# QE-GPU: between performance, correctness and sustainability.

Filippo SPIGA[1,2] <fs395@cam.ac.uk>

[1] High Performance Computing Service, University of Cambridge
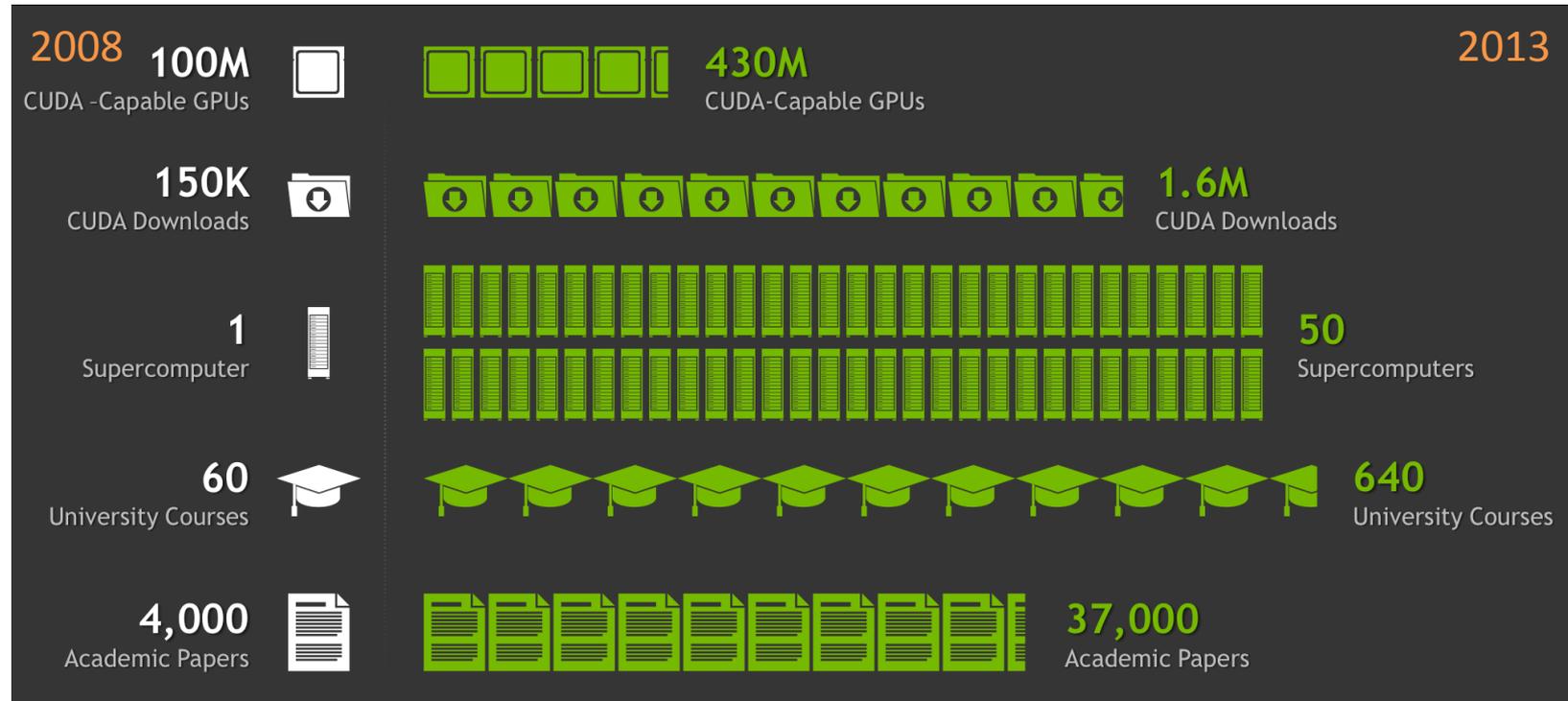
[2] Quantum ESPRESSO Foundation

*«What I cannot compute, I do not understand.»* (adapted from Richard P. Feynman)
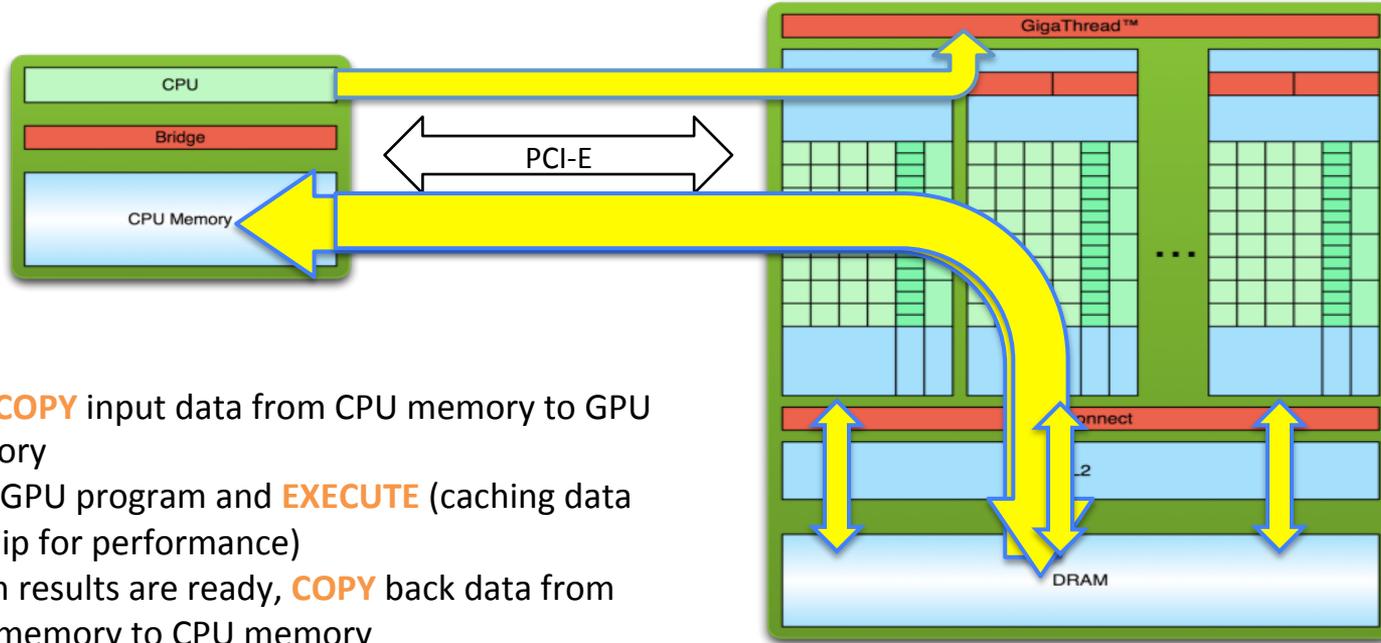
# Outline

- Overview

- The Quantum ESPRESSO project

- PWSCF GPU porting

- Performance evaluation

- Final considerations: performance, correctness and sustainability
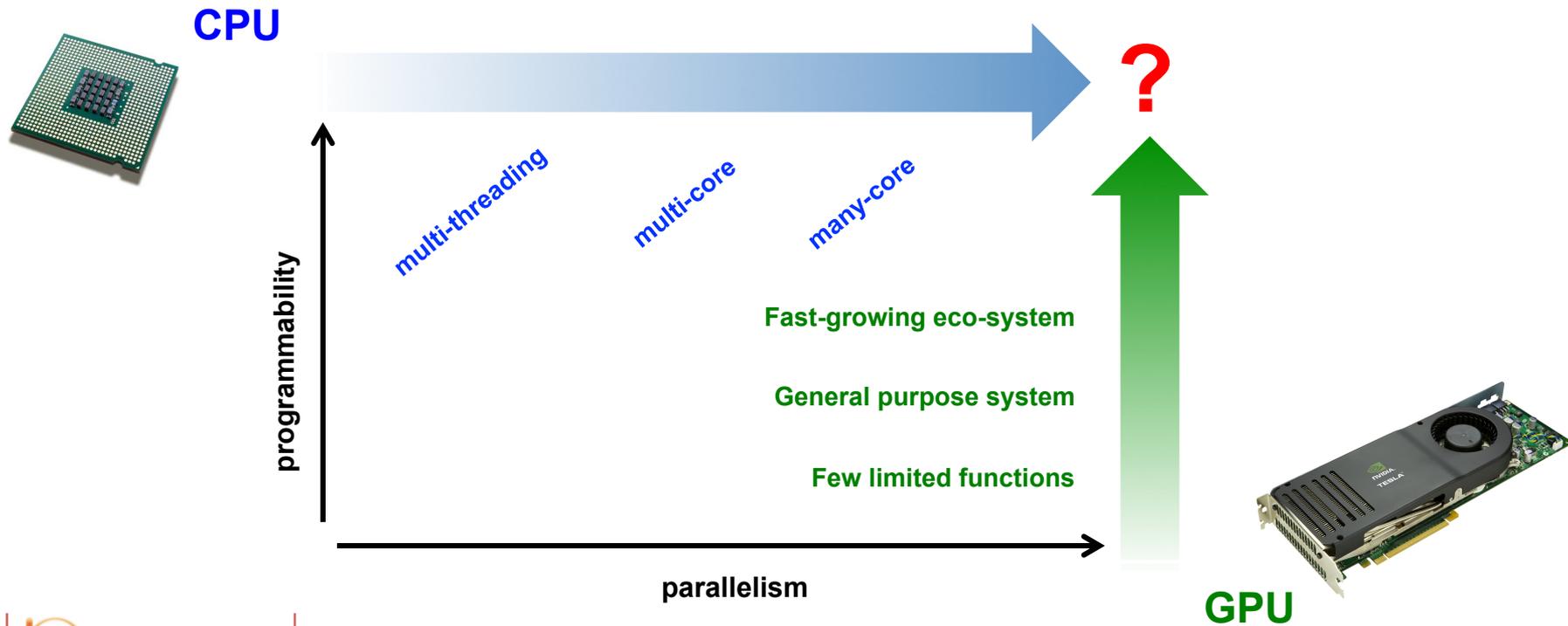
# GPU momentum continues to grow
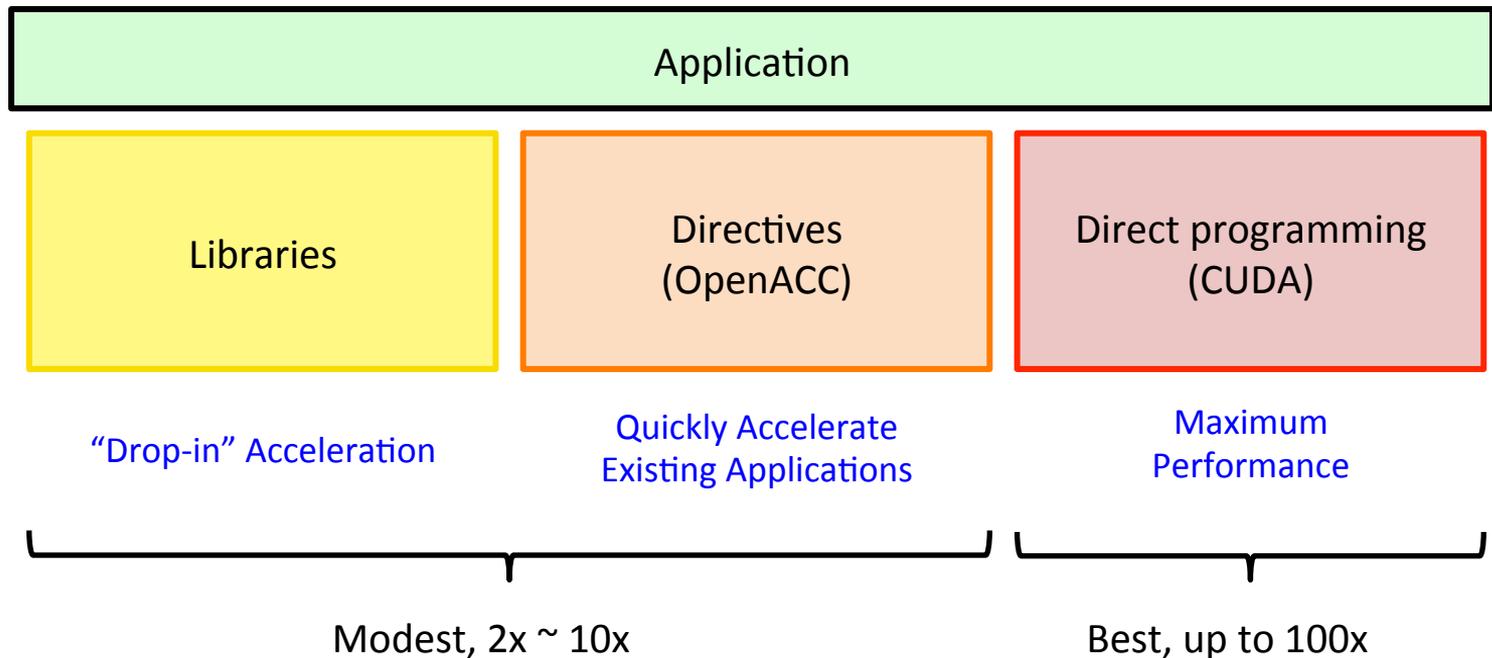
# The *Accelerator* model



1. First **COPY** input data from CPU memory to GPU memory
2. Load GPU program and **EXECUTE** (caching data on chip for performance)
3. When results are ready, **COPY** back data from GPU memory to CPU memory

# Collision or Convergence?



CPU

multi-threading     multi-core     many-core

?

programmability

Fast-growing eco-system

General purpose system

Few limited functions

parallelism

GPU

QUANTUM ESPRESSO FOUNDATION

# 3 ways to accelerate codes using GPU

# QUANTUM ESPRESSO

# What is QUANTUM ESPRESSO?

- QUANTUM ESPRESSO is an integrated suite of computer codes for atomistic simulations based on DFT, pseudo-potentials, and plane waves

- "ESPRESSO" stands for opEn Source Package for Research in Electronic Structure, Simulation, and Optimization

- QUANTUM ESPRESSO is an initiative of SISSA, EPFL, and ICTP, with many partners in Europe and worldwide

- QUANTUM ESPRESSO is free software that can be freely downloaded. Everybody is free to use it and welcome to contribute to its development
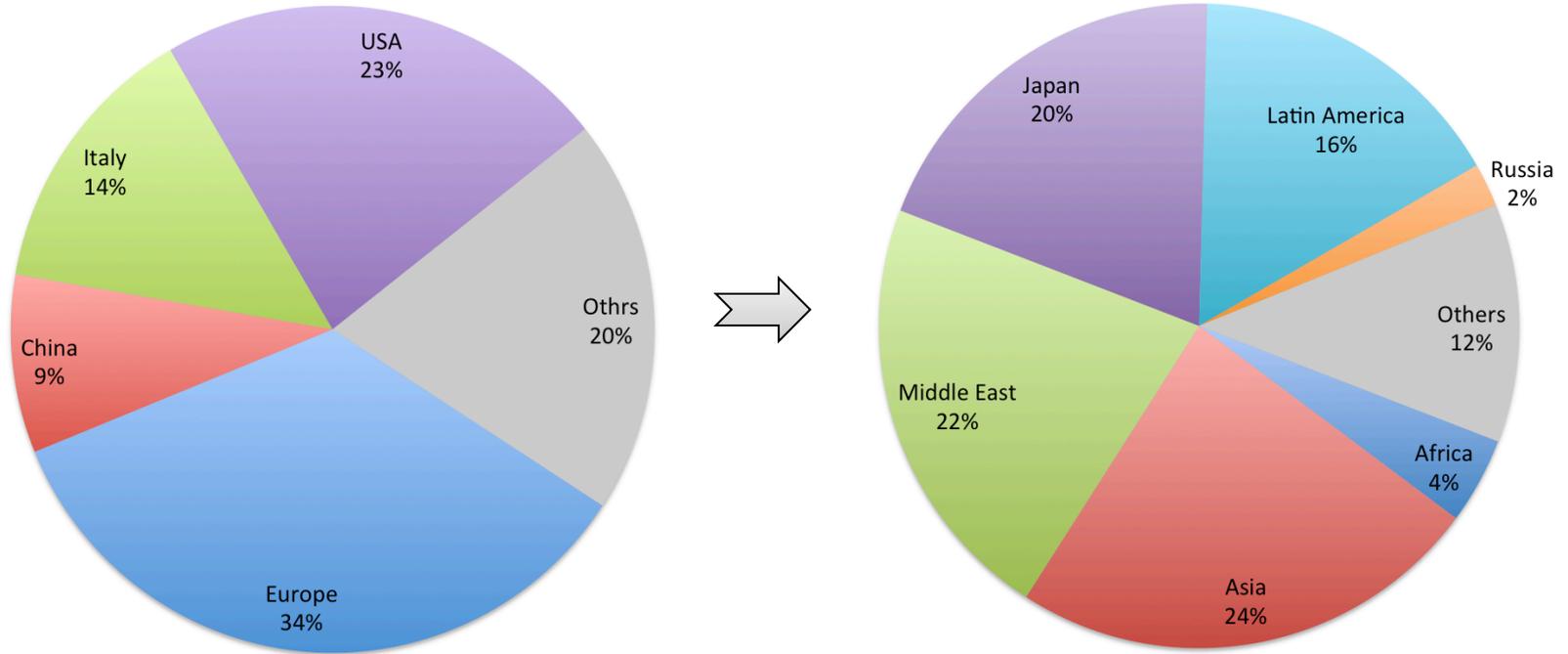
# What can QUANTUM ESPRESSO do?

- norm-conserving as well as ultra-soft and PAW pseudo-potentials
- many different energy functionals, including meta-GGA, DFT+U, and hybrids (van der Waals soon to be available)
- scalar-relativistic as well as fully relativistic (spin-orbit) calculations
- magnetic systems, including non-collinear magnetism
- **ground-state calculations**
  - Kohn-Sham orbitals and energies, total energies and atomic forces
  - finite as well as infinite system
  - any crystal structure or supercell
  - insulators and metals (different schemes of BZ integration)
  - structural optimization (many minimization schemes available)
  - transition states and minimum-energy paths (via NEB or string dynamics) electronic polarization via Berry's phase
  - finite electric fields via saw-tooth potential or electric enthalpy
- **Wannier intepolations**

- **ab-initio molecular dynamics**
  - Car-Parrinello (many ensembles and flavors)
  - Born-Oppenheimer (many ensembles and flavors)
  - QM-MM (interface with LAMMPS)
- **linear response and vibrational dynamics**
  - phonon dispersions, real-space interatomic force constants
  - electron-phonon interactions and superconductivity effective charges and dielectric tensors
  - third-order an-harmonicities and phonon lifetimes
  - infrared and (off-resonance) Raman cross sections
  - thermal properties via the quasi-harmonic approximation
- **electronic excited states**
  - TDDFT for very large systems (both real-time and "turbo-Lanczos")
  - MBPT for very large systems (GW, BSE)

*plus several post processing tools!*

QUANTUM ESPRESSO FOUNDATION

# QUANTUM ESPRESSO in numbers

- 350,000+ lines of FORTRAN/C code

- 46 registered developers

- 1600+ registered users

- 5700+ downloads of the latest 5.x.x version

- 2 web-sites (QUANTUM-ESPRESSO.ORG & QE-FORGE.ORG)

- 1 user mailing-list, 1 developer mailing-list

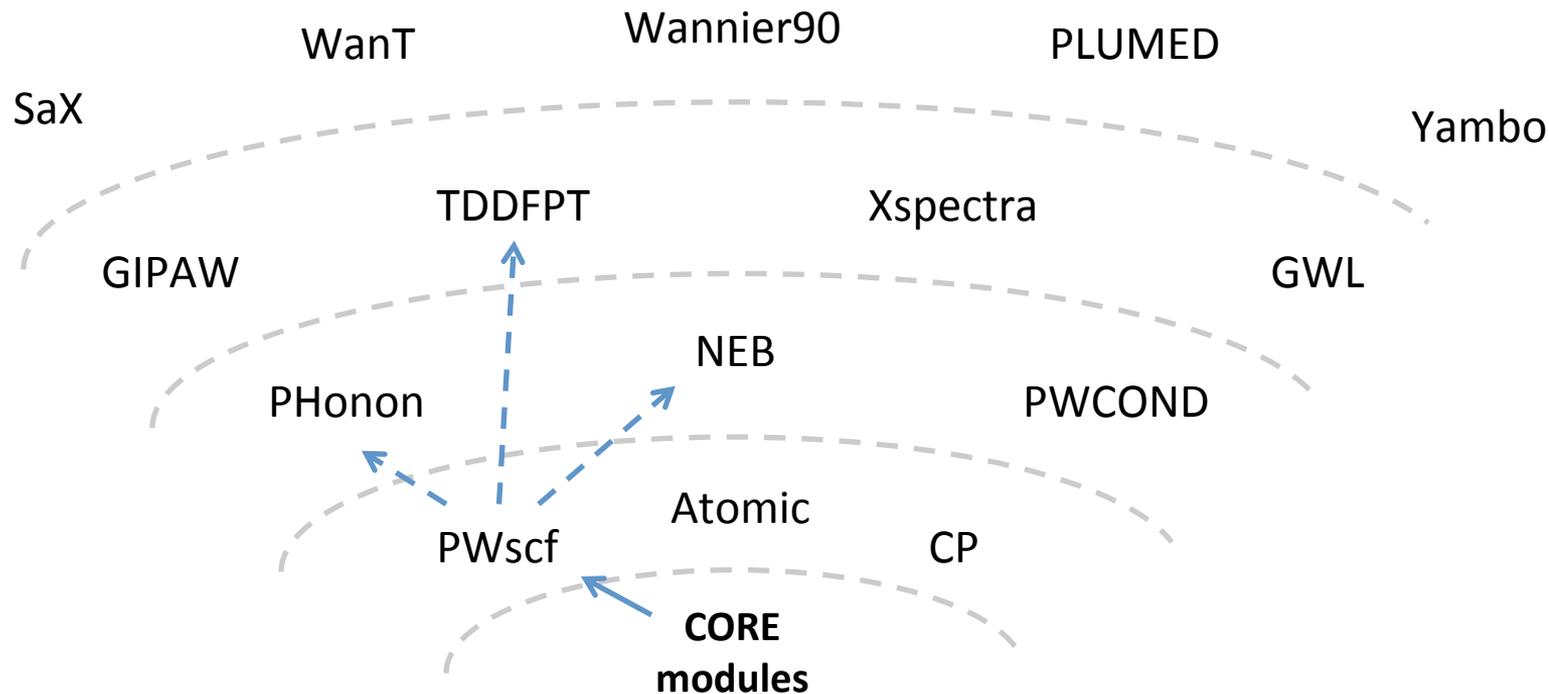- 24 international schools and training courses (1000+ participants)

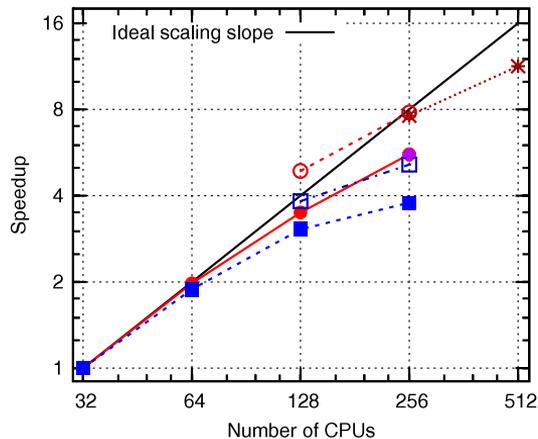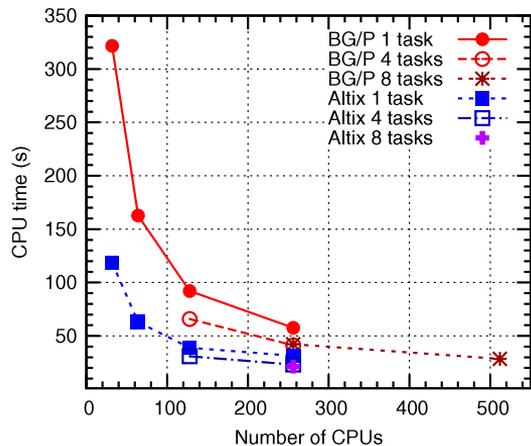# QUANTUM ESPRESSO: a global community

# The development model

- QUANTUM ESPRESSO is **not a monolithic application**, but an integrated ecosystem thriving around a small number of core components developed and maintained by a small number of developers

- the ecosystem is designed so as to be alien-friendly: a number of **third-party QE-compatible applications and add-ons**, often designed to be code-agnostic, are distributed with QE (notable examples include wannier90, Yambo, EPW, WanT, XCrysDen, ...)

- the environment that allows the ecosystem to prosper is provided by the QE-FORGE.ORG platform, **freely available** to researchers and developers from all over the world

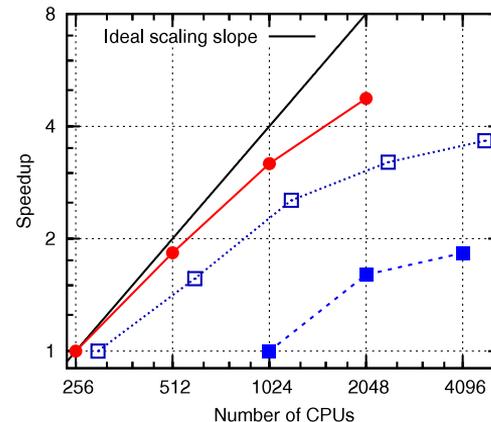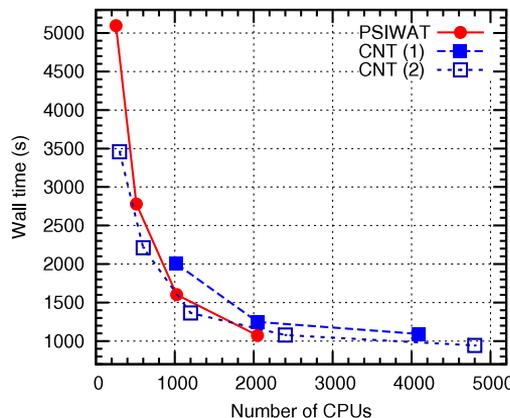# Quantum ESPRESSO package portfolio

# QE Scalability
## (2010)

scalability for < 1000 processors

CP

(*A*β-peptide in water: 838 atoms and 2311 electrons, gamma point)

PWscf

(PSIWAT: 587 atoms, 2552 electrons, 4 k-points)
(CNT: 1532 atoms, 5232 electrons, gamma point)

Scalability for > 1000 processors

# PWSCF GPU PORTING

# Plane Wave Self-Consistency Field (PWSCF)

The solution of the Kohn-Sham requires the diagonalization of the matrix $H_{KS}$ whose matrix elements are

$$\langle \mathbf{k} + \mathbf{G} | \mathbf{T} | \mathbf{k} + \mathbf{G}' \rangle = \frac{\hbar^2}{2m} (\mathbf{k} + \mathbf{G}')^2 \delta_{\mathbf{G},\mathbf{G}'}$$    Kinetic energy
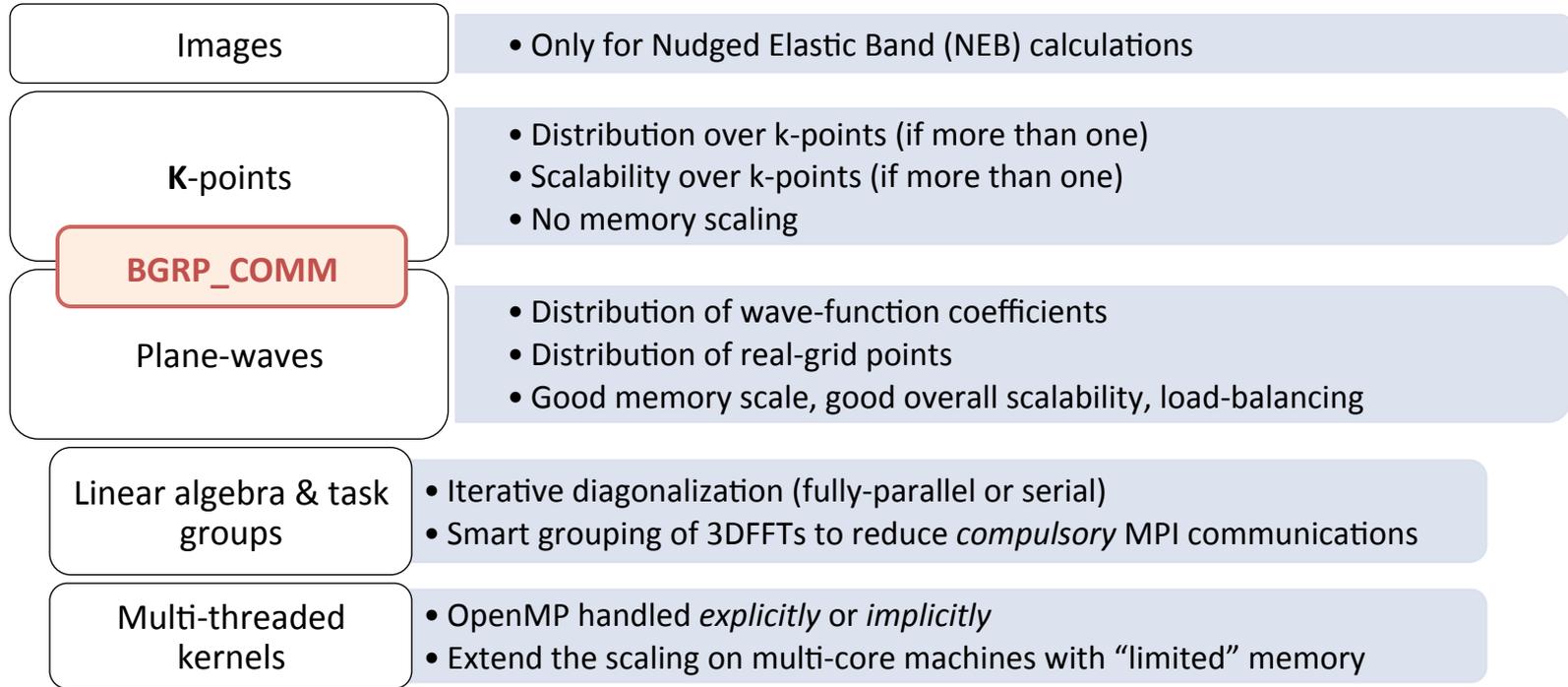
$$\langle \mathbf{k} + \mathbf{G} | V_H | \mathbf{k} + \mathbf{G}' \rangle = V_H(\mathbf{G} - \mathbf{G}') = 4\pi e^2 \frac{n(\mathbf{G} - \mathbf{G}')}{|(\mathbf{G} - \mathbf{G}')|^2}$$    Hartree term
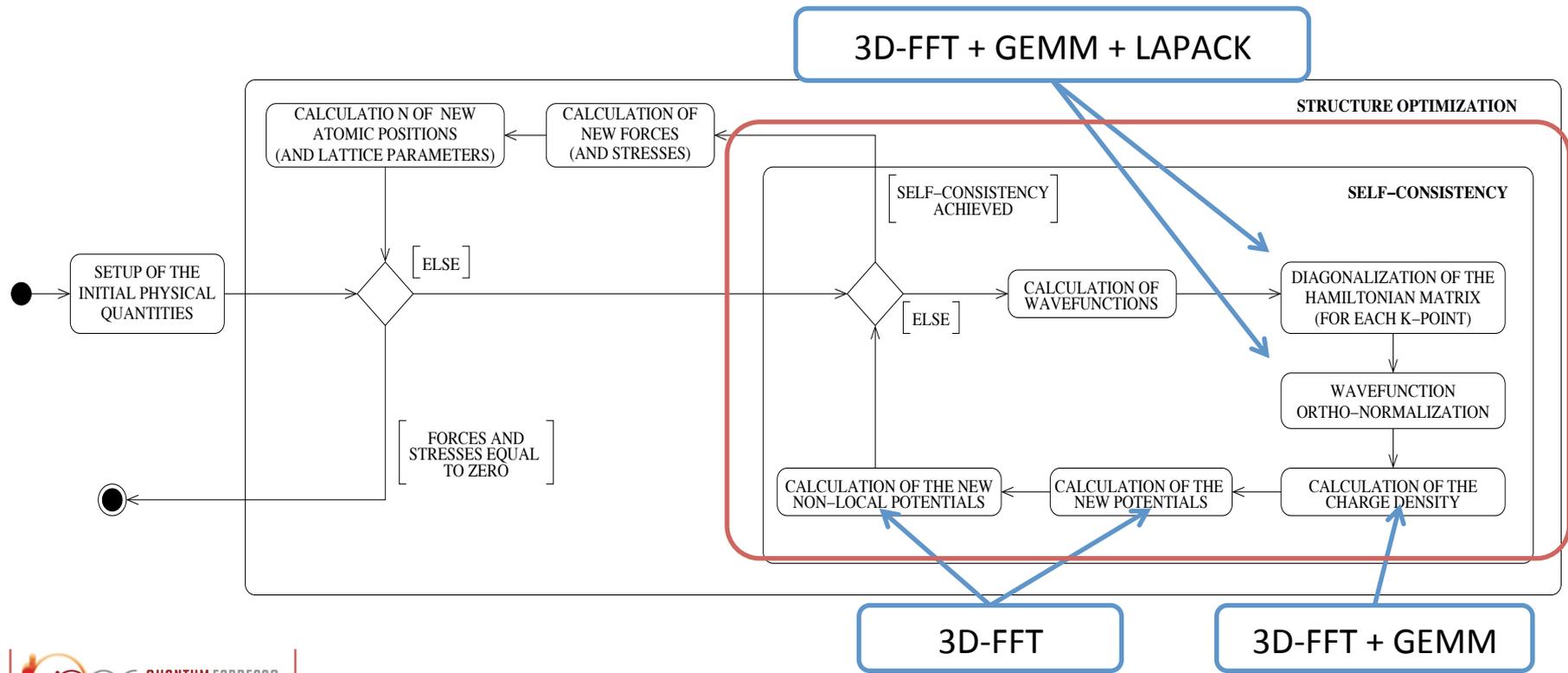
$$\langle \mathbf{k} + \mathbf{G} | V_{xc} | \mathbf{k} + \mathbf{G}' \rangle = FT[V_{xc}(r)]$$    Exchange correlation

# PWscf arallelization hierarchy

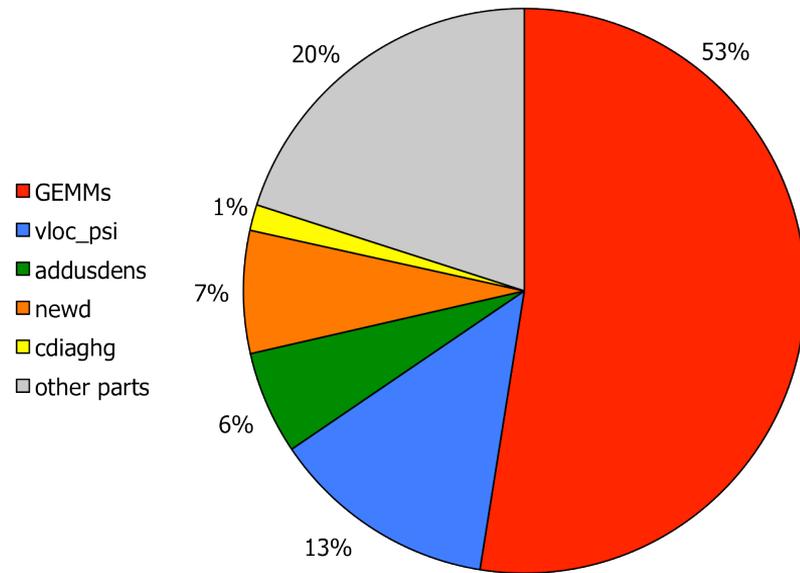| Images | • Only for Nudged Elastic Band (NEB) calculations |

| **K**-points | • Distribution over k-points (if more than one)<br>• Scalability over k-points (if more than one)<br>• No memory scaling |

**BGRP_COMM**

| Plane-waves | • Distribution of wave-function coefficients<br>• Distribution of real-grid points<br>• Good memory scale, good overall scalability, load-balancing |

| Linear algebra & task groups | • Iterative diagonalization (fully-parallel or serial)<br>• Smart grouping of 3DFFTs to reduce *compulsory* MPI communications |

| Multi-threaded kernels | • OpenMP handled *explicitly* or *implicitly*<br>• Extend the scaling on multi-core machines with "limited" memory |

QUANTUM ESPRESSO FOUNDATION

# Simplified PWSCF life-cycle



3D-FFT + GEMM + LAPACK

STRUCTURE OPTIMIZATION

SELF–CONSISTENCY

SETUP OF THE INITIAL PHYSICAL QUANTITIES

CALCULATIO N OF  NEW ATOMIC POSITIONS (AND LATTICE PARAMETERS)

CALCULATION OF NEW FORCES (AND STRESSES)

SELF–CONSISTENCY ACHIEVED

ELSE

ELSE

FORCES AND STRESSES EQUAL TO ZERO

CALCULATION OF WAVEFUNCTIONS

DIAGONALIZATION OF THE HAMILTONIAN MATRIX (FOR EACH K–POINT)

WAVEFUNCTION ORTHO–NORMALIZATION

CALCULATION OF THE NEW NON–LOCAL POTENTIALS

CALCULATION OF THE NEW POTENTIALS

CALCULATION OF THE CHARGE DENSITY

3D-FFT

3D-FFT + GEMM

# QE-GPU and PRACE

## PRACE - *Partnership for Advanced Computing in Europe*

- PRACE Project 1-IP (2010 – 2012)
  - GPU development: PHIGEMM library and PWSCF
  - Key-rule to leverage Quantum ESPRESSO as an EU community code
  - (Better) Parallelization of the GIPAW code (over bands)

- PRACE Project 2-IP (2011 – 2013)
  - Extend the multi-threading support with OpenMP
  - Exploratory of the adoption of OpenACC (GPU with directives)
  - Improvement in the linear algebra and the diagonalization

# A "lucky" starting point

## AUSURF, 112 Au atoms, 2 k-point



Legend:
- GEMMs (red)
- vloc_psi (blue)
- addusdens (green)
- newd (orange)
- cdiaghg (yellow)
- other parts (gray)

1 SCF iteration: GEMMs 53%, vloc_psi 13%, addusdens 6%, newd 7%, cdiaghg 1%, other parts 20%

1 full SCF cycle: GEMMs 26%, vloc_psi 19%, addusdens 9%, newd 8%, cdiaghg 10%, other parts 28%

# GPU developments

- MPI-GPU binding & GPU memory management
- NEWD → CUDA NEWD (multiple kernels combined)
- ADDUSDENS → CUDA ADDUSDENS
- VLOC_PSI → CUDA VLOC_PSI (CUDA kernels + CUFFT)
- BLAS3 *GEMM → PHIGEMM library (CUBLAS)
- (serial) LAPACK → MAGMA library



VLOC_PSI acts over distributed data
NEWD/ADDUSDENS act over local data

# phiGEMM

- Inspired by M. Fatica LINPACK work

- Independent open-source library, BSD license

- GPU+CPU BLAS 3 *GEMM routine

- Manual or "semi-automatic" (SELF-TUNE) split

- Special-K for rectangular matrices

- GEMM→GEMV fallback

- Detailed call-by-call profiling

- Pinned/non-pinned, sync/async

- Support of multi-GPU

web: http://qe-forge.org/projects/phigemm/

# MAGMA: LAPACK for GPU



MAGMA uses **<u>HYBRIDIZATION</u>** methodology based on
- Representing linear algebra algorithms as collections of TASKS and DATA DEPENDENCIES among them
- Properly SCHEDULING the tasks' execution over the multicore and the GPU hardware components

*What does HYBRIDIZATION means?*
- Panels (Level 2 BLAS) are factored on CPU using LAPACK
- Trailing matrix updates (Level 3 BLAS) are done on the GPU using "look-ahead"

# CUDA kernels

- ADDUSDENS
  - compute-bounded kernel
  - Best performance measured: 20x* (Realistic? 9x~10x)
- NEWD
  - compute-bounded kernels
  - Best performance measured: 7.2x* (Realistic? 3x~4x)
- VLOC_PSI
  - memory-bounded kernels
  - Best (serial) performance measured: 9x (Realistic? …)
- All the data is moved to GPU memory only at once
- External loops over atomic species are kept on the CPU side

# The parallel FFT "issue"

There are two "FFT grid" representation in Reciprocal Space: wave functions ($E_{cut}$) and charge density ($4E_{cut}$)

A single 3D-FFT is divided in independent 1D-FFTs

Transform along Z

$E_c$

$4E_c$

y

z — x

~ Nx Ny / 5  FFT along z

Parallel Transpose ~ Nx Ny Nz / (5 Np) data exchanged per PE

Transform along Y

z

y — x

Nx Nz / 2  FFT along y

Transform along X

z

y — x

Ny Nz  FFT along x

0  PE 0     2  PE 2
1  PE 1     3  PE 3

Data are not contiguous and not "trivially" distributed across processors

Zeros are not transformed. Different cut-offs preserve accuracy

# H * psi

```
compute/update H * psi:
```
<span style="color:gray">compute kinetic and non-local term (in G space)</span>

$\qquad$ complexity : $N_i \times (N \times N_g + N_g \times N \times N_p)$

Loop over (not converged) bands:

$\qquad$ FFT (psi) to R space

$\qquad\qquad$ complexity : $N_i \times N_b \times FFT(N_r)$

$\qquad\qquad$ compute V * psi

$\qquad\qquad$ complexity : $N_i \times N_b \times N_r$

$\qquad\qquad$ FFT (V * psi) back to G space

$\qquad\qquad$ complexity : $N_i \times N_b \times FFT(N_r)$

$\qquad\qquad$ <span style="color:gray">compute Vexx:</span>

$\qquad\qquad$ complexity : $N_i \times N_c \times N_q \times N_b \times (5 \times N_r + 2 \times FFT(N_r))$

$N = 2 \times N_b$ (where $N_b$ = number of valence bands) $\qquad$ $N_p$ = number of PP projector

$N_g$ = number of G vectors $\qquad\qquad\qquad\qquad\qquad\qquad$ $N_r$ = size of the 3D FFT grid

$N_i$ = number of Davidson iteration $\qquad\qquad\qquad\qquad$ $N_q$ = number of q-point (may be different from $N_k$)

QUANTUM ESPRESSO FOUNDATION

# Task-group parallelization



```
do i = 1, n
   compute 3D FFT( psi(i) )
end do
```

```
do i = 1, n/groupsize
  merge( psi( i ), psi( i + 1 ) )
  compute groupsize 3D FFT
  (at the same time)
end do
```
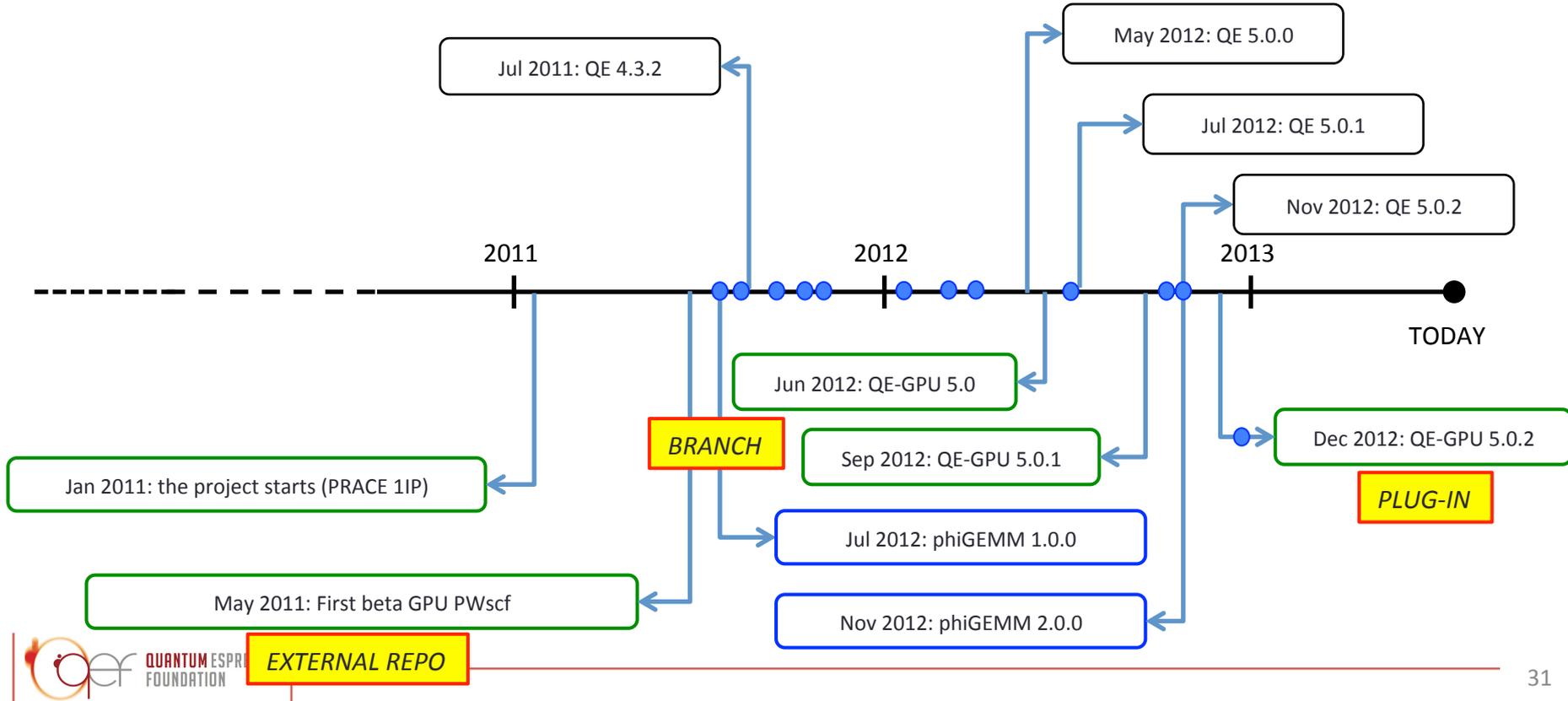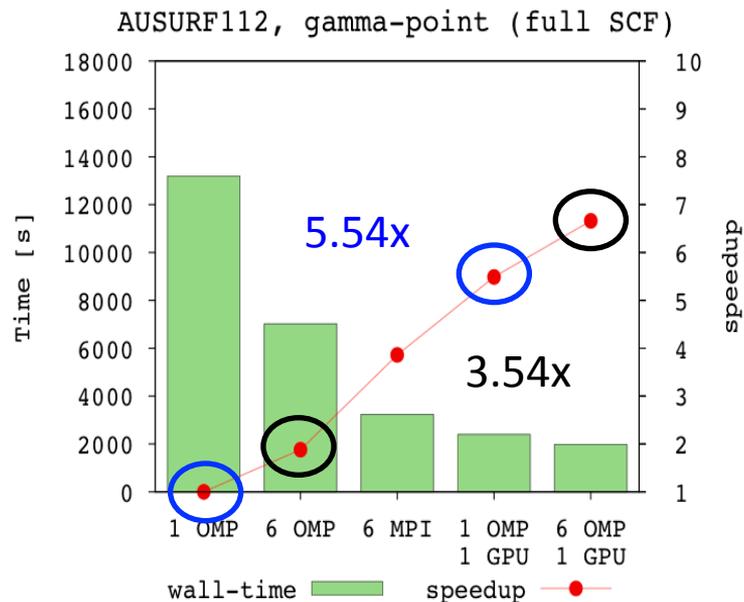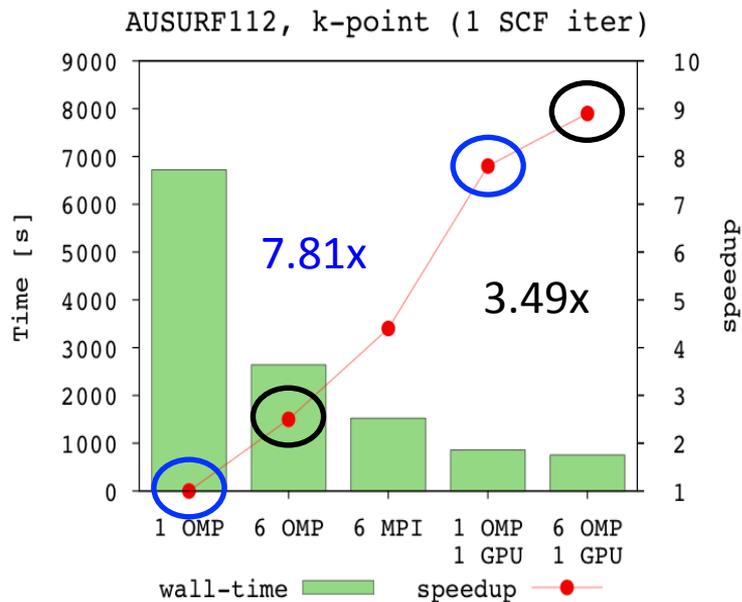
# Parallel VLOC_PSI

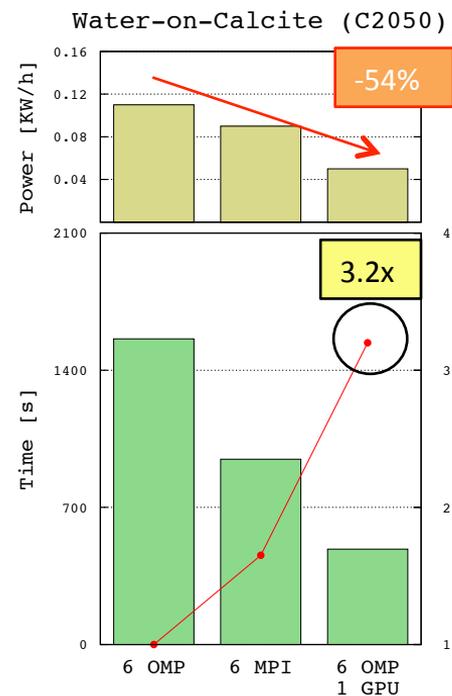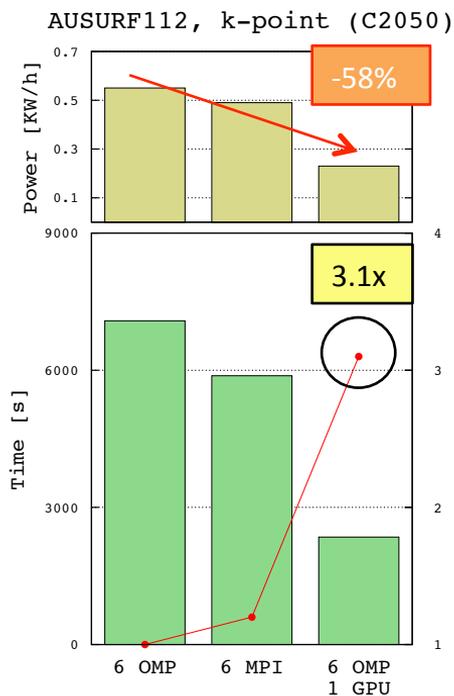# Parallel VLOC_PSI - Limitations

# QE-GPU

# QE-GPU Timeline ( ~2010 – today )
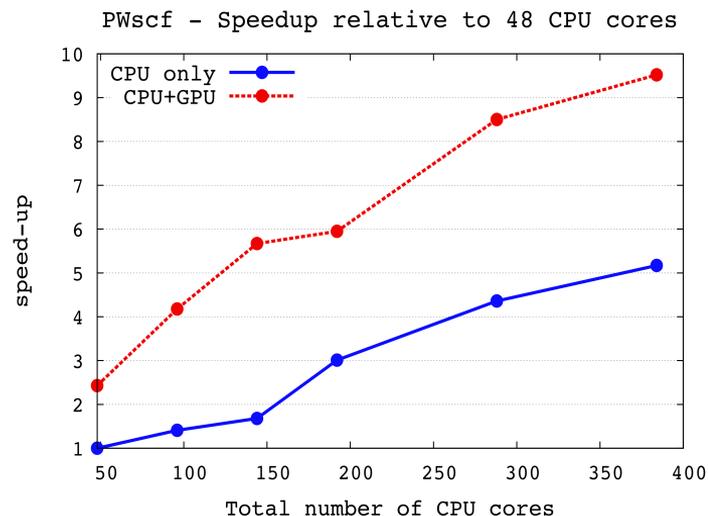
# PERFORMANCE EVALUATION

# AUSURF112, serial

*Tests run early 2012*

# Performance & Power consumption (serial)



Shilu-3 (C2050)

AUSURF112, k-point (C2050)
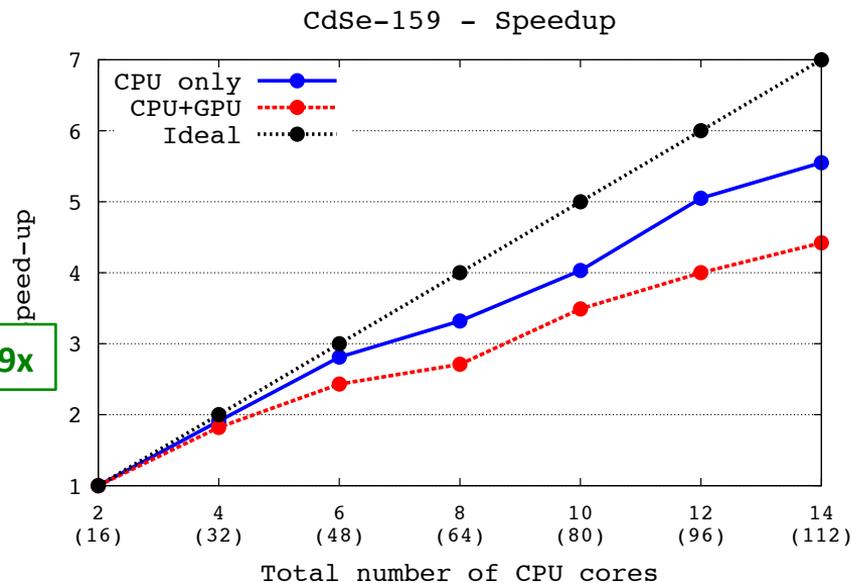
Water-on-Calcite (C2050)

*Tests run early 2012*

# MGST-hex

216 atoms of {Ge, Mn, Te, Sb}, gamma-only (courtesy of *Zhang W.* – RWTH/AACHEN)



PWscf - Walltime of 1 SCF cycle

2.43x

2.96x

3.38x



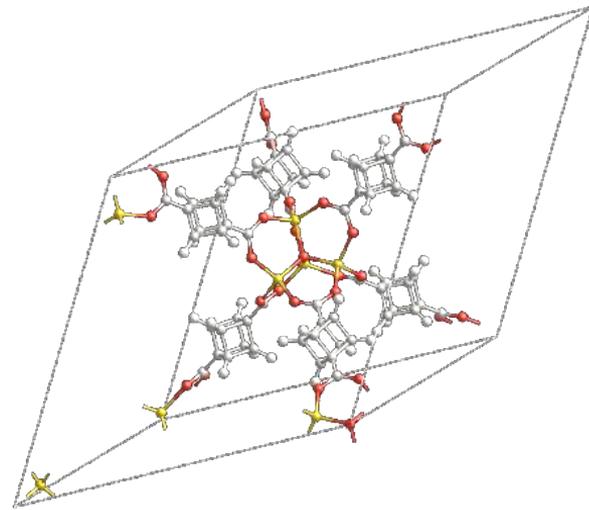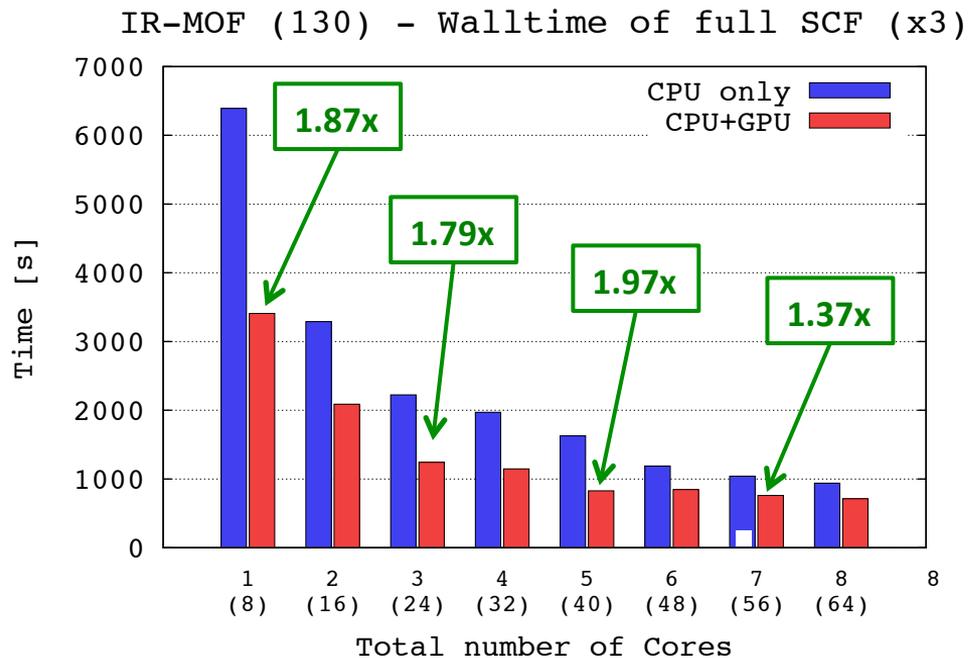PWscf - Speedup relative to 48 CPU cores

*Tests run early 2012*

# CdSe159

159 atoms of {Cd, Se}, gamma-only (courtesy of *Calzolari A.* – CNR/NANO)
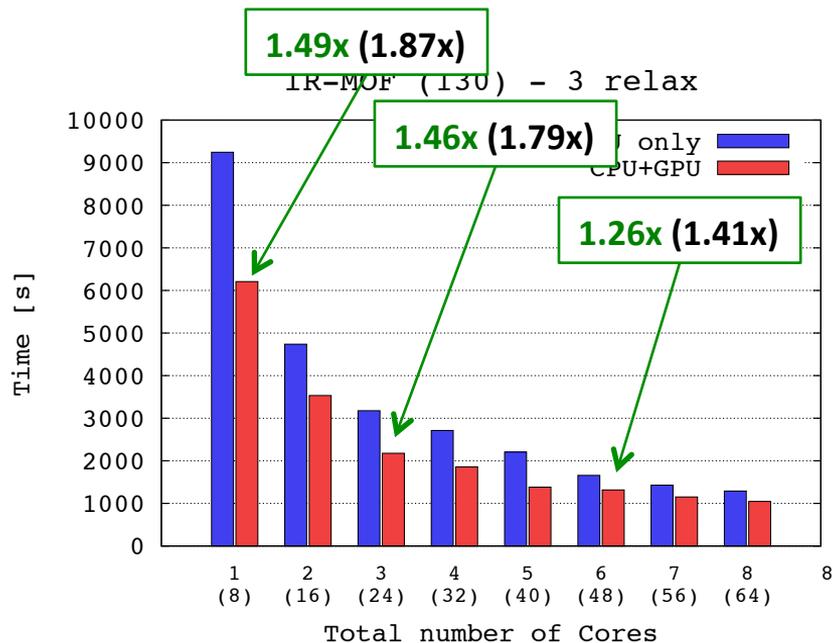
# IRMOF-M11

## 130 atoms of {O, Zn, C, H},  1 K-point (courtesy of *Clima S.* - IMEC)



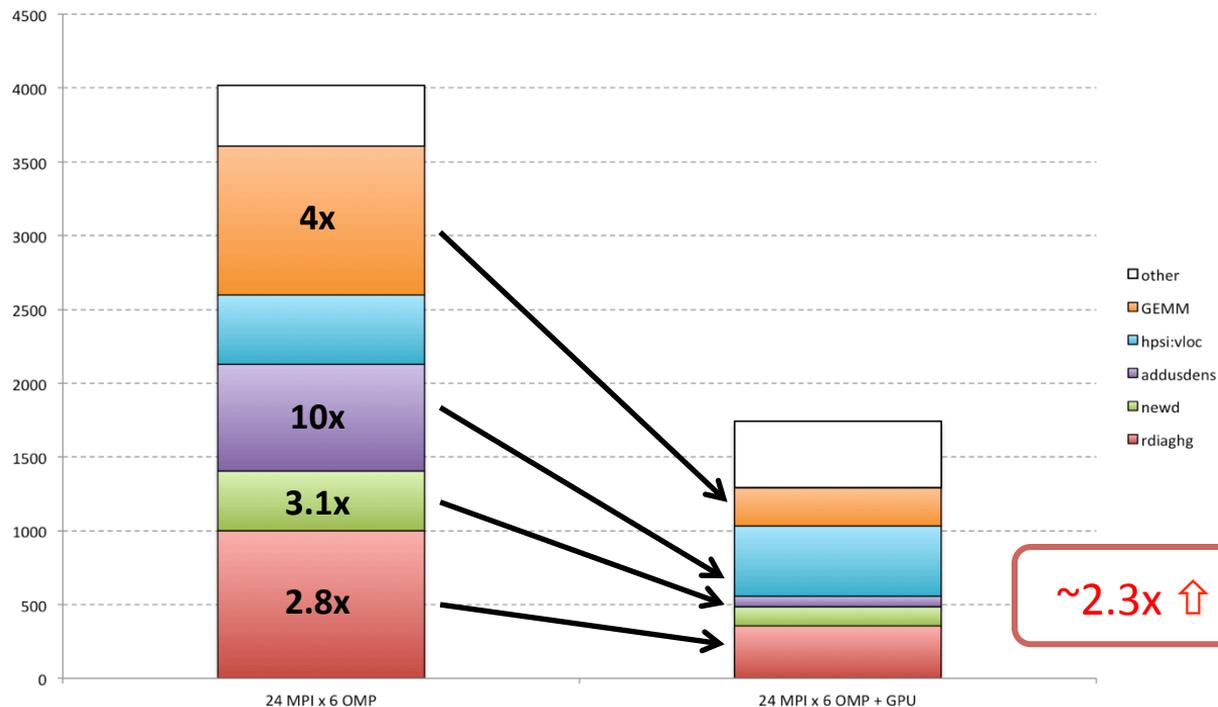*Tests run early 2012*

# Hitting the limit

130 atoms of {O, Zn, C, H},  1 K-point (courtesy of  *Clima S.* - IMEC)

*Tests run early 2012*
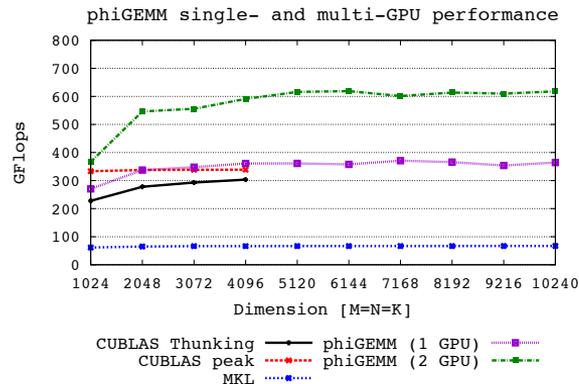
# CdSe489 - PRACE Preparatory Access

489 atoms of {Cd, Se} , 3 SCF steps, 24 MPI x 6 OMP (144 cores) @ PLX (CINECA)

# phiGEMM

- Inspired by M. Fatica LINPACK work
- Independent open-source library, BSD license
- GPU+CPU BLAS 3 *GEMM routine
- Manual or "semi-automatic" (SELF-TUNE) split
- Special-K for rectangular matrices
- GEMM→GEMV fallback
- Detailed call-by-call profiling
- Pinned/non-pinned, sync/async
- Support of multi-GPU

web: http://qe-forge.org/projects/phigemm/

# Hitting the limit

489 atoms of {Cd, Se}, 3(+1/2) SCF steps & STRESS & FORCES, 24 MPI x 6 OMP (144 cores) @ PLX (CINECA)



*Tests run early 2012*

# Cheating using parallelism: USE_3D_FFT

120 atoms of {Bi, Fe, O, LS, Mn}, reduced to 8 k-point (courtesy of *Ferreira R.* – Rio de Janeiro Federal Univ.)

# pool = # MPI processes
( dimension each pool = 1, k-point distributed in "round-robin")

| Computer Nodes | Execution Time [s]<br>#10 self-consistency cycles | Speed-up |
|---|---|---|
| 2 x iDataPlex DX360M3, dual Xeon E5645 6-cores 2.40 GHz (24 cores) | 52057.22 | |
| 2 x iDataPlex DX360M3, dual Xeon E5645 6-cores 2.40 GHz (24 cores) + **4 NVIDIA 2070 (__USE_3D_FFT)** | **10029.1** | **5.2x** |

*Tests run early 2012*

# We are all limited by Amdahl
## (parallel – MGST-hex)



Best CPU+GPU walltime, details [s]

Between 60%~70% of the total wall-time is portable on CUDA with different degree of acceleration

Legend:
- *GEMM
- addusdens
- newd
- rdiaghg
- hpsi:vloc

# FINAL CONSIDERATIONS
## *PERFORMANCE, CORRECTNESS AND SUSTAINABILITY*

# What drive (my) agenda?



Sustainability

Performance

Correctness

# Correctness vs Sustainability

- Bugs can be very tricky and not easy to detect
  - long runs → lot of output; big runs → lot of resources
  - time & reproducibility

- An embedded Unit Test suite is a solution but …

- … an Acceptance Testing Procedure is required too!

- Old-fashion debugging works but …
  - dirty source code
  - plenty of preprocessor macro

QUANTUM ESPRESSO FOUNDATION

# Performance vs Correctness

- GPU are not CPU (*what a surprise!*)

- Third-part libraries hide complexity but also introduce "noise"

- Reduction operations (explicit managed or implicitly inherit) are the nightmare

- Innovative algorithms can improve performances but who does a proper validation?

# Sustainability vs Performance

- CUDA or not-CUDA? (*Shakespearean dilemma*)

- Benchmarking requires time (*again, what a surprise!*)

- The need of certifying QE-GPU benchmarks*

- Guarantee performances is very difficult because...

  - users expect/pretend performance portability

  - too many HW combinations

# OpenACC: yes, no, maybe?

PROs:

- easy learning curve

- fast first deployment

- acceptable performance compromises

CONs:

- Code has to be rewritten to suite the accelerator

- Lack of direct control of the generated GPU code

- memory sharing and competition

and now OpenMP 4.0 RC2 with *Directives for attached accelerators*

# (pseudo-random) Personal thoughts

- Training users is challenging (we kept everything simple but...)

- User support is time consuming

- CUDA is easier than most people think

- Focusing on fine-grain performance ONLY at the final stage

- There is not only the parallel world

- No more old-fashion porting

# The *Optimization Dilemma*

?

Is it <u>worth</u> to spend time/effort optimizing the code for the latest available architecture by using the latest cool set of features provided by the most updated version of CUDA toolkit/driver?
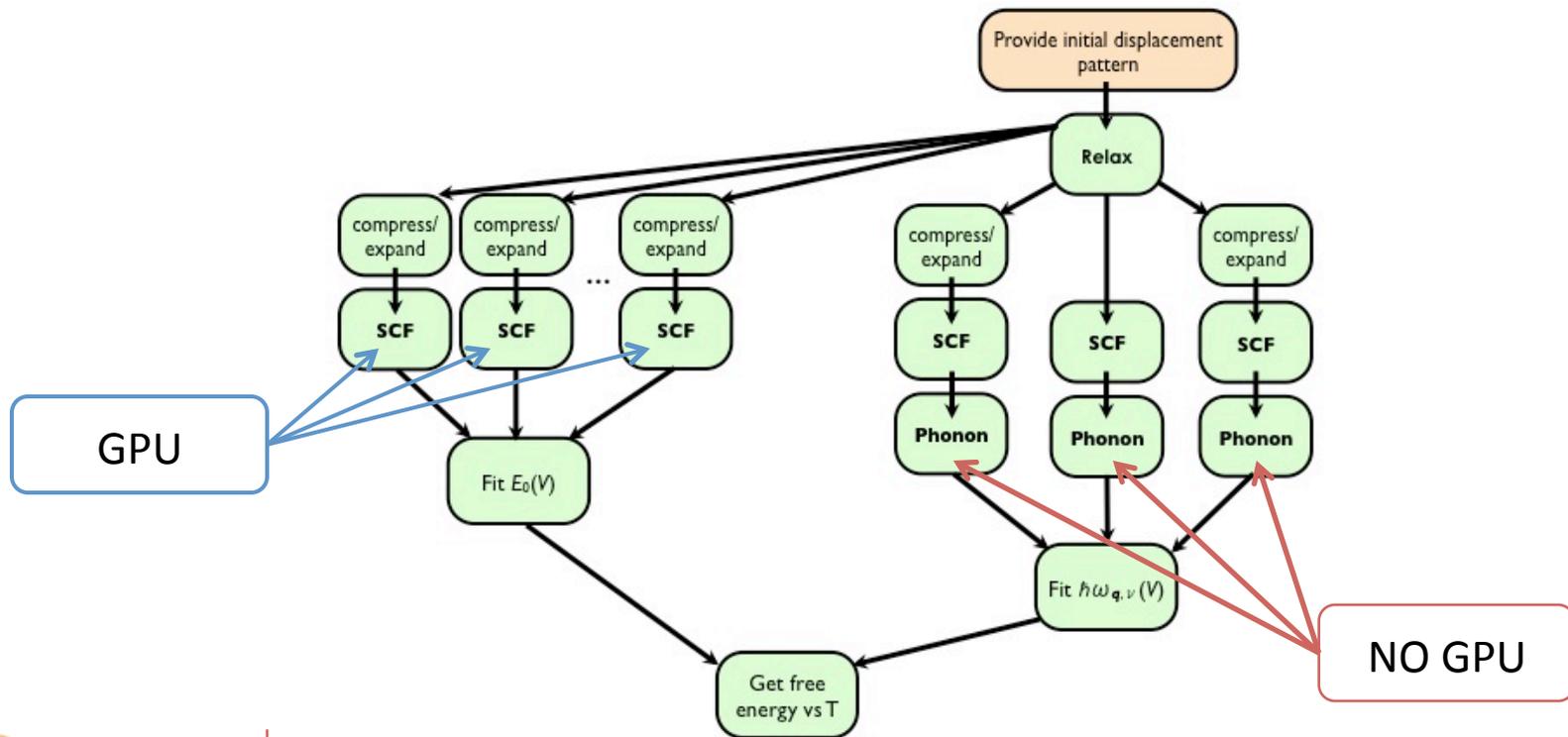
*Who are the users? Are they users or developer? What hardware they have? What is their scientific workflow? What are the scientific challenges they have?  What is their budget? What ROI they expect?*

# GPU-accelerated PWscf, when and why?

- Physical system criteria
  - "crazy" amount of atoms ( $O(1000)$ ) in gamma-only (or 1~2 k-point)
  - reasonable amount of atoms ($O(100)$ )
  - small amount of atoms ($O(10)$ ) but many k-points
  - $O(10){\sim}O(100)$ but with lot of atomic species

- Hardware criteria*
  - workstations equipped with commodity GTX and/or high-end TESLA GPUs
  - reasonable HPC clusters ($O(100)$ nodes, 1:1 ration GPU versus CPU socket )

- Work-flow criteria
  - long SCF energy minimization
  - high throughput

QUANTUM ESPRESSO FOUNDATION

# What scientists do?

# Next developments/priorities

TOP PRIORITY:

- Specific NVIDIA Kepler optimizations → NVIDIA Dev Tech

- non collinear & spin-orbital magnetization

- PHonon using OpenACC

"LESS" PRIORITY:

- first GPU-accelerated CP (special focus on DFORCE)

- Alternatives for current eigen-solvers (ELPA+GPU? Lancroz?)

and... DOCUMENTATION!

# CASTEP and GPU

CASTEP and QE shares...

- similar data distribution (by plane waves, by k-point, by bands, …)
- similar MPI communication patterns (MPI_alltoall, MPI_reduce)
- similar constrains in scalability

Where GPU...

- *Block Davidson solver with density mixing* → eigen-solvers, vloc_psi
- *Support of ultrasoft/norm-conserving pseudopotential* → newd & addusdens
- *Geometry optimization* → stress & forces calculations
- BLAS 3 operations → CUBLAS, phiGEMM

# THANK YOU FOR YOUR ATTENTION!

Links:

- http://www.quantum-espresso.org/
- http://foundation.quantum-espresso.org/
- http://qe-forge.org/gf/project/q-e/
- http://qe-forge.org/gf/project/q-e-gpu/