# We're still having FAWN

David Andersen, Michael Kaminsky*, Michael A. Kozuch*, Padmanabhan Pillali* ,Vijay Vasudevan, Amar Phanishayee, Lawrence Tan, Jason Franklin, Iulian Moraru, Sang Kil Cha, Hyeontaek Lim, Bin Fan, Reinhard Munz, Nathan Wan, Jack Ferris, Hrishikesh Amur**, Wolfgang Richter, Michael Freedman***, Wyatt Lloyd***,  Dong Zhou

Carnegie Mellon University      *Intel Labs Pittsburgh
** Princeton University   *** Georgia Tech

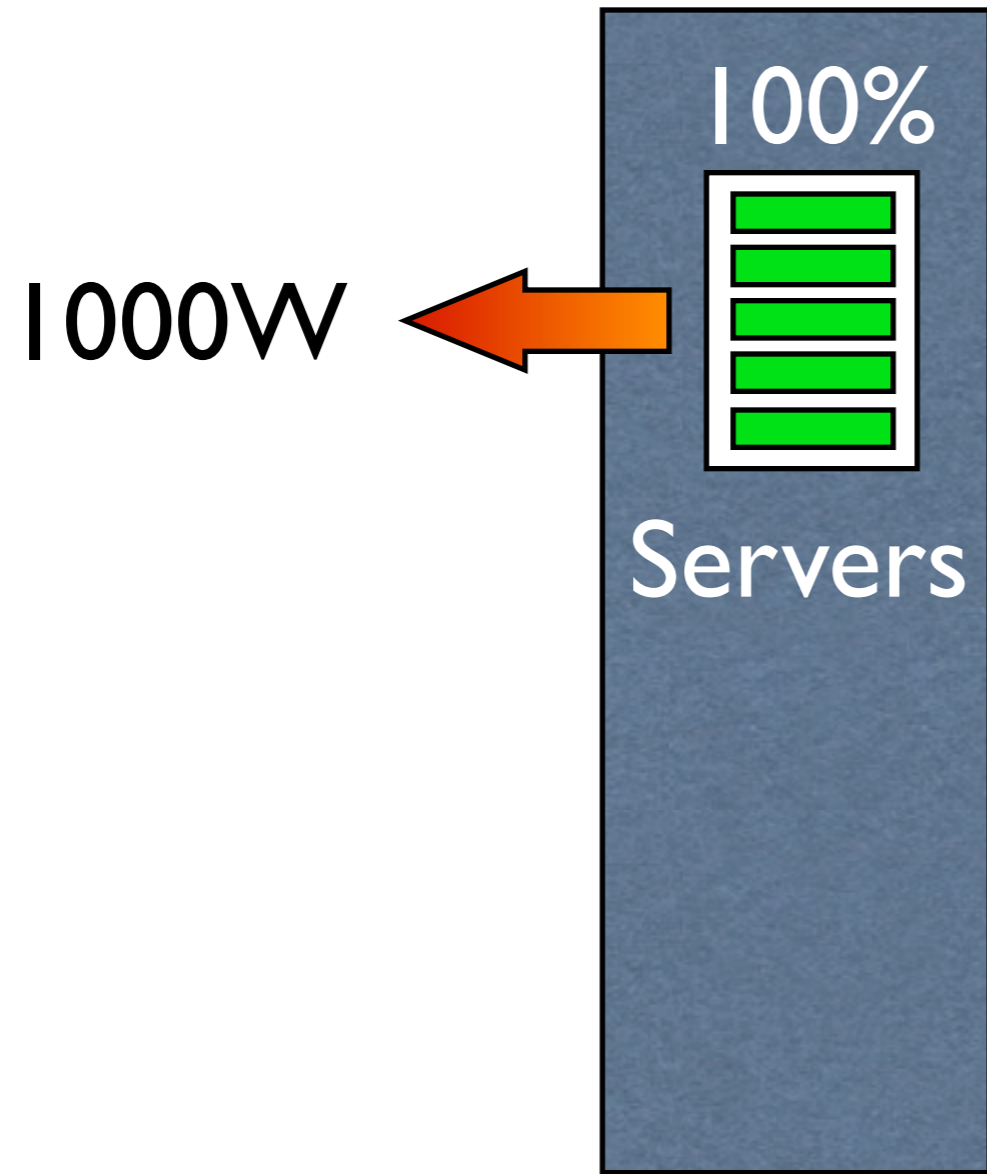# Power limits computing

# Power limits computing

# Power limits computing

# Power limits computing
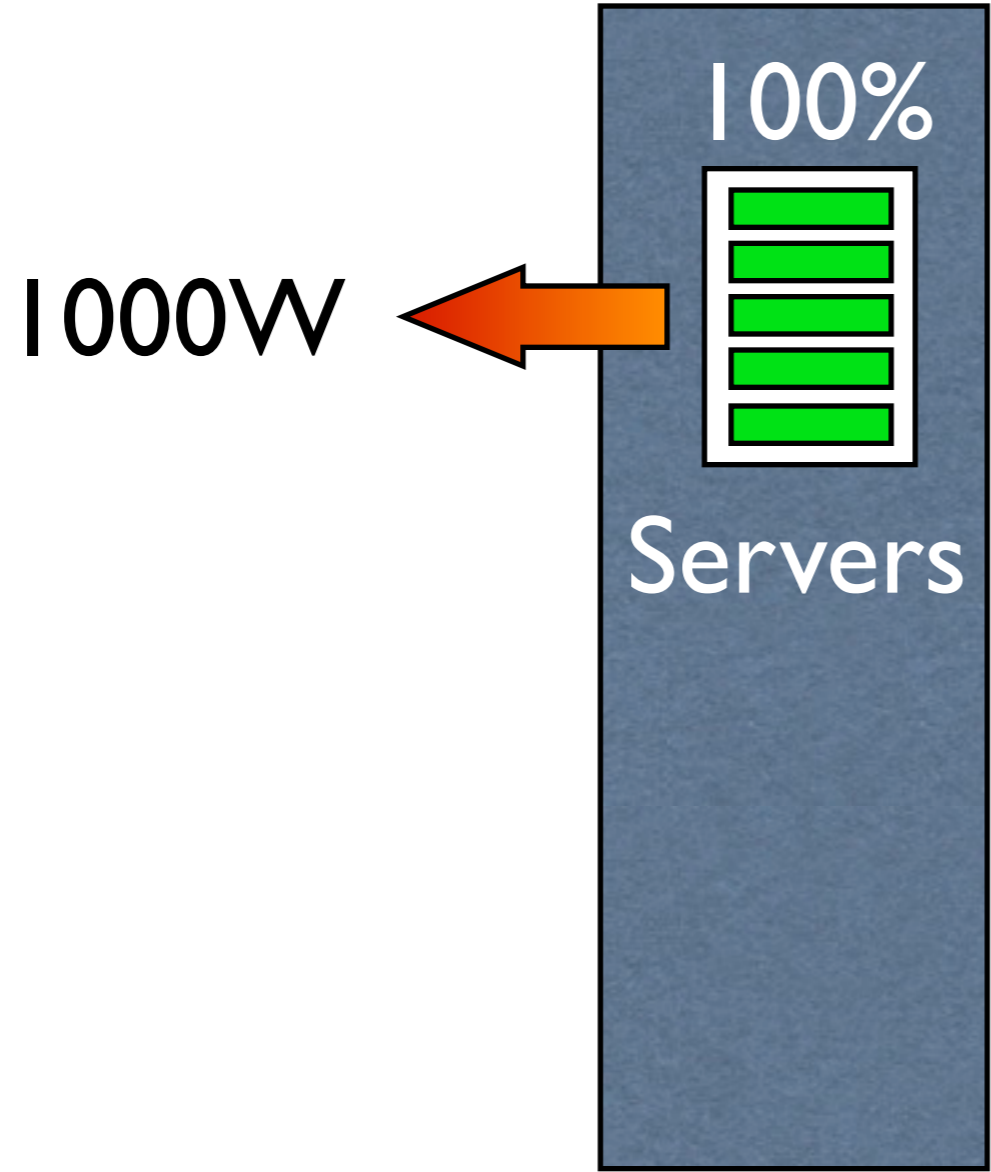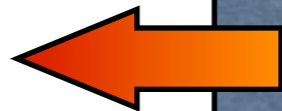
100%

1000W

Servers

2000W

1000W

100%

Servers

**Infrastructure: PUE**
2005: 2–3
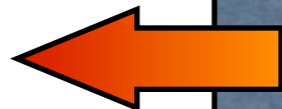2012: ~1.1
*Leave it to industry*

100%

1000W

Servers

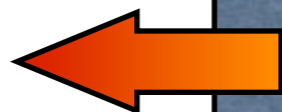Example:

1.6 GHz Atom Z2460

1 GB LPDDR2

# Two pillars

- Gigahertz costs twice:
    - Once for the switching speed
    - Once for the memory wall
- Memory capacity costs (at least) once:
    - Longer buses < efficient

# "Wimpy" Nodes

1.6 GHz Dual-core Atom
**32-160 GB Flash SSD**
**Only 1 GB DRAM!**

*"Each decimal order of magnitude increase in parallelism requires a major redesign and rewrite of parallel code"* - *Kathy Yelick*

The FAWN Quad of Pain

Load Balancing

Parallelization

Bigger Clusters

Wimpy Nodes

Hardware Specificity

Memory Capacity

# It's not just masochism

Moore                                    Dennard



(Figures from Danowitz, Kelley, Mao, Stevenson, and Horowitz:  CPU DB)

*All* systems will face this challenge over time

# FAWN:
# It started
# with a key-value store

# Small record, random access

# Small record, random access

Select name,photo from users where uid=513542;

# Small record, random access



Select name,photo from users where uid=818503;

# Small record, random access



Select name,photo from users where uid=468883;

# Small record, random access



Select wallpost from posts where pid=13821828188;

Select name,photo from users where uid=124111;

# Small record, random access

Select wallpost from posts where pid=89888333522;

Select wallpost from posts where pid=13821828188;

Select name,photo from users where uid=474488;

Select name,photo from users where uid=124566;

Select name,photo from users where uid=124111;

Select wallpost from posts where pid=12314144887;

Select name,photo from users where uid=997788;

Select wallpost from posts where pid=738838402;

Select name,photo from users where uid=357845;

# FAWN-DS

key-value backend store
one node
optimized for wimpy
nodes and flash

Fawn-DS

A cluster-distributed
key-value store
Minimizes work on churn

Fawn-KV

Fawn-DS

Fawn-DS

Fawn-DS

# FAWN-DS FAWN-KV SILT

backend store
hyper-optimized
for low DRAM
and large flash

Fawn-KV

SILT

SILT

SILT

SILT

Cache

SILT

SILT

SILT

Provable load
balancing
using a tiny cache

Parallel, fast, memory-efficient memcached using *optimistic cuckoo hashing*

Cuckoo Cache

SILT

SILT

SILT

# After first victory, moved to Atom+SSD

| | | |
|---|---|---|
| Geode 500Mhz | 256MB | 4GB CF Card ~2k IOPS |
| 6x | 8x | 30-60x |
| Atom 1.6 Ghz single-core | 2GB | 120GB SSD ~60k IOPS |

Fawn-DS

Fawn-DS

Fawn-KV

Fawn-DS

SILT

SILT

SILT

backend store
hyper-optimized
for low DRAM
and large flash

Fawn-KV

# Flash Must be Used Carefully

| Random reads / sec | 48,000 |
| --- | --- |

→ Fast, but not <u>THAT</u> fast

# Flash Must be Used Carefully

| Random reads / sec | 48,000 |
|---|---|

→ Fast, but not <u>THAT</u> fast

| $ / GB | 1.83 |
|---|---|

→ Space is precious

# Flash Must be Used Carefully

| Random reads / sec | 48,000 |
|---|---|

→ Fast, but not THAT fast

| $ / GB | 1.83 |
|---|---|

→ Space is precious

Another long-standing problem:
**random writes** are slow and bad for flash life (wearout)

# Three Metrics to Minimize

**Memory overhead** = Index size per entry

- Ideally 0 bytes/entry (no memory overhead)

# Three Metrics to Minimize

**Memory overhead** = Index size per entry

- Ideally 0 bytes/entry (no memory overhead)

**Read amplification** = Flash reads per query

- Limits **query throughput**
- Ideally 1 (no wasted flash reads)

# Three Metrics to Minimize

**Memory overhead** = Index size per entry

- Ideally 0 bytes/entry (no memory overhead)

**Read amplification** = Flash reads per query

- Limits **query throughput**
- Ideally 1 (no wasted flash reads)

**Write amplification** = Flash writes per entry

- Limits **insert throughput**
- Also reduces **flash life expectancy**
  - Must be small enough for flash to last a few years

SkimpyStash

FAWN-DS
FlashStore
HashCache
BufferHash

Memory efficiency

High performance

# (static) "External Dictionary"

DRAM
Index

Flash
Data

- Prior state of the art: "EPH":  ~3.8 bits/entry
- Ours:  Entropy-coded tries,      ~2.5 bits/entry

- Important considerations:
  - Construction speed;  query speed
  - Aw, it's read-only... [need a system; built it]

Workload: 90% GET (100~ M keys) + 10% PUT

Caveat: Not on wimpies. Still working on reducing CPU cost! :-)

# And now... Load imbalance

- Distributed key-value system

# And now... Load imbalance

- Distributed key-value system

**Atom CPU**

**Backend1**

**FrontEnd**

# And now... Load imbalance

- Distributed key-value system

**Atom CPU**

**SSD**

**Backend1**

**FrontEnd**

# And now... Load imbalance

- Distributed key-value system

# And now… Load imbalance

- Distributed key-value system

**Atom CPU**

**SSD**

**1. get(key)**

FrontEnd

Backend1

Backend2

Back.. 85

Back.. 88

# And now... Load imbalance

- Distributed key-value system

**Atom CPU**

**SSD**

Backend1

1. `get(key)`

FrontEnd

2. `BackendID=hash(key)`

Backend2

Back.. 85

Back.. 88

# And now… Load imbalance

- Distributed key-value system

**Atom CPU**

**SSD**

Backend1

Backend2

**1. get(key)**

FrontEnd

**3. val=lookup(key)**

**2. BackendID=hash(key)**

Back.. 85

Back.. 88

# And now… Load imbalance

- Distributed key-value system

Atom CPU

SSD

**Backend1**

**Backend2**

3. `val=lookup(key)`

**Back.. 85**

**Back.. 88**

1. `get(key)`

**FrontEnd**

4. `return val`

2. `BackendID=hash(key)`

# And now... Load imbalance

- Distributed key-value system

# And now... Load imbalance

- Distributed key-value system
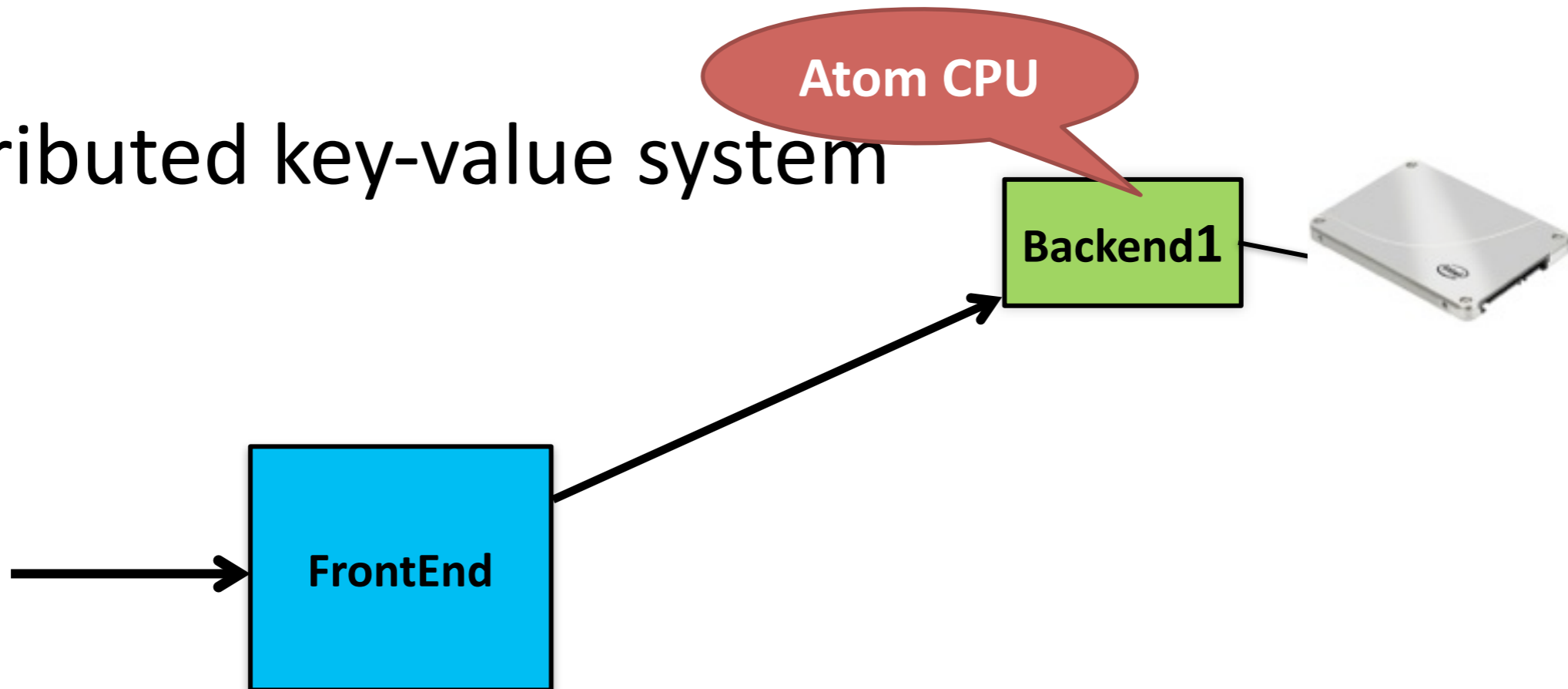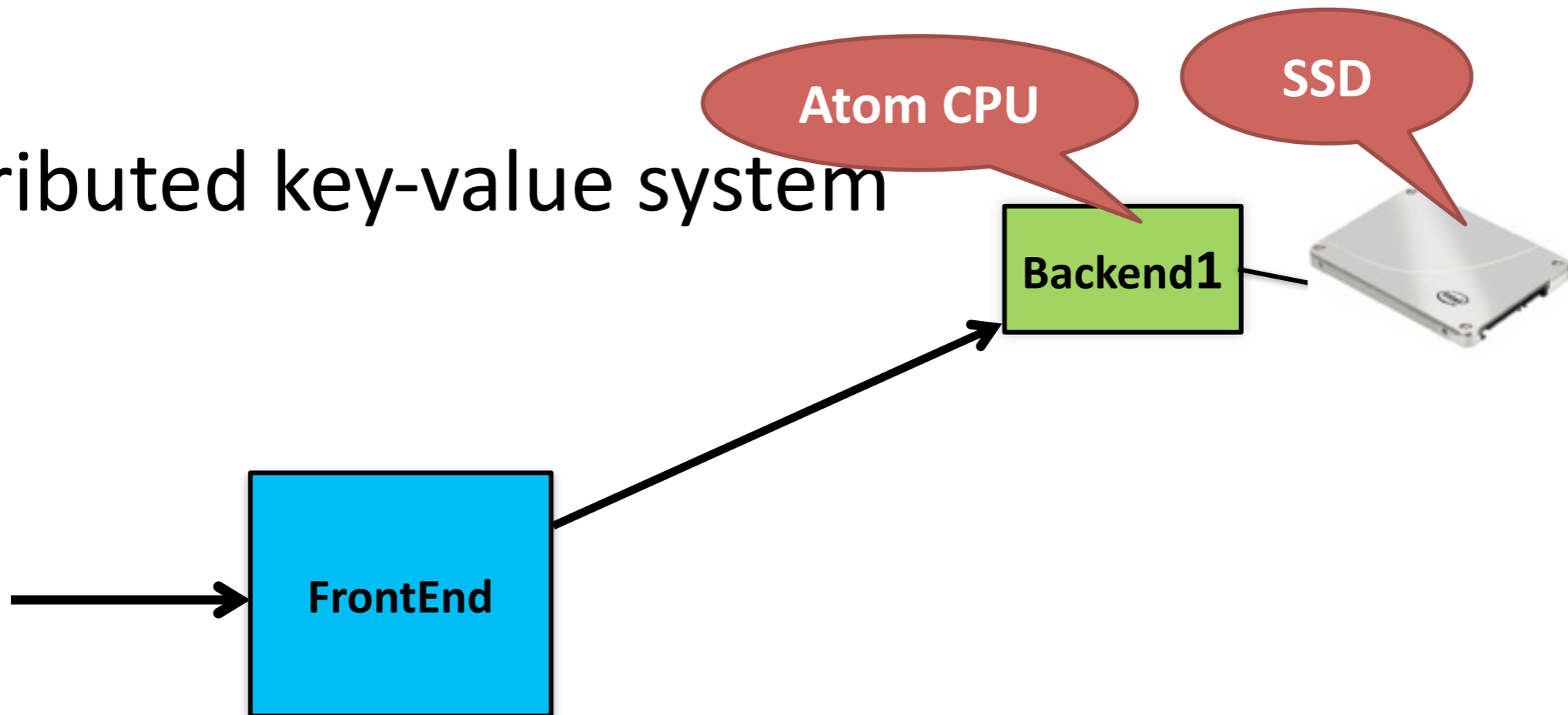


**Atom CPU**

**SSD**
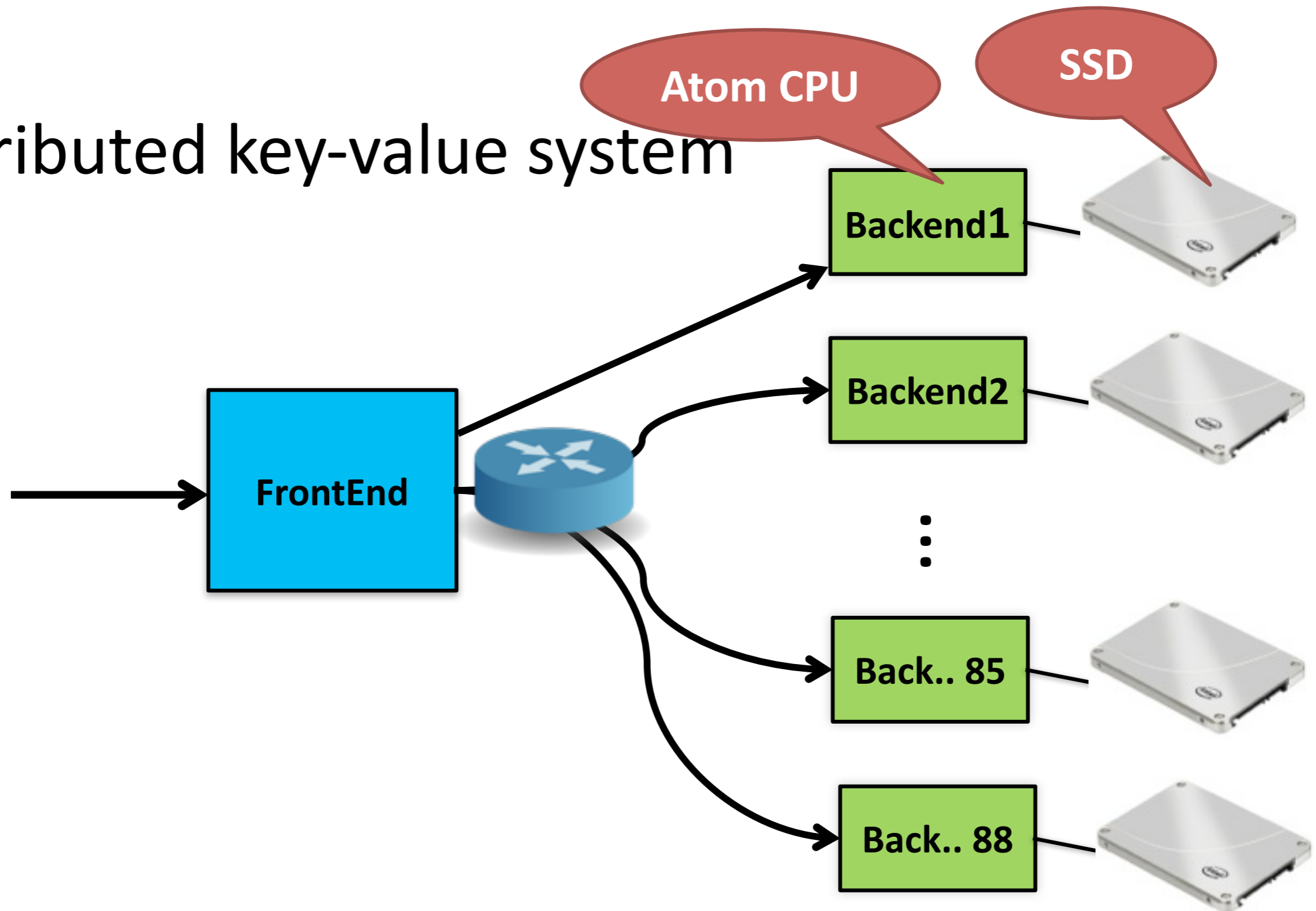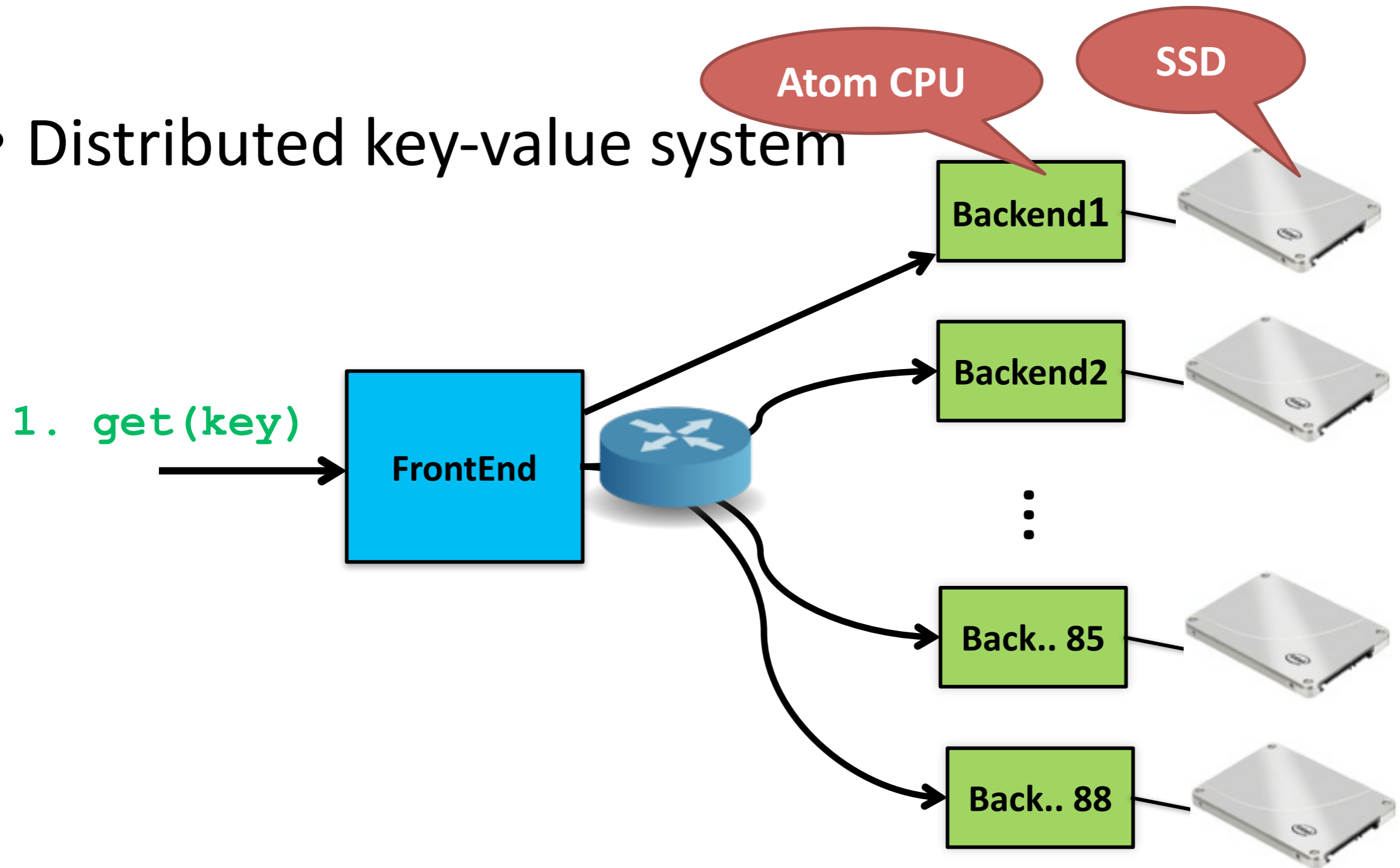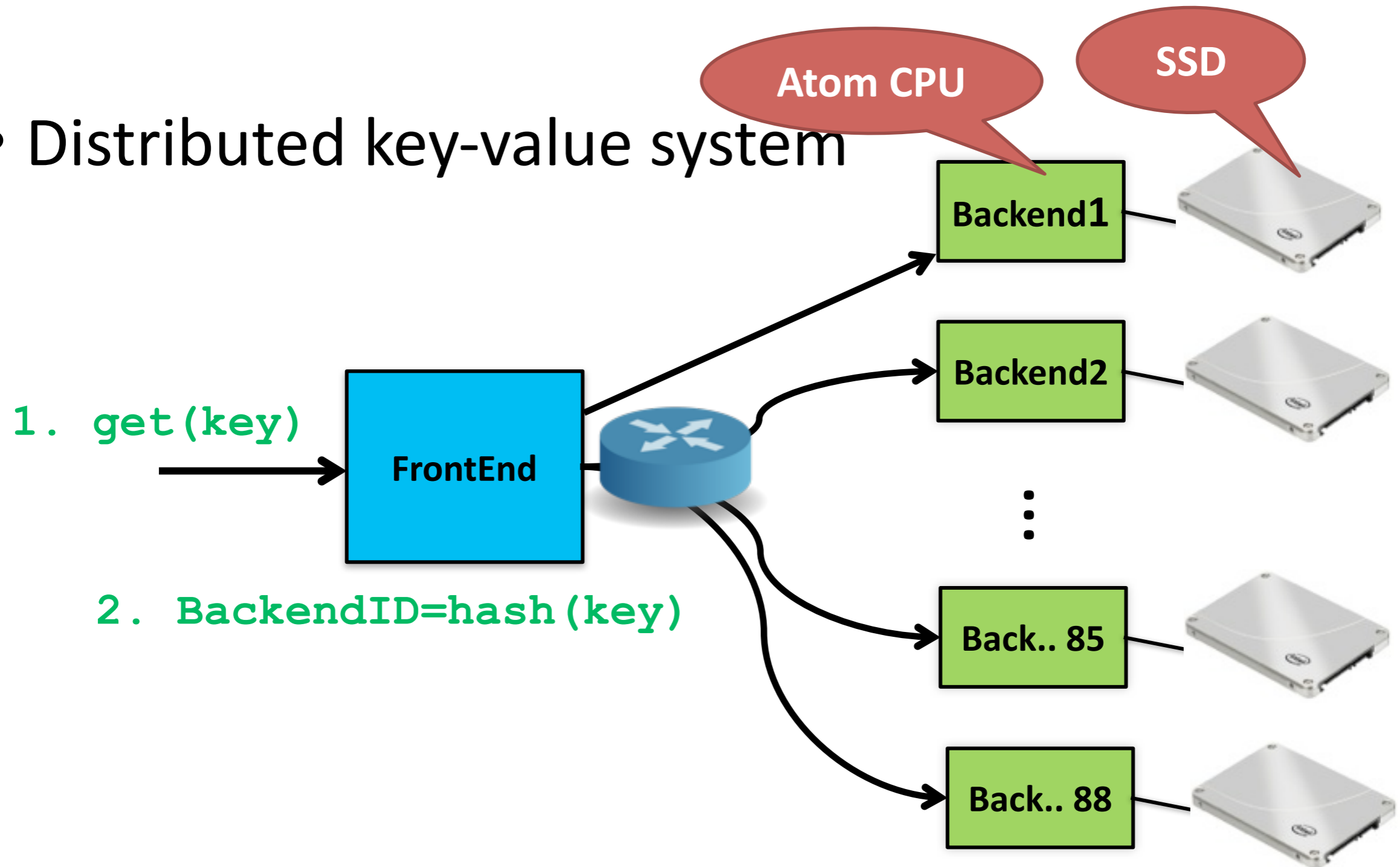
SLA: 850,000 queries/sec

10,000 queries/sec

Backend1

Backend2

Back.. 85

Back.. 88

FrontEnd

1. `get(key)`

4. `return val`

2. `BackendID=hash(key)`

3. `val=lookup(key)`

# Measured tput on FAWN testbed



23

Overall throughput (KQPS) vs n: number of nodes

- uniform
- Zipf (1.01)
- adversarial

*Queries*

FrontEnd

Backend1

Backend2

Backend8

Backend8

How many items to cache?

# small/fast cache is enough!



**Queries**

**FrontEnd**

**cache**

**Backend1**

**Backend2**

⋮

**Backend8**

**Backend8**

We prove that, for *n* nodes
- Only need to cache O(*n log n*) most popular entries
- worst case perf. =
  *(1 - ε)* * n * single node capacity

25

# small/fast cache is enough!

E.g., for 1KB (k,v) pair, 85 nodes,
3MB needed, fitting in CPU L3 cache

*Queries*

**Front**

**cache**

**Backend1**

**Backend2**

**Backend8**

**Backend8**

We prove that, for *n* nodes
- Only need to cache O(*n log n*) most
  popular entries
- worst case perf. =
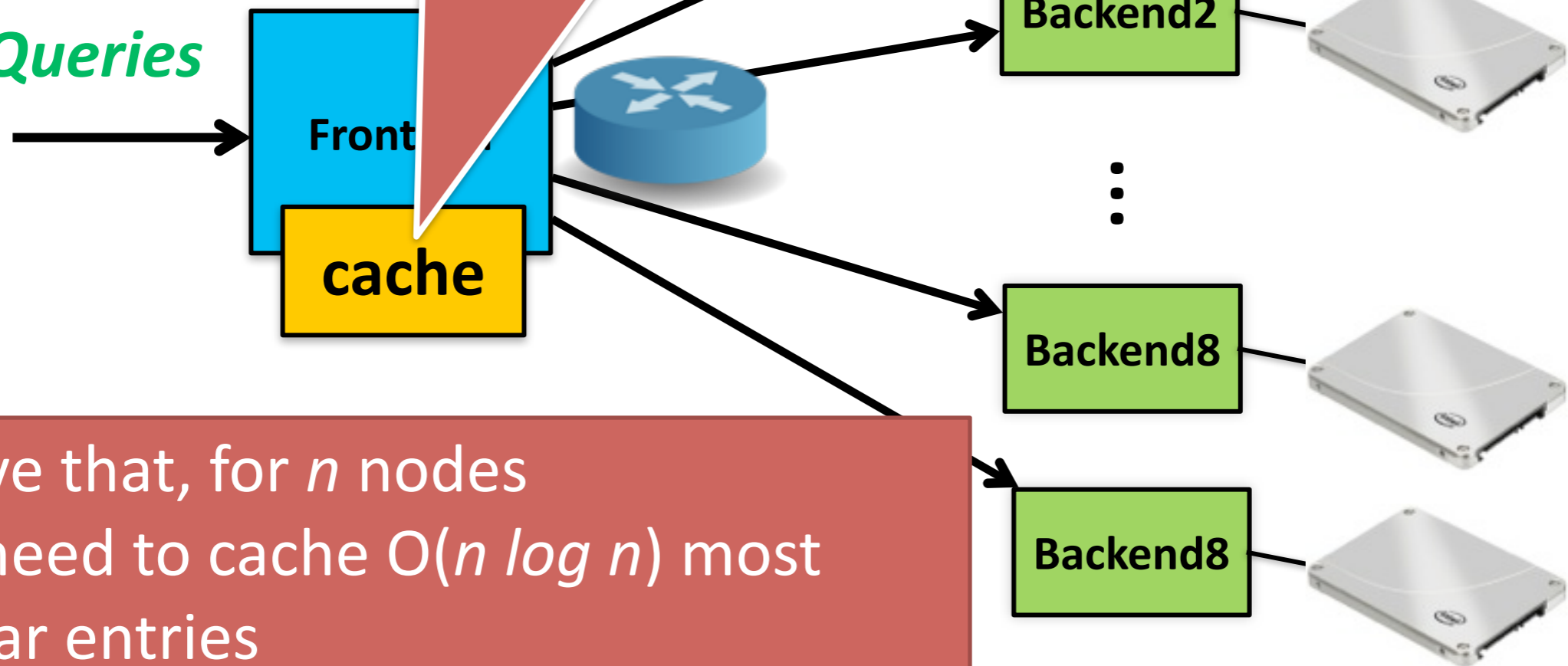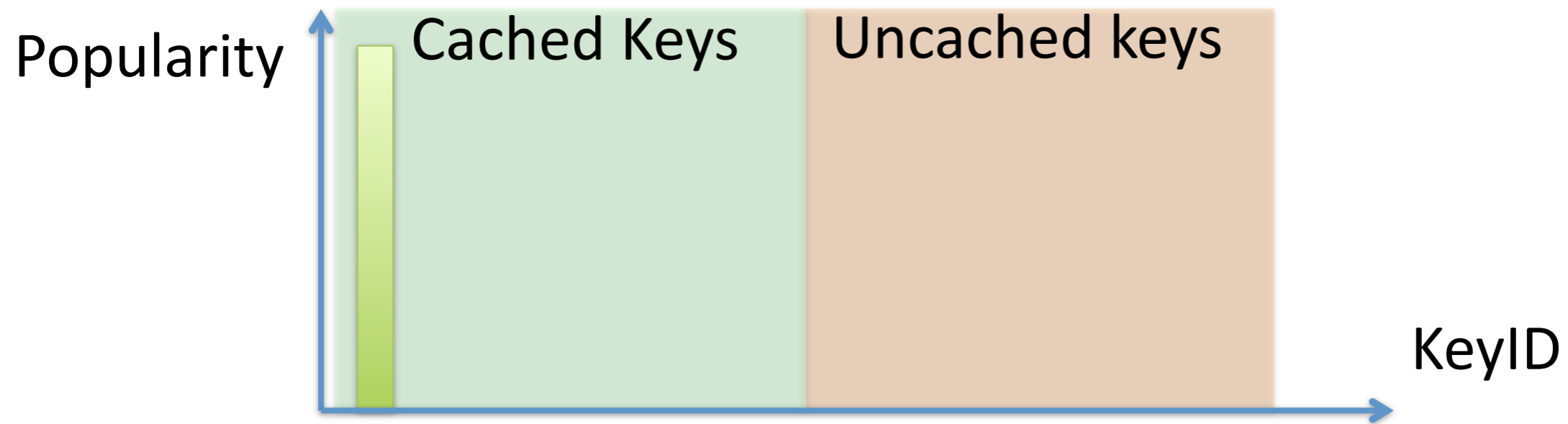*(1 - ε)* * n * single node capacity

# Cache forces near-uniform dist.

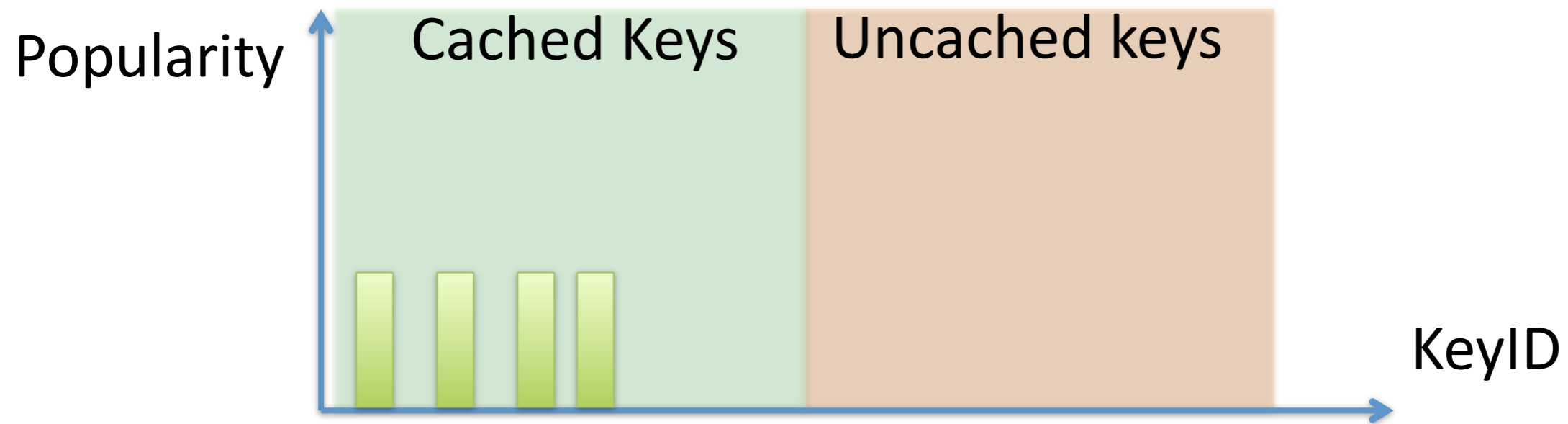# Cache forces near-uniform dist.

# Cache forces near-uniform dist.

# Cache forces near-uniform dist.

# Worst case?  Now best case

# Thus...

FAWN-DS    FAWN-KV    SILT    Small Cache    **Cuckoo**

"Wimpy" servers

"Brawny" server

Cache

SILT

SILT

SILT

FAWN-DS   FAWN-KV   SILT   Small Cache   Cuckoo

"Wimpy" servers

"Brawny" server

Insanely Fast
Cache

SILT

SILT

SILT

Optimistic Cuckoo
Hashing

# Hashing again

- Well known, old technique, nothing new.... right?

# Memory efficiency vs speed

**Linear Probing**

| key | dat |
|-----|-----|
|  |  |
|  |  |
| key3 | dat3 |
|  |  |
|  |  |
| key2 | dat2 |

Wastes 50% of slots

**Chaining**



Key
Dat

Key8
Dat8

Slow: Pointer chasing
Overhead: Pointer space

# Cuckoo Hashing



Item

Uses 93% of slots

**Fast, compact, but ... Single threaded!**

# Optimistic Cuckoo

- Multiple reader, lock-free, single-writer

- Applied in Memcached - *Huge* speedup

- (Talk to Bin @ posters! :)

**MemC3**
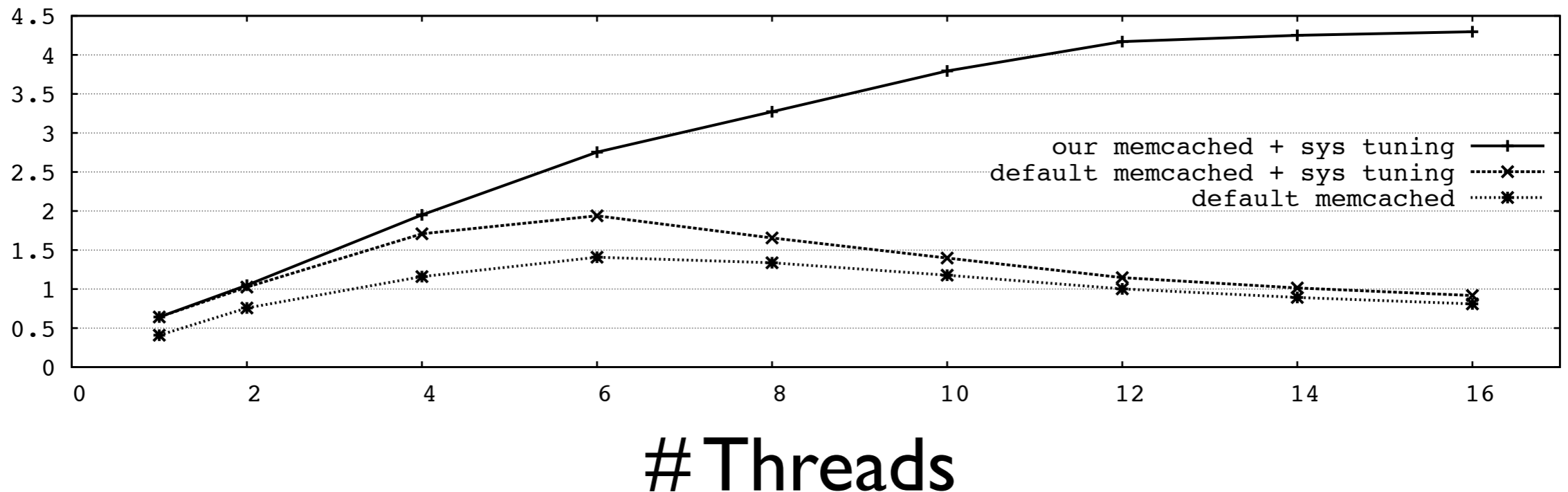


Millions of Reqs/Sec vs # Threads

our memcached + sys tuning
default memcached + sys tuning
default memcached

FAWN-DS    FAWN-KV    SILT    Small Cache    Cuckoo

"Wimpy" servers [FAWN, SOSP 2009]

"Brawny" server



Insanely
Fast Cache

SILT

SILT

SILT

[SILT, SOSP 2011]

O(N log N) ["small cache" socc 2011]

Multi-reader
parallel cuckoo
hashing    [under submission]

Entropy-coded tries [SILT + under submission]

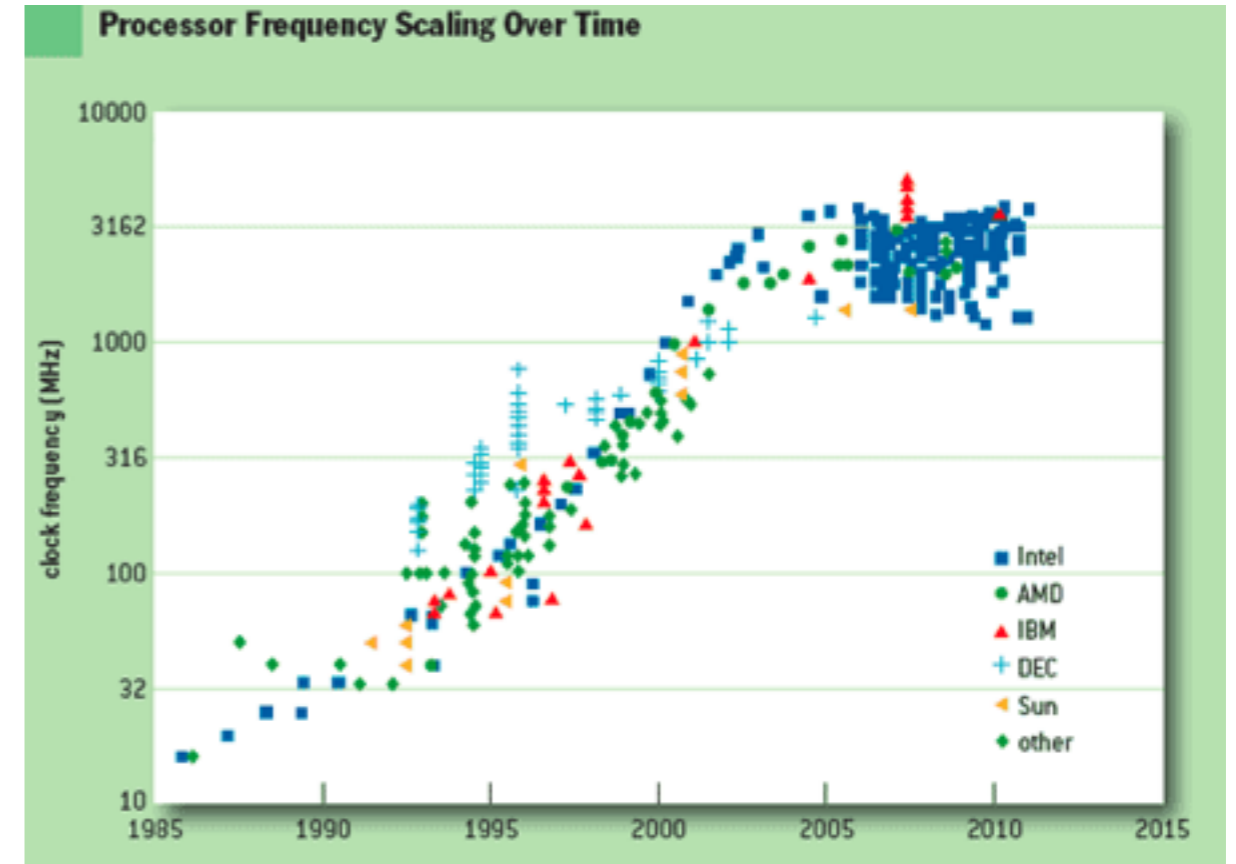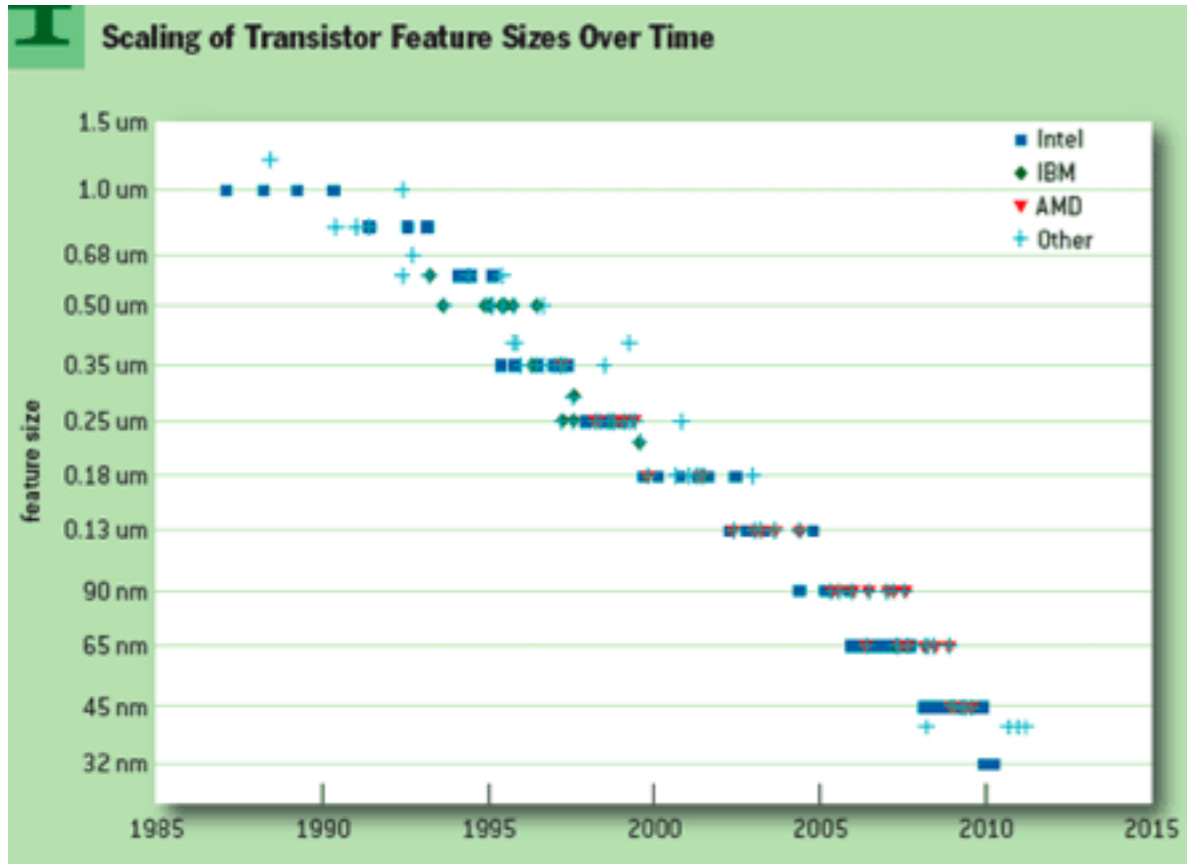Partial-key cuckoo hashing

Cuckoo filter

# Moore

# Dennard

Scaling of Transistor Feature Sizes Over Time

Processor Frequency Scaling Over Time

highly parallel, lower-GHz, (memory-constrained?):

*Architectures, algorithms, and programming*