

NoSQL Roadshow Basel

NEAR REAL TIME PROCESSING OF SOCIAL MEDIA DATA WITH HBASE

Christian Gügi & Jean-Pierre König



- Why Hadoop and HBase?
- Social Media Monitoring
 - Prospective Search and Coprocessors
- Challenges & Lessons Learned
- Resources to get started

- Spin-off of MeMo News AG, the leading provider for Social Media Monitoring & Analytics in Switzerland
- Big Data expert, focused on Hadoop, HBase and Solr
- Objective: Transforming data into insights

NoSQL Roadshow Basel

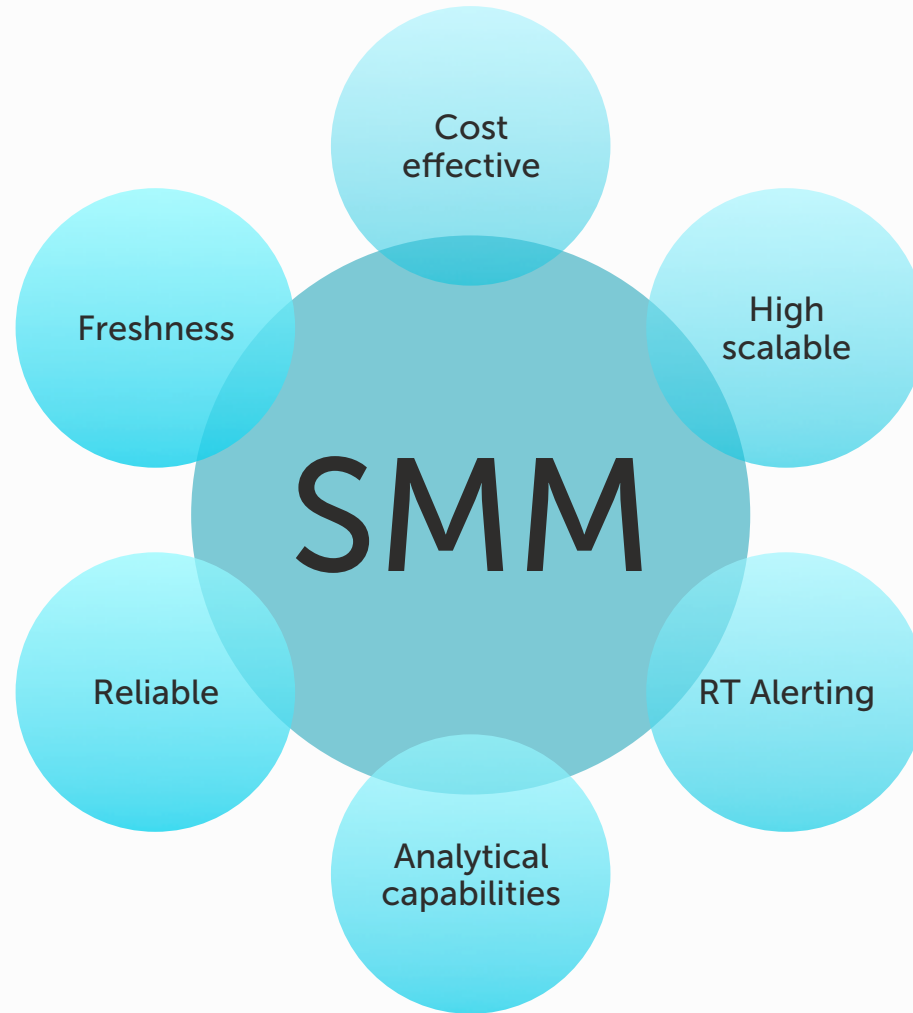
WHY HADOOP
AND HBASE?

Why?



Why Hadoop and HBase?

Social Media Monitoring Process



Why Hadoop and HBase?

Requirements

- HDFS + MapReduce
- Based on Google Papers
- Distributed Storage and Computation Framework
- Affordable Hardware, Free Software
- Significant Adoption

- Non-Relational, Distributed Database
- Column-Oriented
- Multi-Dimensional
- High Availability
- High Performance
- Build on top of HDFS as storage layer

Storage

HBase /HDFS

Search

Solr

Analytics

Hadoop

Mahout

Event mechanism (MQ)

HBase RowLog

Real-time alerting

Prospective search

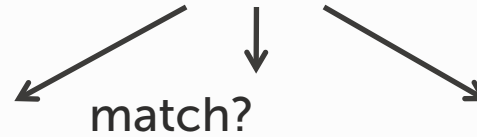
NoSQL Roadshow Basel

SOCIAL MEDIA MONITORING

**WHAT ARE
YOU
LOOKING AT?**



Downloaded Articles



Search Agents



Output



Web-UI

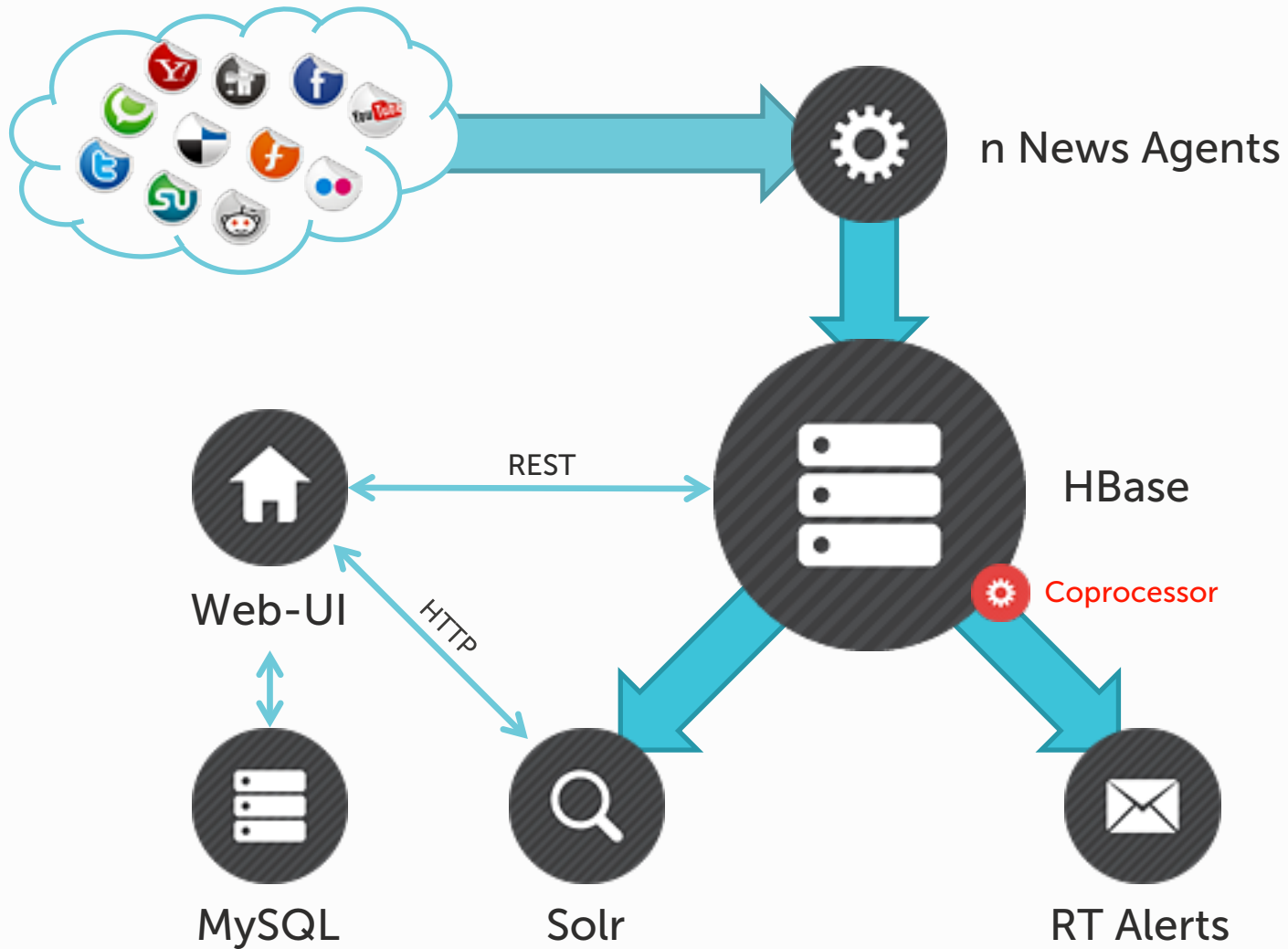


Reports



RT Alerts

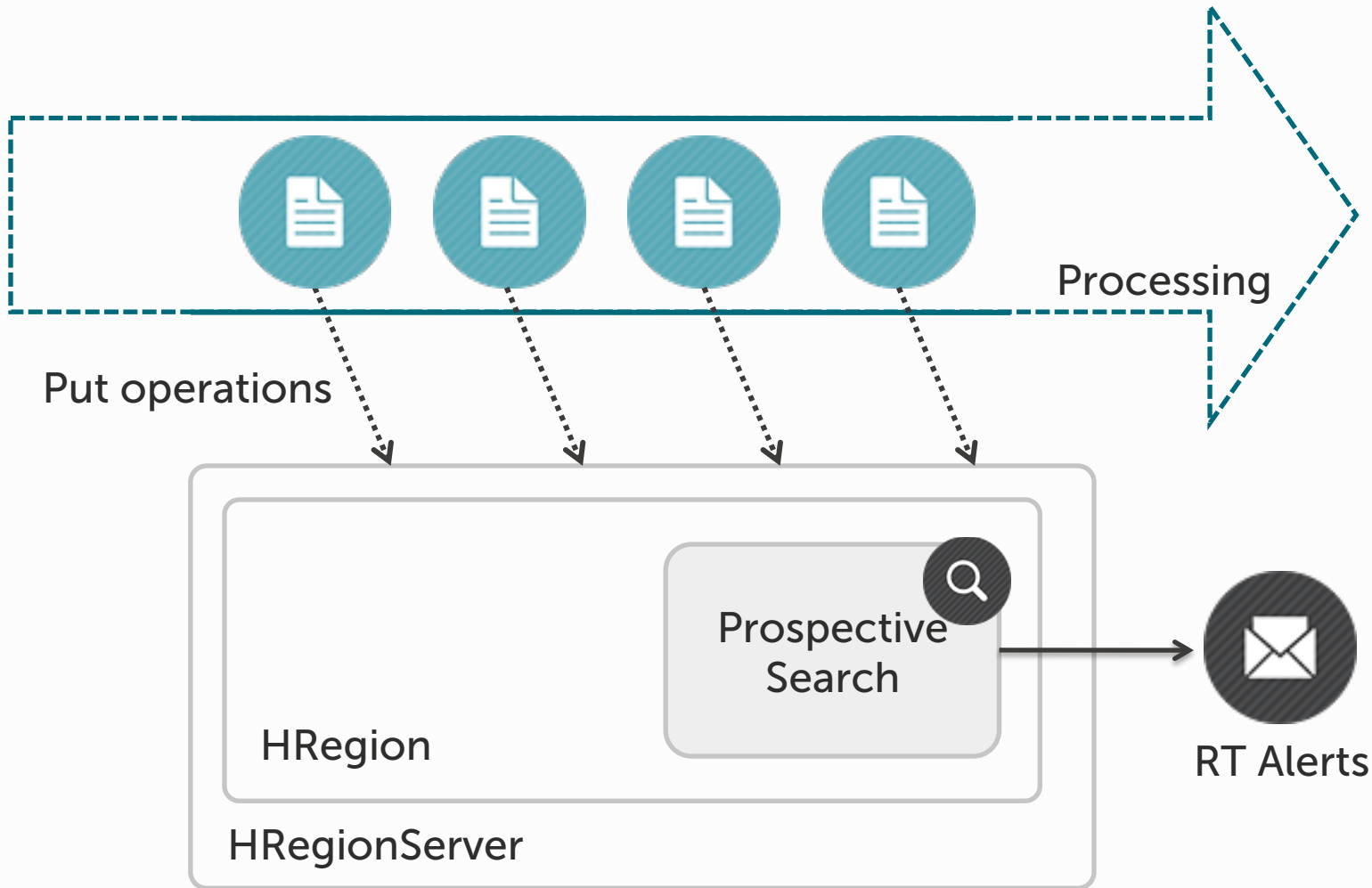
Icons by <http://dryicons.com>



Icons by <http://dryicons.com>

Social Media Monitoring

Solution Architecture



Icons by <http://dryicons.com>

Social Media Monitoring

Prospective Search with Coprocessors

- Monthly growth
 - Index: 200GB
 - 50 Mio. docs/month
 - HBase: 600 GB
 - Raw data, meta data and extracted data
- A few 1000 map-reduce jobs/month

NoSQL Roadshow Basel

CHALLENGES & LESSONS LEARNED

NORWEGIAN WOOD

TEACH US TO OUTGROW

CUPPY / STEIG

The Decline and Fall of
Practically Everybody

In the Realm of a Dying Emperor

MURKAWA

Kenyan

- 1 Benchmarks - workloads
- 2 Supervision
- 3 Keys and shards – Schema design /LG
- 4 Timestamps, the 4th dimension
- 5 Short ColumnFamily names->
- 6 File handles. OS
- 7 JVM Tuning, GC !!!
- 8 Scaling region servers, data locality!
- 9 Automatic vs manual splits, compaction
- 10 Do not use HBase as rock solid in prod
- 11 Forget feuerwehr aktionen, it takes some time
- 12 Use Hbase for a appropriate use case
- 13 Tune and tweak – it's not a project – it's a process
- 14 You need devops in production
- 15 Huge know-how curve, you need to know the hole ecosystem
- 16 Use a distribution, ist packed, tested and supports migration, enterprise grade
- 17 Virtualisierung, Hardware
- 18 Dont struggle to much, there is a good community
- 19 Share your knowledge
- 20 It's early state, many tools around, a few still missing

- Everyone is still learning
- Some issues only appear at scale
 - At scale, nothing works as advertised
- Production cluster configuration
 - Hardware issues
 - Tuning cluster configuration to our work loads
- HBase stability
- Monitoring health of HBase

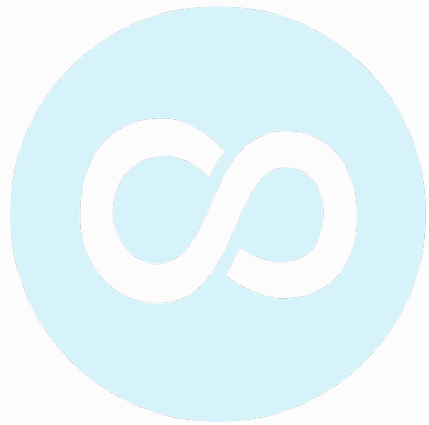
- Do not rely on HBase as frontend storage layer. It's not going to be rock solid
- Don't struggle to much, there is a good community
- Share your knowledge
- It's early stage, many tools around, a few still missing

- Use HBase for an appropriate use case
- Use a distribution, its packed, tested and supports migration, enterprise grade
- Benchmarks – know your workloads & query patterns
 - YCSB
- Schema & Key Design
 - What's queried together should be stored together
- Scaling region servers, data locality!
- Virtualization vs. Real Hardware

- Number of CF < 10
 - Compaction + Flushing I/O intensive
- Short ColumnFamily names
 - HFile index size occupying aloc RAM (storefileindexSize)
- OS file handles
 - ulimit -n 32768
- JVM Tuning, GC !!!
 - HMaster 1024 MB
 - RegionServer 8192 MB
 - -XX:+UseConcMarkSweepGC
 - -XX:+CMSIncrementalMode
- Automatic vs. manual splits
- Be careful with expensive operations in coprocessors
- Play with all the configurations and benchmark for tuning

- Monitoring/Operational tooling is most important
- Forget “emergency actions”, it takes some time
- Tune and tweak – it’s not a project – it’s a process
- You need DevOps in production
- Huge know-how curve, you need to know the whole ecosystem
 - Hadoop, HDFS, MapRed

- <http://hbase.apache.org/book.html>
- <http://www.sentric.ch/blog/best-practice-why-monitoring-hbase-is-important>
- <http://www.sentric.ch/blog/hadoop-overview-of-top-3-distributions>
- <http://www.sentric.ch/blog/hadoop-best-practice-cluster-checklist>
- <http://outerthought.org/blog/465-ot.html>



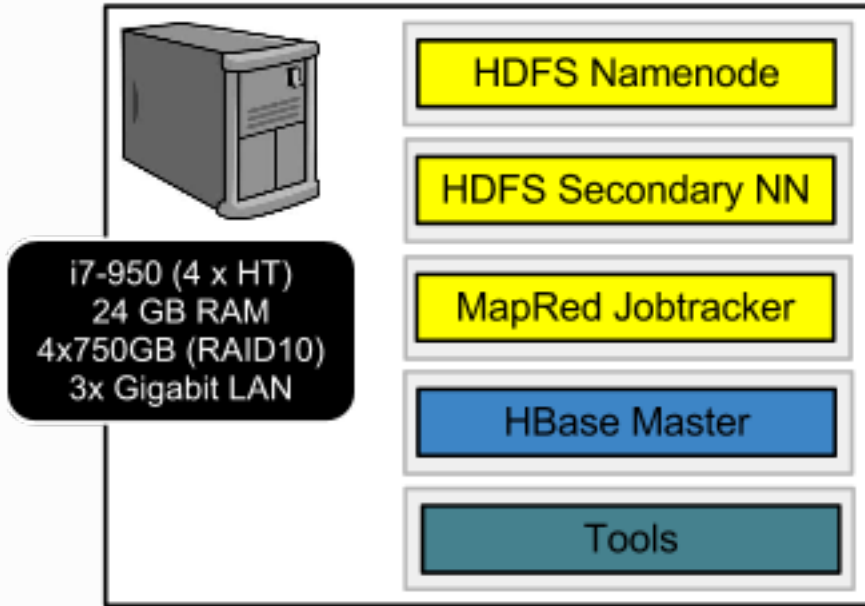
entric

Questions?

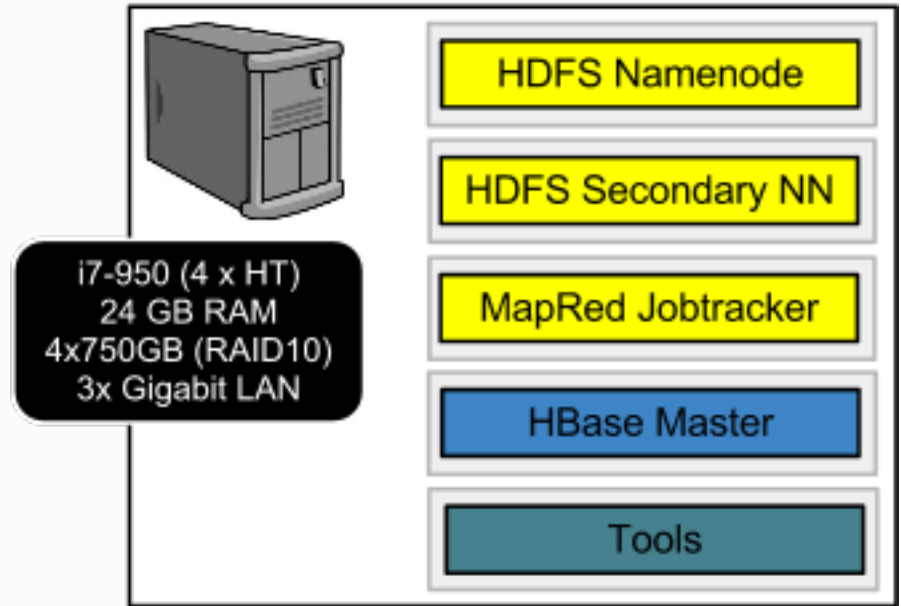
Christian Gügi, christian.guegi@entric.ch

Jean-Pierre König, jean-pierre.koenig@entric.ch

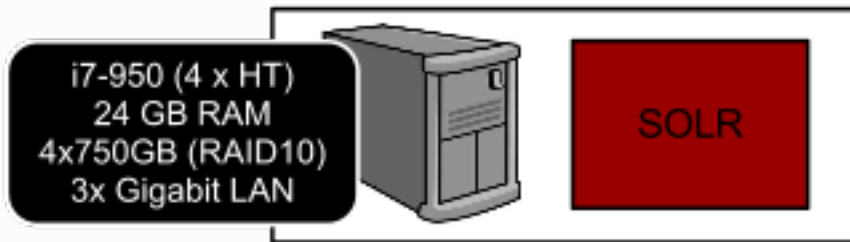
Virtual Master Host 1



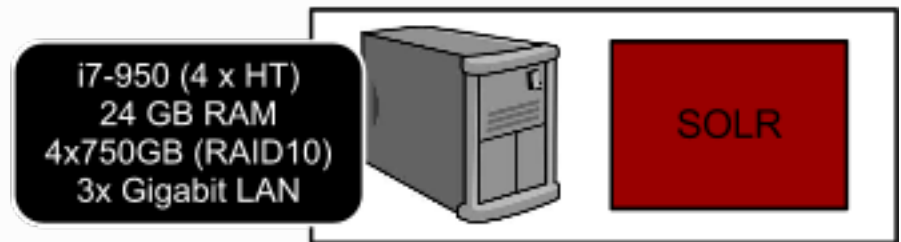
Virtual Master Host 2



Search - 1



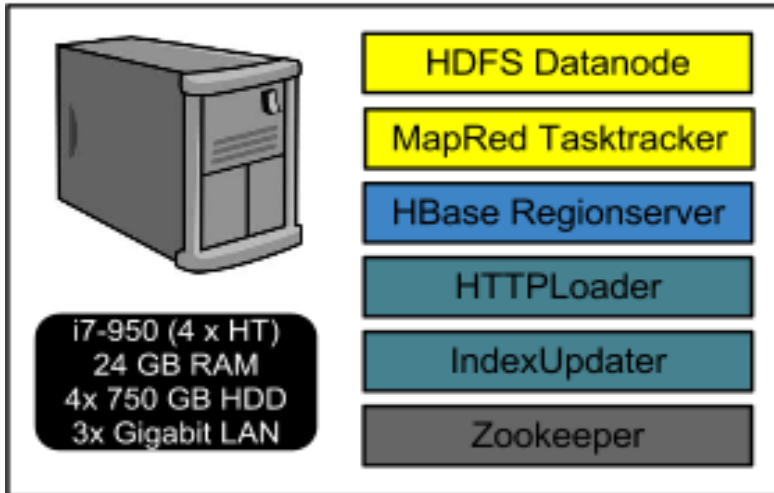
Search - 2



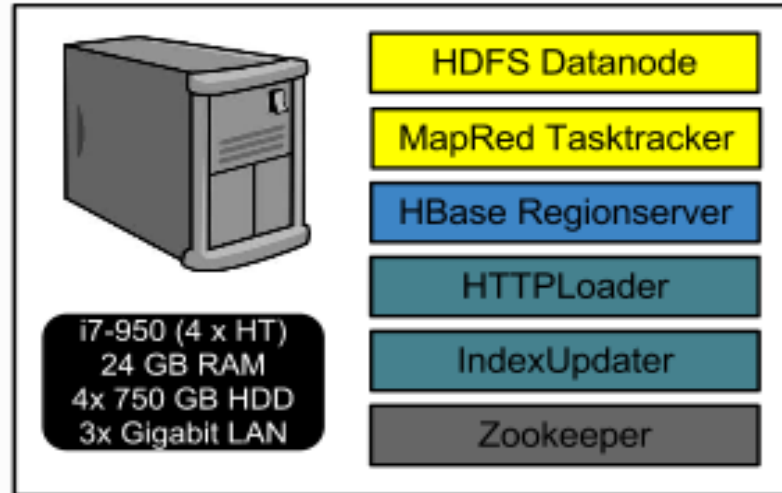
Masters

Cluster

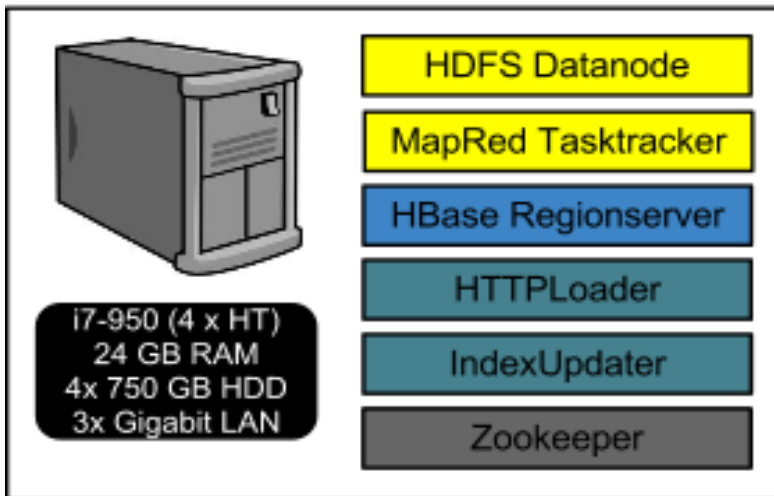
Worker-1



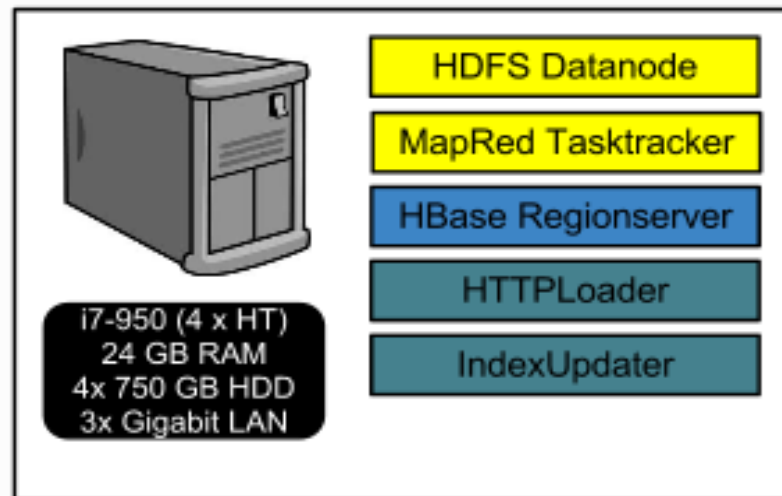
Worker-2



Worker-3



Worker-4



Worker

Cluster