

Berlin Buzzwords 2012

USING HBASE COPROCESSORS TO IMPLEMENT PROSPECTIVE SEARCH

Christian Gügi, Software Architect



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



5. Juni
2012

3

Software Architect
@ sentric

Co-founder and
organizer of the
Swiss HUG

Contact:

christian.guegi@sentric.ch

<http://www.sentric.ch>

@chrisgugi

About me

Berlin
buzz
words
search. share. scale

 sentric

- 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

Berlin Buzzwords 2012

WHY HADOOP
AND HBASE?

NOT.

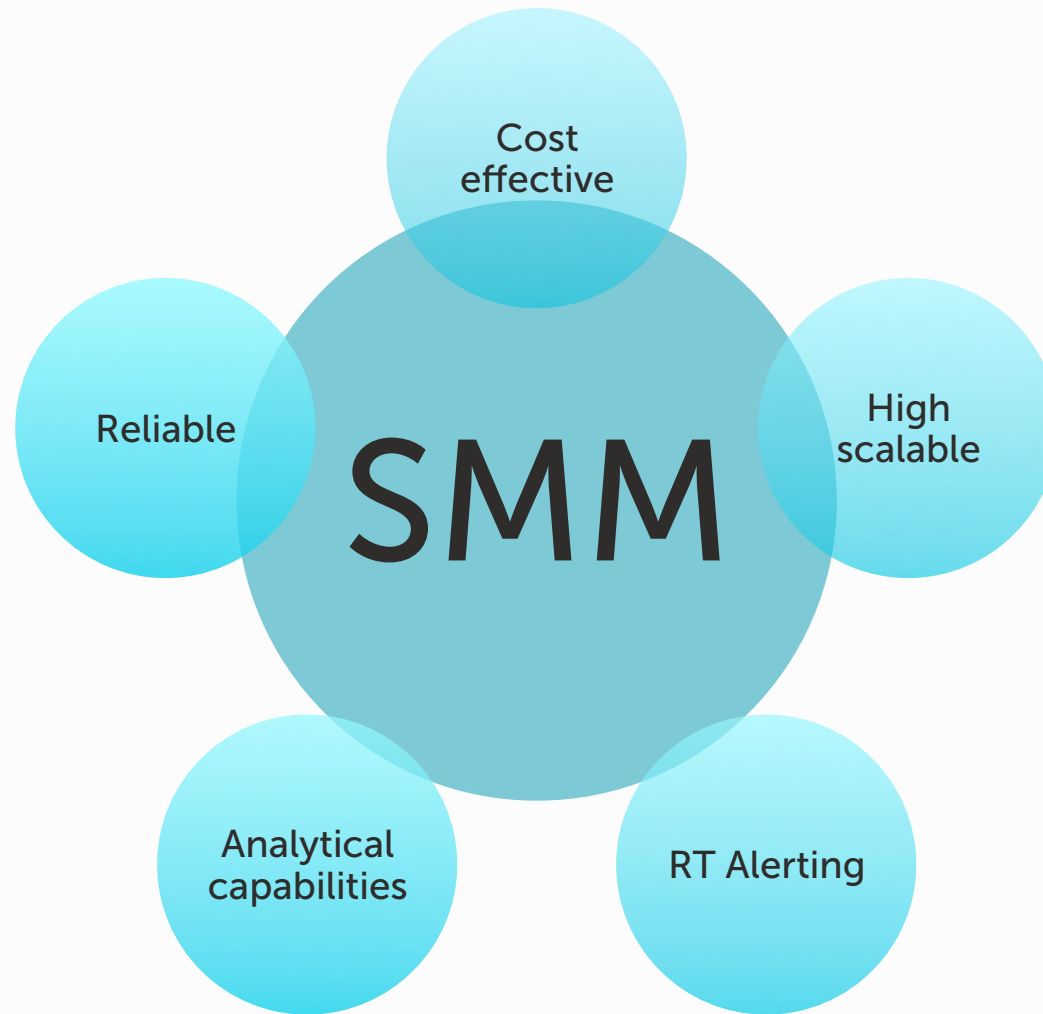


Why Hadoop and HBase?

Social Media Monitoring Process

Berlin
buzz
words
search. share. scale

 sentric



Why Hadoop and HBase?

Requirements

Storage

HBase /HDFS

Search

Solr

Analytics

Hadoop

Mahout

Event mechanism (MQ)

HBase RowLog

Real-time alerting

Prospective search

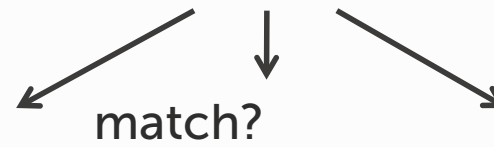
Berlin Buzzwords 2012

SOCIAL MEDIA MONITORING

WHAT ARE
YOU
LOOKING AT?



Downloaded Articles



Search Agents



Output



Web-UI

Reports

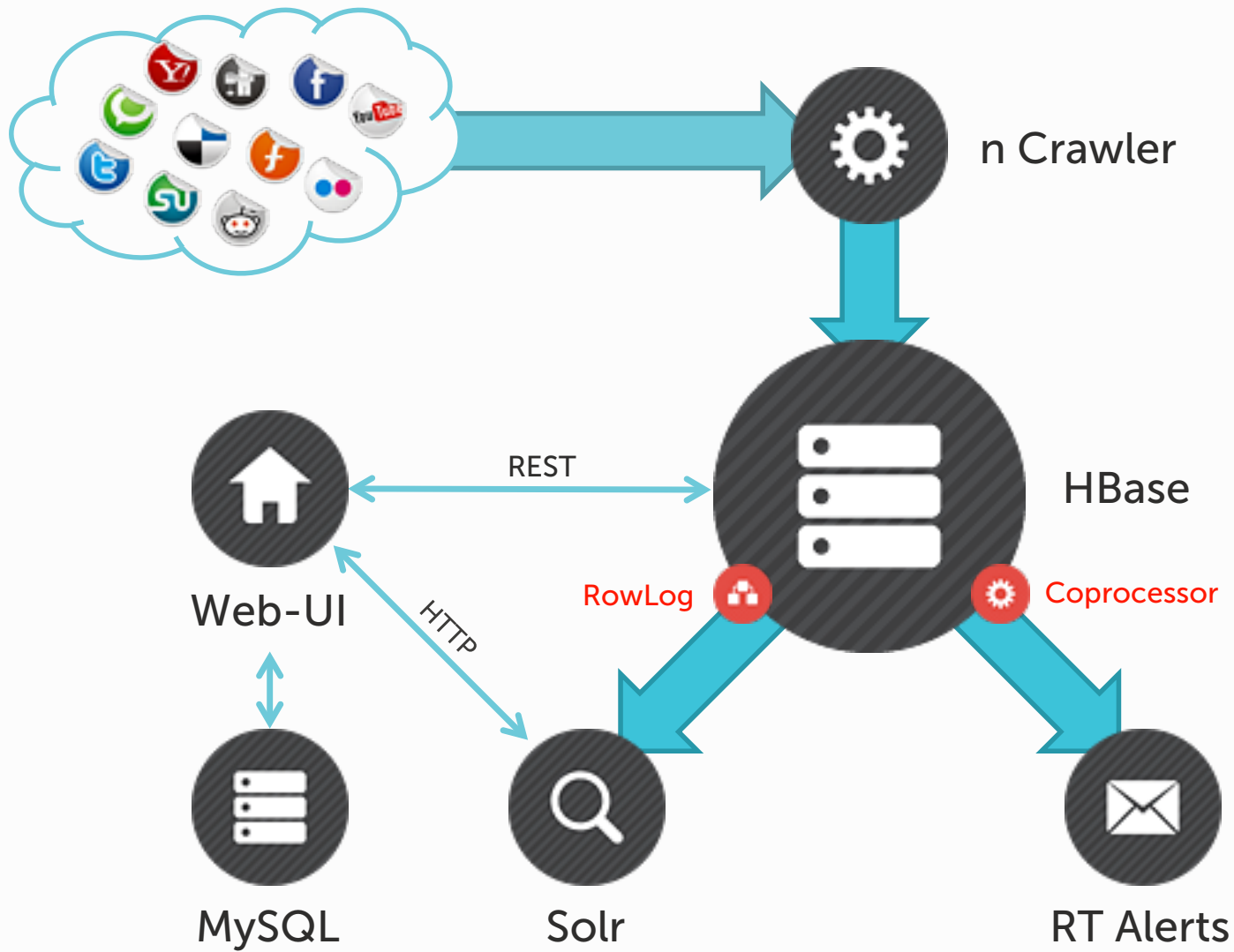
RT Alerts

Icons by <http://dryicons.com>

Social Media Monitoring

Overview

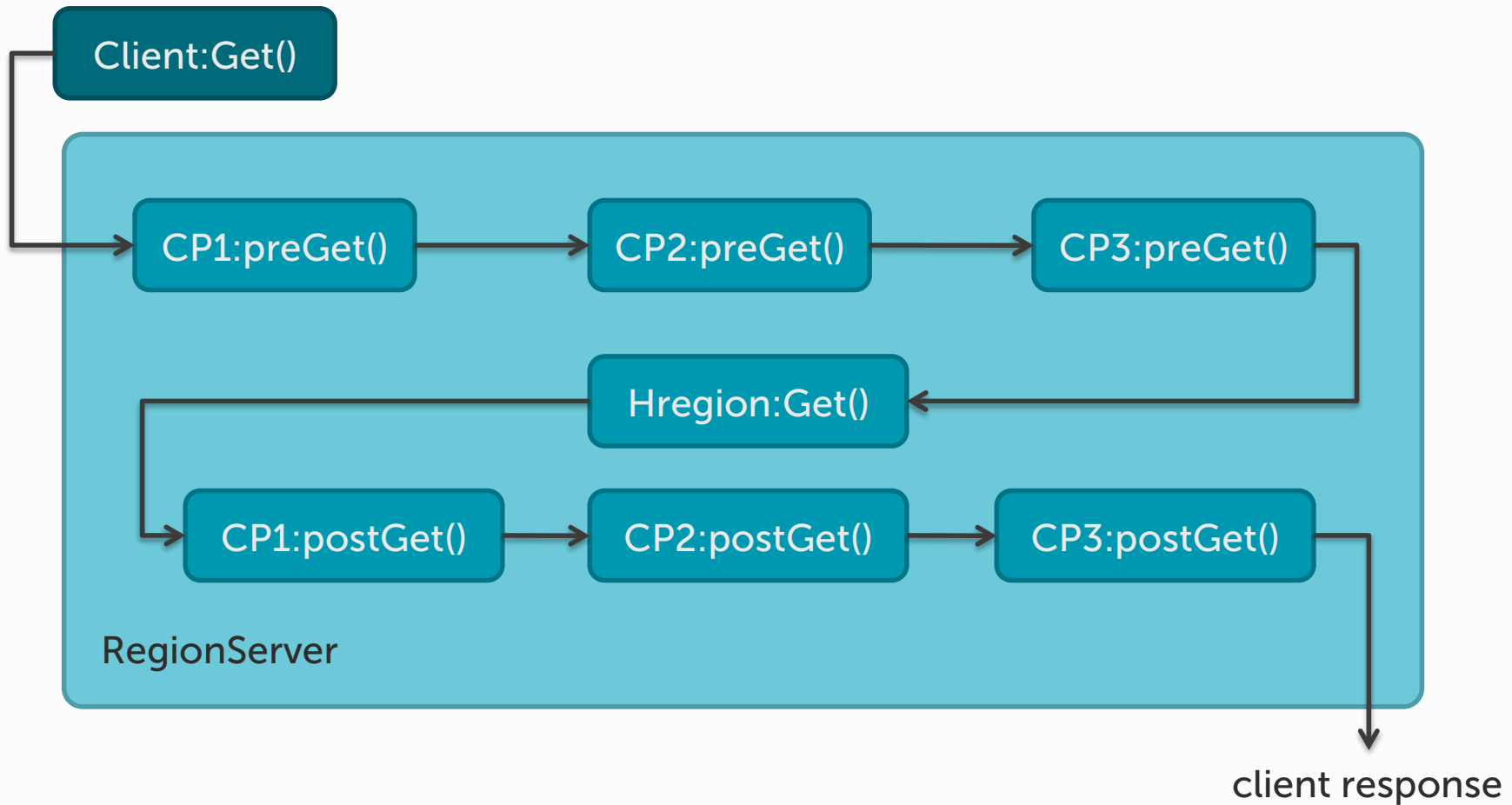




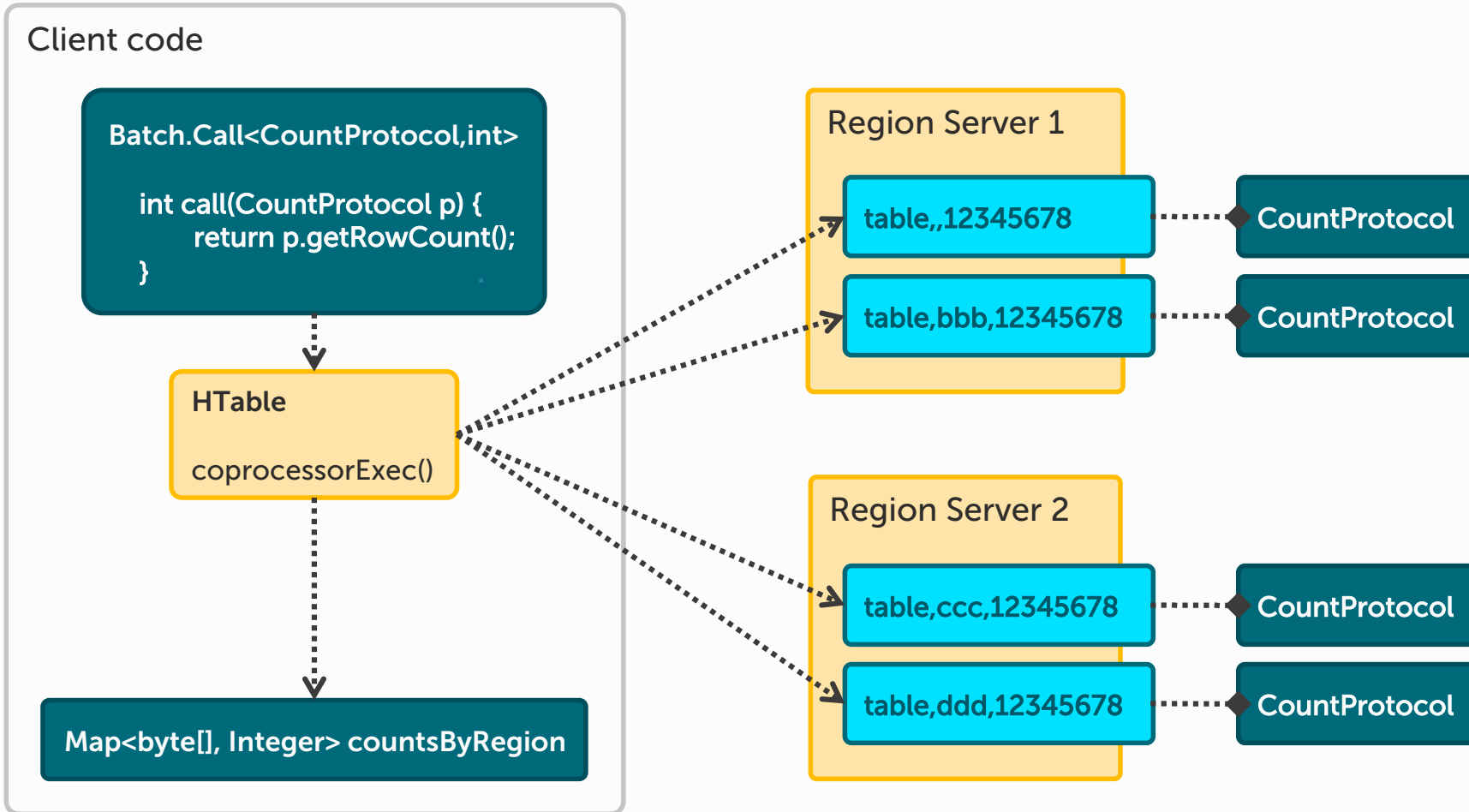
Icons by <http://dryicons.com>

- Inspired by Google Bigtable coprocessors
- HBase version 0.92
- Embed code directly into server processes
- High-level call interface for clients
- Automatic scaling, load balancing, request routing

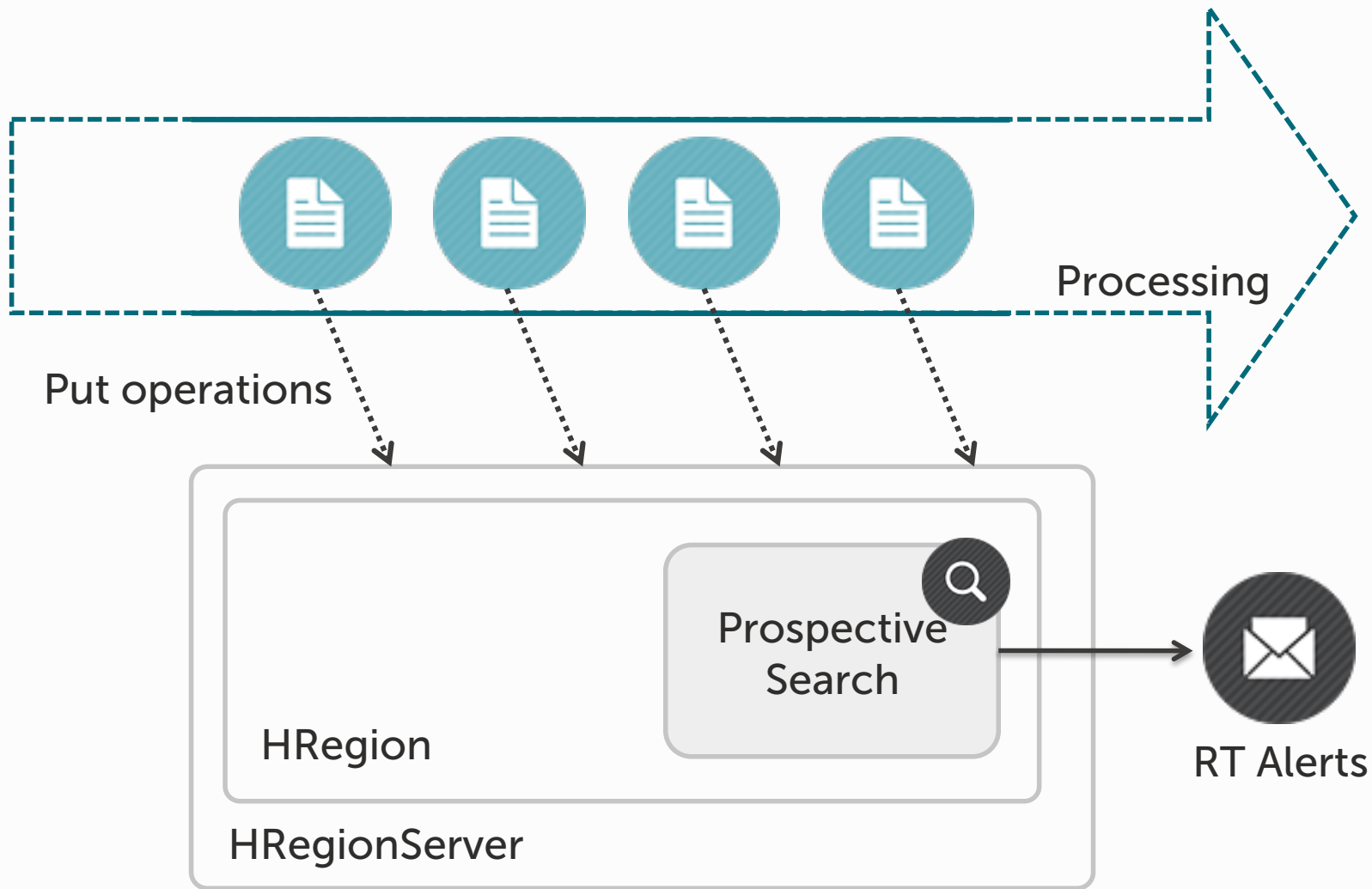
- Like a database trigger
 - Provides event based hooks
- Concrete Implementations
 - RegionObserver
 - CRUD or DML type operations
 - MasterObserver
 - DDL or metadata operations and cluster administration
 - WALObserver
 - Write-ahead-log appending and restoration



- Comparable to stored procedures
 - Custom RPC protocol, used between client and region server
- Loaded in region server
- Client call APIs over single row or a row range
 - Framework translates row keys to region location
 - Parallel execution

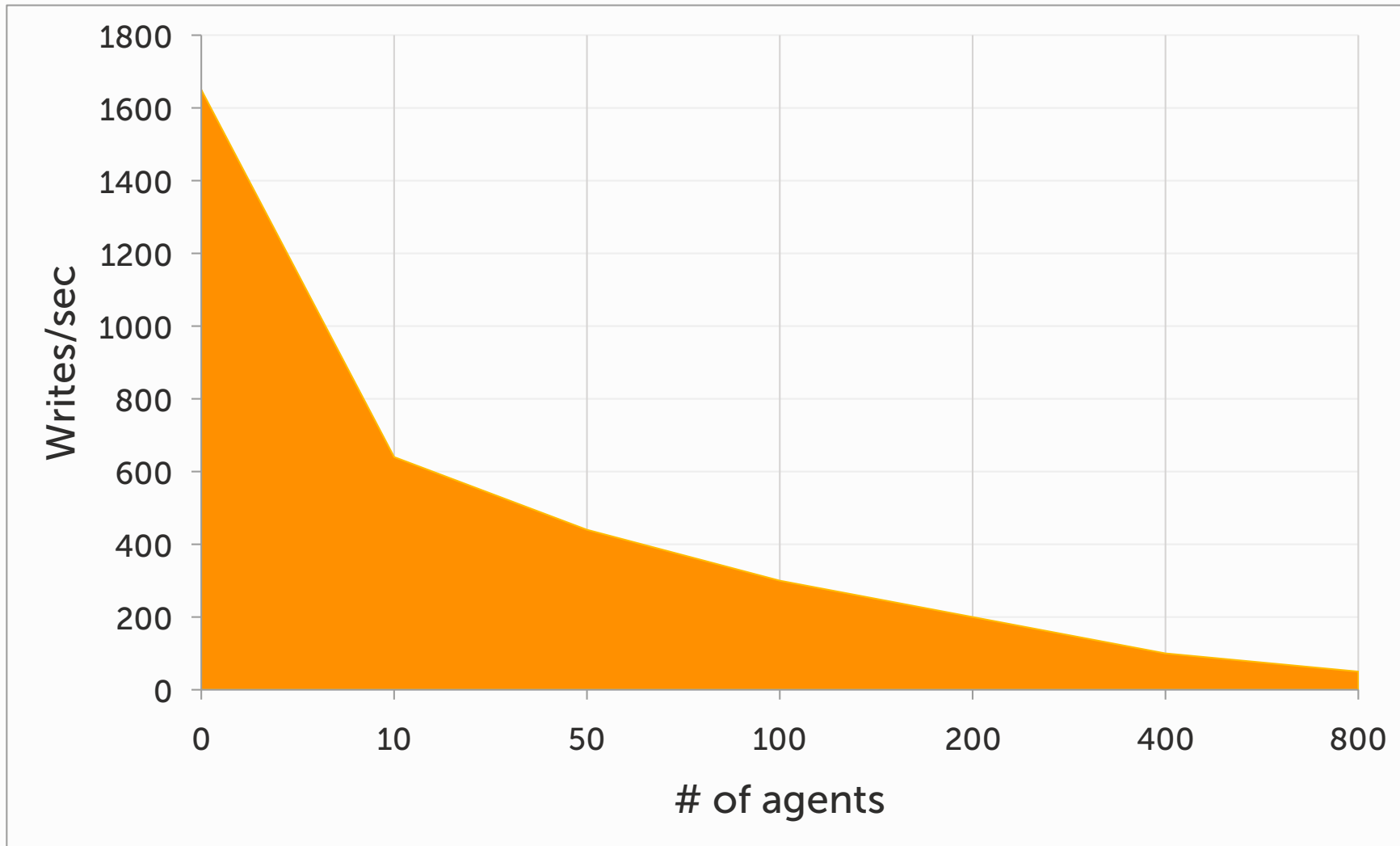



- HBase Security (Version 0.94)
- Aggregate operations avg(), sum()
 - AggregatorProtocol
- HBASE-3529: Embedded search



Icons by <http://dryicons.com>

- Standard, virtualized test cluster:
4RS/DN, 1HM, 1NN, 3ZK
- Test dataset created from 2h of live
index (1GB)
- Drive load on RS/DN





Berlin Buzzwords 2012

CHALLENGES & LESSONS LEARNED

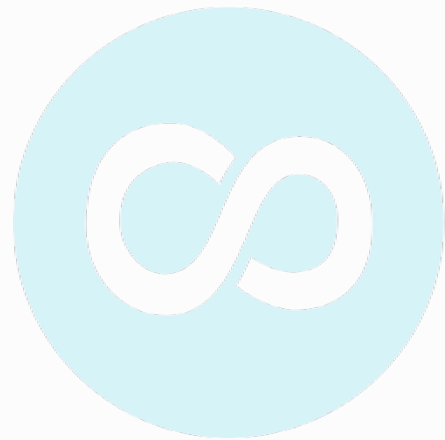
- Everyone is still learning
- Some issues only appear at scale
- Production cluster configuration
 - Hardware issues
 - Tuning cluster configuration to our work loads
- HBase stability
- Monitoring health of HBase

- Be careful with expensive operations in coprocessors
- At scale, nothing works as advertised
- Monitoring/Operational tooling is most important
- Play with all the configurations and benchmark for tuning

- https://blogs.apache.org/hbase/entry/coprocessor_introduction
- <http://hbase.apache.org/apidocs/index.html>
- <http://www.lilyproject.org/lily/about/playground/hbaserowlog.html>
- <http://www.github.com/sentric/HBasePS>

5. Juni
2012

25



sentric

Questions?

Christian Gügi

christian.guegi@sentric.ch

Berlin Buzzwords 2012

Thank you!

Berlin
buzz
words
search. share. scale

 sentric