

data infrastructure @ linkedin

figures

- >2M company pages
- 16 langages
- 4.2B prof. searches

architecture

- mainly on java
- scala
- layered
 - presentation
 - business service tier
 - data service tier
 - data infrastructure

Sid Anand

- linkedin
 - 2116 empl
 - montain view
- web dev
- search network analytics

databus

- oracles stores data
- data replicated in
 - search index
 - graph index
 - stadardization
 - read replicas
- data change events to replicate changes
- Relay (sharded)
 - capture change (from DB)
 - on-line changes bootstrap service
 - generate consistent snapshots and consolidated deltas
 - continuous updates
 - with long running queries
 - circular buffer
 - Guarantee order of changes

oracle

- all user provided data
- tens of physical instances
- tought problem : scaling writes
- scaling read
 - oracle instances
 - DSC token
 - databus
 - voldemort
 - timestamp comparison between
 - master
 - slave

voldemort

- distributed
- persistent
- K/V store influenced by AWS dynamo paper
- several self-healing mecanism
 - read repair
 - hinted handoff
 - anti-entropy repair
 - scans the entire dataset on a node and fixes it
 - to inssure consistency
- failure detection
- fat client : 3 layer going to move on the server
 - repair mecanism
 - failure detector
 - routing
- layered and abstract api's
 - easy to replace a layer
 - easy to extend
- better RO performance than mysql (latency)
- usages
 - RW store
 - RO store
- map reduce to generate
 - people you may know
 - jobs
 -

data infrastructure

- oracle
- voldemore
- expresso
- k/v store
- kafka
- databus
- data replication

Kafka

- messaging
- used for application level data streams
- for
 - end to end user request
 - metrics
- zookeeper
 - topic partition
 - offset management
- scale
 - B of events
 - TB of data
- design choices
 - zero copy
 - no double buffering
 - rely on the OS page cache
 - less presure on mem and GC