

## **Enabling faster material science modeling using the accelerated Quantum ESPRESSO**

### Filippo Spiga (ICHEC)

Ivan Girotto (ICHEC) Carlo Cavazzoni (CINECA)





Co-funded by the Irish Government and the European Union



EUROPEAN REGIONAL DEVELOPMENT FUND



OIDEACHAIS EDUCATION AGUS EOLAÍOCHTA A N D S C I E N C E



Higher Education Authority An tÚdarás um Ard-Oideachas



## Irish Centre for High-End Computing (ICHEC)

- Founded in mid-2005 by SFI and HEA
- ~25 staff member
- Two main systems:
  - SGI Altix ICE 8200EX (~4k cores)
  - Bull Novascale R422-E2, ~ 500 cores + 48 NVIDIA GPUs, National Service Production
- **Objectives:** 
  - Provide computational resources
  - Provide education and training to third-level institutions
  - Tech transfer and consultancy services to develop Irish smart economy





## Outline

- What is Quantum ESPRESSO
- Project Context, Objectives and Goals
- Implementation strategies
- Performance & Power Measurements
- Current Limitations & Best Practices
- Partnership & (Big) Challenges
- Future developments
- Q/A



## What is QUANTUM ESPRESSO?

QUANTUM ESPRESSO is...

- an integrated suite of computer codes for electronic-structure calculations and materials modeling at nano-scale
- based on density-functional theory (DFT), plane waves (PW) basis set and pseudo-potentials (PP)
- a DEMOCRITOS initiative, later joined by ICTP, CINECA Bologna, EPF Lausanne, Princeton University, MIT, Paris VI, Oxford, IJS Ljubljana, **ICHEC**, et al.
- composed of many packages: PWscF, CP, PHONON, ATOMIC, PWCOND, XSPECTRA, GIPAW, GWL, TDDFPT, WANT, ...
- able to run in serial and parallel on several architectures (Linux clusters, IBM systems, CRAY, NEC)
- distributed under the GNU General Public License (GPL)
- supported by a worldwide community of developers and users



## **QUANTUM ESPRESSO in numbers...**

- ~300,000 FORTRAN90/C lines of code (in total, including examples and docs, ~500,000)
- ~300 citations at 2010Q1 (>1,000 today)
- ~1,300 subscribers to the mailing-lists
- ~4,500 mails/year (2010/2011 traffic)
  - new mailing list dedicated to the GPU PWscf, 56 member registered
- 4,000~4,500 downloads from the website per version (SVN? who knows...)
- ~40 active projects on QE-forge (including PHIGEMM and QE-GPU!)
- 252 QE-forge users
- 115,633 QE-forge pages seen in 2010
- 20 international schools/events all over the world



(trends evolve very fast... the number always increasing!)



## **Project Context, Objectives and Goals**

### Who:

- the Irish Centre for High-End Computing within EC-funded PRACE Partnership For Advanced Computing in Europe, 1<sup>st</sup> implementation phase project (FP7/2007-2013 under grant RI-261557) and SFI - Science Foundation Ireland (grant 08/HEC/I1450).
- CINECA & DEMOCRITOS, as technical and consultancy partners

### What:

- accelerate the Plane Wave Self-Consistency Field (PWSCF) code exploiting the NVIDIA GPU capabilities
- target both serial and parallel version, assessing the numerical accuracy and the overall performance
- GPU code maintenance, support to the growing community and package dissemination

### Why:

- PWSCF is one of the most used package of the Quantum ESPRESSO suite. SCF calculations represent the starting point of other type of calculations (PHONON, GIPAW, GWL,...)
- Quantum ESPRESSO is recognized at European level as *community code* and PWscF is part of the PRACE Official Benchmark suite



### **Preliminaries: Kohn-Sham with plane waves**

The solution of the Kohn-Sham equation set requires the diagonalization of the matrix H<sub>KS</sub> whose matrix elements are

Kinetic energy 
$$\langle \mathbf{k} + \mathbf{G} | \mathbf{T} | \mathbf{k} + \mathbf{G} \rangle = \frac{\hbar^2}{2m} (\mathbf{k} + \mathbf{G} )^2 \mathcal{O}_{\mathbf{G},\mathbf{G}^{\sharp}}$$
  
Hartree term  $\langle \mathbf{k} + \mathbf{G} | V_H | \mathbf{k} + \mathbf{G} \rangle = V_H (\mathbf{G} - \mathbf{G}^{\sharp}) = 4\rho e^2 \frac{n(\mathbf{G} - \mathbf{G}^{\sharp})}{|(\mathbf{G} - \mathbf{G}^{\sharp})|^2}$   
Exchange correlation  $\langle \mathbf{k} + \mathbf{G} | V_{xc} | \mathbf{k} + \mathbf{G}^{\sharp} \rangle = FT \notin V_{xc} (\mathbf{r})$ 

FFT is used (or "abused") to move from real to reciprocal space (and vice-versa) in order to "simplify" some operations



### **Simplified PWscF life-cycle**



A <u>REAL</u> simulation for scientific purpose is usually composed by several SCF steps in a global structure optimization loop...



### **Time-consuming steps in PWSCF**

- Calculation of Charge Density
  - FFT (full 3D or set of 1D)
  - matrix-matrix multiplications
- Calculation of Potential
  - FFT (full 3D or set of 1D)
  - operations on real-space grid
- Davidson Iterative Diagonalization (SCF)
  - eigenvalues/eigenvectors
  - FFT (full 3D or set of 1D)
  - matrix-matrix multiplications



Basically most CPU time spent in linear-algebra operations, it is implemented among BLAS, LAPACK and FFT libraries!



## **Starting point: AUSURF112**

<u>Physical description</u>: gold slab surface of 112 Au atoms <u>Technical description</u>: both gamma- and k-point version, can run in serial and parallel. No high scalable. <u>Reference</u>: subset of a big test case (PSIWAT), provided by Arrigo Calzolari (CNR-NANO)





### **Development strategy: ADDUSDENS and NEWD**

- Both routines are more compute-bounded than memory-bounded
  - we achieved up to 19~20-times\* and 7~8-times acceleration
- All the data is moved to GPU memory only at once
  - some data structures do not change during the computation  $\rightarrow$  preload
- External loops over atomic species are kept on the CPU side
- The QVAN2 (computation of the Fourier transformation of Q functions) is not yet CUDA
  - it is memory consuming and we are investigating how to split it
  - if the implementation is not efficient, better to keep it on the CPU (now we *overlap*)



### **Development strategy: PHIGEMM library**

- Project started in 2010 (P. Yang)
- Inspired by M.Fatica LINPACK work
- Independent open-source library, BSD license
- GPU+CPU BLAS 3 \*GEMM routine
- Manual or "semi-automatic" (SELFUNE) workload split
- Special-K for rectangular matrices
- C and FORTRAN interfaces
- Detailed profiling of each call (with a little "hack" in your code)
- CUDA 3.x\* and CUDA 4.x compatible
- Pinned/non-pinned, sync/async
- Support of multi-GPU
- Current version 1.8.3

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







### **Development strategy: VLOC\_PSI**

It combines CUDA kernels and CUFFT/CUBLAS calls

- CUDA kernel (INIT\_PSI) assembles the FFT grid
- CUFFT\_INVERSE is performed (CUFFTEXECZ2Z)
- A CUDA kernel (VEC\_PROD) performs multiplications over the vector
- CUFFT\_FORWARD is performed (CUFFTEXECZ2Z)
- After transformation data need to be scaled (CUBLASZDSCAL)
- CUDA kernel (SAVE\_HPSI) adds the new contribution to the final vector
- $\rightarrow$  variants for gamma-point and k-point
- → assembling kernels are memory-bounded

### Note:

- External loop over bands is kept on CPU side
- External loop over bands is not fixed (it decreases during the iterative SCF process!)



### **Development strategy: VLOC\_PSI serial**



![](_page_14_Picture_0.jpeg)

## The parallel "FFT issue"

![](_page_14_Figure_2.jpeg)

![](_page_15_Picture_0.jpeg)

### **Development strategy: VLOC\_PSI parallel**

![](_page_15_Figure_2.jpeg)

**Overlapping is possible!!** 

![](_page_16_Picture_0.jpeg)

### **Development strategy: memory allocation**

### On GPU side...

- Simple CUDAMALLOC of 85~90% of GPU memory
  - managed 1:1 MPI:GPU or N:1 MPI:GPU
  - if MAGMA is used than more space has to be left un-allocated
- Pointer to DEV\_SCRATCH\_QE is globally accessible
- PHIGEMM (through PHIGEMMINIT) makes use this buffer to perform GEMM operation

### On CPU side...

- CUDAMEMCOPY (few CUDAMEMSET)
- "shifts" are pre-computed considering memory alignment
- Specific data allocations can be either pinned or not
  - BUT pinned memory slow down all the application!
  - No pinned memory by default!

![](_page_17_Picture_0.jpeg)

## **Benchmarking platforms**

Wide set of workstations (and different GPU) & HPC clusters...

- FERMI (ICHEC): assembled workstation
  - CPU: 2 Intel Xeon X5650 (6-core), 24 GByte RAM
  - GPU: 2 {C2050, GTX480, C2075}
  - SW: CUDA 4.1, Intel compilers
  - GEMINI (ICHEC): Dell Power Edge C6145 & C410x
    - CPU: 4 AMD 6136 (8-core), 64 GByte RAM
    - GPU: 4 M2090
    - SW: CUDA 4.x, Intel compilers, PGI (12.x)
  - STONEY (ICHEC): Bull Novascale R422-E2, 24 GPU nodes
    - CPU: 2 Intel Xeon X5560 (4-core), 48 GByte RAM
    - GPU: 2 M2090
    - SW: CUDA 4.0, Intel compilers

- LONGHORN (TACC): Dell XD Cluster, 240 GPU nodes
  - CPU: 2 Intel Xeon 5355 (4-core), 48 GByte RAM
  - GPU: 2 Quadro FX 5800
  - SW: CUDA 4.0, Intel compilers
- CURIE (CEA): Bull GPUs B505 blades, 144 GPU nodes
  - CPU: 2 Intel Westmere (4-core), 24 GByte RAM
  - GPU: 2 M2090
  - SW: CUDA 4.1, Bull MPI stack, Intel compiler
- PLX (CINECA): IBM iDataPlex DX360M3, 264 GPU nodes
  - CPU: 2 Intel Westmere (6-core), 48 GByte RAM
  - GPU: 2 M2070
  - SW: CUDA 4.0, Intel compilers, PGI (11.x)

![](_page_18_Picture_0.jpeg)

### **Power measurement**

(cheap) Prodigit 200M Plug-in Main Power and Energy Monitor

- max voltage 250V
- max current 15A
- max active power 3750 Watts
- kWH accuracy: 30ppm
- kWh displayed accuracy: +/- 0.01

Plug to FERMI...

- measuring operational power absorbed
- the workstation is "equivalent" to 1 PLX node
  - IBM Dataplex nodes are higher energy efficient than a workstation, CPU are equivalent, GPUs absorb comparable amount of power
- qualitative considerations...

![](_page_18_Picture_13.jpeg)

![](_page_19_Picture_0.jpeg)

## **Benchmark philosophy**

Impossible to represent a realistic scenario using only ad-hoc designed benchmarks...

- user inputs that represent on-going investigations
- user inputs that represent challenge (too long to run here)
- user inputs that represent starting point to go through new "science"

Two-side goal was

- provide a feedback to US, evaluating if what we implemented actually cover what a generic input might trigger
- provide a feedback to the USERS, encouraging them to try and support the GPU implementation
- > We received almost 15 contributions after a "Call of Benchmarks"
- We selected 8~9 of them as benchmarks (both for serial and parallel
  - 3 are quite challenging (> 500 atoms, thousands of electrons)

![](_page_20_Picture_0.jpeg)

### Those who have contacted us directly...

![](_page_20_Figure_2.jpeg)

![](_page_21_Picture_0.jpeg)

### AUSURF112, serial (FERMI)

![](_page_21_Figure_2.jpeg)

![](_page_22_Picture_0.jpeg)

### **Performance & Power consumption (serial)**

-58%

3.1x

6 OMP

1 GPU

![](_page_22_Figure_2.jpeg)

![](_page_22_Figure_3.jpeg)

23

![](_page_23_Picture_0.jpeg)

### GeSnTe134.in, parallel (PLX)

#### **1.61x 1.83x** GeSnTe13 Walltime of full SCF 4000 2.33x CPIL on ly: PŪ: 2.69x 3500 3000 2.32x 2500 ິ **1.91x** 2000 Time V 2.04x 1500 V 1000 500 0 8 1.2 1.6 2.4 32 4.4 4 (4.8) (288)(384)(9.6)(144)(192)(528)Total number of Cores

### MGST.hex.rx, parallel (PLX)

![](_page_23_Figure_4.jpeg)

![](_page_24_Picture_0.jpeg)

## IRMOF-M11 (130), parallel (STONEY)

![](_page_24_Figure_2.jpeg)

Very small system... it can even run in SERIAL!

IR-MOF 130, 5 SCF. 2500 6 4.26x 2000 5 speedup 1500 4 3 1000 500 2 0. 12 OMP 12 OMP 12 MPI 4 MPI 4 MPI 3 OMP 1 GPU 3 OMP 2 GPU wall-time speedup

May 17, 2012

ິ

Time

![](_page_25_Picture_0.jpeg)

## IRMOF-M11 (520), parallel (STONEY)

![](_page_25_Figure_2.jpeg)

![](_page_26_Picture_0.jpeg)

## **Hitting the limit**

![](_page_26_Figure_2.jpeg)

![](_page_27_Picture_0.jpeg)

### **Closing the loop...**

![](_page_27_Figure_2.jpeg)

What we learn...

- easy (more or less) but *test-and-try* approach
- difficulties to compile the code using PGI
- (now) no data transfer overlapping
- **BUT** if not heavy computation  $\rightarrow$  big loss of performance

![](_page_28_Picture_0.jpeg)

### **PRACE Preparatory Access**

- <u>Physical interests</u>:
  - prototypical material for optoelectronic applications (e.g. light emitting diodes, solar cells)
  - easy-growth nanoparticles through and chemical processes (colloidal synthesis)
- <u>Numerical challenge</u>:
  - high electrons-to-atoms ratio in pseudo-potential
  - calculations due to the inclusion of Cd<sub>4d</sub> electrons in valence shell
- GPU challenge:
  - accelerate stress/forces calculations using (first) OpenACC
- <u>Collaborations</u>: A. Calzolari (CNR-NANO), C. Cavazzoni (CINECA)
- Funding: project pa0699 (2012), CURIE cluster, 300K hours

![](_page_28_Picture_12.jpeg)

r=10 Å #at =159, #el = 638 1ay 17, 2012

![](_page_28_Figure_14.jpeg)

r=12 A #at =275, #el = 1110

![](_page_28_Picture_16.jpeg)

r=15 Å r=18 Å #at =489, #el = 1938 #at =922, #el = 3668 F.Spiga, GPU PWscf, GTC2012

![](_page_28_Picture_18.jpeg)

X. Peng et al. Nature **404**, 59 (2000).

r=20 Å #at =12149, #el = 4844

> r=25 Å #at =2365, #el = 9370

> > r=30 Å #at =4109, #ela = 16282

![](_page_29_Picture_0.jpeg)

### CdSe-159, parallel (STONEY)

![](_page_29_Figure_2.jpeg)

![](_page_30_Picture_0.jpeg)

### Achievements

Direct effects...

- PWscF has been extended to use GPUs to accelerate both gamma- and k-points calculations
- Performance improvements for a **PRODUCTION** package
  - for serial, an average of 3-/3.5-times (full-socket vs full-socket + GPU)
  - for parallel, an average of >= 2-times (full-node vs full-node + GPUs)
  - no (visible) scalability improvement but faster time-to-solution
- 99% match of numerical consistency between CPU-only and CPU+GPU calculations
- tested on several platforms and several GPUs
- average of 150 downloads per released versions (6 different 0.X releases since July 2011)

Side effects...

- new interesting benchmarks provided by users
- lots of profiling for both CPU and CPU+GPU (CPU load, GPU load, memory occupancy, internal clocks,...)
- few improvements on the CPU code (OpenMP)

![](_page_31_Picture_0.jpeg)

## **Best Practices**

- <u>Scientific case</u>: LSMO-BFO (120 atoms)
- <u>PI</u>: Rodrigo Neumann Barros Ferreira, PhD Student Solid State Physics Department Physics Institute, Rio de Janeiro Federal University
- <u>Description</u>: 1024 electrons, 615 different quantum-mechanical states considered, **40 k-points** for the integration over the Brillouin zone.
- <u>Goal</u>: exploit QE *pool* parallelism and GPU  $\rightarrow$  keep the FFT local by using npool=npocs

Computer Nodes	Execution Time [s] #10 self-consistency cycles	Speed-up	
1 x IBM Power 575, P6 4.7 GHz (32 cores)	20314.59		60%
2 x iDataPlex DX360M3, dual Xeon E5645 6-cores 2.40 GHz (24 cores)	52057.22		user budget
2 x iDataPlex DX360M3, dual Xeon E5645 6-cores 2.40 GHz (24 cores) + <b>4 NVIDIA 2070 (USE_3D_FFT)</b>	10029.1	5.2x (2x)	

![](_page_32_Picture_0.jpeg)

## **Current (known) limitations**

- non-collinear calculations
  - special variant for VLOC\_PSI not yet accelerated
- spin magnetization
  - additional (CPU) code that slow down the GPU performance
- Low number of atoms (<32)
  - if there is not "enough" work to keep CPU busy... GPU will not do better
- High number of k-points
  - systems with high number of k-points (> 8) suffer of performance degradation due to I/O
  - mitigable in principle by using in a smart way the MPI parallelism
- Periodic systems with low number of atoms (<32) + high number of k-points
  - better to physically create a large system with less k-point (if it has sense)
  - there might me exceptions... (i.e. when the cell is big)

![](_page_33_Picture_0.jpeg)

## **Big Challenges & Collaborations**

Petascale computations in mineral physics with the Quantum ESPRESSO

- <u>P.I.</u>: Prof Renata Wentzcovitch (Chemical Engineering & Materials Science, U. of Minnesota)
- Funded by: NSF
- <u>Objectives</u>: investigation of mineral properties must be investigated in a wide range of pressure, temperature, and chemical compositions
- <u>Target machine (facility)</u>: Blue Waters (NCSA)

### Center for software innovation: (R)Evolutionary Materials Development

- <u>P.I.</u>: Prof Marco Buongiorno Nardelli (Physics and Chemistry Departments, U. of North Texas)
- <u>Funded by</u>: DoE
- <u>Objectives</u>: mapping new materials enabled technologies (METs) genome across different length scales to enable accelerated discovery and technological transfer
- <u>Target machine (facility)</u>: Jaguar/Titan (ORNL) + local resources

![](_page_34_Picture_0.jpeg)

### **Future Plans**

Future developments will focus on three independent directions...

- 1. Improve PWscF by engage big scientific challenges with scientists (users drive the development priorities)
- 1. Extend GPU capabilities to other code of the suite (next: CP, PHONON)
- 1. Improve the support for multi-GPU in serial calculations

QE-GPU is an open collaboration:

- repository connected to the main one, QE-GPU is like "an extension" of the QE suite
- contributors can develop GPU code with no impact to the main suite (modularized structure)
- GPU technology exploration and evaluation

![](_page_35_Picture_0.jpeg)

# Thank you for your attention!

No CPUs or GPUs have been damaged during the preparation of this talk (-:

### DOWNLOAD IT AT http://tinyurl.com/gpu-pwscf

Acknowledgments:

Ivan Girotto (ICHEC, now ICTP), Carlo Cavazzoni (CINECA), Paolo Giannozzi (U. of UDine/DEMOCRITOS), Layla Martin-Samos (DEMOCRITOS/U. of Nuova Goriza), Arrigo Calzolari (CNR-NANO), Wei Zhang (RWTH Aachen University), Clima Sergiu (IMEC), Koroteev Victor (NIIC SB RAS), Bhagawan Sahu (Globalfoundries) and many others...

![](_page_35_Picture_6.jpeg)

![](_page_35_Picture_7.jpeg)

Ireland's EU Structural Funds Programmes 2007 - 2013

Co-funded by the Irish Government and the European Union

![](_page_35_Picture_10.jpeg)

EUROPEAN REGIONAL DEVELOPMENT FUND

![](_page_35_Picture_12.jpeg)

GUS EOLAIOCHTA

![](_page_35_Picture_13.jpeg)

Higher Education Authority An tÚdarás um Ard-Oideachas

![](_page_36_Picture_0.jpeg)

## What can QUANTUM ESPRESSO's PWscF do?

- both gamma-point and **k**-point calculation
- both insulators and metals, with various flavors of broadening, or tetrahedra
- any crystal structure or supercell form
- ground-state energy and one-electron (Kohn-Sham) orbitals;
- atomic forces, stresses, and structural optimization;
- molecular dynamics on the ground-state Born-Oppenheimer surface, also with variable cell;
- Nudged Elastic Band (NEB) and Fourier String Method Dynamics (SMD) methods
- norm-conserving PP's in separable form, ultrasoft Vanderbilt PP's, PAW
- almost all flavours of LDA and of gradient-corrected exchange-correlation functionals (PW91, PBE, B88-P86, BLYP,...), DFT+U, exact exchange and a few hybrid functionals (PBE0, B3LYP), TPSS meta-GGA
- spin-polarized, magnetic systems (including non-collinear magnetism and spin-orbit interactions)

(see: http://www.quantum-espresso.org/whatcanqedo.php)

![](_page_37_Picture_0.jpeg)

### **Coding numerical formulas ...**

#### SCF:

compute potential
solve KS eigen-problem
Loop over k-points:
 Davidson iteration / CG iteration:
 compute/update H \* psi:
 compute kinetic and non-local term (in G space)
 Loop over (not converged) bands:
 FFT psi to R space
 compute V \* psi
 FFT V \* psi back to G space

project H in the reduced space (ZGEMM) diagonalize the reduced Hamiltonian: cholesky factorization call to LAPACK/SCALAPACK/MAGMA/PLASMA diagonalization routine

**BACKUP SLIDE** 

 $\hat{H}_{\scriptscriptstyle K\!S}ig| \mathcal{Y}_{ec{k},b}ig
angle$ 

 $\left< \mathbf{y}_{\vec{k},a} \middle| \hat{H}_{KS} \middle| \mathbf{y}_{\vec{k},b} \right>$ 

 $\hat{H}_{KS} \left| \psi_{\vec{k},b} \right\rangle = \varepsilon_{\vec{k},b} \left| \psi_{\vec{k},b} \right\rangle$ 

compute new density loop over k-points: loop over bands: FFT psi to R space accumulate psi charge density symmetrisation

 $n(\vec{r}) = 2 \mathop{\text{a}}_{\vec{k}} \mathop{\text{a}}_{\nu} \left| \mathcal{Y}_{\nu,\vec{k}}(\vec{r}) \right|^2$ 

![](_page_38_Picture_0.jpeg)

## H \* psi (VLOC\_PSI)

```
compute/update H * psi:
    compute kinetic and non-local term (in G space)
        complexity : N_i \times (N \times N_g + N_g \times N \times N_p)
    Loop over (not converged) bands:
    FFT (psi) to R space
        complexity : N_i \times N_b \times FFT(N_r)
        compute V * psi
        complexity : N_i \times N_b \times N_r
    FFT (V * psi) back to G space
        complexity : N_i \times N_b \times FFT(N_r)
    compute Vexx:
        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)  $N_g =$  number of G vectors  $N_i =$  number of Davidson iteration (usually about 10)  $N_p$  = number of PP projector  $N_r$  = size of the 3D FFT grid  $N_q$  = number of q-point (may be different from  $N_k$ )

![](_page_39_Picture_0.jpeg)

### **Pinned** versus **non-Pinned** AUSURF112, k point (ZGEMM), 4MPI x 3 OMP with 2 GPUs, few SCF iterations

![](_page_39_Figure_2.jpeg)

May 17, 2012 BACKUP SLIDE

![](_page_40_Picture_0.jpeg)

### **Development strategy: diagonalization**

![](_page_40_Figure_2.jpeg)

BACKUP SLIDE

![](_page_41_Picture_0.jpeg)

### PHIGEMM : special-K/1

![](_page_41_Figure_2.jpeg)

PHIGEMM does not perform very well, overlapping COMP & COMM does not produce any increment in performance

### $\bigcirc$

Split C = alpha\*op(A)\*op(B) + beta\*C and perform the multiplication on GPU and the sum on CPU

### $\bigcirc$

small block "scheuduled" statically one after the other, overlapping in data movement (pinned buffer)

May 17, 2012

**BACKUP SLIDE** 

![](_page_42_Picture_0.jpeg)

### PHIGEMM : special-K/2

![](_page_42_Figure_2.jpeg)