Motivation
Single Server Caching
FemtoCaching Problem
Cooperative MultiPoint Systems
CoMP Caching Policies
Conclusion

# Caching Policies for Delay Minimization in Small Cell Networks with Joint Transmissions

Guilherme I. Ricardo [1,2]

Giovanni Neglia [2]    Thrasyvoulos Spyropoulos [1]

[1]EURECOM, Communication Systems Department

[2]Inria, Université Côte d'Azur

Motivation
Single Server Caching
FemtoCaching Problem
Cooperative MultiPoint Systems
CoMP Caching Policies
Conclusion

# Agenda

1. **Motivation**

2. **Single Server Caching**

3. **FemtoCaching Problem**

4. **Cooperative MultiPoint Systems**

5. **CoMP Caching Policies**

6. **Conclusion**

Motivation
Single Server Caching
FemtoCaching Problem
Cooperative MultiPoint Systems
CoMP Caching Policies
Conclusion

# Motivation
## Content Distribution Networks

- Scenario : Increasing mobile and cellular data usage.
- Question : How to provide better QoS under such scenario ?
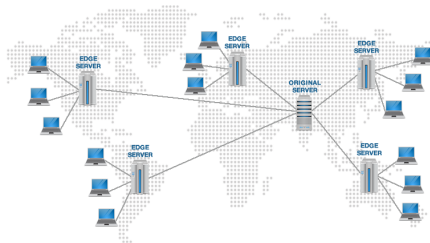- Solution : Content replication closer to final user - Caching !



FIGURE – CDN Multiserver Caching Strategy – Source

Motivation
Single Server Caching
FemtoCaching Problem
Cooperative MultiPoint Systems
CoMP Caching Policies
Conclusion

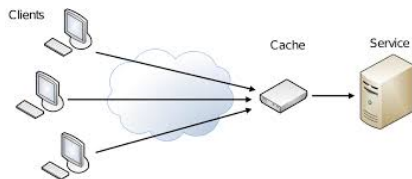# Single Server Caching

Introduction



FIGURE – Single Server Caching – Source

- Problem : What to cache ?
- Performance metric : Hit Ratio
- Popularity is known : Store the most popular contents
- Popularity is unknown/dynamic : Caching algorithms (policies)

Motivation
Single Server Caching
FemtoCaching Problem
Cooperative MultiPoint Systems
CoMP Caching Policies
Conclusion

# Single Server Caching

Policies Examples - Least Frequently Used (LFU)



FIGURE – LFU Caching Policy – Source

Motivation
Single Server Caching
FemtoCaching Problem
Cooperative MultiPoint Systems
CoMP Caching Policies
Conclusion

# Single Server Caching
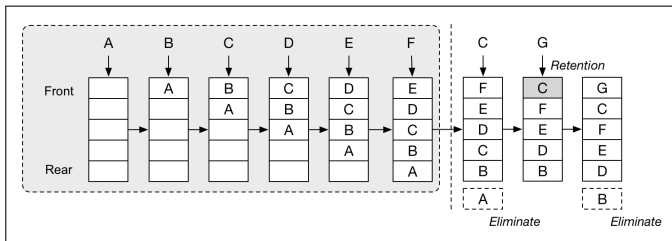## Policies Examples - Least Recently Used (LRU)



FIGURE – LRU Caching Policy – Source

- Variations :
    - $q$LRU – probabilistic insertion, $0 \leq q \leq 1$
    - $k$LRU – multilevel cache, $k = 1, 2, ...$

Motivation
Single Server Caching
FemtoCaching Problem
Cooperative MultiPoint Systems
CoMP Caching Policies
Conclusion

# FemtoCaching Problem

## 5G Heterogeneous Networks Topology



FIGURE – Heterogeneous Network Topology – Source

Motivation
Single Server Caching
FemtoCaching Problem
Cooperative MultiPoint Systems
CoMP Caching Policies
Conclusion

## FemtoCaching Problem
The Optimization Formulation

Let $\mathbf{X}$ be the allocation matrix such that $x_{hf} = 1$ if helper $h$ caches content $f$ and $x_{hf} = 0$ otherwise. The problem is :

$$\underset{x}{\text{maximize}} \quad F(\mathbf{X}) = \frac{1}{U} \sum_{f=1}^{F} p_f \sum_{u=1}^{U} \mathbb{1}_{\{k(u,f)>0\}}$$

$$\text{subject to} \quad \sum_{f=1}^{F} x_{hf} = C, \ h = 1, \ldots, H,$$

where $F$ is the catalog size, $C$ is the cache capacity, $U$ is the number of users, $k(u,f) \triangleq \sum_{h \in \mathcal{H}(u)} x_{hf}$, and $\mathcal{H}(u)$ is the set of helpers covering user $u$.

Motivation
Single Server Caching
FemtoCaching Problem
Cooperative MultiPoint Systems
CoMP Caching Policies
Conclusion

## FemtoCaching Problem
The Offline Solution – Femto (2015)

- NP-Hard Problem (Combinatorial Nature)
- Greedy Algorithm :
    - $F(\mathbf{X})$ is monotone and submodular
    - Constraints form a matroid partition
    - 1/2-Approximation ratio
- Drawbacks : Strong assumptions, e.g.,
    - Centralized intelligence
    - Network topology and popularities are static and known

Motivation
Single Server Caching
FemtoCaching Problem
Cooperative MultiPoint Systems
CoMP Caching Policies
Conclusion

# FemtoCaching Problem
The Online Solutions – Caching Policies

- LRU-One and LRU-All – Giovanidis (2016)
- $q$LRU-Lazy – Neglia (2018)

> ### $q$LRU-Lazy Policy Description
>
> 1. Only the helper that served the file can update its cache; and
>
> 2. It only does so if it is the only one able to actually serve it

Motivation
Single Server Caching
FemtoCaching Problem
Cooperative MultiPoint Systems
CoMP Caching Policies
Conclusion

# Cooperative Multipoint Systems

Joint Transmissions and The Delay Metric

### Definition

The delay $d(u, f, \mathbf{X})$ for user $u$ to download content $f$ under allocation $\mathbf{X}$ is

$$d(u, f, \mathbf{X}) = \begin{cases} d_B + \frac{M}{W \log_2(1+\max_h g_{hu})}, & \text{if cache miss} \\ \frac{M}{W \log_2(1+\sum_h g_{hu} x_{hf})}, & \text{if cache hit,} \end{cases}$$

where $d_B$ is the backhaul delay, $g_{hu}$ is the SNR from $h$ to $u$, $M$ is the file size, and $W$ is the channel bandwidth.

Motivation
Single Server Caching
FemtoCaching Problem
Cooperative MultiPoint Systems
CoMP Caching Policies
Conclusion

# Cooperative Multipoint Systems

The Optimization Problem : Formulation – Tuholukova (2018)

---

### Delay Minimization Problem

$$\underset{x}{\text{minimize}} \quad F(\mathbf{X}) = \frac{1}{U} \sum_{f=1}^{F} p_f \sum_{u=1}^{U} d(u, f, \mathbf{X})$$

$$\text{subject to} \quad \sum_{f=1}^{F} x_{hf} = C, \ h = 1, \dots, H$$

---

### Remark

Submodularity Proof and Greedy Algorithm

Motivation
Single Server Caching
FemtoCaching Problem
Cooperative MultiPoint Systems
CoMP Caching Policies
Conclusion

# Cooperative Multipoint Systems

Hit Ratio → Avg. Delay



FIGURE – Static allocation for different overlapping levels (full rep.)

Motivation
Single Server Caching
FemtoCaching Problem
Cooperative MultiPoint Systems
CoMP Caching Policies
Conclusion

# Cooperative Multipoint Systems

Optimal Allocation : $d_B$ x SNR Bounds

Assumptions :

- Completely overlap
- Homogeneous SNR $(\gamma)$

For a given $\gamma$, if $d_B \geq d_{B,max}$ such that

$$d_{B,\max}(C, H, \alpha, \gamma) \triangleq (HC)^\alpha \frac{M}{W}(\frac{1}{\log_2(1 + \gamma)} - \frac{1}{\log_2(1 + 2\gamma)})$$

then the optimal allocation is full diversity.
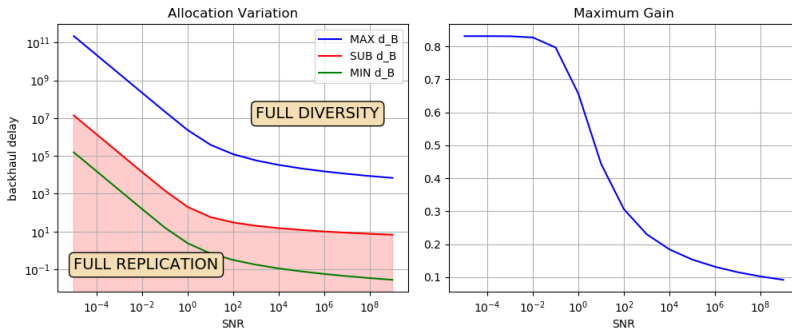For a given $\gamma$, if $d_B \leq d_{B,min}$ such that

$$d_{B,\min}(C, H, \alpha, \gamma) \triangleq \left(\frac{C+1}{C}\right)^\alpha \frac{M}{W} \left(\frac{1}{\log_2(1 + (H-1)\gamma)} - \frac{1}{\log_2(1 + H\gamma)}\right)$$

then the optimal allocation is full replication.

Motivation
Single Server Caching
FemtoCaching Problem
Cooperative MultiPoint Systems
CoMP Caching Policies
Conclusion

# Cooperative Multipoint Systems

## Optimal Allocation : $d_B$ x SNR Bounds, Example

Motivation
Single Server Caching
FemtoCaching Problem
Cooperative MultiPoint Systems
CoMP Caching Policies
Conclusion

## CoMP Caching Algorithms
$q$LRU-$\Delta d$ Policy Notation

- Let $I_u$ be the set of helpers covering user $u$ and $J_{u,f} \subseteq I_u$ be the subset of those helpers caching $f$.

- The marginal gain for adding a copy of file $f$ at helper $h$ is defined as :

$$\Delta d^{(h)}(u, f, \mathbf{X}) \triangleq d(u, f, \mathbf{X} \ominus \mathbf{e}^{(h)}) - d(u, f, \mathbf{X})$$

- Normalizers :

$$\beta \triangleq 1/(\max_{f,h,u,\mathbf{X}} \Delta d^{(h)}(u, f, \mathbf{X}))$$

$$\gamma \triangleq 1/(\max_{f,h,u,\mathbf{X}} \Delta d^{(h)}(u, f, \mathbf{X} \oplus \mathbf{e}^{(h)})).$$

Motivation
Single Server Caching
FemtoCaching Problem
Cooperative MultiPoint Systems
CoMP Caching Policies
Conclusion

# CoMP Caching Algorithms

$q$LRU-$\Delta d$ Policy Introduction

---

### $q$LRU-$\Delta d$ Policy General Description

At every request $(u, f)$, each $h \in I_u$ updates its cache as follows :

- If $h \in J_{u,f}$, reset $f$'s cache position with probability :

$$\rho^{(h)}(u, f, \mathbf{X}) = \beta \cdot \Delta d^{(h)}(u, f, \mathbf{X})$$

- If $h \in I_u \backslash J_{u,f}$, store $f$ to $h$'s cache with probability $q \cdot \sigma^{(h)}(u, f, \mathbf{X})$, where $q \in (0, 1]$ is fixed and

$$\sigma^{(h)}(u, f, \mathbf{X}) = \gamma \cdot \Delta d^{(h)}(u, f, \mathbf{X} \oplus \mathbf{e}^{(h)})$$

Motivation
Single Server Caching
FemtoCaching Problem
Cooperative MultiPoint Systems
CoMP Caching Policies
Conclusion

# CoMP Caching Algorithms

$q$LRU-$\Delta d$ Policy Introduction

---

### $q$LRU-$\Delta d$ Policy Algorithmic Description
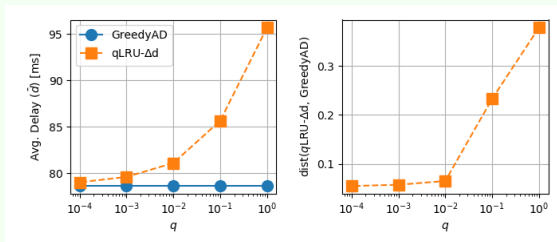
**Input:** $I_u, J_{u,f}$, and $g_{h',u}, \forall h' \in I_u$
**for** $h \in I_u$ **do**
    **if** $h \in J_{u,f}$ **then**
        Move $f$ to the front with prob. $\rho^{(h)}$
    **else**
        Evict file in the cache's last position;
        Insert $f$ with prob. $q \cdot \sigma^{(h)}(u, \mathbf{X}_f)$.
    **end**
**end**

Motivation
Single Server Caching
FemtoCaching Problem
Cooperative MultiPoint Systems
CoMP Caching Policies
Conclusion

# CoMP Caching Algorithms

$q$LRU-$\Delta d$ Policy Introduction

## Remark – Ricardo (2020)

Under IRM, Che's, and Exponentialization approximations, a network of $q$LRU-$\Delta d$ caches converges to a locally-optimal caching configuration when $q \to 0$.

Motivation
Single Server Caching
FemtoCaching Problem
Cooperative MultiPoint Systems
CoMP Caching Policies
Conclusion

# CoMP Caching Algorithms
2LRU-$\Delta d$ Policy Notation

- IRM $\neq$ Real request process (Temporal locality)
- Each helper deploys a 2-levels cache : the physical cache storing the actual file and the virtual cache storing files' metadata (i.e., ID)
- Let $I_u$ be the set of helpers covering user $u$ and let $J_{u,f}, \hat{J}_{u,f} \subseteq I_u$ be the subsets of those helpers storing $f$ at the physical cache and at the virtual cache, respectively.

Motivation
Single Server Caching
FemtoCaching Problem
Cooperative MultiPoint Systems
CoMP Caching Policies
Conclusion

# CoMP Caching Algorithms

2LRU-$\Delta d$ Policy Introduction

---

## 2LRU-$\Delta d$ Policy General Description

At every request $(u, f)$, each $h \in I_u$ updates its cache as follows :

- If $h \in \hat{J}_{u,f}$, move $f$'s ID to the front of $h$'s virtual cache and,
  - if $h \in J_{u,f}$, move $f$ to the front of $h$'s physical cache with prob. $\rho^{(h)}(u, f, \mathbf{X})$ ;
  - else, evict the file in the physical cache's last position and insert $f$.
- If $h \notin \hat{J}_{u,f}$, with prob. $q \cdot \sigma^{(h)}(u, f, \mathbf{X})$, evict the ID in $h$'s virtual cache's last position and insert $f$'s ID
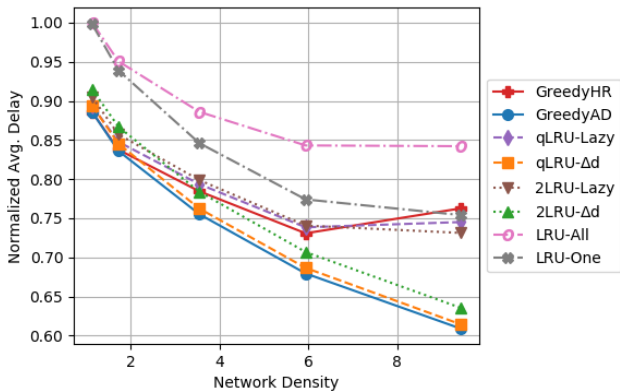
---

Motivation
Single Server Caching
FemtoCaching Problem
Cooperative MultiPoint Systems
CoMP Caching Policies
Conclusion

# CoMP Caching Algorithms

## 2LRU-$\Delta d$ Policy Introduction

---

**2LRU-$\Delta d$ Policy Algorithmic Description**

---

**Input:** $I_u, J_{u,f}, \hat{J}_{u,f}$, and $g_{h',u}, \forall h' \in I_u$

**for** $h \in I_u$ **do**

    **if** $h \in \hat{J}_{u,f}$ **then**

        Move $f$'s ID to the front of the virtual cache;

        **if** $h \in J_{u,f}$ **then**

            Move $f$ to the front of the physical cache
            with prob. $\rho^{(h)}$

        **else**

            Evict file in physical cache's last position;
            Insert $f$.

        **end**

    **else**

        Evict file's ID in virtual cache's last position;
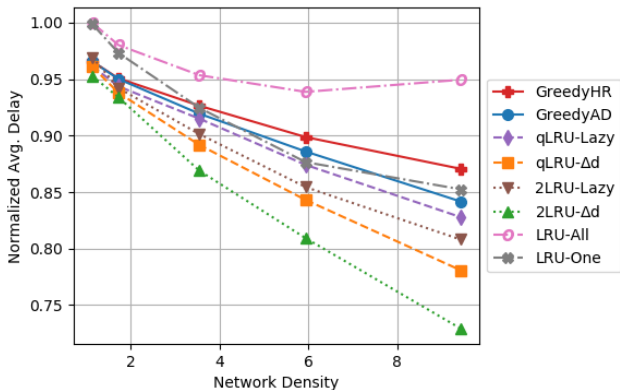        Insert $f$'s ID with prob. $\sigma^{(h)}$.

    **end**

**end**

---

Motivation
Single Server Caching
FemtoCaching Problem
Cooperative MultiPoint Systems
CoMP Caching Policies
Conclusion

# CoMP Caching Algorithms

## Numerical Results – IRM, Homogeneous SNR

Motivation
Single Server Caching
FemtoCaching Problem
Cooperative MultiPoint Systems
CoMP Caching Policies
Conclusion

# CoMP Caching Algorithms

Numerical Results – Real, Homogeneous SNR

Motivation
Single Server Caching
FemtoCaching Problem
Cooperative MultiPoint Systems
CoMP Caching Policies
Conclusion

# CoMP Caching Algorithms

Numerical Results – Real, Heterogeneous SNR

Motivation
Single Server Caching
FemtoCaching Problem
Cooperative MultiPoint Systems
CoMP Caching Policies
Conclusion

## Conclusion and Future Works

- Conclusions
  - Delay cost function under CoMP provides different allocation with potentially better download rates
  - $q$LRU-$\Delta d$ Policy outperforms other Hit Ratio dynamic policies for synthetic requests
- Future Work
  - Finish Real Traces Experiments
  - Greedy Algorithm with pair of files
  - Finish Algorithm

Thank You!