



Big Time-series Data

Overview of related projects and activities at CIPSI

Tomasz Wiktor Wlodarczyk

Associate Professor

Center for IP-based Service Innovation

University of Stavanger

tomasz.w.wlodarczyk@uis.no

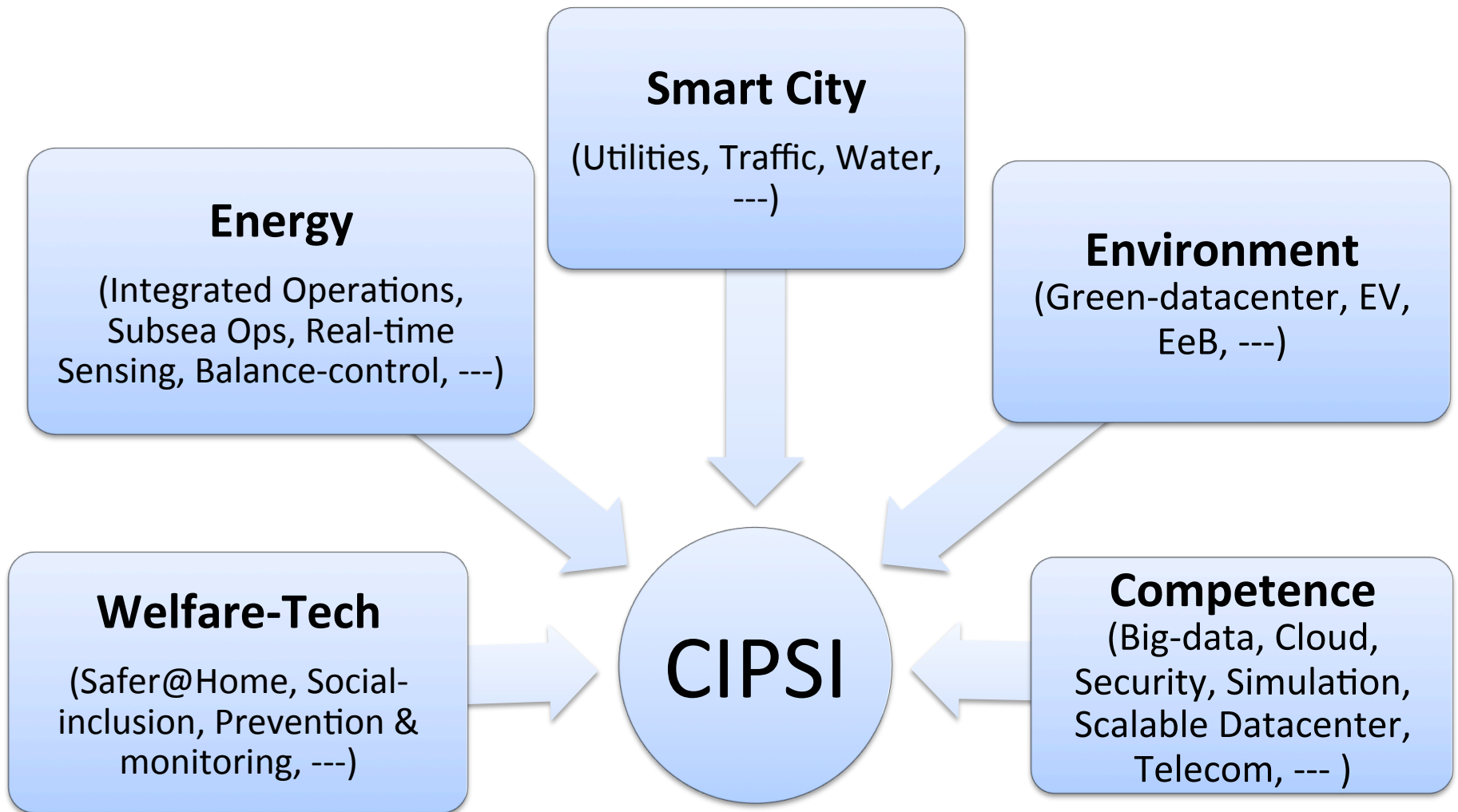
Overview

- Several applied ICT activities revolving around **Big Data**
- Many have significant focus on **time series**
- **R2Time** is a library for time series analysis stored in OpenTSDB over a Hadoop cluster based on extended RHIPE library
- **Data Intensive Systems course** developed jointly with Purdue University and AMD Research
- **CBOK of Data Intensive Science** with Research Data Alliance

Overview

- CIPSI
- SEEDS
- Uninett
- R2Time
- A4Cloud
- Safer@Home
- Data Intensive System course

Innovation on Smart ICT Services



CIPSI

- Faculty
 - Chunming Rong
 - Tomasz Wiktor Wlodarczyk
- Researchers and Post Docs
 - Rui Maximo Esteves
- PhD students
 - Antorweep Chakravorty
 - Aryan TaheriMonfared
 - Bikash Agarwal
 - Ali Abassi
- ~4 Master Students

Overview

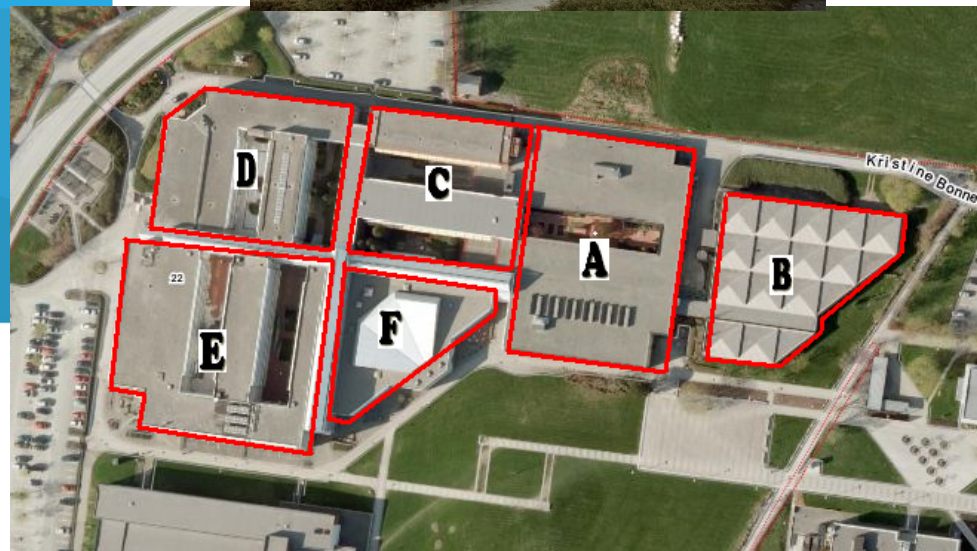
- CIPSI
- **SEEDS**
- Uninett
- R2Time
- A4Cloud
- Safer@Home
- Data Intensive System course



Self learning Energy Efficient buildDings and open Spaces



<http://seeds-fp7.eu>



UiS Demo Space

SEEDS

Self learning Energy Efficient buildings and open Spaces

ENTER

Projects Details:

- Project Acronym : SEEDS
- Grant Agreement No. 285150
- Framework Programme: FP7
- Start Date : 01 - 09 - 2011
- End Date : 31 - 08 - 2014
- Duration : 36 months
- Project Cost : 4.081.646 euros
- Project Funding : 2.898.966 euros
- Project Status : Execution



This website has been produced with funding received from the European Community's Seventh Framework Programme under Grant Agreement N° 285150.



This website is the property of the SEEDS consortium and shall not be distributed or reproduced without the formal approval of the SEEDS General Assembly.

EeB-ICT-2011.6.4 ICT for energy-efficient buildings and spaces of public use

VICIONE

Devices Global Rooms

Freitag 25. Jan 12:39

7° SW 2

Freitag 5° 8° Samstag 5° 13° Sonntag 3° 11° Montag 1° 12°

- Room 1
- Room 2
- Room 3
- Room 4
- Room 5
- Room 6
- Room 7
- Room 8
- Room 9
- Room 10

Room 1

Room 2

Room 3

Room 4

Room 5

Room 6

Room 7

Room 8

Room 9

Room 10

Set Comfort Temperature Maximum Delta

2.0 °C

Occupancy

Set Number of Persons by Time

00:00	08:00	12:00	16:00	20:00
0	1	9	5	0

Solar Radiation Intensity

Set Solar Radiation Intensity by Time

00:00	08:00	12:00	16:00	20:00
1,00 W/m ²	0,06 W/m ²	0,07 W/m ²	0,10 W/m ²	0,09 W/m ²

Data Logging

Current Temperature

Energy Consumption

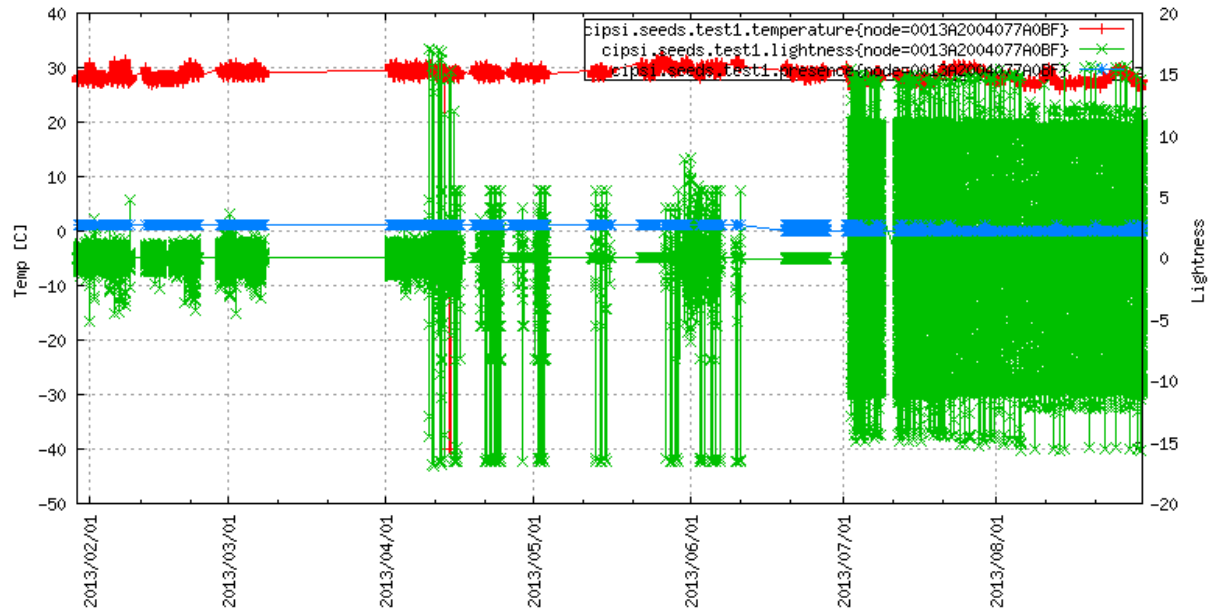
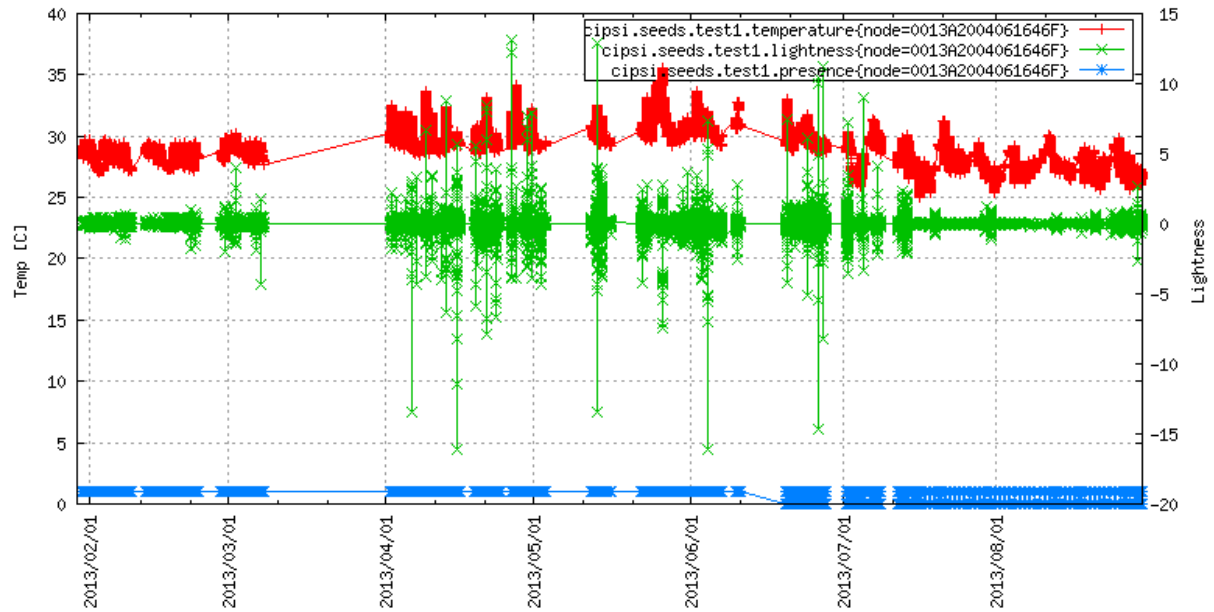
SEEDS



OpenTSDB

- How to store large amount of sensor data
 - frequently sampled (often every minute)
 - store forever with original precision
 - make available for analysis
- Traditional solutions, relational databases or local time series databases either cannot handle so much data or require significant workarounds
- OpenTSDB is a distributed, scalable Time Series Database (TSDB) written on top of HBase.
- HBase is a database built on top of HDFS
- HDFS is Hadoop Distributed File System
- OpenTSDB provides standardized and efficient schema and API for time-series data

SEEDS



SEEDS

- Pre-demo phase 7 sensor nodes (from Jan 2013)
- Demo phase >100 sensor nodes (starting soon)
- 8 sensors per node
- Archive data from all sensors and enable post demo evaluation
- Enable self-learning and optimization of energy saving algorithms
- Possibly archive all predictions made by optimization algorithm for later analysis

Overview

- CIPSI
- SEEDS
- **Uninett**
- R2Time
- A4Cloud
- Safer@Home
- Data Intensive System course

Uninett

Traffic Type	Statistics/day		
	Avg	Max	Min
Distinct Source IPs	987104	4740760	122266
Distinct Source IPs and Source ports	6083640	13188647	844898
Distinct Destination IPs	1613040	2488893	420686
Distinct Destination IPs and Destination ports	7010330	16379274	1113095
Distinct Bidirectional flows	10683200	21454096	1829854
NetFlow records	21962800	44036078	4373665

Uninett

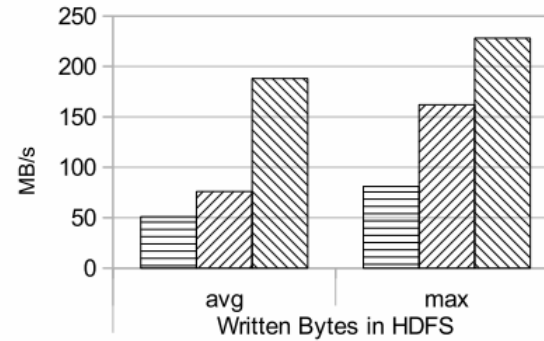
Table	Row Key					Query Type
T1	[sa]	[sp]	[da]	[dp]	[1 - ts]	Extended queries
T2	[da]	[dp]	[sa]	[sp]	[1 - ts]	
T3	[sa]	[da]	[sp]	[dp]	[1 - ts]	Source-Destination address queries
T4	[da]	[sa]	[dp]	[sp]	[1 - ts]	Source-Destination address queries
T5	[sp]	[sa]	[da]	[dp]	[1 - ts]	Service server discovery queries
T6	[dp]	[da]	[sa]	[sp]	[1 - ts]	Service server discovery queries
T7	[sp]	[da]	[sa]	[dp]	[1 - ts]	Service client discovery queries
T8	[dp]	[sa]	[da]	[sp]	[1 - ts]	Service client discovery queries

	Est. # records	Storage for T1	Storage for T2-T8	Storage for OpenTSDB	Total
Single Record	1	$(37 * 23) + (133) \sim 1KB$	$7tables * 23B = 161B$	$5metrics * 2B = 10B$	$\sim 1KB$
Daily Import	$\sim 20million$	$1KB * 20 * 10^6 = 20GB$	$161B * 20 * 10^6 \sim 3GB$	$10B * 20 * 10^6 = 200MB$	$\sim 23GB$
Initial Import	$20m * 150days \sim 3 * 10^9$	$1KB * 3 * 10^9 = 3TB$	$161B * 3 * 10^9 \sim 500GB$	$10B * 3 * 10^9 = 30GB$	$\sim 3.5TB$

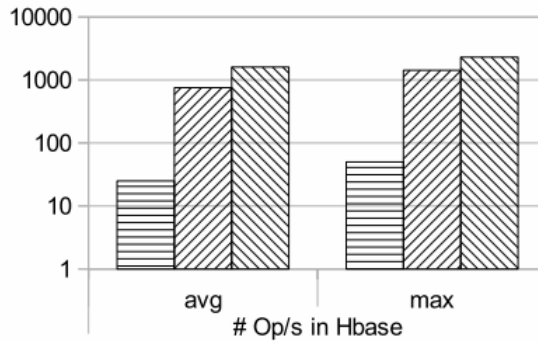
Uninett



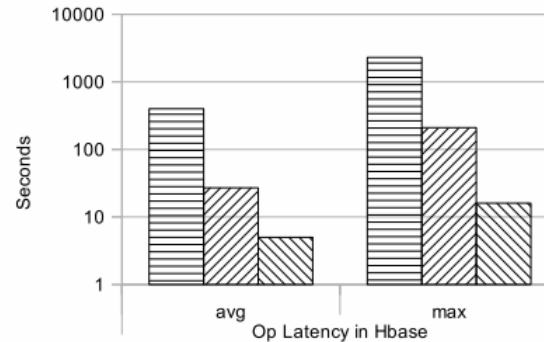
(a) Jobs finishing time



(b) HDFS IO



(c) Number of operations per seconds in HBase



(d) Operation latency in HBase

Overview

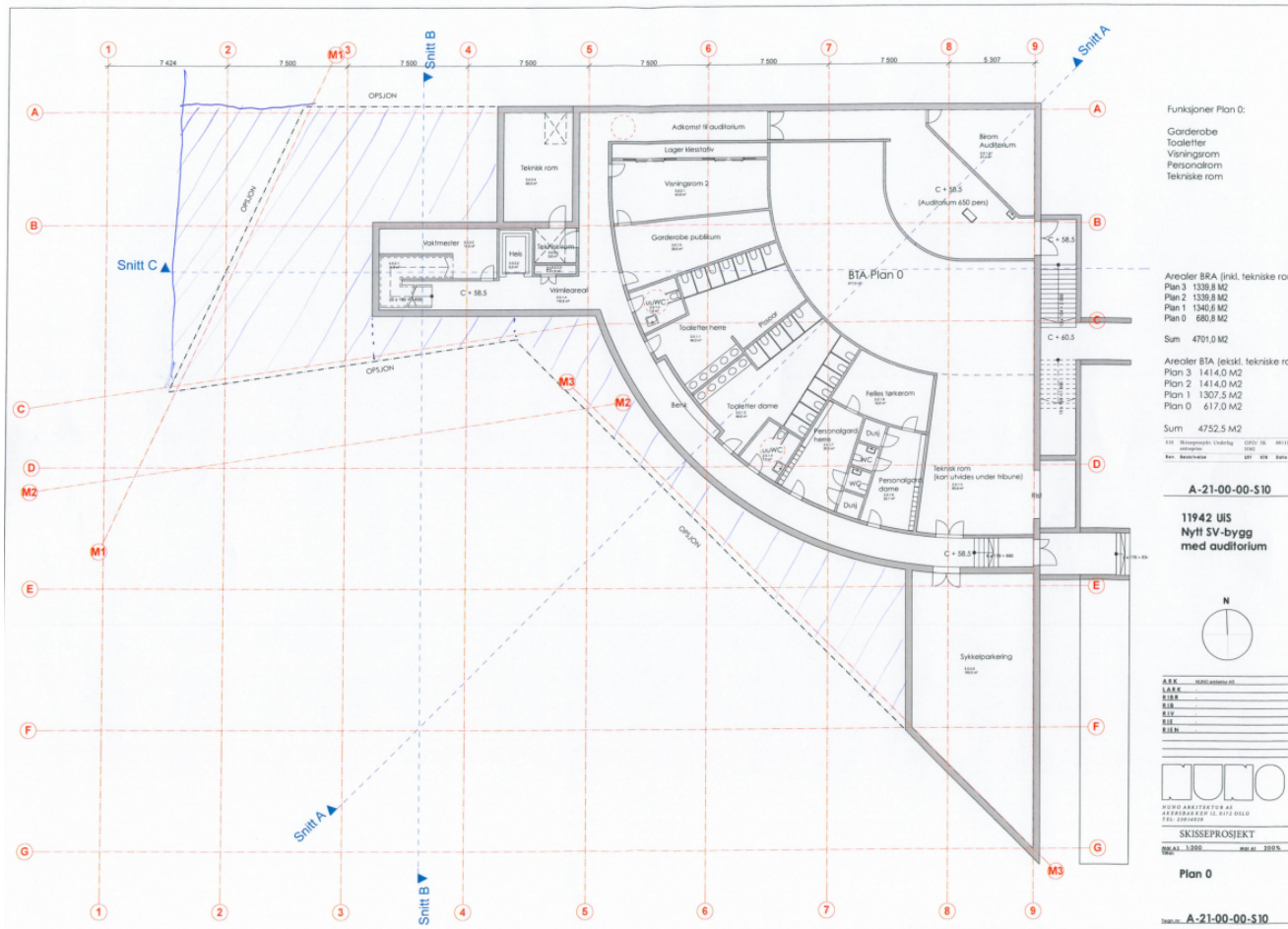
- CIPSI
- SEEDS
- Uninett
- **R2Time**
- A4Cloud
- Safer@Home
- Data Intensive System course

Data Cluster Infrastructure



39 nodes
1 TB RAM
300 TB Disk

Infrastructure extension – Data Center



- 200 – 300 m²
- Central water cooling
- Ca. 20 racks, 10 operational initially
- Construction to start in Fall 2013
- Expected finish Fall 2015
- Off-site backup and test site at Green Mountain Data Centre through Lyse in 2013/2014

Infrastructure extension – Off-site Backup



Infrastructure extension – Off-site Backup



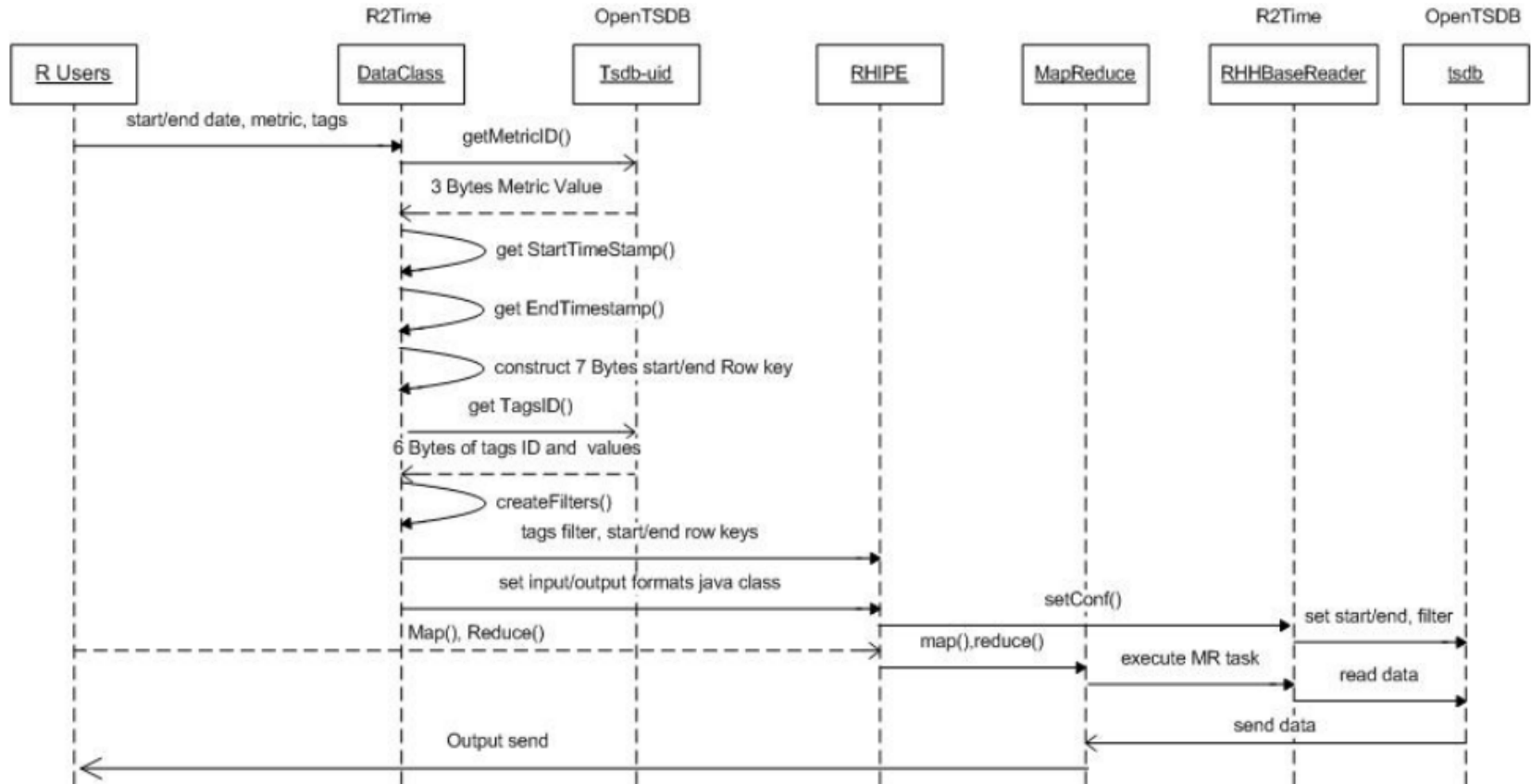
Beyond storage

- OpenTSDb provides storage and basic visualization capabilities
- There are no analytic capabilities
- Currently the only way to analyze data is to export it out e.g. to R (v. opentsdbR from Berkeley)
- It is not suitable for large analysis
- Data should be analyzed in place (move computation not data)

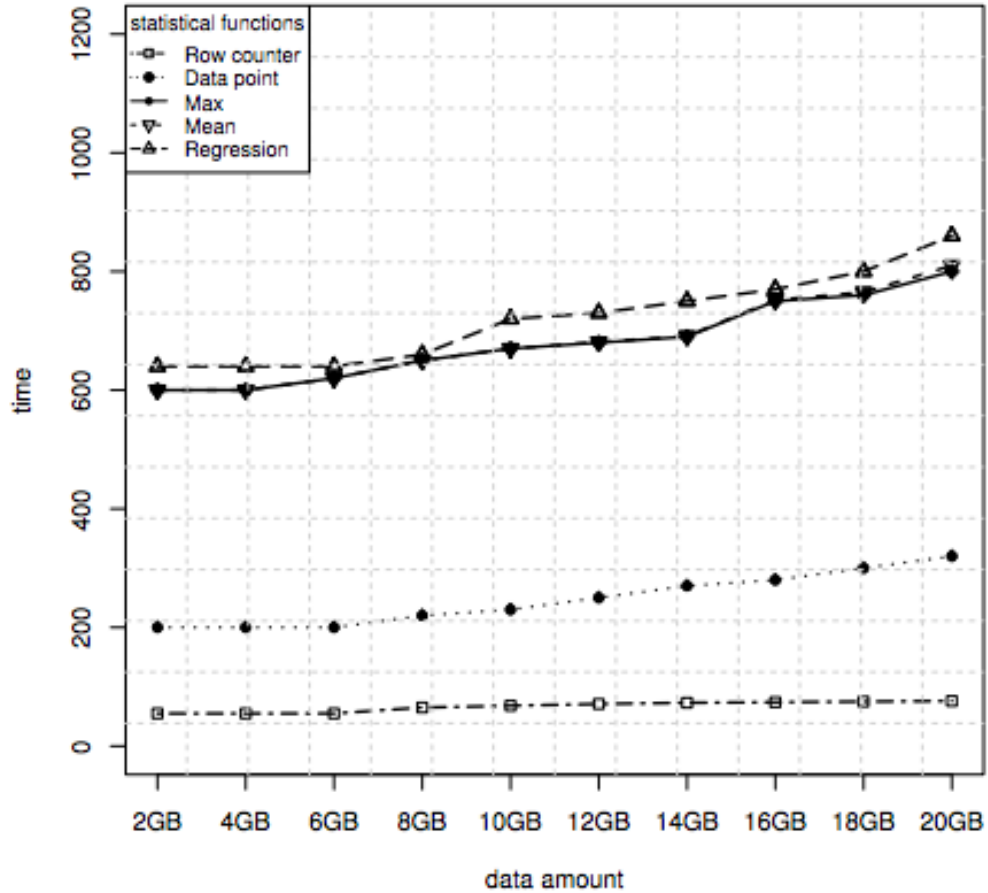
Beyond storage

- RHIPE is one of two major libraries to run R over Hadoop
- RHIPE support for HBase is limited, it is mostly focused on HDFS
- OpenTSDB key structure is complex
- We created R2Time library
 - Fixes some of the RHIPE HBase problems
 - Provides support to handle OpenTSDB key structure, use provides only timestamps
 - Improves performance by ensuring uniform data distribution over non-uniform time samples

R2Time



R2Time



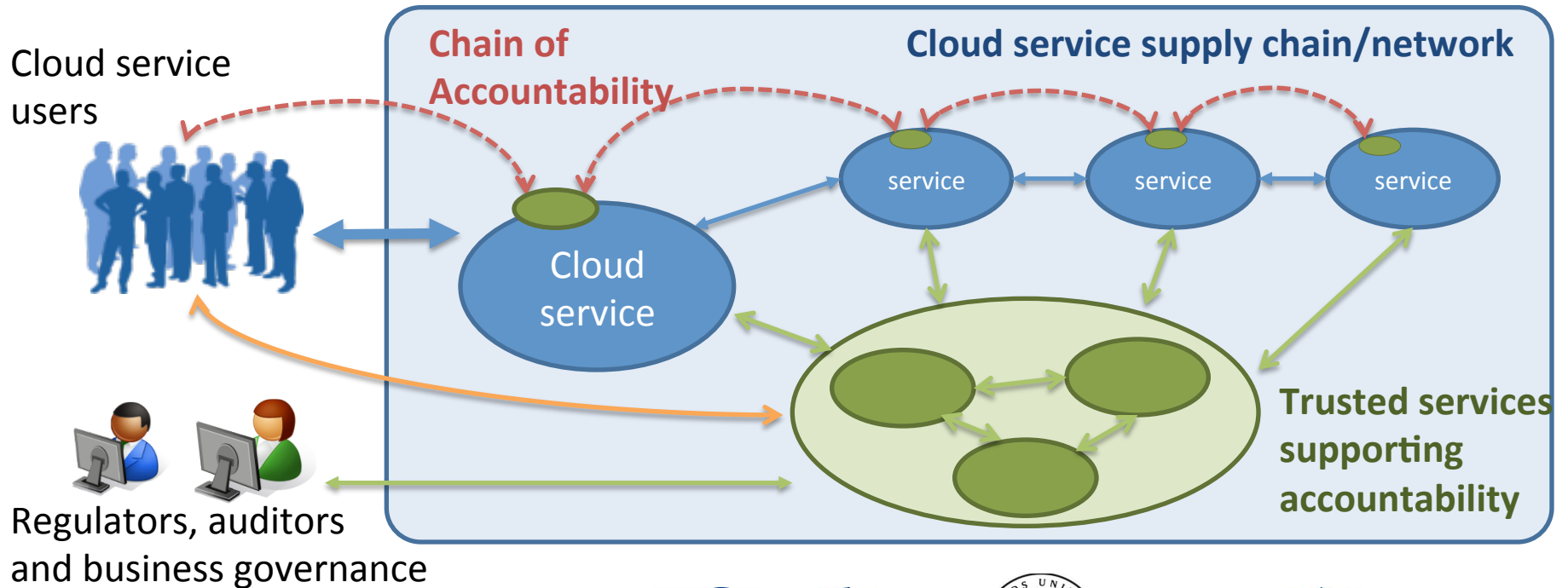
R2Time – next steps

- Move from alpha to beta – standardize and simplify API
- Implement typical analytic functions – so that users without MapReduce knowledge can also use the library
- Integrate typical functions with OpenTSDB interface
- Open source

Overview

- CIPSI
- SEEDS
- Uninett
- R2Time
- **A4Cloud**
- Safer@Home
- Data Intensive System course

A4Cloud – Accountability for Cloud and other Future Internet Services



A4Cloud

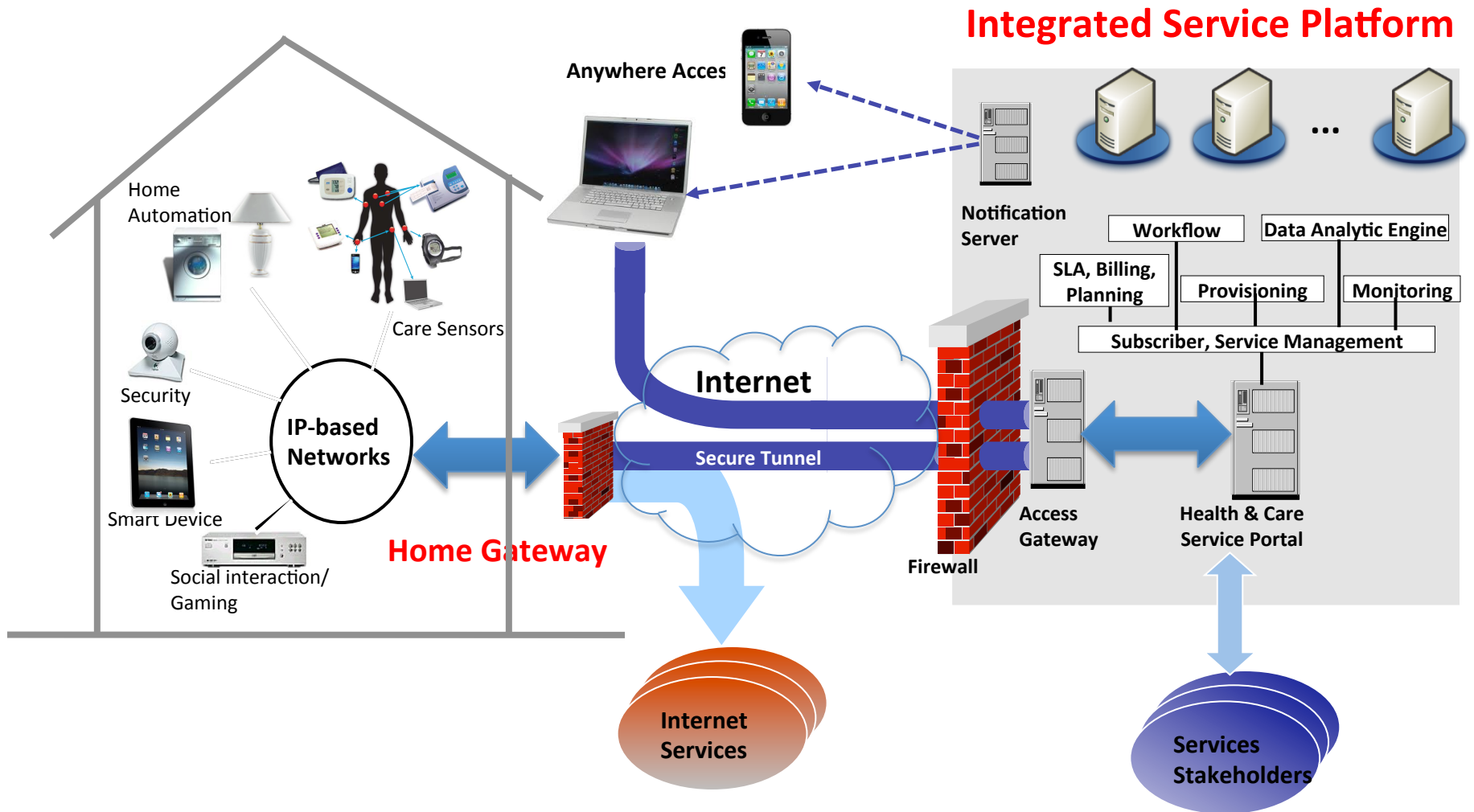
- Framework for evidence collection
- Collect, store and process log data to ensure compliance with policies
- Time correlation of events
- Long-term archiving (min 3-5 years) for legal purposes

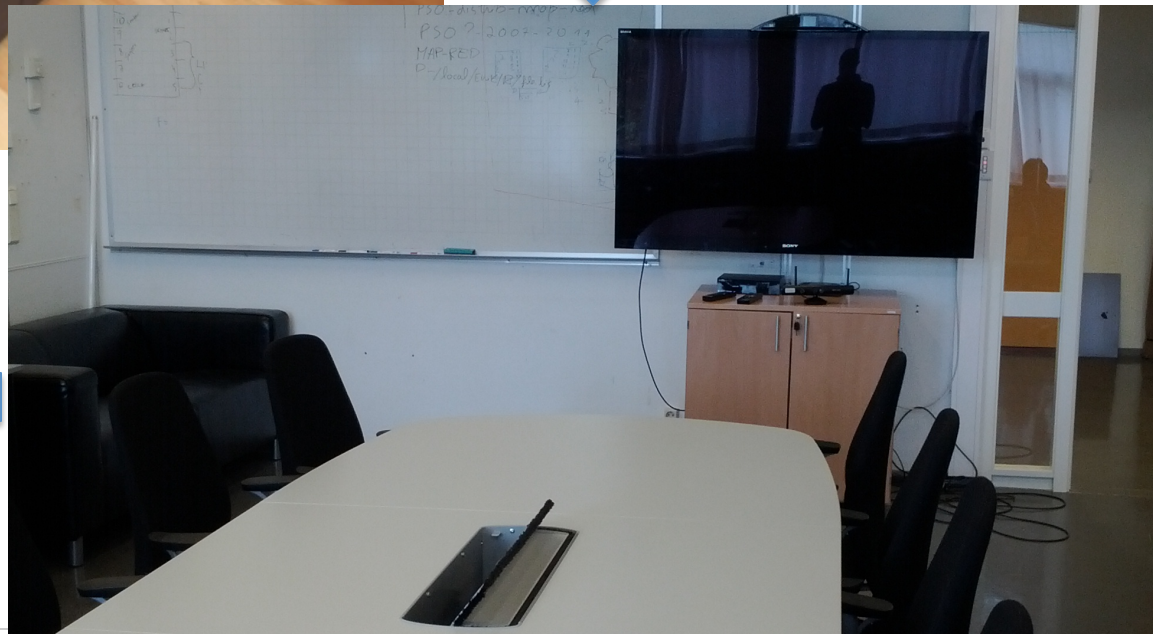
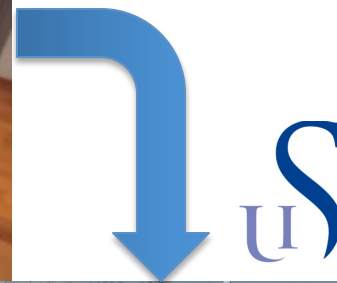
Overview

- CIPSI
- SEEDS
- Uninett
- R2Time
- A4Cloud
- **Safer@Home**
- Data Intensive System course

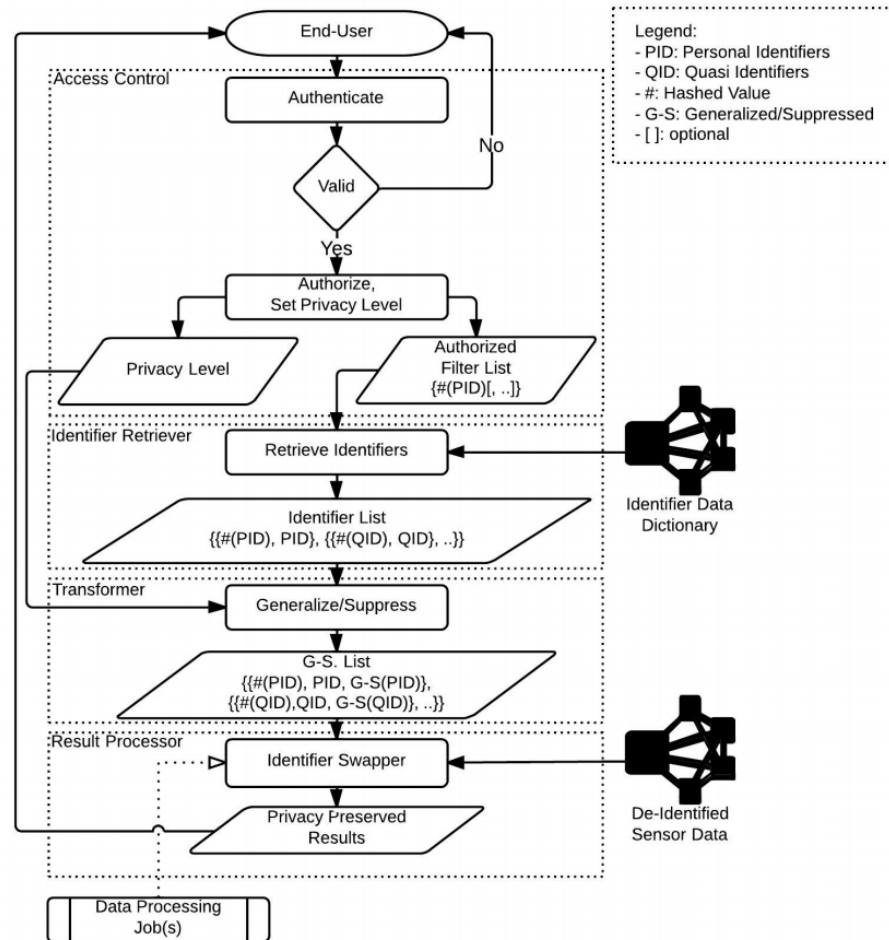
Safer@Home (supported by NFR)

Smart System to Support Safer Independent Living & Social Interaction for Elderly at Home



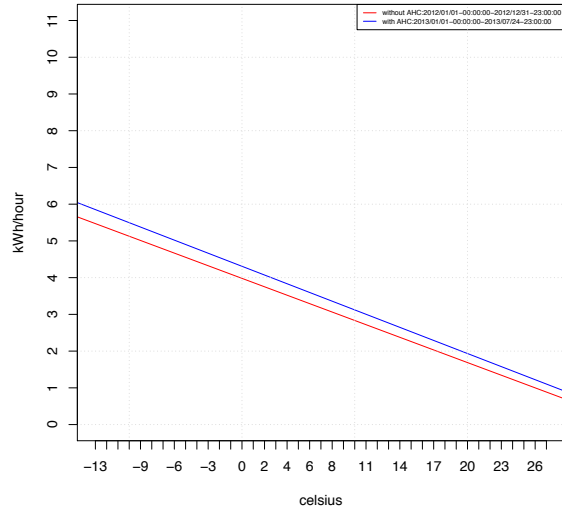


Safer@Home - Privacy

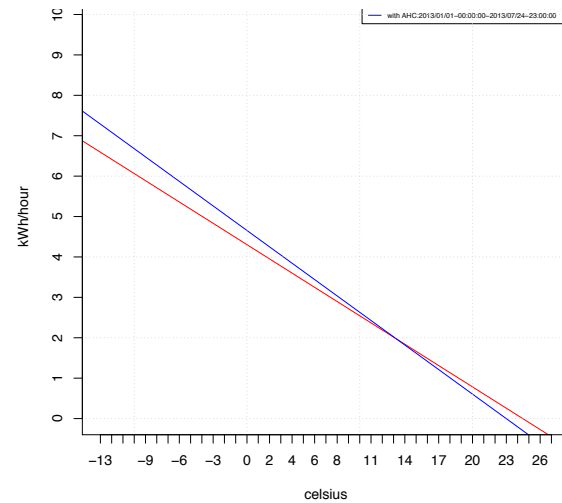
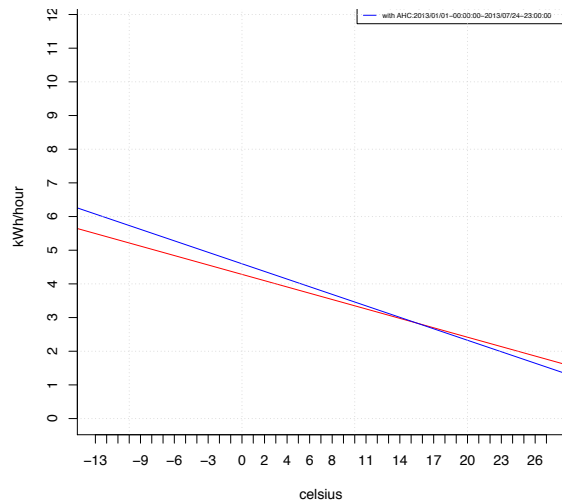
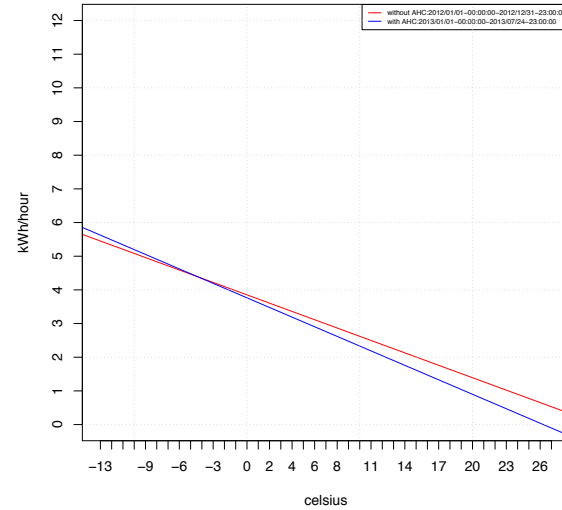


Safer@Home - Energy Saving

Linear Regression for meterid:1034409 and outside temperature



Linear Regression for meterid:1038227 and outside temperature



Safer@Home - Student housing



Overview

- CIPSI
- SEEDS
- Uninett
- R2Time
- A4Cloud
- Safer@Home
- **Data Intensive System course**

Strategic collaboration on Big Data Analytics and Communication



Universitetet
i Stavanger

- SIU sponsored North American Program (2012 – 2016)
 - staff and student exchanges
 - joint curriculum development, teaching and student supervision
 - further funding collaboration

Data Intensive System course

- Available courses are typically based on repackaged content
- Available materials are focused on
 - engineering aspect (most of Hadoop related books)
 - particular algorithmic problems (statistics, machine learning, ...)
- There is lack of
 - coherent and stable learning objectives
 - balance between theory and practice

Data Intensive System course

- 10 ECTS, 1/3 semester load
- Master (early graduate) level
 - Considering later Bachelor/undergrad level
- Co-taught simultaneously at Purdue, UiS and AMD Research
- Created from scratch (avoiding repackaged content)
- Started Fall 2013
- Alternative Short Course (3-5 week long)
 - Focused on summer/winter schools or industrial needs
 - First one given in June 2013
 - ~6-8h a day, lecture + lab
- Collaboration with Research Data Alliance on Common Body of Knowledge within Data Intensive Science

Overview

- Several applied ICT activities revolving around **Big Data**
- Many have significant focus on **time series**
- **R2Time** is a library for time series analysis stored in OpenTSDB over a Hadoop cluster based on extended RHIPE library
- **Data Intensive Systems course** developed jointly with Purdue University and AMD Research
- **CBOK of Data Intensive Science** with Research Data Alliance



<http://2013.cloudcom.org>

• 2009



• 2010



• 2011



• 2012



 **IEEE CloudCom 2013**
Bristol, UK, Dec. 2-5, 2013
<http://2013.cloudcom.org>

