

# OPEN

Compute Engineering Workshop

March 9, 2015

San Jose



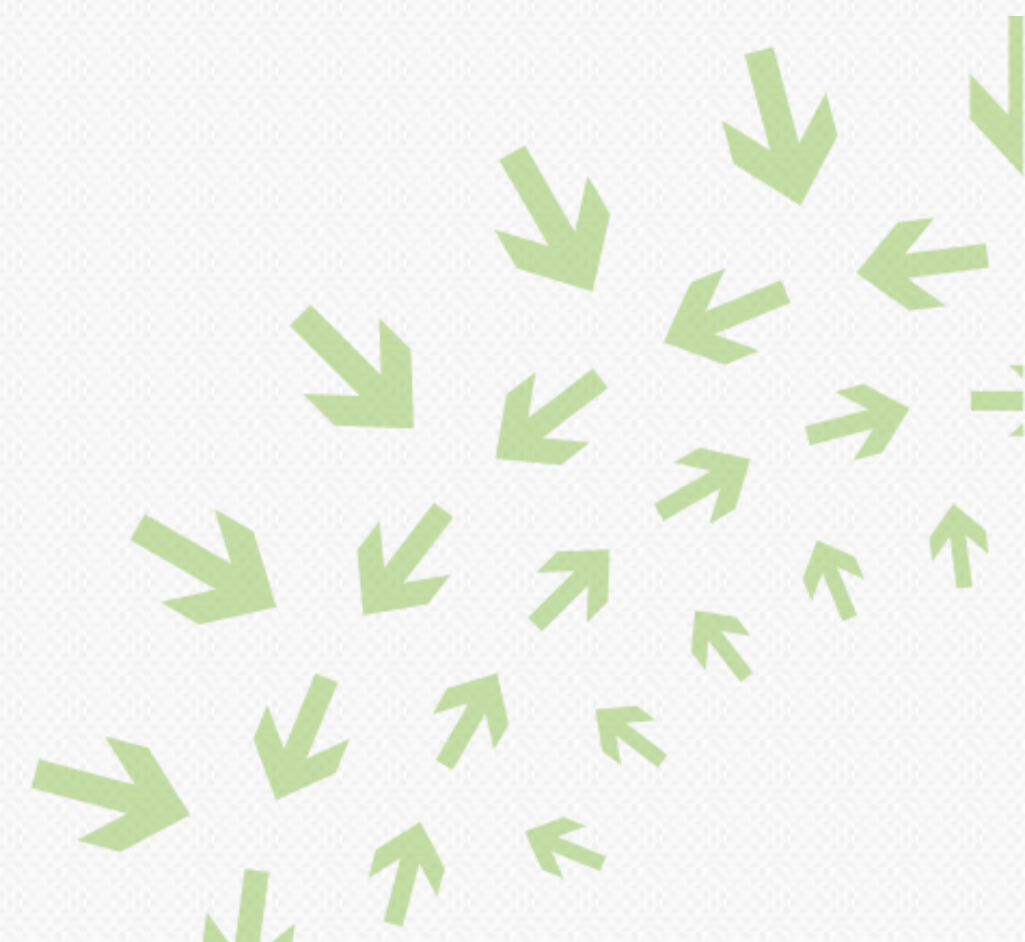
# Express Fabric within an Open Compute Project and hybrid architecture based on x86/ARM

## The DaaP project

Jean-Marie Verdun

Horizon Computing Solutions

President



# Project background

- Started in June 2014 with the intend to study new generation fabric based on PCI-Express Gen 3
- Released under OCP license
- Specifications and Implementations are open
- This project is designed in collaboration with OCP community members in Europe, Thales and Horizon Computing Solutions





# Data center as a PCB

- Run the PCB up to 35 degrees celsius ambient using air cooling or even more using immersive cooling
- Share any expensive features like I/O boards or management between multiple compute nodes
- Remove as much as possible cabling which does represent up to 70% of datacenter failures
- High resiliency capabilities through CPU count increase

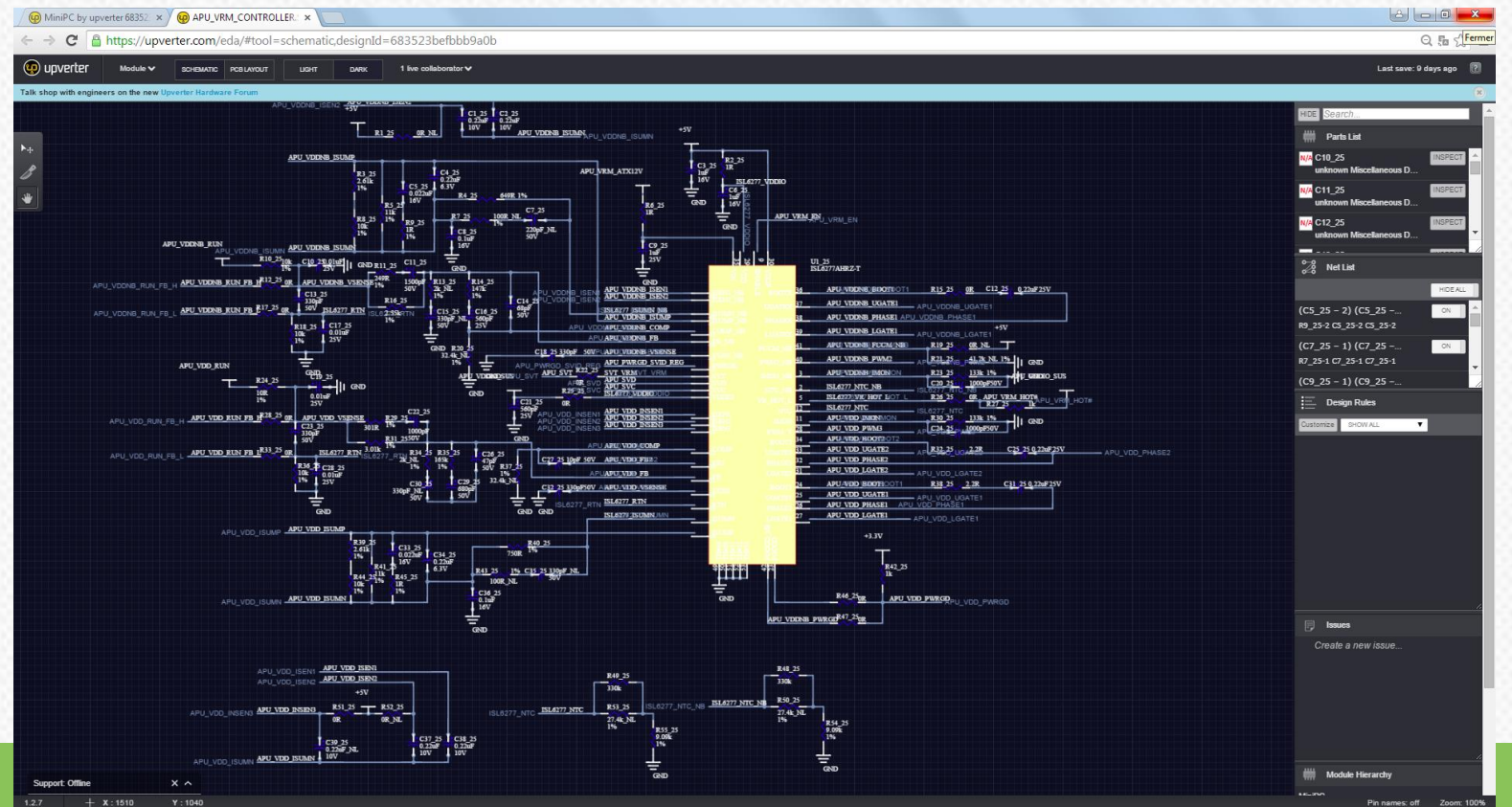
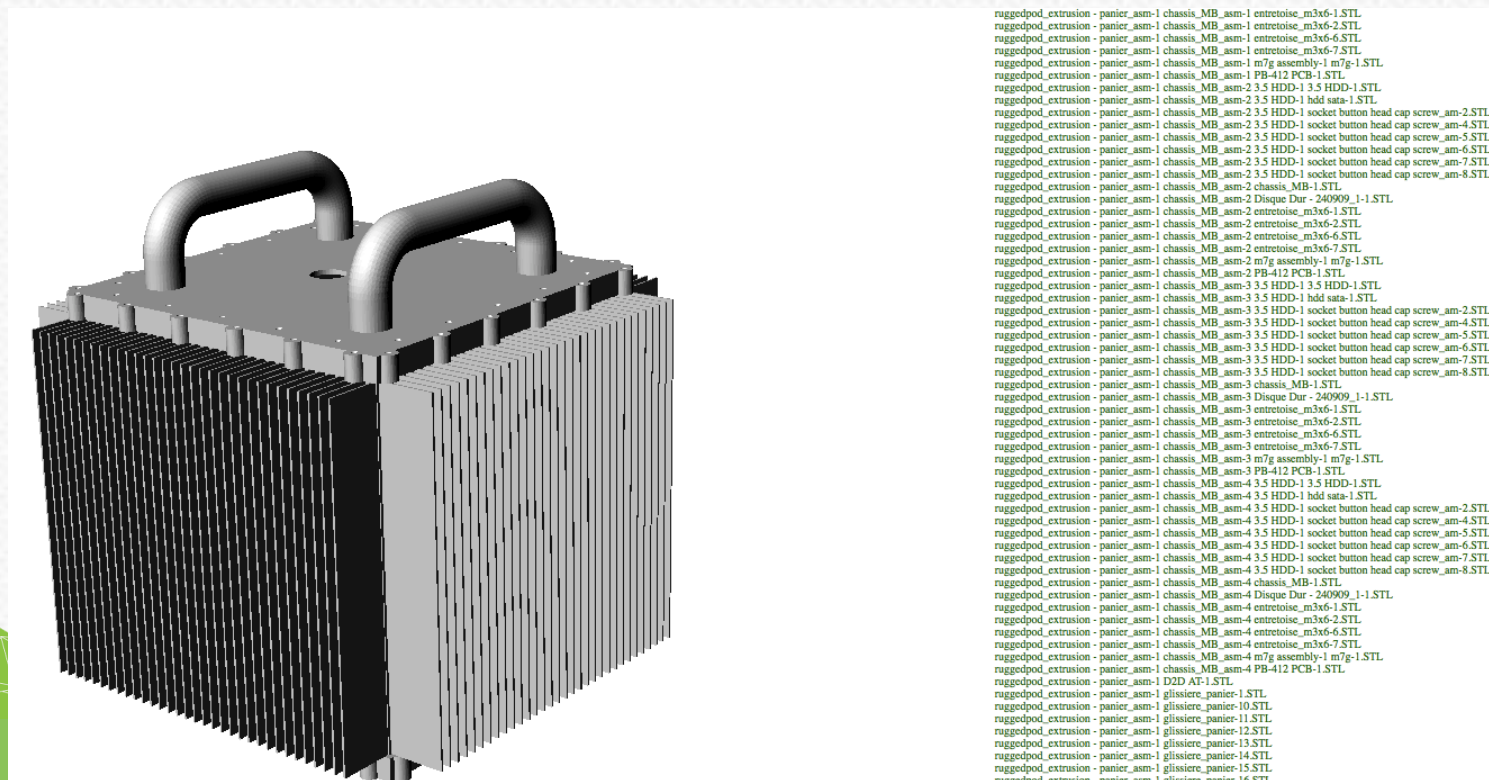
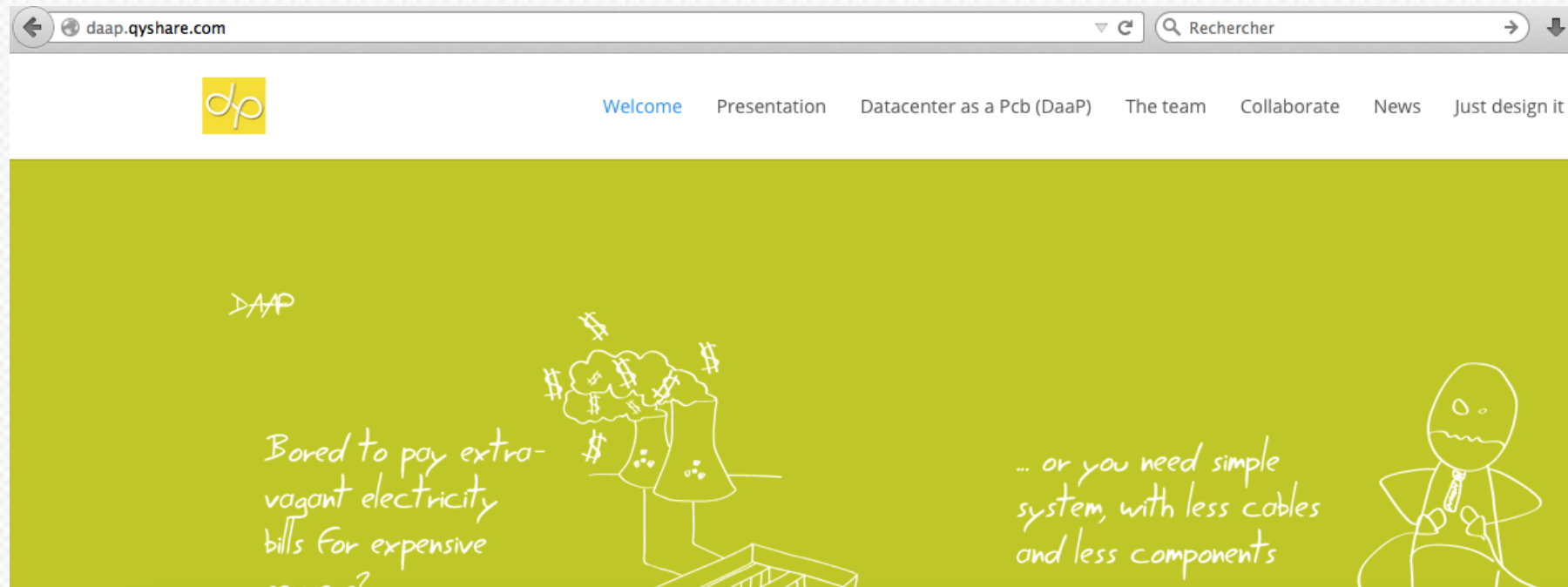
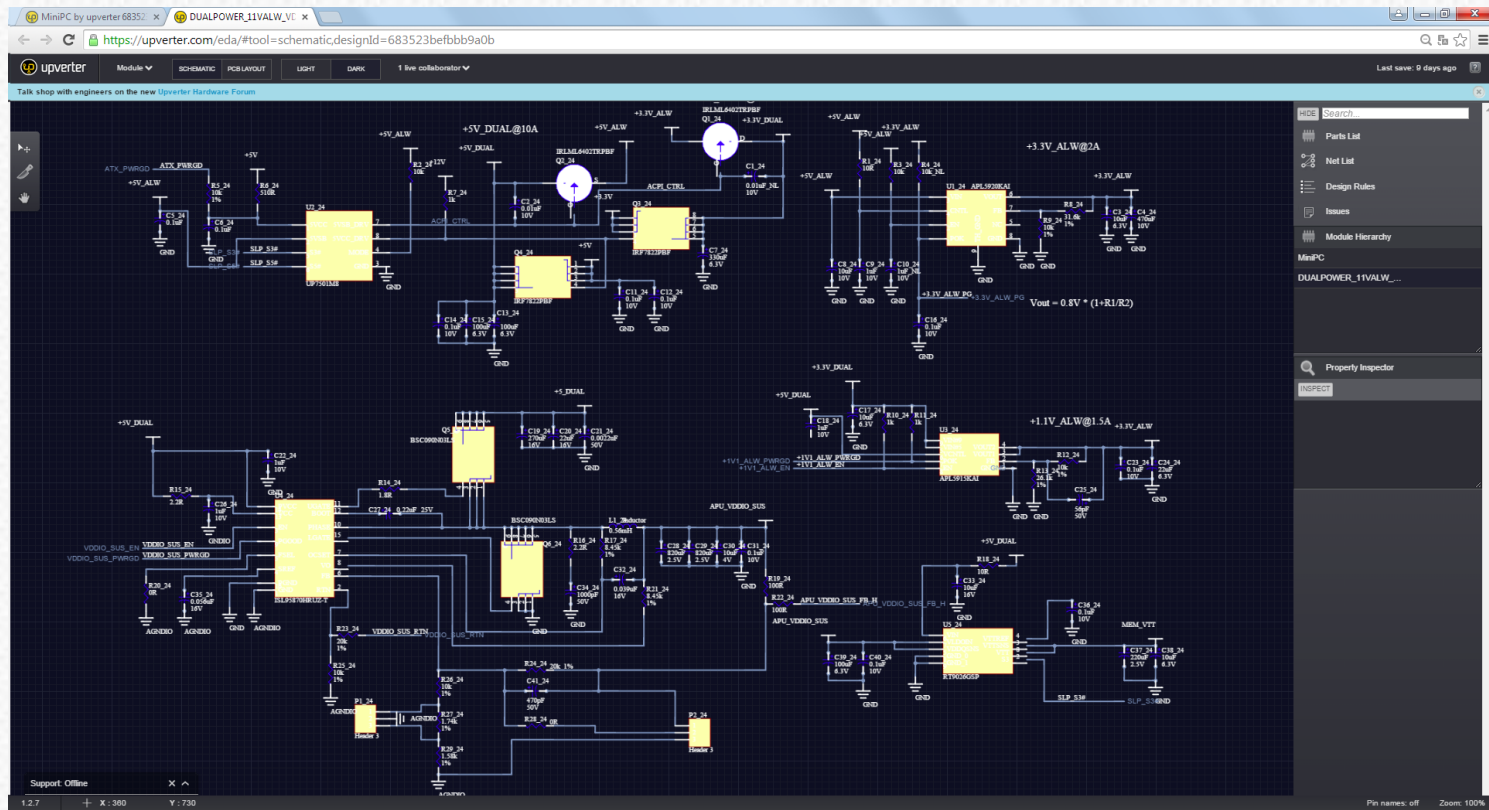


# How to work together ?

- We build up from ground a VM that is providing:
- Project Management based on OpenProject
- Email to team members and calendar based on SOGo
- Web server with private content creation for team exchange
- Cloud storage for data exchange
- SSO to the various tool
- 3D WebGL player for mechanical
- Need an example: <http://daap.qyshare.com> or <http://ruggedpod.qyshare.com>



# Participate ?



# Up to 35 degrees

## Control $W/mm^2$ and use high $T_j$ components

- Xeon or Opteron are low  $T_{case}$  chips around 65 C
- $dT$  with ambient is crucial and low thermal resistance heatsink are required
- Mobile chips have a much higher  $T_j$  which may vary between 95 to 105 C.
- Increase the number of PWM stages from VRM.





# Shared I/O

## PCIe Fabric

- Any modern chip is coming up with a PCIe interface
- Cost of the PCIe interface is 0 \$US against 300 \$US for a 40Gb/s NIC
- PLX technology is currently developing an SR-IOV with multi root support PCIe switch which might provide tremendous results.





# Reliability and resiliency

## Solder down everything (except RAM)

- CPU and chipset have a high reliability rate
- CPU socket requires higher footprint space than BGA solutions. Use BGA to increase node count.
- RAM are second cause of failure within datacenter



# Management

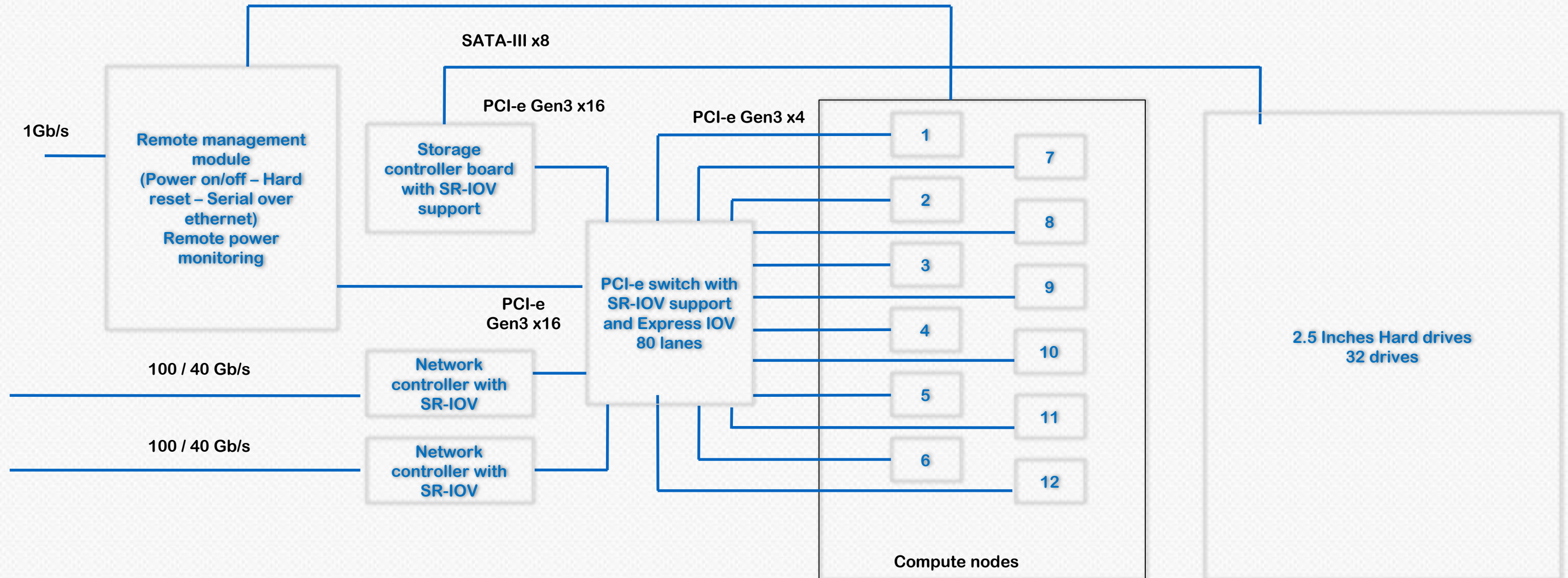
## Get rid of 1 BMC per node

- Traditional management systems requires 1 BMC per system board
- Moving management task out of system board might simplify
  - Node design and improve density
  - Management board design
- Adapt a 2 states management for compute node



# Block diagram

- ✓ First generation Compute node based on **AMD APU**
- ✓ Second generation Compute node based on **AMD HeroFalcon ARM** and **APM XGene**

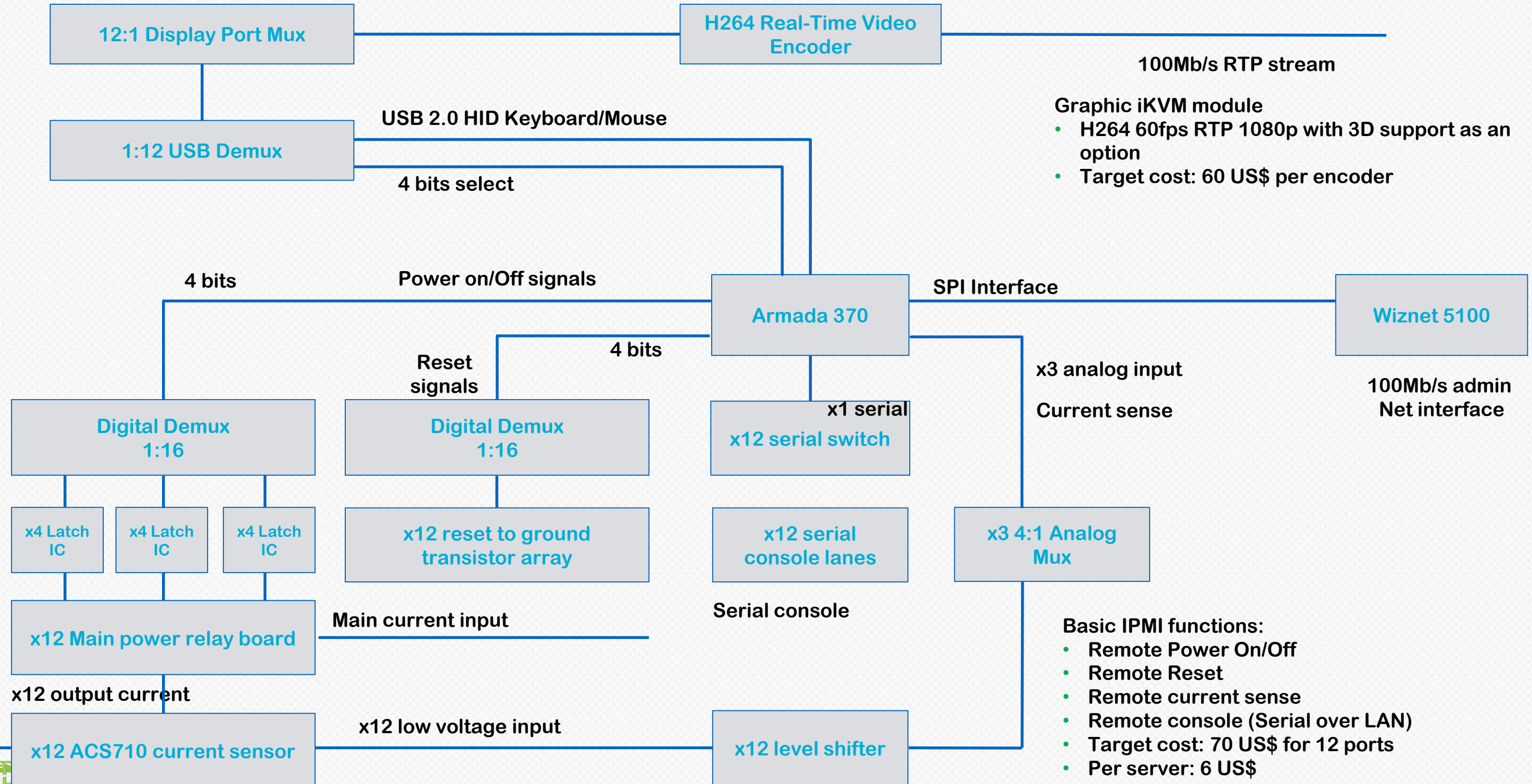


- ✓ 48 cores in 1 U
- ✓ 384 GB ECC main memory
- ✓ 48Gb/s local drive bandwidth
- ✓ 64 TB local storage
- 32 Gb/s interconnect per node
- 384Gb/s local fabric
- 80 Gb/s network output





# Remote management module



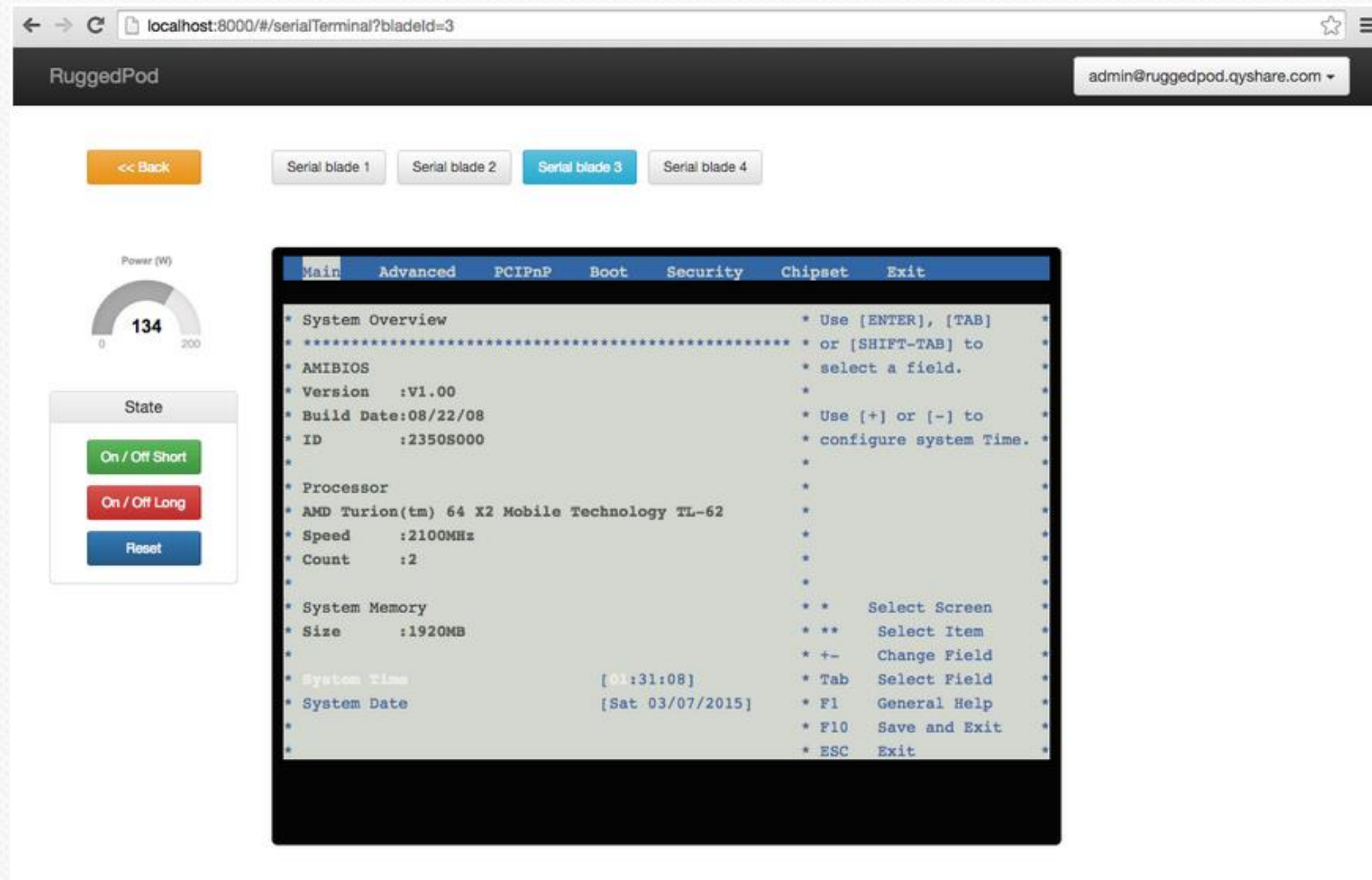
# Firmware

## Go OpenSource for flexibility and customer added value

- System firmware (traditionally called BIOS) are proprietary firmware which lacks of innovation. DaaP will adopt Open Source firmware
- What might be coming from OpenSource ?
  - BIOS through coreboot
  - Management node firmware through Rest API
  - Network boot firmware (PXE removal)
  - PCIe Fabric configuration firmware



# What it looks like ?



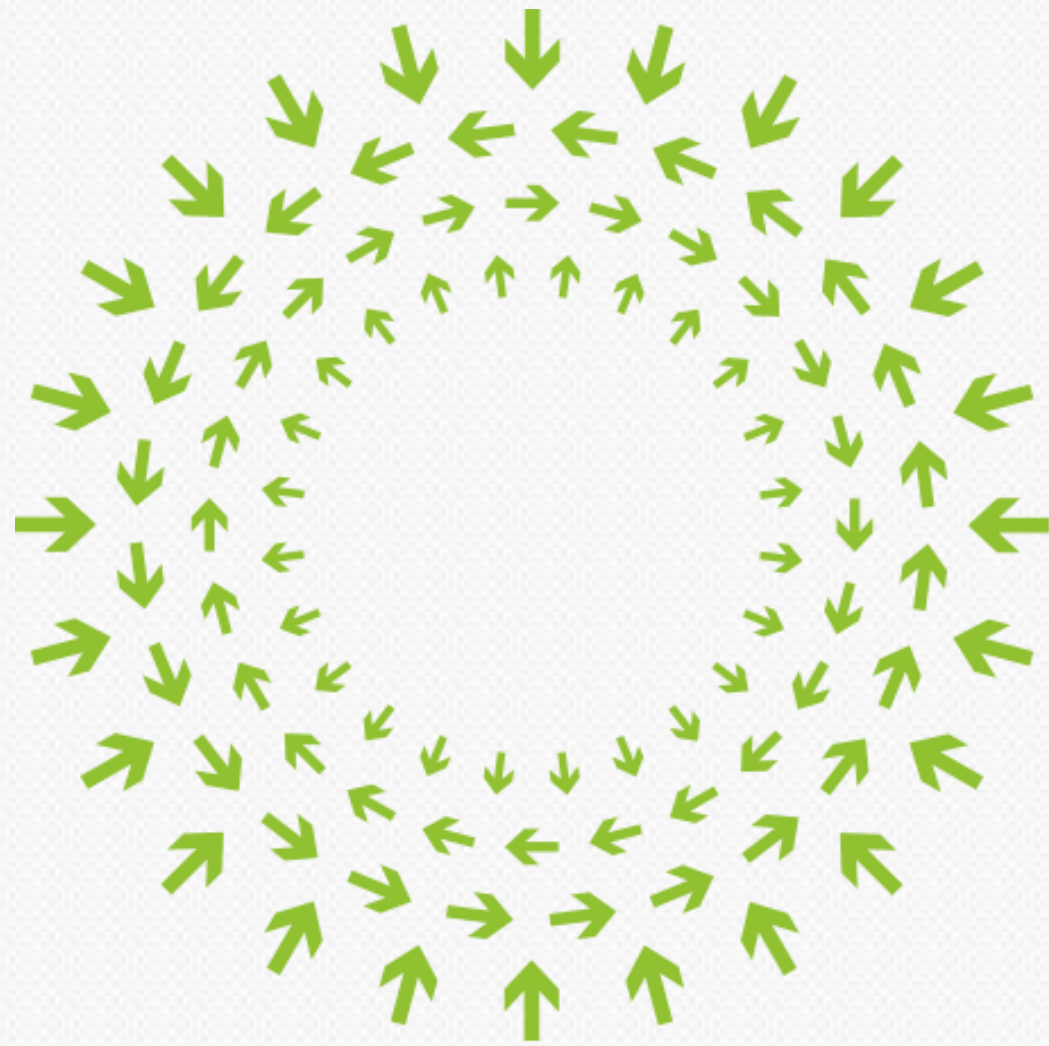


# Estimated production cost

All Cost in \$US	Quantities per board		Total
Management Module (without iKVM)	70	1	70
Compute module			
1 CPU Quad Core 4 Ghz 2MB L2	110	12	1320
32 GB Main memory	200	12	2400
Compute board and "accessories" (like SB)	70	12	840
PCI-e Switch	400	1	400
40 Gb/s NIC	300	2	600
Storage HSA with SR-IOV	300	1	300
PCB	400	1	400
Mechanical	200	1	200
Total			6530

<b>Per Server (\$US)</b>	<b>544</b>
Total Power (Watts) full load	720
Per server (Watts) full load	60





# OPEN

Compute Engineering Workshop

March 9, 2015

San Jose