



Prediction-based Redundant Data Elimination with Content Overhearing in Wireless Networks

Haiying Shent[†], Shenghua He^{*}, Lei Yu[‡] and **Ankur Sarker[†]**

[†]Department of Computer Science, University of Virginia

^{*}Department of Computer Science and Engineering, Washington University in St. Louis

[‡]College of Computing, Georgia Institute of Technology

Outline

- Introduction
- System Design
- Performance Evaluation
- Conclusion

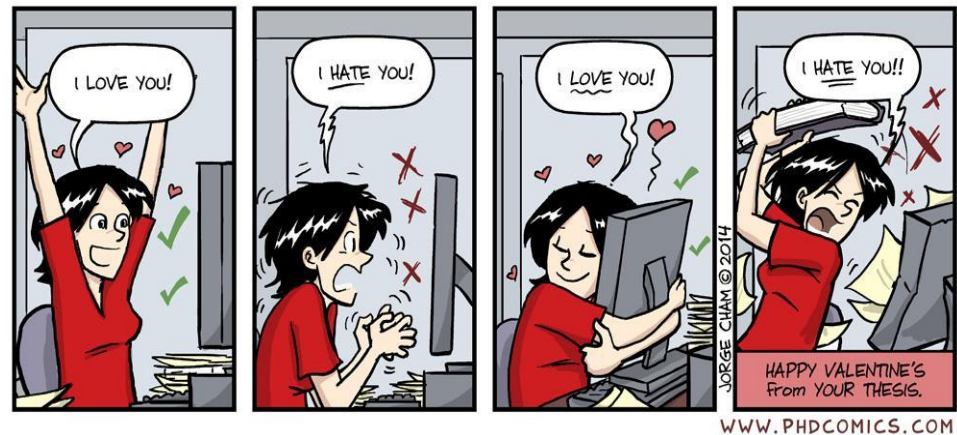
Outline

- **Introduction**
- System Design
- Performance Evaluation
- Conclusion

Introduction

Internet Traffic

- ❑ How Netflix, YouTube, Hulu, and Amazon became the Internet.
- ❑ File sharing is another major source.
- ❑ 30% of the Internet is just a **Copy** of itself.



Redundancy

Introduction

Redundant data elimination

Suppressing duplicated data transmission using **Redundancy Elimination** techniques:

1. Packet-based RE
2. Content-based RE



Redundancy Elimination

Related work

Wired network

1. WAN optimization
 - ❑ [Riverbed Networks 2013] [Juniper Network 2014]
2. Server to client RE
 - ❑ [Agarwal et al., NSDI 2010] [Hua et al., Infocom 2014]
3. Prediction-based RE
 - ❑ [Zohar et al., Sigcomm 2011] [Yu et al., ICNP 2012]

Wireless network

1. Packet-based RE
 - ❑ [Hua et al., Infocom 2015] [Sanadhya et al., Mobicom 2012]
2. Content-based RE
 - ❑ [Dogar et al., Mobicom 2008] [Afanasyev et al., NSDI 2008]

Introduction

Problems in previous methods

1. Caches at sender and receiver would be outdated
 - ❑ Disrupts RE's correctness and degrade its performance
2. Overhearing probability estimation is difficult
 - ❑ Consequently degrades the performance of RE
 - ❑ Causes significant communication cost and complex coordination among nodes

Solution: Prediction-based Redundancy Elimination

1. The receiver stores the received and overheard data stream in a chain of chunks.
2. It compares the chunks of the incoming packet with the stored chunk chains in the cache.
3. The receiver sends to the sender future data predictions that include the hashes of chunks on the chain.

Introduction

Challenges of RE

1. Identifying duplicate chunks of hundreds of bytes at sub-packet level and work on lower-bandwidth wireless links.
2. Transmission cost of predictions has to be considered.
3. There are possibilities of missing data.

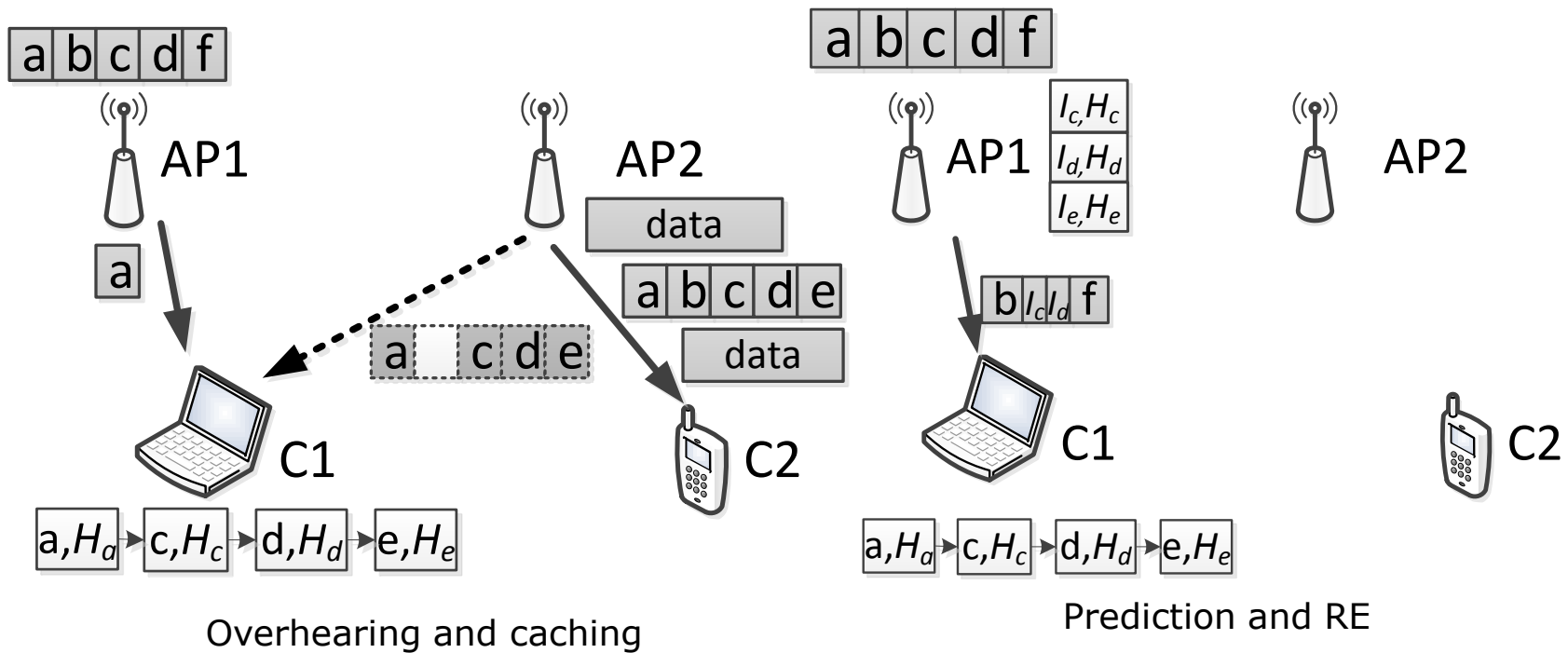
Our method: PRECO

1. Redundancy elimination technique
 - ❑ An effective, efficient and scalable solution for content-overhearing based IP-layer RE over wireless links.
2. Chunking and caching
 - ❑ To divide the payload into evenly distributed chain of chunks.
3. Adaptive prediction algorithm
 - ❑ To improve prediction accuracy and reduce the prediction overhead.

Outline

- Introduction
- **System Design**
- Performance Evaluation
- Conclusion

System Design Overview



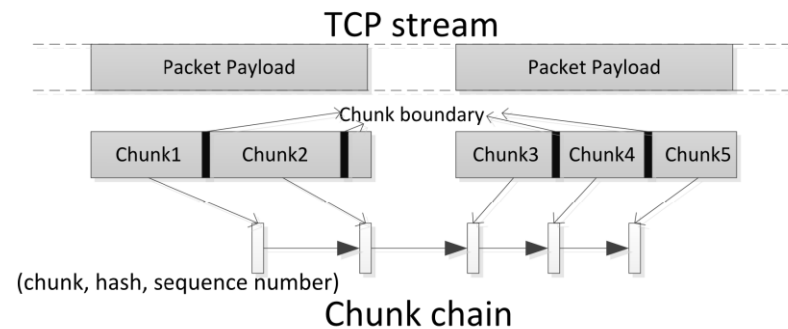
System Design

Chunking and caching

The overall procedures divide into two steps:

- ❑ Chunking algorithm: PRECO divides the payload of a packet into several chunks. MAXP [1], a content-based chunking algorithm is used to define chunk boundary.

- ❑ Caching received and overheard chunks: Nodes overheard TCP streams and stored as chunks based on stream IDs (src, dst, src port, dst port).



[1] Anand, Ashok, et al. "Redundancy in network traffic: findings and implications." *ACM SIGMETRICS Performance Evaluation Review* 37.1 (2009).

System Design

Prediction-based RE

Prediction of receive packets works as follow:

- ❑ PRECO determines one matching chunk as a prediction anchor.
- ❑ Chooses a chain of chunks for prediction based on highest matching length.
- ❑ Prediction chunk is chosen based on the matching degree.
- ❑ Virtual chunk is created based on matching degree.
- ❑ For a received packet, one chain of chunks is selected as prediction based on virtual chunk.

Prediction transmission and shim decoding:

- ❑ A prediction windows is used to increase the efficiency.
- ❑ Receiver sends the chunk prediction in a prediction message.
- ❑ Upon receiving, the sender stores the prediction in cache.
- ❑ For an outgoing packet, the sender performs chunking using same algorithm and insert shim into the packet.
- ❑ Once receiving a packet containing shim, receiver finds the shim from sender.

System Design

Adaptive prediction algorithm

Size of the prediction window, W

$$W = \begin{cases} (1 + R(P_A))W_0, & \text{if } N_B - (N_A + |P_A|) < d_T \\ W_0, & \text{if } N_B - (N_A + |P_A|) \geq d_T \end{cases}$$

where

P_A is prediction chunk chain based on Anchor A

$R(P_A)$ is hit ratio of prediction PA

W_0 is initial prediction window size

N_A is the next expected byte sequence number based on A

N_B is the next expected byte sequence number based on B

d_T prediction distance threshold

System Design

Adaptive prediction algorithm

Size of the prediction window, W

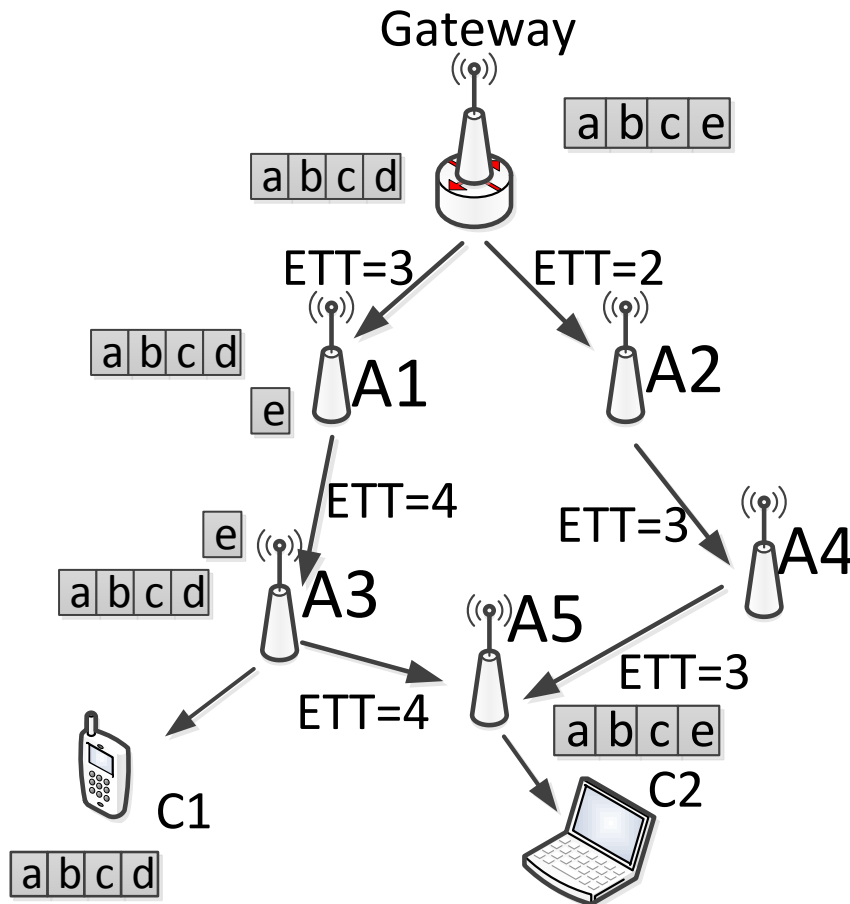
$$W = \begin{cases} (1 + R(P_A))W_o, & \text{if } N_B - (N_A + |P_A|) < d_T \\ W_o, & \text{if } N_B - (N_A + |P_A|) \geq d_T \end{cases}$$

When making predictions with a new prediction anchor:

- ❑ the algorithm first computes the hit ratio of previous prediction, and then
- ❑ accordingly adjusts the prediction window.

System Design

Redundancy-aware source routing



ETT Estimation:

- Using "A1-A3-A5": $3+3+16 = 23$
- Using "A2-A4-A5": $8+12+12 = 32$

System Design

Redundancy-aware source routing

The redundancy-aware source routing:

- ❑ Routing metric: Redundancy Estimated Transmission Time

$$RETT = ETX \times \frac{s(1-\alpha)}{B}$$

where, ETX is Expected Transmission Count [1]

S is average packet size

B is the bandwidth

α is average redundancy ratio

- ❑ Routing Protocol:
 - ❑ Compute RETT matrix for all links
 - ❑ Apply Dijkstra's shortest path algorithm to find path route with lowest RETT
 - ❑ Without any overhearing consideration

[1] De Couto, Douglas SJ, et al. "A high-throughput path metric for multi-hop wireless routing." *Wireless Networks* 11.4 (2005).

Outline

- Introduction
- System Design
- **Performance Evaluation**
- Conclusion

Experiment

Simulation settings

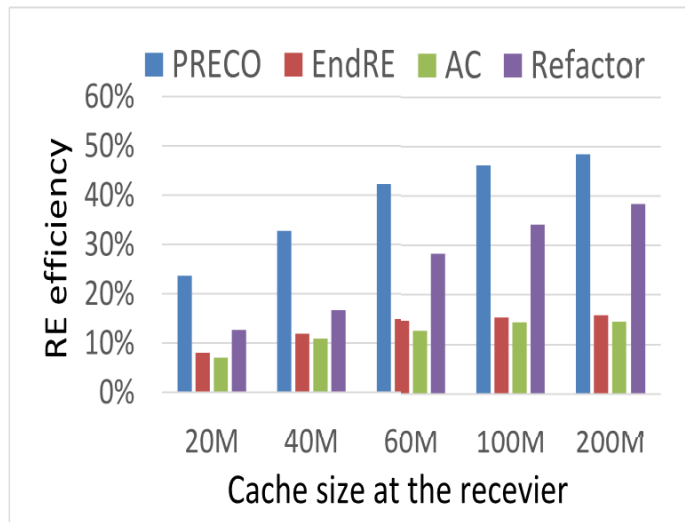
1. Real trace collected using YouTube app over 3.5 GB data
 - Two smartphones (Iphone 6 and Xiao Mi 3) connected with one laptop (Lenovo T420 Windows 10 machine)
 - 2 different videos with similar contents, watched twice
 - 60 minutes a day for 7 day
 - Captured packets using Wireshark
2. 2 scenarios-
 - One AP and one client without content overhearing
 - Two Aps and two clients, with content overhearing

Compared methods

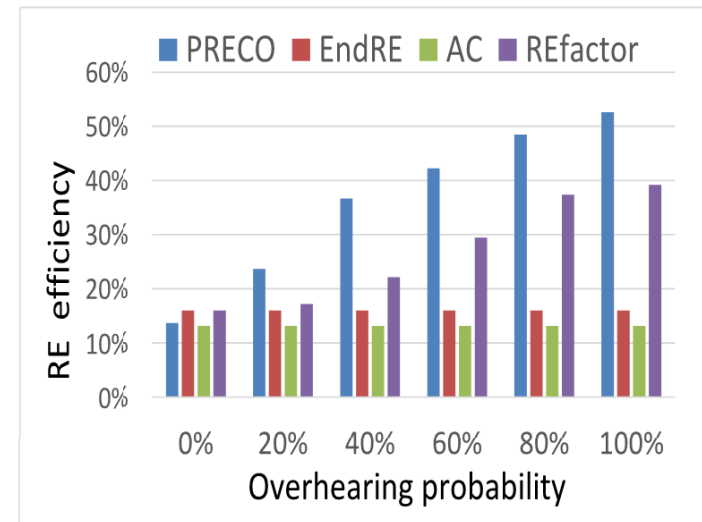
1. EndRE: new finger print technique for end users
2. Asymmetric Caching (AC): RE operations based on feedback cache
3. REfactor: finer-granularity redundancy at the sub-packet level with content overhearing

Experiment

RE efficiency



(a) Different receiver's cache size



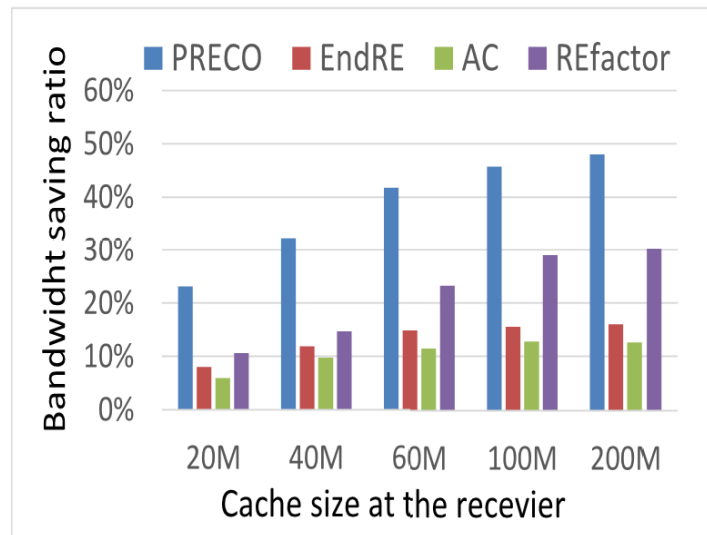
(b) Different overhearing probability

Observation: RE efficiency follows PRECO>Refactor>EndRE>AC

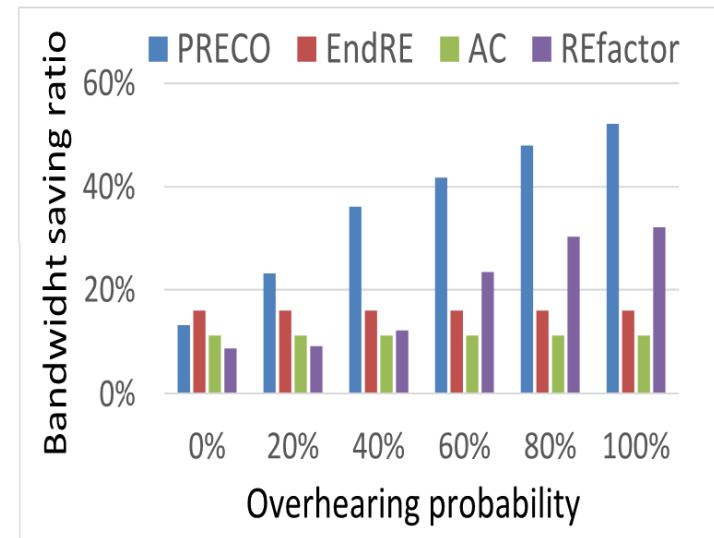
Reason: EndRE and AC do not support overhearing; Refactor uses overhearing probability estimation.

Experiment

Content overhearing



(a) Different receiver's cache size



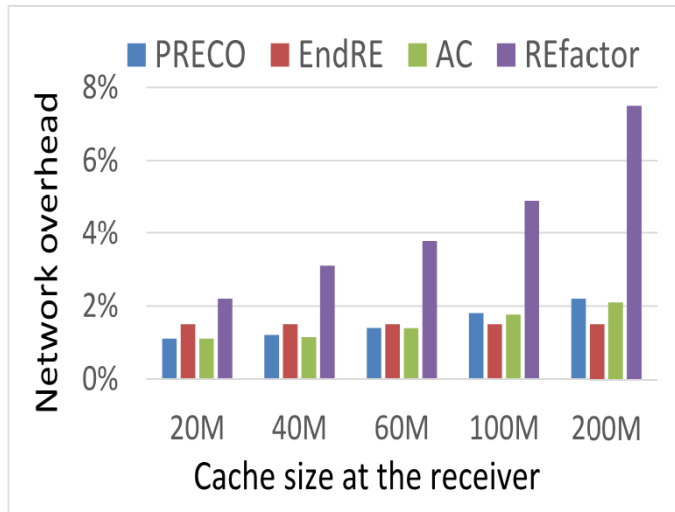
(b) Different overhearing probability

Observation: PRECO>REfactor>EndRE>AC

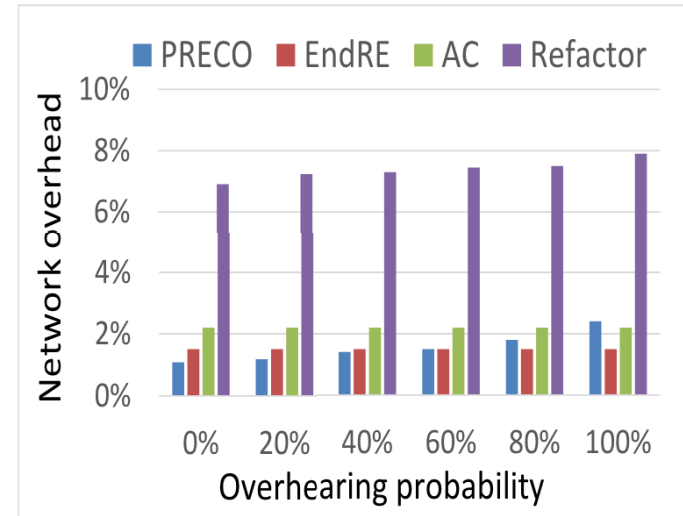
Reason: The bandwidth saving is mainly caused by RE efficiency

Experiment

Network overhead



(a) Different receiver's cache size



(b) Different overhearing probability

Observation: Refactor has the highest network overhead among all these RE methods

Reason: The overhearing probability estimation results larger amount of network overhead

Experiment

Simulation settings

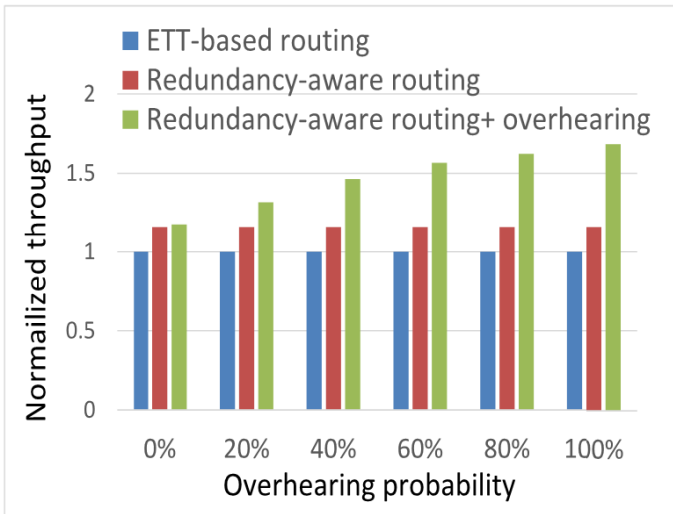
1. we deployed a mesh network with 5 rows and 5 columns in total
 1. First node is the gateway
 2. Three clients associated in three distinct nodes
 3. the overhearing coverage to 1
 4. The average data rate for each link varies from 800Kps to 1200Kps
2. The gateway send traces to two different clients

Compared methods

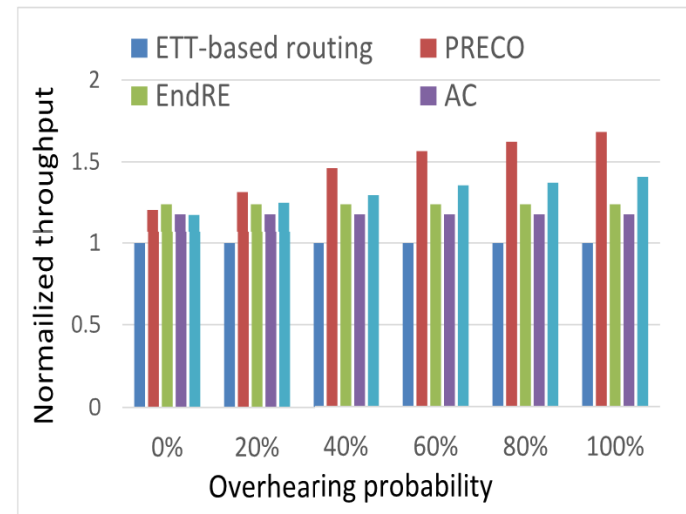
1. ETT-based routing: The gateway determines the optimal route to a receiver using ETT metric, and no network-wide PRECO deployment to perform RE.
2. Redundancy-aware routing without content overhearing: network-wide PRECO is deployed without content overhearing.
3. Redundancy-aware routing with content overhearing: network-wide PRECO with content overhearing.

Experiment

Redundancy-aware routing



(a) Different overhearing probability



(b) Different overhearing probability

Observation: Redundancy-aware routing with content overhearing produces more throughput

Reason: The gateway steers the traffic through the nodes with high redundancy

Outline

- Introduction
- System Design
- Performance Evaluation
- **Conclusion**

Conclusions

1. we propose a prediction-based IP-layer RE method with content overheard named PRECO for wireless networks.
2. We propose novel prediction algorithms that allow PRECO to effectively improve prediction accuracy and overall bandwidth saving.
3. Trace-driven simulation results show that PRECO provides significant performance benefits in comparison with other RE methods.

Future work

Further take into account efficiently learn the overhead data streams of all nodes for route determination in mesh networks.

Thank you!
Questions & Comments?

Ankur Sarker

as4mz@Virginia.edu

Ph.D. Candidate

Pervasive Communication Laboratory

University of Virginia