

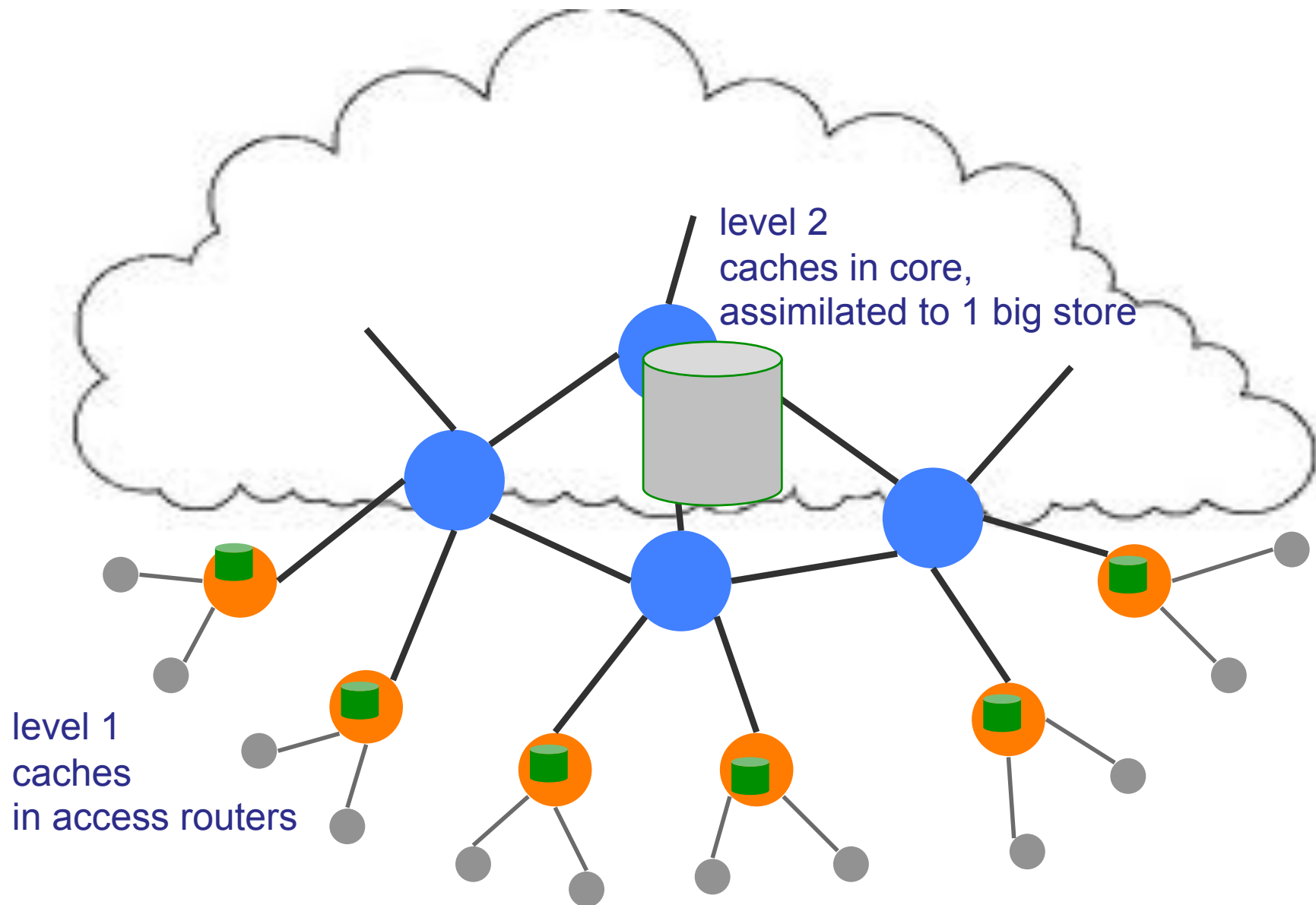
Impact of traffic mix on caching performance

Christine Fricker, Philippe Robert, Jim Roberts, Nada Sbihi

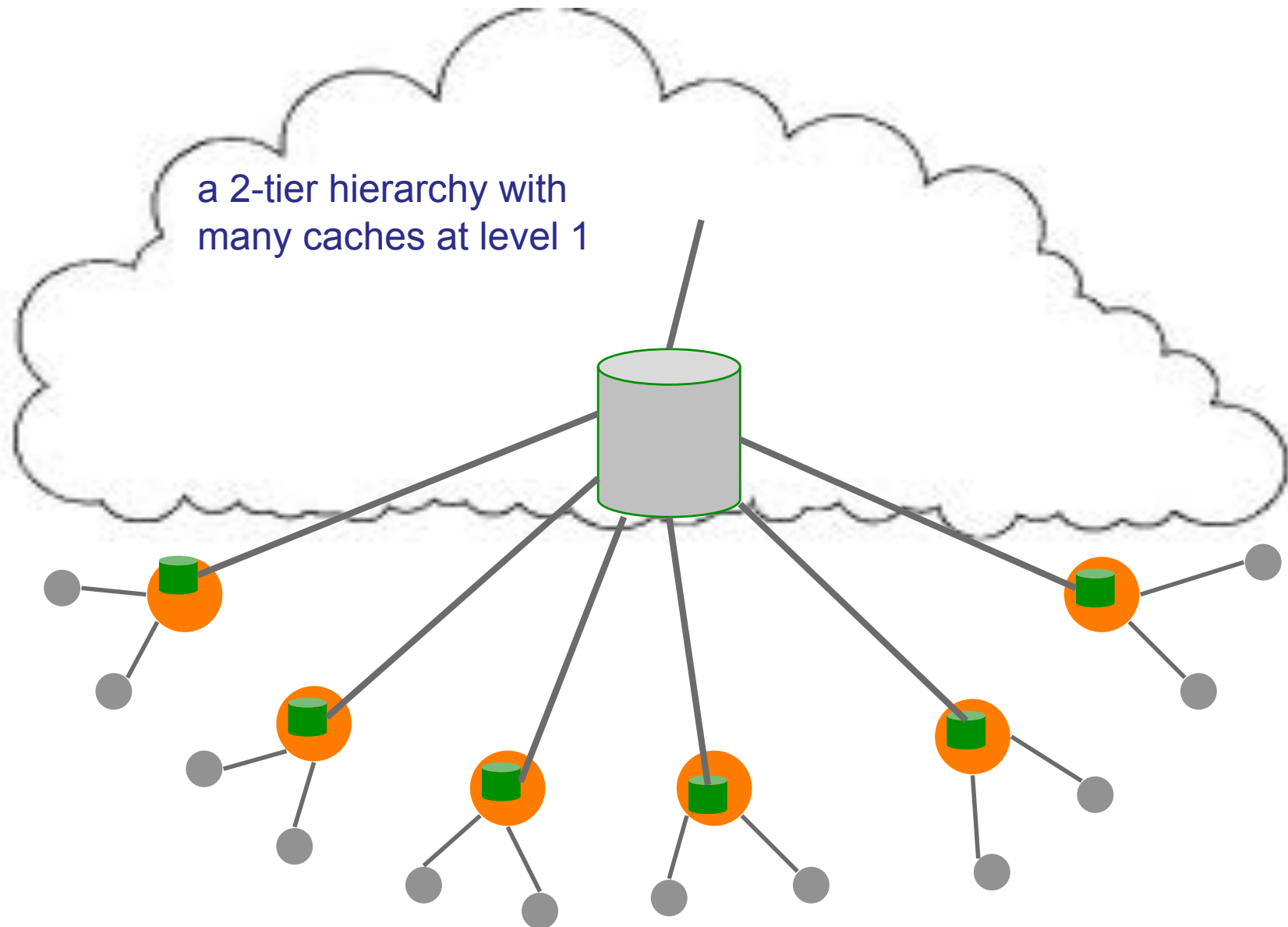
INRIA/RAP

Evaluation Seminar
Rungis, 22 March 2012

A network of caches



A network of caches



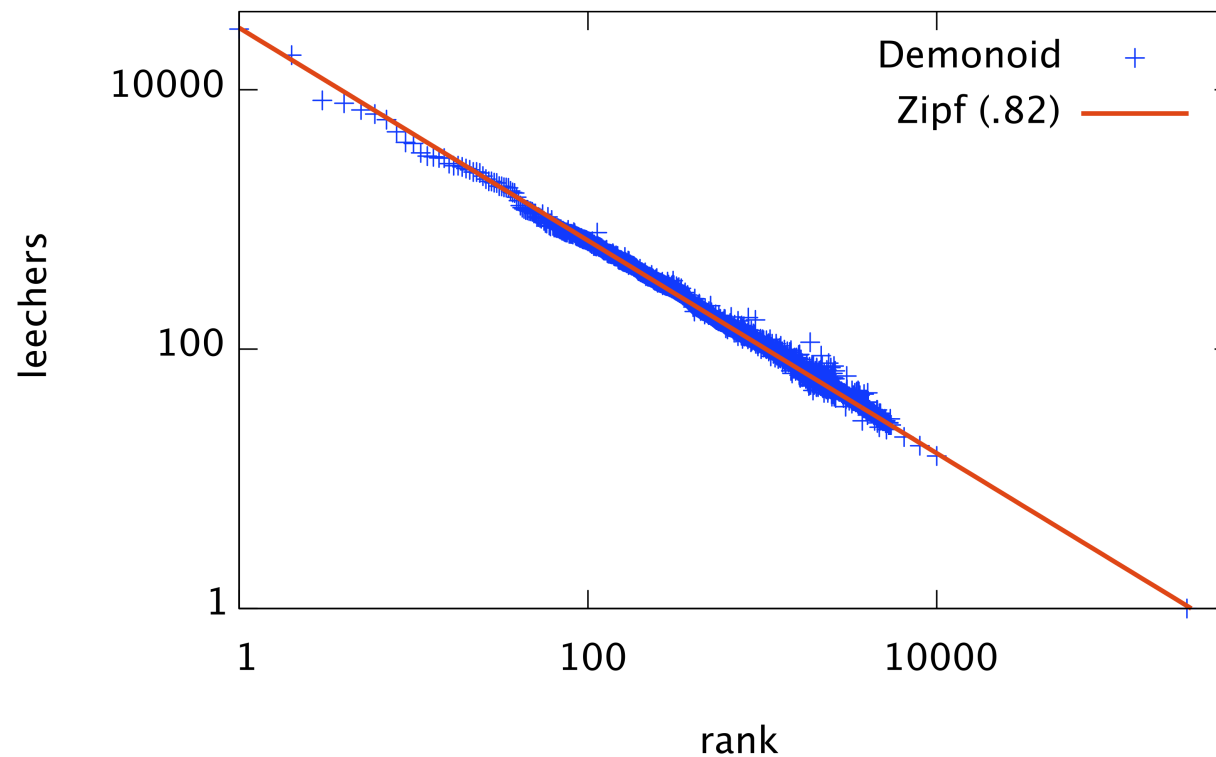
Content mix

- Cisco: "96% of traffic is content transfer"
- web, file sharing, user generated content, video on demand
- billions of objects, petabytes of content!

	share	objects	size	volume
web	.18	10^{11}	10 KB	1 PB
file sharing	.36	10^5	10 GB	1 PB
UGC	.23	10^8	10 MB	1 PB
VoD	.23	10^4	100 MB	1 TB

Content popularity

- Zipf law popularity
 - request rate for n^{th} most popular object, $q(n) \propto 1/n^\alpha$
 - eg, for 270000 torrents on Demonoid.me, $\alpha = 0.82$



Content popularity

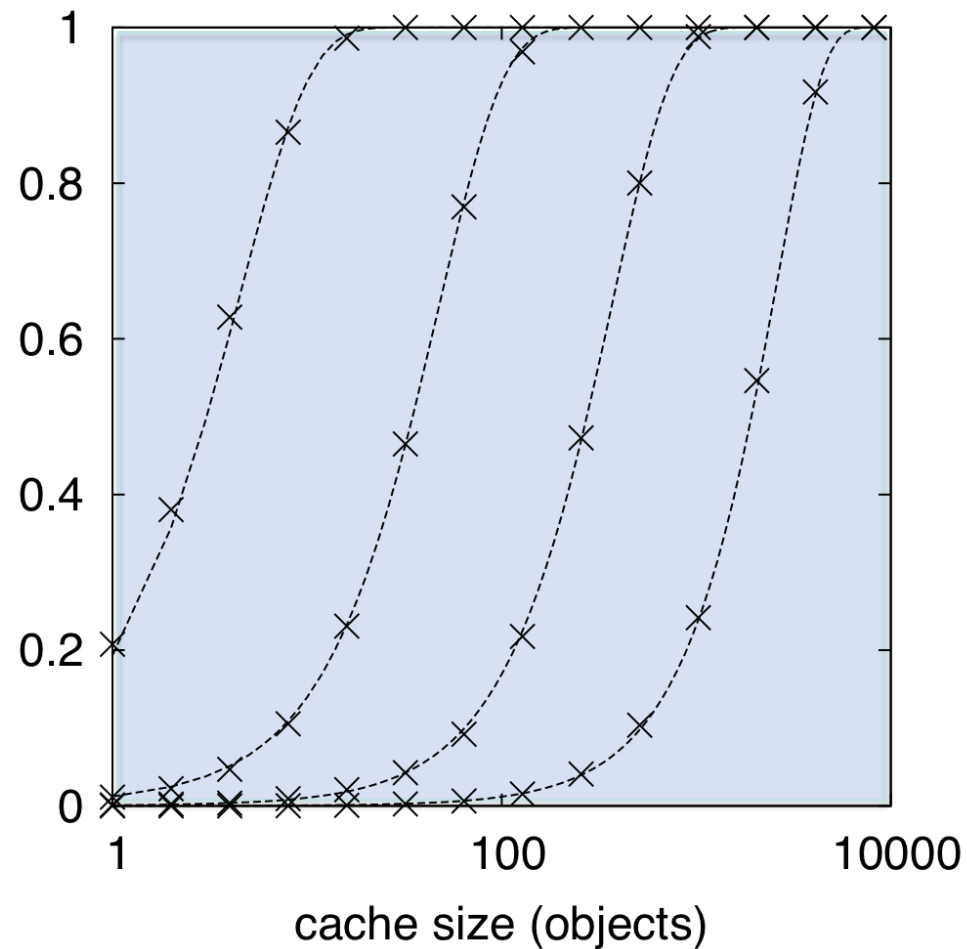
- Zipf law popularity
 - request rate for n^{th} most popular object, $q(n) \propto 1/n^\alpha$
 - eg, for 270000 torrents on Demonoid.me, $\alpha = 0.82$
- web page popularity Zipf(a), $.64 < \alpha < .83$
- file sharing Zipf(a), $.75 < \alpha < .82$
- user generated content Zipf(a), $.56 < \alpha < .88$
- video on demand Zipf(a), $.65 < \alpha < 1.2$, Weibull, ...
- for our evaluations, we suppose
 - Zipf(.8) for web, file sharing, UGC
 - Zipf(.8) or Zipf(1.2) for VoD

Calculating LRU hit rates (Che, Tung and Wang, 2002)

- cache size C objects, popularity of object n is $\propto q(n)$
- assume "independent reference model" or, equivalently, Poisson request arrivals at rate $q(n)$ for object n
- "characteristic time" T_C is time for C different objects to be requested
- assume random variable T_C is approximately deterministic, $T_C \sim t_C$
- then, hit rate for object n is $h(n) = 1 - \exp\{-q(n)t_C\}$
- now, $C = \sum_n \mathbf{1}\{\text{object } n \text{ is in cache}\}$
- taking expectations, $C = \sum_n h(n) = \sum_n (1 - \exp\{-q(n)t_C\})$
- solving numerically for t_C yields $h(n)$

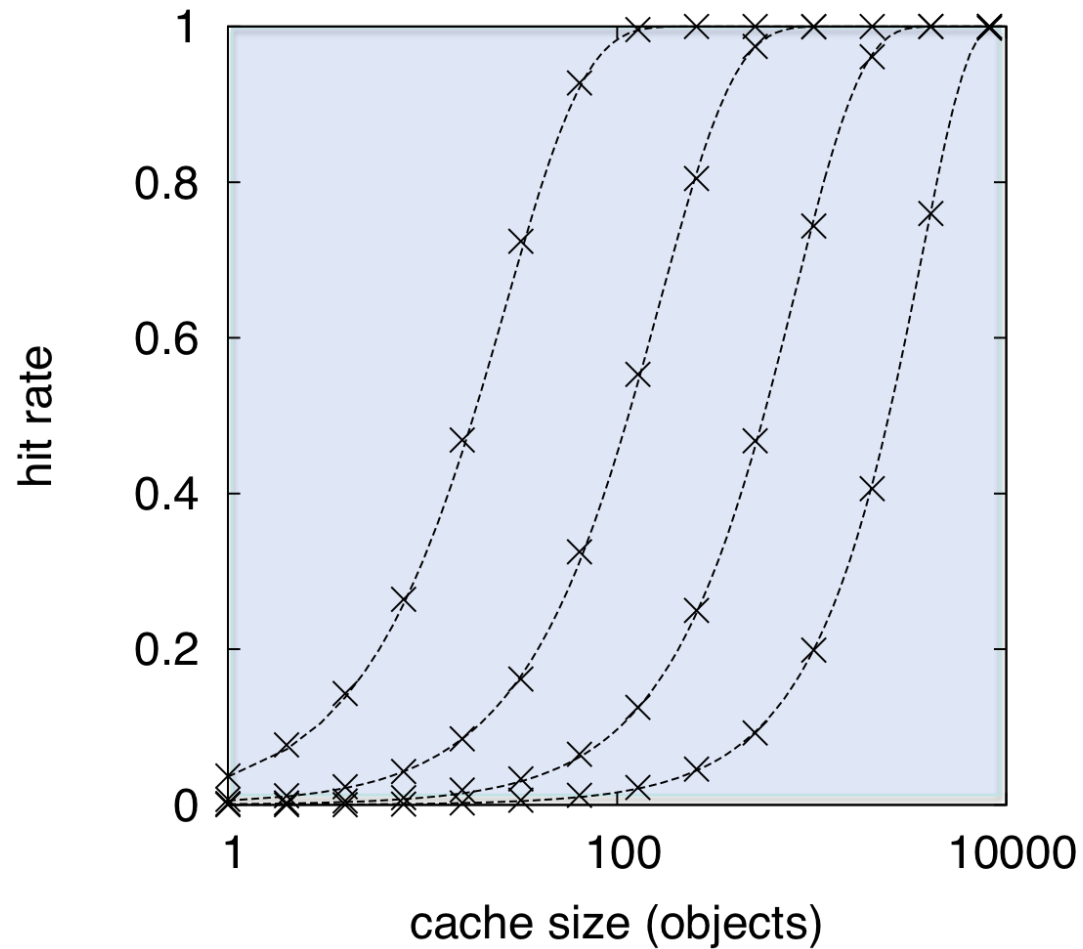
Accuracy of Che approximation

- 10000 objects, Zipf(1.2) popularity



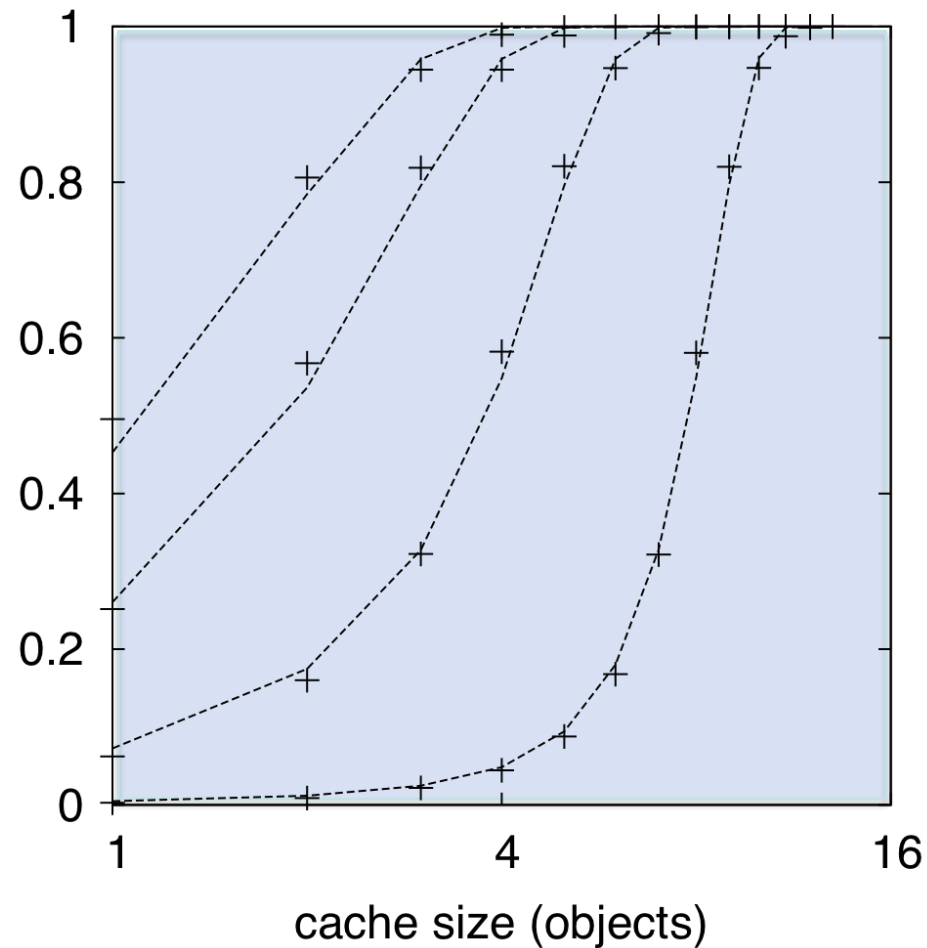
Accuracy of Che approximation

- 10000 objects, Zipf(.8) popularity



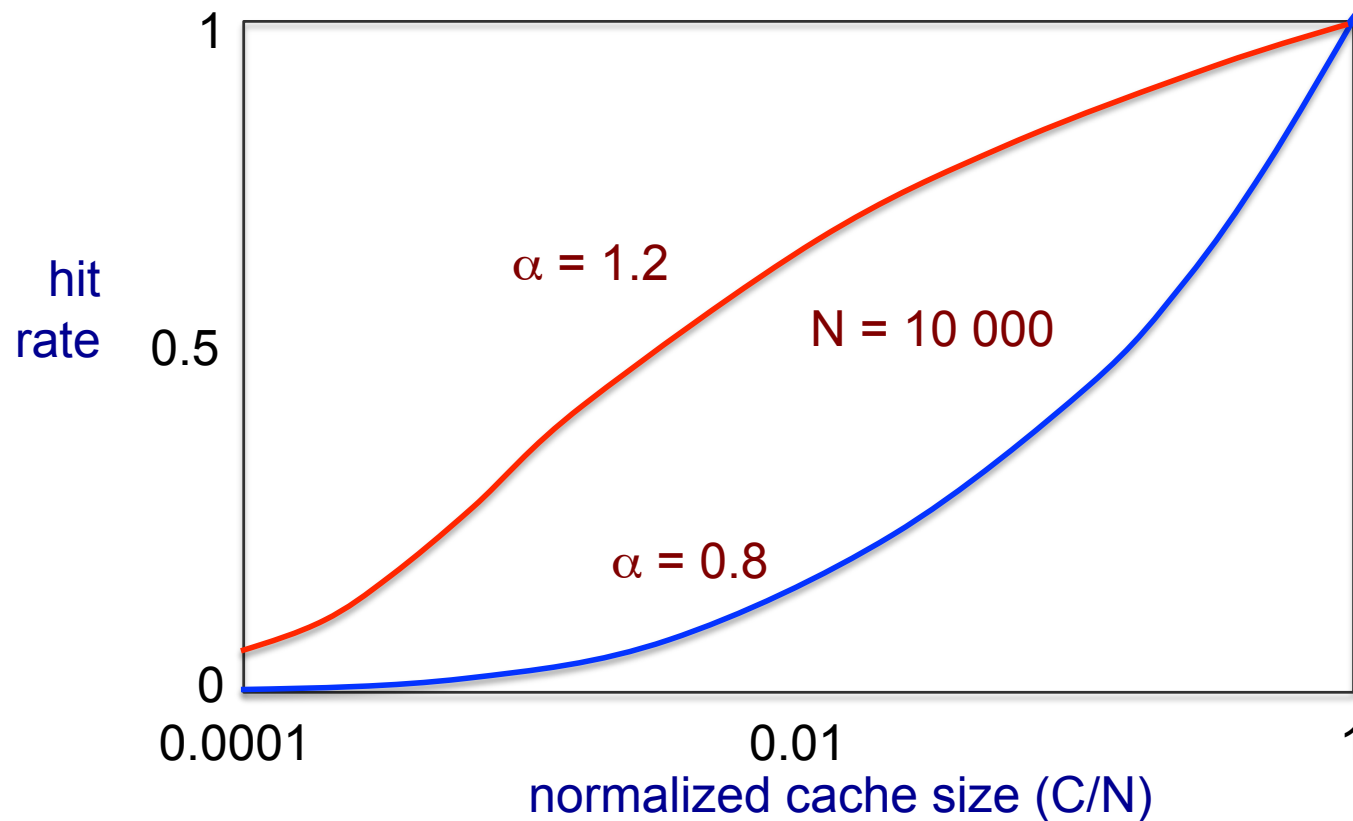
Accuracy of Che approximation

- 16 objects, geometric(.5) popularity



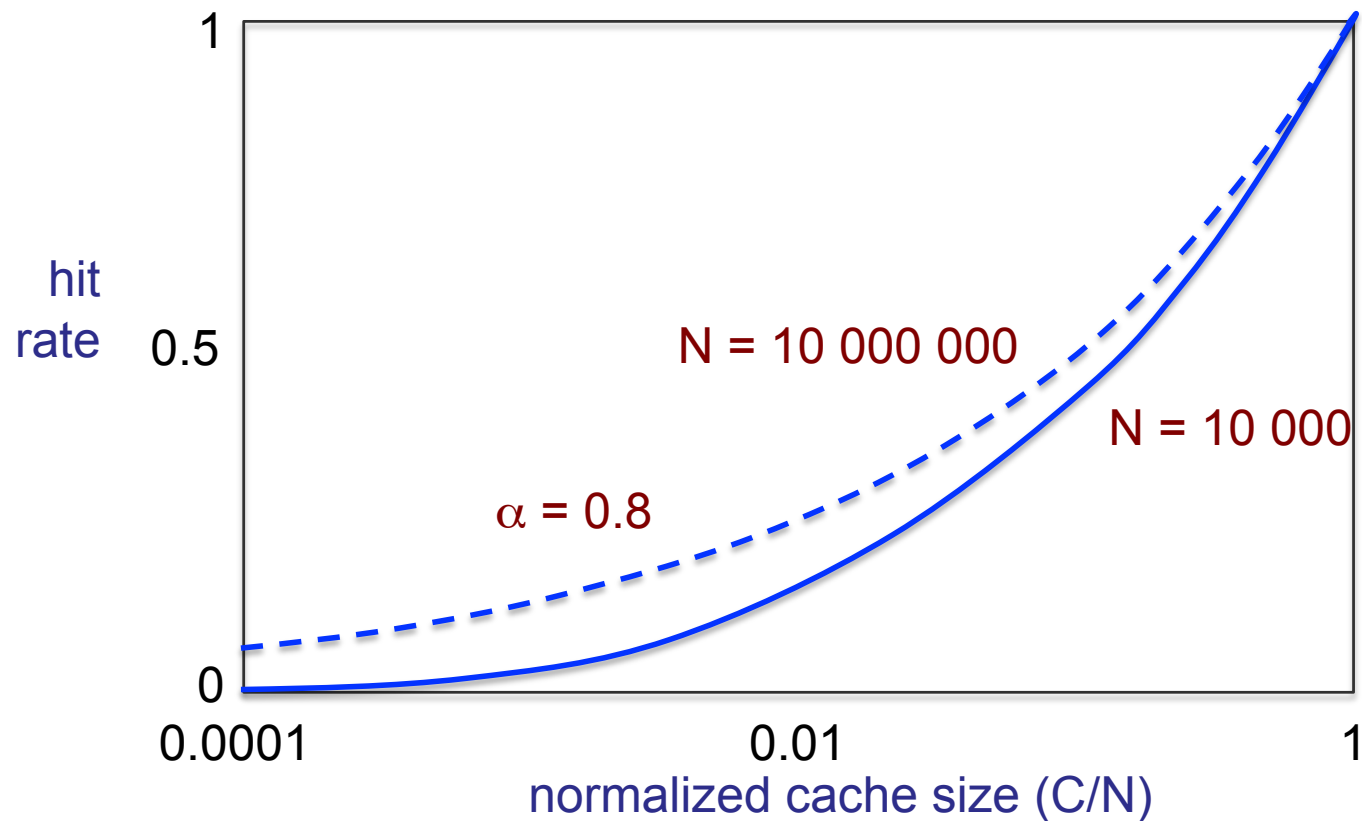
LRU - single cache performance

- strong impact of Zipf parameter α
- strong impact of object population N



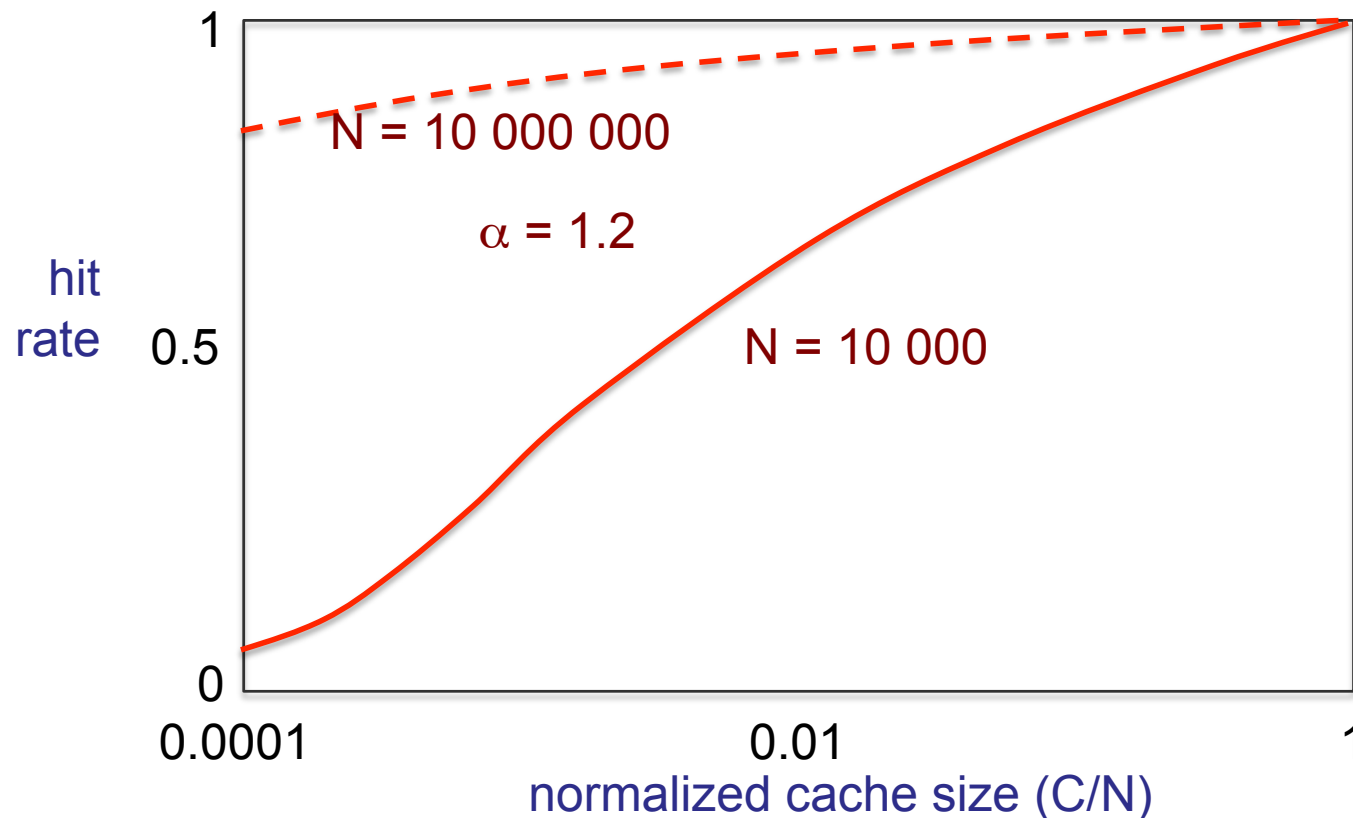
LRU - single cache performance

- strong impact of Zipf parameter α
- strong impact of object population N



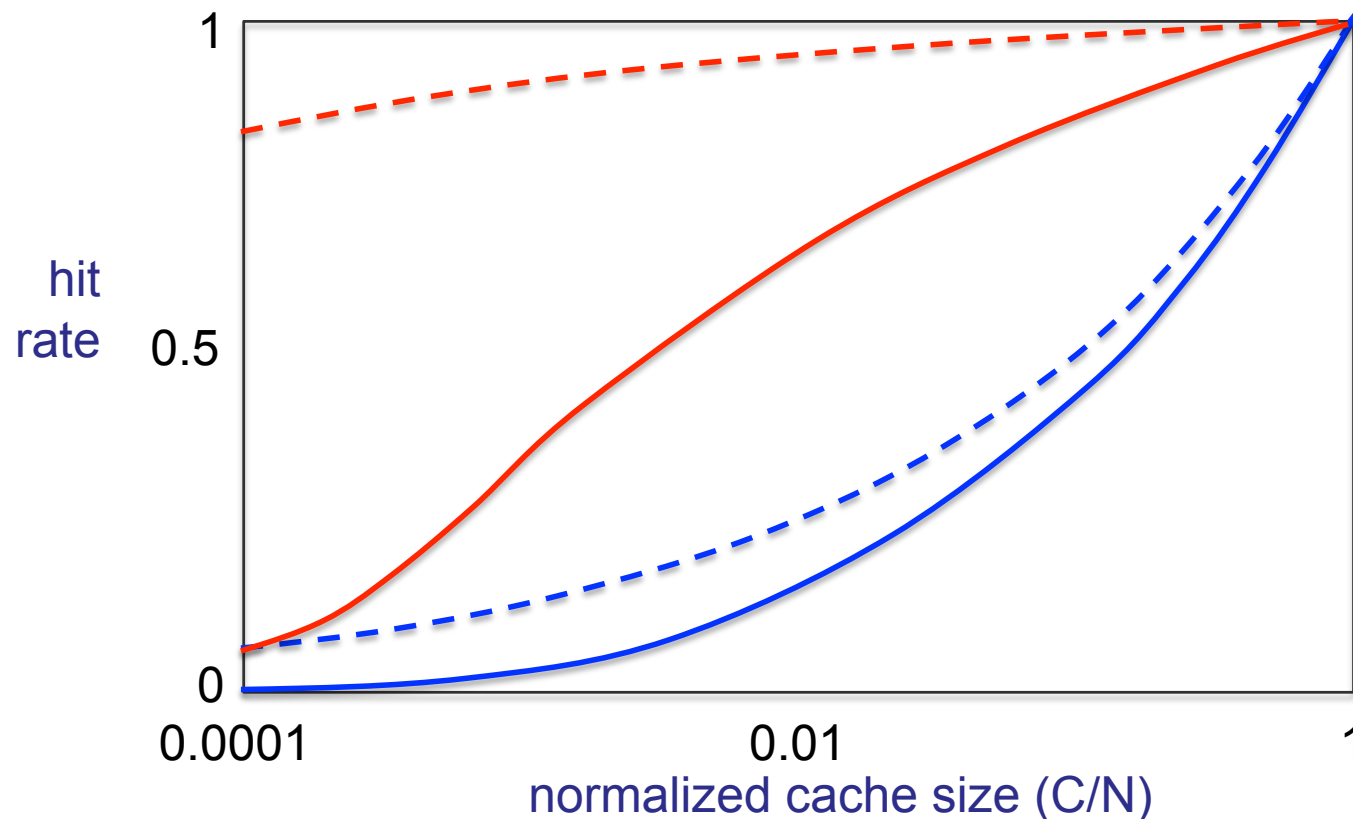
LRU - single cache performance

- strong impact of Zipf parameter α
- strong impact of object population N



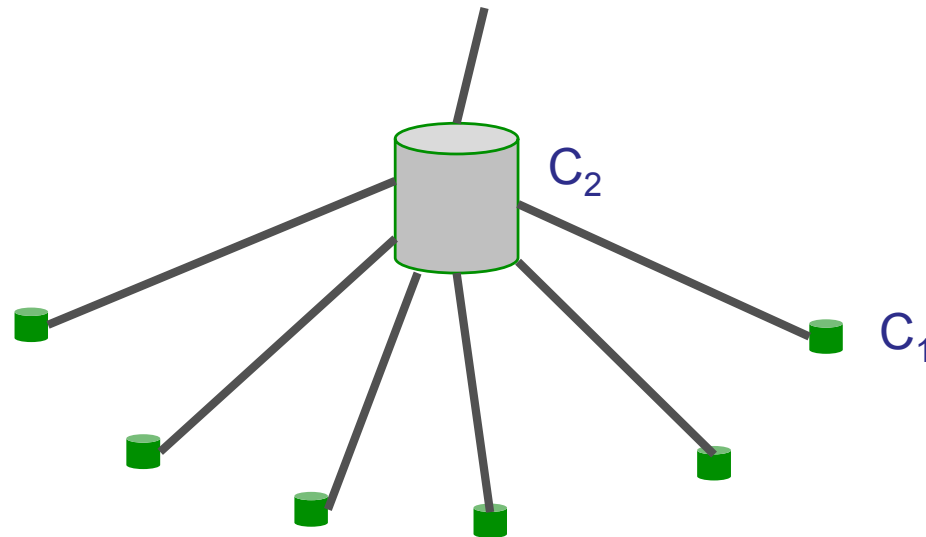
LRU - single cache performance

- strong impact of Zipf parameter α
- strong impact of object population N



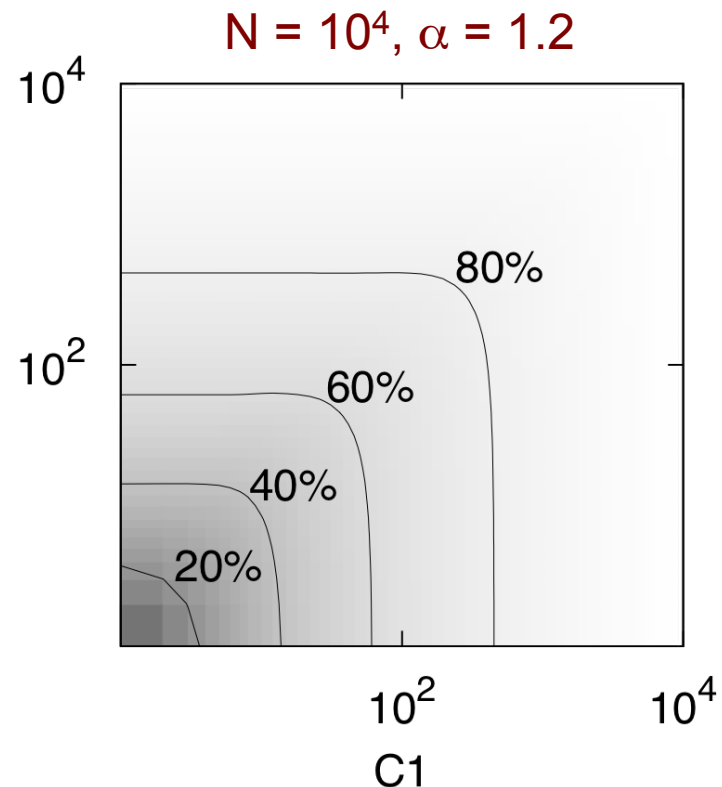
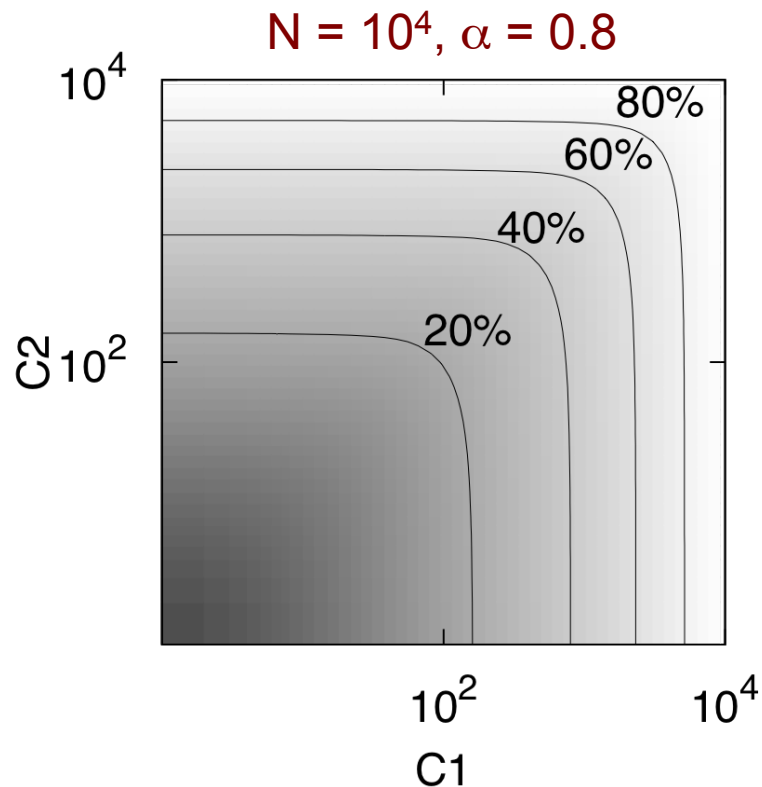
Overall hit rate in 2-layer hierarchy

- assume independent reference model at level 2 and independent cache occupancies (since there are many level 1 caches)
- popularity at level 2
 - $q'(n) = q(n) (1-h(n))$
- apply Che approximation to derive hit rates $h'(n)$
- overall hit rate = $\sum q(n) (h(n) + (1-h(n)) h'(n)) / \sum q(n)$



LRU hit rate as function of C_1, C_2 for $N = 10000$, Zipf(α) popularity

- overall hit rate depends essentially on $C_1 + C_2$
- strong impact of Zipf parameter
 - eg, for VoD content

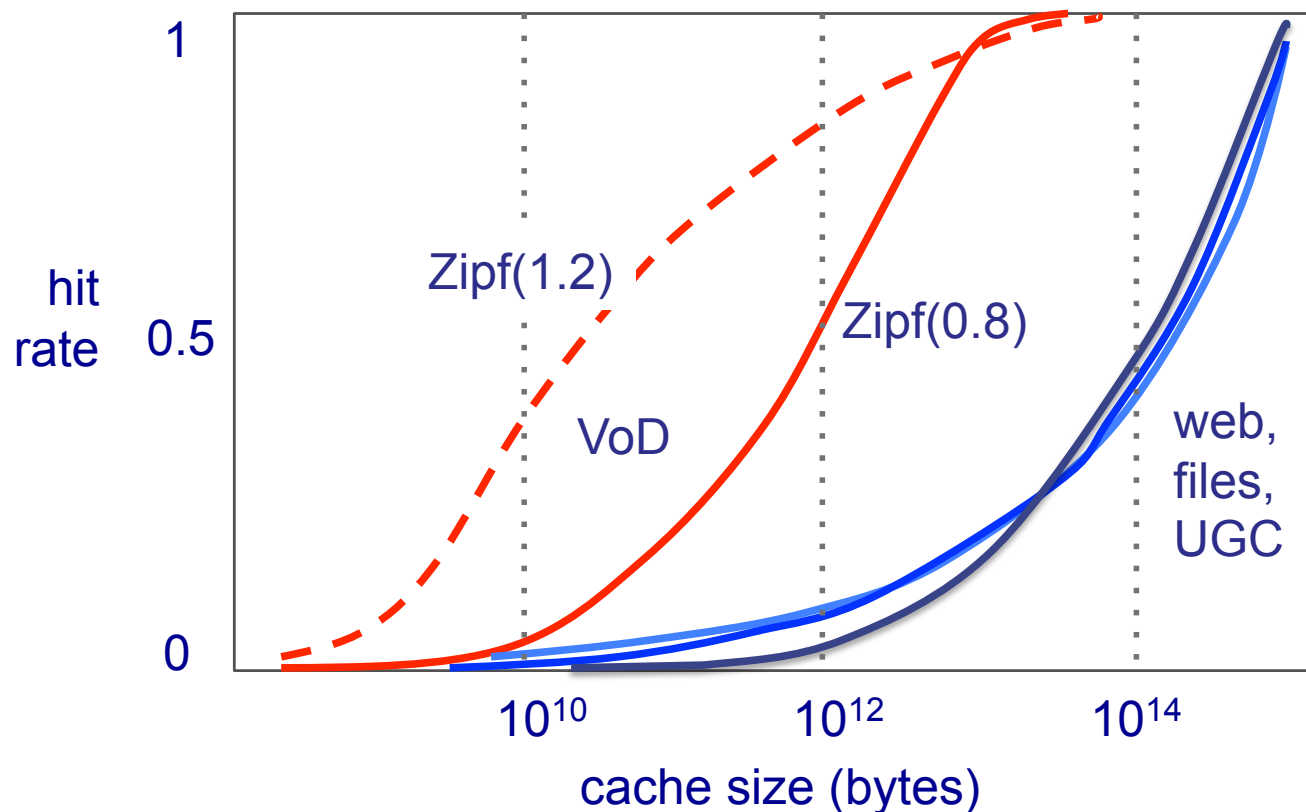


Generalized Che approximation for single cache

- cache size C bytes, popularity of object n of type i is $q_i(n)$, size of this object is $\theta_i(n)$
- assume "independent reference model"
- "critical time" T_C is time for different objects of total size C to be requested, assume $T_C \sim t_C$
- then, hit rate for type i object n is $h_i(n) \approx 1 - \exp\{-q_i(n)t_C\}$
- now, $C = \sum_i \sum_n 1 \{\text{type } i \text{ object } n \text{ present}\} \theta_i(n)$
- taking expectations, $C = \sum_i \sum_n h_i(n) \theta_i(n)$
- solving for t_C yields the $h_i(n)$

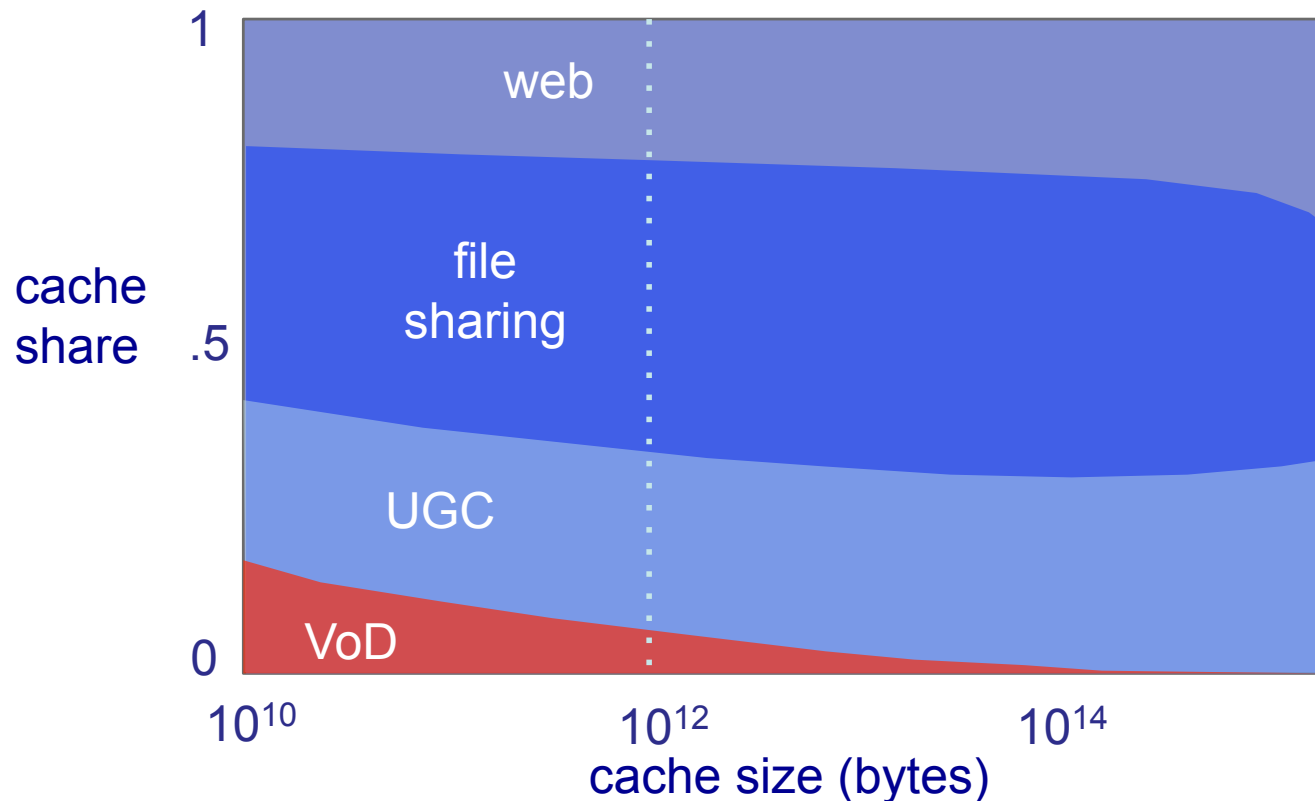
Per content type hit rates

- for web, file sharing and UGC, we need $O(10^{14})$ cache size
- for VoD, $O(10^{12})$ cache size yields high hit rates



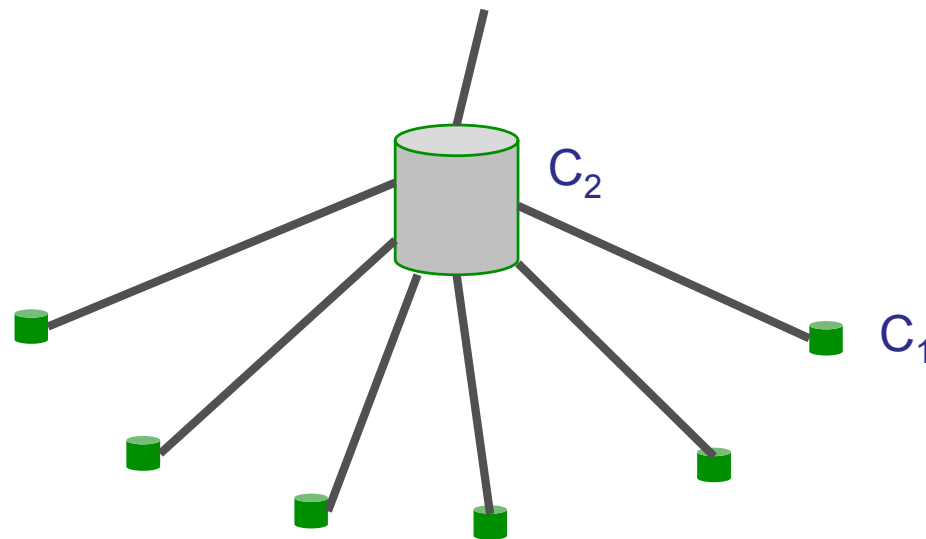
Cache occupancy by type of content - VoD popularity is Zipf(1.2)

- disproportionate cache use, e.g., for 1 TB cache:
 - 95% share for web, files and UGC for only 5% hit rate
 - VoD hit rate is 85% but only uses 5% of cache capacity



Overall hit rates in network

- adapt Che approximation as before to account for
 - mix of types and cache size in bytes
 - popularity filtering at first level cache
- yielding overall per-type hit rate as function of cache sizes in bytes, C_1 and C_2

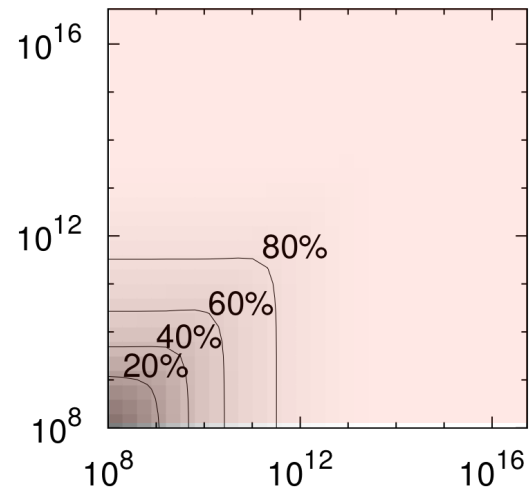
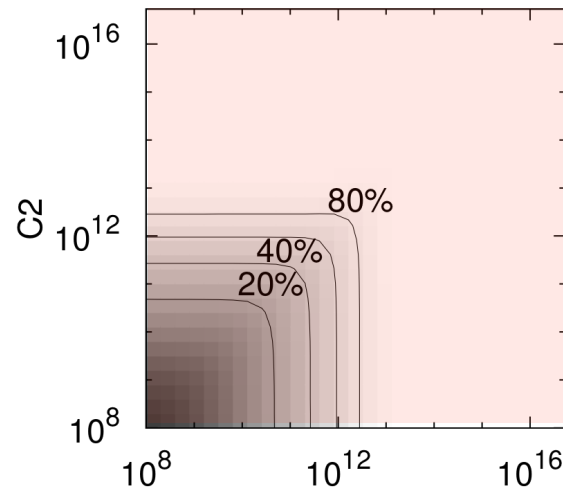


Overall hit rates

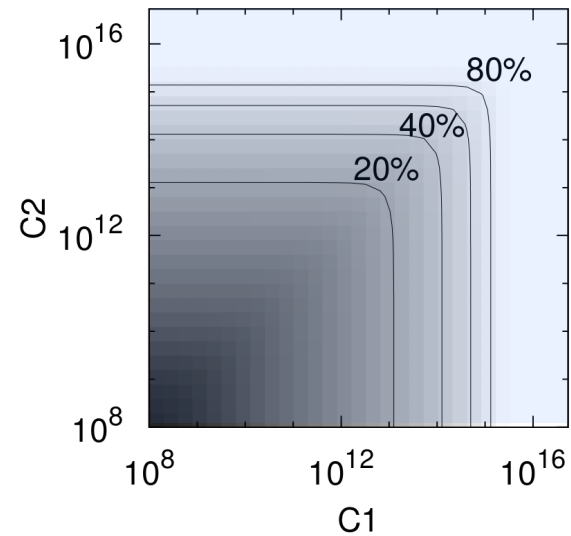
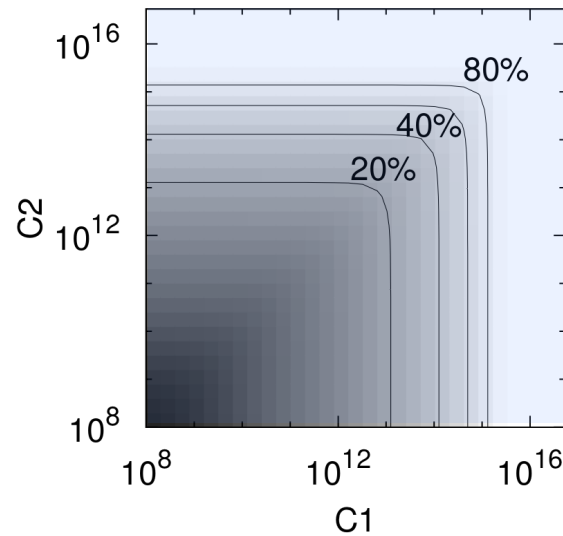
VoD popularity: Zipf(0.8)

Zipf(1.2)

VoD



web,
files,
UGC



Is ubiquitous caching really a good idea?

- level 1 cache cannot be very big, $O(10^{12})$ bytes, say
 - for cheap storage in a large number of access routers
- significant VoD hit rate but heavy pollution from other types
- it would be more effective to specialize level 1 for VoD
- eg, assume $C_1 = 1$ TB, $C_2 = 100$ TB, VoD popularity = Zipf(0.8)

level 1	level 1 hit rate	overall hit rate
shared	17%	50%
VoD only	23%	58%

Conclusions

- cache performance depends on real traffic characteristics
 - Zipf(α) popularity with $\alpha < 1$, large volumes $\sum N_i \theta_i = O(1 \text{ PB})$
- the Che approximation is a versatile tool for LRU performance
 - accounting for mixed content, filtered popularity
 - see Fricker et al., 2012, <http://arxiv.org/pdf/1202.3974v1>
- we need a very big cache to significantly reduce traffic due to web, file sharing and UGC
 - ~100 TB for a 50% hit rate
- a smaller cache is sufficient for VoD
 - ~1 TB for 100% hit rate if level 1 is dedicated to VoD
- the level 2 cache is in fact a network of coordinated caches
 - optimal design is the aim of our ongoing research