# Energy-Efficient Data Transfers in Radio Astronomy with Software UDP RDMA

Third Workshop on Innovating the Network for Data-Intensive Science, INDIS16

Przemek Lenkiewicz, Researcher@IBM Netherlands
Bernard Metzler, Researcher@IBM Zurich Research Lab
Chris Broekema, Researcher@ASTRON Netherlands Institute for RadioAstronomy

ASTRON
Netherlands Institute for Radio Astronomy
IBM

# Table of contents/Agenda template

**Radio-astronomy & The Square Kilometre Array**

**Data Transport in Radio Astronomy**

**Our Solution and Experiments**

**Conclusions and Future Work**

**The DOME project**

IBM ®
Netherlands

IBM ®
Zürich Research Lab

ASTRON
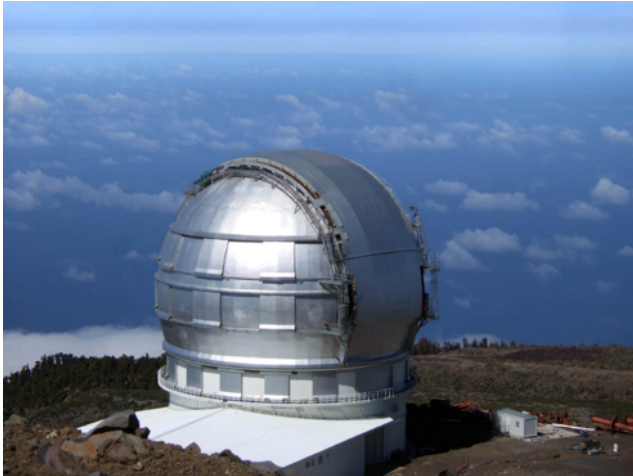Netherlands Institute for Radio Astronomy

13 November 2016

# Radio astronomy &
# The Square Kilometre Array

A brief introduction

# Astronomy
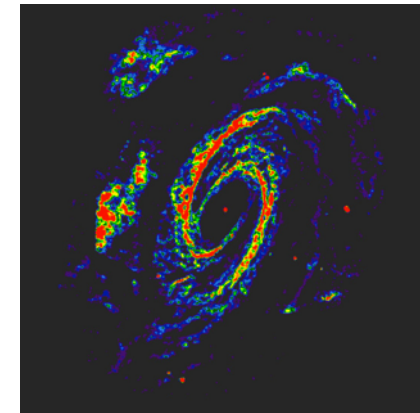
- Lenses, mirrors, sensors
- Light
- Picture of object

- Array of antennas and/or dishes
- Radio frequencies
- Map of radio sources

**Gran Telescopio CANARIAS**



**Low-Frequency Array (LOFAR)**





**The M33 Galaxy**



**The M81 Galaxy**

13 November 2016

# The Square Kilometre Array



SKA1 MID - the SKA's mid-frequency instrument

The Square Kilometre Array (SKA) will be the world's largest radio telescope, revolutionising our understanding of the Universe. The SKA will be built in two phases - SKA1 and SKA2 - starting in 2018, with SKA1 representing a fraction of the full SKA. SKA1 will include two instruments - SKA1 MID and SKA1 LOW - observing the Universe at different frequencies.

Location: South Africa

Frequency range: 350 MHz to 14 GHz

~200 dishes (including 64 MeerKAT dishes)

Total collecting area: 33,000m² or 126 tennis courts

Maximum distance between dishes: 150km

Total raw data output:
2 terabytes per second
62 exabytes per year

x340,000

Enough to fill 340,000 average laptops with content every day

Compared to the JVLA, the current best similar instrument in the world:
4x the resolution
5x more sensitive
60x the survey speed

www.skatelescope.org  Square Kilometre Array  @SKA_telescope  The Square Kilometre Array

SKA1 LOW - the SKA's low-frequency instrument

The Square Kilometre Array (SKA) will be the world's largest radio telescope, revolutionising our understanding of the Universe. The SKA will be built in two phases - SKA1 and SKA2 - starting in 2018, with SKA1 representing a fraction of the full SKA. SKA1 will include two instruments - SKA1 MID and SKA1 LOW - observing the Universe at different frequencies.

Location: Australia

Frequency range: 50 MHz to 350 MHz

~130,000 antennas spread between 500 stations

Total collecting area: 0.4km²

Maximum distance between stations: 65km

Total raw data output:
157 terabytes per second
4.9 zettabytes per year

Enough to fill up 35,000 DVDs every second

5x the estimated global internet traffic in 2015 (source: Cisco)

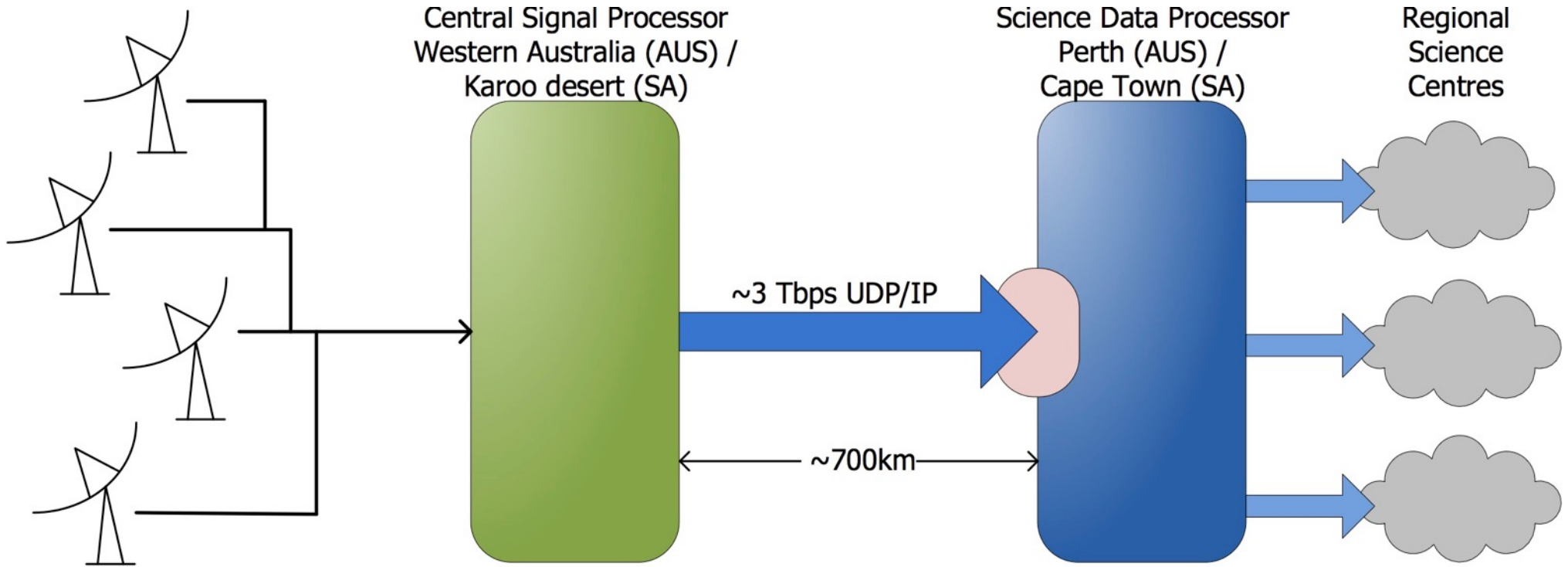Compared to LOFAR Netherlands, the current best similar instrument in the world
25% better resolution
8x more sensitive
135x the survey speed

www.skatelescope.org  Square Kilometre Array  @SKA_telescope  The Square Kilometre Array

# Radio astronomy data transport

# SKA telescope data flow



Central Signal Processor
Western Australia (AUS) /
Karoo desert (SA)

Science Data Processor
Perth (AUS) /
Cape Town (SA)

Regional
Science
Centres

~3 Tbps UDP/IP

~700km

13 November 2016

# SKA Phase 1 in numbers
## (italics are derived and/or speculative)

| | SKA1 MID | SKA1 LOW |
|---|---|---|
| Location | Karoo, South Africa | Western Australia |
| Number of receivers | 197 (133 SKA + 64 MeerKAT) | 131.072 (512 st x 256 el) |
| Receiver diameter | 15 m (13,5 m MeerKAT) | 35 m (station) |
| Maximum baseline | 150 km | 65 km |
| Frequency channels | 65.536 | 65.536 |
| SDP input bandwidth | 3,1 Tbps | 3,1 Tbps |
| *Req'd Compute capacity[*]* | *20-72 PFLOPS* | *16-41,5 PFLOPS* |
| *Archive growth rate* | *10 – 100 Gbps (50yr life)* | *25 – 100 Gbps (50yr life)* |
| *SDP Energy budget* | *<5MW* | *<5MW* |

[*]These are sustained PFLOPS, computational (in)efficiency not included
Cost cap for the first phase of SKA: € 650 M (2014)
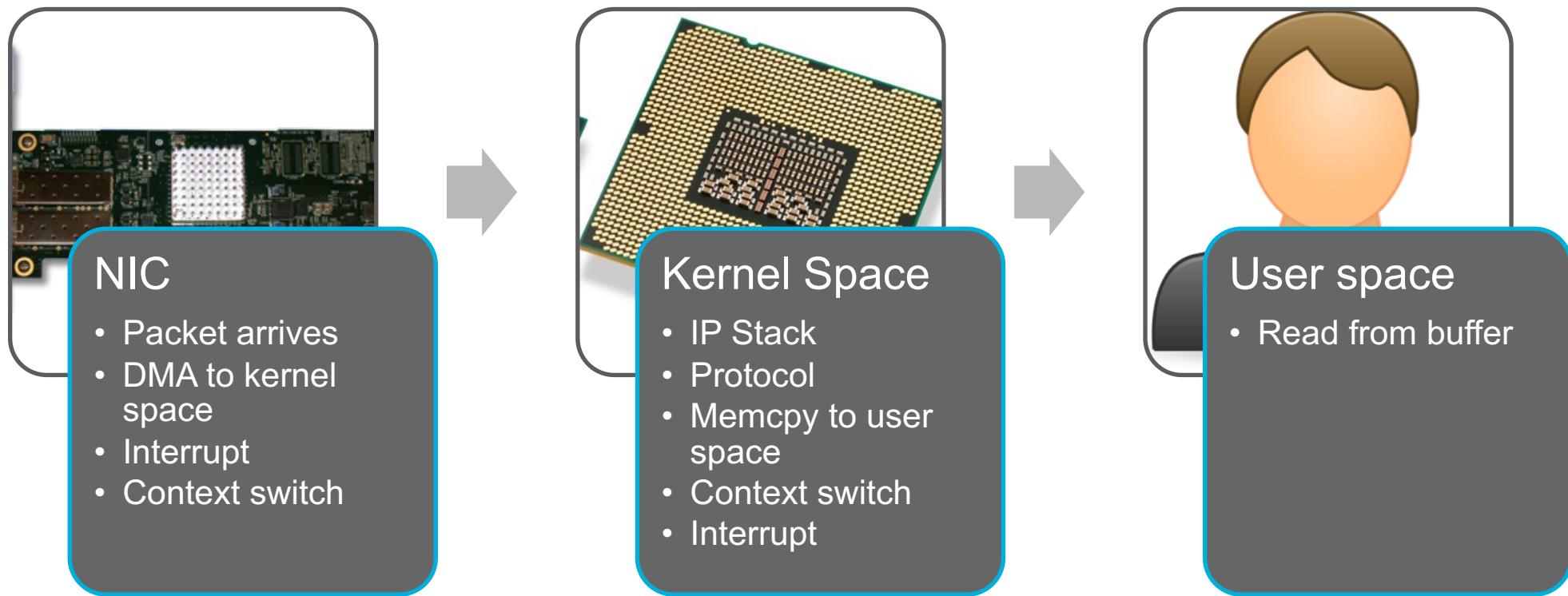
13 November 2016

# Compare: Top 500 development

## Projected Performance Development



- Efficiency of radio astronomy algorithms: ~10%

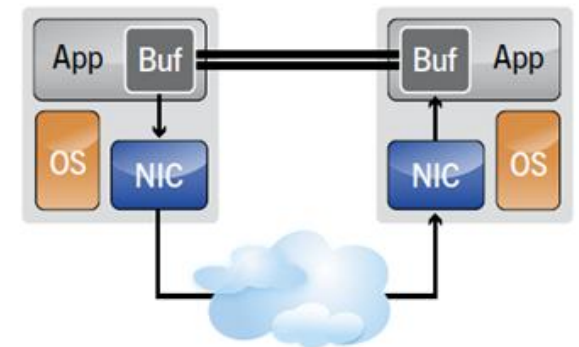- Effective required compute capacity: hundreds of PFLOPS

Souce: top 500 lists

# Receive data through hierarchical OS

**NIC**
- Packet arrives
- DMA to kernel space
- Interrupt
- Context switch

**Kernel Space**
- IP Stack
- Protocol
- Memcpy to user space
- Context switch
- Interrupt

**User space**
- Read from buffer

      13 November 2016

# Requirements for astronomical data transport service

- Very high data rates – Terabits per second per instrument

- Almost entirely uni-directional traffic

- UDP/IP over Ethernet

- Prioritizing bandwidth over latency

- Desire for very high energy efficiency
  - Receiving end crucial!

- Full reliability is not crucial, some data loss is tolerable

        13 November 2016

# Approach - RDMA

- Moving data from user space memory of one machine to that of another
- No involvement of host operating system
- Memory buffers registered with the local RDMA-capable network adapter (RNIC) and usually pinned to local physical memory
- Fully asynchronous to allow overlapping communication and computation

13 November 2016

# RDMA in radioastronomy

- We looked at:
  - RoCE
  - iWARP
  - Infiniband

- Reliable connection only

- Short range

     13 November 2016

# Approach – SoftiWARP UDP

- SofiWARP (SIW)

- An open source software implementation of the iWARP protocol suite

- Developed at the IBM Zurich Research Lab and available from GitHub
  – https://github.com/zrlio/softiwarp

- Exports the OpenFabrics RDMA API to both user space and kernel space applications

- Fully compatible with hardware iWARP RNICs

- Utilizes kernel sockets for efficient communication and less data touching

- In DOME:
  – Development of UDP transport layer

SoftiWARP

      13 November 2016

# Experiments and measurements

# Power sensor for PCIe card slot

- ARDUINO board
- Current flow sensors
- PCIe riser card



          13 November 2016

# Chelsio T5



Chelsio T5 Power consumption
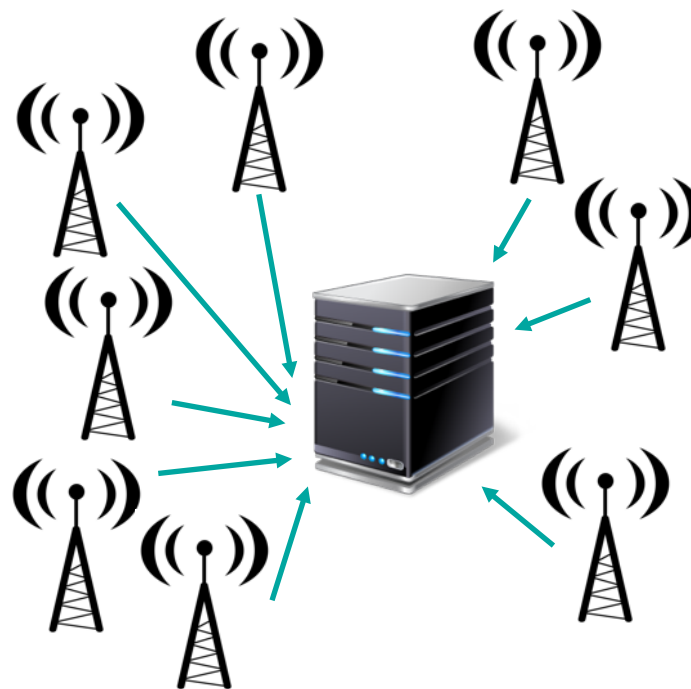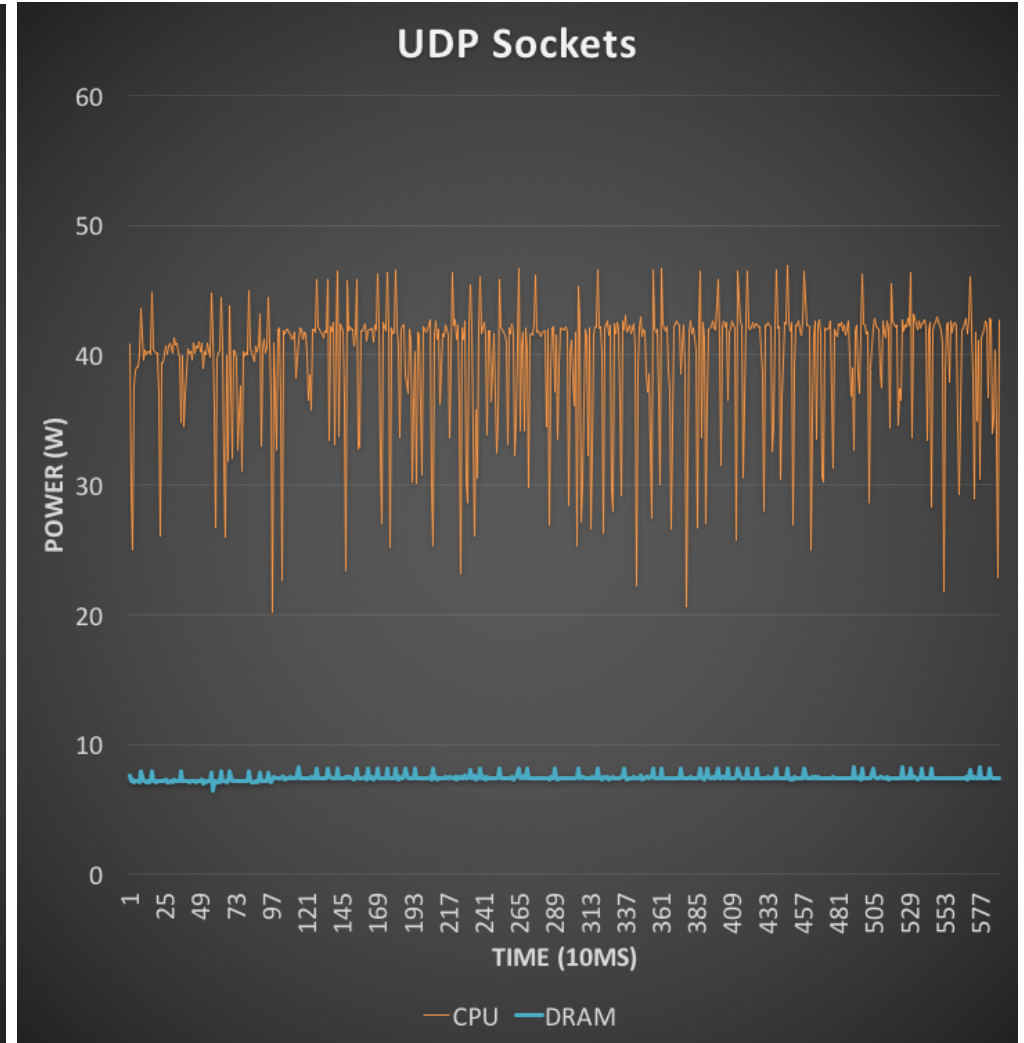
13 November 2016

# Power consumption measurements
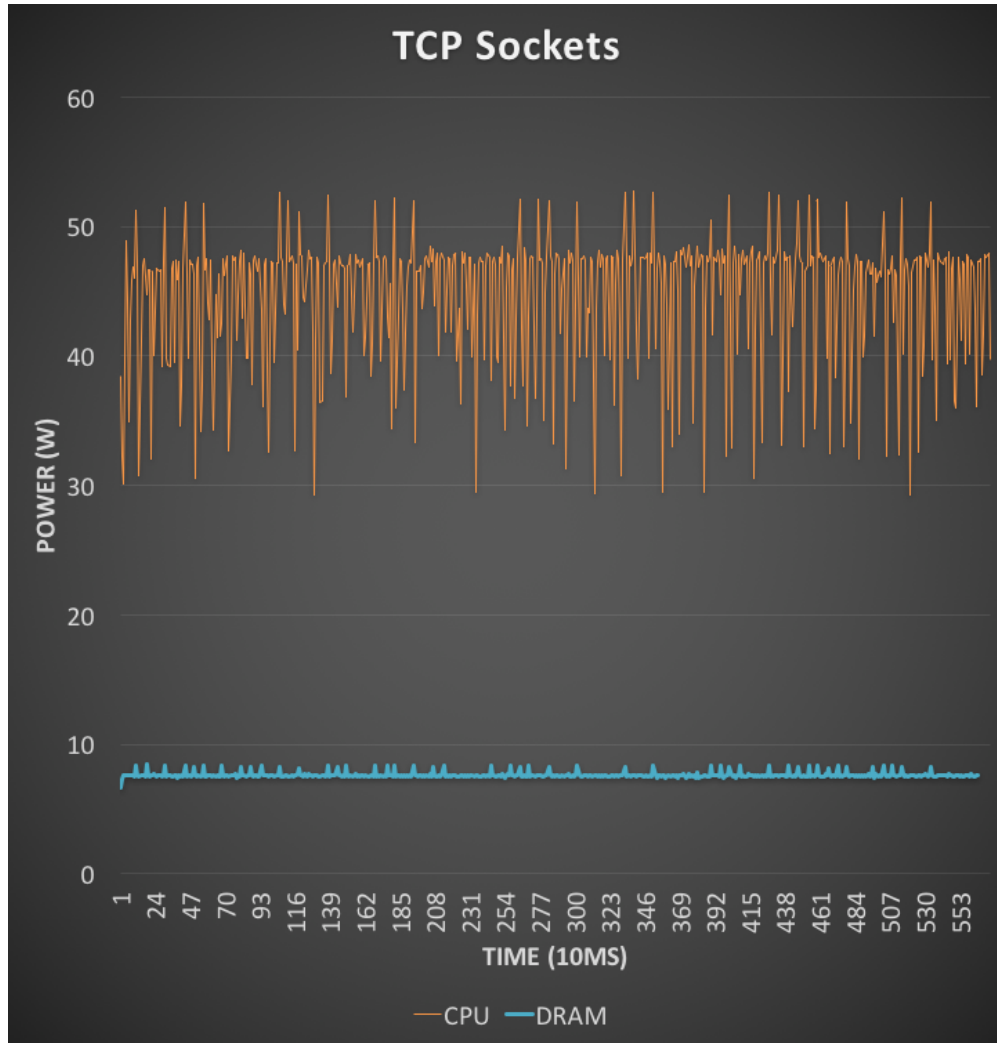
- CPU and DRAM power consumption

- RAPL for power readings

- Radio-astronomy traffic simulation

- Netperf benchmarking tool

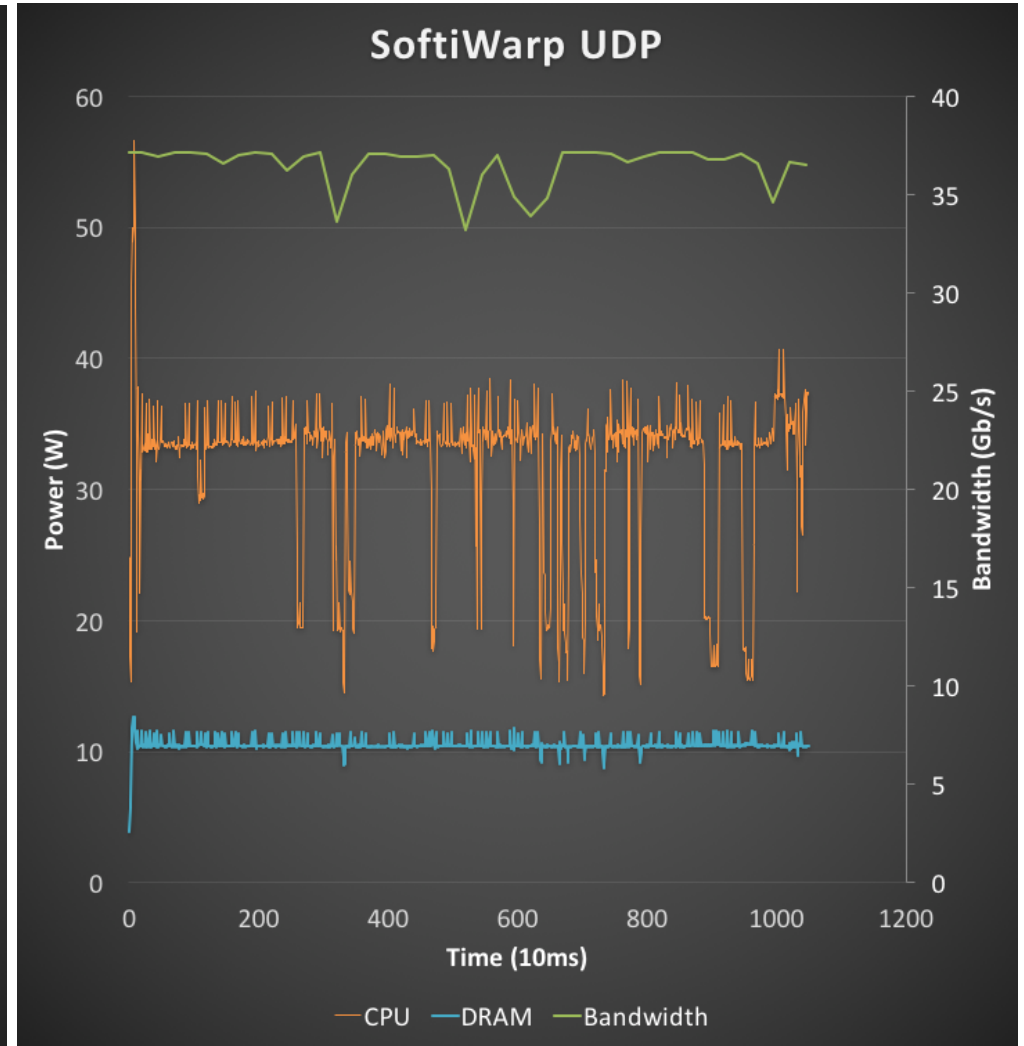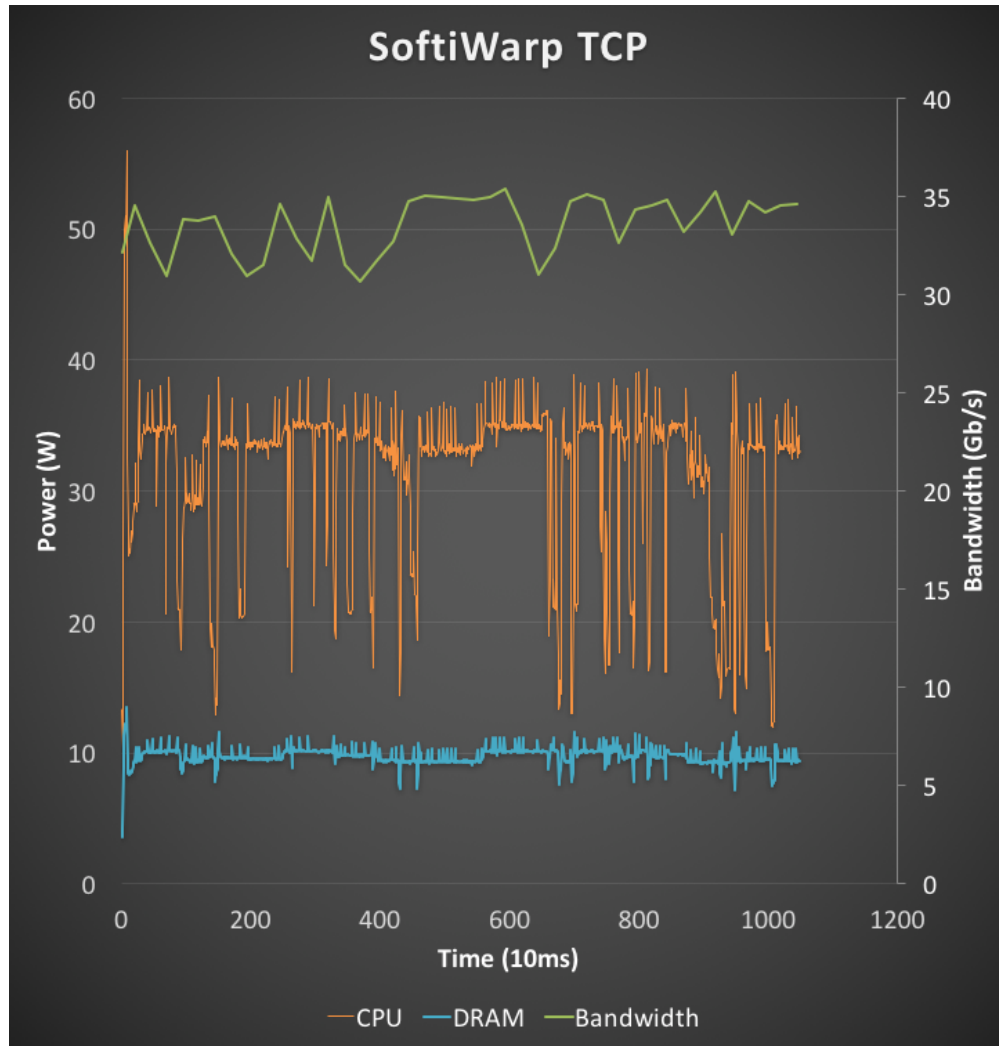          13 November 2016

# Radio-astronomy data flow

- Mimic the data flow from LOFAR

- Emulate the data produced by a LOFAR Remote Station Processing (RSP) board.

- UDP/IP data stream, ~760 Mb/s, packets of 8 kB.

- 50 Data streams received by a single CPU core

- Our emulator supports:
  - TCP Sockets
  - UDP Sockets
  - Softiwarp TCP
  - Softiwarp UDP

# Radio-astronomy data flow



 13 November 2016

# Radio-astronomy data flow



          13 November 2016

# Power measurements, no offloading

| | | | Power (W) | | | |
|---|---|---|---|---|---|---|
| | | BW (Gb/s) | Send cpu | Send dram | Recv cpu | Recv dram |
| TCP | Sock: | 38.17 | 24.13 | 7.23 | 24.46 | 6.99 |
| UDP | Sock: | 39.55 | 24.14 | 9.59 | 23.35 | 7.09 |
| SIWTCP | RDMA Read | 35.66 | 25.6 | 4.8 | 24.64 | 5.36 |
| SIWTCP | RDMA Write | 23.83 | 22.65 | 6.18 | 11.48 | 5.53 |
| SIWUDP | RDMA Read | 39.17 | 21.55 | 6.64 | 21.7 | 5.3 |
| SIWUDP | RDMA Write | 38.33 | 26.58 | 8.03 | 13.34 | 6.98 |

     13 November 2016

# Power efficiency comparison
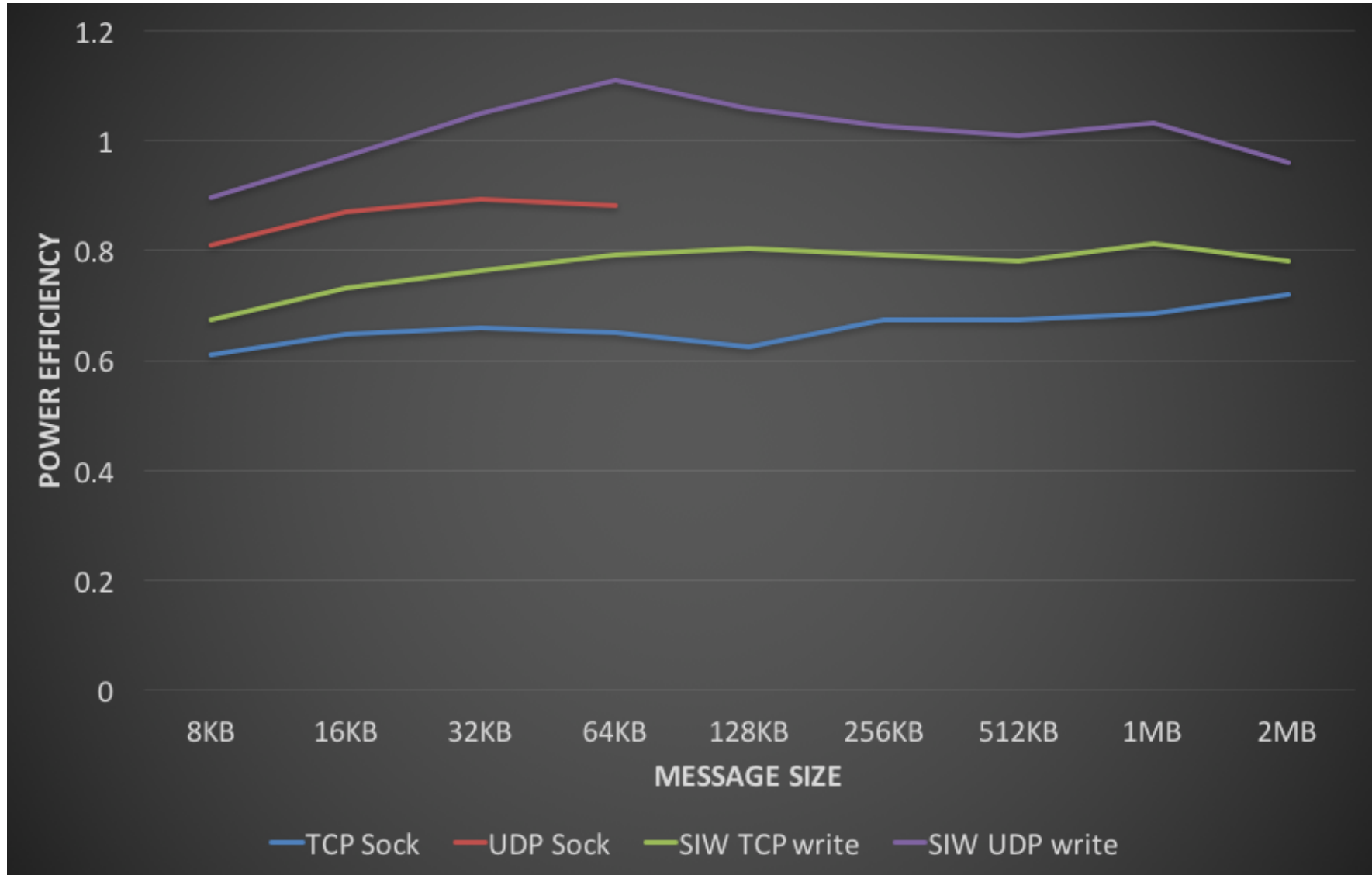
- We perform tests of all protocols and vary the message size
- Our power efficiency metric:

$$\frac{Bandwidth(Gb/s)}{PowerCons(Watt)}$$

13 November 2016

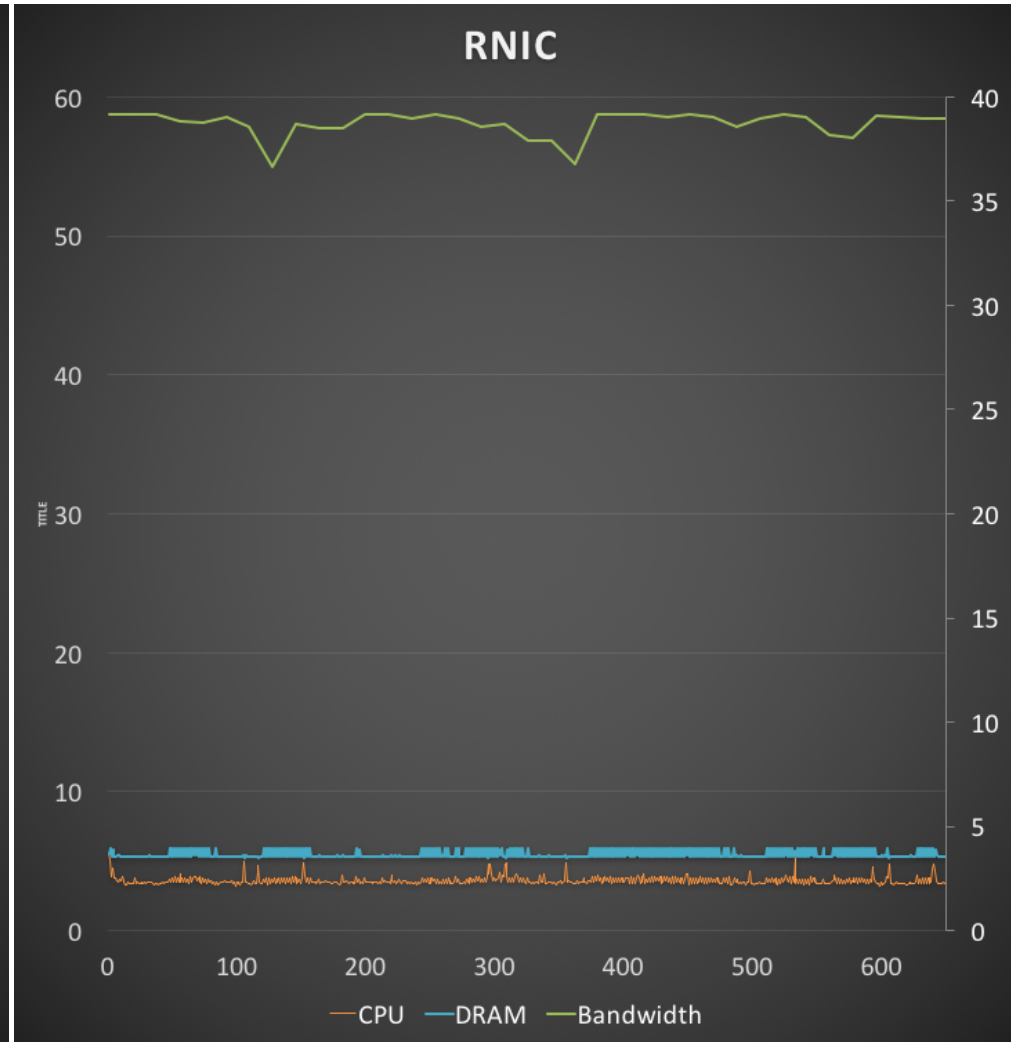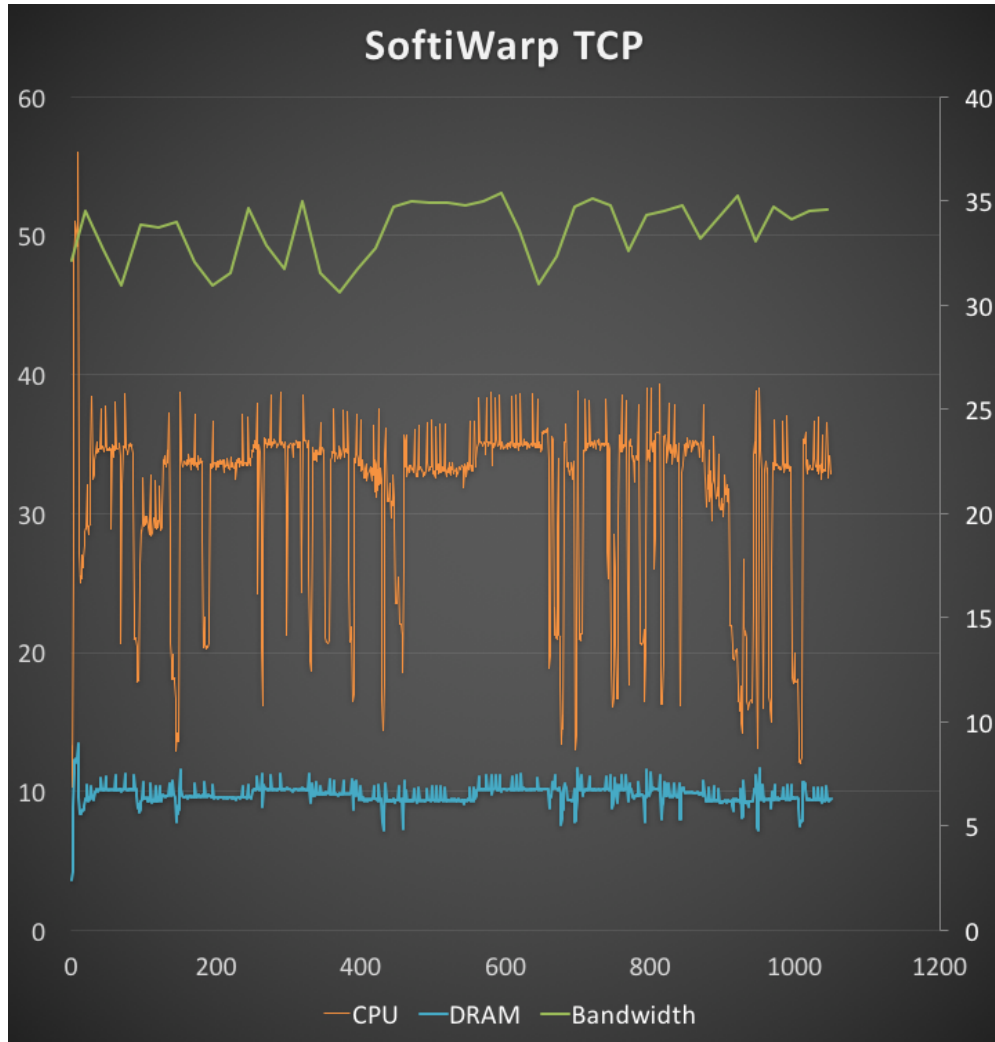# Bandwidth and Power Consumption



     13 November 2016

# Calculated power efficiency (normalized power consumption)

# Radio-astronomy data flow

13 November 2016

# Performance in case of packet loss

# Bandwidth versus packet loss

Introduced packet loss (%)

BANDWIDTH GB/s

—TCP Sockets —SIW TCP read —SIW UDP read —SIW TCP write —SIW UDP write —UDP Sockets

     13 November 2016

# Conclusions and future work

- An unreliable datagram-based iWARP protocol suits the requirements for radio-astronomy data transport service

- Software prototype looks promising

- Best power efficiency with hardware product (FPGA-based the most viable solution)

- Further work directions
  - Investigate the use of flash storage and big data frameworks (Spark, Hadoop)

     13 November 2016

# Thank you!

[http://www.dome-exascale.nl](http://www.dome-exascale.nl)

[lenkiewicz@nl.ibm.com](mailto:lenkiewicz@nl.ibm.com)

     13 November 2016