

# Fundamental Limits of Coded Caching: From Uncoded Prefetching to Coded Prefetching

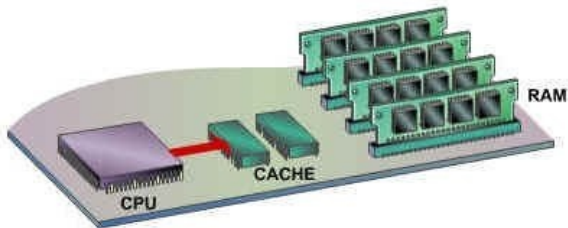
Kai Zhang and Chao Tian

Texas A&M University

ISIT 2018

# Coded Caching and Its Applications

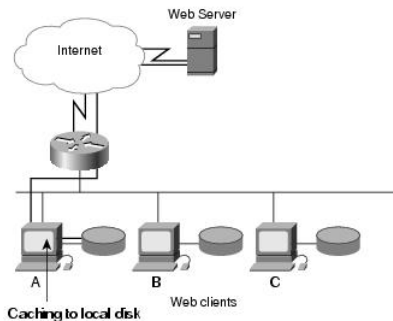
- ▶ A data management strategy to reduce delay during peak-traffic time,
- ▶ Single user setting, e.g. on-CPU caches, RAM in computers,
- ▶ Multiple user setting, e.g. networked system.



When the CPU needs data, it looks first in cache memory.

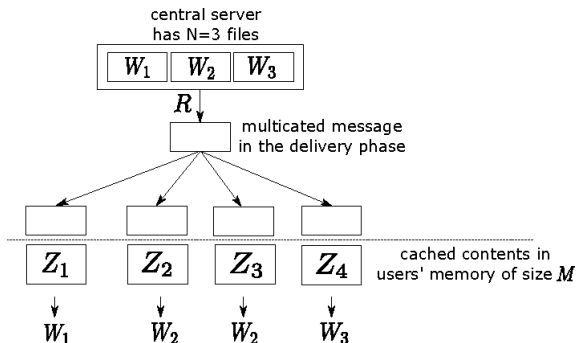
# Coded Caching and Its Applications

- ▶ A data management strategy to reduce delay during peak-traffic time,
- ▶ Single user setting, e.g. on-CPU caches, RAM in computers,
- ▶ Multiple user setting, e.g. networked system.



# An Information-theoretic Formulation

- ▶  $N$  files of same size,  $K$  users, each with a local cache of size  $M$ ,
- ▶ Prefetching phase: users fill their cache without knowing the requests,
- ▶ Delivery phase: each user requests a single file, server multicasts information to accommodate the requests.



# Existing Schemes Using Uncoded Prefetching

Maddah-Ali & Niesen<sup>1</sup>, improved by Yu et al<sup>2</sup>, proved being optimal,

► Prefetching strategy:

- Partition each file into  $\binom{K}{t}$  segments, each associated with a cardinality- $t$  subset,
- Each user caches  $\binom{K-1}{t-1}$  segments,
- E.g.:  $(N, K) = (3, 4)$ ,  $t = 2$ ,

User 1	$A_{12}$	$A_{13}$	$A_{14}$	$B_{12}$	$B_{13}$	$B_{14}$	$C_{12}$	$C_{13}$	$C_{14}$
User 2	$A_{12}$	$A_{23}$	$A_{24}$	$B_{12}$	$B_{23}$	$B_{24}$	$C_{12}$	$C_{23}$	$C_{24}$
User 3	$A_{13}$	$A_{23}$	$A_{34}$	$B_{13}$	$B_{23}$	$B_{34}$	$C_{13}$	$C_{23}$	$C_{34}$
User 4	$A_{14}$	$A_{24}$	$A_{34}$	$B_{14}$	$B_{24}$	$B_{34}$	$C_{14}$	$C_{24}$	$C_{34}$

► Delivery strategy

- For any  $(t + 1)$  users subset  $\mathcal{B}$ , send the XOR of requested segments.

---

<sup>1</sup>Maddah-Ali M A, Niesen U. Fundamental limits of caching[J]. IEEE Transactions on Information Theory, 2014, 60(5): 2856-2867.

<sup>2</sup>Yu Q, Maddah-Ali M A, Avestimehr A S. The exact rate-memory tradeoff for caching with uncoded prefetching[J]. IEEE Transactions on Information Theory, 2017.

# Existing Schemes Using Uncoded Prefetching

- ▶ Server transmits:

$$\mathcal{B} = \{1, 2, 3\} : A_{23} + A_{13} + B_{12}, \mathcal{B} = \{1, 2, 4\} : A_{24} + A_{14} + C_{12}$$

$$\mathcal{B} = \{1, 3, 4\} : A_{34} + B_{14} + C_{13}, \mathcal{B} = \{2, 3, 4\} : A_{34} + B_{24} + C_{23}$$

Using a brief moment, let's look at the decoding steps for the 1st user, when demand is (A,A,B,C):

# Existing Schemes Using Uncoded Prefetching

- Demand (A,A,B,C), server transmits:

$$\begin{aligned} &A_{23} + A_{13} + B_{12}, \quad A_{24} + A_{14} + C_{12} \\ &A_{34} + B_{14} + C_{13}, \quad A_{34} + B_{24} + C_{23} \end{aligned}$$

User 1	$A_{12}$	$A_{13}$	$A_{14}$	$B_{12}$	$B_{13}$	$B_{14}$	$C_{12}$	$C_{13}$	$C_{14}$
User 2	$A_{12}$	$A_{23}$	$A_{24}$	$B_{12}$	$B_{23}$	$B_{24}$	$C_{12}$	$C_{23}$	$C_{24}$
User 3	$A_{13}$	$A_{23}$	$A_{34}$	$B_{13}$	$B_{23}$	$B_{34}$	$C_{13}$	$C_{23}$	$C_{34}$
User 4	$A_{14}$	$A_{24}$	$A_{34}$	$B_{14}$	$B_{24}$	$B_{34}$	$C_{14}$	$C_{24}$	$C_{34}$

Step 1: using  $A_{13}$  and  $B_{12}$  to decode  $A_{23}$ ;

# Existing Schemes Using Uncoded Prefetching

- Demand (A,A,B,C), server transmits:

$$\begin{aligned} &A_{23} + A_{13} + B_{12}, \quad A_{24} + A_{14} + C_{12} \\ &A_{34} + B_{14} + C_{13}, \quad A_{34} + B_{24} + C_{23} \end{aligned}$$

User 1	$A_{12}$	$A_{13}$	$A_{14}$	$B_{12}$	$B_{13}$	$B_{14}$	$C_{12}$	$C_{13}$	$C_{14}$
User 2	$A_{12}$	$A_{23}$	$A_{24}$	$B_{12}$	$B_{23}$	$B_{24}$	$C_{12}$	$C_{23}$	$C_{24}$
User 3	$A_{13}$	$A_{23}$	$A_{34}$	$B_{13}$	$B_{23}$	$B_{34}$	$C_{13}$	$C_{23}$	$C_{34}$
User 4	$A_{14}$	$A_{24}$	$A_{34}$	$B_{14}$	$B_{24}$	$B_{34}$	$C_{14}$	$C_{24}$	$C_{34}$

Step 2: using  $A_{14}$  and  $C_{12}$  to decode  $A_{24}$ ;



# Existing Schemes Using Uncoded Prefetching

- Demand (A,A,B,C), server transmits

$$A_{23} + A_{13} + B_{12}, \quad A_{24} + A_{14} + C_{12}$$

$$A_{34} + B_{14} + C_{13}, \quad A_{34} + B_{24} + C_{23}$$

User 1	$A_{12}$	$A_{13}$	$A_{14}$	$B_{12}$	$B_{13}$	$B_{14}$	$C_{12}$	$C_{13}$	$C_{14}$
User 2	$A_{12}$	$A_{23}$	$A_{24}$	$B_{12}$	$B_{23}$	$B_{24}$	$C_{12}$	$C_{23}$	$C_{24}$
User 3	$A_{13}$	$A_{23}$	$A_{34}$	$B_{13}$	$B_{23}$	$B_{34}$	$C_{13}$	$C_{23}$	$C_{34}$
User 4	$A_{14}$	$A_{24}$	$A_{34}$	$B_{14}$	$B_{24}$	$B_{34}$	$C_{14}$	$C_{24}$	$C_{34}$

Step 3: using  $B_{14}$  and  $C_{13}$  to decode  $A_{34}$ ;

User 1 now decodes the entire file A.

# Existing Schemes Using Coded Prefetching

Tian-Chen scheme<sup>1</sup>,

► Prefetching strategy:

- A file is partitioned into  $\binom{K}{t}$  segments, each user allocated  $\binom{K-1}{t-1}$  segments,
- User encodes file segments using a rank metric code and stores the parities,
- E.g.:  $(N, K) = (3, 4)$ ,  $t = 2$ , user 1 caches

$$\begin{pmatrix} \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \end{pmatrix}_{(5 \times 9)} \cdot$$
$$(A_{12}, A_{13}, A_{14}, B_{12}, B_{13}, B_{14}, C_{12}, C_{13}, C_{14})^T$$

In the example, a  $(14, 9)$  systematic rank metric code is used,  
user 1 stores the 5 parities.

---

<sup>1</sup>Tian C, Chen J. Caching and delivery via interference elimination[C]. Information Theory (ISIT), 2016 IEEE International Symposium on. IEEE, 2016: 830-834.

# Existing Schemes Using Coded Prefetching

- ▶ Delivery strategy: Linear combination of segments from one file.
  - E.g.: demand (A,A,B,C), transmits 9 (coded) symbols:

$$A_{34}, B_{12}, B_{14}, B_{24}, C_{12}, C_{13}, C_{23}, A_{13} + A_{23}, A_{14} + A_{24}$$

- Each symbol serves two roles:
  - 1) content delivery for some users;
  - 2) help resolve coded symbols for some other users.

To be clear, let's check user 1,

# Existing Schemes Using Coded Prefetching

- Demand (A,A,B,C), server transmits:

$$A_{34}, B_{12}, B_{14}, B_{24}, C_{12}, C_{13}, C_{23}, A_{13} + A_{23}, A_{14} + A_{24}$$

$$\begin{pmatrix} \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ & & & & & & & & \cdot \\ & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \end{pmatrix} \quad \cdot$$

$$(A_{12}, A_{13}, A_{14}, B_{12}, B_{13}, B_{14}, C_{12}, C_{13}, C_{14})^T$$

Step 1: user 1 collects 4 symbols  $B_{12}, B_{14}, C_{12}, C_{13}$ ;

Together with 5 (coded) symbols in his cache, he can decode all 9 symbols.

# Existing Schemes Using Coded Prefetching

- Demand (A,A,B,C), server transmits:

$$A_{34}, B_{12}, B_{14}, B_{24}, C_{12}, C_{13}, C_{23}, A_{13} + A_{23}, A_{14} + A_{24}$$

Step 2: user 1 already resolved all symbols in the cache, as following:

User 1	$A_{12}$	$A_{13}$	$A_{14}$	$B_{12}$	$B_{13}$	$B_{14}$	$C_{12}$	$C_{13}$	$C_{14}$
--------	----------	----------	----------	----------	----------	----------	----------	----------	----------

He then collects  $A_{13} + A_{23}, A_{14} + A_{24}$ , he can decode  $A_{23}, A_{24}$ ;

# Existing Schemes Using Coded Prefetching

- Demand (A,A,B,C), server transmits:

$$A_{34}, B_{12}, B_{14}, B_{24}, C_{12}, C_{13}, C_{23}, A_{13} + A_{23}, A_{14} + A_{24}$$

Step 3: user 1 collects  $A_{34}$ , now he decodes the entire file A;

# Two Quite Different Schemes

- ▶ Uncoded prefetching vs. coded prefetching,
- ▶ Binary code vs. non-binary code,
- ▶ Encoding/decoding using simple combinatorics vs. sophisticated coding techniques (rank metric codes),
- ▶ Uncoded prefetching performs better at high cache memory regime vs. coded prefetching performs better at low cache memory regime.

*Largely unrelated...*

*Are there any connections?*

# Some Backgrounds

- ▶ Linearized Polynomial: A degree- $q^{P-1}$  linearized polynomial in the finite field  $\mathbb{F}_{q^m}$

$$f(x) = \sum_{i=1,2,\dots,P} v_i x^{q^{i-1}}, v_i \in \mathbb{F}_{q^m} \quad (1)$$

can be uniquely identified from evaluations at any  $P$  points  $x = \theta_i \in \mathbb{F}_{q^m}, i = 1, 2, \dots, P$ , that are linearly independent over  $\mathbb{F}_q$ .



# Some Backgrounds

## Lemma

Let  $f(x)$  be a degree- $q^{P-1}$  linearized polynomial in  $\mathbb{F}_{q^m}$ , and  $\theta_i \in \mathbb{F}_{q^m}, i = 1, 2, \dots, P_o$ , be linearly independent over  $\mathbb{F}_q$ . Let  $G$  be a  $P_o \times P$  full rank (rank  $P$ ) matrix with entries in  $\mathbb{F}_q$ , then  $f(x)$  can be uniquely identified from

$$[f(\theta_1), f(\theta_2), \dots, f(\theta_{P_o})] \cdot G. \quad (2)$$

- $(v_1, \dots, v_P)$  are information symbols to be encoded;  
 $[f(\theta_1), \dots, f(\theta_{P_o})]$  are coded symbols;
- A  $(P_o, P)$  MDS code in terms of rank metric.

## A Hidden Connection

Continuing example  $(N, K) = (3, 4)$ ,  $t = 2$ , demand = (A,A,B,C):

- ▶ Yu et al scheme transmits (4 symbols):

$$\begin{array}{ll} A_{23} + A_{13} + B_{12}, & A_{24} + A_{14} + C_{12} \\ A_{34} + B_{14} + C_{13}, & A_{34} + B_{24} + C_{23} \end{array}$$

Separate different files & remove repeated transmissions,



- ▶ Tian-Chen scheme (9 symbols):

$$\begin{array}{llll} A_{23} + A_{13}, & B_{12}, & A_{24} + A_{14}, & C_{12} \\ A_{34}, & B_{14}, & C_{13}, & B_{24}, & C_{23} \end{array}$$

Instead of fully decomposing a transmission, what if we partially decompose them?

# An Alternative Scheme

- ▶ A partial decomposition: deliver 5 coded symbols:

- Demand (A,A,B,C), transmits

$$\begin{array}{ll} A_{23} + A_{13} + B_{12}, & A_{24} + A_{14} + C_{12} \\ A_{34}, & B_{14} + C_{13}, \quad B_{24} + C_{23} \end{array}$$

- Demand (A,A,B,B), transmits

$$\begin{array}{ll} A_{23} + A_{13} + B_{12}, & A_{24} + A_{14} + B_{12} \\ A_{34}, & B_{14} + B_{13}, \quad B_{24} + B_{23} \end{array}$$

# An Alternative Scheme

- ▶ A partial decomposition: deliver 5 coded symbols:

- Demand (A,A,A,C), transmits

$$\begin{array}{ll} A_{23} + A_{13} + A_{12}, & A_{24} + A_{14} + C_{12} \\ A_{34}, & A_{14} + C_{13}, \quad A_{24} + C_{23} \end{array}$$

- Demand (A,A,A,A), transmits

$$\begin{array}{ll} A_{23} + A_{13} + A_{12}, & A_{24} + A_{14} + A_{12} \\ A_{34}, & A_{14} + A_{13}, \quad A_{24} + A_{23} \end{array}$$

# An Alternative Scheme

- ▶ Each user only needs to cache 8 linear combinations,
- ▶ Every transmission could have different ways to be partially decomposed, i.e. *decomposition patterns*,
- ▶ Above is a good choice of decomposing, in fact produces a new corner point at  $(M, R) = (4/3, 5/6)$ .

We will give a generalized scheme and code examples.

# Partial Decomposition and A New Code Example

To facilitate understanding, we introduce some notions:

- ▶ **Demand vector**  $\mathbf{d} = \{d_1, \dots, d_K\}$  denotes the demands by  $K$  users,  $d_k \in [1 : N]$ ,
- ▶ Each transmission in Yu scheme can be represented as

$$\bigoplus_{k \in \mathcal{B}} W_{d_k, \mathcal{B} \setminus k}, \quad |\mathcal{B}| = t + 1,$$

a **transmission type**  $\mathbf{t} = (t_1, \dots, t_N)$  is associated with a transmission where  $t_n$  denotes the number of users in subset  $|\mathcal{B}|$  requesting file  $W_n$ . Set  $\mathcal{T}_{\mathbf{d}}^{(t)}$  denotes all valid transmission types for a demand  $\mathbf{d}$ .

# Partial Decomposition and A New Code Example

- ▶ A **partial decomposition**  $\mathcal{P}_{\mathbf{t},\mathbf{d}}$  of transmission type  $\mathbf{t}$  is a partition of  $\text{supp}(\mathbf{t})$ . The full set of decomposition patterns is  $\mathcal{P}_{\mathbf{d}}^{(t)} = \{\mathcal{P}_{\mathbf{t},\mathbf{d}} | \mathbf{t} \in \mathcal{T}_{\mathbf{d}}^{(t)}\}$ .
- ▶ The collection of all possible  $\mathcal{P}_{\mathbf{t},\mathbf{d}}$  is  $\mathfrak{P}_{\mathbf{t},\mathbf{d}}$ .
- ▶ E.g.:  $(N, K) = (3, 4)$ ,  $t = 2$ ,  $\mathbf{d} = (1, 1, 2, 3)$ ,  
 $\mathcal{T}_{\mathbf{d}}^{(t)} = \{(2, 1, 0), (2, 0, 1), (1, 1, 1)\}$ , for  $\mathbf{t} = (1, 1, 1)$ :  
$$\begin{aligned} \mathcal{B} = \{1, 3, 4\} : A_{34} + B_{14} + C_{13} &\xrightarrow{\text{decompose}} \{A_{34}, B_{14} + C_{13}\} \\ \mathcal{B} = \{2, 3, 4\} : A_{34} + B_{24} + C_{23} &\{A_{34}, B_{24} + C_{23}\} \end{aligned}$$

Decomposition pattern  $\mathcal{P}_{\mathbf{t},\mathbf{d}} = \mathcal{P}_{(1,1,1)(1,1,2,3)} = \{\{1\}, \{2, 3\}\}$ .

# Partial Decomposition and A New Code Example

Partial decomposition could produce new memory-rate points:

- ▶ For each demand  $\mathbf{d}$ , using multiple coding instances,
- ▶ Fix a decomposition pattern  $\mathcal{P}_{\mathbf{t},\mathbf{d}}$  for each  $\mathbf{t}$ ,
- ▶ Set an auxiliary variable  $\alpha_{\mathbf{d},\mathcal{P}_{\mathbf{d}}^{(t)}}$  for each decomposition pattern, s.t.

$$\sum_{\mathcal{P}_{\mathbf{d}}^{(t)} \in \mathfrak{P}_{\mathbf{t},\mathbf{d}}} \alpha_{\mathbf{d},\mathcal{P}_{\mathbf{d}}^{(t)}} = 1, \quad \mathbf{d} \in \mathcal{D},$$
$$1 \geq \alpha_{\mathbf{d},\mathcal{P}_{\mathbf{d}}^{(t)}} \geq 0, \quad \mathbf{d} \in \mathcal{D}, \quad \mathcal{P}_{\mathbf{d}}^{(t)} \in \mathfrak{P}_{\mathbf{t},\mathbf{d}}.$$

- $\alpha_{\mathbf{d},\mathcal{P}_{\mathbf{d}}^{(t)}}$ : how much fraction of instances using a decomposition pattern;
- Intuitively same as “time sharing”, but not equivalent.



# Partial Decomposition and A New Code Example

## ► Prefetching:

- Partition each file into  $r\binom{K}{t}$  segments ( $r$  instances each of  $\binom{K}{t}$  segments), each to be in  $\mathbb{F}_{2^m}$  for sufficiently large  $m$ ,
- Each user is allocated  $P = rN\binom{K-1}{t-1}$  symbols; caches  $P_o - P$  linear combinations:

$$P_o - P = rM'_r\binom{K}{t}$$

( $M'_r$  to be defined later)

- Using  $(P_o, P)$  systematic rank metric code to encode  $P$  symbols, caches the parities;
- $m \geq P_o$  suffices for existence of rank metric code.

# Partial Decomposition and A New Code Example

► Delivery:

- For demand  $\mathbf{d} = (A,A,B,C)=(1,1,2,3)$ , let  $\alpha_{\mathbf{d},\mathcal{P}_d^{(t)}} = 1$ ,  $\mathcal{P}_d^{(t)} = \mathcal{P}_{(1,1,2,3)}^{(2)}$  are:

$$\mathcal{P}_{(2,1,0),(1,1,2,3)} = \{\{1, 2\}\} = \{\{A, B\}\} : A_{23} + A_{13} + B_{12},$$

$$\mathcal{P}_{(2,0,1),(1,1,2,3)} = \{\{1, 3\}\} = \{\{A, C\}\} : A_{24} + A_{14} + C_{12},$$

$$\mathcal{P}_{(1,1,1),(1,1,2,3)} = \{\{1\}, \{2, 3\}\} = \{\{A\}, \{B, C\}\} : \\ A_{34}, \quad B_{14} + C_{13}, \quad B_{24} + C_{23}$$

Thus  $R_{\mathbf{d},\mathcal{P}_d^{(t)}} = 5$ ,  $M_{\mathbf{d},\mathcal{P}_d^{(t)},k} = 8$ , for  $k \in [1 : 4]$ .

# Partial Decomposition and A New Code Example

► Delivery:

- For demand  $\mathbf{d} = (A,A,B,B)=(1,1,2,2)$ , two decomposition patterns:

1st pattern:  $\mathcal{P}_{\mathbf{d}}^{(t)}$  without any decomposition, let

$$\alpha_{\mathbf{d}, \mathcal{P}_{\mathbf{d}}^{(t)}} = 1/2,$$

$$A_{23}^{(1)} + A_{13}^{(1)} + B_{12}^{(1)}, \quad A_{24}^{(1)} + A_{14}^{(1)} + B_{12}^{(1)}$$

$$A_{34}^{(1)} + B_{14}^{(1)} + B_{13}^{(1)}, \quad A_{34}^{(1)} + B_{24}^{(1)} + B_{23}^{(1)}$$

For this pattern,  $R_{\mathbf{d}, \mathcal{P}_{\mathbf{d}}^{(t)}} = 4, M_{\mathbf{d}, \mathcal{P}_{\mathbf{d}}^{(t)}, k} = 9, k \in [1 : 4]$ .

# Partial Decomposition and A New Code Example

► Delivery:

- For demand  $\mathbf{d} = (A,A,B,B)=(1,1,2,2)$ , two decomposition patterns:

2nd pattern: the following decompositions  $\mathcal{P}_d^{(t)}$ :

$$\mathcal{P}_{(2,1,0),(1,1,2,2)} = \{\{1\}, \{2\}\} = \{\{A\}, \{B\}\} : A_{23}^{(2)} + A_{13}^{(2)}, B_{12}^{(2)}, \\ A_{24}^{(2)} + A_{14}^{(2)}, B_{12}^{(2)};$$

$$\mathcal{P}_{(1,2,0),(1,1,2,2)} = \{\{1\}, \{2\}\} = \{\{A\}, \{B\}\} : A_{34}^{(2)}, B_{14}^{(2)} + B_{13}^{(2)}, \\ A_{34}^{(2)}, B_{24}^{(2)} + B_{23}^{(2)};$$

For this pattern,  $R_{\mathbf{d}, \mathcal{P}_d^{(t)}} = 6$ ,  $M_{\mathbf{d}, \mathcal{P}_d^{(t)}, k} = 7$ , for  $k \in [1 : 4]$ .

# Partial Decomposition and A New Code Example

► Delivery:

- For demand  $\mathbf{d} = (A,A,A,A)=(1,1,1,1)$ , two decomposition patterns:

1st pattern:  $\mathcal{P}_d^{(t)}$  without any decomposition, let

$$\alpha_{\mathbf{d}, \mathcal{P}_d^{(t)}} = 5/6,$$

$$\begin{aligned} &A_{23}^{(i)} + A_{13}^{(i)} + A_{12}^{(i)}, & A_{24}^{(i)} + A_{14}^{(i)} + A_{12}^{(i)} \\ &A_{34}^{(i)} + A_{14}^{(i)} + A_{13}^{(i)}, & \text{for } i = 1, 2, 3, 4, 5. \end{aligned}$$

For this pattern,  $R_{\mathbf{d}, \mathcal{P}_d^{(t)}} = 3, M_{\mathbf{d}, \mathcal{P}_d^{(t)}, k} = 9, k \in [1 : 4]$ .

# Partial Decomposition and A New Code Example

► Delivery:

- For demand  $\mathbf{d} = (A,A,A,A)=(1,1,1,1)$ , two decomposition patterns:

2nd pattern: special uncoded decomposition pattern  $\check{\mathcal{P}}_{\mathbf{d}}^{(t)}$ , let  $\alpha_{\mathbf{d}, \check{\mathcal{P}}_{\mathbf{d}}^{(t)}} = 1/6$ ,

$$A_{23}^{(6)}, A_{13}^{(6)}, A_{24}^{(6)}, A_{14}^{(6)}, A_{12}^{(6)}, A_{34}^{(6)} \\ B_{23}^{(6)}, B_{13}^{(6)}, B_{24}^{(6)}, B_{14}^{(6)}, B_{12}^{(6)}, B_{34}^{(6)}$$

For this pattern,  $R_{\mathbf{d}, \check{\mathcal{P}}_{\mathbf{d}}^{(t)}} = 12, M_{\mathbf{d}, \check{\mathcal{P}}_{\mathbf{d}}^{(t)}, k} = 3, k \in [1 : 4]$ .

# Partial Decomposition and A New Code Example

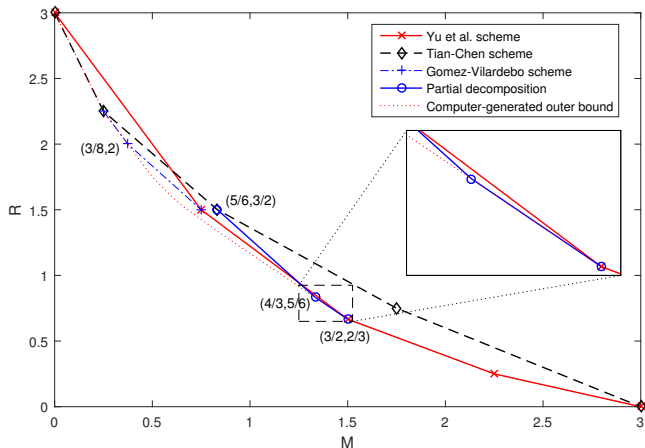
- ▶ The above example uses  $r = 6$  instances of code,
- ▶ Of all 6 instances,  $r_{d, \mathcal{P}_d^{(t)}}$  of them adopt decomposition  $\mathcal{P}_d^{(t)}$ ,
- ▶ Memory-rate pair

$$(M, R) = \frac{1}{r^{(K)}_t} \left( \sum_{\mathcal{P}_d^{(t)}} r_{d, \mathcal{P}_d^{(t)}} M_{d, \mathcal{P}_d^{(t)}, k}, \sum_{\mathcal{P}_d^{(t)}} r_{d, \mathcal{P}_d^{(t)}} R_{d, \mathcal{P}_d^{(t)}} \right)$$

would be

$$(M, R) = (4/3, 5/6).$$

# A New Information-Theoretic Inner Bound



**Figure:** A new information-theoretic inner-bound for the example caching system  $(N, K) = (3, 4)$ ,  $t = 2$ .



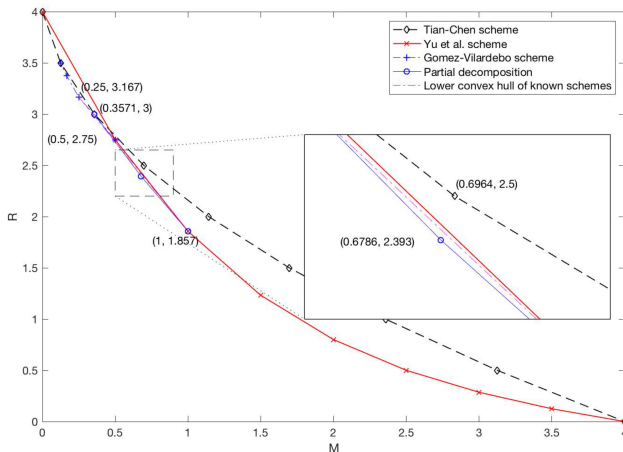
# A New Information-Theoretic Inner Bound

- ▶ For some demands, special uncoded transmission needed to keep the decomposition rule;
- ▶ As the parameters increase, number of decomposition patterns for each  $\mathbf{t}$  grows;
- ▶ The best decompositions for all  $\mathbf{t}$ ? Exhaustive search, practically impossible to calculate by hand;
- ▶ We can use a computer-aided approach.

Before the details, we show a new corner point for caching system  $(N, K) = (4, 8)$ ,  $t = 2$ .

This time it is calculated using a computer program.

# A New Information-Theoretic Inner Bound



**Figure:** A new information-theoretic inner-bound for the example caching system  $(N, K) = (4, 8)$ ,  $t = 2$ .

# A Linear Programming Framework

The inner-bound of the new coding scheme can be calculated using linear programming:

- Define rate region  $\mathcal{R}^{(t)}$ : collection of memory-rate pairs  $(M, R)$  such that exists set of  $\{\alpha_{\mathbf{d}, \mathcal{P}_d^{(t)}}\}$  satisfying

$$\sum_{\mathcal{P}_d^{(t)} \in \mathfrak{P}_{t,d}} \alpha_{\mathbf{d}, \mathcal{P}_d^{(t)}} = 1, \quad \mathbf{d} \in \mathcal{D},$$

$$1 \geq \alpha_{\mathbf{d}, \mathcal{P}_d^{(t)}} \geq 0, \quad \mathbf{d} \in \mathcal{D}, \quad \mathcal{P}_d^{(t)} \in \mathfrak{P}_{t,d},$$

$$\sum_{\mathcal{P}_d^{(t)} \in \mathfrak{P}_{t,d}} \alpha_{\mathbf{d}, \mathcal{P}_d^{(t)}} R_{\mathbf{d}, \mathcal{P}_d^{(t)}} \leq R \binom{K}{t}, \quad \mathbf{d} \in \mathcal{D},$$

$$\sum_{\mathcal{P}_d^{(t)} \in \mathfrak{P}_{t,d}} \alpha_{\mathbf{d}, \mathcal{P}_d^{(t)}} M_{\mathbf{d}, \mathcal{P}_d^{(t)}} \leq M \binom{K}{t}, \quad \mathbf{d} \in \mathcal{D}, k \in [1 : K].$$

# A Linear Programming Framework

The inner-bound of the new coding scheme can be calculated using linear programming:

- Further define

$$(M'_d, R'_d) = \frac{1}{r_t^{(K)}} \left( \max_{k \in [1:K]} \sum_{\mathcal{P}_d^{(t)}} r_{d, \mathcal{P}_d^{(t)}} M_{d, \mathcal{P}_d^{(t)}, k}, \sum_{\mathcal{P}_d^{(t)}} r_{d, \mathcal{P}_d^{(t)}} R_{d, \mathcal{P}_d^{(t)}} \right)$$

and

$$M'_r \triangleq \max_{d \in \mathcal{D}} M'_d, \quad R'_r \triangleq \max_{d \in \mathcal{D}} R'_d.$$

# A Linear Programming Approach

The inner-bound of the new coding scheme can be calculated using linear programming:

- ▶ The number of instances  $r$  can be chosen s.t. exists  $\{r_{\mathbf{d}, \mathcal{P}_d^{(t)}}\}$

$$\left| \frac{r_{\mathbf{d}, \mathcal{P}_d^{(t)}}}{r} - \alpha_{\mathbf{d}, \mathcal{P}_d^{(t)}} \right| \leq \epsilon.$$

By making  $r$  large and choosing appropriate  $\{r_{\mathbf{d}, \mathcal{P}_d^{(t)}}\}$  s.t.  $\epsilon \geq 0$  being arbitrarily small. We have

$$\lim_{r \rightarrow \infty} (M'_r, R'_r) = (M, R),$$

it will be the effective memory-rate pair of the new code.

# Conclusion

- ▶ A single scheme unifying two general classes of schemes (uncoded and coded),
- ▶ The Yu et al scheme and Tian-Chen scheme are two extreme points of the new scheme,
- ▶ The notion of Transmission type plays an important role and is overlooked before this work,
- ▶ The performance improvement is not quite large, and mostly reside in the middle  $M$  regime, which are not surprise.

# Conclusion

- ▶ The new coding scheme can not incorporate the coding scheme of Gómez-Vilardebó<sup>1</sup>, a more generalized code may exist?
- ▶ The complexity: certain decomposing patterns are obviously bad choice; Simplifying the proposed scheme appears worthwhile.

---

<sup>1</sup>J. Gómez-Vilardebó, "Fundamental limits of caching: Improved bounds with coded prefetching," *arXiv:1612.09071*, Dec. 2016.