



# Evaluation of low power Systems on Chip for scientific computing

Roberto ALFIERI, Enrico CALORE, Daniele CESINI, Andrea FERRARO,  
Michele MICHELOTTO, **Lucia MORGANTI**, Fabio SCHIFANO - INFN

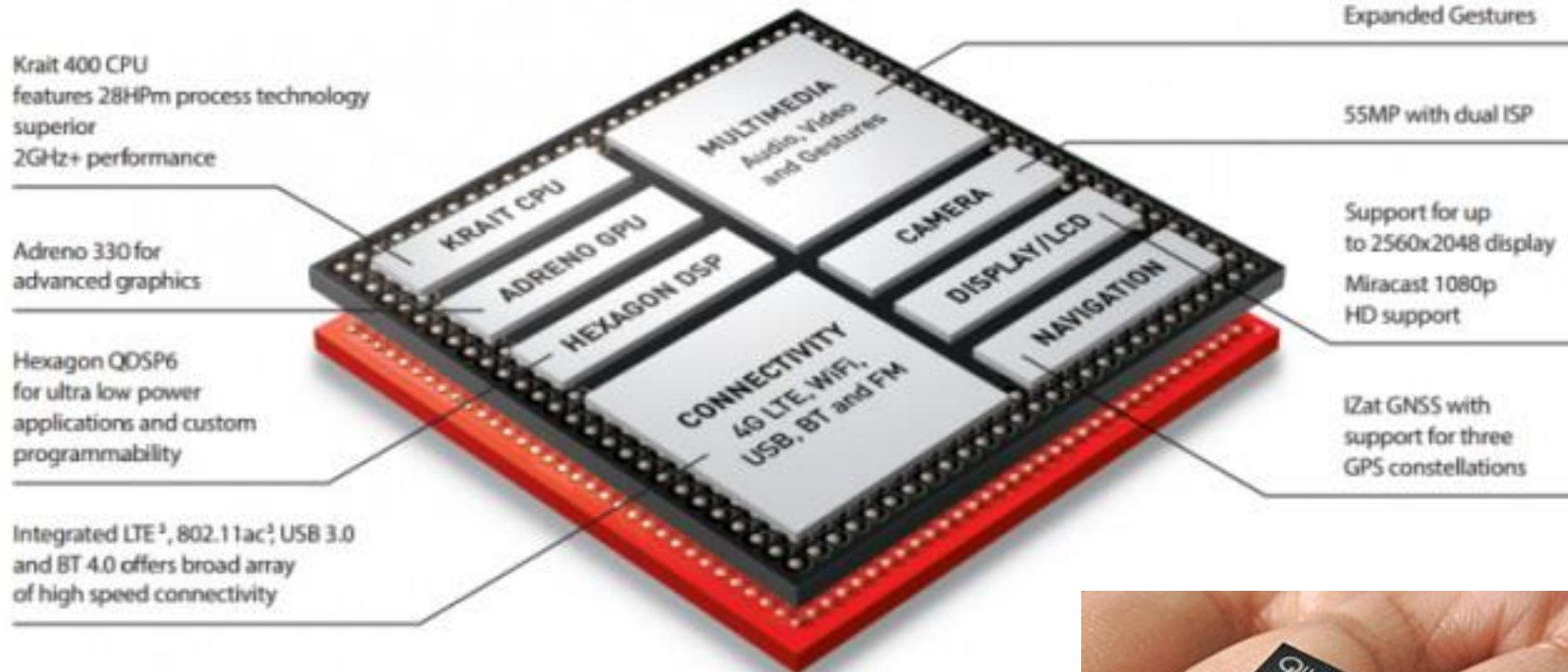


# Outline

- SoCs: an overview
- Computing on SoCs: limitations and scientific applications
- Low power from Intel: preliminary tests
- Conclusions

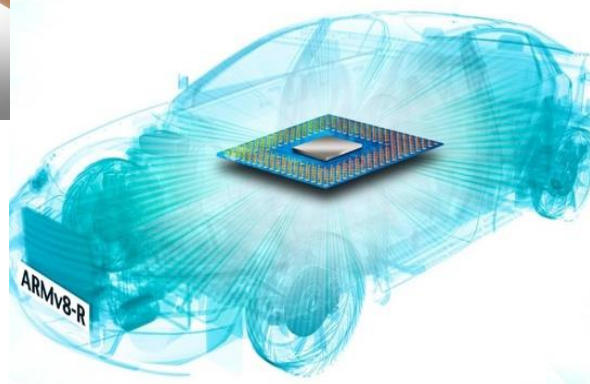
# Systems on Chip (SoCs)

## Qualcomm Snapdragon 800 Processors



# SoCs: where do we find them?

- Mobile
- Embedded
- Wearable
- IoT



# + Ok, but then... an iPhone cluster?

- NO, we are not thinking of building an iPhone cluster
- We want to use these processors in a standard computing centre configuration
  - Rack mounted
  - Linux powered
  - Running scientific application mostly in a batch environment
- ... Use development boards...



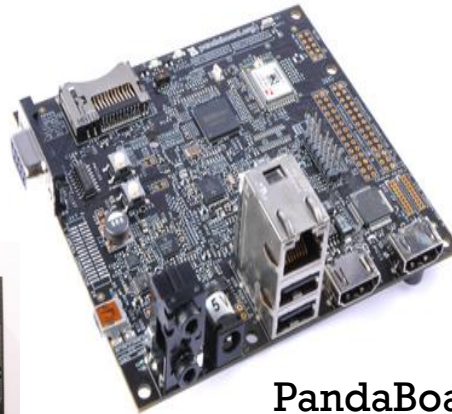
# Plenty of nice boards



WandBoard



Rock2 Board



PandaBoard



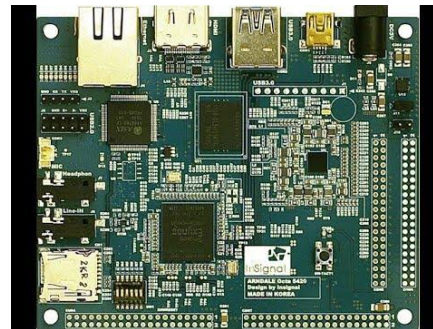
DragonBoard



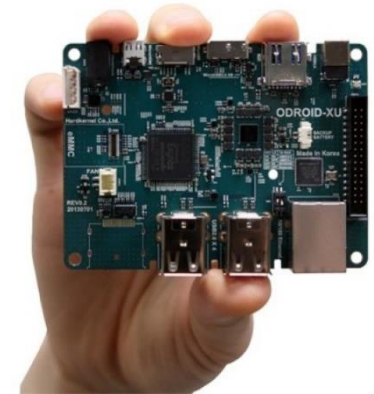
SabreBoard



CubieBoard



Arndale OCTA Board



Odroid-XU

[http://elinux.org/Development\\_Platforms](http://elinux.org/Development_Platforms)

■ ... and counting...

# Some specs

<b>BOARD</b>	<b>Model</b>	<b>ARM IP</b>	<b>GPU/DSP</b>	<b>GFLOPS (GPU)</b>	<b>GFLOPS (CPU+GPU)</b>	<b>Eth</b>
<b>FREESCALE (Embedded SoC)</b> SABRE Board	<b>Freescale</b> i.MX6Q	<b>ARM</b> A9(4)	<b>Vivante</b> GC2100	19.2	25	1Gb
<b>ARNDALE (Mobile SoC)</b> Octa Board	<b>Samsung</b> Exynos 5420	<b>ARM</b> A15(4) A7(4)	<b>ARM</b> Mali-T628 MP6	110	115	10/100
<b>HARDKERNEL (Mobile SoC)</b> Odroid-XU-E	<b>Samsung</b> Exynos 5410	<b>ARM</b> A15(4) A7(4)	<b>Imagination Technologies</b> PowerVR SGX544MP3	51.1	65	10/100
<b>HARDKERNEL (Mobile SoC)</b> Odroid-XU3	<b>Samsung</b> Exynos 5422	<b>ARM</b> A15(4) A7(4) <b>(HMP)</b>	<b>ARM</b> Mali-T628 MP6	110	130	10/100
<b>INTRINSIC (Mobile SoC)</b> DragonBoard	<b>Qualcomm</b> Snapdragon 800	<b>Qualcomm</b> Krait(4)	<b>Qualcomm</b> Adreno 330	130	145	1Gb
<b>TI (Embedded SoC)</b> EVMK2H	<b>TI Keystone</b> 66AK2H14	<b>ARM</b> A15(2)	<b>TI</b> MS320C66x	189	210	1Gb (10Gb)



# NVIDIA JETSON K1



**TEGRA K1**  
192-core  
Kepler-Class Chip

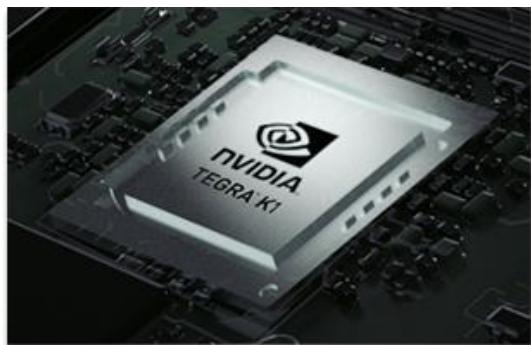
Quad A15 CPUs  
32-bit  
3-way Superscalar  
Up to 2.3GHz  
32K+32K L1\$



- First **ARM+CUDA programmable** GPU-accelerated Linux development board
- 4 cores ARM A15 CPU
- 192 cores NVIDIA GPU → 300 GFLOPS (peak sp)
- ... for less than 200 Euros
- Only 32bit



# GPU acceleration in K1



4 core ARM A15 ~ 18 GFLOPS  
 Kepler SMX1 192 core ~ 300 GFLOPS (sp)  
 ~ 15 W  
 ~ **21 GFLOPS /W**



2 x (E5-2673v2 (IvyBridge) 8 core)  
 ~ 200 GFLOPS (dp)  
 220 W



NVIDIA TESLA K40 2880 core  
 ~ 1400 GFLOPS (dp) ~ 4300 GFLOPS (sp)  
 235 W

**CPU+GPU ~ 3 GFLOPS/W dp**  
**~ 9 GFLOPS/W sp**

# + How can we program SoCs (in a Linux environment)?

- GCC+OpenMP+MPI available for ARM architectures
- CUDA available only on the Jetson K1
  - Computing capability 3.2 (vs 3.5)
- OpenCL for the GPU
  - If you are lucky enough to find working drivers
- GCC5+OpenMP4 tests ongoing...

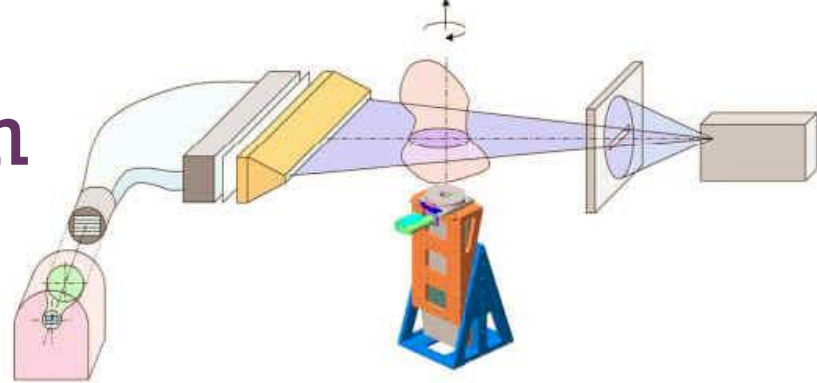


# Limitations

- Commodity SoCs and development boards have a number of limitations:
  - 32 bit (looking forward to 64 bit)
  - Small RAM size in the boards (O(2GB))
    - However modern SoCs can address 40bit
  - No ECC memory
  - Reliability to be tested for production environment
  - Slow connections (10/100Mb eth) in many cases
    - Ethernet via USB in some boards
  - HW bugs



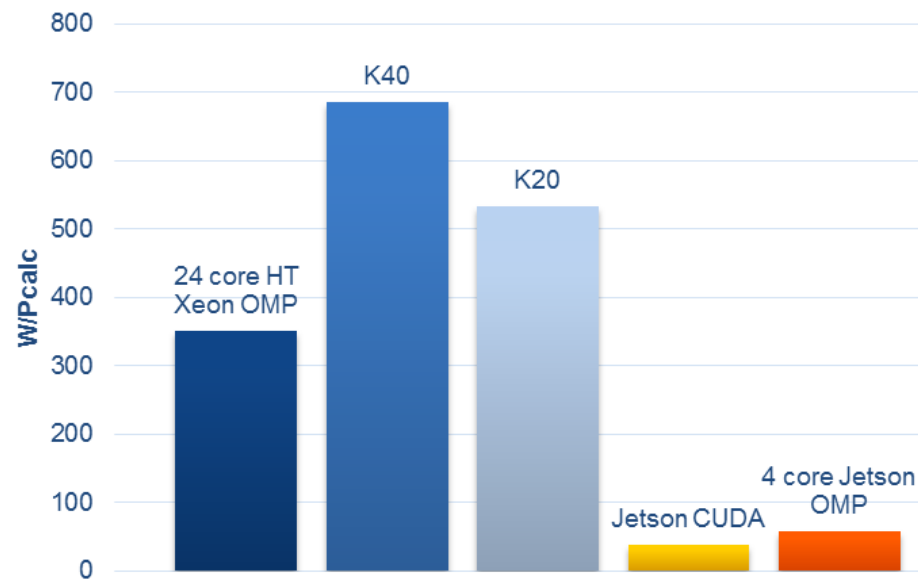
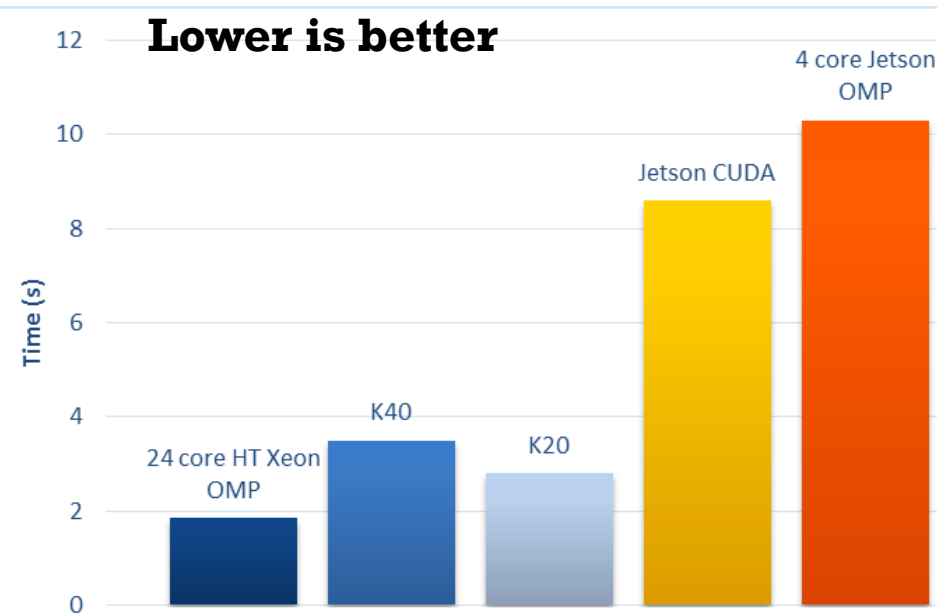
# Filtered Backprojection



CPU and GPU

1024 x 1024 pixel

**Lower is better**



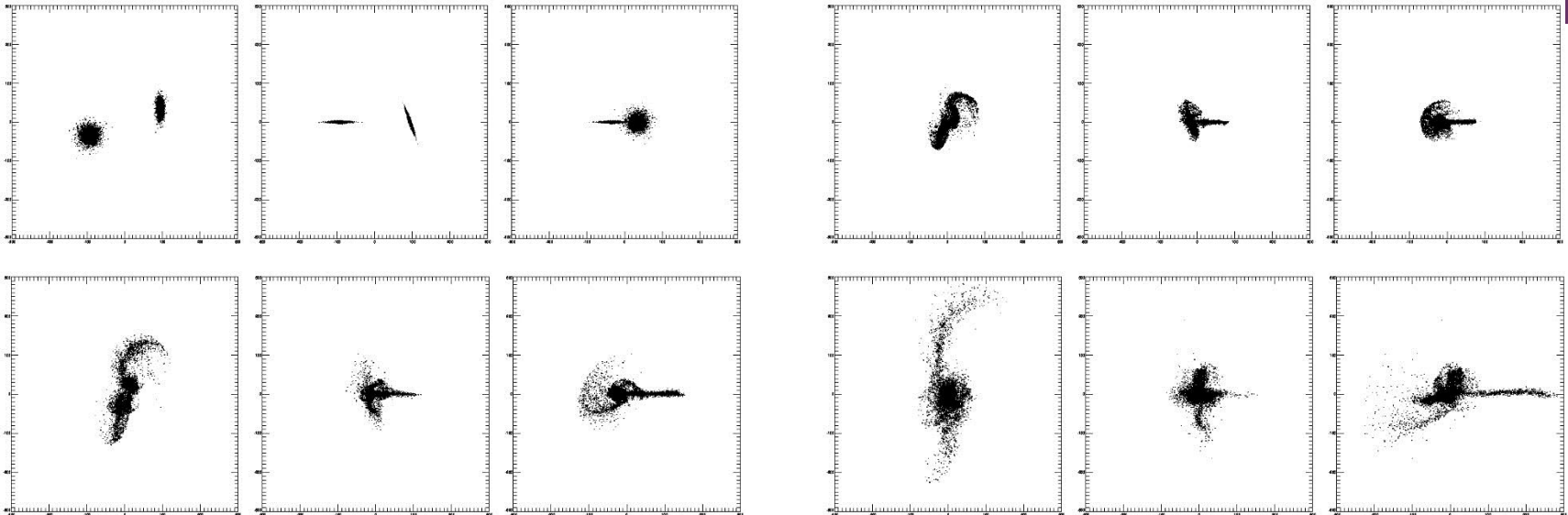
On (2xE5-2620+K20): 3241 slices reconstructed in 1 h: **350 Wh**

On **5** x Jetson-K1: 3840 slices reconstructed in 1 h: **41 Wh**

# Merging galaxies with Gadget-2

MPI code for cosmological simulations by Volker Springel

<http://www.mpa-garching.mpg.de/gadget/>



- 60000 disk and halo particles,  
Barnes Hut N-body simulation

- Jetson-K1      **10.9 W**

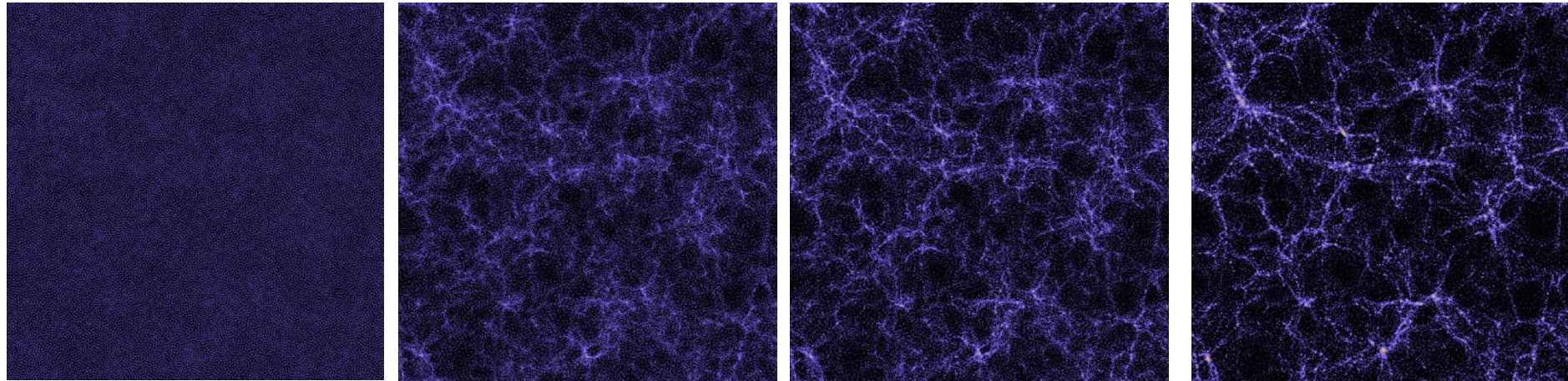
- Xeon            **~220 W**

**Time (s)**

N cores	Jetson K1	Xeon	MIC
4	950,00	446,43	4871,21
8	961,85	290,96	3017,73
16		229,47	1926,00

# + Evolving the universe with Gadget-2

$128^3$  particles in a cube of 250 Mpc side, 13.6 Gyr



- Jetson-K1 about 2.6X slower than Xeon using 4 CPU cores, and about 4.4X slower using 8 CPU cores

- Jetson-K1      **12.1 W**
- Xeon            **~220 W**

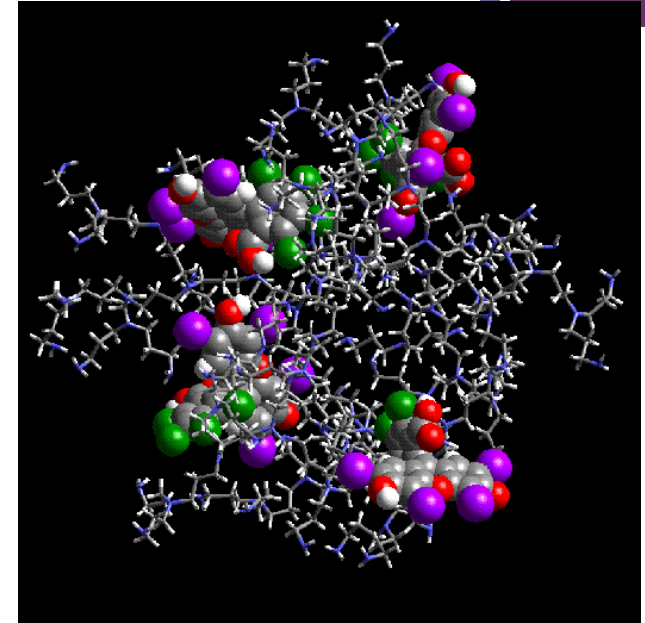
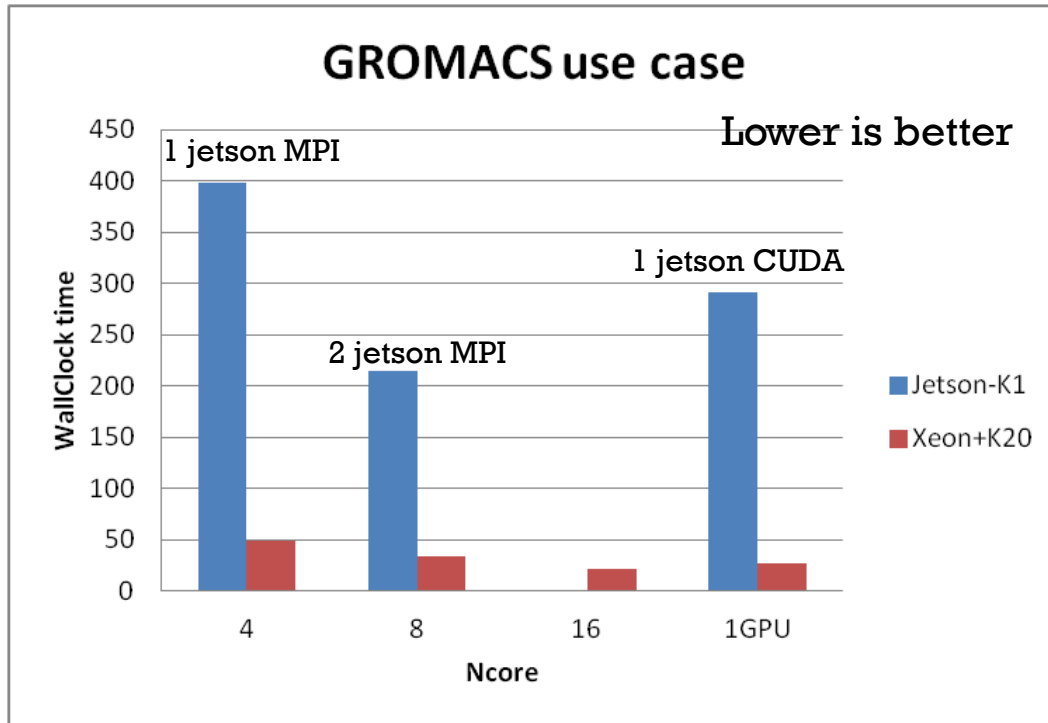
**Time (s)**

N cores	Jetson K1	Xeon	MIC
4	11649,30	4445,58	
8	10563,76	2383,3	
16		1344,78	10608,96

# Molecular Dynamics on Jetson-K1

CPU and GPU

**Parallel application for CPU and GPU:  
real life use case with GROMACS**

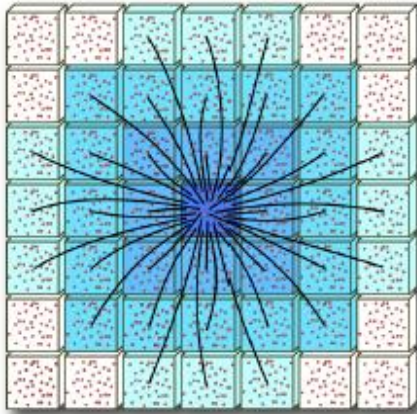


- Jetson-K1 about 10X slower using the same number of CPU cores
- Jetson-K1 about 10X slower using the GPU (vs. an NVIDIA Tesla K20)
  - Jetson-K1 13.5 W
  - Xeon+K20 ~320 W

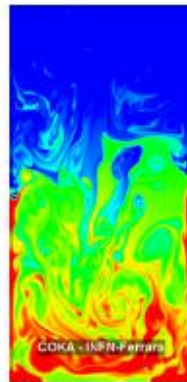
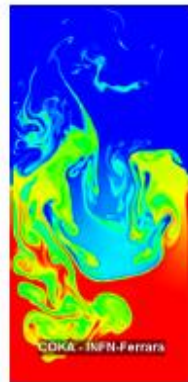
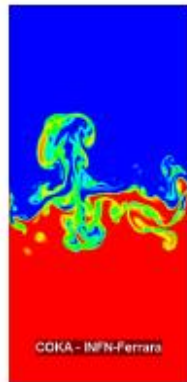
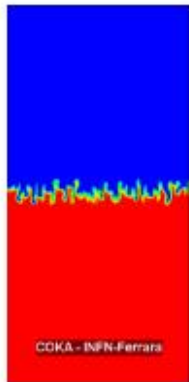
# Lattice Boltzmann on the Tegra K1

GPU only

## Lattice Boltzmann Methods: D2Q37



(\*) Schifano et al.  
*A portable OpenCL  
Lattice Boltzmann code  
for multi- And many-core  
processor architectures,*  
Proc. Comp. Sci. 29, 2014



## Performance comparison with K20

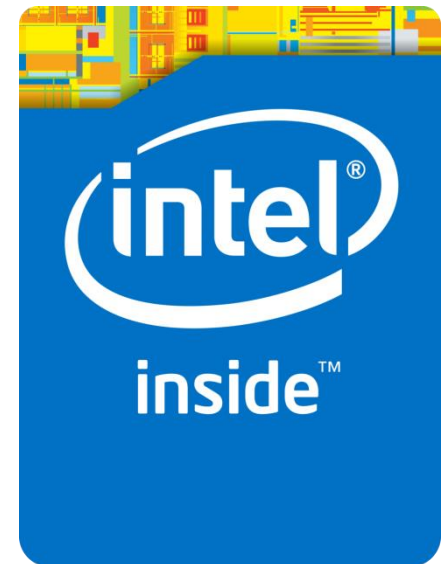
Propag. (MLUPS)		Collide (MLUPS)	
Tegra	K20	Tegra	K20
17.8	256.0	1.6	89.4

- Porting easier than expected
- Performance under investigation



# + Only ARM based SoCs? And Intel?

- INTEL produce SoCs
  - Probably you have one in your laptop
- Some of them are low power
- Already 64bit
- Integrated GPU
  - CILK++ programmable
  - OpenCL programmable

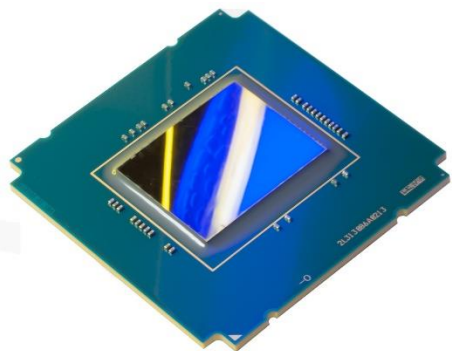


# Some low power from Intel

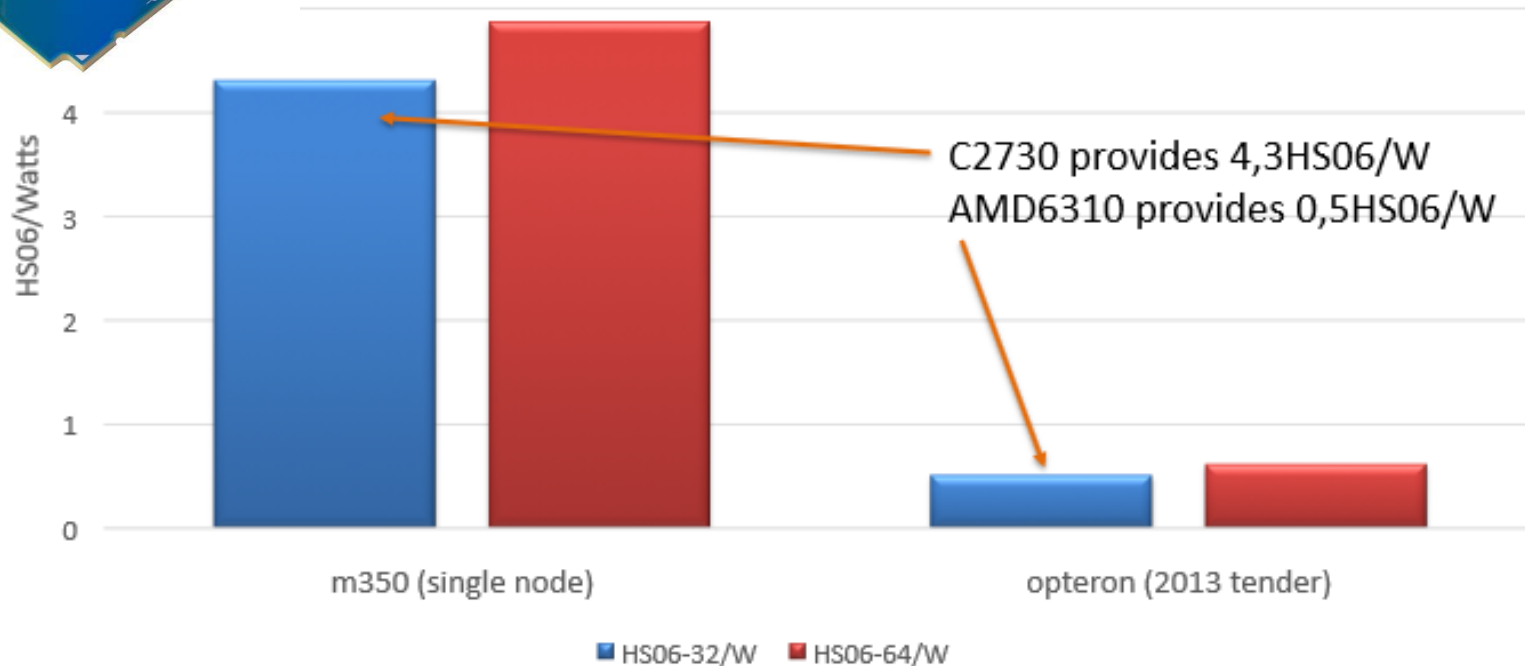
Nome prodotto	Intel® Atom™ Processor E3845 (2M Cache, 1.91 GHz)	Intel® Core™ M-5Y71 Processor (4M Cache, up to 2.90 GHz)	Intel® Core™ i7-4578U Processor (4M Cache, up to 3.50 GHz)	Intel® Core™ i7-4702EC Processor (8M Cache, up to 2.00 GHz)	Intel® Atom™ Processor C2730 (4M Cache, 1.70 GHz)
Nome in codice	Bay Trail	Broadwell	Haswell	Haswell	Avoton
<b>Informazioni di base</b>					
Stato	Launched	Launched	Launched	Launched	Launched
Data di lancio	Q4'13	Q4'14	Q3'14	Q1'14	Q3'13
Numero di processore	E3845	5Y71	i7-4578U	i7-4702EC	C2730
Cache	2 MB L2 Cache	4 MB	4 MB	8 MB Intel® Smart Cache	4 MB
Set di istruzioni	64-bit	64-bit	64-bit	64-bit	64-bit
Opzioni integrate disponibili	Yes	No	No	Yes	No
Litografia	22 nm	14 nm	22 nm	22 nm	22 nm
Prezzo consigliato per il cliente	TRAY: \$52.00	TRAY: \$281.00	TRAY: \$426.00	TRAY: \$459.00	TRAY: \$150.00
Datasheet	<a href="#">Link</a>	<a href="#">Link</a>	<a href="#">Link</a>	<a href="#">Link</a>	<a href="#">Link</a>
Privi di minerali provenienti da zone di conflitto	Yes	Yes	Yes	Yes	
Estensioni set di istruzioni		AVX, SSE	SSE 4.1/4.2, AVX 2.0	SSE 4.1/4.2, AVX 2.0	
Tipo di Bus			DMI2	DMI	
Bus di sistema			5 GT/s	5 GT/s	
<b>Prestazioni</b>					
Numero di core	4	2	2	4	8
Numero di thread	4	4	4	8	8
Frequenza base del processore	1.91 GHz	1.2 GHz	3 GHz	2 GHz	1.7 GHz
TDP	10 W	4,5 W	28 W	27 W	12 W
<b>Specifiche della grafica</b>					
Grafica del processore †	Intel® HD Graphics	Intel® HD Graphics 5300	Intel® Iris™ Graphics 5100		
Frequenza di base grafica	542 MHz	300 MHz	200 MHz		
Frequenza di burst della grafica	792 MHz				
Intel® Quick Sync Video	Yes	Yes	Yes		
Numero massimo di schermi supportati †	2	3	3		
Frequenza massima grafica		900 MHz	1.2 GHz		



# AVOTON on HP Moonshot - HS06



HS06/W: m350 vs 2013 tender



(data from A.Chierici@HEPIX

<https://indico.cern.ch/event/305362/session/2/contribution/22/material/slides/0.pdf> )

# + Conclusions

- Mobile and embedded low power SoCs are becoming attractive for scientific computing
  - In particular if you manage to extract power from the GPU
  - For selected applications
    - Image processing
    - No high RAM/RAM bandwidth requirements
- They still have many limitations for a production environment
  - 32bit, no ECC, bugs, system stability, etc..
  - (BUT we used development boards - not server grade machines)
- NVIDIA K1 in our tests was the most powerful ARM based SoC
  - Easy to install and use
  - Easy to port CUDA based applications
- Intel has interesting low power SoCs
  - Avoton has a high HS06/W ratio
- Looking forward to development boards based on 64bit SoCs with an ARM CPU on board

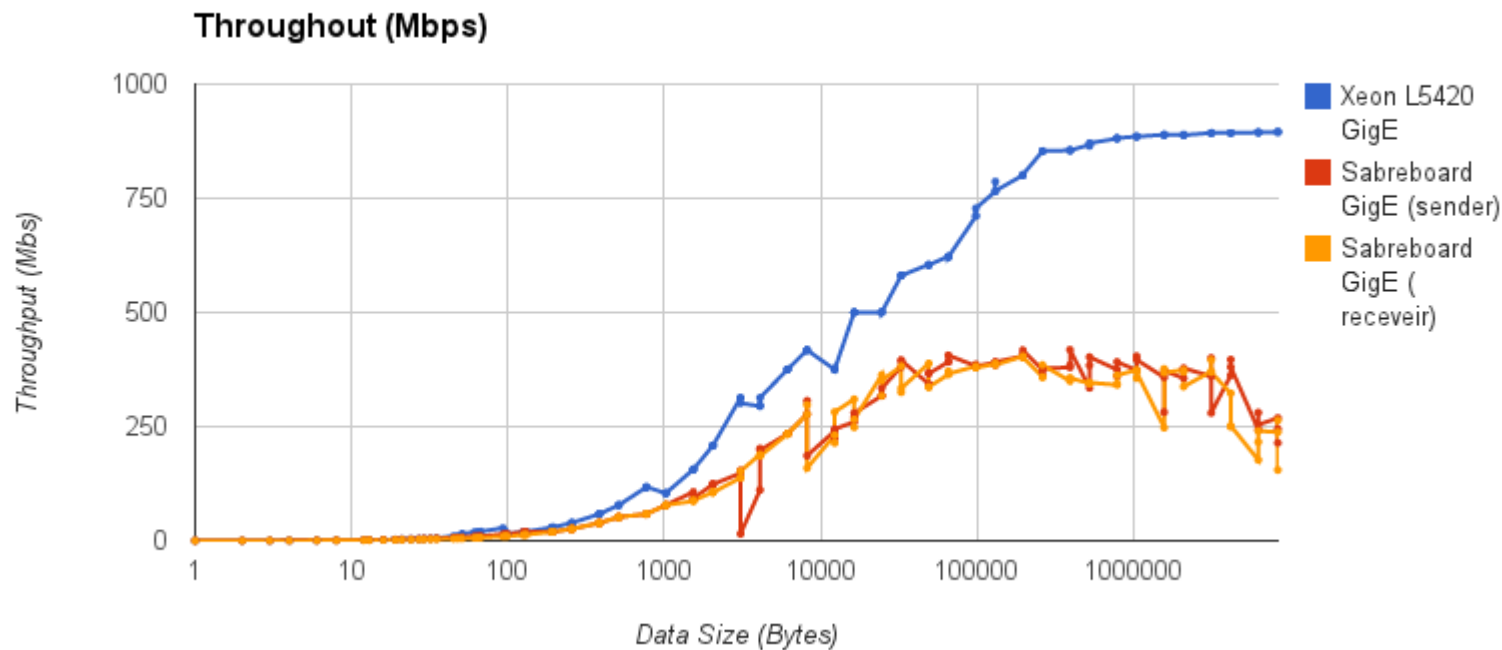


# References

- <http://www.cosa-project.it>
- <http://montblanc-project.eu>
- <https://indico.cern.ch/event/320819/session/3/contribution/30/material/slides/0.pptx>
- <https://indico.cern.ch/event/305362/session/2/contribution/22/material/slides/0.pdf>



# Gbit Ethernet



While latency was comparable to a server class 1Gb ethernet card (50/75 us)

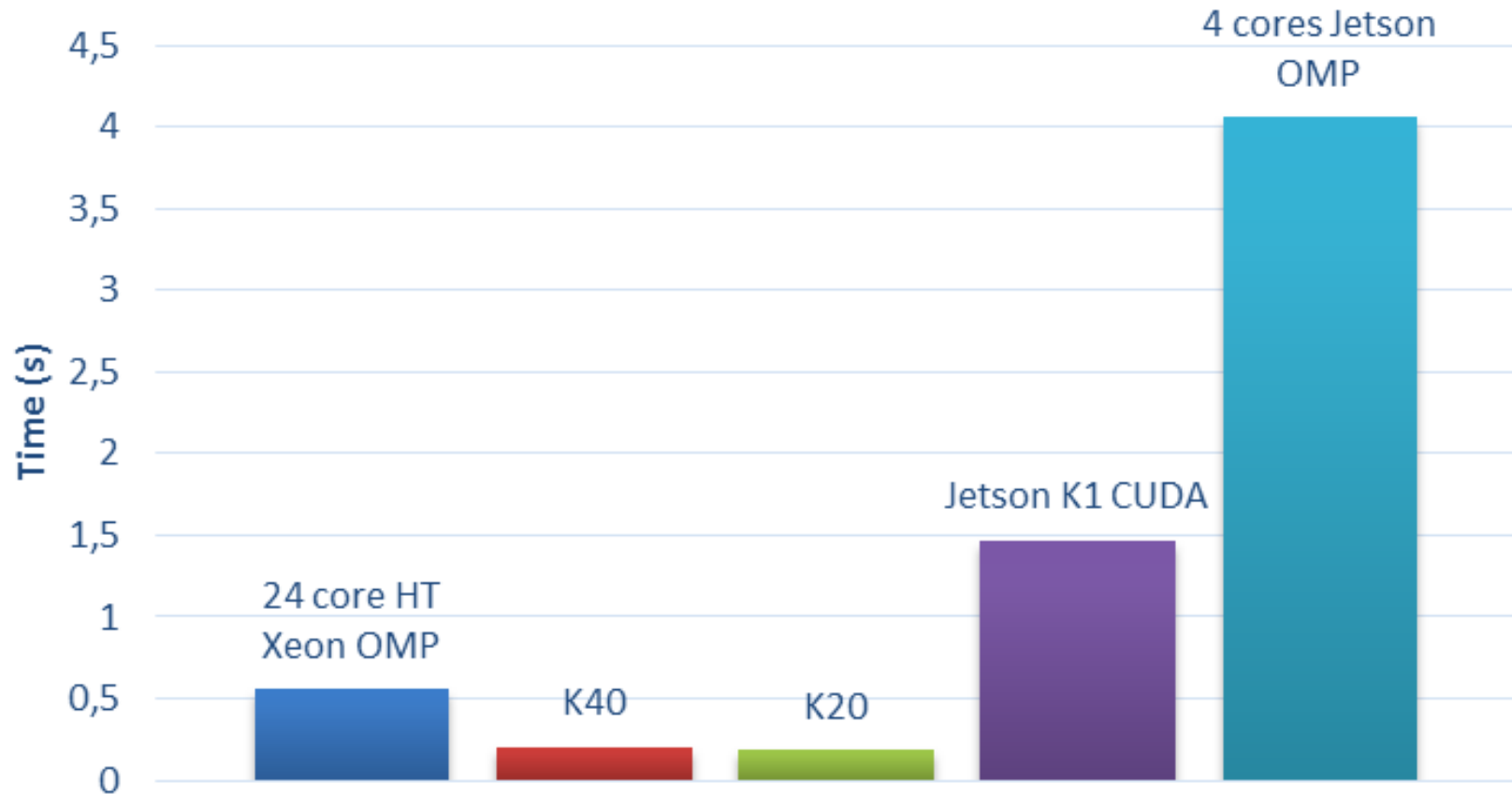


# FFT on CPU and GPU

fftw3 on CPUs  
cuFFT on GPUs

Lower is better

## FFT (Array size $2^{25}$ )



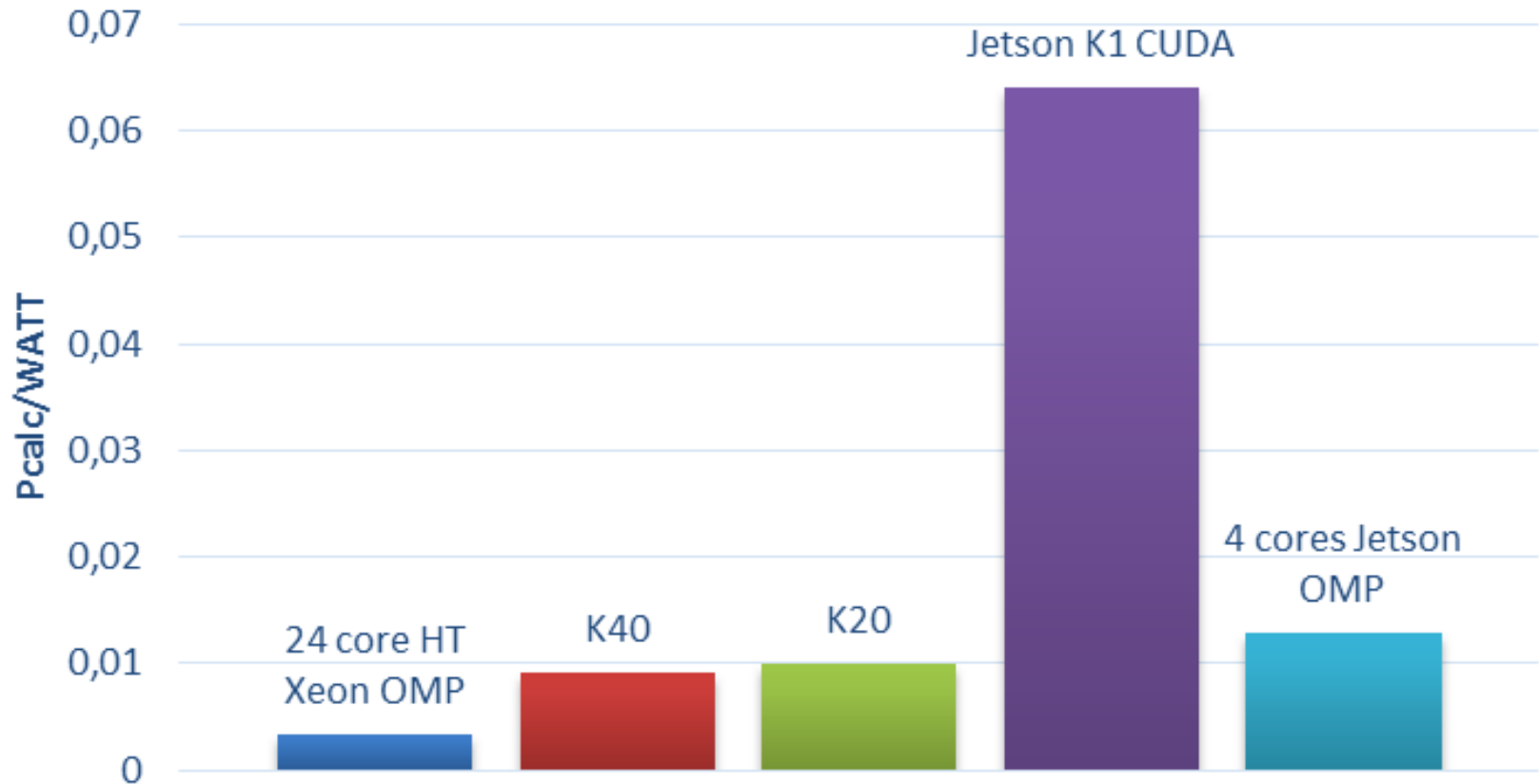


# FFT on CPU and GPU

fftw3 on CPUs  
cuFFT on GPUs

Higher is better

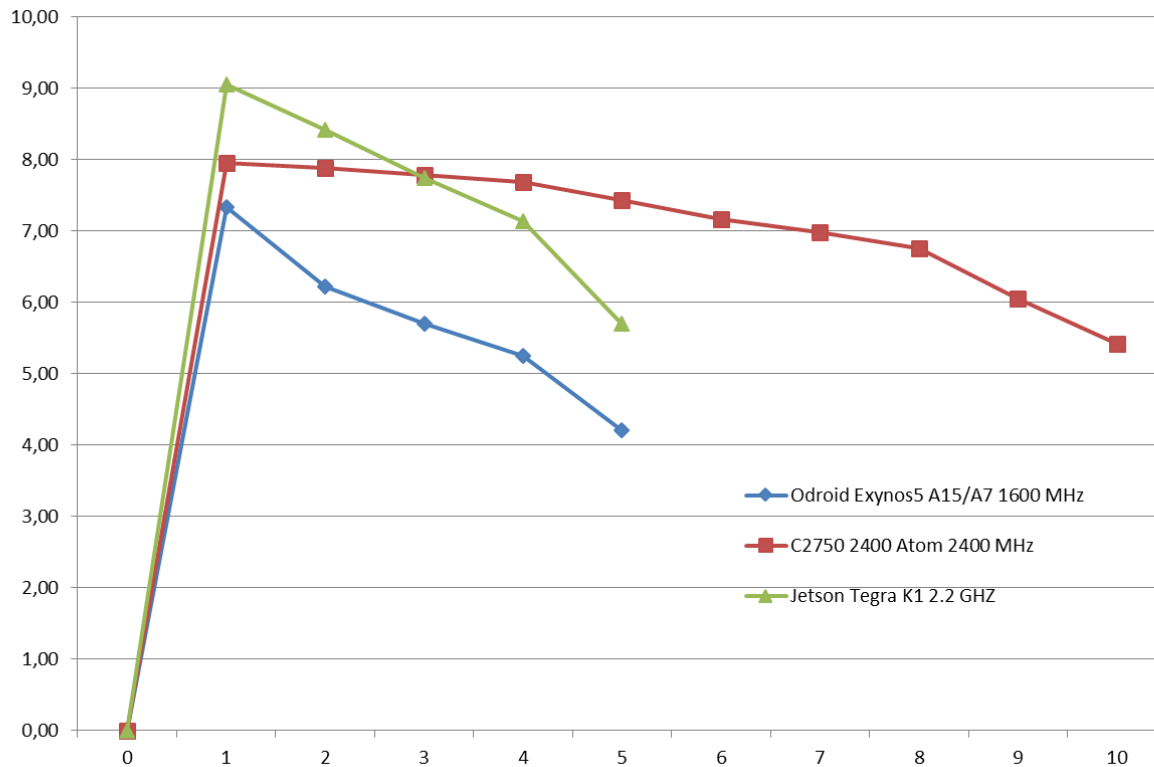
## FFT (Array size $2^{25}$ )





# + HS06

## HS06 on Exynos5, TegraK1 and Atom C2750 – Per core loaded



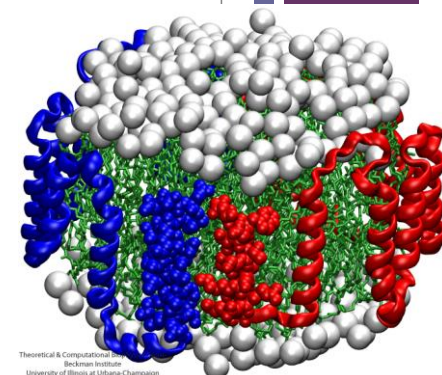
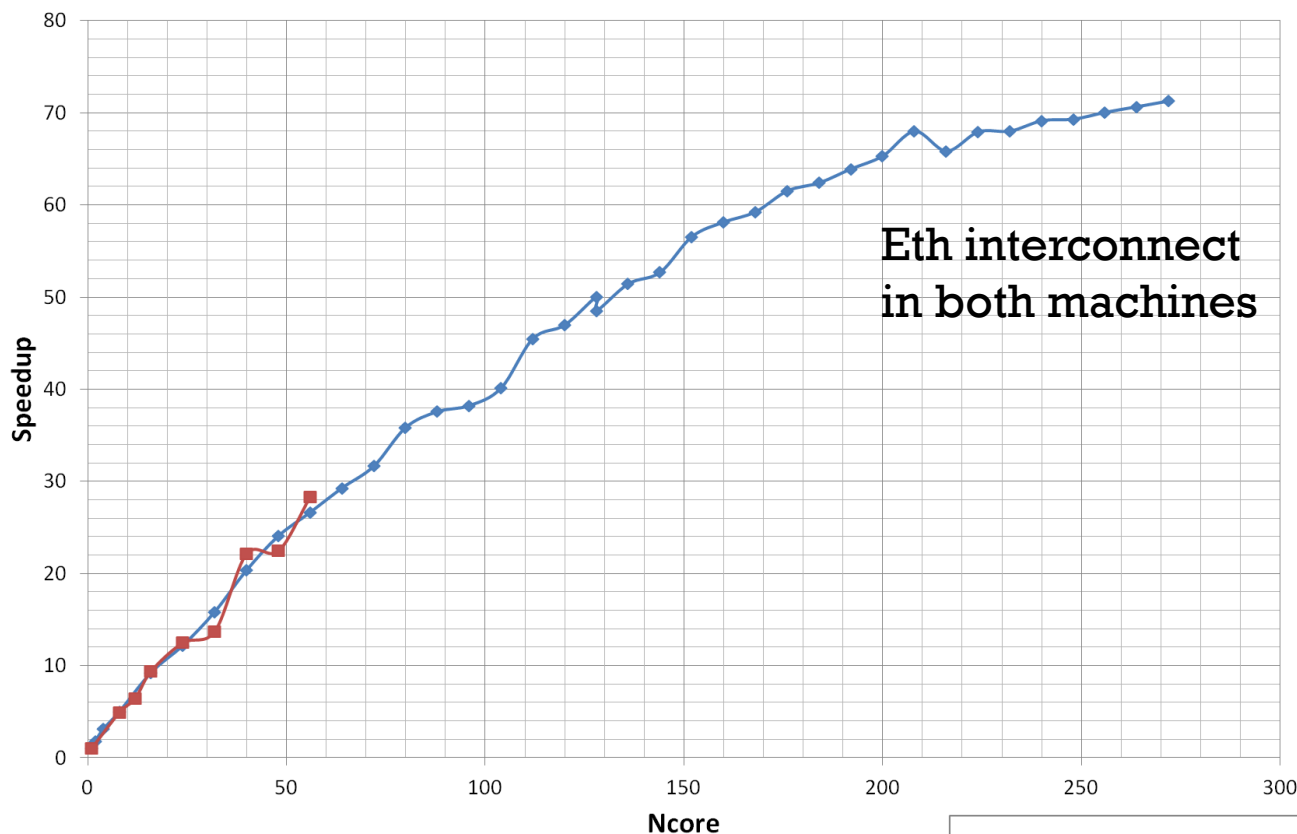
(data from M.Michelotto@HEPIX 2014

<https://indico.cern.ch/event/320819/session/3/contribution/30/material/slides/0.pptx> )



# Test on Intel AVOTON

### NAMD APOA1 Test

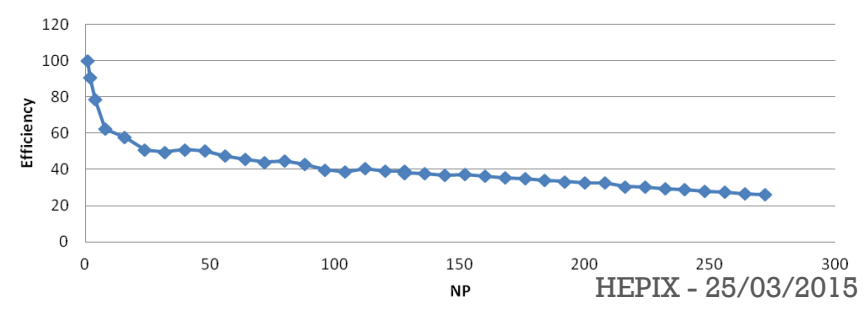


Theoretical & Computational  
Beckman Institute  
University of Illinois at Urbana-Champaign

- Avoton C2730@1.7GHz
- Xeon E5410@2.33GHz

N.B. – Comparison with an old Penryn Xeon CPU

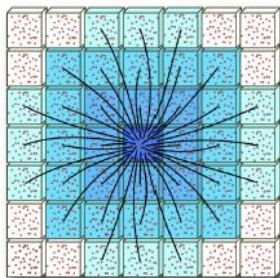
### NAMD APOA1 efficiency@avoton C2730 @1.7GHz



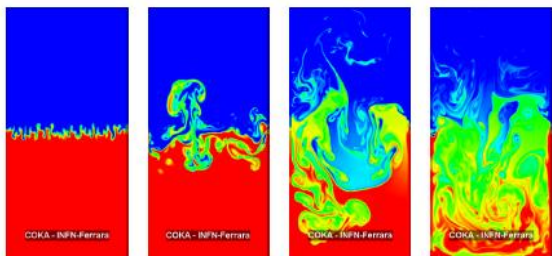
# Lattice Boltzmann on the Tegra K1

GPU only

## Lattice Boltzmann Methods: D2Q37



(\*) Schifano et al.  
*A portable OpenCL  
 Lattice Boltzmann code  
 for multi- And many-  
 core processor  
 architectures, Proc.  
 Comp. Sci. 29, 2014*



## Performance comparison with K20

Propag. (MLUPS)		Collide (MLUPS)	
Tegra	K20	Tegra	K20
17.8	256.0	1.6	89.4

...moving towards energy aware metrics...

Block Size	Propag. [A]	Collide [A]	Propag. [ms]	Collide [ms]
256	0.46	0.36	<b>29.1</b>	<b>340</b>
128	0.46	<b>0.30</b>	29.3	411
64	0.43	0.31	29.4	398
32	<b>0.29</b>	0.31	34.8	398

