

# Software and Computing Domain: WP2.6.2 Computing Hardware Architecture

**Duncan Hall, Chris Broekema**  
**WP2 Meeting**  
**2010 October 29**



- WP2.6.2 Hardware options for SKA computing
- Calibration and imaging requirements for SKA1
- Forecasts of COTS hardware capabilities
- Forecasts of COTS power requirements
- Data input – output challenges

## WP2.6.2 Hardware options for SKA computing

- Calibration and imaging requirements for SKA1
- Forecasts of COTS hardware capabilities
- Forecasts of COTS power requirements
- Data input – output challenges

Contributions to date; Current Status	Challenges to be addressed
<ul style="list-style-type: none"> <li>•A spreadsheet model identifying key cost driver parameters and likely pessimistic-optimistic ranges has been published on the S&amp;C Domain wiki</li> </ul>	<ul style="list-style-type: none"> <li>•The required number of floating point operations per u-v sample (“flops per float”) will likely range from 100,000 to 400,000+</li> <li>•This leads to requirements for hundreds of petaflop/s to ~1 exaflop/s for SKA1 as defined in Memo 125</li> <li>•The “flops” metric is only one measure of HPC performance; disk input-output data rates, cache memory access speeds and sizes, and hardware reliability can be just as important</li> </ul>

# CPG Memo 3 (2009-11-6) confirms requirements for extreme scale computing:

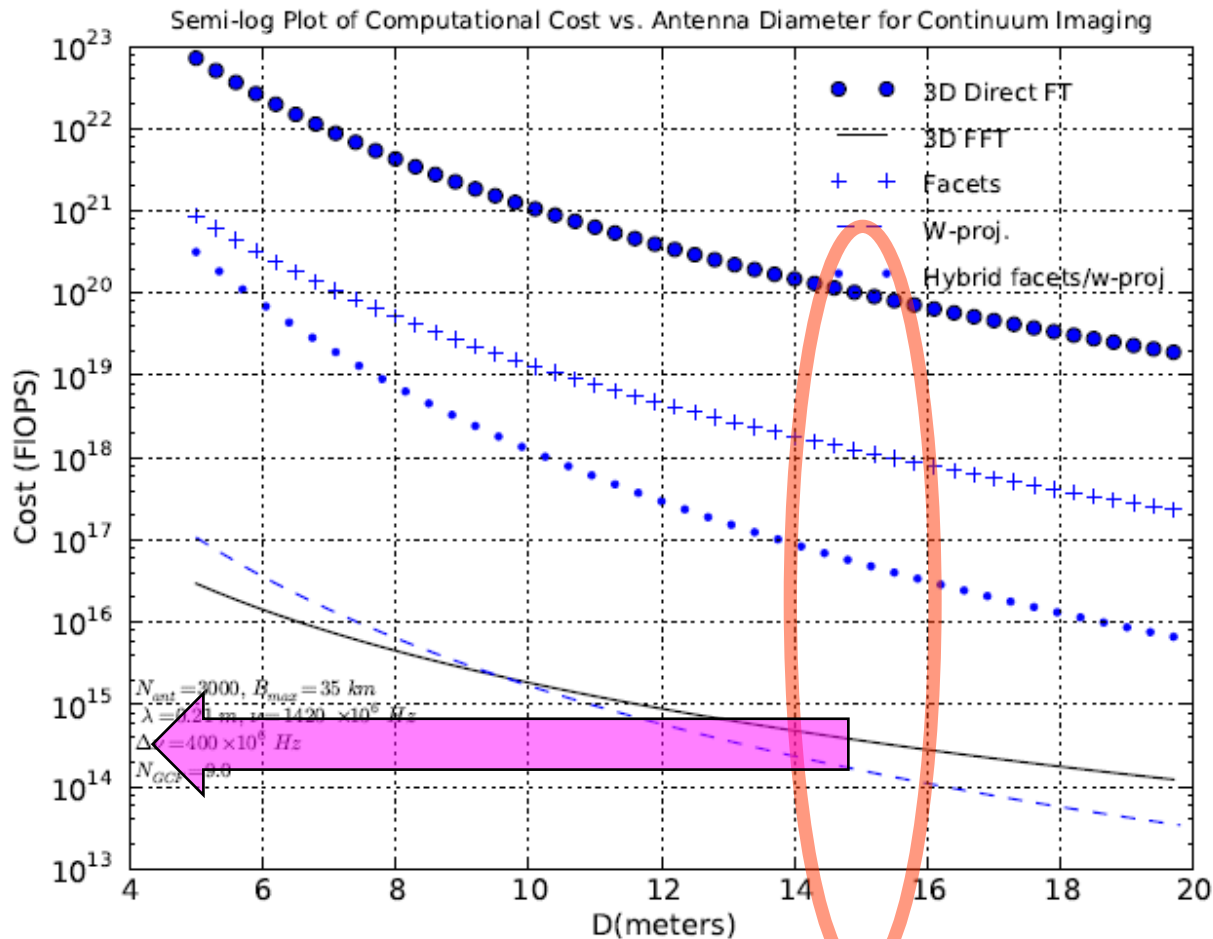
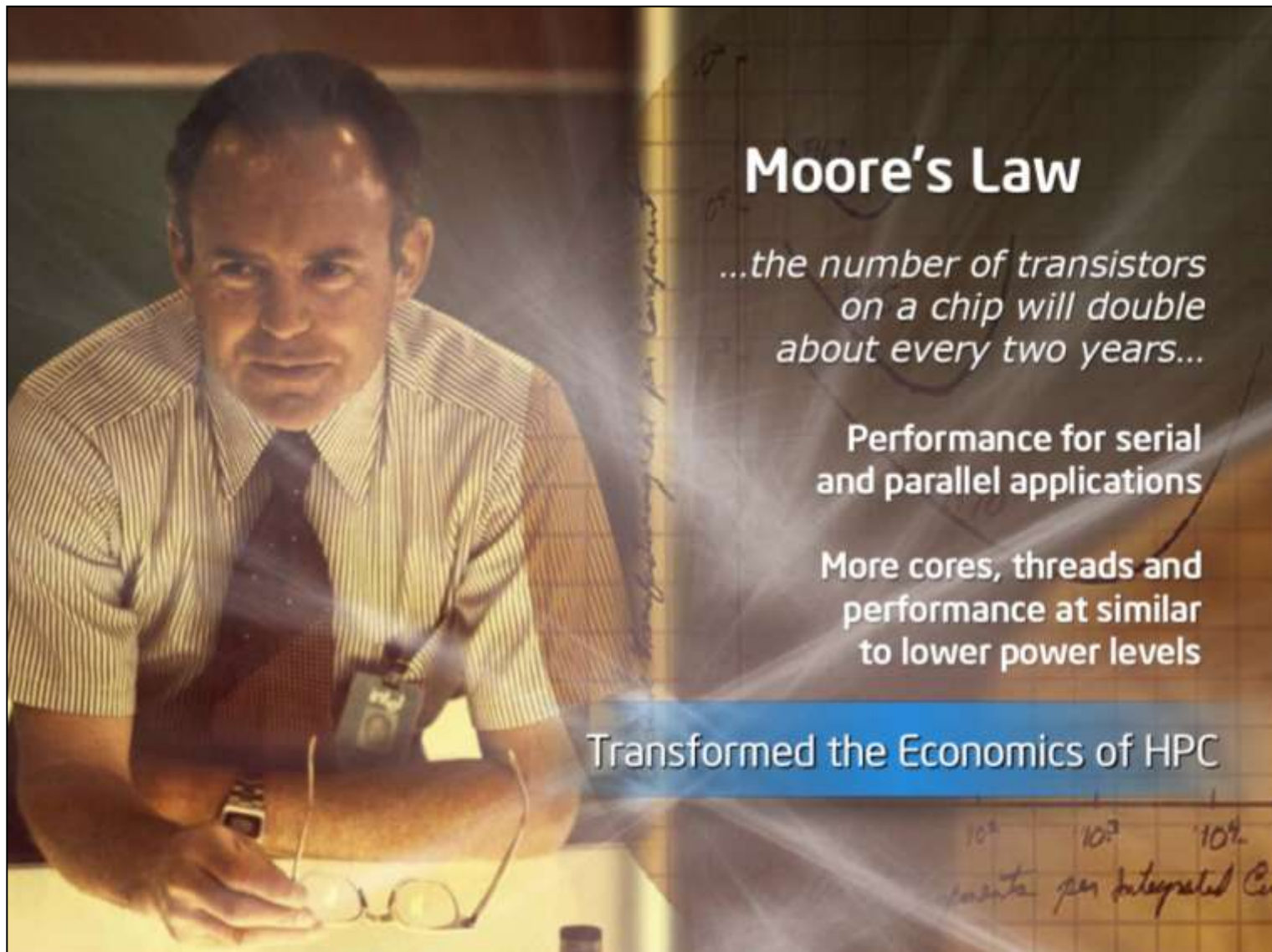


Figure 1: Semi-log y plots of computational costs (without consideration of deconvolution and parallel computing efficiency  $\eta$ ) vs. antenna diameter  $D$  for continuum imaging for the 3-D direct FT, 3-D FFT, facets, w-projection, and hybrid facets/w projection imaging algorithms.

# Intel asserts that Moore's 1965 observation will continue to hold to ~2020 ...



## Moore's Law

*...the number of transistors  
on a chip will double  
about every two years...*

**Performance for serial  
and parallel applications**

**More cores, threads and  
performance at similar  
to lower power levels**

**Transformed the Economics of HPC**

*10<sup>2</sup> 10<sup>3</sup> 10<sup>4</sup>  
transistors per Integrated Cir*

... but the advent of multiple cores per chip has major implications for software at extreme scale ...

# Cores Per Die

Past



Present



Future



## Amdahl's Laws

Gene Amdahl (1965): **Laws for a balanced system**

- i. Parallelism: max speedup is  $S/(S+P)$
- ii. **One bit of IO/sec per instruction/sec (BW)**
- iii. One byte of memory per one instruction/sec (MEM)

Modern multi-core systems move farther  
away from Amdahl's Laws  
(Bell, Gray and Szalay 2006)





Challenges to be addressed	Work to be done; Milestones; Risks
<ul style="list-style-type: none"> <li>•The required number of floating point operations per u-v sample (“flops per float”) will likely range from 100,000 to 400,000+</li> <li>•This leads to requirements for hundreds of petaflop/s to ~1 exaflop/s for SKA1 as defined in Memo 125</li> <li>•The “flops” metric is only one measure of HPC performance; disk input-output data rates, cache memory access speeds and sizes, and hardware reliability can be just as important</li> </ul>	<ul style="list-style-type: none"> <li>•Opportunities to improve performances over that of current codes must be explored – in WP2.6.3 work</li> <li>•Opportunities to shift computing load from general purpose von Neumann architecture to special purpose architectures – using for example hardware accelerators – must be explored as a follow on from WP2.6.3 work</li> </ul>

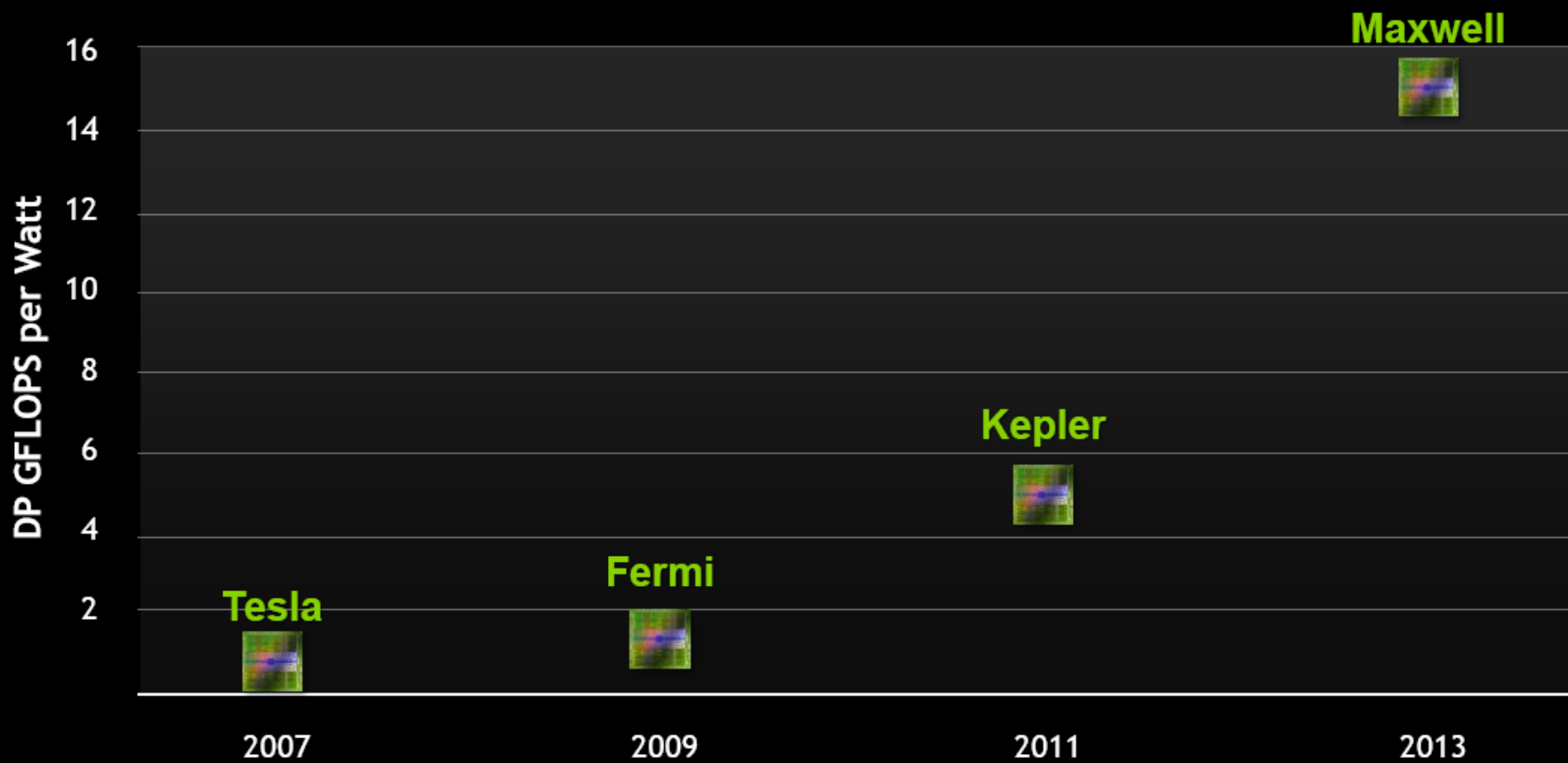
- WP2.6.2 Hardware options for SKA computing
  - Calibration and imaging requirements for SKA1
  - Forecasts of COTS hardware capabilities
  - Forecasts of COTS power requirements
  - Data input – output challenges



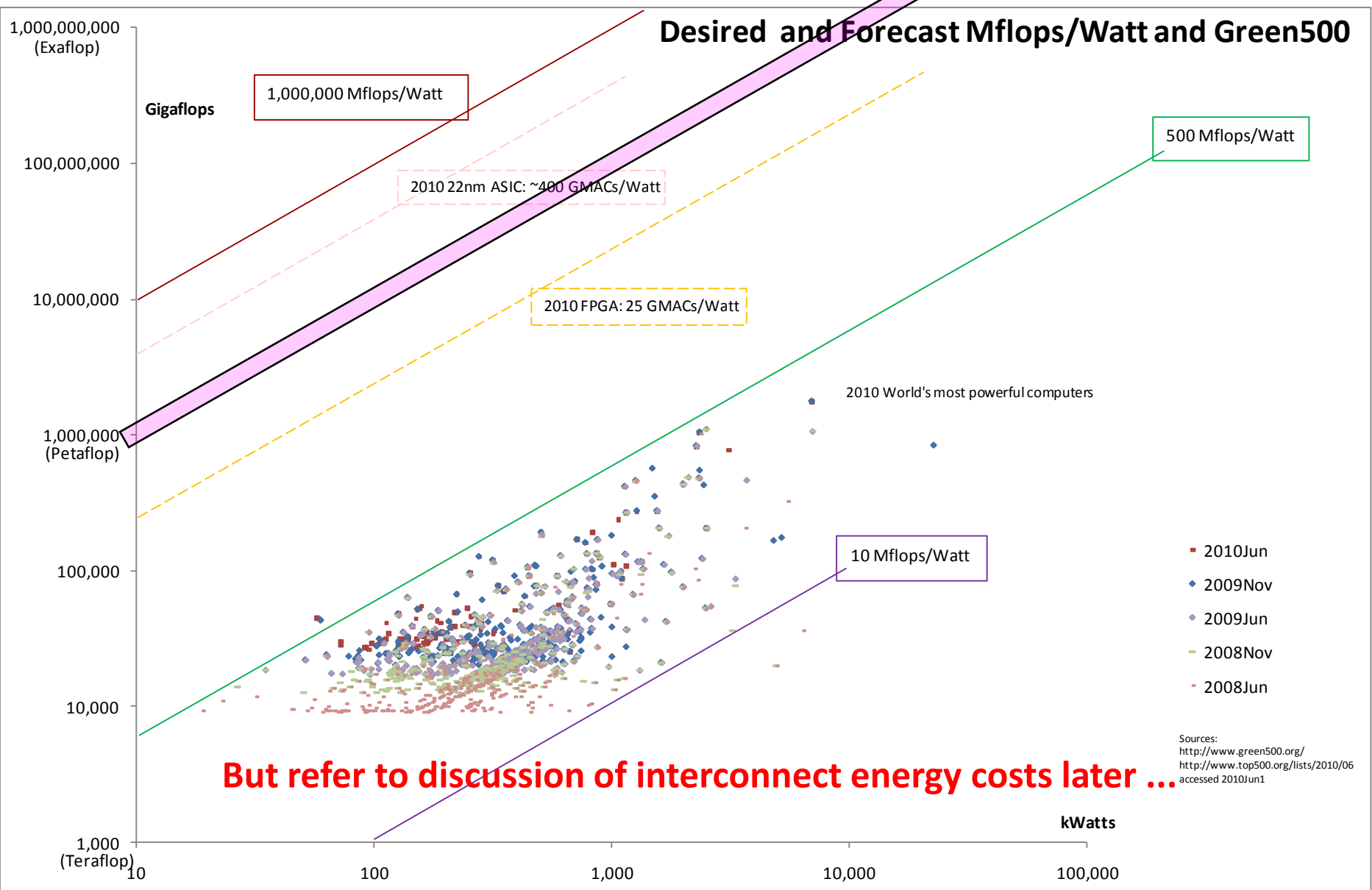
## WP2.6.2 Forecasts of COTS hardware capabilities and power requirements

Contributions to date; Current Status	Challenges to be addressed
<ul style="list-style-type: none"><li>•Recent published research on Commercial Off The Shelf (COTS) hardware capabilities and power requirements has been posted on the S&amp;C Domain wiki</li><li>•Forecasts of performance from vendors of compute elements (e.g. Intel and NVIDIA) have been posted on the S&amp;C Domain wiki</li></ul>	<ul style="list-style-type: none"><li>•Other than essentially one-off “icon” installations, COTS High Performance Computers (HPCs) with capacities of ~1 exaflop/s are not likely to be commercially available until ~2020</li><li>•The typical cost of high end HPC hardware is ~€100 million; other infrastructure such as persistent storage is additional</li><li>•The stretch target power consumption for exaflop/s class HPCs given to US vendors is ~20 MW for the HPC alone</li></ul>

# CUDA GPU Roadmap



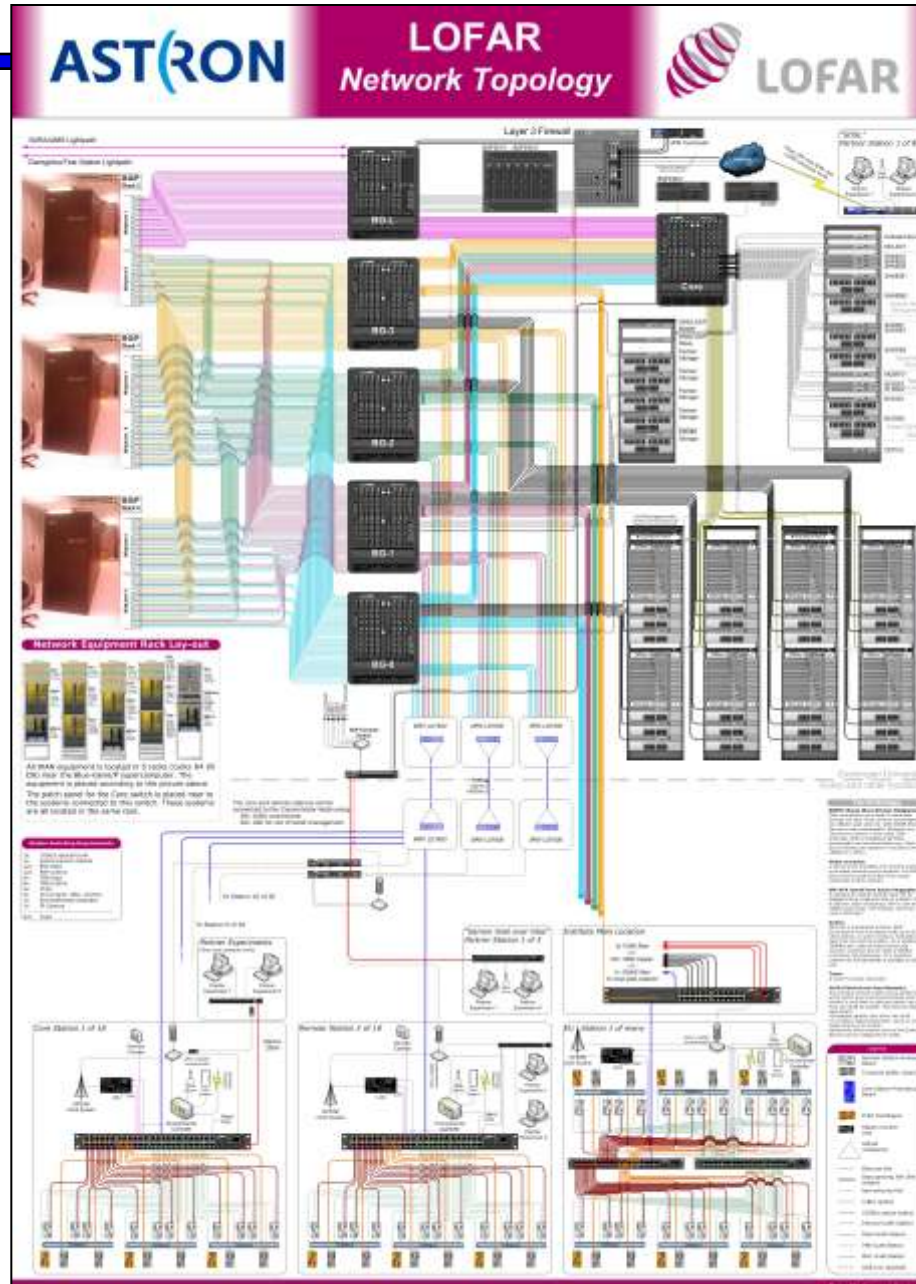
It is not impossible that the equivalent of GPU-powered boards operating at 100+ SP Gflop/s per watt could deliver 1 EF for ~10 MW



**But refer to discussion of interconnect energy costs later ...**

Challenges to be addressed	Work to be done; Milestones; Risks
<ul style="list-style-type: none"> <li>• Other than essentially one-off “icon” installations, COTS High Performance Computers (HPCs) with capacities of ~1 exaflop/s are not likely to be commercially available until ~2020</li> <li>• The typical cost of high end HPC hardware is ~€100 million; other infrastructure such as persistent storage is additional</li> <li>• The stretch target power consumption for exaflop/s class HPCs given to US vendors is ~20 MW for the HPC alone</li> </ul>	<ul style="list-style-type: none"> <li>• Characterisation of all sources of power consumption – including interconnection, power conditioning, persistent storage and environmental control – is required in order to estimate the total power costs</li> <li>• Significant work is required to assess the feasibility of HPC using heterogeneous architectures, e.g. learn from LOFAR, ASKAP and Single Digital Backend</li> <li>• Interim results of work to be delivered before CoDR in September 2011</li> <li>• Little downside risk: work is research</li> </ul>

- 
- WP2.6.2 Hardware options for SKA computing
    - Calibration and imaging requirements for SKA1
    - Forecasts of COTS hardware capabilities
    - Forecasts of COTS power requirements
    - Data input – output challenges



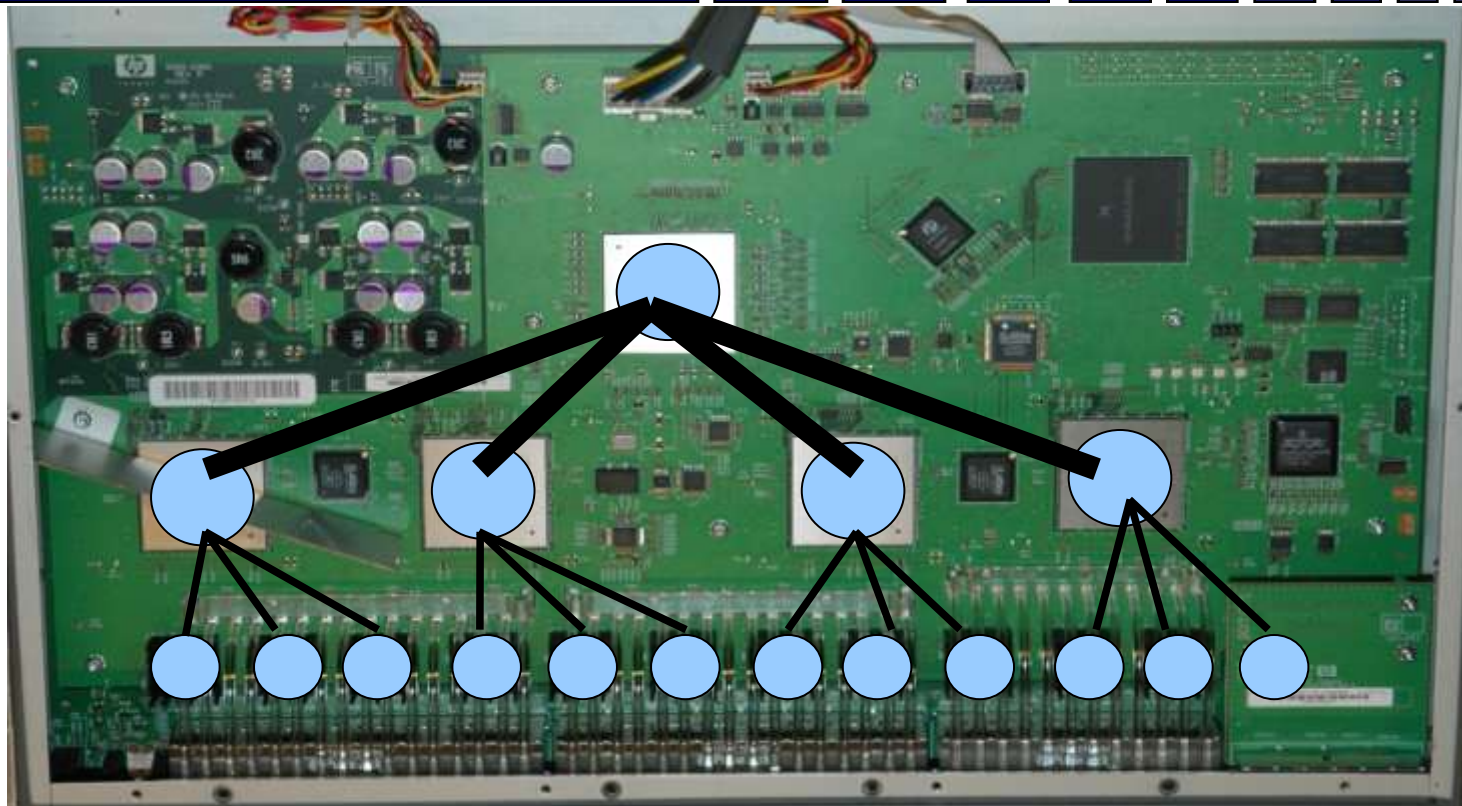


Contributions to date; Current Status	Challenges to be addressed
<ul style="list-style-type: none"> <li>•ASTRON experience: for LOFAR data streaming, the default IBM Blue Gene P input I/O configuration was more than 50% too slow</li> <li>•Various optimisations were required to improve performance</li> <li>•Need to consider I/O on all levels:               <ul style="list-style-type: none"> <li>•System design, operating system, communications stacks, application</li> </ul> </li> <li>•Current software not optimised for throughput:               <ul style="list-style-type: none"> <li>•Mostly optimised for connections and stability</li> <li>•Much of the required work needs to be done in-house</li> <li>•At least some expert knowledge is required, e.g. Linux kernels on I/O nodes</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>•Cost of I/O is non-linear with scale:               <ul style="list-style-type: none"> <li>•<u>Number of active switch chips</u>: scales faster than port count because of fat-tree configuration; doesn't scale beyond a few thousand ports due to limited port count per component with much worse scaling from there on</li> <li>•<u>Procurement cost</u>: added costs of interconnect components; additional complexity; and high bandwidth COTS cards are just not available beyond a few thousand ports</li> <li>•<u>Energy</u>: line energy cost is ~linear with distance, but need to add costs of noise handling algorithms and interconnect active components</li> </ul> </li> </ul>



A 48 port GbE switch, representative of a hardware-efficient universal fat-tree network design:

- Loads of chips needed already
- Non-blocking performance not guaranteed



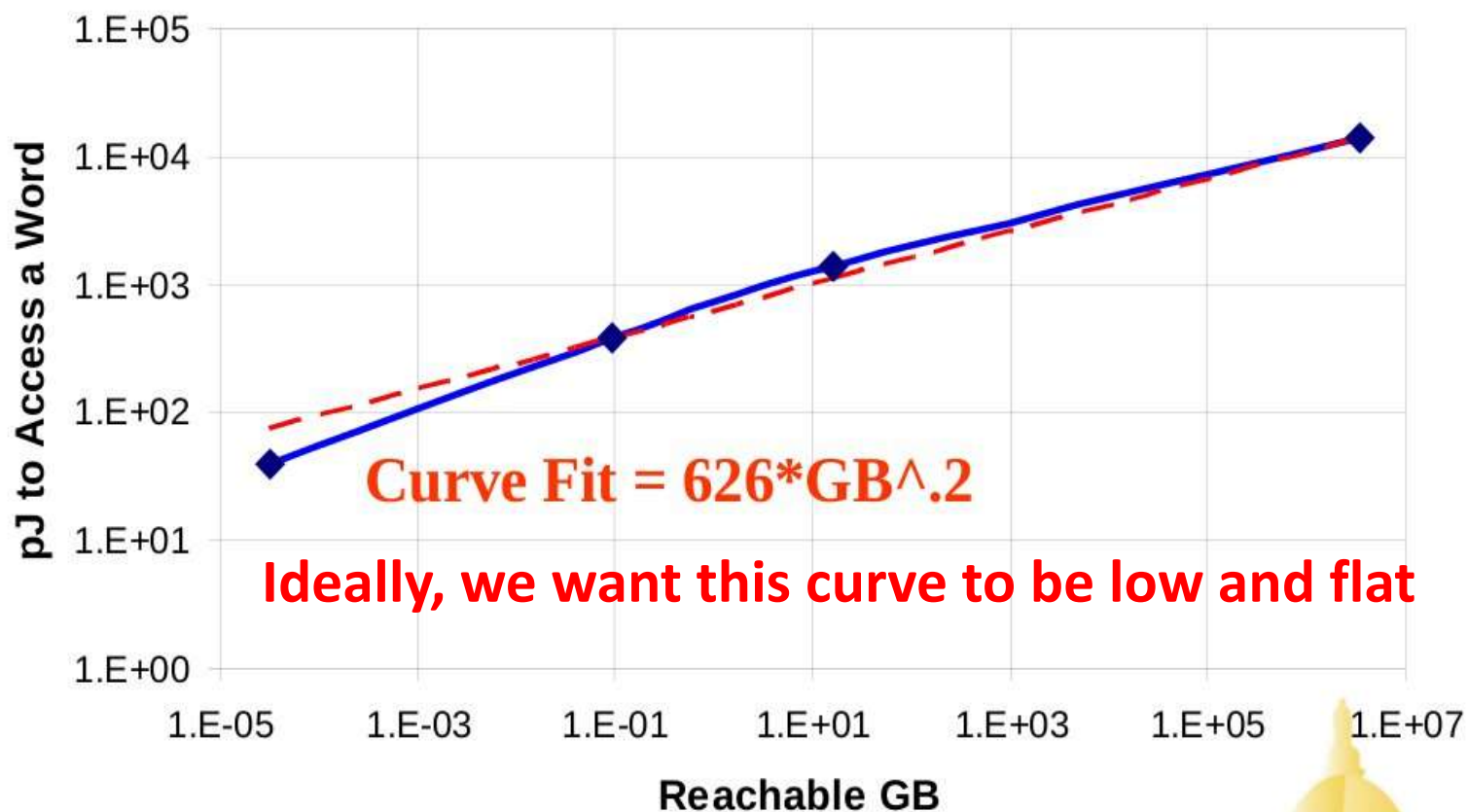
**This design will not scale indefinitely.**

**For N ports per switch chip:**

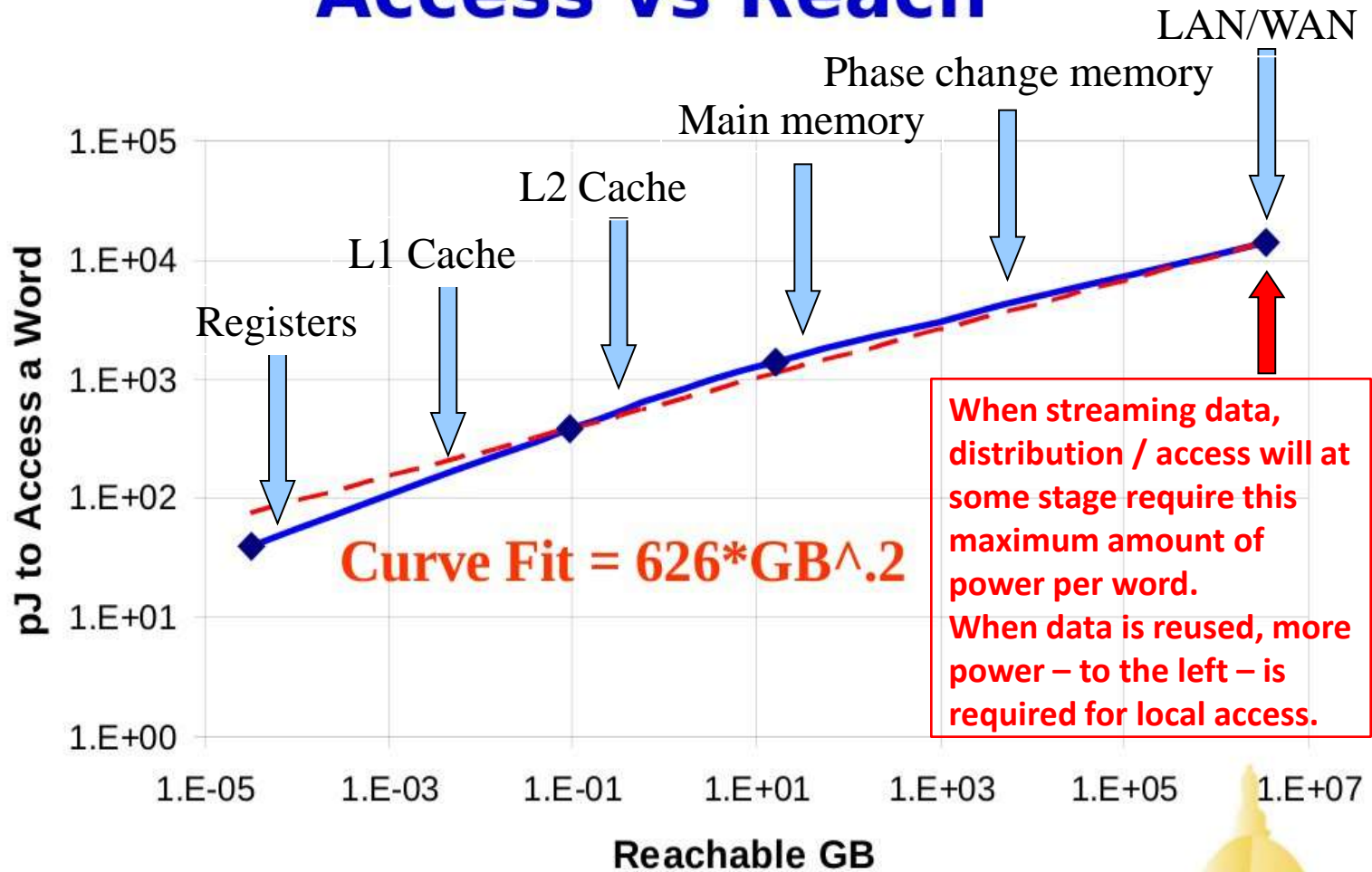
**Max port count =  $N^2 / 2$**

**when using a single type of switch chip; this uses three types**

# Access vs Reach

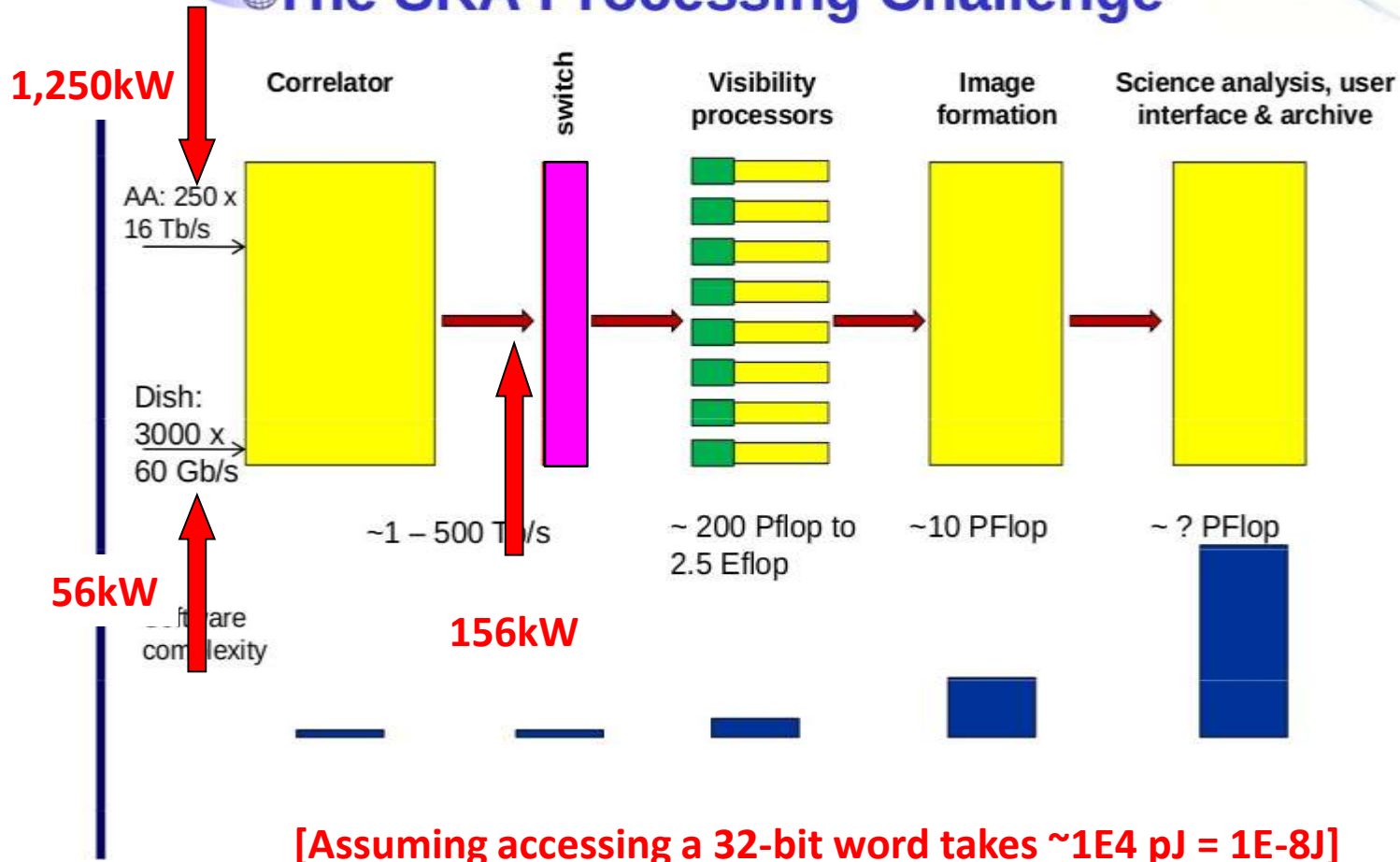


# Access vs Reach





# The SKA Processing Challenge



---

## **Moving data costs a tremendous amount of energy**

Hundreds of kilowatts to megawatts for SKA

Moving data over hundreds or thousands of km will cost even more energy

## **Switching costs a lot of energy too**

A non-blocking thousands-of-ports switch will consume power at a staggering rate

## **This will probably not scale well**

N-dimensional torus may be feasible; but is not currently COTS

## **In addition, getting data into various machines is challenging**

see LOFAR case

## **Limiting the port count can significantly reduce cost**

By avoiding low efficiency network topologies

Should investigate design of an over-subscribed network

## **Finally, need to reduce data flows to keep SKA affordable**

Otherwise moving bits will end up dominating the power budget

# Key messages:

---

- **Sustained COTS exascale is coming: ~2020**
  - US and other government-funded initiatives for e.g. energy and climate modelling and related research
  - Substantial change in hardware architectures will drive change in software: e.g. multiple threading across millions of cores (IESP)
- **But, SKA1's requirements will push the 2020 envelope of COTS hardware capabilities**
- **Purpose built hardware solutions for SKA1 will also be subject to the challenges of computation at exascale:**
  - Amdahl's laws
  - Energy for interconnection
  - Reliability