

### LSMs, DFSs, and Other Acronyms: Revisiting Storage and Layering for Big Data Management



Michael Carey Information Systems Group CS Department UC Irvine

#AsterixDB

### **Rough Topical Plan**

- Background and motivation (quick!)
- Big Data storage landscape (satellite view <sup>(C)</sup>)
  - Two points of view (plus cloudy skies)
- AsterixDB: our next-generation BDMS
  - What it does (in a nutshell)
  - What we do about storage
- Storage research plans and Q&A





#### **A Few Presenter Cautions**





#### HPTS Has Some Great Debates...

- *Debate #1:* The **TP Architecture Wars** in the early days of HPTS (late 1980's?)
  - TP Heavy: Transaction monitors (middleware)
  - TP Lite: Stored procedures ( $\rightarrow$  one less tier)



#### A DB History Lesson: **DIRECT**



http://isg.ics.uci.edu

#### Debate #2: The Shared What? Wars



http://isg.ics.uci.edu

### Big Data in the **Database** World

- Enterprise data warehouses
  - 1980's: Shared-nothing parallel DBMSs
  - 2000's: Enter new players (Netezza, Aster Data, DATAllegro, Greenplum, Vertica, ParAccel, ...)
- Scalable OLTP

UCIRVINE

http://isg.ics.uci.edu

– 1980's: Tandem's NonStop SQL





#### Notes:

- One storage manager per machine
- Upper layers orchestrate query execution
- One way in/out: through the SQL door

### Later: Big Data in the Systems World

- Out to index and query the Web, Google laid a new foundation in the early 2000's
  - Google File System (GFS): Files spanning many machines with 3-way replication
  - MapReduce (MR): "Parallel programming for dummies" (UDFs + parallel framework)
  - Yahoo!, FB, et al read the papers
    - HDFS and Hadoop MapReduce
    - Declarative HLLs: Pig, Hive, ...
    - HLLs now heavily preferred to MR
  - Also key-value stores ("NoSQL")
    - Social sites, online games, ...
    - BigTable/HBase,
       Dynamo/Cassandra,
       MongoDB, ...





#### Remember History...? (DIRECT)





#### **One More Bit of History**

#### **Operating System Support for Database Management**

Michael Stonebraker University of California, Berkeley

> Communications of the ACM

July 1981 Volume 24 Number 7

#### 7. Conclusions

The bottom line is that operating system services in many existing systems are either too slow or inappropriate. Current DBMSs usually provide their own and make little or no use of those offered by the operating system. It is important that future operating system designers become more sensitive to DBMS needs.

ting system services are examined oplicability to support of database lese services include buffer pool om; scheduling, process managecommunication; and consistency



#### Plus: Today's Big Data Tangle



### AsterixDB: "One Size Fits a Bunch"

# Semistructured Data Management Asteri World of Parallel Hadoop & Friends Database Systems

#### **BDMS** Desiderata:

- Flexible data model
- Efficient runtime
- Full query capability
- Cost proportional to task at hand (!)
- Designed for continuous data ingestion
- Support today's "Big Data data types"



(Note: This work began in 2009.)

### ASTERIX Data Model (ADM)

}

create dataverse TinySocial; use dataverse TinySocial;

- create type MugshotUserType as {
  - id: int32,
  - alias: string,
  - name: string,
  - user-since: datetime,
  - address: {
    - street: string,
    - city: string,
    - state: string,
    - zip: string,
    - country: string
  - },

friend-ids: {{ int32 }},

employment: [EmploymentType]

create type EmploymentType as open {
 organization-name: string,
 start-date: date,
 end-date: date?

create dataset MugshotUsers(MugshotUserType) primary key id;

*Note:* We store and manage datasets...

#### Highlights include:

- JSON++ based data model
- Rich type support (spatial, temporal, ...)
- Records, lists, bags
- Open vs. closed types

#### ASTERIX Data Model (ADM)



#### Highlights include:

- JSON++ based data model
- Rich type support (spatial, temporal, ...)
- Records, lists, bags
- Open vs. closed types



### **Other Data Management Features**

create index msUserSinceIdx on MugshotUsers(user-since); create index msTimestampIdx on MugshotMessages(timestamp); create index msAuthorIdx on MugshotMessages(author-id) type btree; create index msSenderLocIndex on MugshotMessages(sender-location) type rtree; create index msMessageIdx on MugshotMessages(message) type keyword;

#### create type AccessLogType as closed

(("path"="{hostname}://{path}"), ("format"="delimited-text"), ("delimiter"="|"));

create feed socket\_feed using socket\_adaptor
 (("sockets"="{address}:{port}"), ("addressType"="IP"),
 ("type-name"="MugshotMessageType"), ("format"="adm"));
connect feed socket\_feed to dataset MugshotMessages;



## ASTERIX Query Language (AQL)

• *Ex:* Identify active users and group/count them by country:

with \$end := current-datetime()
with \$start := \$end - duration("P30D")
from \$user in dataset MugshotUsers
where some \$logrecord in dataset AccessLog
satisfies \$user.alias = \$logrecord.user
and datetime(\$logrecord.time) >= \$start
and datetime(\$logrecord.time) <= \$end
group by \$country := \$user.address.country keeping \$user
select {</pre>

"country" : \$country,
"active users" : count(\$user)

# http://isg.ics.uci.edu

#### AQL highlights:

- Lots of other features (see website!)
- Spatial predicates and aggregation
- Set-similarity matching
- And plans for more...

#### AsterixDB System Overview



UCIRVINE http://isg.ics.uci.edu

#### Local LSM-Based Storage & Indexes



### **Distributed Storage in AsterixDB**

- Hash-partitioned, shared-nothing, local drives
  - Partitioning based on primary key (hashing)
  - Secondary indexes local to, and consistent with, corresponding primary partitions (all LSM-based)
- Also offer external dataset feature (for HDFS)
  - Multiple (Hive) formats, secondary index support
  - Index partitions co-located with data (if possible)
  - Developed for space and "IT comfort" reasons



### Data Replication in AsterixDB (WIP)



#### Chained Declustering



(synchronous, recovery-only copies kept)





### Hedging Our Bets

- We're currently porting our LSM-based storage system to also work on *top* of HDFS (and YARN)
  - Might somehow feel more "comforting" (and/or "environmentally friendly") to Big Data IT shops
  - Another path to replication and high availability
- Interesting experiments lie ahead!
  - Revisit Stonebraker-like OS issues (modern version)
  - Bake-off: Distributed record management vs. DFS
  - Just how well does HDFS do *w.r.t.* locality of writes?



#### What About the Cloud?

- Computing may be elastic, but data is not...!
  - Native storage  $\rightarrow$  hard to expand & contract
  - Seems to demand a shared-disk-like approach based on cloud storage facilities?
- Experimentation is needed
  - E.g., AWS storage or Google persistent disks
  - Performance implications seem pretty interesting to explore...





#### For More AsterixDB Info...

NSF project page (UC Irvine and UC Riverside):

• <u>http://asterixdb.ics.uci.edu</u>

Apache AsterixDB (Incubating) project page:

• <a href="https://asterixdb.incubator.apache.org/">https://asterixdb.incubator.apache.org/</a>



#### I Asked the Questions, So....

• Got any answers?





#AsterixDB