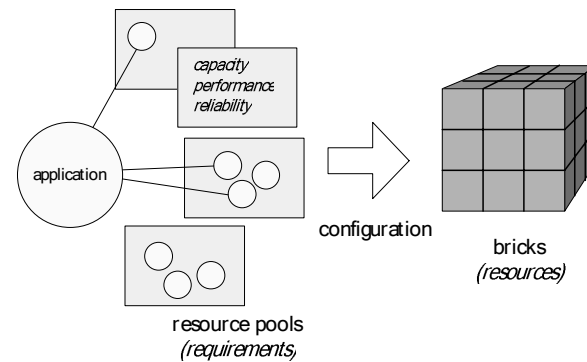*Zygaria: storage performance
as a managed resource*

*Richard Golding, Theodore Wong, Caixue Lin,
and Ralph Becker-Szendy*
*IBM Almaden Research Center*
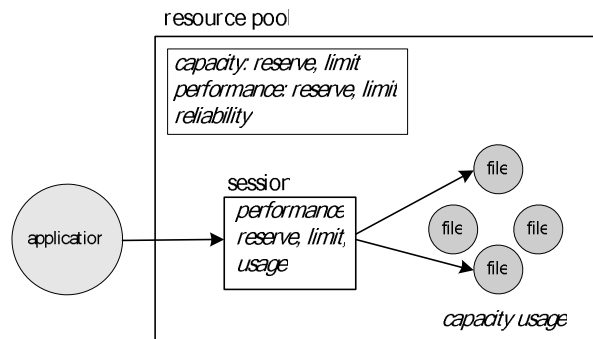*25 August 2005*

1

---

## Kybos: high-level context



configuration

bricks
*(resources)*

resource pools
*(requirements)*

- **Transparency matters**

---

## Resource pool model

resource pool



*capacity: reserve, limit
performance: reserve, limit
reliability*

session

*performance:
reserve, limit,
usage*

application

file

file    file

file

*capacity usage*

- **Goals: enforce reserve and limit; fairly share
unreserved resources**

---

## Virtualizing a resource pool



resource pool

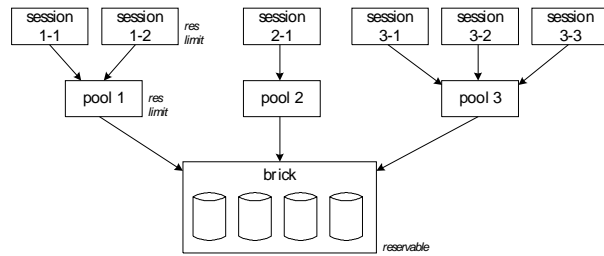*allocation pool*

session

file

bricks

*subsession*

*physical
object*

- **Back virtual entities by physical entities on bricks**
- **Achieve global resource control using local mechanisms**
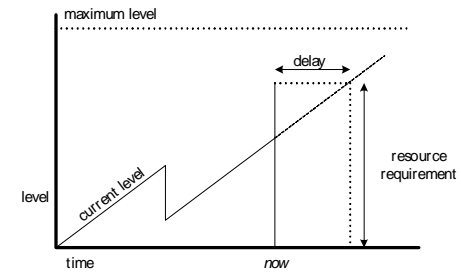
## Local performance control



- Pool/session admission by simple allocation
- Sum of session reserves <= pool
- Limits can be arbitrary; smaller limit applies
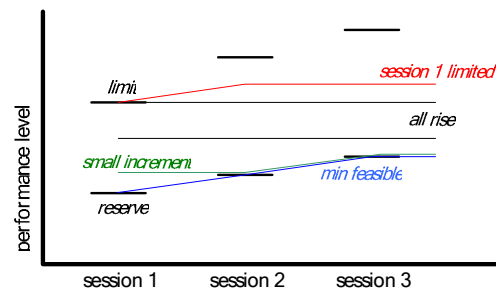- *Prefer to treat local storage as a black box*

## Token buckets



- Usage relative to reference rate
- Reserve: deadline is when level = IO size
- Limit: must wait until level >= IO size
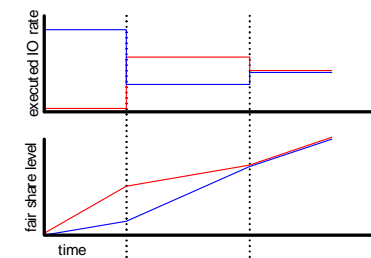
## Fair share: definition



- What defines fair? Over what time period?
- Using "water level" model -- many others possible
- Actually handle fairness first between pools, then sessions within pools
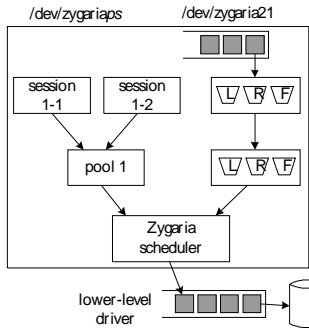
## Fair share: time frame



- Current: long term view, total amount behind (in bytes)
- Alternatives: cap difference in bytes, cap time window
- Pick session (pool) with highest fair share level; can also do a lottery

## Zygaria: algorithm

/dev/zygaria·ps    /dev/zygaria21

session 1-1    session 1-2
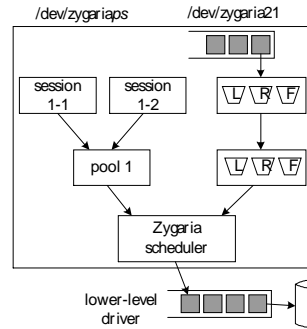
pool 1

Zygaria scheduler

lower-level driver

- Token buckets + EDF + slack stealing
- Compute IO *release times*
  - Later of session and pool limit
- Then compute IO *deadlines*
  - If in past, send now
  - Pool past deadline => send earliest session
- If slack and lower-level queue < $q$ IOs:
  - Pick pool with highest fair share
  - Pick highest session in that pool
- Repeat until nothing more to do
- Trigger on IO arrival, completion, deadline timer

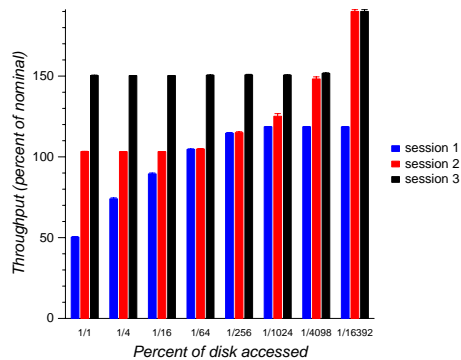25 August 2005          Zygaria          9

## Zygaria: implementation

/dev/zygaria·ps    /dev/zygaria21

session 1-1    session 1-2

pool 1

Zygaria scheduler

lower-level driver

- Loadable driver, Linux 2.6.11
- 2116 lines of C
- $q$=10 outstanding IOs
- Standard Linux IDE driver

- 1.2 GHz Pentium III
- 120 GB, 7200 rpm, IDE

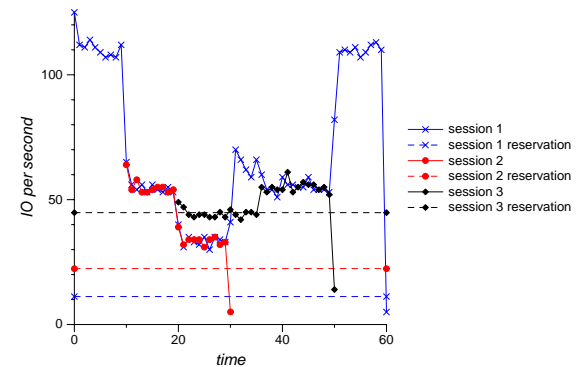25 August 2005          Zygaria          10

## Basic results: sharing



*Throughput (percent of nominal)* vs *Percent of disk accessed*

session 1, session 2, session 3

1/1  1/4  1/16  1/64  1/256  1/1024  1/4098  1/16392

25 August 2005          Zygaria          11

## Basic results: over time



*IO per second* vs *time*

session 1, session 1 reservation, session 2, session 2 reservation, session 3, session 3 reservation

- Disk can do 112 uniformly-distributed random IOps
- All sessions use closed/20 outstanding/0 think time generators, random 1KB requests uniform over whole disk

25 August 2005          Zygaria          12

3

## Overhead

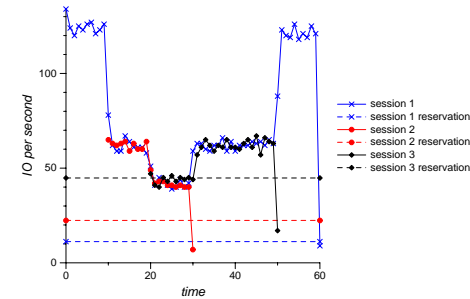| Configuration | Total CPU (%) | Throughput (IOPS) |
|---|---|---|
| No Zygaria | 1.11 | 145 |
| 1 pool, 1 sess, $q$=10 | 1.07 | 112 |
| 10 pool, 10 sess, $q$=10 | 1.24 | 112 |
| 1 pool, 1 sess, $q$=100 | 0.88 | 145 |
| 10 pool, 10 sess, $q$=100 | 0.78 | 145 |

- CPU overhead is in the noise
- Expect proportional to number of active sessions/pools
- Throughput not so good – fewer IOs for head scheduling
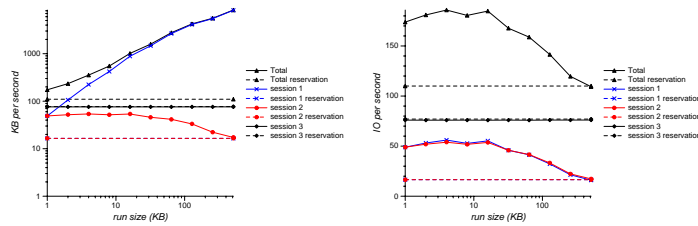
## Adding batching



- Send IOs in batches of up to 10 from sessions
- More aggressive about slack: adapt $q$ based on reservation of active sessions, observed disk performance
- Higher total throughput
- Retain performance for closed/1 workloads

## Combining bandwidth and IO rate



- Use IO *runs* per second; run is up to $r$ KB sequential requests
    - Less bias toward random traffic, handle multiple IO sizes
    - Just group IO requests in a session's queue
- Parameter $r$ manually set; needs to be global for comparability between systems

## Related work

- IO and soft realtime scheduling
    - Media-oriented IO scheduling: Clockwise, Cello
    - Zygaria provides looser scheduling than traditional SRT
- Hierarchical resource allocation
    - UCSC hierarchical disk sharing, Q-RAM, HLS, DQM
- Fair share scheduling
    - Lottery scheduling, YFQ
- Façade, SLEDrunner
    - Adaptive mechanisms, focused on latency
    - Built around an EDF scheduler
    - Zygaria provides stronger guarantees, better control transparency

## Future directions

- Many variations on current algorithm
  - Fair sharing using lottery, other usage estimators
  - Equal-increment fairness, instead of water-level
- Fit Zygaria into a storage system
- Caching
  - Above or below Zygaria?
  - Cache anonymizes traffic
  - Influence load coming into scheduler
- Network flow control
  - Throttling often best done at client
  - Connect Zygaria throttling to protocol
  - Current case: 30K processor BG/L system
- Connect Zygaria to *resource capabilities*

25 August 2005        Zygaria        17

## Conclusions

- Global resource control based on local enforcement
- Zygaria algorithm was simple to implement
- Provides:
  - Reserve and limit enforcement per pool, session
  - Fair sharing
  - Isolation between sessions (applications)
- Good performance requires request batching and aggressive slack use

25 August 2005        Zygaria        18

## Contact

- Richard Golding (rgolding@us.ibm.com)
- Theodore Wong (theowong@us.ibm.com)

25 August 2005        Zygaria        19