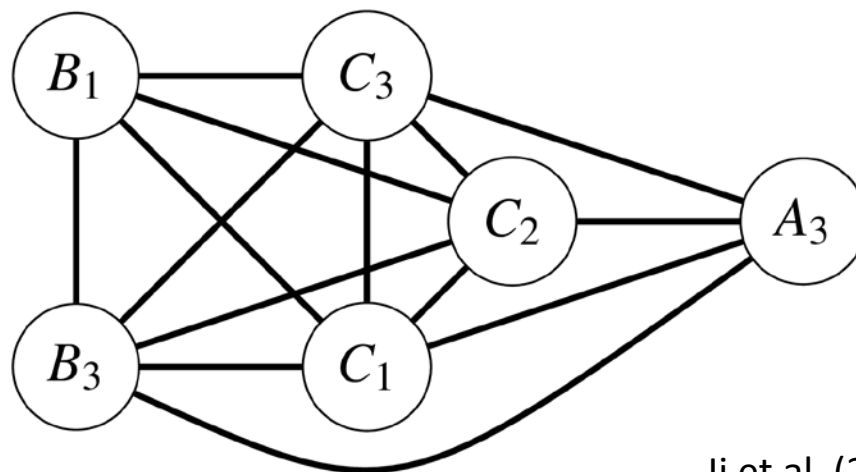
- Replicate if multiple requests of a subfile

Edge if

- Not same identity (cannot connect subfile to itself)
- Request(er) not among users caching the other vertex
 - see (A_3, B_1)

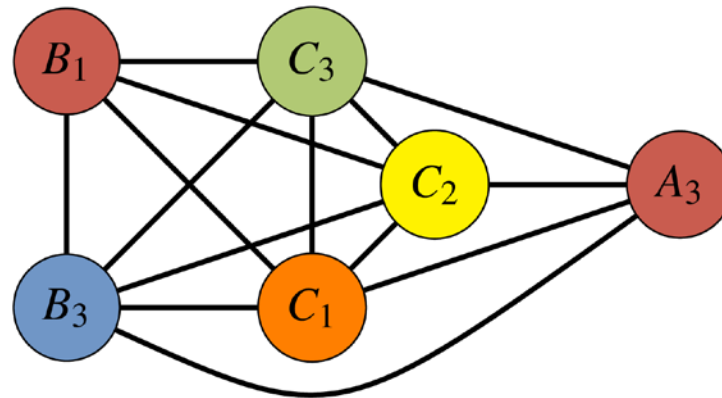
Requests: user1 $\rightarrow A$, user2 $\rightarrow B$, user3 $\rightarrow C$

Queried parts: $Q = \{A_3, B_1, B_3, C_1, C_2, C_3\}$



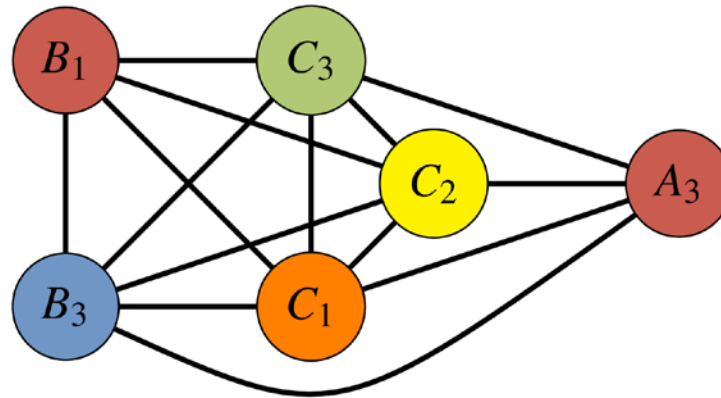
Non-uniform Example: Graph Coloring for XORs

Connected vertices must have different colors



Non-uniform Example: Graph Coloring for XORs

Connected vertices must have different colors

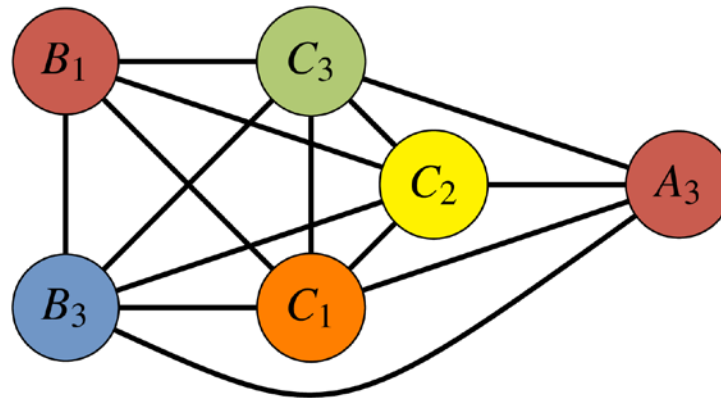


Transmission



Non-uniform Example: Graph Coloring for XORs

Connected vertices must have different colors



Transmission



$$T(\gamma) = 5/3$$

Gain

$$\frac{|Q|}{\chi(H_{C,Q})} \quad (\chi \text{ is chromatic number})$$

Calculation:

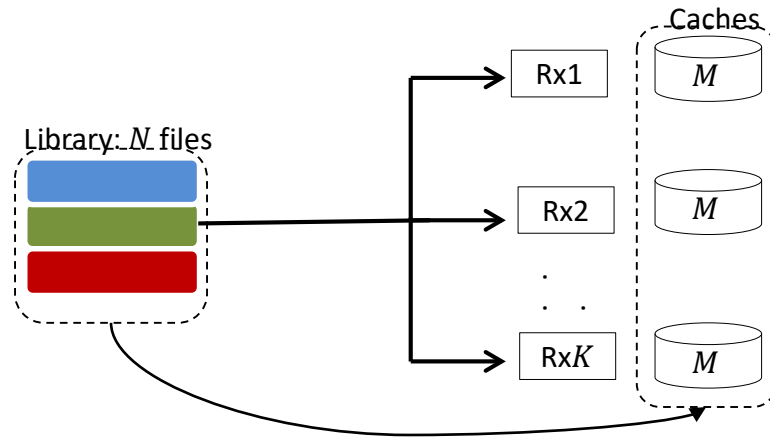
$$\begin{aligned} \frac{|Q|}{\chi(H_{C,Q})} &= \frac{K(1-\gamma)}{T} \\ &= \frac{3\left(1-\frac{1}{3}\right)}{5/3} = \frac{6}{5} \end{aligned}$$

The Optimization Problem:

Optimize how you cache and send each bit

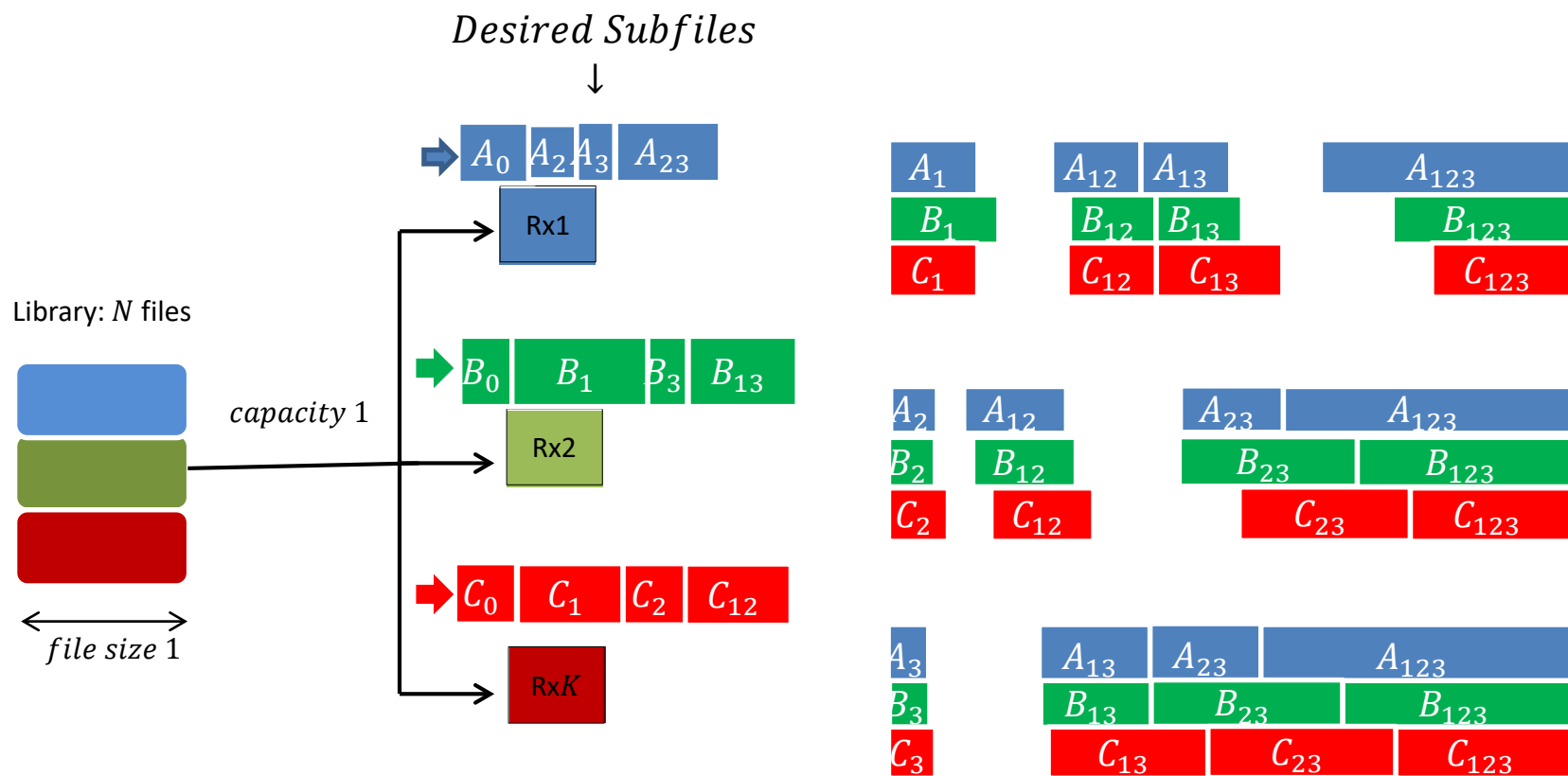
Another (ongoing) Success of
Information Theory

The Optimization Problem (Caching Placement Part)



$$T^*(M) = \min_{\text{schemes } \chi} \max_{\text{demands } \mathbf{d}} T(\mathbf{d}, \chi, M)$$

The Optimization Problem (Delivery Part)



Transmission

Scheme F $\rightarrow F_1(A_{P_1}, B_{Q_1}, C_{R_1}) \quad F_2(A_{P_2}, B_{Q_2}, C_{R_2}) \quad \dots \quad F_{T_F}(A_{P_T}, B_{Q_T}, C_{R_T})$

Scheme Φ $\rightarrow \Phi_1(A_{P_1}, B_{Q_1}, C_{R_1}) \quad \Phi_2(A_{P_2}, B_{Q_2}, C_{R_2}) \quad \dots \quad \Phi_{T_\Phi}(A_{P_T}, B_{Q_T}, C_{R_T})$

Scheme \mathcal{K} $\rightarrow \mathcal{K}_1(A_{P_1}, B_{Q_1}, C_{R_1}) \quad \mathcal{K}_2(A_{P_2}, B_{Q_2}, C_{R_2}) \quad \dots \quad \mathcal{K}_{T_{\mathcal{K}}}(A_{P_T}, B_{Q_T}, C_{R_T})$

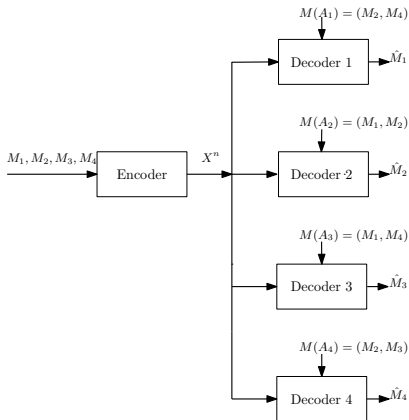
⋮

Tool: Coded Caching \rightarrow Index Coding \rightarrow Graphs

- This is an index coding problem

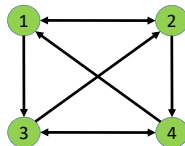
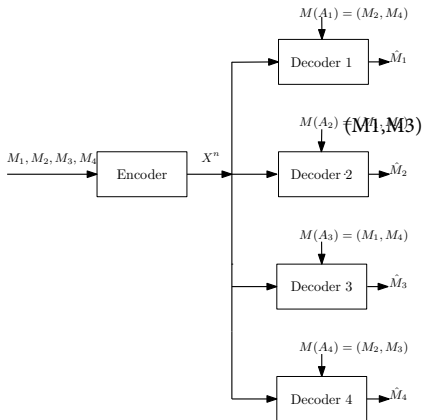
Tool: Coded Caching \rightarrow Index Coding \rightarrow Graphs

- This is an index coding problem
- Example: 4-message index coding problem



Tool: Coded Caching \rightarrow Index Coding \rightarrow Graphs

- This is an index coding problem
- Example: 4-message index coding problem



Tool: Maximum Acyclic Subgraph Bound

Theorem

If G is acyclic, then optimal delay is

$$T^*(G) = |G|.$$

¹Arbabjolfaei et al., 2013.

Tool: Maximum Acyclic Subgraph Bound

Theorem

If G is acyclic, then optimal delay is

$$T^*(G) = |G|.$$



Figure: As if no side information (i.e., TDMA is best you can do)

¹Arbabjolfaei et al., 2013.

Tool: Maximum Acyclic Subgraph Bound

Theorem

If G is acyclic, then optimal delay is

$$T^*(G) = |G|.$$



Figure: As if no side information (i.e., TDMA is best you can do)

Theorem

For any problem corresponding to an arbitrary G , then

$$T \geq \sum_{j \in J} |M_j|$$

for all acyclic subgraphs $J \subset G$.

¹

¹Arbabjolfaei et al., 2013.

Step: Create Graph for One Caching Problem

- Get \mathbf{d} (e.g. $\mathbf{d} = (1, 2, 3)$)

Step: Create Graph for One Caching Problem

- Get \mathbf{d} (e.g. $\mathbf{d} = (1, 2, 3)$)
- General split F_i to $F_{i,\mathcal{W}}$
 - $\mathcal{W} \in 2^{[3]} = \{\emptyset, \{1\}, \{2\}, \{3\}, \{1, 2\}, \{1, 3\}, \{2, 3\}, \{1, 2, 3\}\}$

Step: Create Graph for One Caching Problem

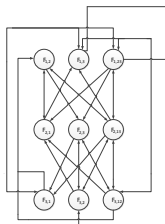
- Get \mathbf{d} (e.g. $\mathbf{d} = (1, 2, 3)$)
- General split F_i to $F_{i,\mathcal{W}}$
 - $\mathcal{W} \in 2^{[3]} = \{\emptyset, \{1\}, \{2\}, \{3\}, \{1, 2\}, \{1, 3\}, \{2, 3\}, \{1, 2, 3\}\}$
- Each node = desired sub-file

Step: Create Graph for One Caching Problem

- Get \mathbf{d} (e.g. $\mathbf{d} = (1, 2, 3)$)
- General split F_i to $F_{i,\mathcal{W}}$
 - $\mathcal{W} \in 2^{[3]} = \{\emptyset, \{1\}, \{2\}, \{3\}, \{1, 2\}, \{1, 3\}, \{2, 3\}, \{1, 2, 3\}\}$
- Each node = desired sub-file
- Row j for user j

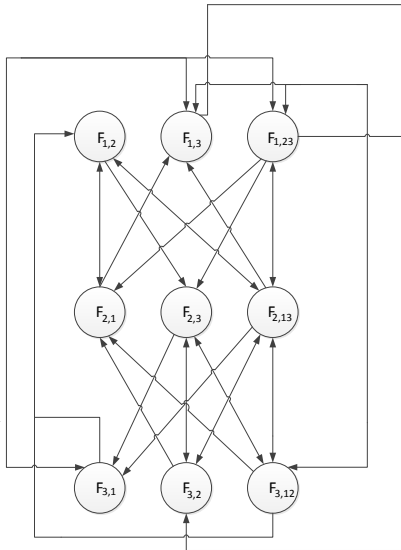
Step: Create Graph for One Caching Problem

- Get \mathbf{d} (e.g. $\mathbf{d} = (1, 2, 3)$)
- General split F_i to $F_{i,\mathcal{W}}$
 - $\mathcal{W} \in 2^{[3]} = \{\emptyset, \{1\}, \{2\}, \{3\}, \{1, 2\}, \{1, 3\}, \{2, 3\}, \{1, 2, 3\}\}$
- Each node = desired sub-file
- Row j for user j
- Edge from i to j if user j knows sub-file i



- K. Wan, D. Tuninetti and P. Piantanida, "On the Optimality of Uncoded Cache Placement", 2016

Step: Create Graph for One Caching Problem



¹Image source: Wan et al. 2017

Step: Create Maximal Acyclic Subgraph - Arrows Down

F_{d_1, \mathcal{W}_1} for all $\mathcal{W}_1 \subseteq [1 : 3] \setminus \{1\}$

F_{d_2, \mathcal{W}_2} for all $\mathcal{W}_2 \subseteq [1 : 3] \setminus \{1, 2\}$

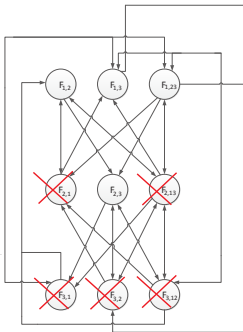
F_{d_3, \mathcal{W}_3} for all $\mathcal{W}_3 \subseteq [1 : 3] \setminus \{1, 2, 3\} = \emptyset$

Step: Create Maximal Acyclic Subgraph - Arrows Down

F_{d_1, \mathcal{W}_1} for all $\mathcal{W}_1 \subseteq [1 : 3] \setminus \{1\}$

F_{d_2, \mathcal{W}_2} for all $\mathcal{W}_2 \subseteq [1 : 3] \setminus \{1, 2\}$

F_{d_3, \mathcal{W}_3} for all $\mathcal{W}_3 \subseteq [1 : 3] \setminus \{1, 2, 3\} = \emptyset$

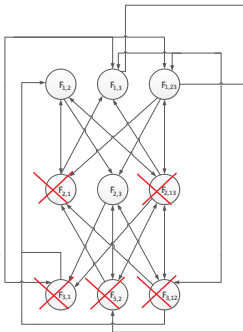


Step: Create Maximal Acyclic Subgraph - Arrows Down

F_{d_1, \mathcal{W}_1} for all $\mathcal{W}_1 \subseteq [1 : 3] \setminus \{1\}$

F_{d_2, \mathcal{W}_2} for all $\mathcal{W}_2 \subseteq [1 : 3] \setminus \{1, 2\}$

F_{d_3, \mathcal{W}_3} for all $\mathcal{W}_3 \subseteq [1 : 3] \setminus \{1, 2, 3\} = \emptyset$



$$T_u(M) \geq |F_{1,\emptyset}| + |F_{1,2}| + |F_{1,3}| + |F_{1,23}| + |F_{2,\emptyset}| + |F_{2,3}| + |F_{3,\emptyset}|$$

Step: More Acyclic Subgraphs - Permute Users ($\pi = (1, 3, 2)$)

$$d_{\pi_1} = d_1 = 1; \mathcal{W}_1 \subseteq [1 : 3] \setminus \{\pi_1\} = \{2, 3\}$$

$$d_{\pi_2} = d_3 = 3; \mathcal{W}_1 \subseteq [1 : 3] \setminus \{\pi_1, \pi_2\} = \{2\}$$

$$d_{\pi_3} = d_2 = 2; \mathcal{W}_3 \subseteq [1 : 3] \setminus \{\pi_1, \pi_2, \pi_3\} = \emptyset$$

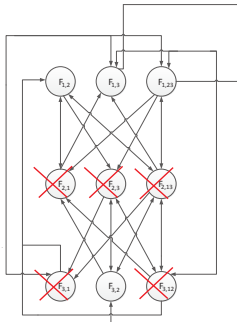
Step: More Acyclic Subgraphs - Permute Users ($\pi = (1, 3, 2)$)

$$d_{\pi_1} = d_1 = 1; \mathcal{W}_1 \subseteq [1 : 3] \setminus \{\pi_1\} = \{2, 3\}$$

$$d_{\pi_2} = d_3 = 3; \mathcal{W}_2 \subseteq [1 : 3] \setminus \{\pi_1, \pi_2\} = \{2\}$$

$$d_{\pi_3} = d_2 = 2; \mathcal{W}_3 \subseteq [1 : 3] \setminus \{\pi_1, \pi_2, \pi_3\} = \emptyset$$

$$T_u(M) \geq |F_{1,\emptyset}| + |F_{1,2}| + |F_{1,3}| + |F_{1,23}| + |F_{3,2}| + |F_{3,\emptyset}| + |F_{2,\emptyset}|$$



Step: More Graphs for Symmetry - Add Acyclic Subgraphs

- $[3!]$ possible demand vectors
- $[3!]$ possible permutations
- sum $[3!]^2$ possible bounds

Step: More Graphs for Symmetry - Add Acyclic Subgraphs

- $[3!]$ possible demand vectors
- $[3!]$ possible permutations
- sum $[3!]^2$ possible bounds

$$T_u(M) \geq |F_{1,\emptyset}| + |F_{1,2}| + |F_{1,3}| + |F_{1,23}| + |F_{2,\emptyset}| + |F_{2,3}| + |F_{3,\emptyset}|$$

$$T_u(M) \geq |F_{1,\emptyset}| + |F_{1,2}| + |F_{1,3}| + |F_{1,23}| + |F_{3,2}| + |F_{3,\emptyset}| + |F_{2,\emptyset}|$$

$$\vdots$$

Step: More Graphs for Symmetry - Add Acyclic Subgraphs

- $[3!]$ possible demand vectors
- $[3!]$ possible permutations
- sum $[3!]^2$ possible bounds

$$T_u(M) \geq |F_{1,\emptyset}| + |F_{1,2}| + |F_{1,3}| + |F_{1,23}| + |F_{2,\emptyset}| + |F_{2,3}| + |F_{3,\emptyset}|$$

$$T_u(M) \geq |F_{1,\emptyset}| + |F_{1,2}| + |F_{1,3}| + |F_{1,23}| + |F_{3,2}| + |F_{3,\emptyset}| + |F_{2,\emptyset}|$$

\vdots

$$T_u(M)(3!)^2 \geq \sum_{\mathbf{d}} \sum_{\pi} \sum_{j \in [3]} \sum_{\mathcal{W}_j \in 2^{[3]} : \mathcal{W}_j \setminus \{\pi_1, \dots, \pi_j\}} |F_{d_{\pi_j}, \mathcal{W}_j}|$$

Step: Exploit Symmetry - Counting Arguments

- Each sub-file that is cached in $|\mathcal{W}| = t$ caches, appears an equal number of times

Step: Exploit Symmetry - Counting Arguments

- Each sub-file that is cached in $|\mathcal{W}| = t$ caches, appears an equal number of times
- Allows for nice transition, from

$$T_u(M) \geq \sum_{j \in [3]} \sum_{\mathcal{W}_j \in 2^{[3]} : \mathcal{W}_j \setminus \{\pi_1, \dots, \pi_j\}} |F_{d_{\pi_j}, \mathcal{W}_j}|$$

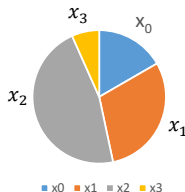
Step: Exploit Symmetry - Counting Arguments

- Each sub-file that is cached in $|\mathcal{W}| = t$ caches, appears an equal number of times
- Allows for nice transition, from

$$T_u(M) \geq \sum_{j \in [3]} \sum_{\mathcal{W}_j \in 2^{[3]}: \mathcal{W}_j \setminus \{\pi_1, \dots, \pi_j\}} |F_{d_{\pi_j}, \mathcal{W}_j}|$$

- to

$$T_u(M) \geq \sum_{i \in [0:3]} x_i \frac{1 - i/3}{1 + i} = x_0 + \frac{1}{3}x_1 + \frac{1}{9}x_2 + 0x_3 \quad (1)$$



Step: Final optimization

- Optimize

$$x_0 + \frac{1}{3} \cdot x_1 + \frac{1}{9} \cdot x_2 + 0 \cdot x_3$$

Step: Final optimization

- Optimize

$$x_0 + \frac{1}{3} \cdot x_1 + \frac{1}{9} \cdot x_2 + 0 \cdot x_3$$

- Under total file size constraint:

$$x_0 + x_1 + x_2 + x_3 = 3$$

- total cache size constraint:

$$0 \cdot x_0 + 1 \cdot x_1 + 2x_2 + 3 \cdot x_3 \leq 3M$$

Step: Final optimization

- Optimize

$$x_0 + \frac{1}{3} \cdot x_1 + \frac{1}{9} \cdot x_2 + 0 \cdot x_3$$

- Under total file size constraint:

$$x_0 + x_1 + x_2 + x_3 = 3$$

- total cache size constraint:

$$0 \cdot x_0 + 1 \cdot x_1 + 2x_2 + 3 \cdot x_3 \leq 3M$$

- Specific solution (*See also Yu, Maddah-Ali, Avestimehr*):

$$x_2 = 3 = 100\%$$



Step: Final optimization

- Optimize

$$x_0 + \frac{1}{3} \cdot x_1 + \frac{1}{9} \cdot x_2 + 0 \cdot x_3$$

- Under total file size constraint:

$$x_0 + x_1 + x_2 + x_3 = 3$$

- total cache size constraint:

$$0 \cdot x_0 + 1 \cdot x_1 + 2x_2 + 3 \cdot x_3 \leq 3M$$

- Specific solution (*See also Yu, Maddah-Ali, Avestimehr*):

$$x_2 = 3 = 100\%$$



- Final answer:

$$T^* \geq \frac{K - t}{1 + t} = \frac{K(1 - \gamma)}{1 + K\gamma}$$

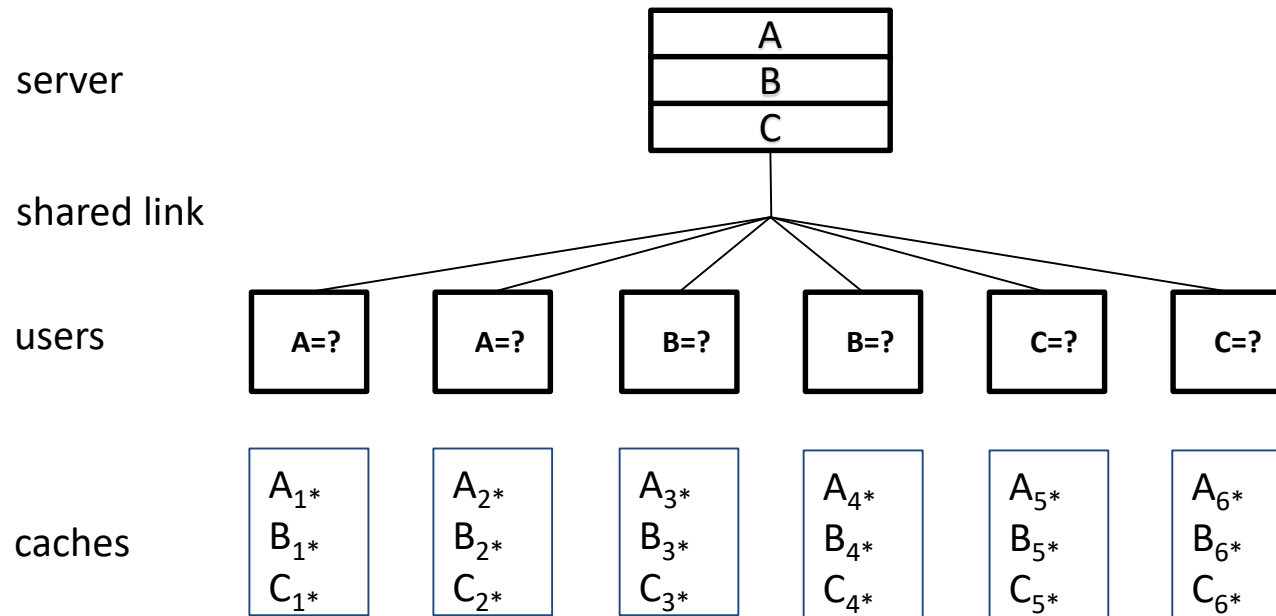
Toward average delay: $N < K$

- Optimality result extends to case of $N < K$

$$T^* = \frac{\binom{K}{K\gamma + 1} - \binom{K - \min(K, N)}{K\gamma + 1}}{\binom{K}{K\gamma}}$$

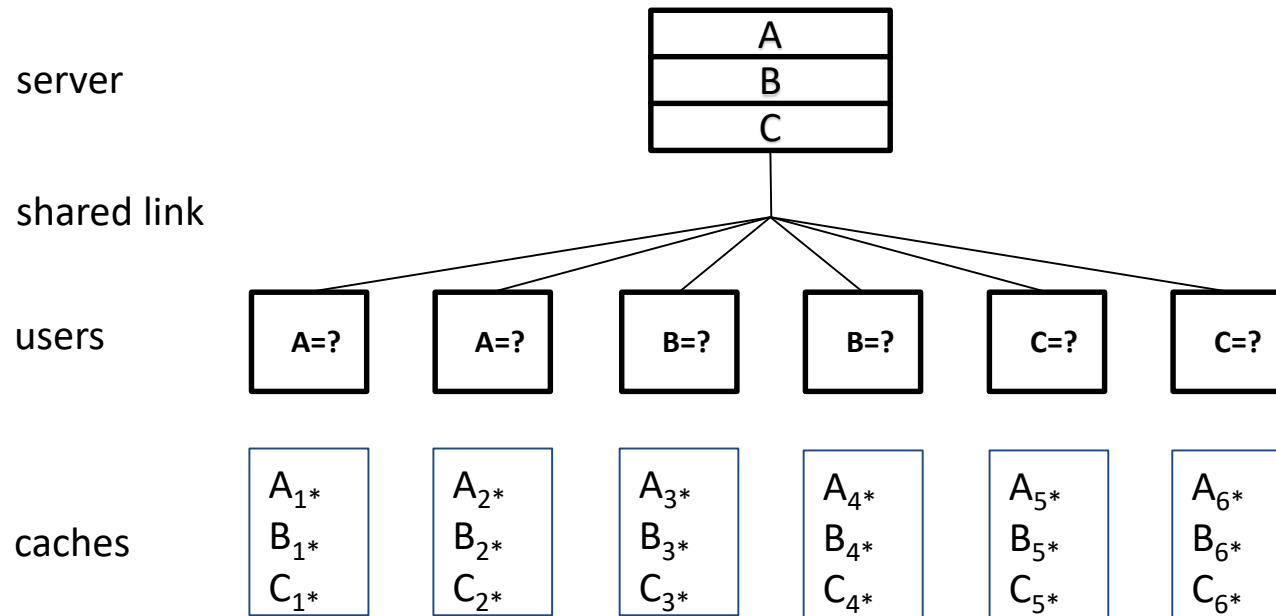
- As well as to average-delay case for uniform distribution
 - Single stream scenario
- Proof based on new scheme for $N < K$

Toward average delay: $N < K$



- $K = 6, N = 3, M = 1$

Toward average delay: $N < K$

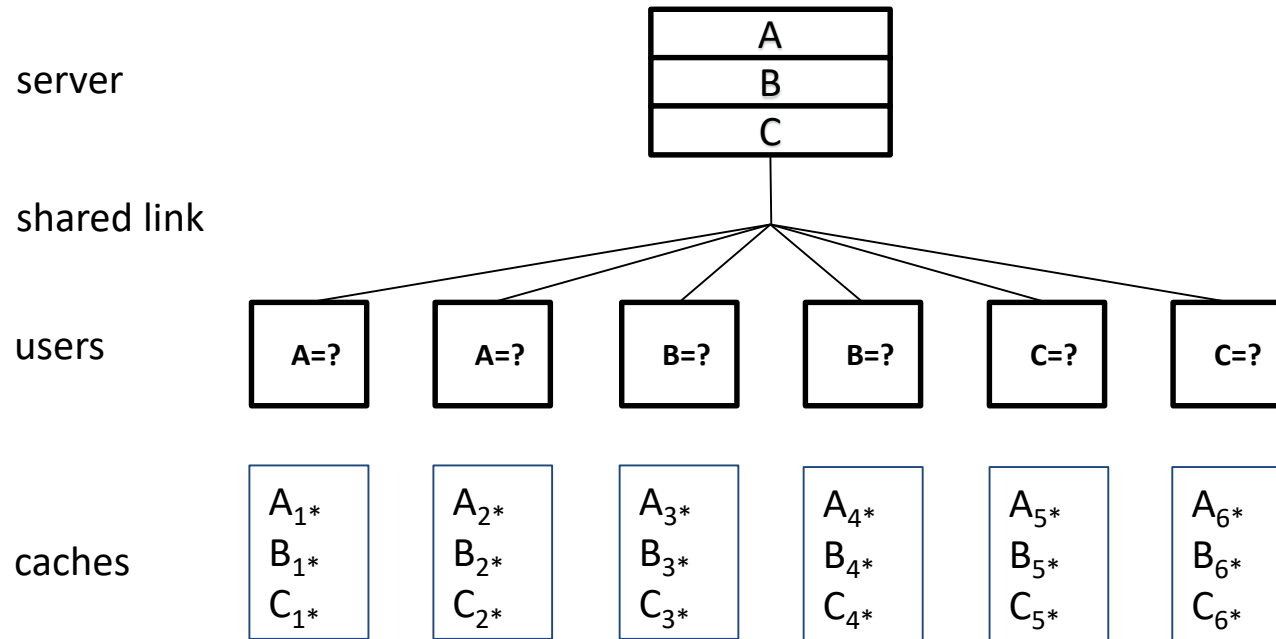


- $K = 6, N = 3, M = 1$
- Subpacketization $\binom{6}{2} = 15$
 $A_{12}, A_{13}, A_{14}, \dots, A_{56},$
 $B_{12}, \dots, B_{56}, C_{12}, \dots, C_{56}.$

- MN Placement
- MN Delivery:

Broadcast $\binom{6}{3} = 20$ XORs $\Rightarrow T = \frac{20}{15}.$

Toward average delay: $N < K$



- $K = 6, N = 3, M = 1$

- Subpacketization $\binom{6}{2} = 15$

$$A_{12}, A_{13}, A_{14}, \dots, A_{56}, \\ B_{12}, \dots, B_{56}, C_{12}, \dots, C_{56}.$$

- MN Placement

- MN Delivery:

$$\text{Broadcast } \binom{6}{3} = 20 \text{ XORs} \Rightarrow T = \frac{20}{15}.$$

- Modified scheme can skip 1 XOR

$$\Rightarrow T = \frac{19}{15} \text{ Optimal}$$

Image-source and Example: Yu et al.

Toward average delay: $N < K$

- Ordinarily would send 20 XORs

- $X_{123} = A_{23} \oplus A_{13} \oplus \mathbf{B_{12}}$

- $X_{124} = A_{24} \oplus A_{14} \oplus \mathbf{B_{12}}$

- $X_{125} = A_{25} \oplus A_{15} \oplus \mathbf{C_{12}}$

- $X_{126} = A_{46} \oplus A_{26} \oplus \mathbf{C_{12}}$

-

- $X_{146} = \mathbf{A_{46}} \oplus B_{16} \oplus C_{14}$

-

- $X_{246} = \mathbf{A_{46}} \oplus B_{26} \oplus C_{24}$

-

- $X_{456} = B_{56} \oplus C_{46} \oplus C_{45}$

Toward average delay: $N < K$

- Ordinarily would send 20 XORs

- $X_{123} = A_{23} \oplus A_{13} \oplus \mathbf{B_{12}}$

- $X_{124} = A_{24} \oplus A_{14} \oplus \mathbf{B_{12}}$

- $X_{125} = A_{25} \oplus A_{15} \oplus \mathbf{C_{12}}$

- $X_{126} = A_{46} \oplus A_{26} \oplus \mathbf{C_{12}}$

-

- $X_{146} = \mathbf{A_{46}} \oplus B_{16} \oplus C_{14}$

-

- $X_{246} = \mathbf{A_{46}} \oplus B_{26} \oplus C_{24}$

-

- $X_{456} = B_{56} \oplus C_{46} \oplus C_{45}$

Toward average delay: $N < K$

- Ordinarily would send 20 XORs

- $X_{123} = A_{23} \oplus A_{13} \oplus \mathbf{B_{12}}$

- $X_{124} = A_{24} \oplus A_{14} \oplus \mathbf{B_{12}}$

- $X_{125} = A_{25} \oplus A_{15} \oplus \mathbf{C_{12}}$

- $X_{126} = A_{46} \oplus A_{26} \oplus \mathbf{C_{12}}$

-

- $X_{146} = \mathbf{A_{46}} \oplus B_{16} \oplus C_{14}$

-

- ~~$X_{246} = A_{46} \oplus B_{26} \oplus C_{24}$~~

-

- $X_{456} = B_{56} \oplus C_{46} \oplus C_{45}$

Toward average delay: $N < K$

- Skip $\mathbf{X}_{246} = \mathbf{A}_{46} \oplus \mathbf{B}_{26} \oplus \mathbf{C}_{24}$
- Focus on user 2 (needs \mathbf{A}_{46})

Toward average delay: $N < K$

- Skip $\mathbf{X}_{246} = \mathbf{A}_{46} \oplus \mathbf{B}_{26} \oplus \mathbf{C}_{24}$
- Focus on user 2 (needs \mathbf{A}_{46})
- Have transmitted:
 - $\mathbf{X}_{146} = \mathbf{A}_{46} \oplus \mathbf{B}_{16} \oplus \mathbf{C}_{14}$
 - $\mathbf{X}_{145} = \mathbf{A}_{45} \oplus \mathbf{B}_{15} \oplus \mathbf{C}_{14}$
 - $\mathbf{X}_{136} = \mathbf{A}_{36} \oplus \mathbf{B}_{16} \oplus \mathbf{C}_{13}$
 - $\mathbf{X}_{135} = \mathbf{A}_{35} \oplus \mathbf{B}_{15} \oplus \mathbf{C}_{13}$

Toward average delay: $N < K$

- Skip $\mathbf{X}_{246} = A_{46} \oplus B_{26} \oplus C_{24}$
- Focus on user 2 (needs \mathbf{A}_{46})
- Have transmitted:
 - $X_{146} = A_{46} \oplus B_{16} \oplus C_{14}$
 - $X_{145} = A_{45} \oplus B_{15} \oplus C_{14}$
 - $X_{136} = A_{36} \oplus B_{16} \oplus C_{13}$
 - $X_{135} = A_{35} \oplus B_{15} \oplus C_{13}$
- Adding up:
 - $X_{146} \oplus X_{145} \oplus X_{136} \oplus X_{135} = A_{46} \oplus \mathbf{A}_{45} \oplus \mathbf{A}_{36} \oplus \mathbf{A}_{35}$

Toward average delay: $N < K$

- Skip $\mathbf{X}_{246} = \mathbf{A}_{46} \oplus \mathbf{B}_{26} \oplus \mathbf{C}_{24}$
- Focus on user 2 (needs \mathbf{A}_{46})
- Have transmitted:
 - $\mathbf{X}_{146} = \mathbf{A}_{46} \oplus \mathbf{B}_{16} \oplus \mathbf{C}_{14}$
 - $\mathbf{X}_{145} = \mathbf{A}_{45} \oplus \mathbf{B}_{15} \oplus \mathbf{C}_{14}$
 - $\mathbf{X}_{136} = \mathbf{A}_{36} \oplus \mathbf{B}_{16} \oplus \mathbf{C}_{13}$
 - $\mathbf{X}_{135} = \mathbf{A}_{35} \oplus \mathbf{B}_{15} \oplus \mathbf{C}_{13}$
- Adding up:
 - $\mathbf{X}_{146} \oplus \mathbf{X}_{145} \oplus \mathbf{X}_{136} \oplus \mathbf{X}_{135} = \mathbf{A}_{46} \oplus \mathbf{A}_{45} \oplus \mathbf{A}_{36} \oplus \mathbf{A}_{35}$
- $\mathbf{A}_{45} \leftarrow \mathbf{X}_{245}$
- $\mathbf{A}_{36} \leftarrow \mathbf{X}_{236}$
- $\mathbf{A}_{35} \leftarrow \mathbf{X}_{235}$

Toward average delay: $N < K$

- Skip $\mathbf{X}_{246} = \mathbf{A}_{46} \oplus \mathbf{B}_{26} \oplus \mathbf{C}_{24}$
- Focus on user 2 (needs \mathbf{A}_{46})
- Have transmitted:
 - $\mathbf{X}_{146} = \mathbf{A}_{46} \oplus \mathbf{B}_{16} \oplus \mathbf{C}_{14}$
 - $\mathbf{X}_{145} = \mathbf{A}_{45} \oplus \mathbf{B}_{15} \oplus \mathbf{C}_{14}$
 - $\mathbf{X}_{136} = \mathbf{A}_{36} \oplus \mathbf{B}_{16} \oplus \mathbf{C}_{13}$
 - $\mathbf{X}_{135} = \mathbf{A}_{35} \oplus \mathbf{B}_{15} \oplus \mathbf{C}_{13}$
- Adding up:
 - $\mathbf{X}_{146} \oplus \mathbf{X}_{145} \oplus \mathbf{X}_{136} \oplus \mathbf{X}_{135} = \mathbf{A}_{46} \oplus \mathbf{A}_{45} \oplus \mathbf{A}_{36} \oplus \mathbf{A}_{35}$
- $\mathbf{A}_{45} \leftarrow \mathbf{X}_{245} = \mathbf{A}_{45} \oplus \mathbf{B}_{25} \oplus \mathbf{C}_{24}$
- $\mathbf{A}_{36} \leftarrow \mathbf{X}_{236} = \mathbf{A}_{36} \oplus \mathbf{B}_{26} \oplus \mathbf{C}_{23}$
- $\mathbf{A}_{35} \leftarrow \mathbf{X}_{235} = \mathbf{A}_{35} \oplus \mathbf{B}_{25} \oplus \mathbf{C}_{23}$
- \Rightarrow Decode $\mathbf{A}_{\{4,6\}}$ (\Rightarrow similarly \mathbf{B}_{26} for user 4, \mathbf{C}_{24} for user 6).

Extensions wide open

- Done for very special case of the single-stream BC

Extensions wide open

- Done for very special case of the single-stream BC
- Approach requires new 'tricks' when setting changes
 - Shared caches (Hachem, Karamchandrani et Diggavi)
 - Multiple file demands (Wei-Ulukus)
 - Multi-layer files (high-def, low-def) (Yang-Gündüz)

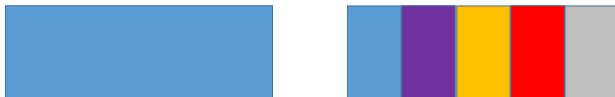
Extensions wide open

- Done for very special case of the single-stream BC
- Approach requires new 'tricks' when setting changes
 - Shared caches (Hachem, Karamchandrani et Diggavi)
 - Multiple file demands (Wei-Ulukus)
 - Multi-layer files (high-def, low-def) (Yang-Gündüz)
- Need methods needed that preserve some asymmetry



Extensions wide open

- Done for very special case of the single-stream BC
- Approach requires new 'tricks' when setting changes
 - Shared caches (Hachem, Karamchandrani et Diggavi)
 - Multiple file demands (Wei-Ulukus)
 - Multi-layer files (high-def, low-def) (Yang-Gündüz)
- Need methods needed that preserve some asymmetry



- Need methods that preserve topology
- Need methods reflecting multiple-senders/multiple-antennas.