From DRAM to SSDs

Challenges with caching @ Facebook

Agenda:

- Introduce 1.
- 2.
- 3.

Sathya Gunasekar Software engineer

acheLib hybrid caching Hybrid caching challenges Call to action

FACEBOOK Infrastructure

Caching is important for social media

Reuse

Users share and consume content through overlapping social circles

Recency

Users interact with recently shared content

Scale & Efficiency

Efficiency is important to scale products to billions of users

Evolution of FB infrastructure

	2005 - 2010	2010 - 2015
Users	10s of millions	100s of million
Product	likes, comments, profile	+ timeline, messaging, videos, photos, location
Infrastructure	Memcached MySQL	fb-Memcached, TAO MySQL, CDN, Haystack

2015 - today

billions

facebook

hundreds of services

Diverse caching needs @ Facebook

Application Domain

4

General purpose KV cache

Social graph cache

CDN caches

Ranking/ML systems

Developer infrastructure

Analytics systems

Architecture

Look-aside

Read & write through

Write back

Should we generalize or specialize ?

Data format

HW

Blob

Structured

Semi-structured

DRAM

SSD

CacheLib – Embedded C++ cache engine

Aggregation point for innovation

CacheLib solves:

- What to cache ? Heuristics 1)
- How to cache ? Storage design 2)

CacheBench – workload benchmarking tool

nost



CacheLib @ Facebook



70+ services, XXX K+ instances, XX PB DRAM, XXX PB NVM

Open sourced by Facebook

www.cachelib.org



CacheLib: Item

7

Unit of object managed in the cache

- Sequence of bytes (<4MB) identified by a key (<255 bytes)
- Exposed through an ItemHandle (shared_ptr<Item>)



CacheLib: API

```
ItemHandle allocate(key, size, ttl)
void* Item::getMemory()
ItemHandle insertOrReplace(handle)
ItemHandle find(key)
bool remove(key)
```



Features

variety of cache policies
memory pools & isolation
resource adaptiveness
zero-copy access
high concurrency
structured data cache

• hybrid cache

Hardware trade-offs

Explore cost and power tradeoffs within a cache system

- marginal cost of an additional hit goes up
 - Hot objects in DRAM
 - Warm objects in NVM



Hybrid Cache goals

Fluid usage of storage mediums based on cost-sensitivity

- DRAM, NVM technologies etc.

Portability of caching applications

- Hide the complexity of hardware technologies behind the API



CacheLib API

DRAM cache NVM cache

Hybrid caches @ Facebook



Social graph



Look-aside key-value cache



CDN

30GB DRAM, 1TB SSD

40GB DRAM, 1TB SSD

30–120 GB DRAM, 2– 8 TB SSD



HDD Storage

10GB DRAM, 400GB SSD

Hybrid Cache efficiency in practice

256GB DRAM -> 60GB DRAM + 2TB SSD

- 10x more cache capacity
- 50% reduction in misses
- 25% reduction power and cost of ownership

Hybrid cache: Challenges

DRAM

mutable bytes small working set infinite R/W BW low & steady latency infinite endurance

m la im ate

NVM

- mutable blocks
- large working set
- limited R/W BW
- latency is complicated
- limited endurance



Hybrid cache: Design

- Portability: Applications always see a DRAM cache
- Items are allocated in DRAM and then evicted to NVM
 - Evicted items can be rejected by NVM admission policy



Latency: accessing Items through find()

- Portability: getMemory() blocks until Item is in DRAM
- Async interface ItemHandle
 - isReady() is Item in DRAM ?
 - wait() wait until Item is in DRAM
 - SemiFuture compatibility



Navy NVM storage engine



Small Item engine design

Accessing billions of objects per TB

- With single IO
- Low DRAM overhead

Design

- Set-associative index (no DRAM)
- Optional DRAM bloom filters
- FIFO eviction (no DRAM)
- No space amp, but large write-amp





Due to endurance constraints, caches can't write all the objects to SSD -> poor miss ratios

Reading an object



RAM (III) Read flash page, if present in filter.



Improving small object caching

Log structured storage ?



Collaboration with Carnegie Mellon University, SOSP'21

To be compute and IO efficient for lookups, log structured cache needs a full DRAM index.

	Set-Associative Cache	Log-Structured Cache
ead	Low	High
ation	High	Low

Kangaroo: Realizing the best of both worlds

Kangaroo: SOSP'21 Design



2)Periodically flush objects from KLog to KSet

3)Move all objects in KLog that map to the same bucket in KSet

Summary: call to action

Hybrid caching is gaining popularity, lots of new challenges

Caching small objects on block devices

- With lower DRAM overheads
- With lower endurance overheads
- With more IO efficiency

New hardware technologies on the horizon





collaborate among industry & academia

acheLibffers an OSS platform to collaborate and explore solutions

www.cachelib.org

FACEBOOK GØ O

