# Personalization

# Want to scale? Think P2P

**Anne-Marie Kermarrec**

# A cry for personalization

# Why is personalization so difficult?

- Huge volume of data: small portion of interest

- Dynamic interests

- Interesting stuff does not come always from friends

- Classical notification systems do not filter enough or too much

Scalable personalization infrastructures
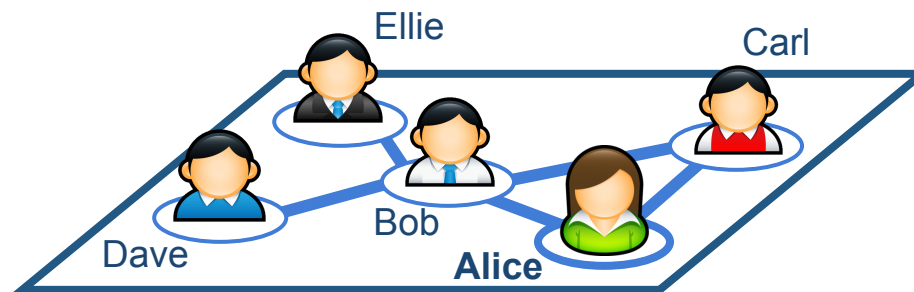
# KNN computation over large data

**Basic building block for many applications**

- Similarity search

- Machine learning

- Data mining

- Image processing

- **Collaborative filtering**

# KNN-based user-centric collaborative filtering

Provide each user with her k closest neighbors

(Users owns a profile, the system has its favorite similarity metric)



Use this topology for

- personalized notifications

- recommendation

# Dealing with truly big data

**Want to scale?  Think P2P**

# Do not look exhaustively

# The key to scalability in KNN graph construction

**Consider a partial set of candidates**

**Sampling-based approach**

# P2P KNN graph construction

Which nodes are close? $\Longrightarrow$ Similarity metric

How to discover them? $\Longrightarrow$ Sampling

# Which nodes are close?

**Model**

*U(sers) × I(tems) (items)*

*Profile(u) = vector of liked/shared/viewed items*

**Cosine similarity metric**

$$Similarity\ (n,p) = \frac{n.p}{\|n\|\,\|p\|}$$

**Jaccard metric**

$$Jaccard(n,p)\ = \left| \frac{n \cap p}{n \cup p} \right|$$

Minimal information: **no tag, no user's input, generic**

# How to discover them: Gossip-based computing
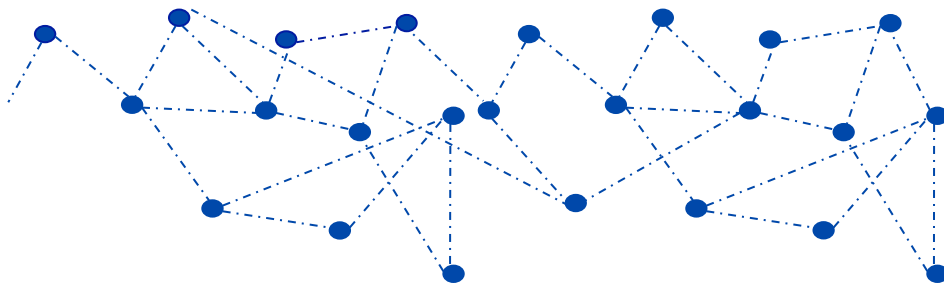
Each node maintains a set of
neighbors (c entries)

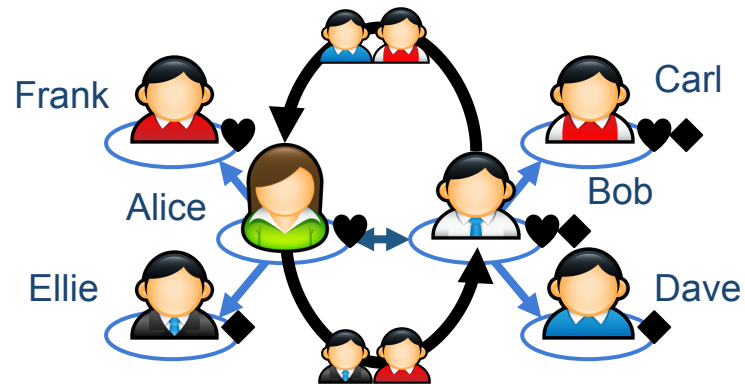Peer exchange

Shuffle

P      Q

Result ➜ random graph

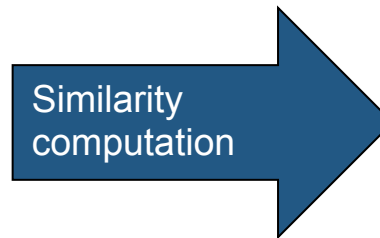Highly resilient against churn, partition

Small diameter

[JGKVV, ACM TOCS 2007]
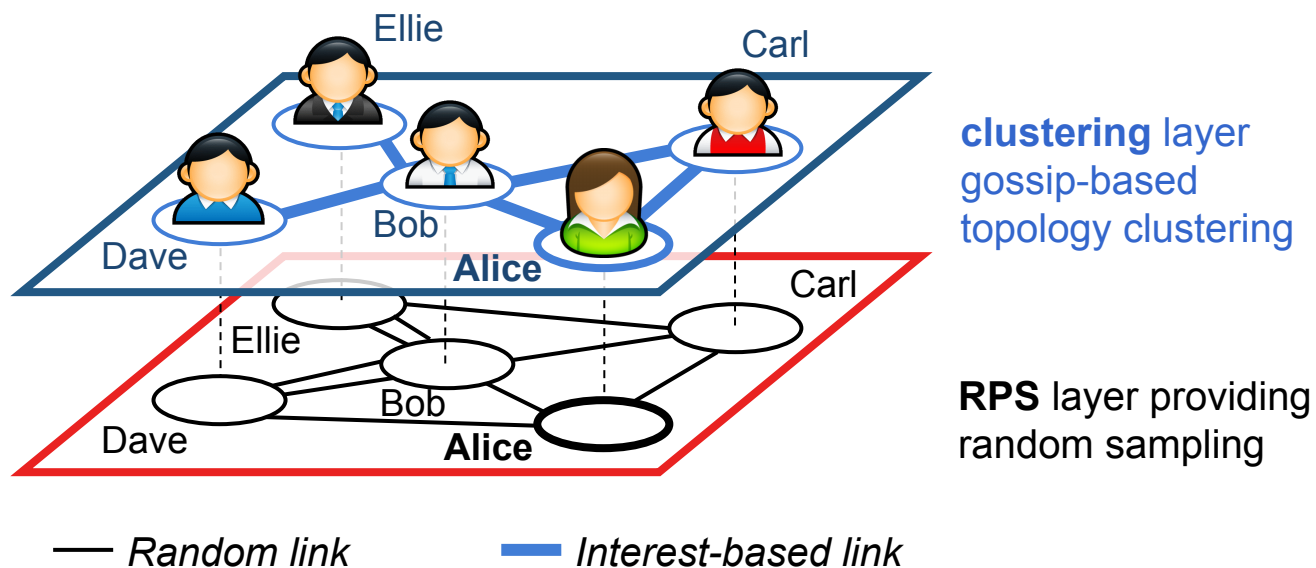
*Inria*

# KNN construction



1. exchange of neighbors lists

Similarity computation

2. neighborhood optimization

# Decentralized KNN selection



clustering layer
gossip-based
topology clustering

RPS layer providing
random sampling
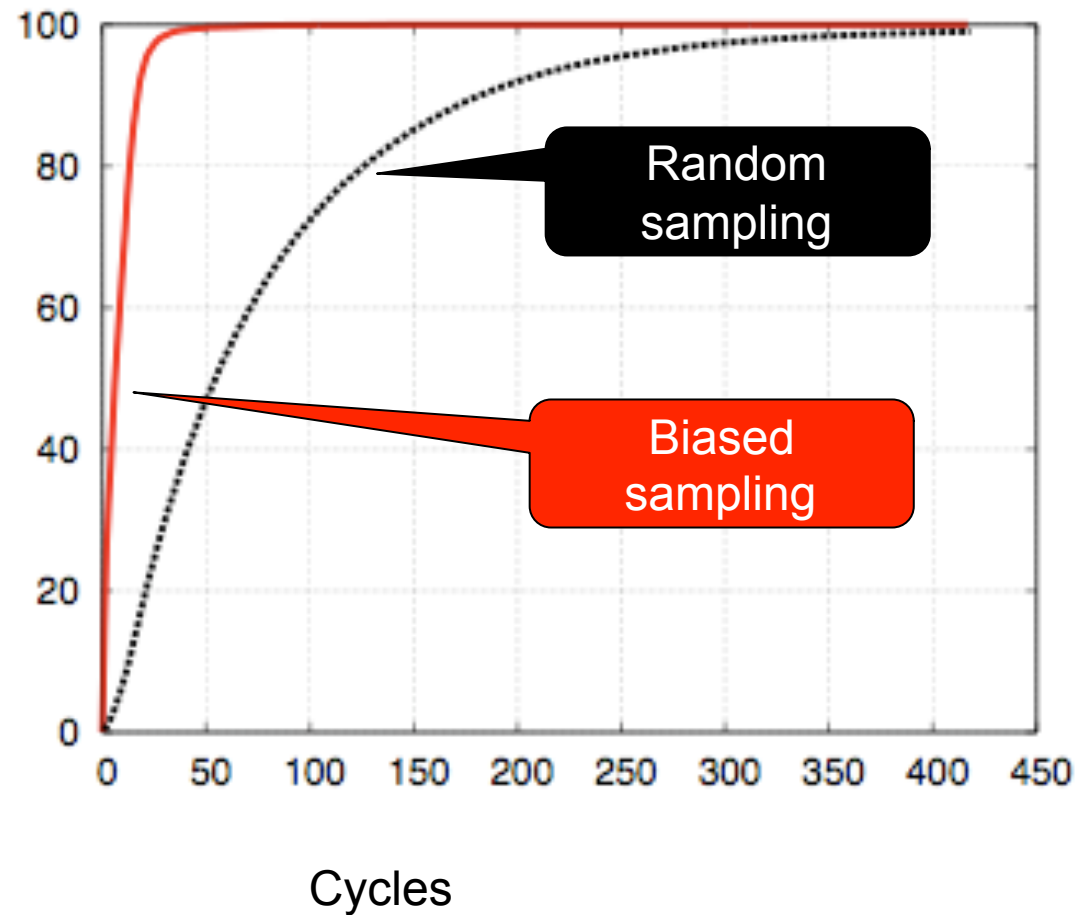
— *Random link*          — *Interest-based link*

[FGKL Middleware 2010]

# Convergence

c current
neighbors versus
the c closest

**Cycles**

# Applications

- **Decentralized news recommendation** [BFGJK, IPDPS 2013]

- Top-K [BGKL, ACM TODS 2011] [BGK, ACM TOIT 2014]

- Geo recommendation [BKKT, ICDCS 2012]

# DECENTRALIZED NEWS RECOMMENDER

# Notification is taking over

# An implicit notification system

# based on collaborative filtering

# WhatsUp in a nutshell



Generate profile

KNN selection

User Opinion

Influence dissemination

Influence topology

Dissemination

Tune dissemination

# Dissemination: orientation and amplification

Orientation: **to whom**?

Amplification: **to how many**?

**Exploit**: Forward To friends

**Explore**: Forward to random users

Increase Fanout (Log(n))

Decrease Fanout (1)

# WhatsUp in action on the survey (480 users)

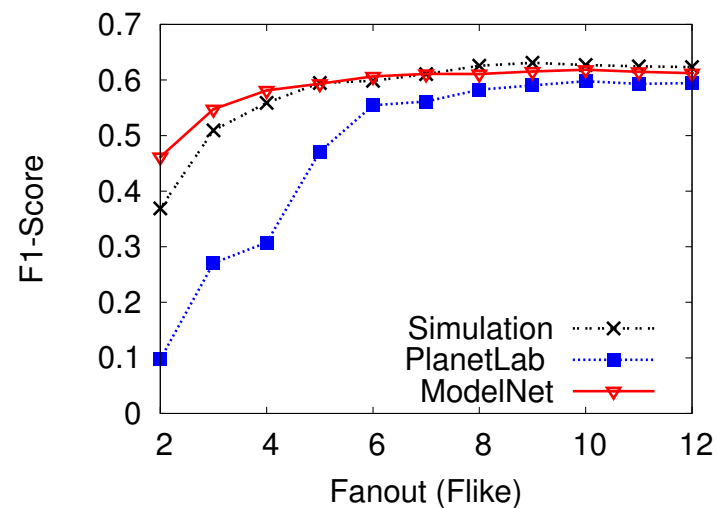|  | Precision | Recall | F1-Score | Messages |
|---|---|---|---|---|
| Gossip (f=4) | 0.34 | 0.99 | 0.51 | 2.3 M |
| Cosine-CF | 0.50 | 0.65 | 0.57 | 5,9k |
| **Whatsup (f=10)** | **0.471** | **0.83** | **0.60** | **2,4k** |

# Orientation (survey)

News items received through a dislike forward

| Number of dislikes | 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|
| Fraction of liked news | 54% | 31% | 10% | 3% | 2% |

# WhatsUp versus Pub/Sub

| Approach | Precision | Recall | F1-Score |
|----------|-----------|--------|----------|
| Pub/Sub  | 0.40      | 1.0    | 0.58     |
| WhatsUp  | 0.47      | 0.83   | 0.60     |

# WhatsUp versus cascading

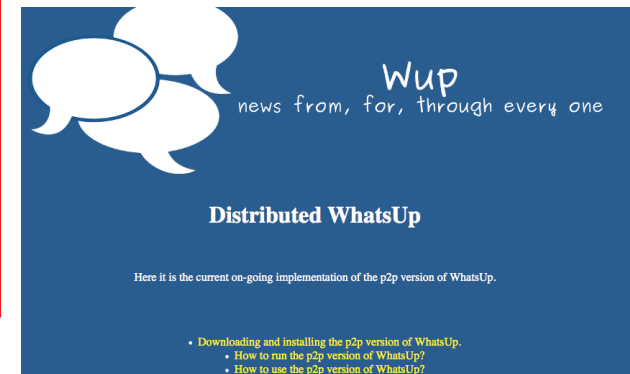| Approach | Precision | Recall | F1-Score |
|----------|-----------|--------|----------|
| Cascading | 0.57 | 0.09 | 0.16 |
| WhatsUp | 0.56 | 0.57 | 0.57 |

**Take away message**

Personalization is needed

Decentralization is healthy

Gossip-based computing is one (the) way to go



**Privacy matters**

- Obfuscation
- Anonymous routing
- Threshold protocol
- Differentially private protocol
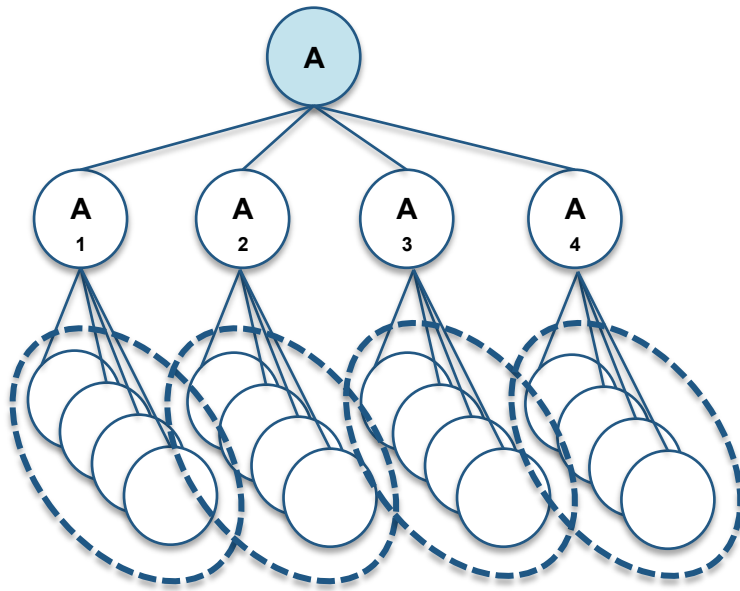- Landmark-based protocol

STRONGER GARANTEES

**http://131.254.213.98:8080/wup/**

**Operational prototype**
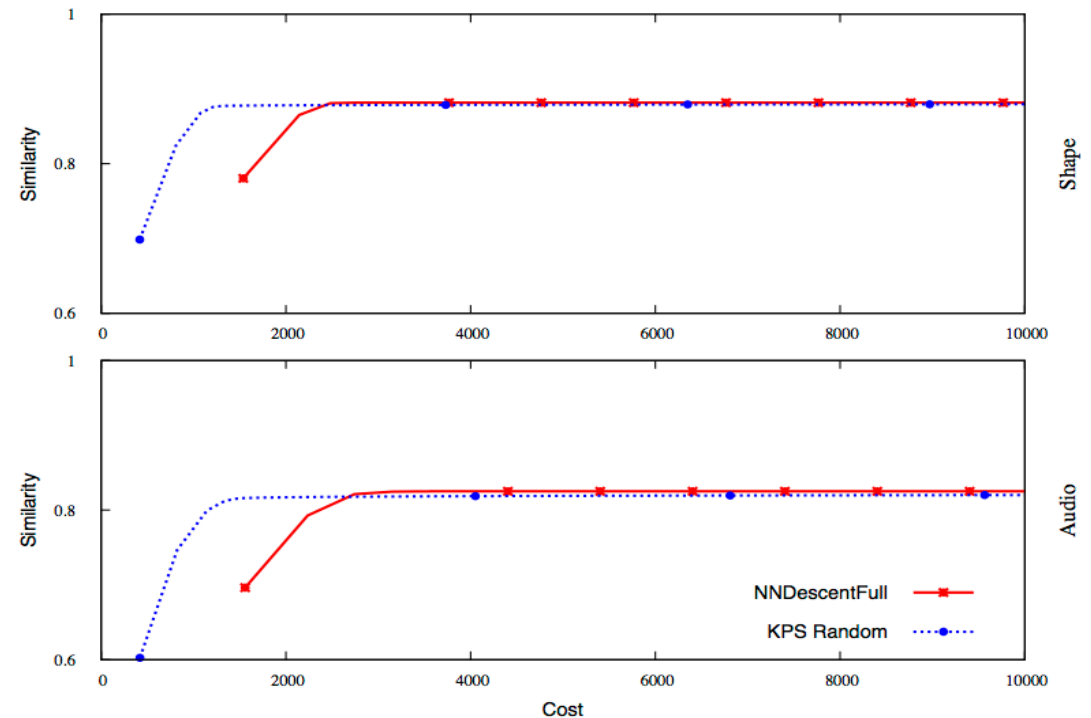
**Tested on 500 users @ TrentoRise last year**

**TRY IT** ☺

# For those who are afraid of P2P

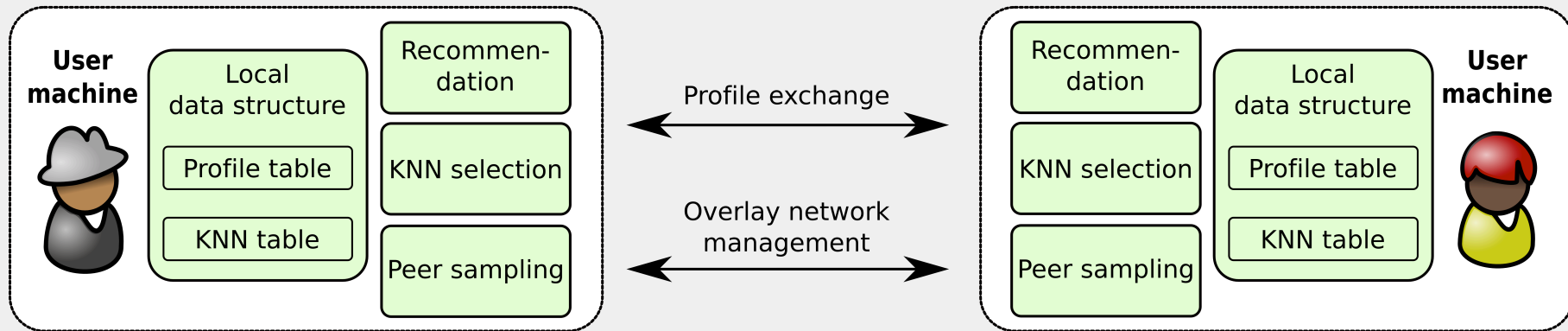# Turns out to be an effective centralized algorithm too.



Candidate set: neighbors of neighbors
+ Random candidates for dynamics

Comparison with [Dong&al, 2012]
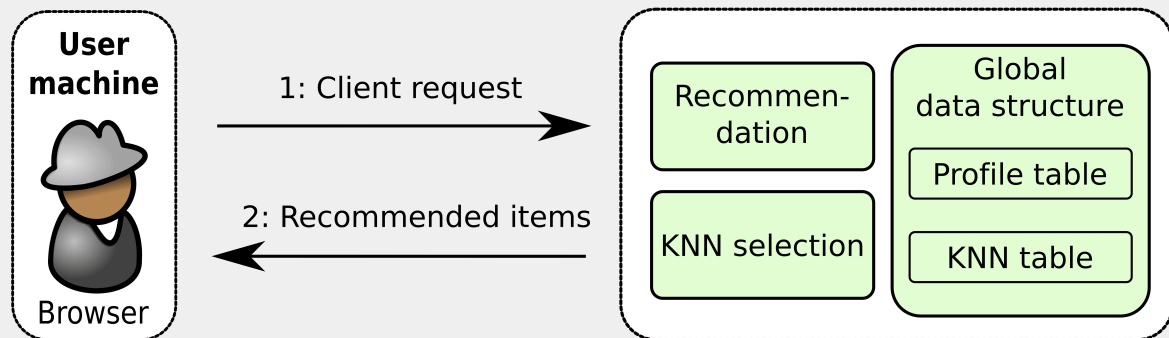
# Hybrid recommendation engine

# Decentralized approach

| User machine | Local data structure | Recommen-dation |
|---|---|---|
| | Profile table | KNN selection |
| | KNN table | Peer sampling |

Profile exchange

Overlay network management

| Recommen-dation | Local data structure | User machine |
|---|---|---|
| KNN selection | Profile table | |
| Peer sampling | KNN table | |

# Data structures

| Profile table | |
|---|---|
| uid | P(uid) = {list of iid} |

| KNN table | |
|---|---|
| uid | Knn(uid) ={list of uid} |

# Centralized approach

User machine

Browser

1: Client request

2: Recommended items

| Recommen-dation | Global data structure |
|---|---|
| | Profile table |
| KNN selection | KNN table |

# Cost

| Dataset | Users | Items | Ratings |
|---|---|---|---|
| MovieLens1 | 943 | 1700 | 100,000 |
| MovieLens2 | 6,040 | 4000 | 1,000,000 |
| MovieLens3 | 69,878 | 10,000 | 10,000,000 |
| Digg | 59,167 | 7724 | 782,807 |

CRec ☐    Mahout ■

ClusMahout ■    Exhaustive ■



Grid 5000 implementation

# HyRec: Taking the best of both worlds

**User machine**

Recommen-dation

KNN selection

Browser

1: Client request →

← 2: Candidate set

3: Update KNN →

**Server**

Personalization orchestrator

Sampler

Global data structure

Profile table

KNN table
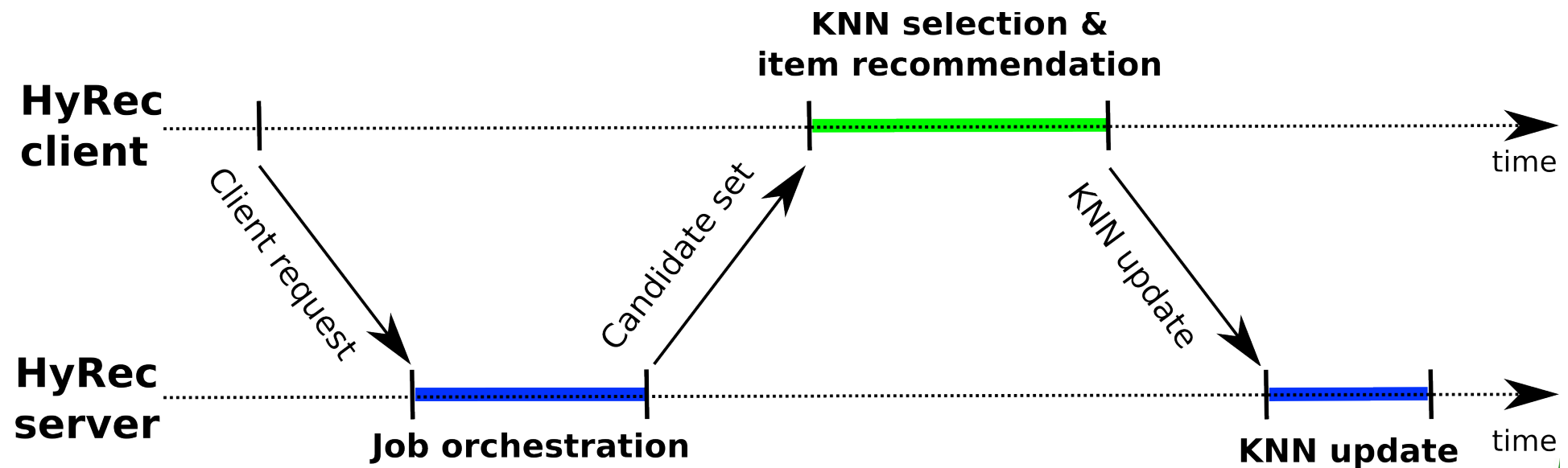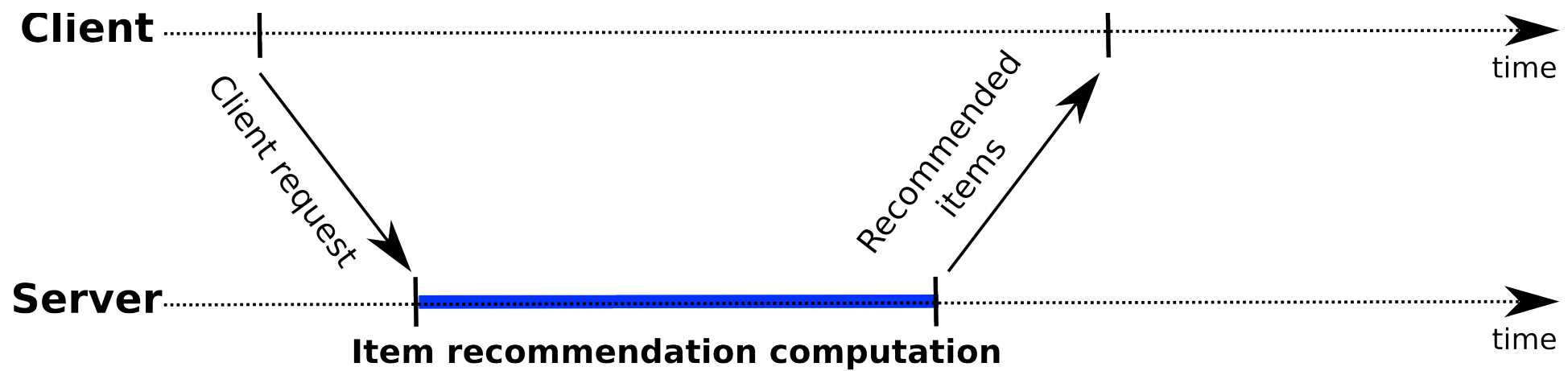
Online KNN selection

Candidate set (k): $k^2$ users for quick convergence, k random (biased) for dynamics
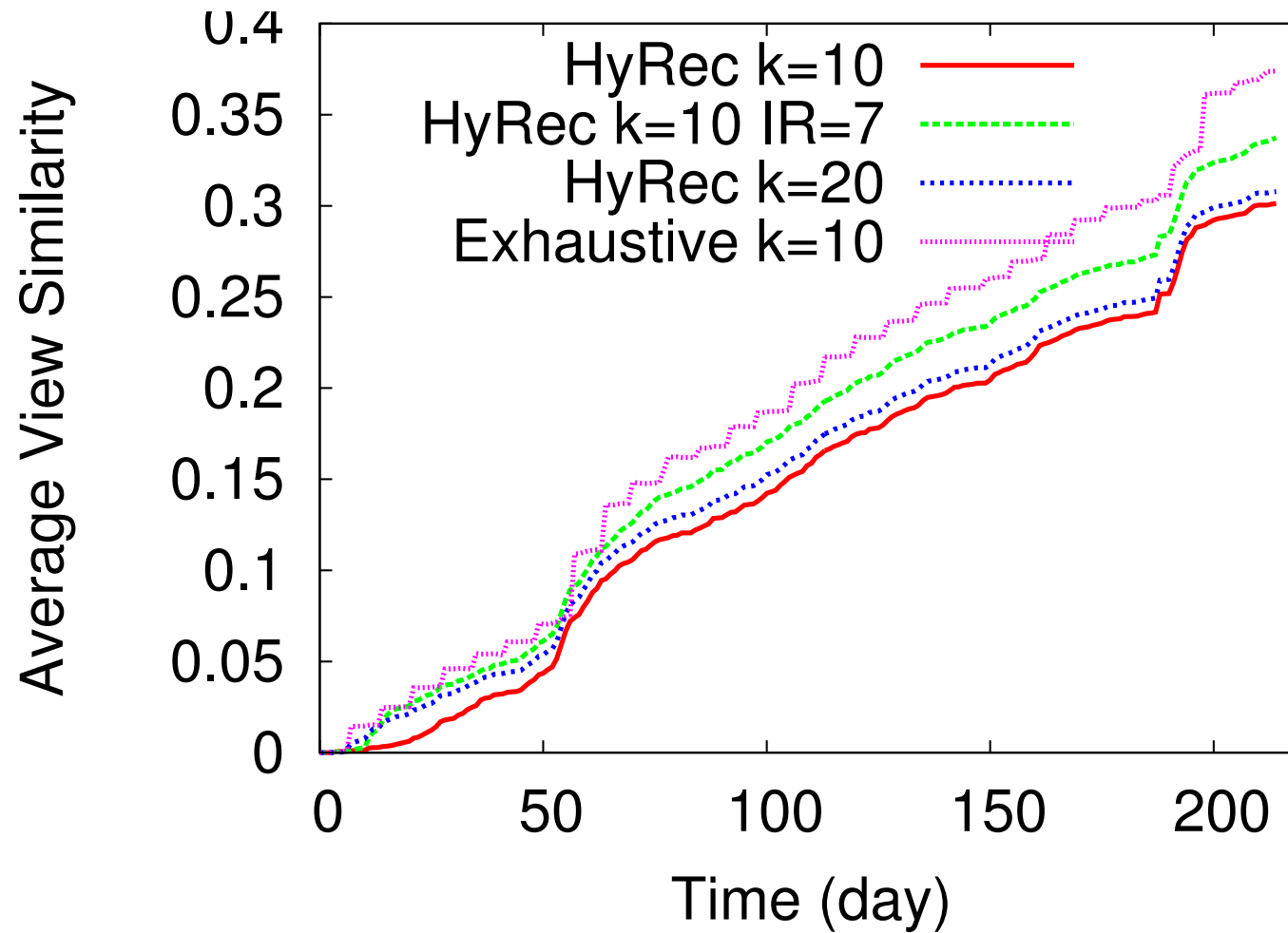
No data stored at the client

Recommendation: R most popular items
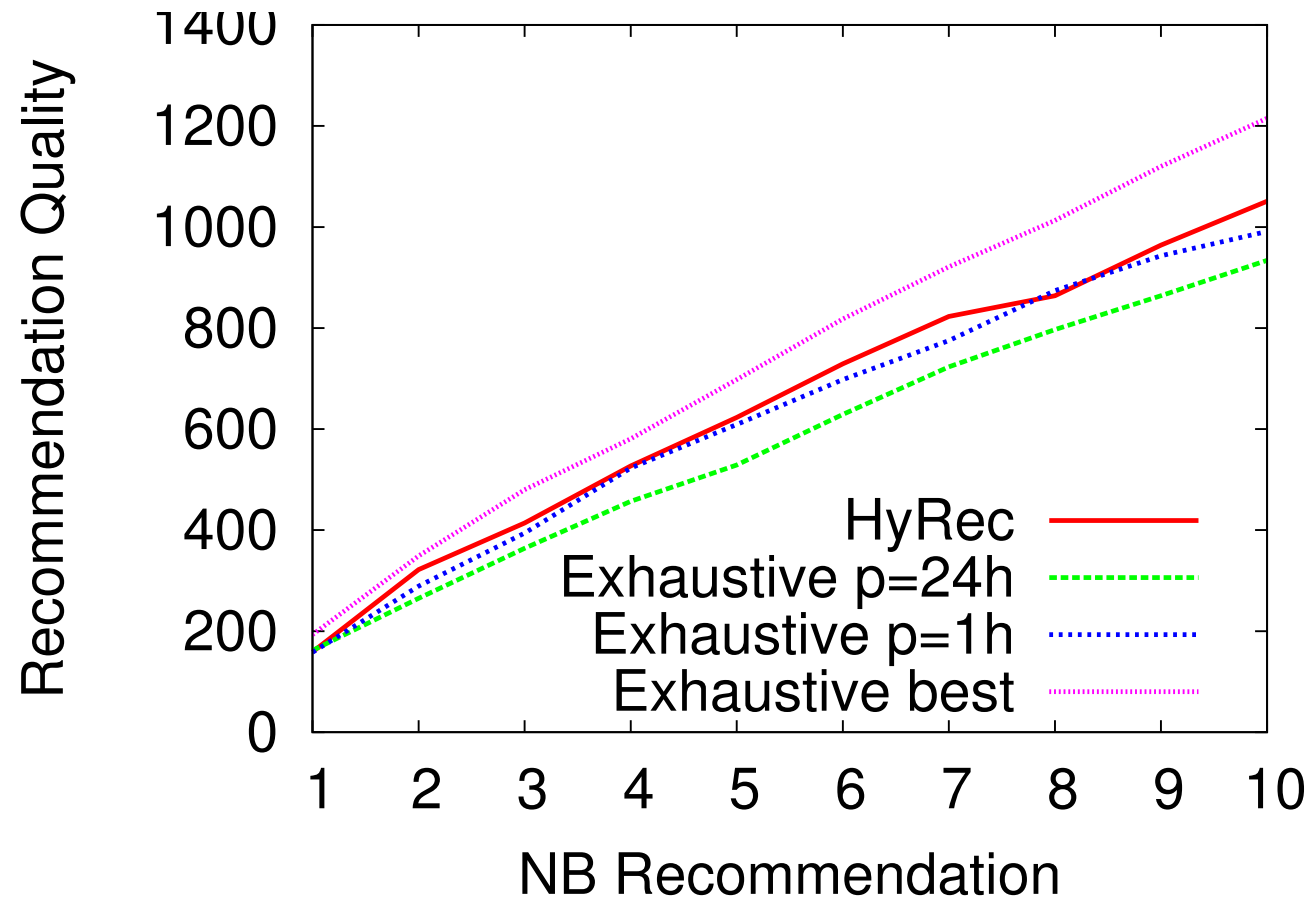
HyRec client: Javascript (widget) running in the browser

**Client** · · · | · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · | · · · · ▶
time

Client request

Recommended items

**Server** · · · | ━━━━━━━━━━━━━━━━━━━━━━━━ | · · · ▶
time

**Item recommendation computation**

**KNN selection & item recommendation**

**HyRec client** · · · | · · · · · · · · · · · · · · · · · · · | ━━━━━━━━━━ | · · · ▶
time

Client request

Candidate set

KNN update

**HyRec server** · · · | ━━━━━━━ | · · · · · · · · · · · · · · · · | ━━━━ | ▶
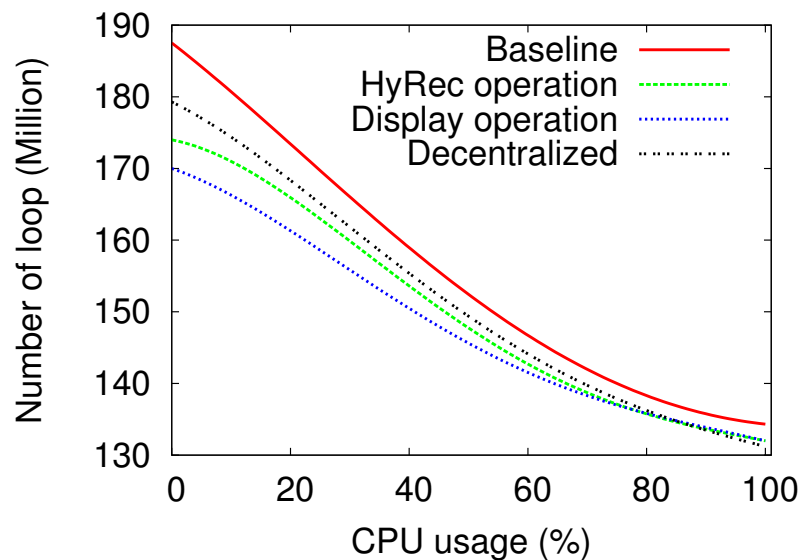time

**Job orchestration**     **KNN update**

# View similarity (MovieLens)
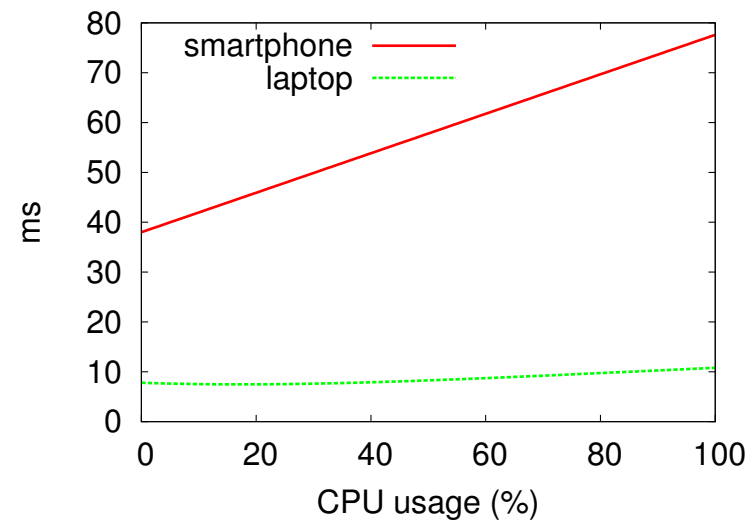
# Recommendation quality

# HyRec versus the client load
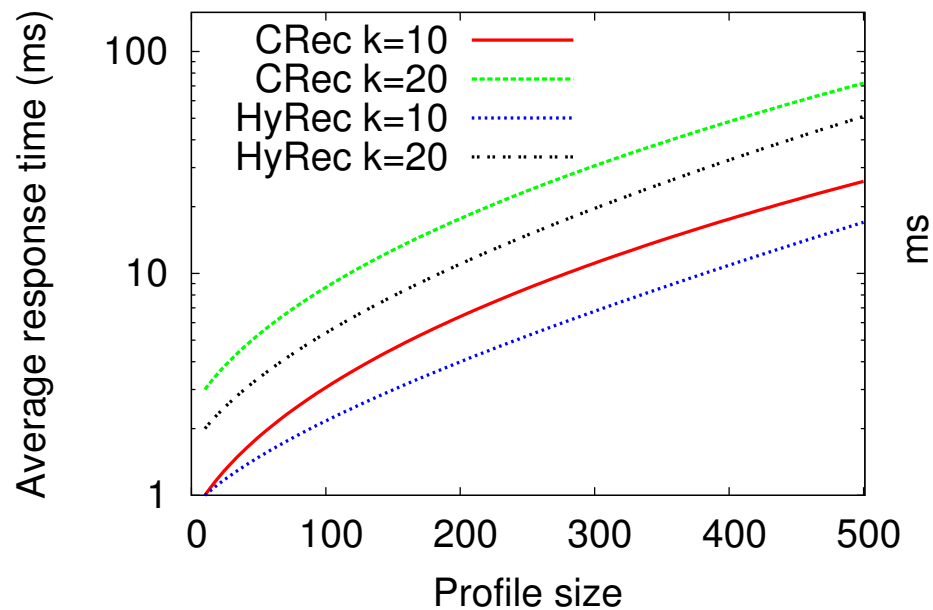


Impact of HyRec

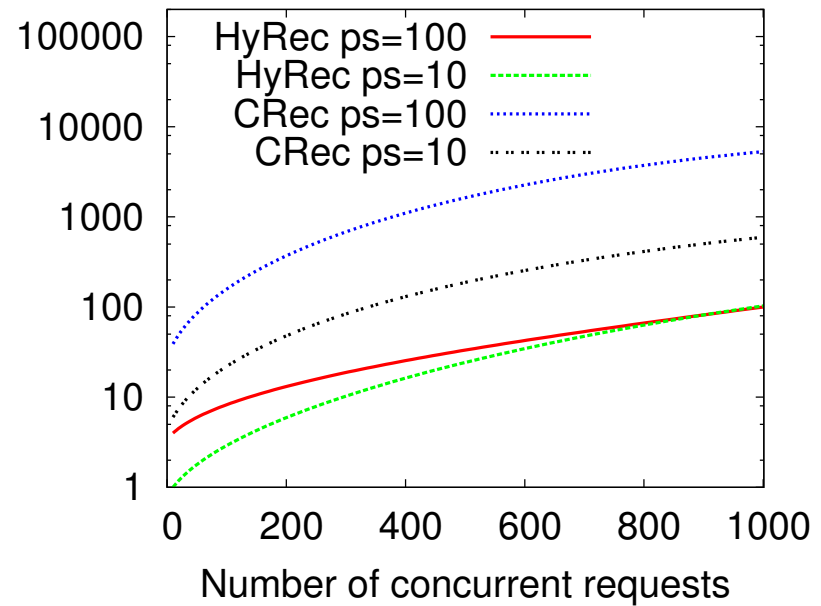Negligible disruption of HyRec

Impact of the client load

50% load
<60ms on smartphone
<10ms on laptop

# HyRec versus a centralized recommender



Impact of the profile size

Impact of the request stress

# Take away message

P2P design is crucial

Leveraging clients machine has a significant impact  on scalability

Enable any content provider to implement personalization

# To take away

Personalization is crucial

P2P in a design mindset

Thank you