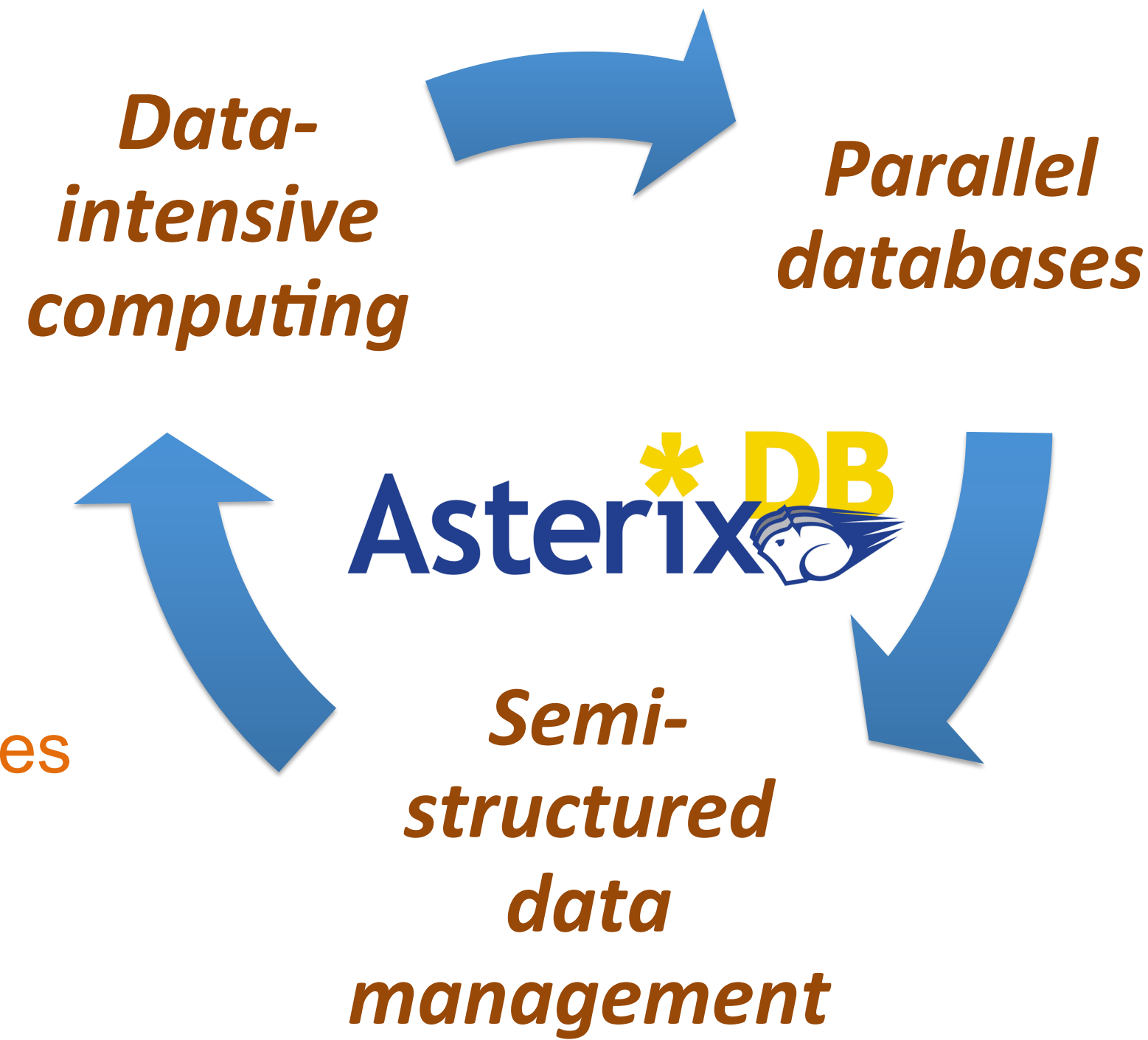


### At a Glance

- Open Source (Apache Incubating!) **Big Data Management System**
- Native support for **data ingestion** (data feeds)
- Runs on large **commodity clusters**
- Designed for mass quantities of **semi-structured data**
- Highly scalable **storage and index** management
- Native support for **rich data types and operations** (e.g., spatial & temporal data)
- Native support for **similarity queries**



### AQL (AsterixDB Query Language) + ADM

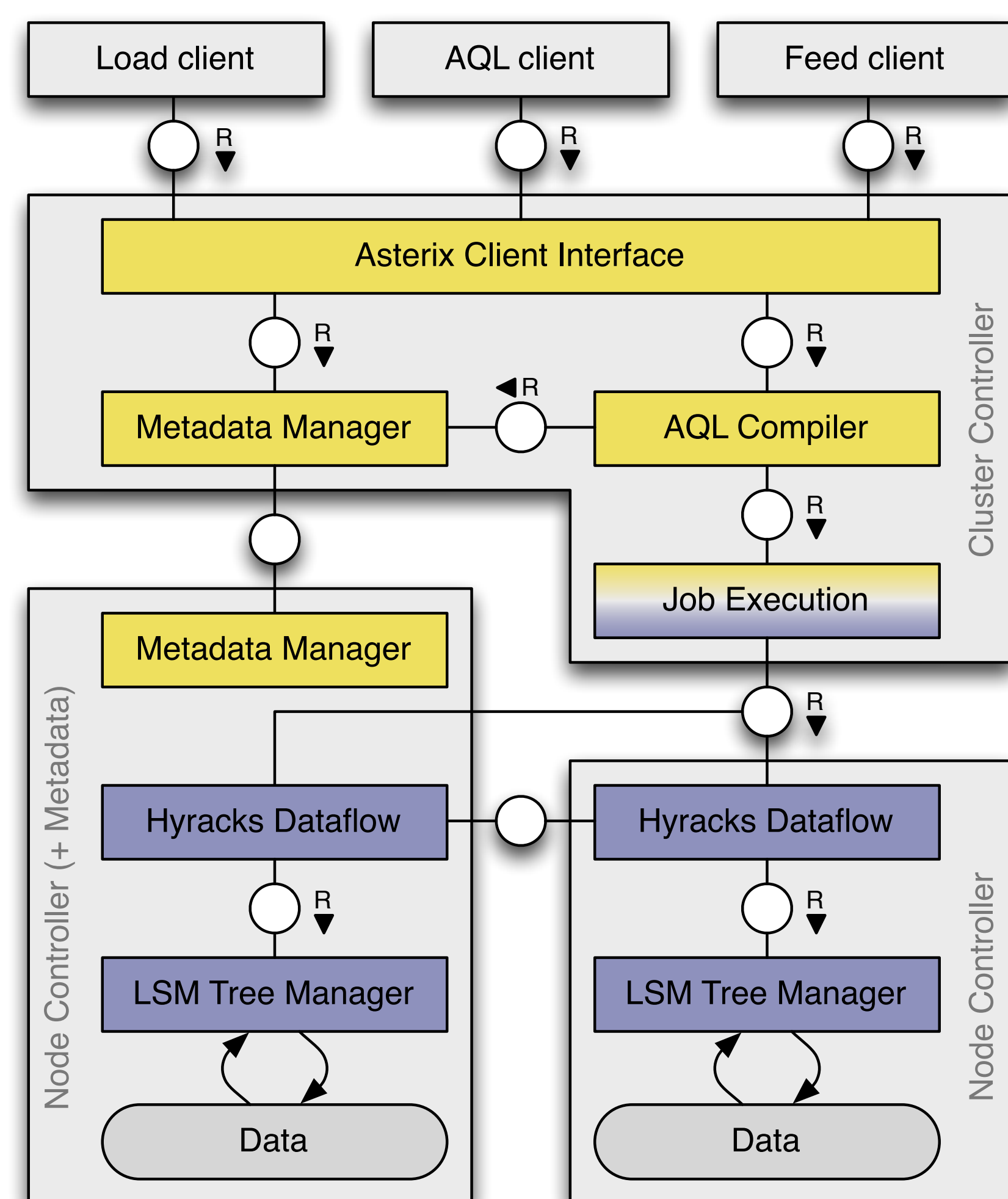
- ADM is a **superset of JSON** with a richer set of types (e.g., bags, spatial data, temporal data, text) and **optional schemas**
- AQL is a **powerful query language** for semi-structured data, influenced by the best parts of W3C's XQuery

• *Ex: List the user name and messages sent by those users who joined the Mugshot social network in a certain time frame:*

```
for $user in dataset MugshotUsers
where $user.user-since >= datetime('2010-07-22T00:00:00')
and $user.user-since <= datetime('2012-07-29T23:59:59')
return {
  "uname" : $user.name,
  "messages" :
    for $message in dataset MugshotMessages
    where $message.author-id = $user.id
    return $message.message
};
```

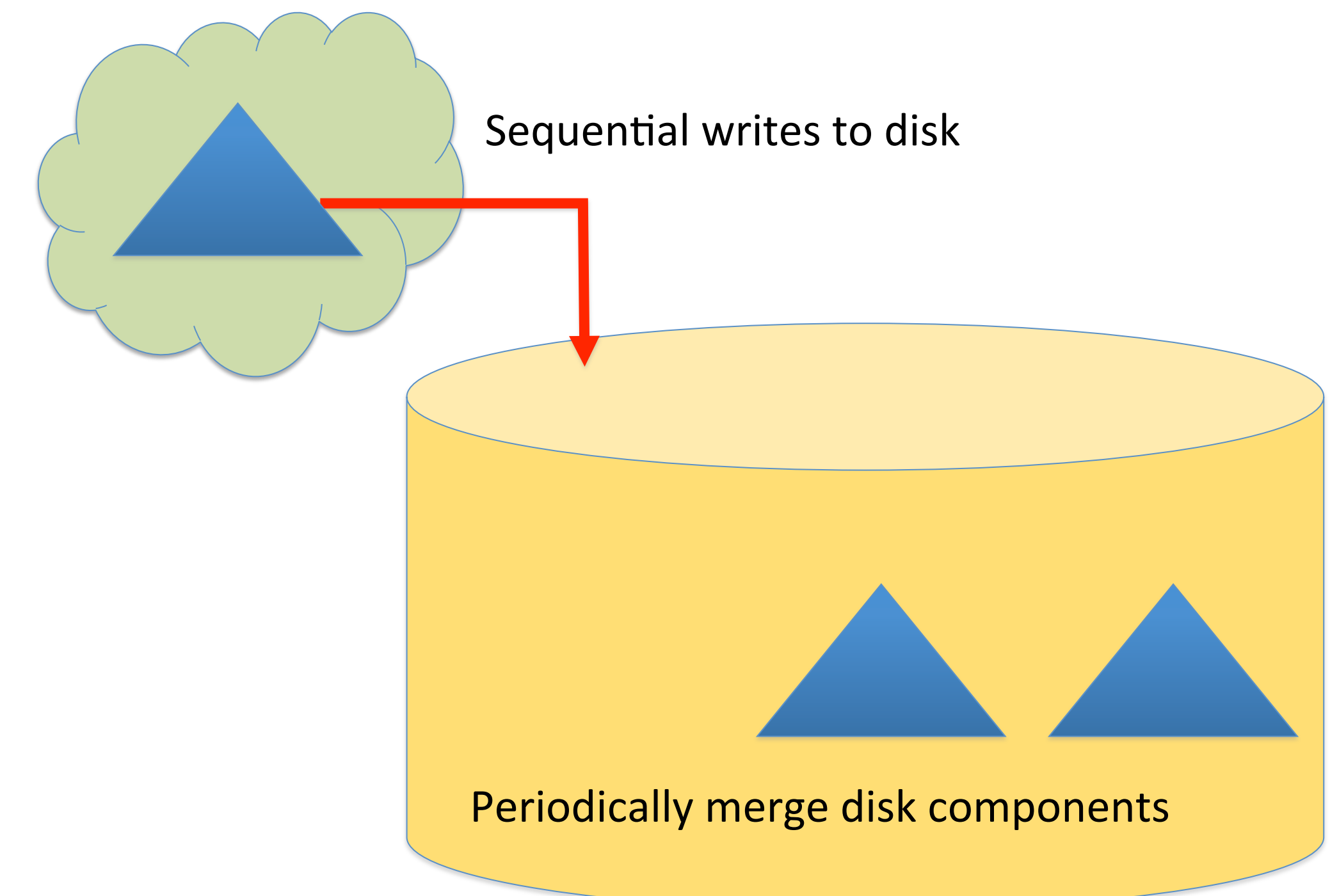
### System Architecture

- Uses the **Hyracks** data parallel platform as its runtime engine
- Shared-nothing** storage
- Built for **commodity clusters**
- AQL uses **Algebricks** to optimize queries
- Each Node Controller stores a **partitioned portion** of the data stored in each index



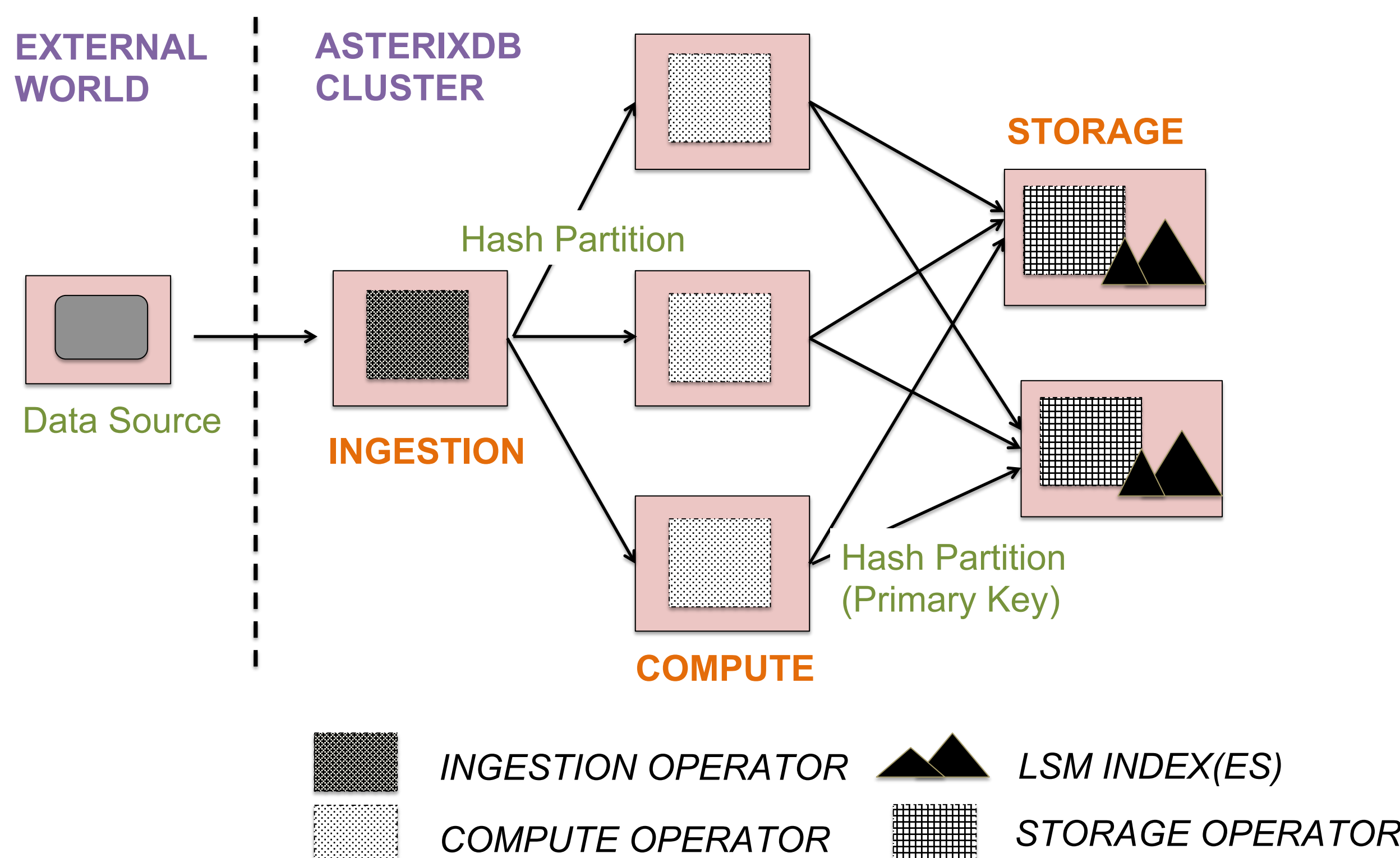
### LSM-Based Storage and Indexing

- Log-Structured Merge** technology to support high data ingestion rates
- All primary and secondary index types are **LSM-ified** (B+Tree, R-Tree, Inverted Indexes, ...)



### Data Feeds

- Integrated system for **intake of data** from live sources (e.g., Twitter)
- Allows for **user-defined computation** to be performed on ingestion



### Competitive Performance

	Users	Messages	Tweets
Asterix (Schema)	192	120	330
Asterix (KeyOnly)	360	240	600
Syst-X	290	100	495
Hive	38	12	25
Mongo	240	215	478

Table 2: Dataset sizes (in GB)

Batch Size	Asterix Schema	Asterix KeyOnly	Syst-X	Mongo
1	0.091	0.093	0.040	0.035
20	0.010	0.011	0.026	0.024

Table 4: Average insert time per record (in sec)

	Asterix Schema	Asterix KeyOnly	Syst-X	Hive	Mongo
Rec Lookup	0.03	0.03	0.12	(379.11)	0.02
Range Scan	79.47	148.15	148.33	11717.18	175.84
— with IX	0.10	0.10	4.90	(11717.18)	0.05
Sel-Join (Sm)	78.03	96.76	55.01	333.56	66.46
— with IX	0.51	0.55	2.13	(333.56)	0.62
Sel2-Join (Lg)	79.62	99.73	56.65	350.92	273.52
— with IX	2.24	2.32	10.59	(350.92)	14.97
Sel2-Join (Sm)	79.06	97.82	55.81	340.02	66.45
— with IX	0.50	0.52	2.62	(340.02)	0.61
Sel2-Join (Lg)	80.18	101.24	56.10	394.11	313.17
— with IX	2.32	2.32	10.70	(394.11)	15.28
Agg (Sm)	128.66	232.30	130.64	83.18	400.97
— with IX	0.16	0.17	0.14	(83.18)	0.19
Agg (Lg)	128.71	232.41	132.19	94.11	401
— with IX	5.53	5.55	4.67	(94.11)	8.34
Grp-Aggr (Sm)	130.20	232.77	131.18	127.85	398.27
— with IX	0.45	0.46	0.17	(127.85)	0.20
Grp-Aggr (Lg)	130.62	234.10	133.02	140.21	400.10
— with IX	5.96	5.91	4.72	(140.21)	9.03

Table 3: Average query response time (in sec)