



**Universit  Pierre et Marie Curie – Sorbonne Universit s
Paris, France**

LIP6/UPMC/CNRS/INRIA
Orange Labs/France Telecom

Predicting Popularity and Adapting Replication of Internet Videos for High-Quality Delivery

(with Hermes)

by

**Guthemberg Silvestre, S bastien Monnet, David Buffoni,
and Pierre Sens**

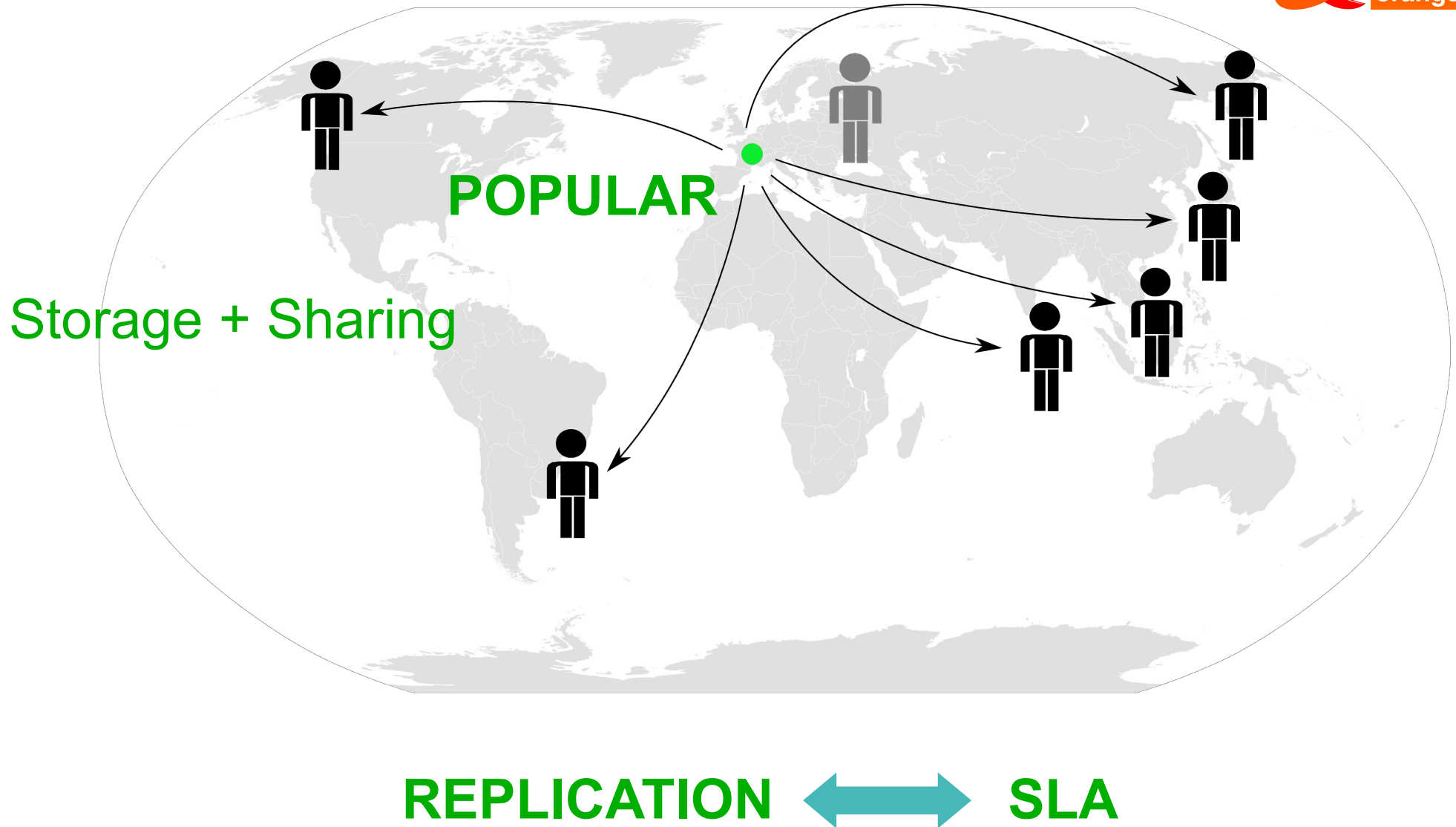
CIFRE partner:



Orange Labs

ICPADS 2013, Seoul, South Korea
18 December 2013

Context



State-of-the-art

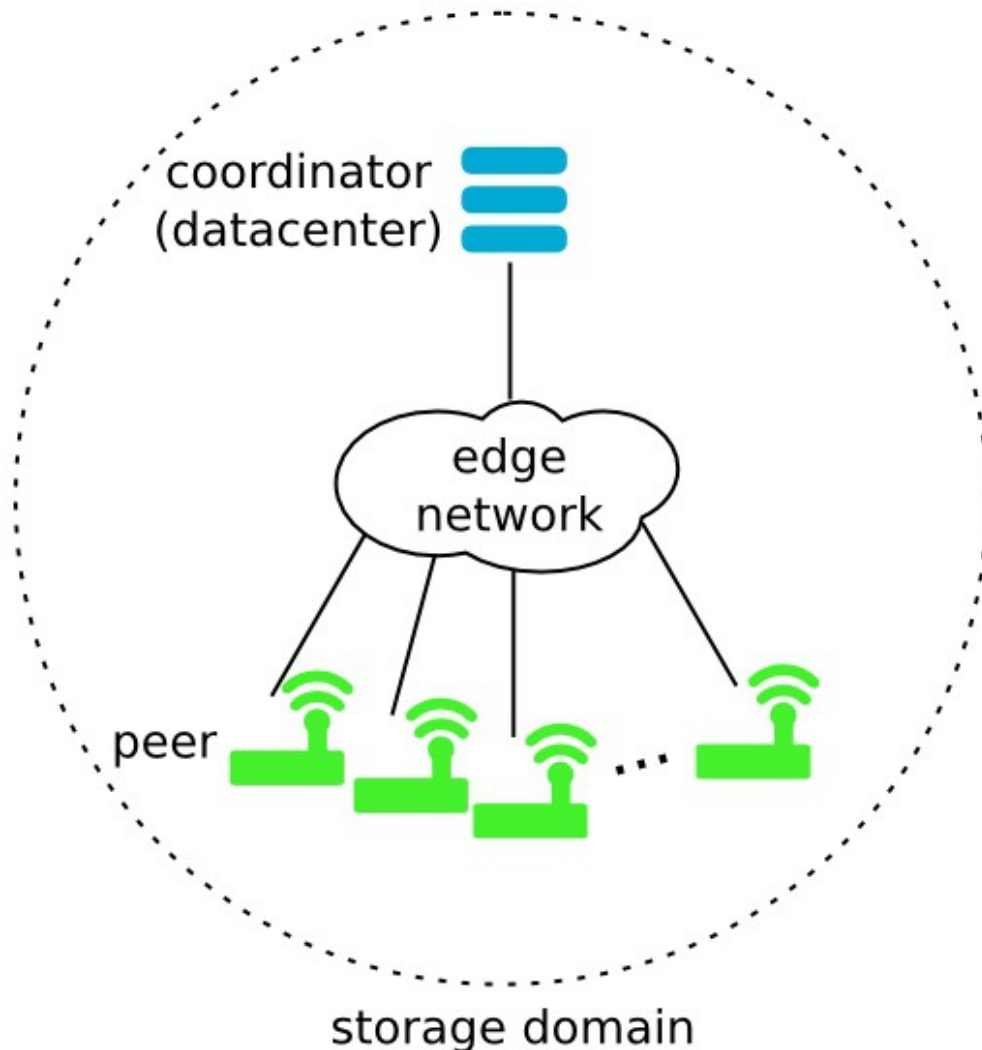
- CDN architectures
 - Infrastructure-based (datacenter): Akamai, Google
 - Hybrid design: NaDa [FP7'11], Echos[SIGCOMM'08]
- SLA
 - Poor content availability: uptime (Amazon S3)
 - Very high content availability: deadline-aware approaches (D3[SIGCOMM'11])
- Content replication
 - Uniform and fixed
 - Adaptive: non-collaborative cache(LRU), EAD [IEEE TPDS'10], Skute [ACM SoCC'10]

Open issues

- Where do we place clients' objects?
- How do we handle edge network devices for object-based storage systems?
- **How many replicas per object should the system create?**
- **How could we prevent SLA violations and optimize edge resources utilization?**

Caju: a content delivery system for edge networks

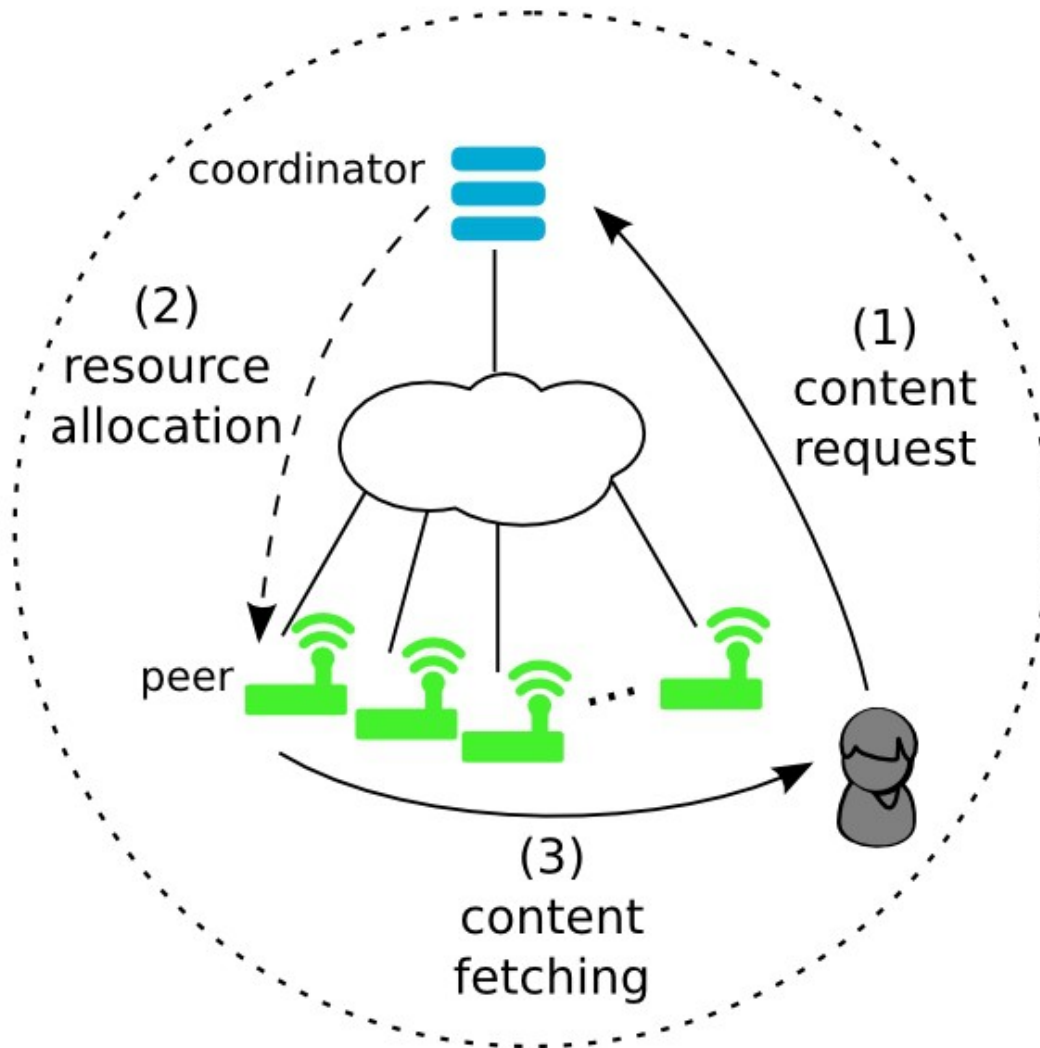
[EUROPAR/BDMC'12]



- Hybrid architecture organized in storage domains
- Two classes of devices: coordinator and peer
- P2P communication, chunks, multi-sourcing ...

Caju: a content delivery system for edge networks

[EUROPAR/BDMC'12]



- Customers are connected to the system through a peer, and assigned to a storage domain
- According to peers' requests (Create, Read, Delete), system performs Replication properly
- Service Level Agreements

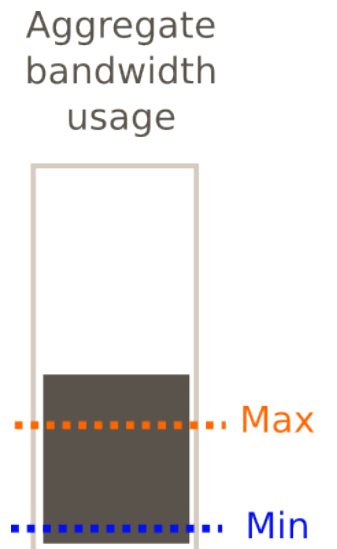
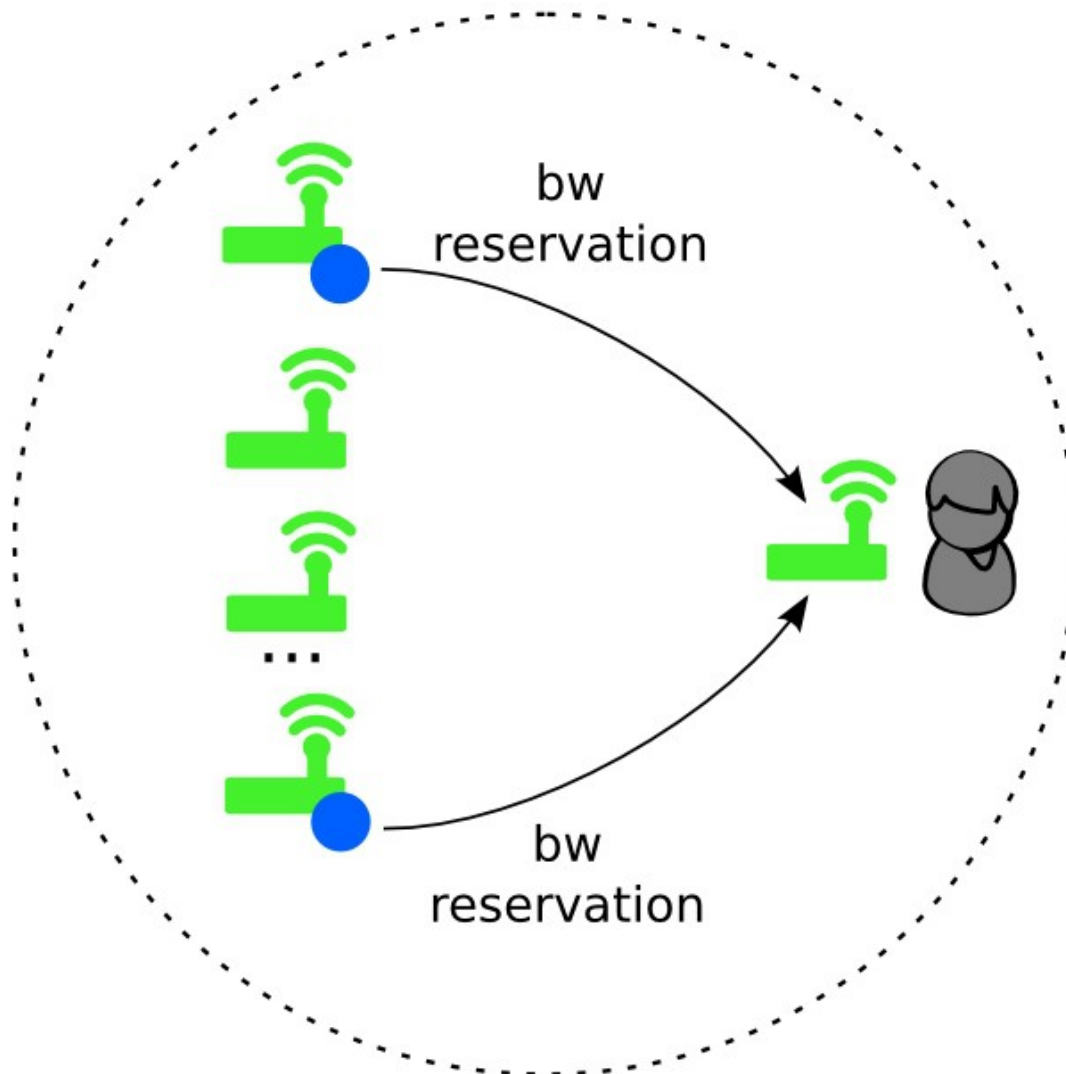
AREN: Adaptive replication scheme for edge networks

[ICPADS'12]

- The performance goals are twofold:
 - (1) prevent SLA violations
 - (2) reduce the usage of edge resources
- Combine bandwidth reservation and collaborative caching using thresholds

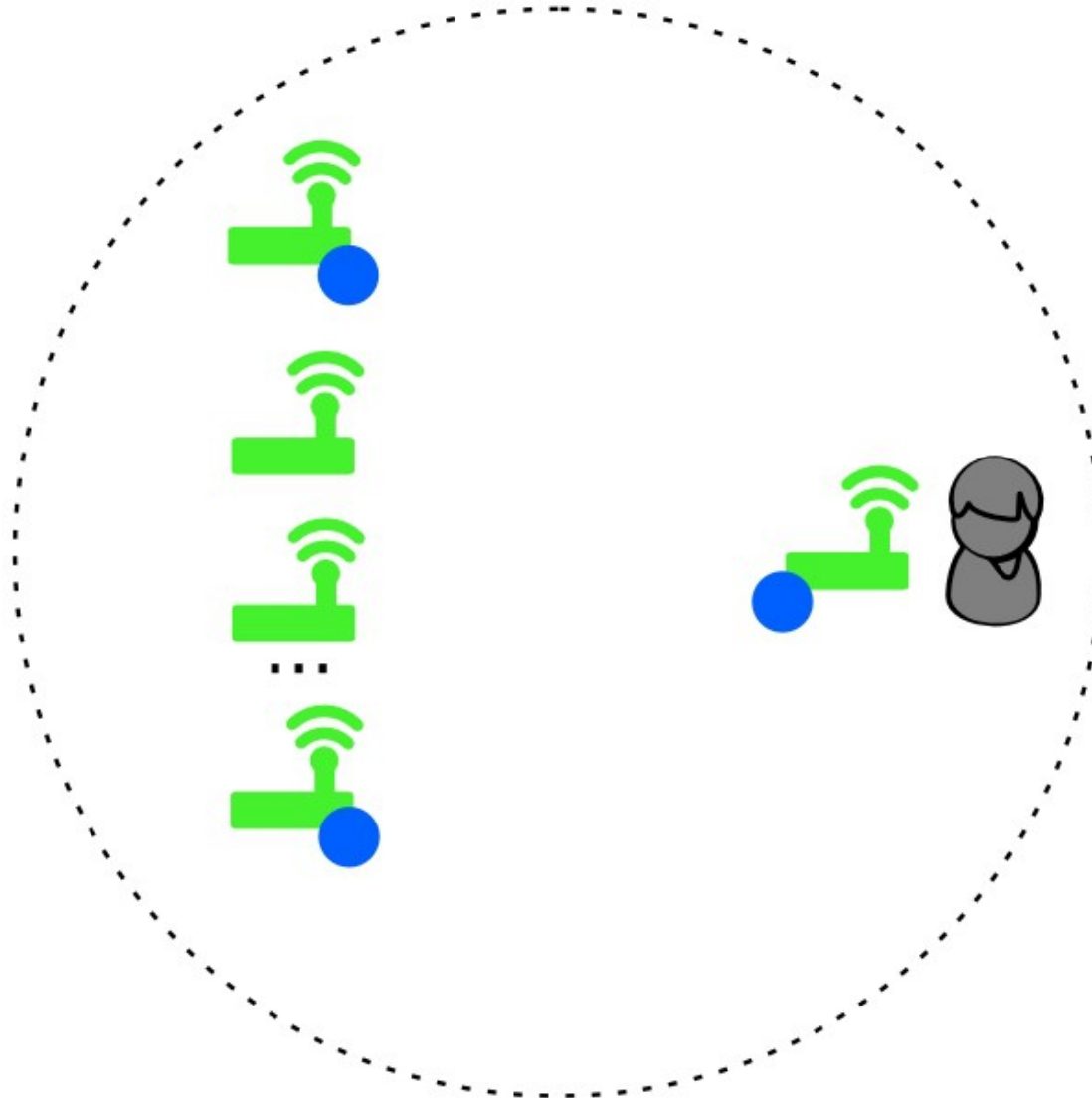
AREN: bandwidth reservation

[ICPADS'12]



AREN: collaborative cache

[ICPADS'12]



AREN replication scheme

- AREN, a novel adaptive replication scheme for cloud storage in edge networks
 - Enforces strict SLA contracts efficiently
 - Improves resource allocation
- AREN tracks bandwidth reservation on edge nodes for operating collaborative caching mechanism

AREN replication scheme

- AREN, a novel adaptive replication scheme for cloud storage in edge networks
 - Enforces strict SLA contracts efficiently
 - Improves resource allocation
- AREN tracks bandwidth reservation on edge nodes for operating collaborative caching mechanism

Issues: AREN relies on bandwidth reservation and thresholds

AREN replication scheme

- AREN, a novel adaptive replication scheme for cloud storage in edge networks
 - Enforces strict SLA contracts efficiently
 - Improves resource allocation
- AREN tracks bandwidth reservation on edge nodes for operating collaborative caching mechanism

Issues: AREN relies on bandwidth reservation and thresholds

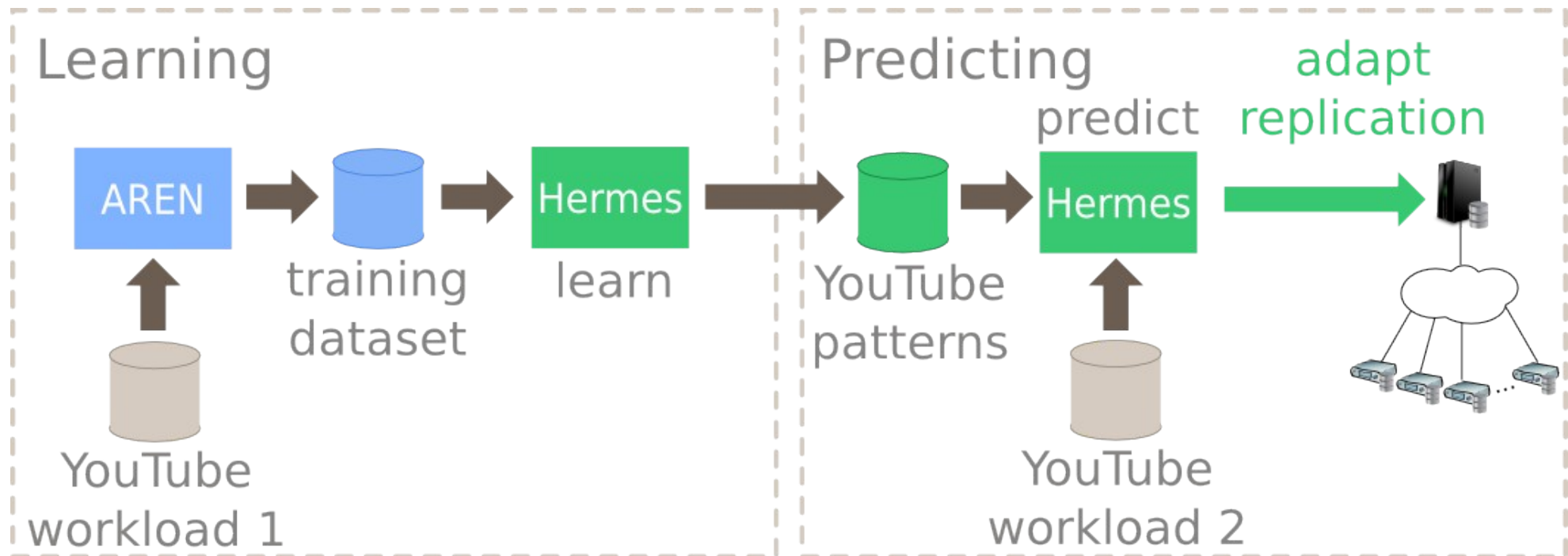
We cope with these issues by making predictions about Internet content demand for adapting replication accordingly

Hermes: predicting popularity and adapting replication

[ICPADS'13]

- Hermes is an adaptive replication scheme based on accurate predictions about the popularity of Internet videos
- Two-step approach (for each view request):
 - (1) Popularity classifier for distinguishing between non-popular and popular contents
 - (2) Replication classifier for maintaining replication of popular content: ***increasing, decreasing, and keeping***

Predicting Content Demand



Framework

Predicting Content Demand

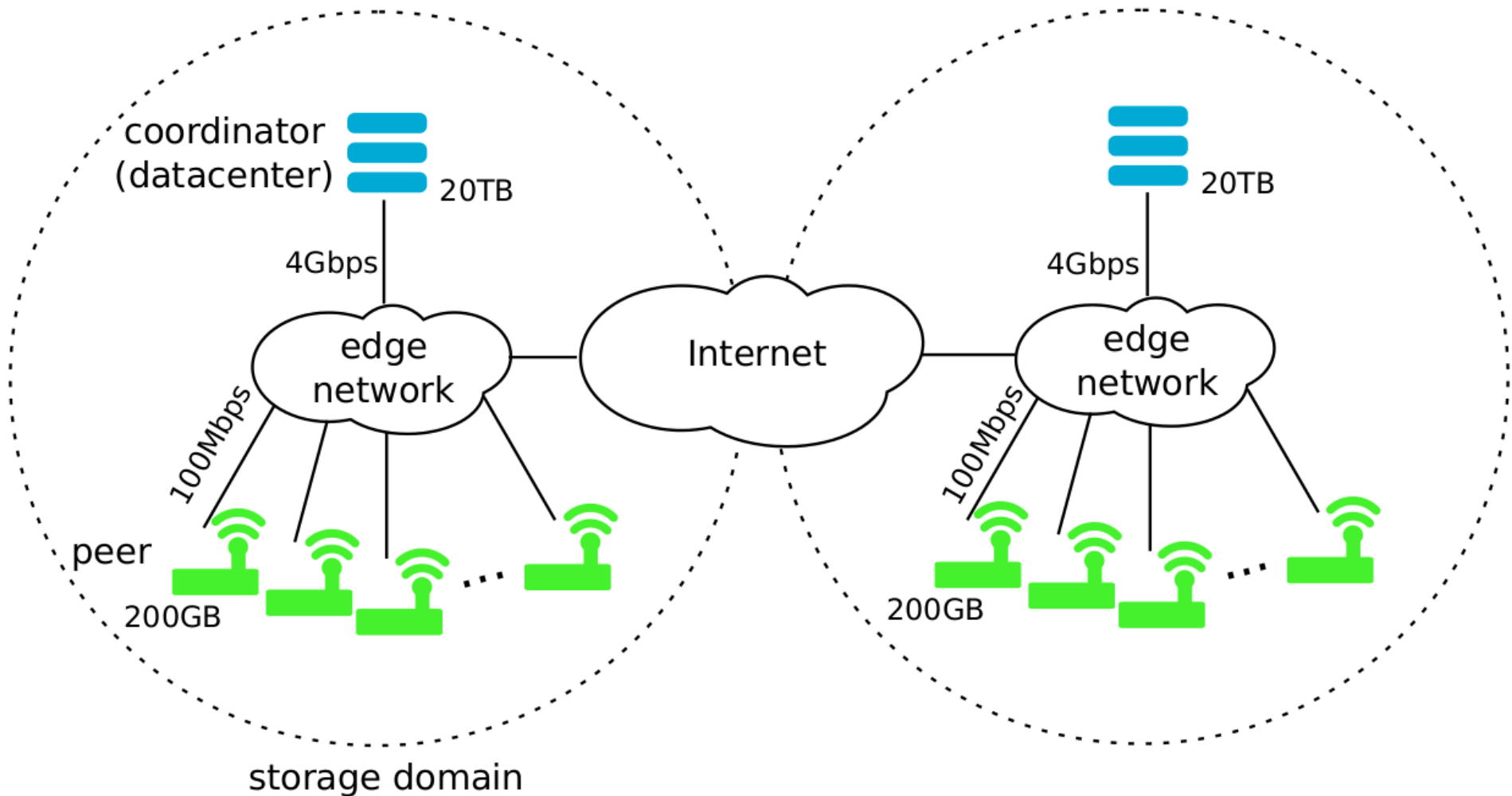
Measurements

- Predictions from 10 measurements of request arrival process:
 - content size
 - network availability
 - network usage (load)
 - # of consumers
 - # of replicas
 - inter-arrival time between requests (delta)
 - aggregate number of downloads
 - mean of time between requests (mtbr)
 - life time
 - average bandwidth

Evaluation Scenario

- PeerSim component: deadline-aware transport mechanism based on data flow
[MOSPAS/HPCS'13]
[PeerSim user code (<http://peersim.sourceforge.net/#code>)]
- Metrics: SLA violations, storage, and network usage
- Compare to: non-collaborative caching and AREN

Evaluation Scenario

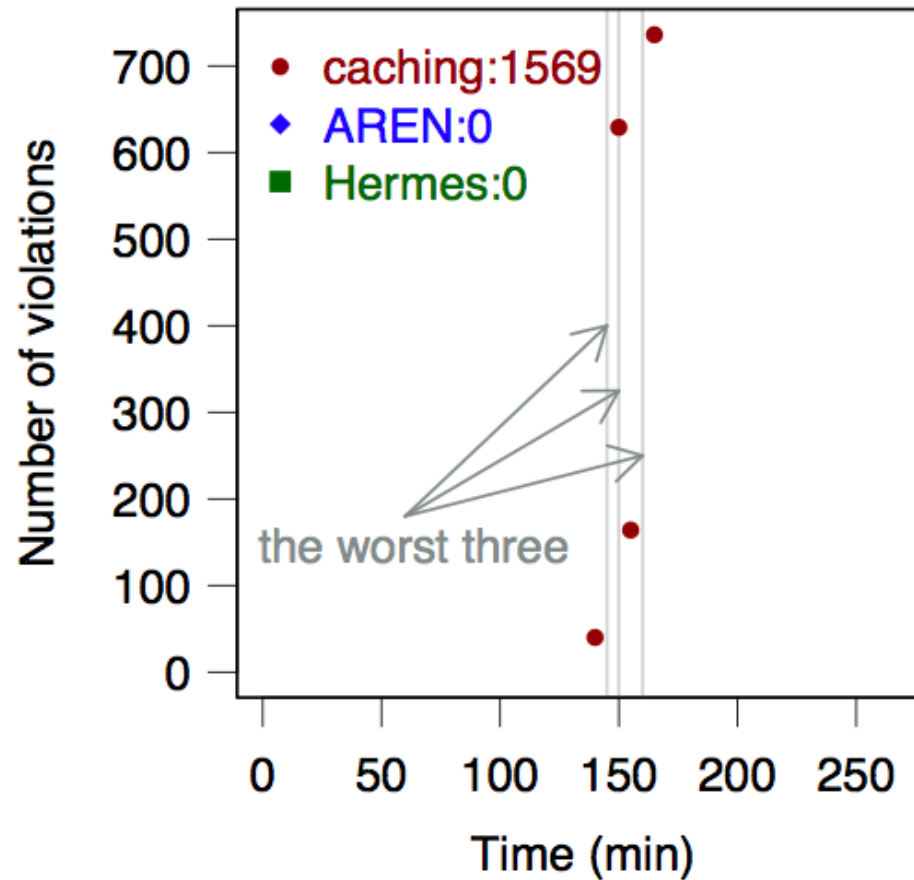


Evaluation Scenario

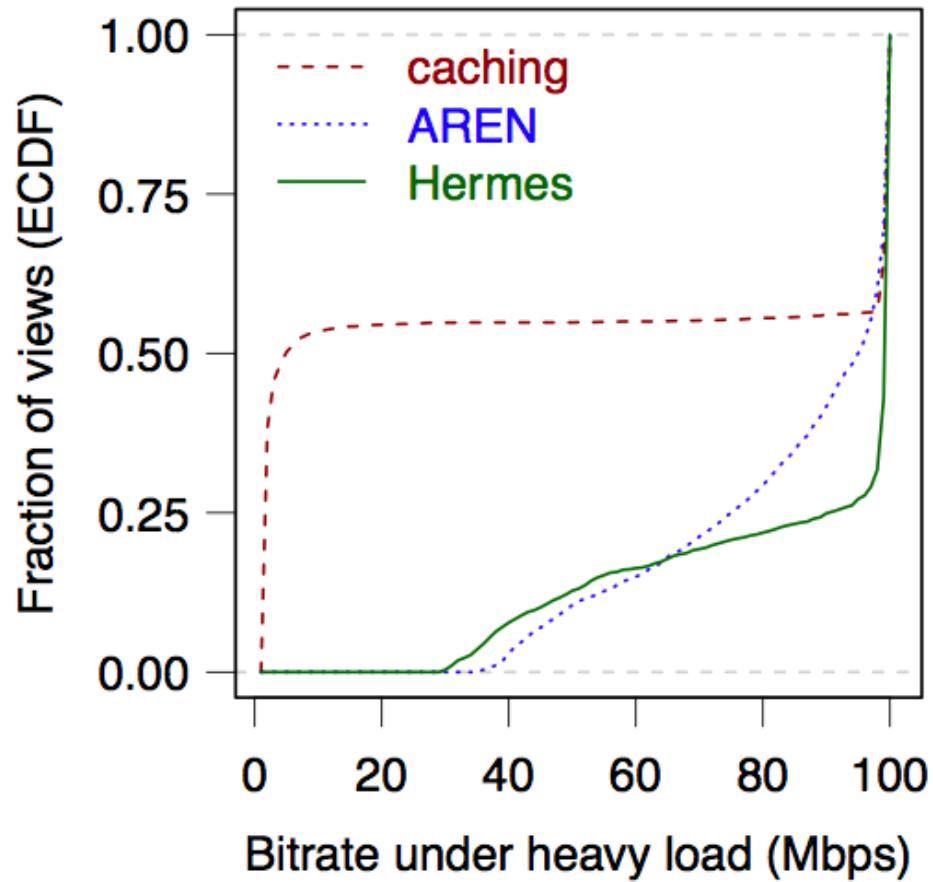
Workload	
Requests per user	uniform
Experiment duration	4 hours
Mean requests per second	100
Requests fractions	5% of creations, 95% of views
Object size (follows Pareto)	shape=3, from 13MB to 1.6GB
Video popularity (Zipf-Mandelbrot)	shape=0.8, cutoff=# of videos
Videos' creation (Poisson)	λ =creations per second
Popularity growth from YouTube traces	21827 distinct patterns

- YouTube traces [Figueiredo et al., WSDM'11]
- SLA definitions for highly available contents
 - Customer-oriented bitrate metric (28MB/s)

Hermes: preventing SLA violations

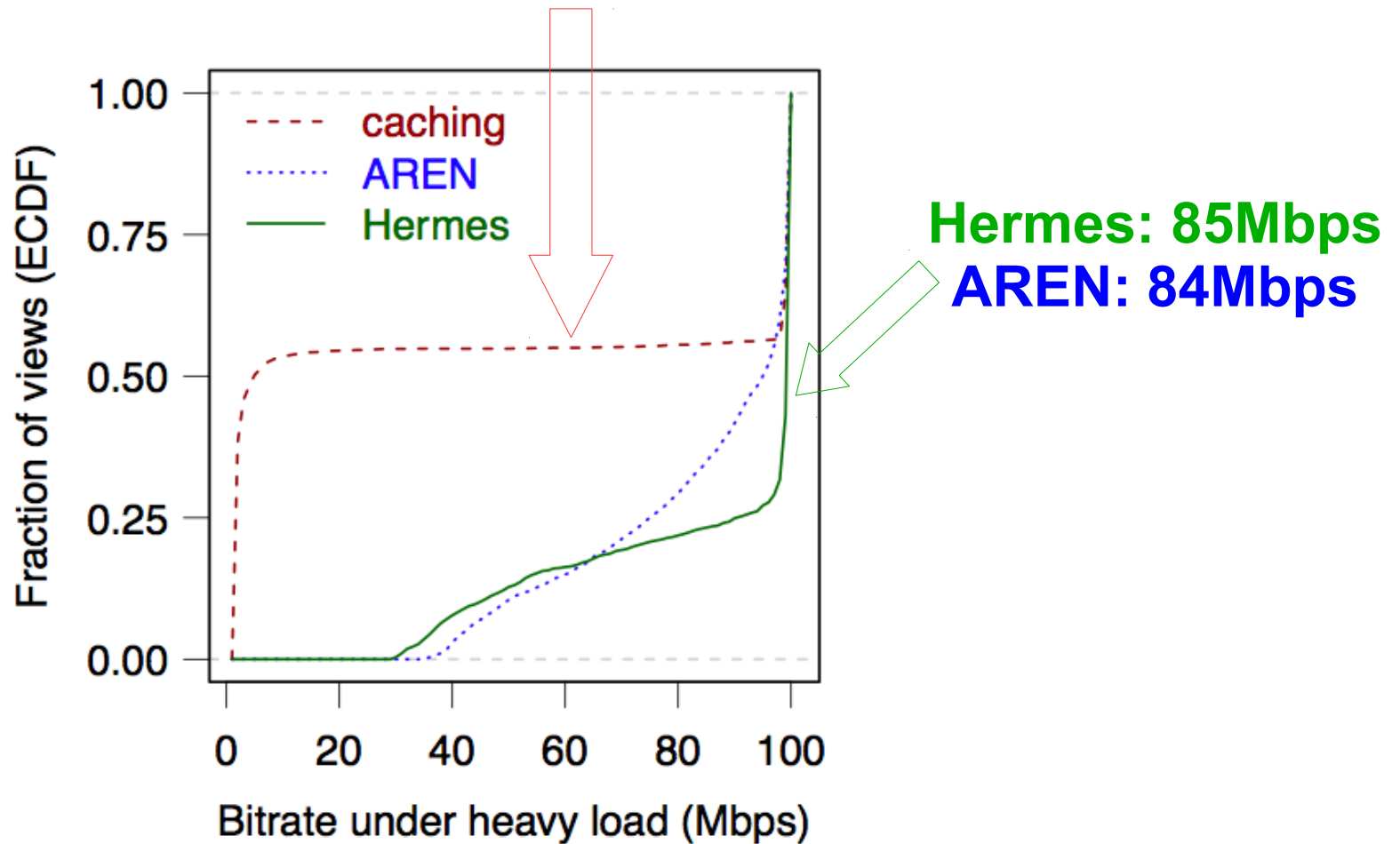


Hermes: improving bitrate provision

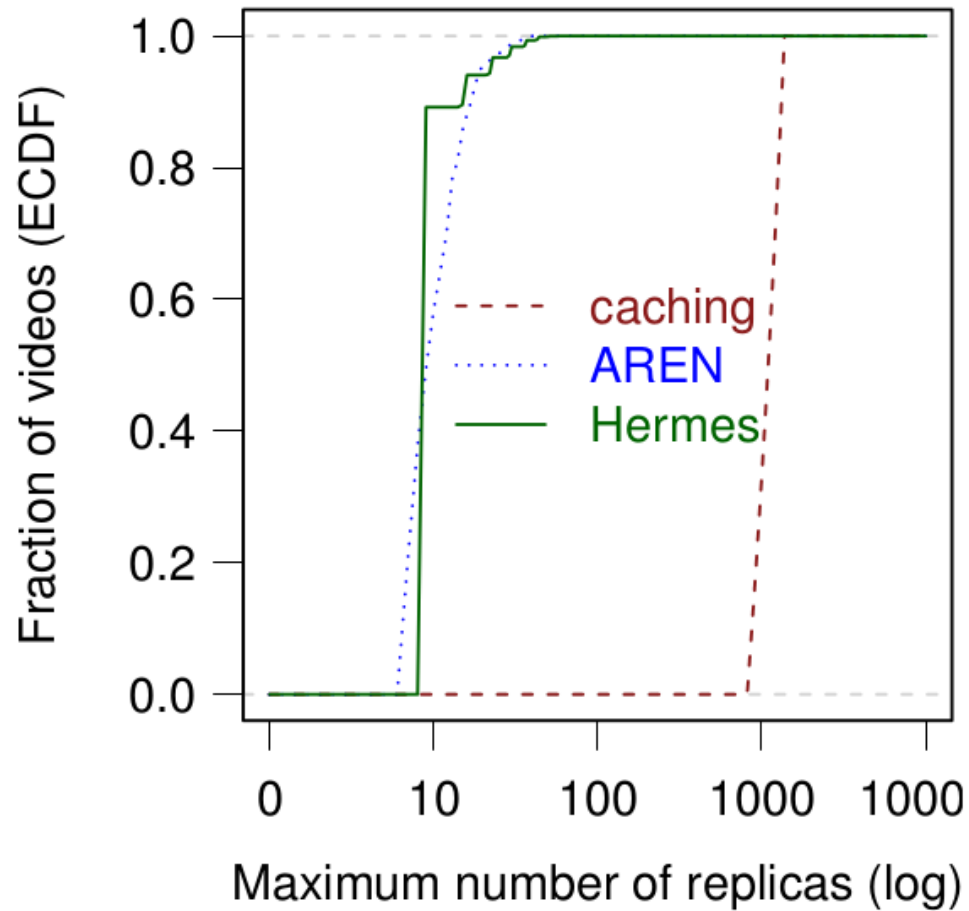


Hermes: improving bitrate provision

caching: 45Mbps

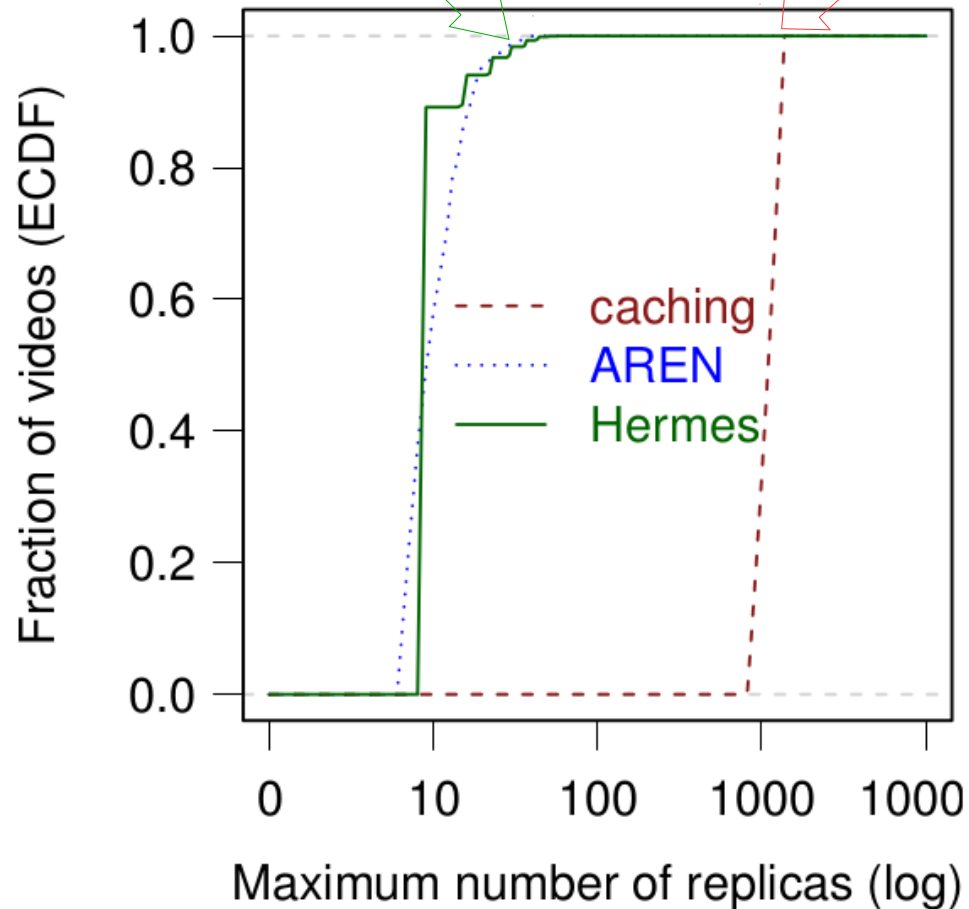


Hermes: reducing the number of replicas

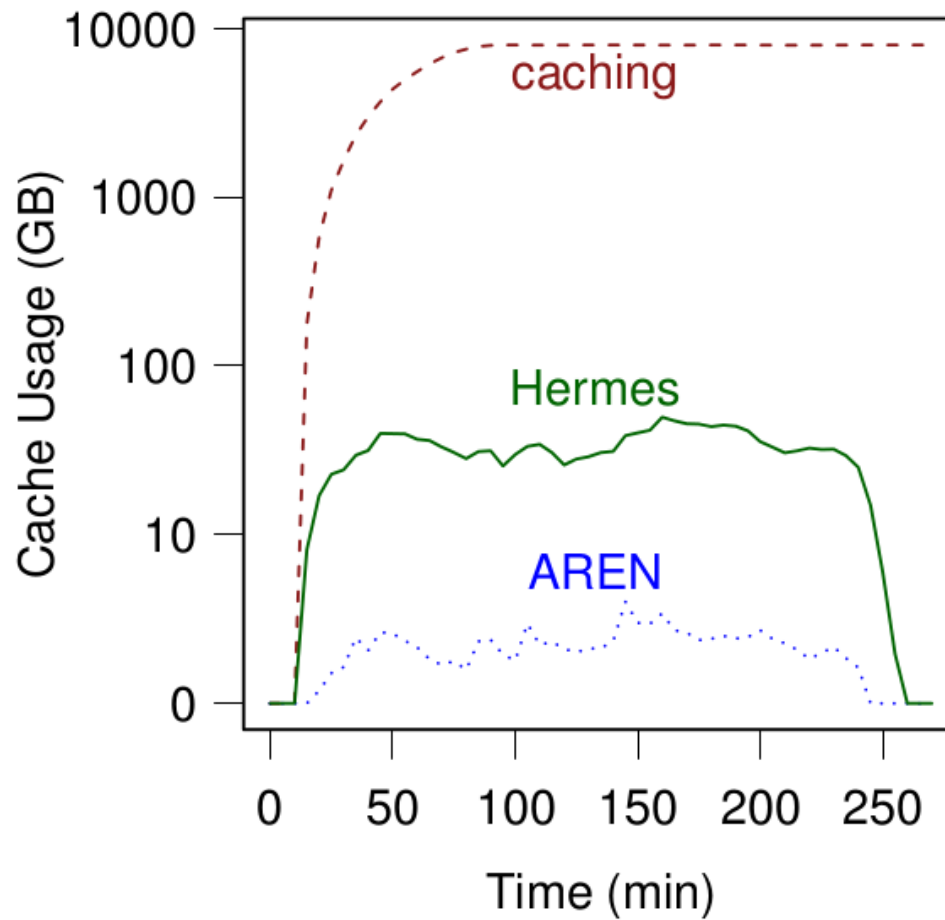


Hermes: reducing the number of replicas

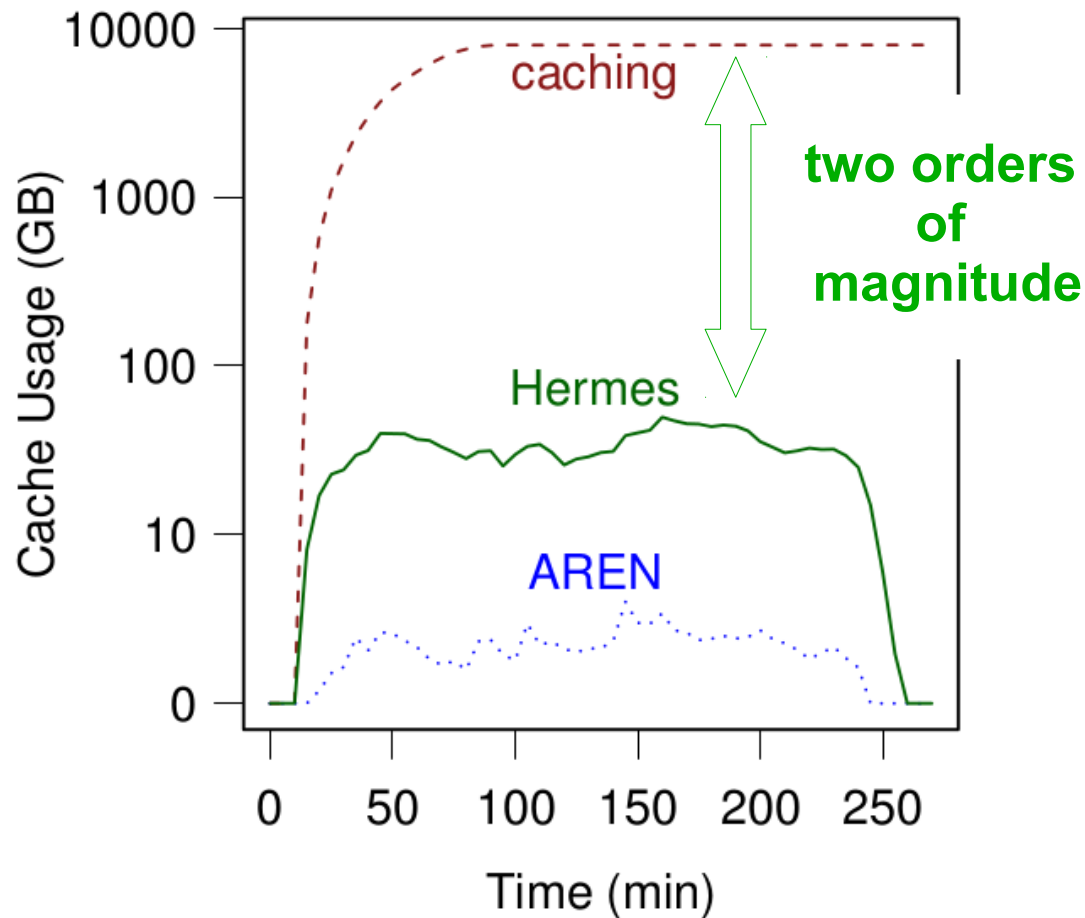
AREN: 39
Hermes: 58 **caching: 1377**



Hermes: saving storage



Hermes: saving storage



Final considerations

- Hermes, an adaptive replication for distributing videos in content delivery networks
 - Adapts content replication based on popularity predictions
 - Prevent SLA violations
 - Improves resource allocation
- Future work
 - Evaluate our scheme through a proof-of-concept prototype

For further information:
<http://guthemberg.co.nr>

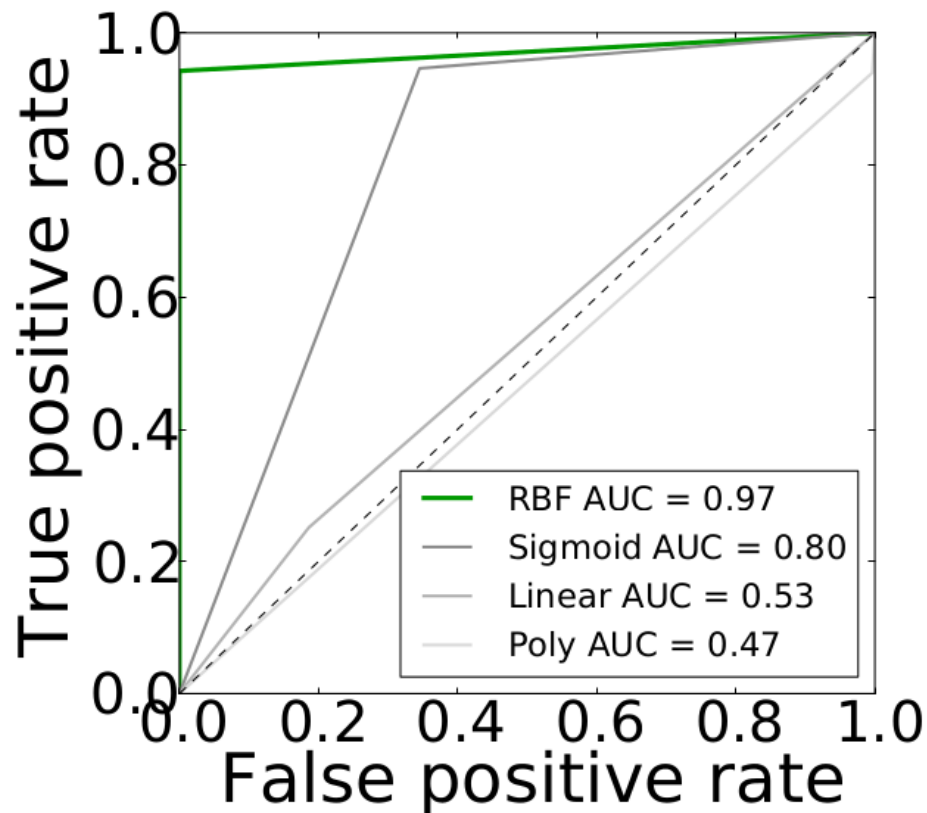
Backup slides

Hermes: predicting popularity and adapting replication

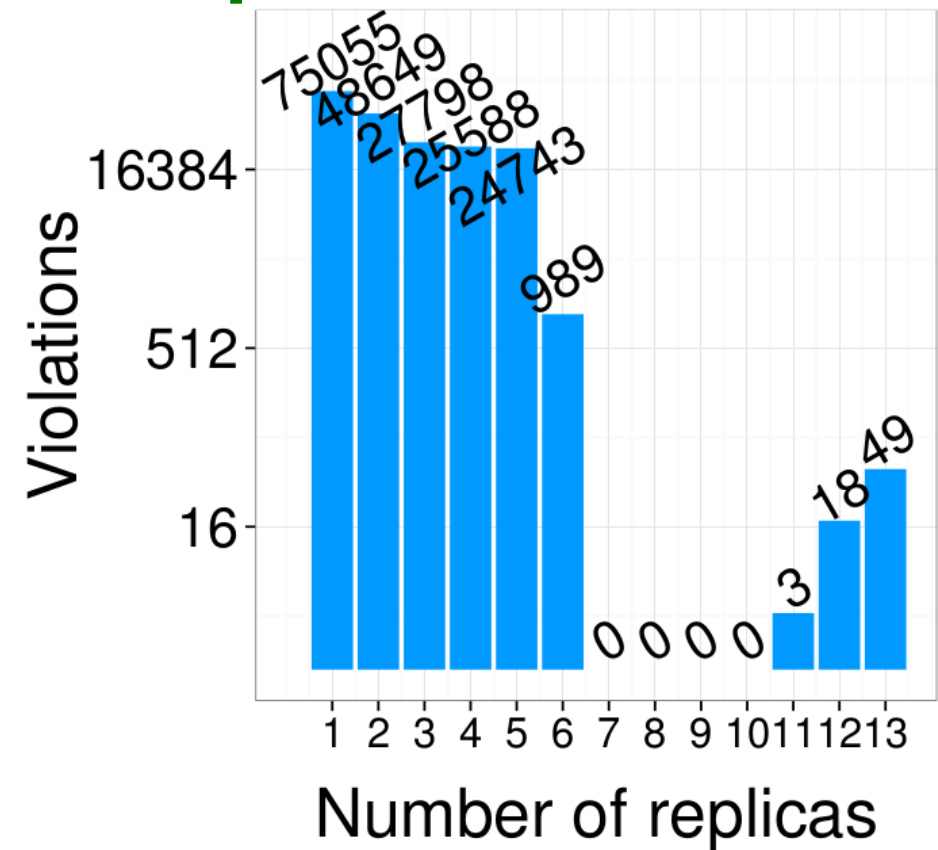
[ICPADS'13]

- Tuning Hermes

Popularity Classifier



Replication Classifier



Perspectives

- Improve learning dataset (traces, logs, testbed measurements...)
- Improve Caju design (like a P2P system)
- Evaluate our schemes through a proof-of-concept prototype
- Consider Hybrid/Mobile CDNs