

How to handle a DDOS which makes you happy (as it comes from your customers)

Michal Barla
FIIT STU
luigisbox.com

luigisbox.com

- Site Search Analytics
- Which queries ends with no results?
- Which queries have low CTR?
- Which queries have wrong results ordering?
- What are people looking for on your site?
- What typos people do when searching for your content?
- A/B testing

luigisbox.com - solution

- Web site administrator inserts our javascript
- Provides few HTML5 annotations
- Collected information ends on our servers
 - Dashboard, diagnostics

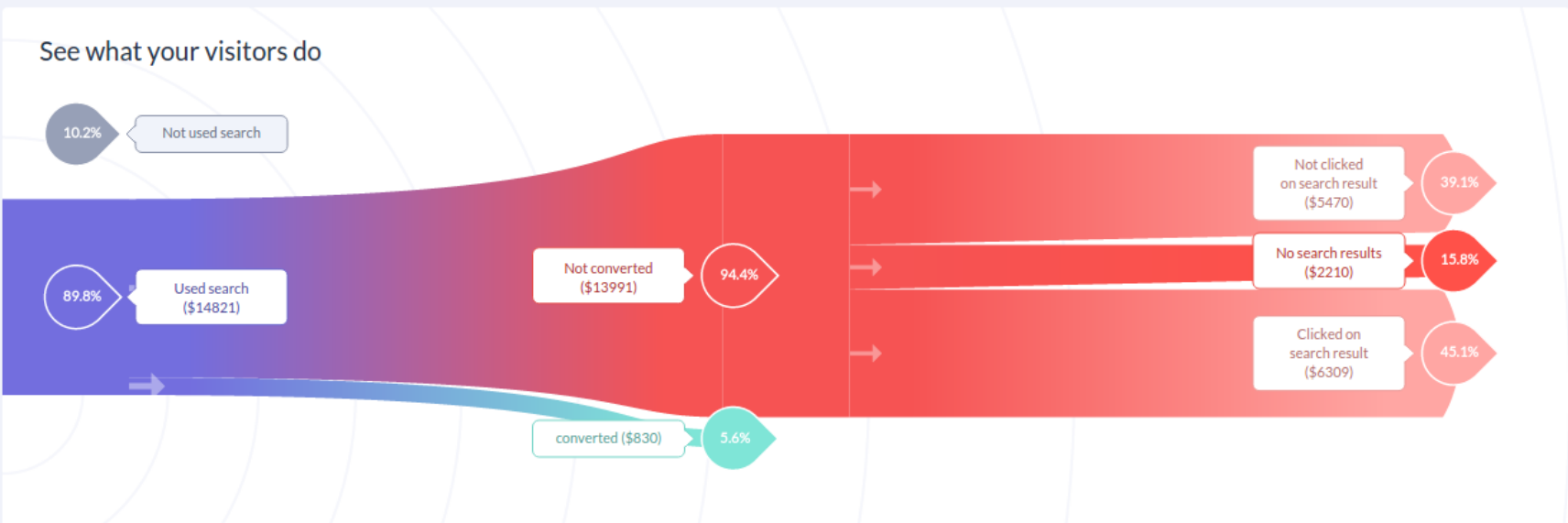
luigisbox.com

MY BOOKSTORE
mybookstore.com

- Search Analytics
- Profiles
- Account
- Help center
- Settings

Dashboard Show data for 7 days

USER PROFILES COLLECTED 17,341	LOST OPPORTUNITY 13,991	SEARCHES PERFORMED 38,982	CONVERSIONS 2,183
--	-----------------------------------	-------------------------------------	-----------------------------



SEARCHES WITHOUT RESULTS

in last 7 days

1	four	64	\$832	Fix it
2	after	52	\$676	Fix it
3	delivery options	46	\$598	

TRENDING SEARCHES

in last 7 days

1	jo nesbo	640	\$6235.0	Fix it
2	harry potter	590	\$7360	Fix it
3	game of thrones	592	\$2201	

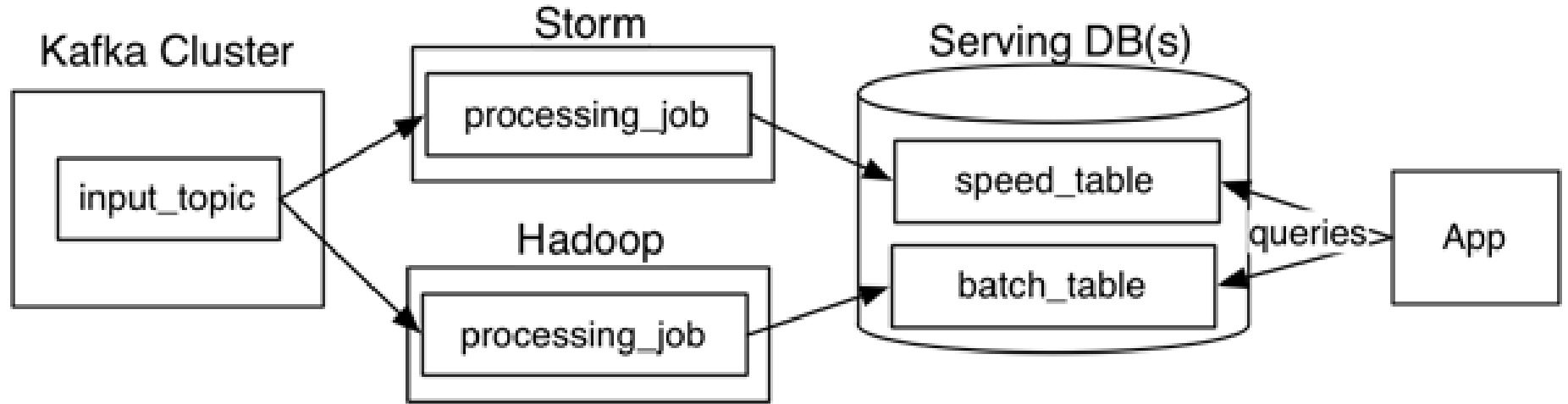
ALL TIPS

- Warning:** You are losing visitors on "four" query due to no search results.
- You are getting very high traffic from

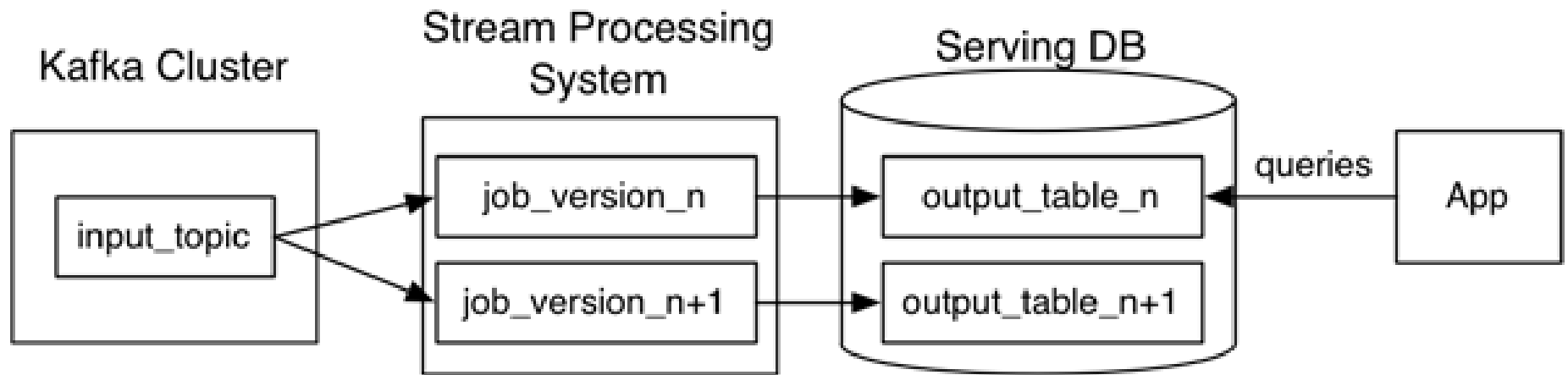
luigisbox.com

- You must handle traffic from customers of your customer
 - Which can be a huge e-shop
- Your failure must be isolated from an end-user
- You do not want to loose data during an outage
- Quick asynchronous processing

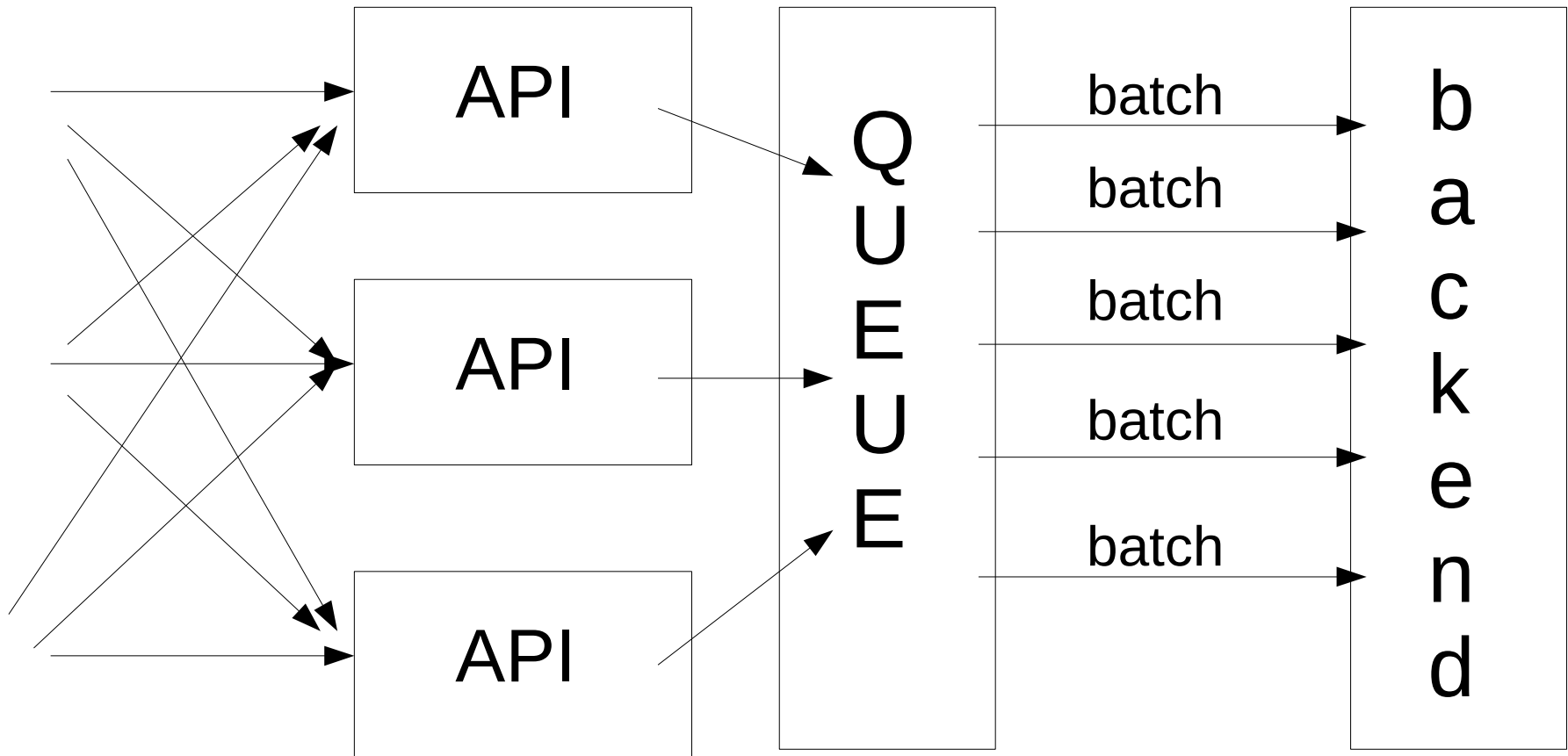
Lambda architecture



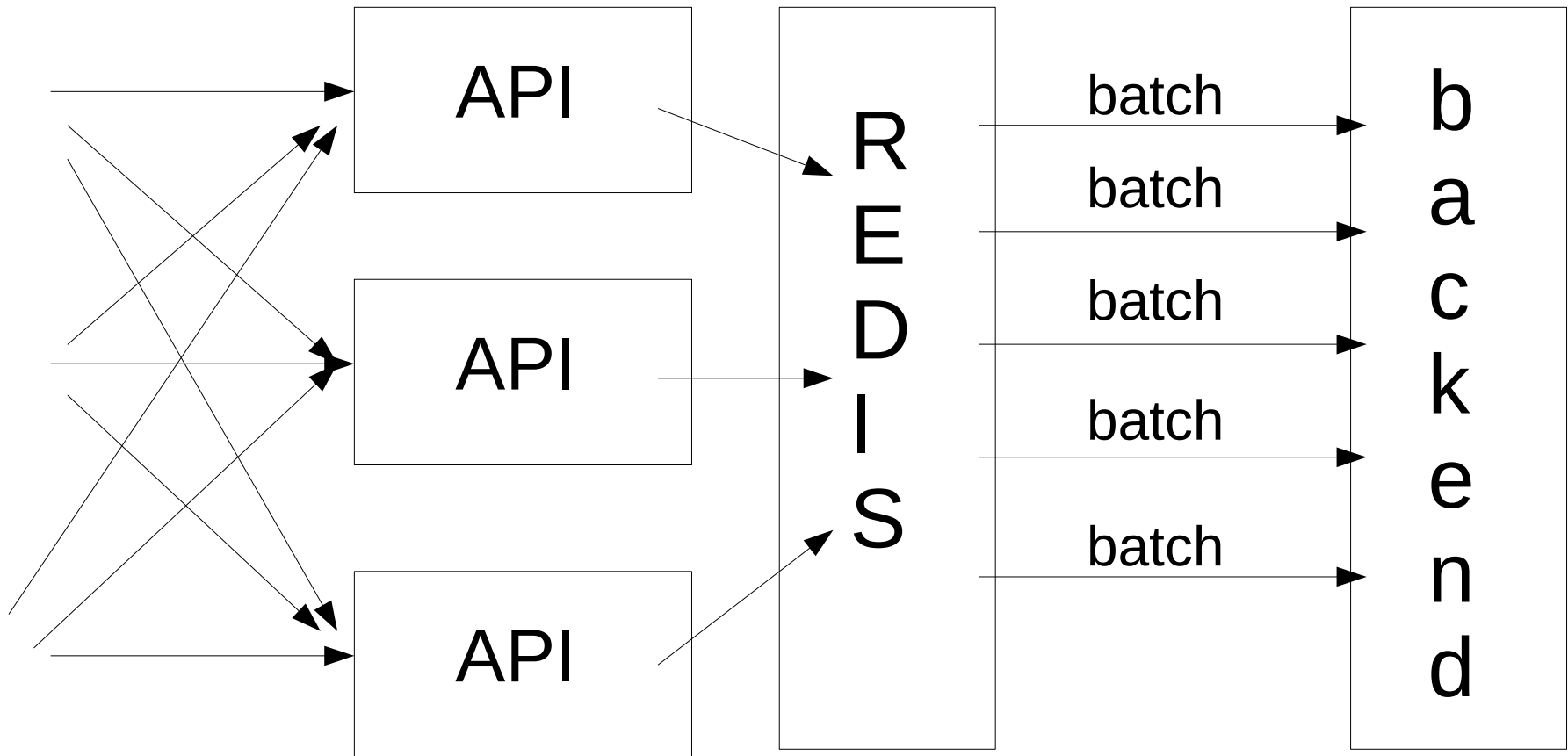
Kappa architecture



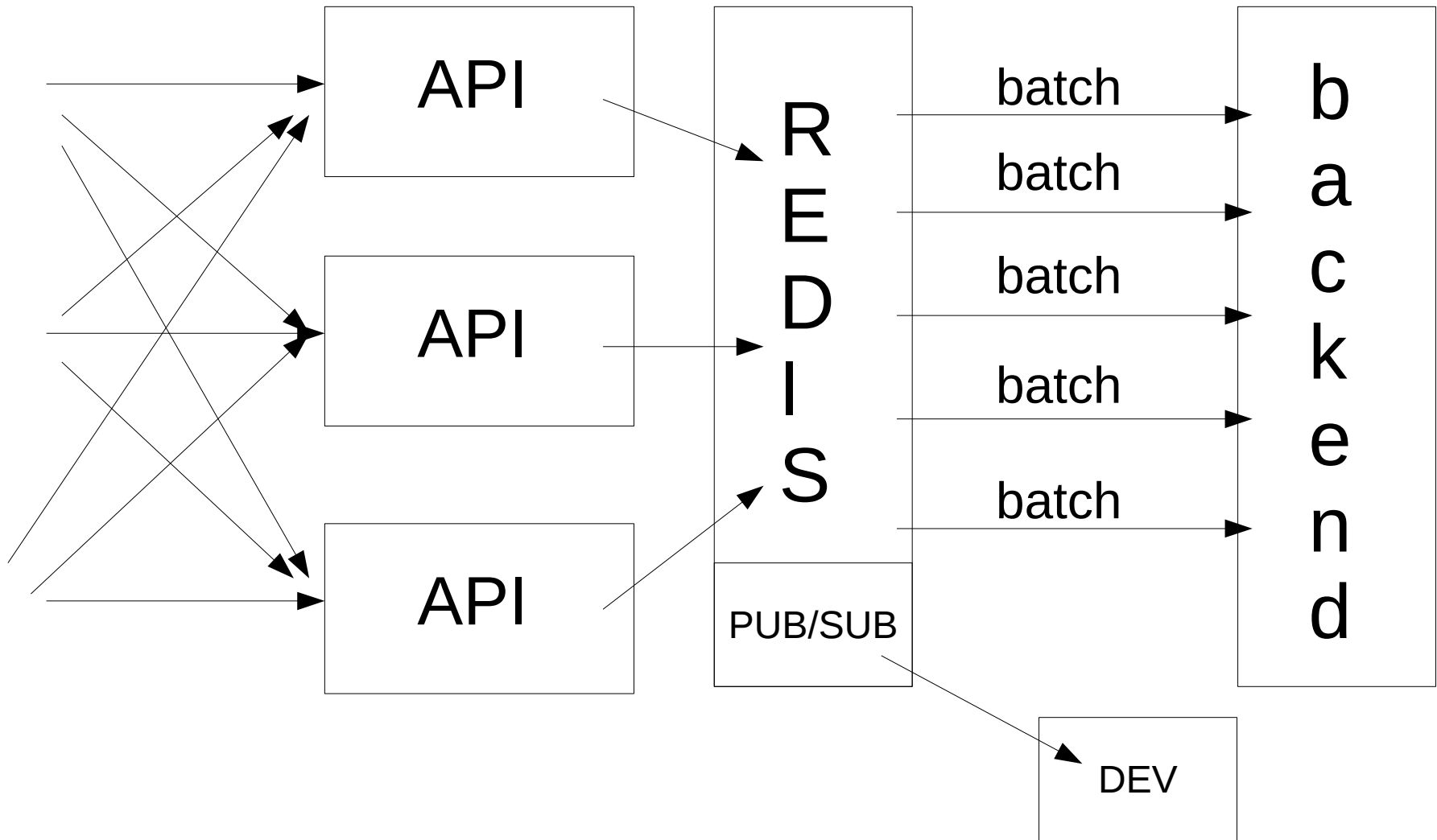
Architecture



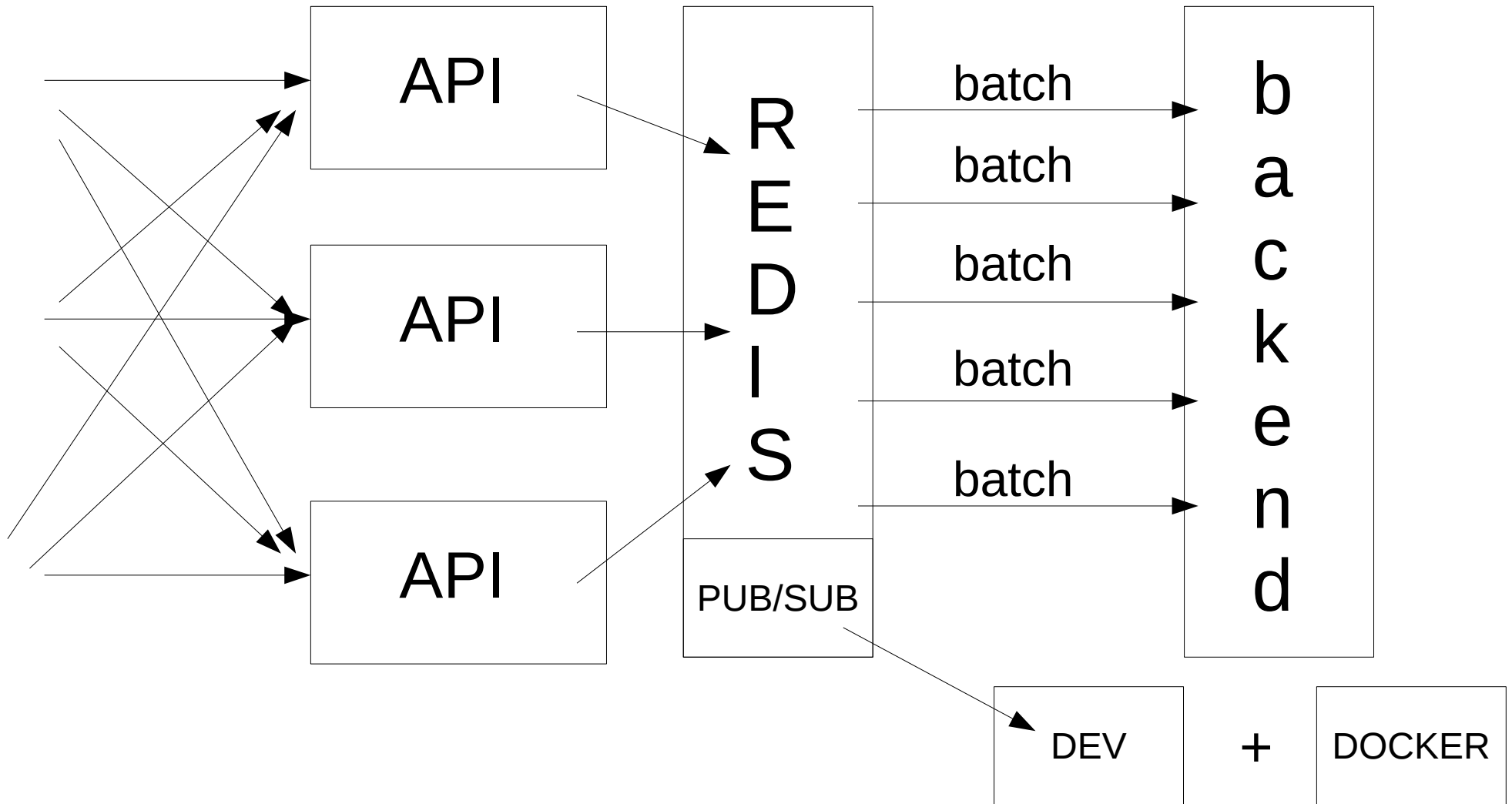
Architecture



Architecture



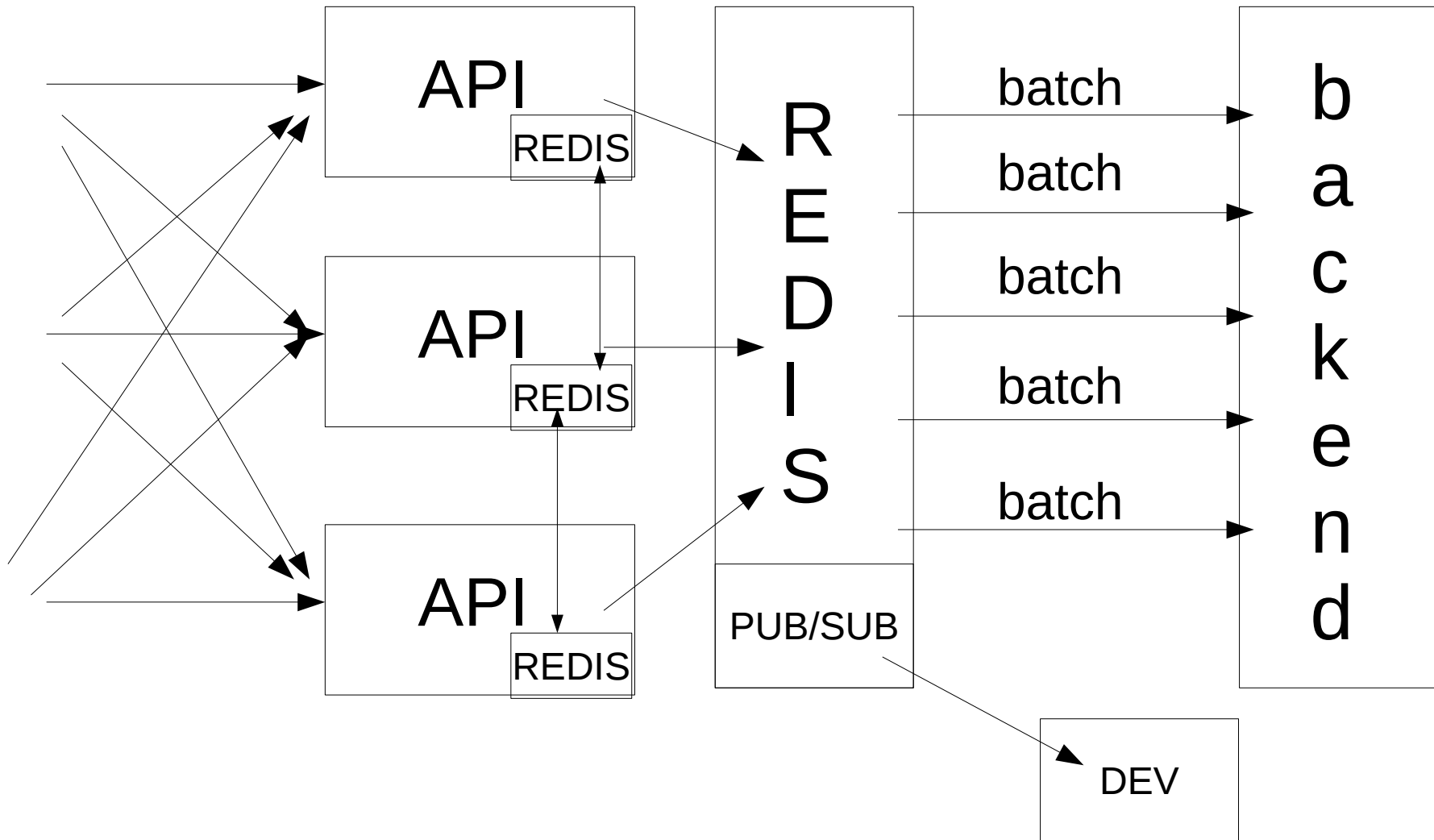
Architecture



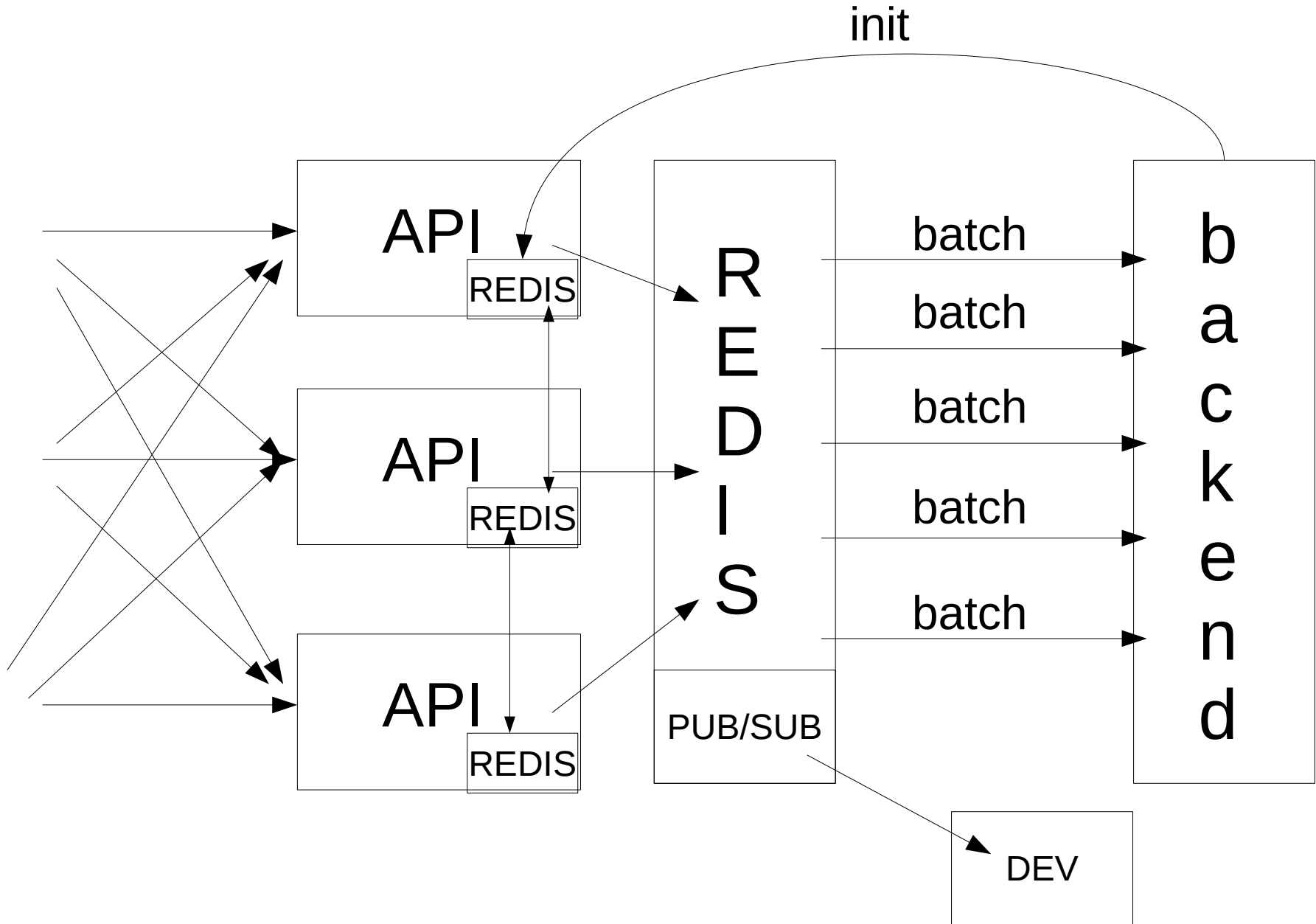
Additional requirements

- Collect every pageview, every interaction
- Quick fix of most serious search-related problems
 - Wrong results ordering
 - Add a result
 - Remove a result
- Site adaptation for a specific customer segment

Architecture



Architecture



luigisbox.com backend

- Many COUNT DISTINCT users/sessions/events + GROUP BY query/content + ORDER BY
- Custom criteria for events and their cardinality
 - GROUP BY HAVING
- A custom time interval for metrics
- Elasticsearch and its Aggregations framework
 - Parent/child a nested documents
- Approximations of distinct counts
 - LinearCounting
 - HyperLogLog

HyperLogLog

- Coin tossing
 - Longest run of consecutive heads
- If the longest run of leading zeroes is N , than cardinality estimation is 2^N
- Hash
 - First k bits identifies a register
 - A register stores longest run of leading zeroes from the rest of the hash
- Estimate is a harmonic mean of all registers
- **They can be combined together!!**
- **postgres-hll**

Druid.io

- OLAP on top of events
- Roll-up on data ingestion
 - Granularity of a milisecond
 - But I need to know my aggregations beforehand

Hadoop and HDFS

- HDFS is the right durable deep storage (?)
 - Until your NameNode runs out of disk space!
- It is complicated
 - YARN, ZooKeeper, Oozie, Nagios, Ganglia, ...
- It is slow...?
 - SQL on Hadoop starts to be very interesting

Apache Drill

- Open source implementation of Google's Dremel paper
 - BigQuery
- ANSI SQL on top of
 - CSV
 - JSON
 - **Parquet !**
 - Hive
 - ...all this **in HDFS**
 - JOIN among different formats
- Massively parallel
- JDBC/ODBC

Greenplum

- Paralel Postgres
 - 8.2 :(
- Master + segments
 - Specialized query planner
 - External tables (HDFS, Apache Parquet)
- Pretty mature technology
 - documentation
 - Instrumentation

Summary

- polyglot persistence
- queues
- pub/sub replication of production data !
- A good approximation instead of exact result
- SQL on top of Big Data is getting faster and faster...