



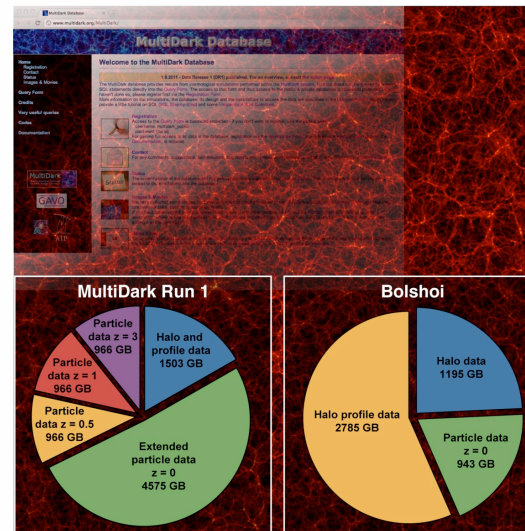
# PaQu - Getting the most out of MySQL with distributed queries

Adrian M. Partl  
Leibniz-Institut für Astrophysik Potsdam (AIP)



[www.multidark.org](http://www.multidark.org)

- Total row count:  
 $2.34 \cdot 10^{11}$
- Most queries 100 sec
- significant amount of queries > 1000 sec
  
- Full table scan for largest tables:  
~ 30 minutes





## Problems while building the MultiDark DB

- Data ingest time:  
Need to convert binary to ASCII CSV format  
(highly inefficient)
- Data transformation:  
Computing values after ingest slow - best during ingest
- Data indexing:  
Index on particle data ( $\sim 10^{10}$  particles) around one week
- Data retrieval times slow on full table scans ( $\sim 30$  min):  
cannot build index for every query
- Spatial queries in 3D hard, impossible in 6D  
nearest neighbour search also inefficient



AIP

## Why RDBMS?

- SQL - it took long time for the community to adopt SQL (we think this is the main problem with NoSQL)
- proven, widely available, large user base
- good for structured data
  
- Problems:
  - Built for different purposes (business, web, ...)  
result sets usually small - mostly in memory solutions
  - parallelisation of data / sharding
  - can be expensive



## Our vision:

- Open source DB solution for scientific purposes:  
A one size fits all solution built by the community for the community
- Developments at AIP:
  - DB independent ingestion library and data transformation tool (*DBIngestor and AsciiIngest*)
  - MySQL job queue (*mysql\_query\_queue*)
  - MySQL sharding solution for scientific queries (*PaQu*)
  - Future:
    - MySQL plugins for data analysis, spatial queries and indexing
    - MySQL storage engine plugins for simulation raw data



## PaQu - HPQ for MySQL

- Full table scan of large table: ~ 30 minutes
- From HPC we know:  
*Distribute your problem over N nodes and you are N times faster (if you are lucky)*
  - i.e. on 10 nodes, full table scan of large table:  
3 minutes!
- From Big Data hype we know:  
*Use Hadoop and everything is fine*
- We want SQL! We want “Hadoop” for MySQL!



AIP

# PaQu - HPQ for MySQL

## Spider engine:

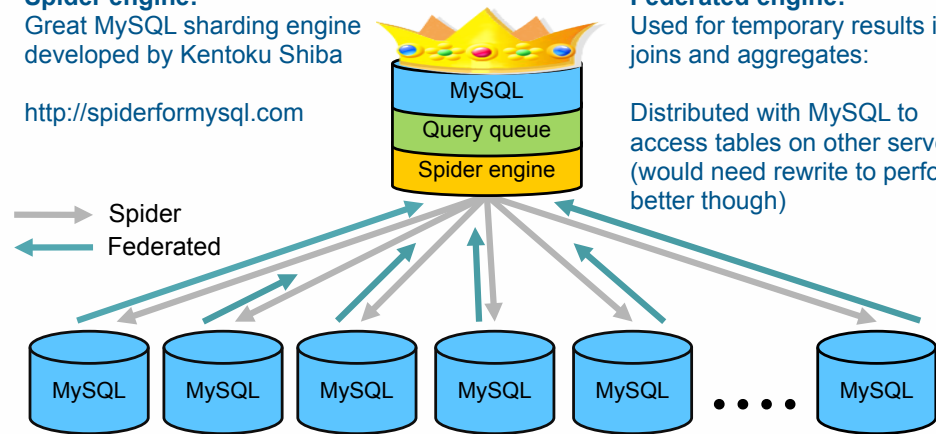
Great MySQL sharding engine developed by Kentoku Shiba

<http://spiderformysql.com>

## Federated engine:

Used for temporary results in joins and aggregates:

Distributed with MySQL to access tables on other servers. (would need rewrite to perform better though)



HPQ = high performance querying



AIP

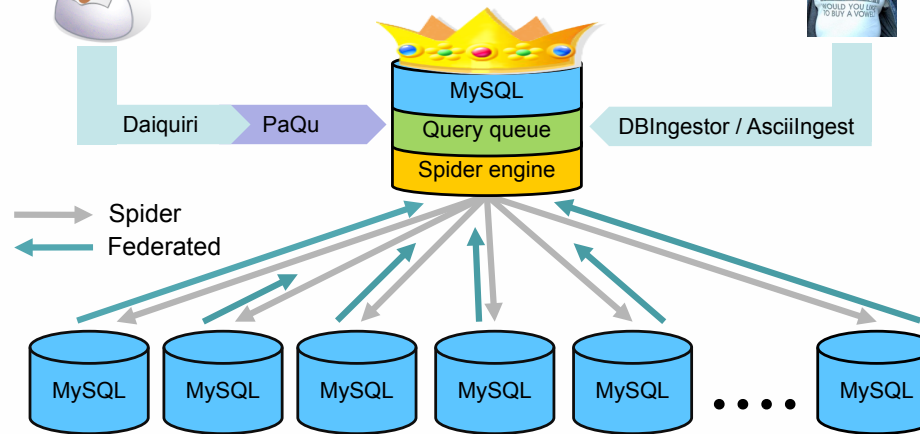
# PaQu - Part of the Daiquiri world



User



Admin

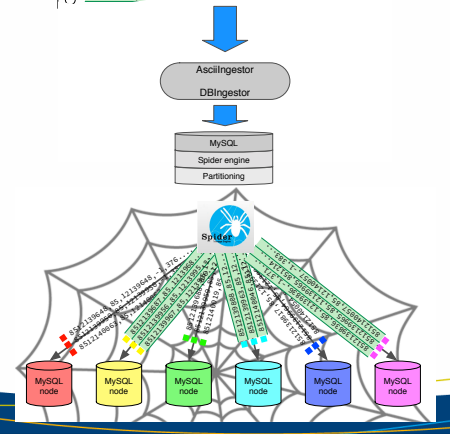






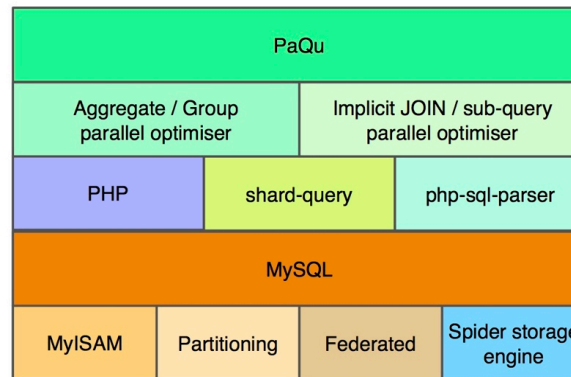
# Sharding

○	bdnid, snapshot, NFNcat, hostFlag, s, y, z, vs, vy, vz, mp, RvYr, #To...	○
●	8512139648_85_12139648_-1_376_6241_422_5894_822_0686_-28...	○
●	8512139687_85_12139687_-1_215_6037_495_9996_824_5079_-27...	○
○	8512139648_85_12139648_8512140063_386_5749_459_3805_824...	○
○	8512139888_85_12139888_-1_386_2757_459_364_823_0633_-200...	○
○	8512139817_85_12139817_-1_379_7038_461_651_824_95_-310_5...	○
○	8512139836_85_12139836_-1_371_6952_468_1314_821_6387_-26...	○
○	8512139958_85_12139958_-1_383_4951_458_8577_826_043_-213...	○
○	8512139956_85_12139956_-1_380_4813_460_9308_825_48_-230...	○
○	8512139959_85_12139959_-1_389_6255_462_5074_821_4885_-28...	○
○	8512139963_85_12139963_-1_369_4372_467_9669_820_5311_-38...	○
○	8512139964_85_12139964_-1_376_2421_467_8033_829_4519_-38...	○
○	8512139965_85_12139965_8512140359_874_2386_468_974_822_7...	○
○	8512139965_85_12139965_-1_389_1884_463_884_821_517_-245...	○
○	8512139967_85_12139967_-1_391_9354_465_8879_821_3747_-26...	○
○	8512140019_85_12140019_-1_362_9375_424_1447_821_813_-265...	○
○	8512140066_85_12140066_-1_377_2039_447_7321_820_3369_-43...	○
○	8512140054_85_12140054_-1_372_1338_454_1272_821_2715_-32...	○
○	8512140057_85_12140057_-1_383_8878_460_9308_825_48_-230...	○
○	8512140053_85_12140053_-1_377_2039_447_7321_820_3369_-43...	○
○	8512140064_85_12140064_-1_377_2039_447_7321_820_3369_-43...	○





# PaQu architecture



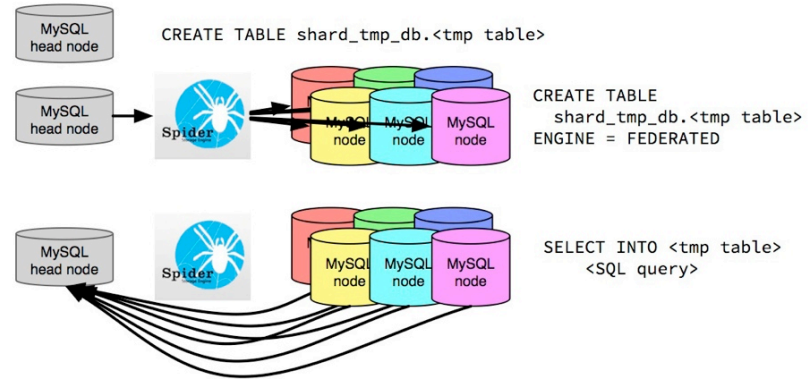
PaQu consists of two components:

- Basic distributed querying operations
- Automatic query paralleliser



# Basic distributed operations

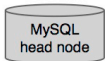
**CALL paquExec( <SQL query>, <tmp table> )**





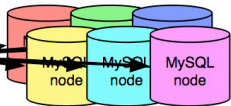
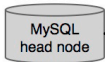
# Basic distributed operations

CALL paquLinkTmp( <tmp table> )



col1	col2
27389	-2,3
57657	3,23
324576	1

Table residing on head node

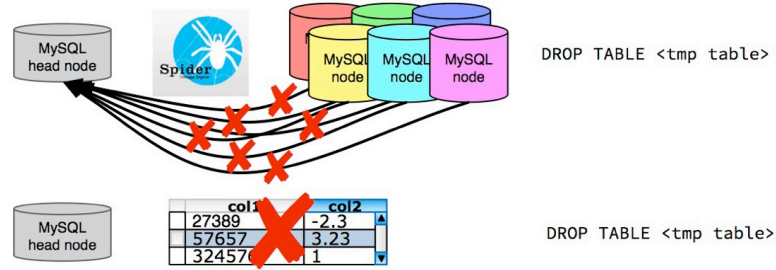


CREATE TABLE  
shard\_tmp\_db.<tmp table>  
ENGINE = FEDERATED



# Basic distributed operations

CALL paquDropTmp( <tmp table> )





AIP

## Query paralleliser - example

```
SELECT  snapnum,  
        AVG(mass)  
FROM    MDR1.FOF  
GROUP BY snapnum;
```



```
SELECT  a.snapnum, SUM(a.mass)/SUM(a.mass)  
FROM  
  (SELECT  snapnum, SUM(mass) as sum,  
           COUNT(mass) as cnt  
   FROM    MDR1.FOF  
   GROUP BY snapnum) as a  
GROUP BY a.snapnum;
```



```
CALL paquExec(  
  "SELECT snapnum AS `snapnum`, COUNT(mass) AS `cnt_avg(mass)`,  
    SUM(mass) AS `sum_avg(mass)`  
  FROM MDR1.FOF GROUP BY snapnum", "aggregation_tmp_8896713");  
  
USE spider_tmp_shard; SET @i=0;  
CREATE TABLE multidark_user_admin.`/*@GEN_RES_TABLE_HERE*/` ENGINE=MyISAM  
  SELECT DISTINCT @i:=@i+1 AS `row_id`, `snapnum`,  
    (SUM(`sum_avg(mass)`) / SUM(`cnt_avg(mass)`)) AS `avg(mass)`  
  FROM `aggregation_tmp_8896713` GROUP BY `snapnum` ;  
  
CALL paquDropTmp("aggregation_tmp_8896713");
```



# Automatic query paralleliser

## Implicit Joins:

```
SELECT  a.*, b.*, c.*
FROM    a, b, c
WHERE   b=2 AND
        b.id=c.b_id AND
        a.id=b.a_id;
```



```
SELECT  a.*, tmp.*
FROM    a,
        (SELECT  b.*, c.*
         FROM    c,
                 (SELECT  b.*
                  FROM    b
                   WHERE  b=2) as b
         WHERE   b.id=c.b_id)
WHERE   a.id=tmp.b.a_ids
```

## Aggregates:

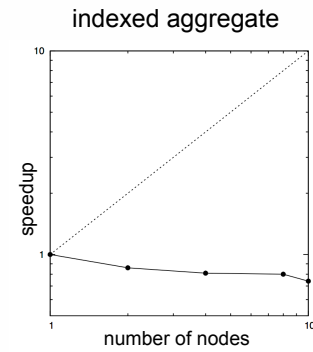
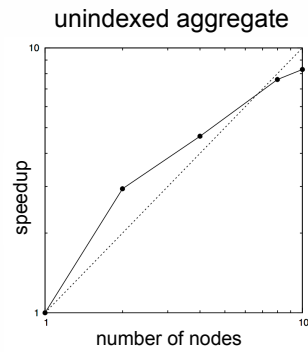
```
SELECT  a.bar, AVG(a.foo)
FROM    a
GROUP BY a.bar;
```



```
SELECT  a.bar,
        SUM(a.sum)/SUM(a.cnt)
FROM    (SELECT  a.bar as bar,
                 SUM(a.foo) as sum,
                 COUNT(a.foo) as cnt
         FROM    a
         GROUP BY a.bar) as a
GROUP BY a.bar;
```



## Performance results



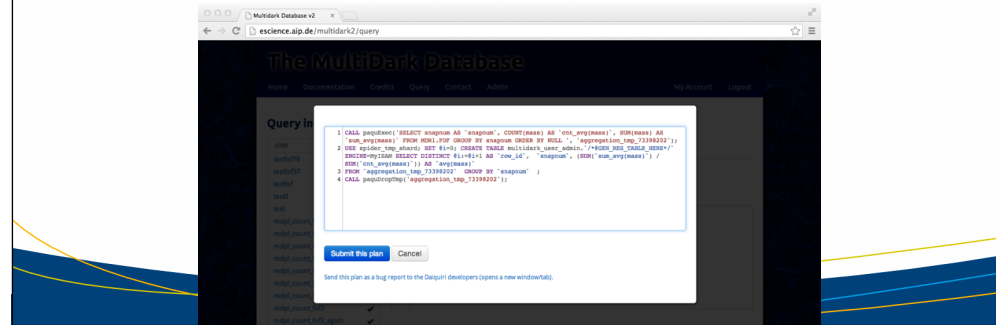
- Strong correlation with hardware setup:
  - Cache sizes, size of data files (smaller is better / partitioning?), network and I/O performance





# MultiDark v2.0

- PaQu covers most known use cases of MultiDark
- All data from MultiDark now on 10 DB nodes
- PaQu fully integrated with Daiquiri
- <http://escience.aip.de/multidark2>





## Conclusions

PaQu brings parallel querying to MySQL.

PaQu relies on open source.

PaQu is open source!

Download our DB developments here:

<https://github.com/adrpar>

Try PaQu yourself and compare with the original:

<http://escience.aip.de/multidark2>