

Introducing Service-level Awareness in Clouds

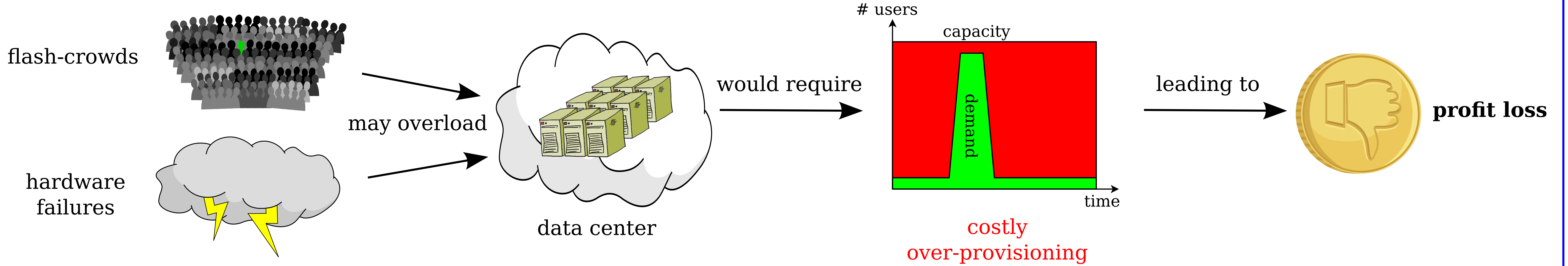
Cristian Klein¹, Martina Maggio², Karl-Erik Årzén², Francisco Hernández-Rodríguez¹

¹Umeå University, Sweden, ²Lund University, Sweden

cristian.klein@cs.umu.se, martina.maggio@control.lth.se, karlerik@control.lth.se, francisco@cs.umu.se

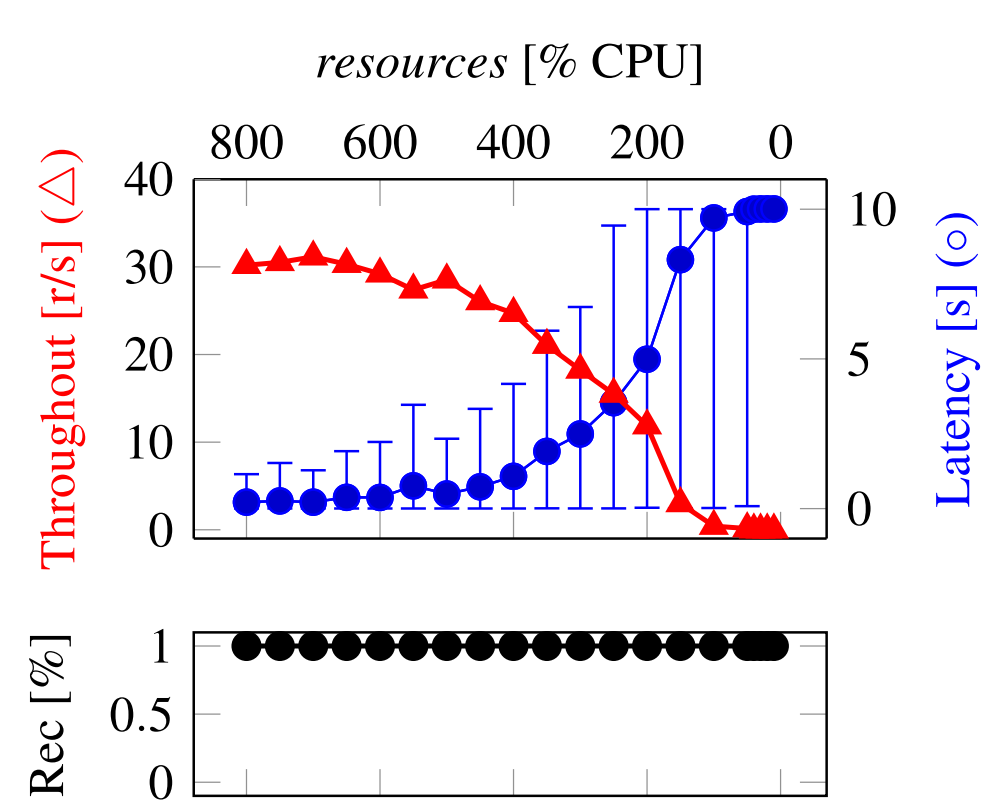
<https://github.com/cristiklein/cloudish>

Problem: unexpected events

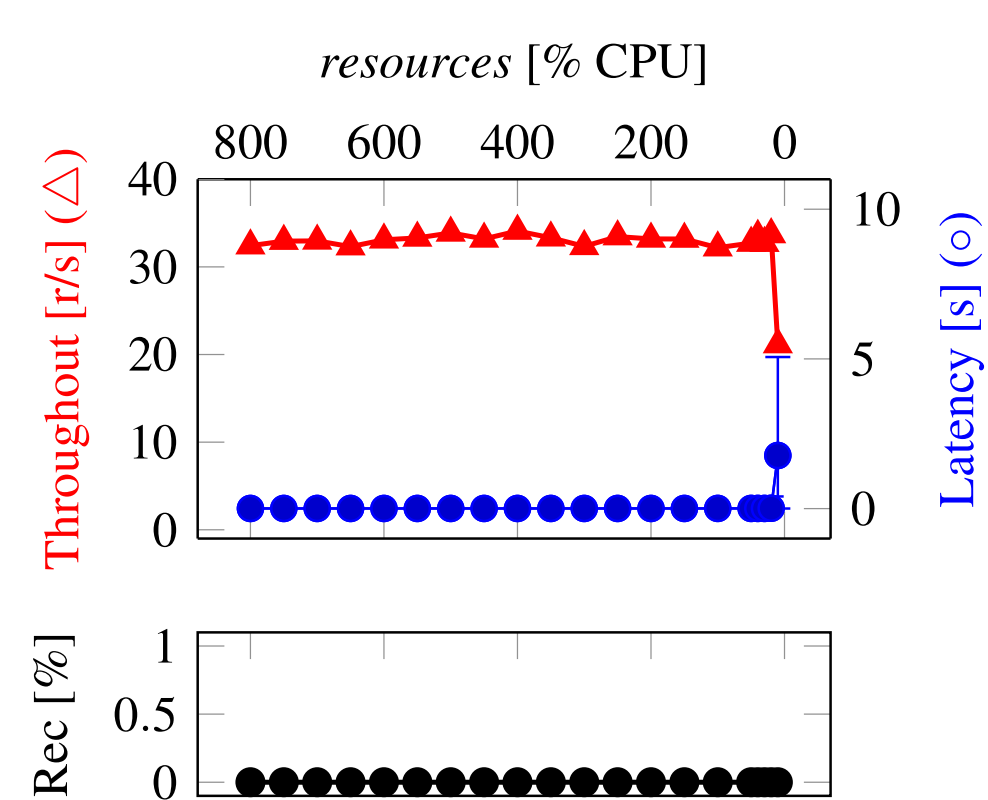


Idea: reduce **service-level** to avoid data center overload

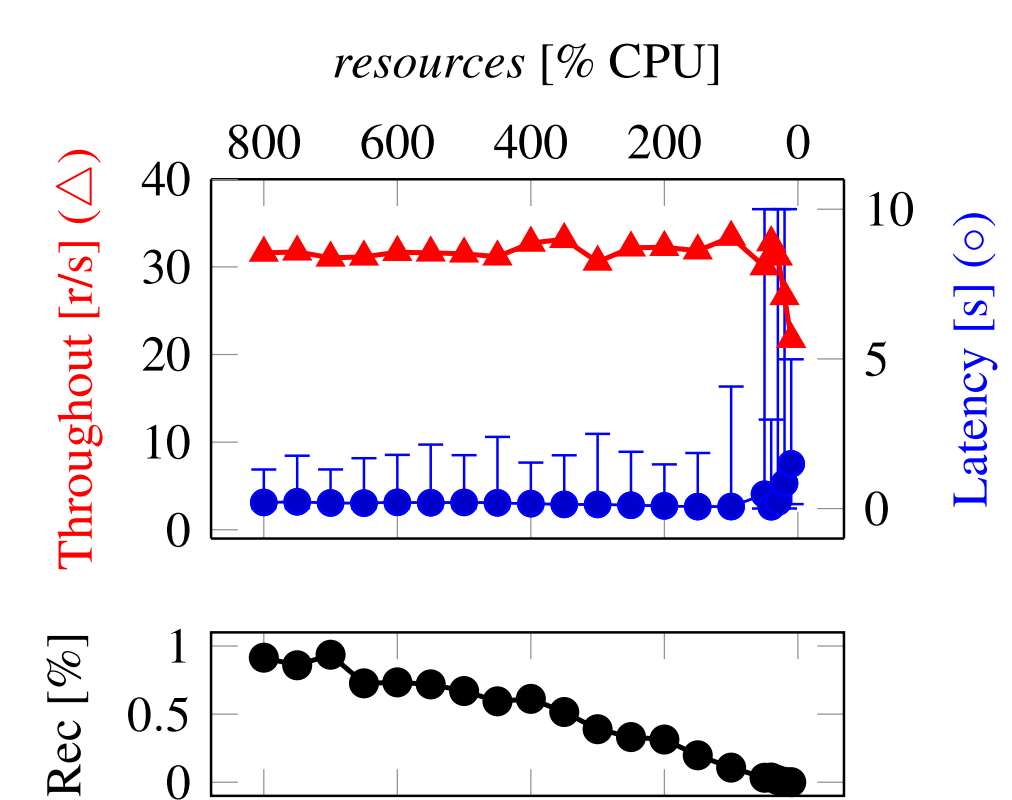
Downgrade **user experience** to reduce capacity requirements and serve more users



RUBiS with full recommendations

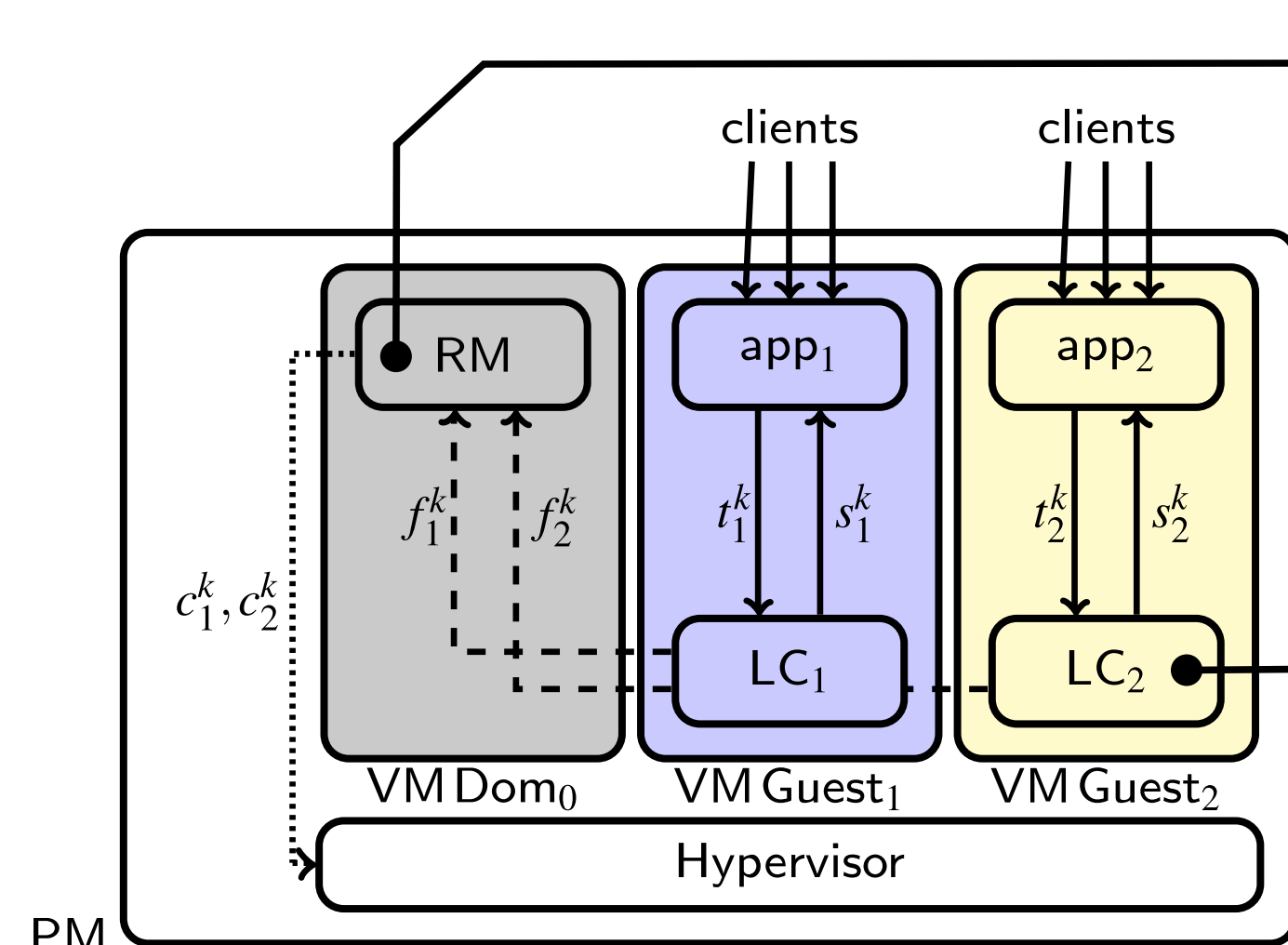
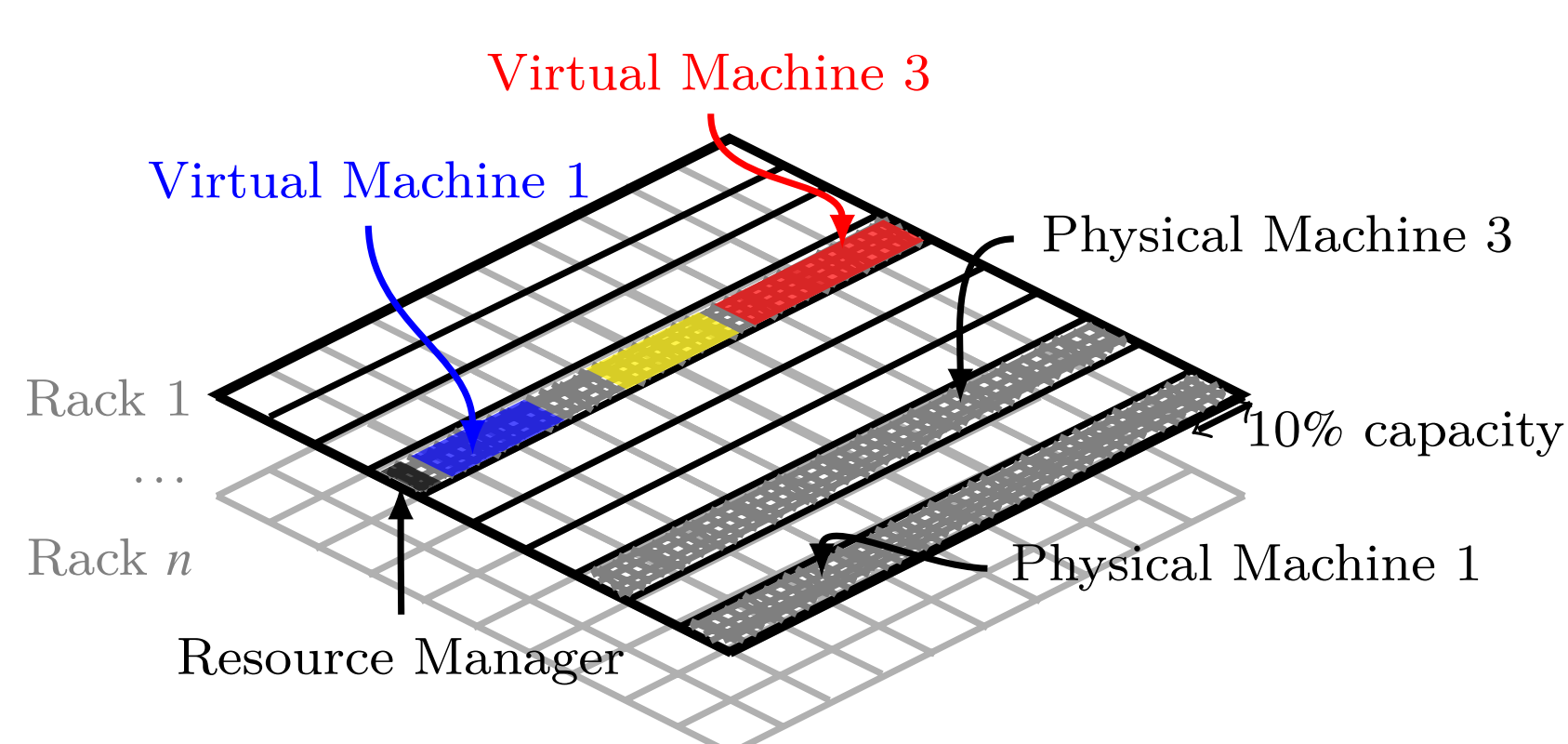


RUBiS with no recommendations



RUBiS with adaptive recommendations

Our Proposal



Resource Manager (RM): designed using **game theory**

$$c_i^{k+1} = c_i^k - \epsilon_{rm} (f_i^k - c_i^k \cdot \sum_i f_i^k)$$

matching value sent by app. i at time index k

capacity allocation of app. i at time index $k+1$

Local Controller (LC): designed using **control theory**

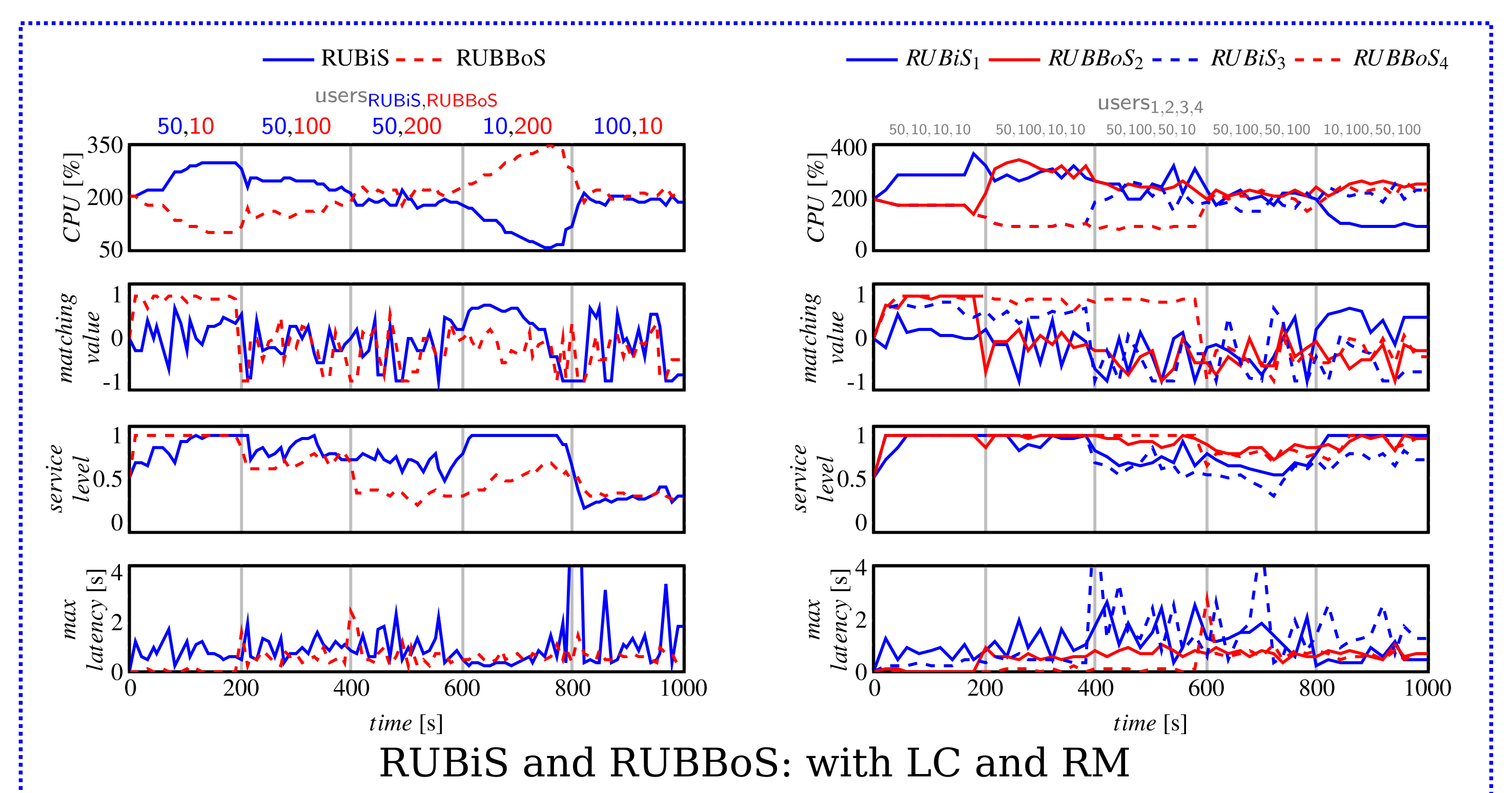
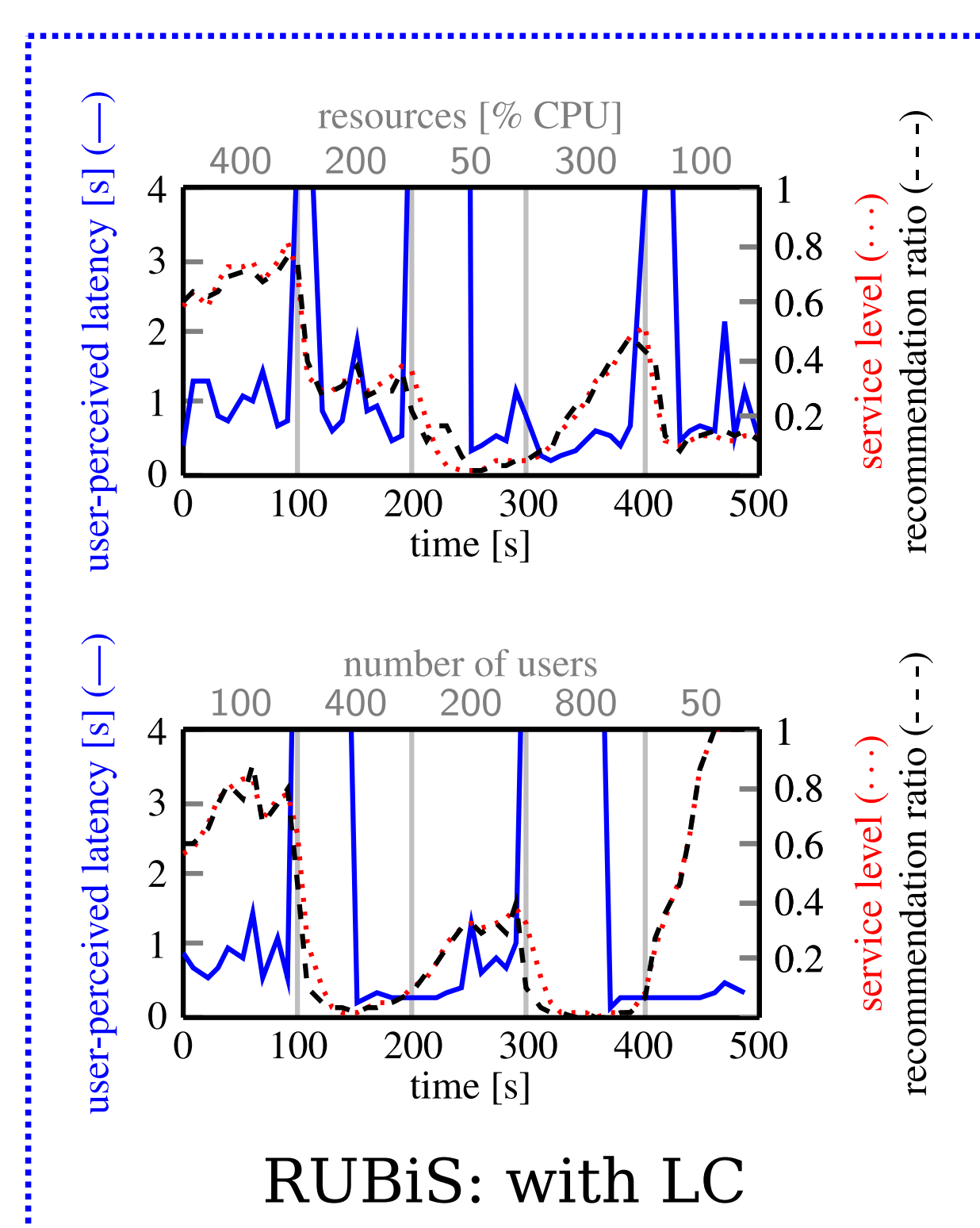
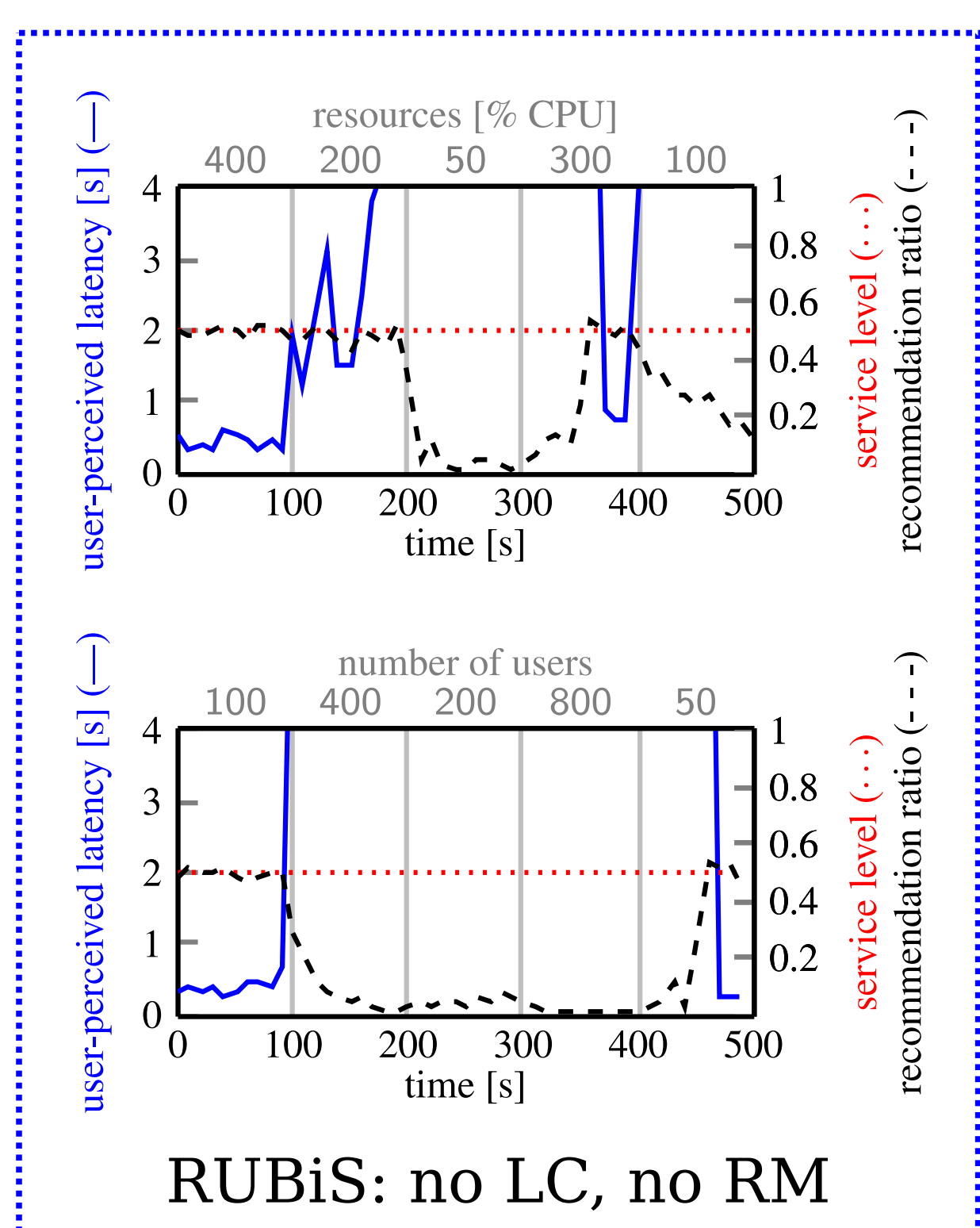
$$s_i^{k+1} = s_i^k + \frac{1-p_1}{\alpha} \cdot (t_i^k - \bar{t}_i) \quad f_i^k = 1 - \frac{t_i^k}{\bar{t}_i}$$

service-level

target latency

measured latency

Results



Conclusions

- clouds can more robustly **withstand capacity shortages**
 - flash-crowds
 - hardware failures
- applications can reduce their service level (turn off optional features)
- infrastructure can rebalance resources among hosted applications

Perspectives

- combining with other mechanism
 - migration
 - horizontal scaling
 - overbooking
- devising a billing model

