QUANTUM ESPRESSO is an integrated suite of **Open-Source** computer codes for **electronic-structure** calculations and **materials modeling** at the nanoscale.

It is based on *density functional theory, plane waves, and pseudopotentials.*
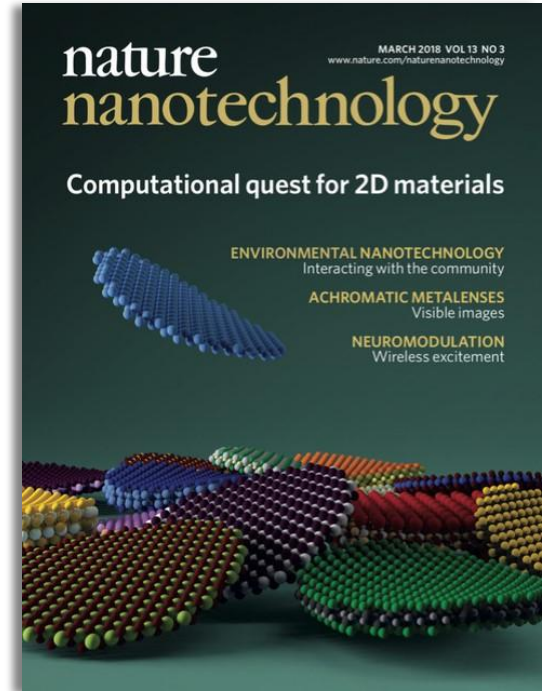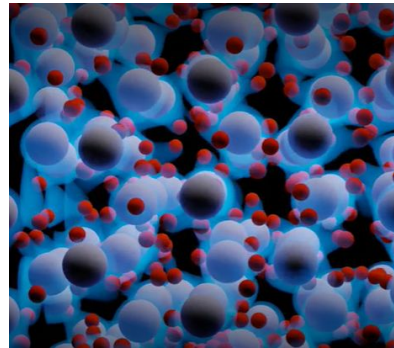
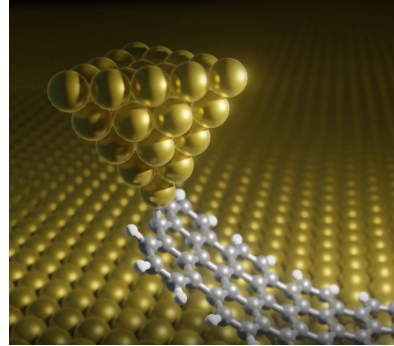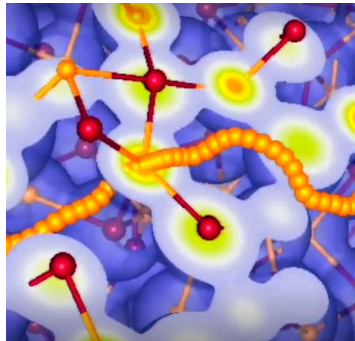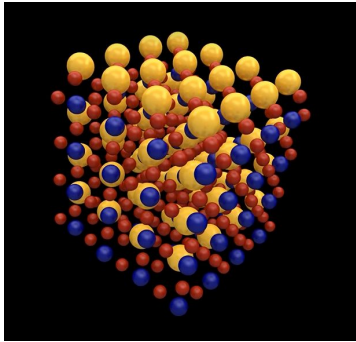QUANTUM ESPRESSO is a community code.

# Outline

- **Quantum mechanics** for **materials**;

- the QUANTUM **ESPRESSO** (**QE**) suite;

- **porting strategy** and constraints;

- **multiple standards** in QE;

- toward a portable **FFTXlib**, library for 3D FFTs;

- results and state of the art of **QE performance**;

- summary and outlook

# Quantum ESPRESSO

AB INITIO QUANTUM MECHANICS

no input parameters for material modeling

reduces costs, accelerates discoveries



*Two-dimensional materials from high-throughput computational exfoliation of experimentally known compounds*, Nature Nanotechnology **13**, 246 (2018). doi:10.1038/s41565-017-0035-5

# Quantum mechanics for materials

# Quantum mechanics for materials
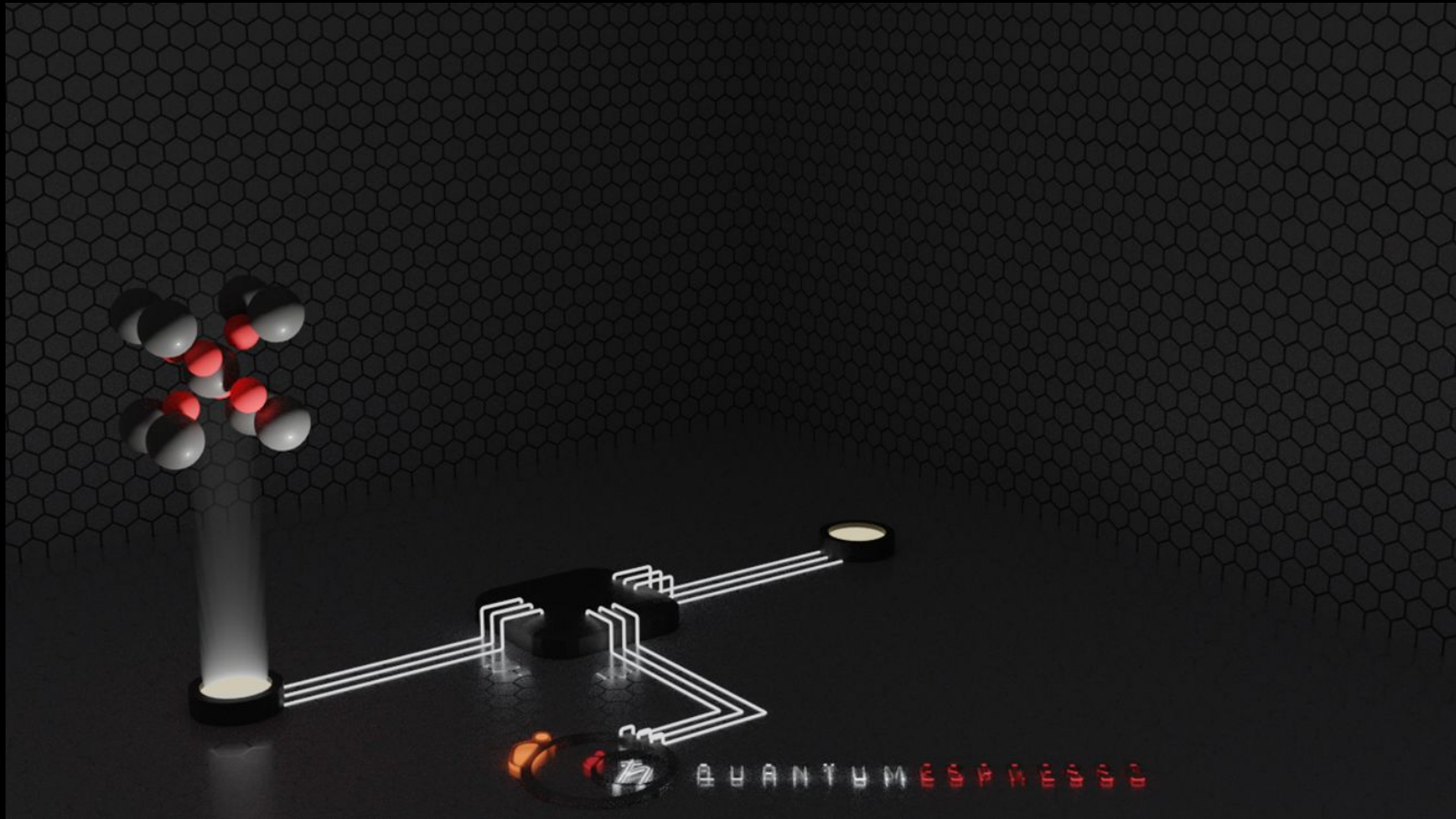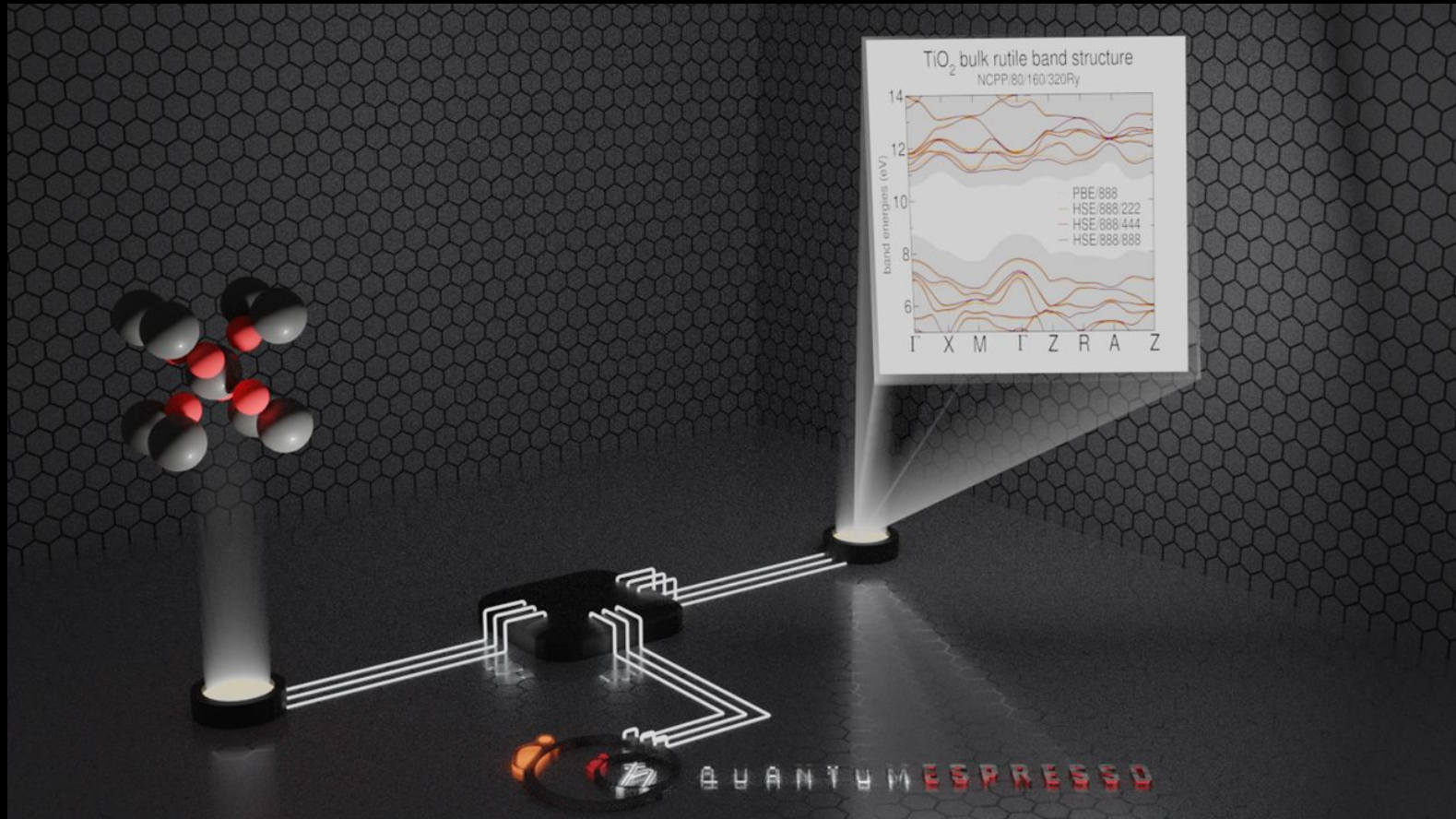
TiO$_2$ bulk rutile band structure
NCPP/80/160/320Ry

PBE/888
HSE/888/222
HSE/888/444
HSE/888/888

band energies (eV)

Γ  X  M  Γ  Z  R  A  Z

QUANTUMESPRESSO

EELS spectrum of Si(100) - p2x1
range 0.0 - 5.0 eV with 51 sampling frequencies

LIGHTHOUSE CODES

DOMAIN EXPERTS & CODE DEVELOPERS

HPC EXPERTS & DATA CENTRES

TECHNOLOGY & CO-DESIGN PARTNERS

Coe for HPC applications in material science

exploit **frontier HPC**
for material science research in strong link with **scientific communities**

CODE PORTING

CO-DESIGN

HPC ECOSYSTEM

# ICSC National Research Centre

for High Performance Computing, Big Data and Quantum Computing

## Flagship codes

Geographic distribution of the authors of the articles citing the main reference articles of Quantum ESPRESSO

**Quantum ESPRESSO** is an **open initiative** involving a large community of developers and contributors from different **regions of the world**



Citations per year

4000

3000

2000

1000

0

2010 2013 2016 2019 2022

Data provided by the courtesy of the Quantum ESPRESSO foundation

# The Quantum ESPRESSO suite

**35000+ download** of the code from the website in 2022, mostly from Europe, USA, India and China

**OTHER**
**23,7%**

**NATIONAL LAB**
**4,9%**

**INDUSTRY**
**6,3%**

**ACADEMIA**
**65,1%**

Geographic distribution and main professional fields of people who have downloaded QE from the website since the beginning of 2022

0 — 1.649

# The Quantum ESPRESSO suite

**DENSITY FUNCTIONAL THEORY**

$$\left[-\frac{\hbar^2}{2m}\nabla^2 + V_s(\mathbf{r})\right]\varphi_i(\mathbf{r}) = \varepsilon_i\varphi_i(\mathbf{r}).$$

**PLANE WAVES & PSEUDOPOTENTIAL**

$$\varphi_\alpha(\mathbf{r}) = \frac{1}{\sqrt{\Omega}}\exp[iG_\alpha \cdot \mathbf{r}]$$

**DUAL SPACE TECHNIQUE**

$$\psi(\mathbf{r}) \rightarrow \psi(\mathbf{k}) \rightarrow \psi(\mathbf{r})$$

# The Quantum ESPRESSO suite

**DENSITY FUNCTIONAL THEORY**

$$\left[ -\frac{\hbar^2}{2m}\nabla^2 + V_{\mathbf{s}}(\mathbf{r}) \right] \varphi_i(\mathbf{r}) = \varepsilon_i\varphi_i(\mathbf{r}).$$
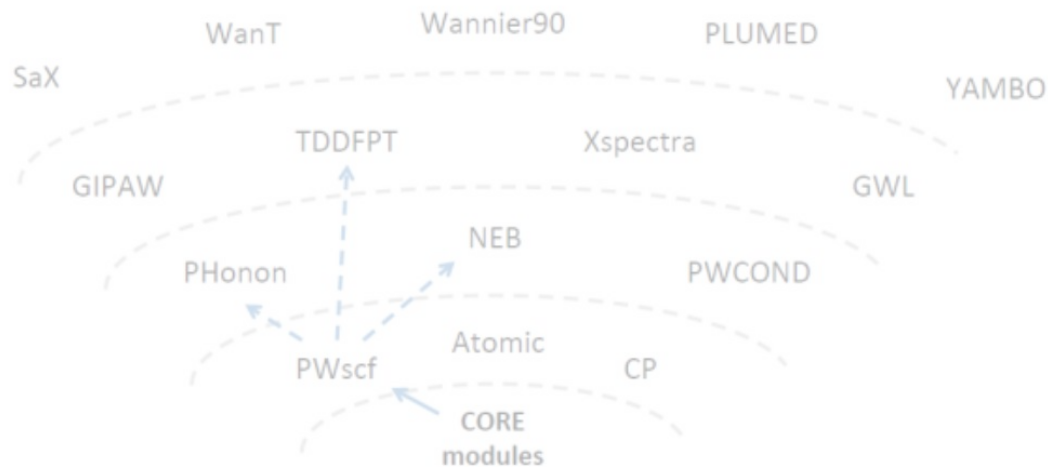
**PLANE WAVES & PSEUDOPOTENTIAL**

$$\varphi_\alpha(\mathbf{r}) = \frac{1}{\sqrt{\Omega}}\exp[iG_\alpha \cdot \mathbf{r}]$$

**DUAL SPACE TECHNIQUE**

$$\psi(\mathbf{r}) \rightarrow \psi(\mathbf{k}) \rightarrow \psi(\mathbf{r})$$

WanT    Wannier90    PLUMED

SaX    YAMBO

TDDFPT    Xspectra

GIPAW    GWL

NEB

PHonon    PWCOND

Atomic

PWscf    CP

CORE modules

# Data distribution



$N_{PW}$ (**G** vectors)

$$\psi_{ik}(\mathbf{r}) = \sum_{\mathbf{G}}^{N_{PW}} C_{\mathbf{G},ik} \frac{e^{i((\mathbf{G}+\mathbf{k})\cdot\mathbf{r})}}{\sqrt{\Omega}}$$

$i = 1, ..., N_b$

$k = 1, ..., N_k$

# The parallel scheme of PWscf

# Porting strategy

# Towards a portable GPU version

The transition from CUDA to Openacc



DIRECTIVE-BASED
PROGRAMMING MODELS

MAINTAINABLE

PORTABLE

SINGLE SOURCE CODE

# Towards a portable GPU version

Modularity supports interoperability and new programming models

Modularity supports interoperability and new programming models

Modularity supports interoperability and new programming models

Modularity supports interoperability and new programming models

# On the porting roadmap

- J. Chem. Phys. **152,** 154105 (2020)

- J. Chem. Theory Comput. **19**, 6992 (2023)



- Until v 6.8;

- from v 7.0;

- under development;

- current goal.

# OMP porting of QE



Basic features:

- **loop** offloading;

- **global variables** offloading and pinning;

- manage different **backends** (linear algebra and FFTs);

- **streams** and/or **tasks** (for async batched FFTs).

*Ported*:

- standard **FFTs** (cpu driver);
- **KS_Solver** (except diagonalization);
- Interfaces for **mathematical libraries**;
- qe instrumentation routines (**rocprof**) have been added.

*To be ported*:

- diagonalization (zhegv);
- **batched FFTs**;
- **Hubbard**, forces, stress;
- codes other than PW.

Multiple standards in QE

## CUF only

**Host to Device**

```
if ( use_gpu ) then
  arg_d = arg
endif
```

**Routine calls**

```
if ( use_gpu ) then
  call abc( arg_d )
else
  call abc( arg )
endif
```

**Interfaces**

```
interface abc
  subroutine abc_cpu( v )
  subroutine abc_gpu( v_d )
end interface
```

# Offload

|  | CUF only | CUF interfaces<br>OpenACC parent code |
|---|---|---|
| **Host to Device** | ```fortran
if ( use_gpu ) then
  arg_d = arg
endif
``` | `!$acc update device(arg)` |
| **Routine calls** | ```fortran
if ( use_gpu ) then
  call abc( arg_d )
else
  call abc( arg )
endif
``` | ```fortran
!$acc host_data use_device(arg)
call abc( arg )
!$acc end host_data
``` |
| **Interfaces** | ```fortran
interface abc
  subroutine abc_cpu( v )
  subroutine abc_gpu( v_d )
end interface
``` | |

# Offload

| | CUF only | CUF interfaces OpenACC parent code | OpenACC only |
|---|---|---|---|
| **Host to Device** | `if ( use_gpu ) then`<br>`  arg_d = arg`<br>`endif` | `!$acc update device(arg)` | |
| **Routine calls** | `if ( use_gpu ) then`<br>`  call abc( arg_d )`<br>`else`<br>`  call abc( arg )`<br>`endif` | `!$acc host_data use_device(arg)`<br>`call abc( arg )`<br>`!$acc end host_data` | `call abc_acc( arg )` |
| **Interfaces** | `interface abc`<br>`  subroutine abc_cpu( v )`<br>`  subroutine abc_gpu( v_d )`<br>`end interface` | | `subroutine abc_acc(v)` |

# Offload

| CUF only | CUF interfaces OpenACC parent code | OpenACC only | OpenACC + OpenMP5 |
|---|---|---|---|

**Host to Device**

CUF only:
```
if ( use_gpu ) then
  arg_d = arg
endif
```

CUF interfaces OpenACC parent code:
```
!$acc update device(arg)
```

OpenACC + OpenMP5:
```
!$acc update device(arg)
!$omp target update to(arg)
```

**Routine calls**

CUF only:
```
if ( use_gpu ) then
  call abc( arg_d )
else
  call abc( arg )
endif
```

CUF interfaces OpenACC parent code:
```
!$acc host_data use_device(arg)
call abc( arg )
!$acc end host_data
```

OpenACC only:
```
call abc_acc( arg )
```

OpenACC + OpenMP5:
```
#if def __OPENACC
  call abc_acc( arg )
#elif def __OPENMP
  call abc_omp( arg )
#endif
```
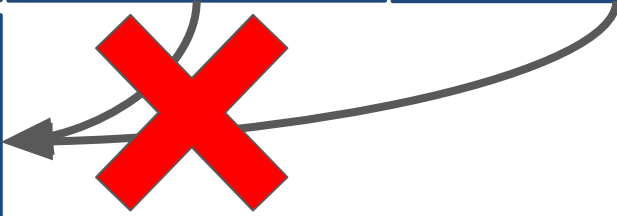
**Interfaces**

```
interface abc
  subroutine abc_cpu( v )
  subroutine abc_gpu( v_d )
end interface
```

# Offload

| | CUF only | CUF interfaces OpenACC parent code | OpenACC only | OpenACC + OpenMP5 |
|---|---|---|---|---|
| **Host to Device** | `if ( use_gpu ) then`<br>`  arg_d = arg`<br>`endif` | `!$acc update device(arg)` | | `!$acc update device(arg)`<br>`!$omp target update to(arg)` |
| **Routine calls** | `if ( use_gpu ) then`<br>`  call abc( arg_d )`<br>`else`<br>`  call abc( arg )`<br>`endif` | `!$acc host_data use_device(arg)`<br>`call abc( arg )`<br>`!$acc end host_data` | `call abc_acc( arg )` | `#if def __OPENACC`<br>`  call abc_acc( arg )`<br>`#elif def __OPENMP`<br>`  call abc_omp( arg )`<br>`#endif` |
| **Interfaces** | `interface abc`<br>`  subroutine abc_cpu( v )`<br>`  subroutine abc_gpu( v_d )`<br>`end interface` | | `subroutine abc_acc(v)` | `subroutine abc_acc( v )`<br>`subroutine abc_omp( v )` |

# Offload

| | CUF only | CUF interfaces OpenACC parent code | OpenACC only | OpenACC + OpenMP5 |
|---|---|---|---|---|
| **Host to Device** | `if ( use_gpu ) then`<br>`  arg_d = arg`<br>`endif` | `!$acc update device(arg)` | | `!$acc update device(arg)`<br>`!$omp target update to(arg)` |
| **Routine calls** | `if ( use_gpu ) then`<br>`  call abc( arg_d )`<br>`else`<br>`  call abc( arg )`<br>`endif` | `!$acc host_data use_device(arg)`<br>`call abc( arg )`<br>`!$acc end host_data` | `call abc( arg, offload )` | `call abc( arg, offload )` |
| **Interfaces** | `interface abc`<br>`  subroutine abc_cpu( v )`<br>`  subroutine abc_gpu( v_d )`<br>`end interface` | | `interface abc`<br>`  subroutine abc_cpu( v, offload )`<br>`  subroutine abc_acc( v, offload )`<br>`  subroutine abc_omp( v, offload )`<br>`end interface` | |

Courtesy of Ivan Carnimeo (SISSA)

# Offload

| OpenACC + OpenMP5 |
|---|
| !$acc update device(arg)<br>!$omp target update to(arg) |
| **call** abc( arg, **offload** ) |

→

| **devXlib** + OpenACC +OpenMP5 |
|---|
| **call** **devXlib**( arg, **offload** ) |
| **call** abc( arg, **offload** ) |

**Host to Device**

**Routine calls**

( ... )

**Interfaces**

```
interface abc
  subroutine abc_cpu( v, offload )
  subroutine abc_acc( v, offload )
  subroutine abc_omp( v, offload )
end interface
```

The Yambo group in Modena is developing a portable library (**devXlib**) to manage porting to multiplatform heterogeneous architectures

Main developers:
A .Ferretti (CNR-NANO)
N. Spallanzani (CNR-NANO)
G. Rossi (Intel)

Courtesy of Ivan Carnimeo (SISSA)

# Wrappers instead of interfaces

```fortran
!-----------------------------------------------------------------
SUBROUTINE wave_r2g( f_in, f_out, dfft, igk, howmany_set, omp_mod )
  !-----------------------------------------------------------------
  !! Wave function FFT from R to G-space.
  !
  USE fft_helper_subroutines,  ONLY: fftx_psi2c_gamma, fftx_psi2c_k
#if defined(__OPENMP_GPU)
  USE fft_helper_subroutines,  ONLY: fftx_psi2c_gamma_omp, fftx_psi2c_k_omp
#endif
  USE control_flags,           ONLY: many_fft
  !
  IMPLICIT NONE
  !

  ...

  !
  omp_offload = .FALSE.
  omp_map     = .FALSE.
#if defined(__OPENMP_GPU)
  IF (PRESENT(omp_mod)) THEN
    omp_offload = omp_mod>=0 ! run FFT on device (data already mapped)
    omp_map     = omp_mod>=1 ! map data and run FFT on device
  ENDIF
#endif
  !
```

```fortran
  !
  !$acc host_data use_device(f_in)
  IF (PRESENT(howmany_set)) THEN
    IF(omp_offload) THEN
#if defined (__OPENMP_GPU)
      IF(omp_map) THEN
        !$omp target data map(tofrom:f_in)
        CALL fwfft_y_omp( 'Wave', f_in, dfft, howmany=howmany_set(3) )
        !$omp end target data
      ELSE
        CALL fwfft_y_omp( 'Wave', f_in, dfft, howmany=howmany_set(3) )
      ENDIF
#endif
    ELSE
      CALL fwfft( 'Wave', f_in, dfft, howmany=howmany_set(3) )
    ENDIF
    !
  ELSE
    IF(omp_offload) THEN
#if defined (__OPENMP_GPU)
      IF(omp_map) THEN
        !$omp target data map(tofrom:f_in)
        CALL fwfft_y_omp( 'Wave', f_in, dfft )
        !$omp end target data
      ELSE
        CALL fwfft_y_omp( 'Wave', f_in, dfft )
      ENDIF
#endif
    ELSE
      CALL fwfft( 'Wave', f_in, dfft )
    ENDIF
  ENDIF
  !$acc end host_data
  !
  IF (gamma_only) THEN
```

# Some notes

- We need **both offloaded and non-offloaded** low level routines (e.g. FFTXlib, LAXlib) at the same time;

- we use **wrappers with offloading switch** to sort CPU and GPU low level library calls;

- **duplication** of **low level routines** still necessary (avoidable? In the future?);

- *omp target* for GPU **protected** from openACC and from CPU *omp.*

# Backends

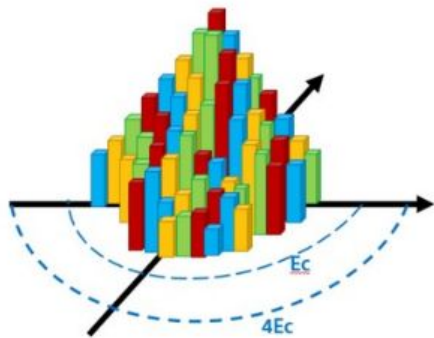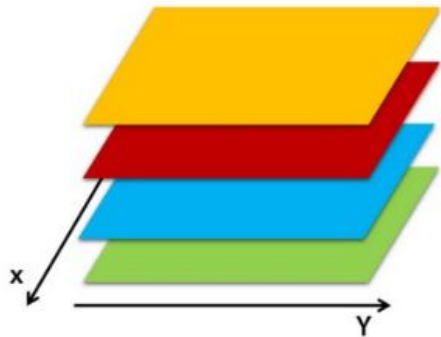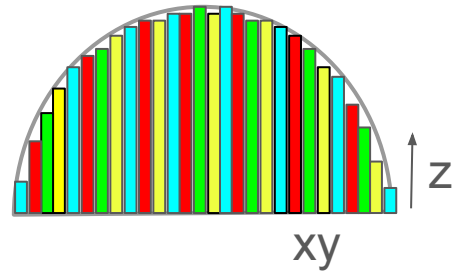| | Explicit streams | gpu/cpu interface | No need c_bind |
|---|---|---|---|
| **Linear Algebra** | | | |
| cuBlas | ✓ | ✓ | ✓ |
| rocBlas | ✓ | ✗ | ✓ |
| oneMKL | ✗ | ✗ | ✓ |
| **Fourier transforms** | | | |
| cuFFT | ✓ | ✓ | ✓ |
| hipFFT | ✓ | ✗ | ✓ |
| oneMKL | ✗ | ✗ | ✓ |

```fortran
SUBROUTINE MYDGEMM2( TRANSA, TRANSB, M, N, K, ALPHA, A, LDA, B, LDB, &
                     BETA, C, LDC, OMP_OFFLOAD )
#if defined(__CUDA)
    use cudafor
    use cublas
#elif defined(__OPENMP_GPU)
#if defined(__ONEMKL)
    use onemkl_blas_gpu
#endif
#if defined(__ROCBLAS)
    use rocblas_utils
#endif
#endif
    CHARACTER*1, INTENT(IN) :: TRANSA, TRANSB
    INTEGER, INTENT(IN) :: M, N, K, LDA, LDB, LDC
    DOUBLE PRECISION, INTENT(IN) ::   ALPHA, BETA
    DOUBLE PRECISION :: A( LDA, * ), B( LDB, * ), C( LDC, * )
    LOGICAL, INTENT(IN) :: OMP_OFFLOAD
#if defined(__CUDA)
    attributes(device) :: A, B, C
    CALL cublasdgemm( TRANSA, TRANSB, M, N, K, ALPHA, A, LDA, B, LDB, &
                      BETA, C, LDC)
#elif defined(__ONEMKL)
    IF (OMP_OFFLOAD) THEN
      !$omp target variant dispatch use_device_ptr(A, B, C)
      CALL dgemm( TRANSA, TRANSB, M, N, K, ALPHA, A, LDA, B, LDB, BETA, &
                  C, LDC)
      !$omp end target variant dispatch
    ELSE
      CALL dgemm( TRANSA, TRANSB, M, N, K, ALPHA, A, LDA, B, LDB, BETA, &
                  C, LDC)
    ENDIF
#elif defined(__ROCBLAS)
    IF (OMP_OFFLOAD) CALL rocblas_dgemm( TRANSA, TRANSB, M, N, K, ALPHA, &
                                         A, LDA, B, LDB, BETA, C, LDC)
    IF (.NOT. OMP_OFFLOAD) CALL dgemm( TRANSA, TRANSB, M, N, K, ALPHA, A, &
                                       LDA, B, LDB, BETA, C, LDC)
#else
    CALL dgemm(TRANSA, TRANSB, M, N, K, ALPHA, A, LDA, B, LDB, BETA, C, LDC)
#endif

END SUBROUTINE MYDGEMM2
```
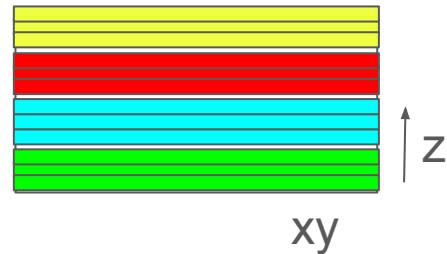
G-space
(sticks)

R-space
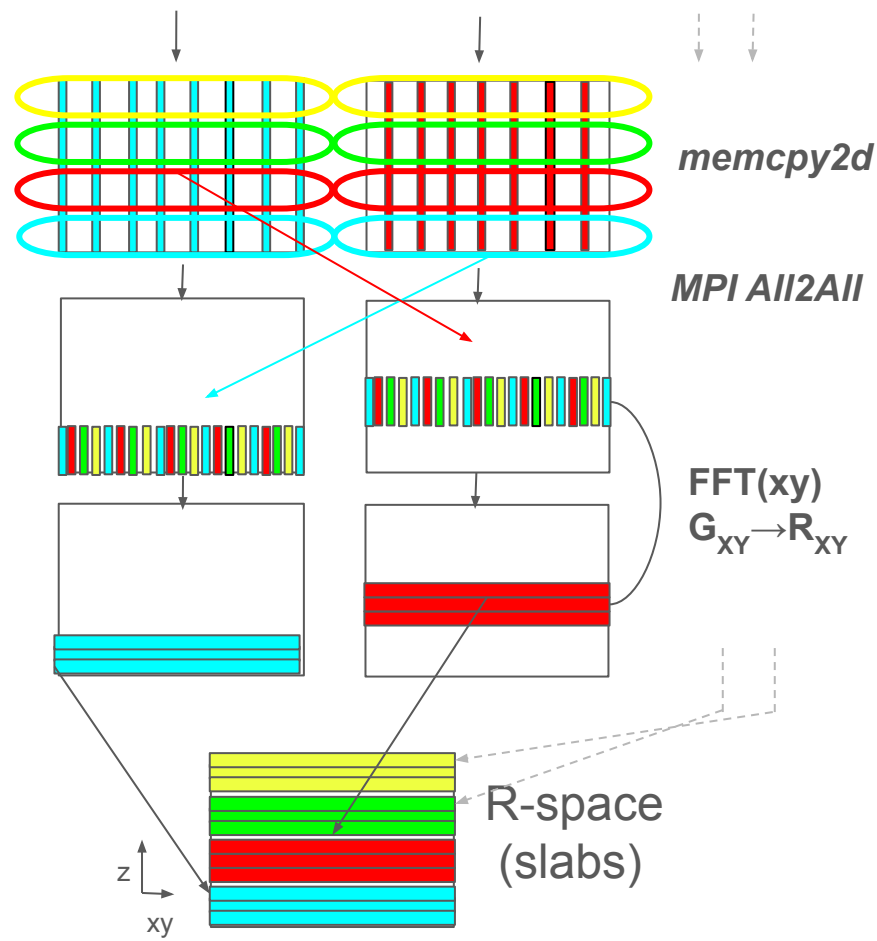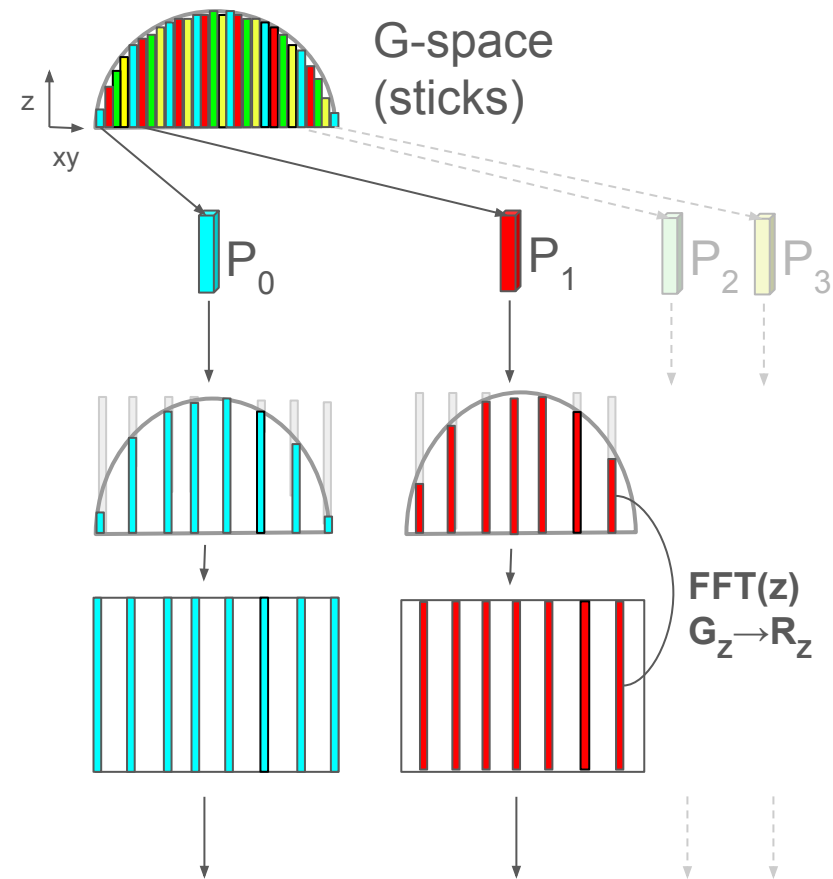(slabs)

# FFTXlib: slab decomposition



G-space
(sticks)

$P_0$  $P_1$  $P_2$  $P_3$

FFT(z)
$G_Z \rightarrow R_Z$

memcpy2d

MPI All2All

FFT(xy)
$G_{XY} \rightarrow R_{XY}$

R-space
(slabs)
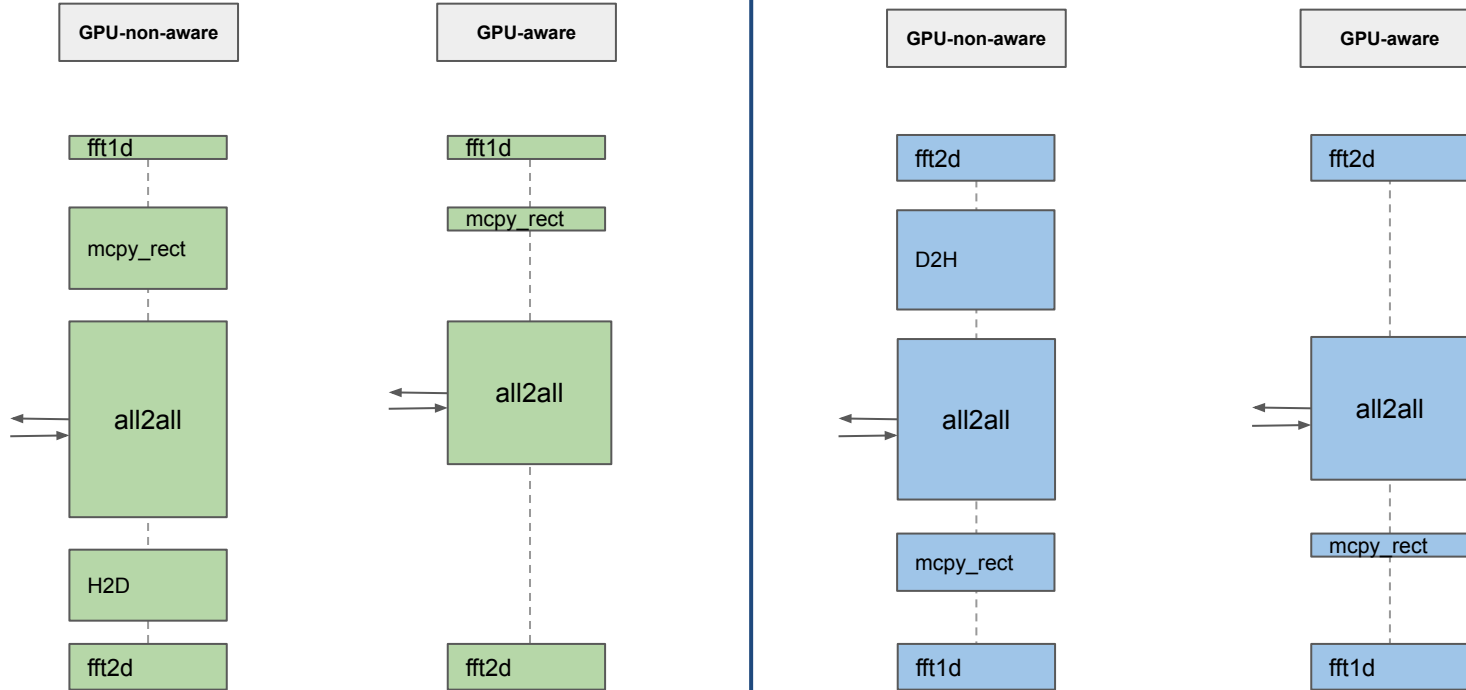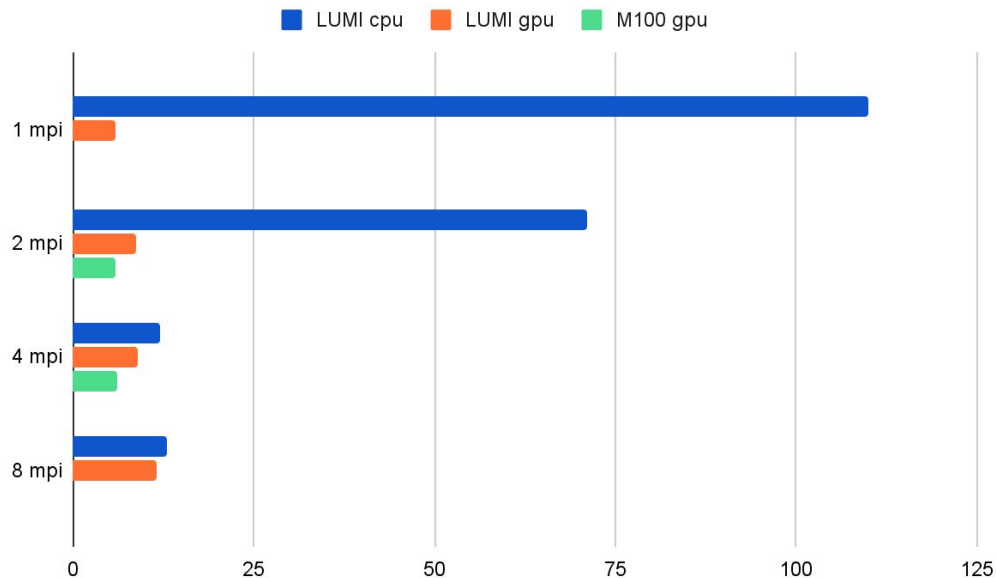
# FFTXlib: standard flow
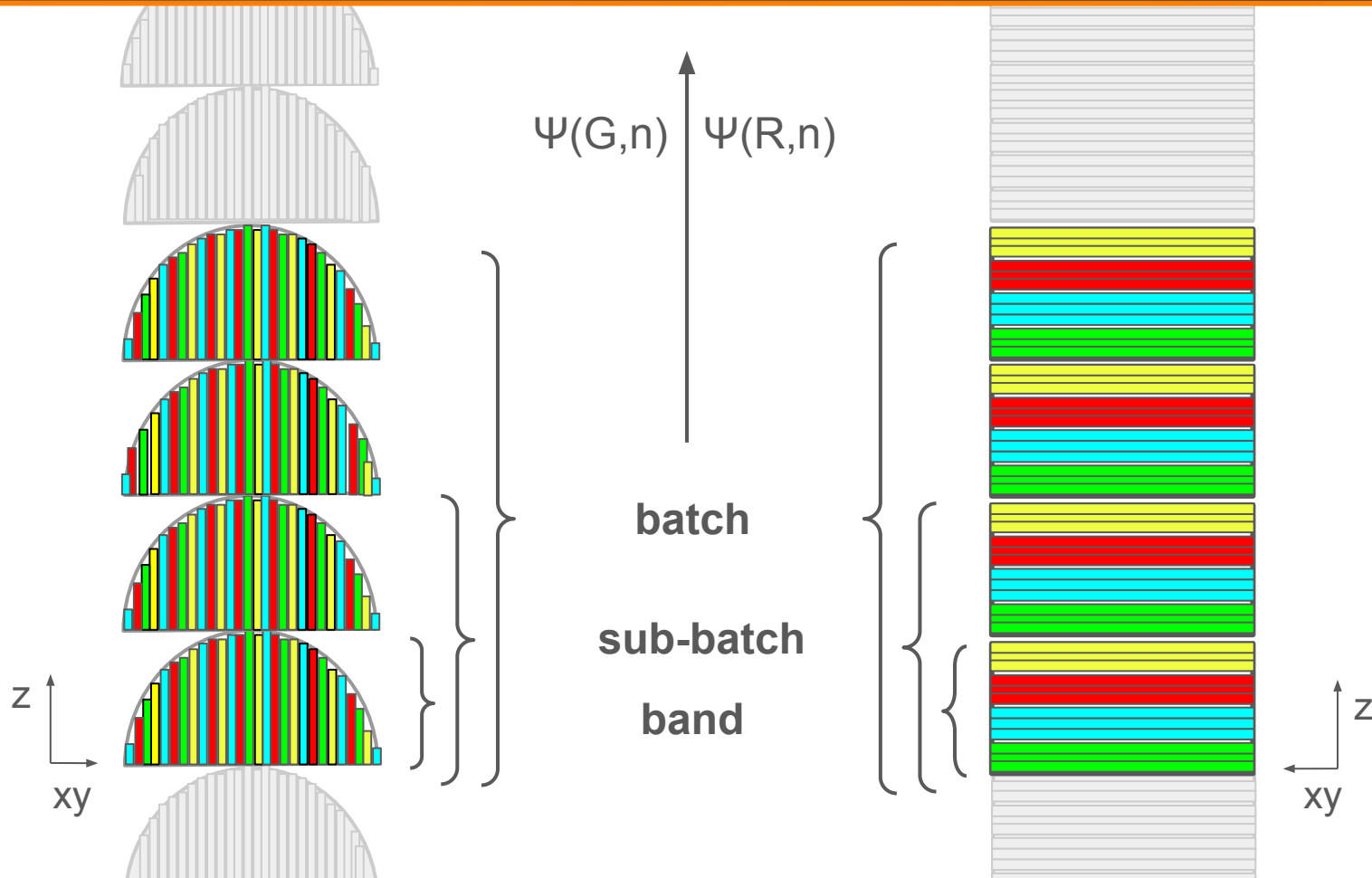
# FFTXlib: basic porting

Reference benchmark: **Ausurf 112** (1 scf step). **FFTXlib** calls, **preliminary** results:
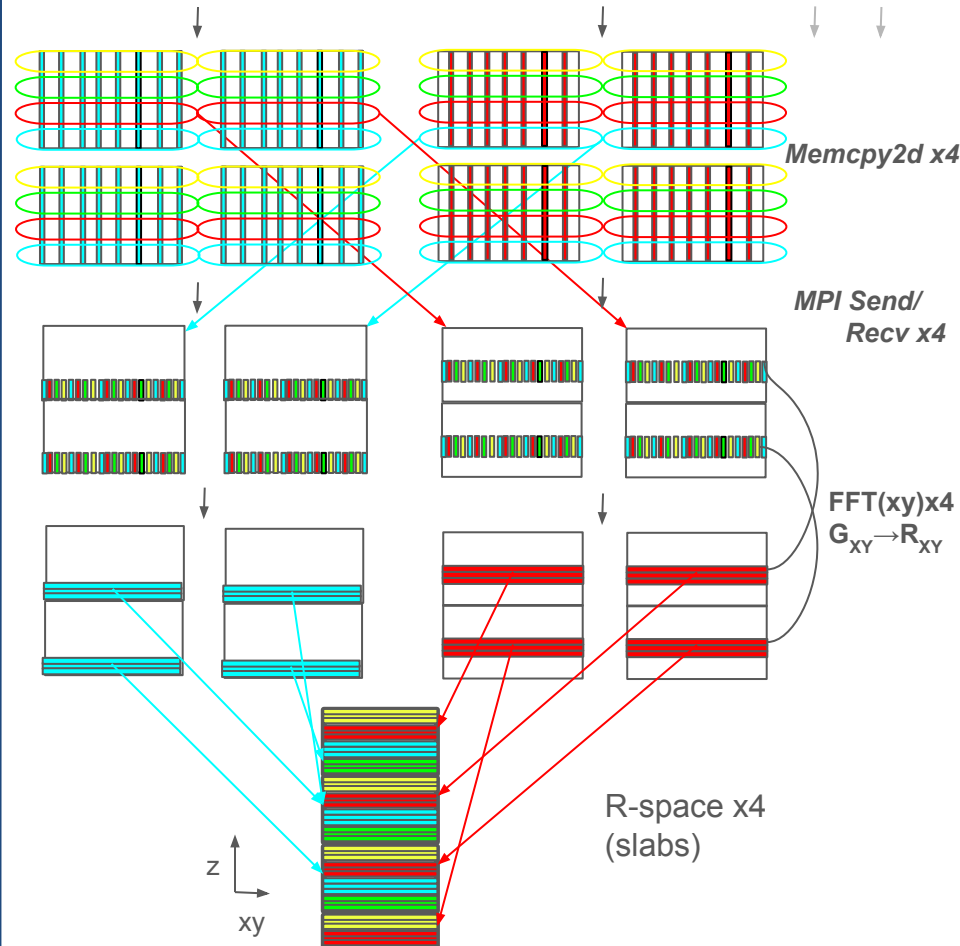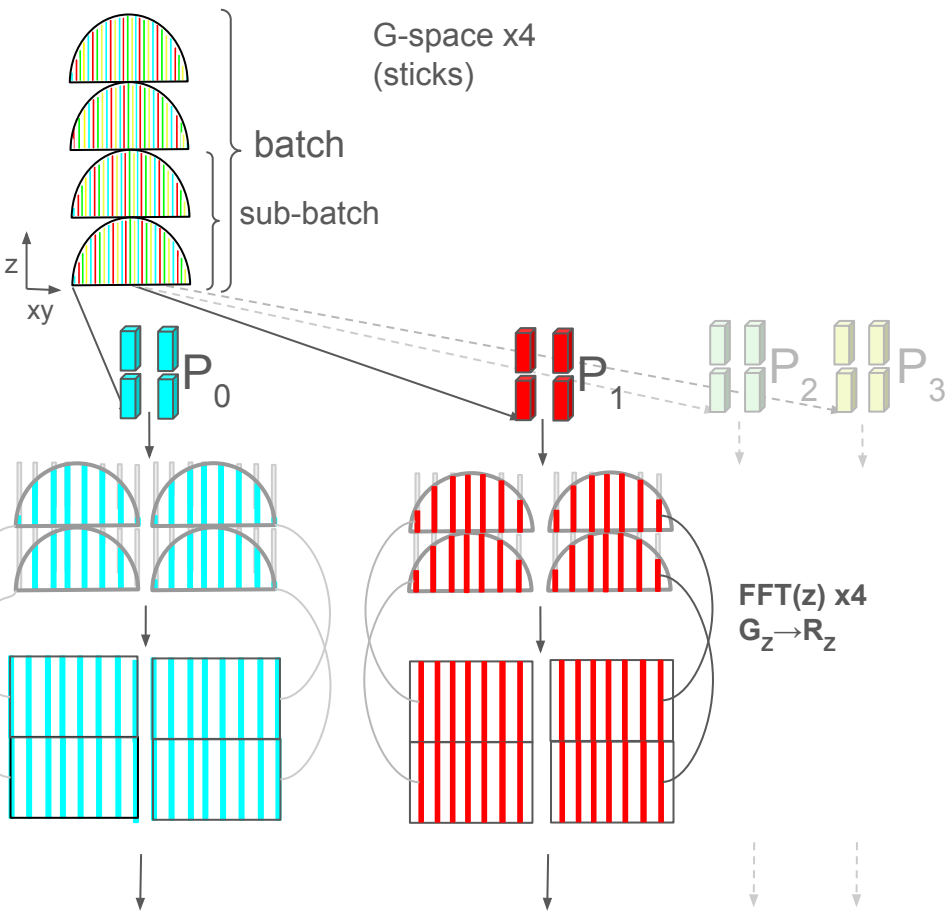


3d FFTs with SLAB decomposition (standard case):

- reference runs: M100 (**V100 gpus**)

- **overall match** between LUMI and M100;

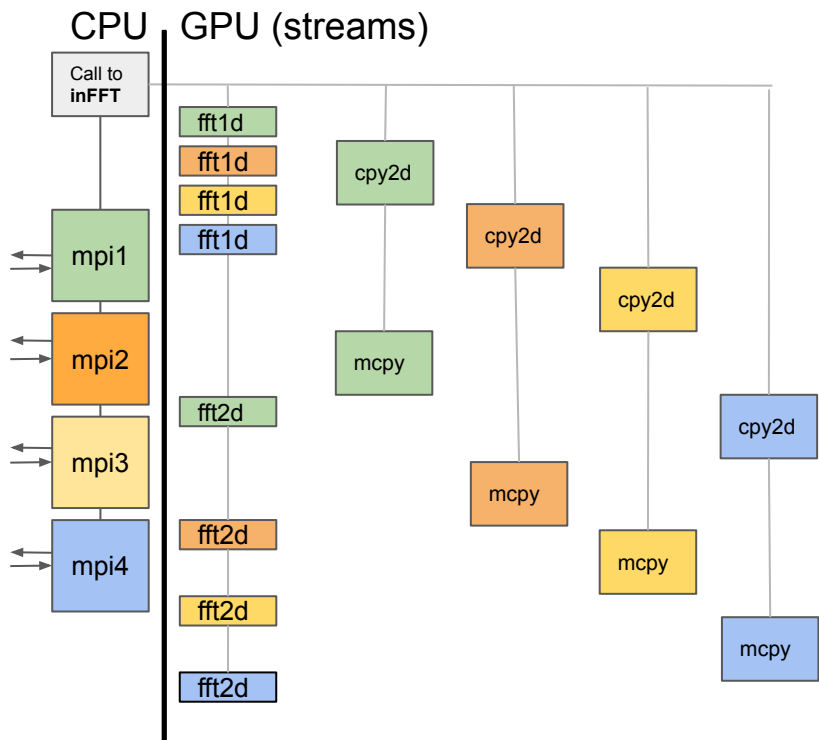- **H2D-D2H** part of the FFT looks a bit slower on LUMI side (still under investigation).

$\Psi(G,n)$ | $\Psi(R,n)$

batch

sub-batch

band

G-space x4
(sticks)

batch

sub-batch

z
xy

$P_0$

$P_1$

$P_2$

$P_3$

FFT(z) x4
$G_z \rightarrow R_z$

Memcpy2d x4

MPI Send/
Recv x4

FFT(xy)x4
$G_{XY} \rightarrow R_{XY}$

R-space x4
(slabs)

z
xy

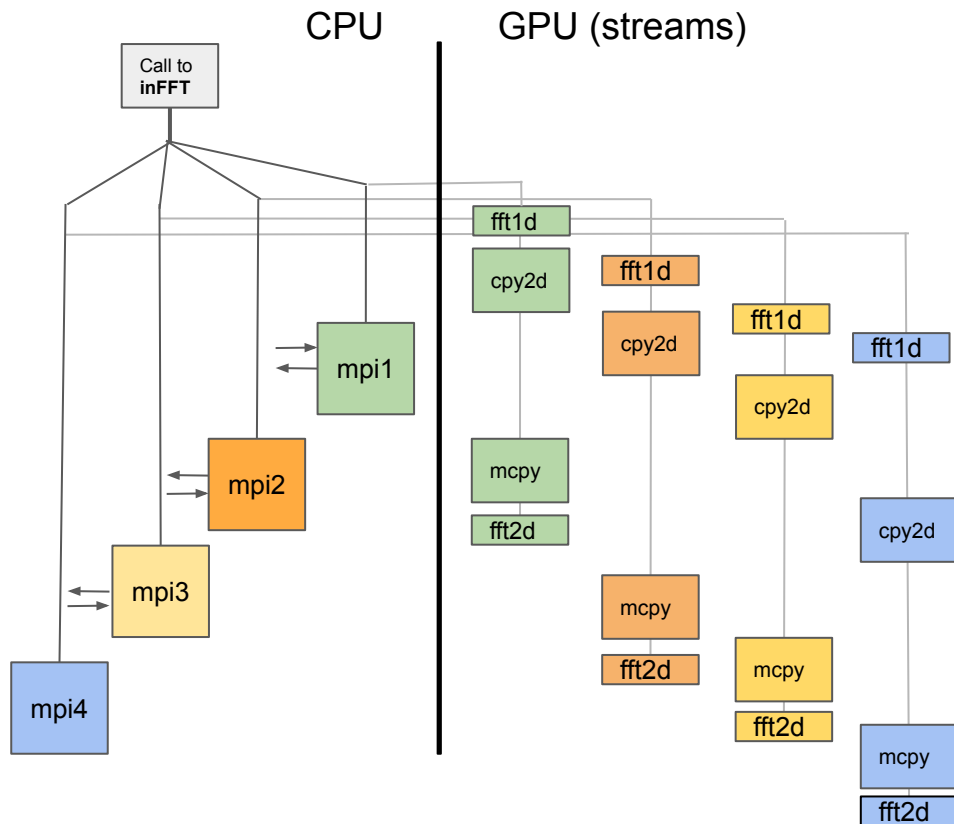# Batched FFTs - CUDA/HIP streams



- **Batched** 3d-FFT of the **wave-function**;

- the input array divided in **4 batches** (on bands);

- 1 stream for **FFTs**, 4 streams for **data movements**;

- 4 **async mpi** communications (ISEND, IRECV).

  Notes:

- **non-asynchronous memcpy**;

- memcpy operations **D2H/H2D** much more time consuming than FFT calls;

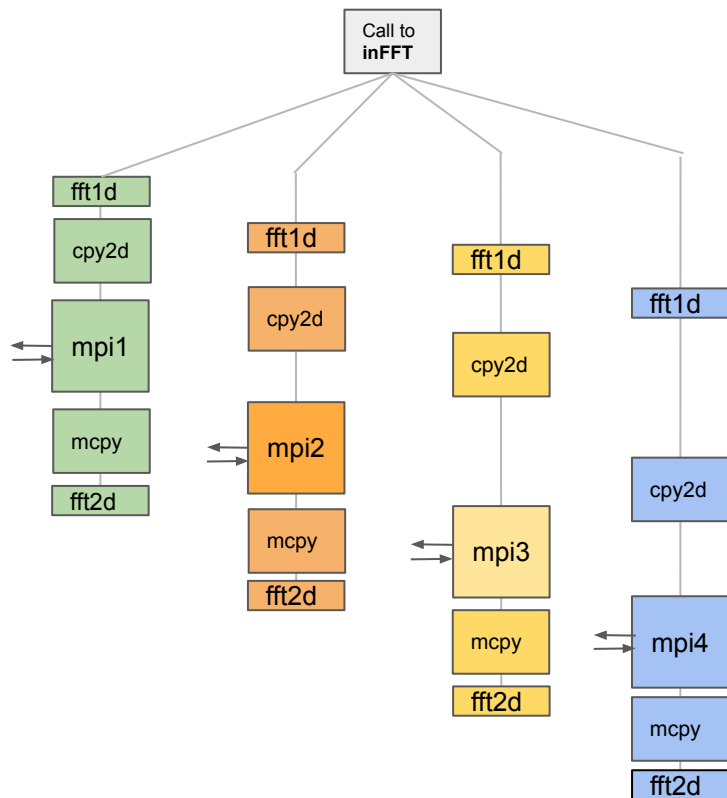- memcpy operations **D2D** same order of FFT calls.

# Batched FFTs - OMP tasks + CUDA/HIP streams



- Need to set up pure OMP porting of batched FFTs for the Intel[®] side;

- setting up a starting scheme by using **omp task** with hip streams and **detach** clause;

```
!$omp parallel
!$omp single

DO j = 0, batchsize-1, dfft%subbatchsize
  currsize = min(dfft%subbatchsize, batchsize - j)
  !
  !$omp task firstprivate(j,currsize) private(i) shared(ptr_callback) detach(event)
  !
  DO i = 0, currsize - 1
    CALL cft_1z_omp( f((j+i)*dfft_nnr + 1:), sticks(me_p), n3, nx3, isgn, &
                     aux(j*dfft_nnr + i*ncpx*nx3 +1:), &
                                   stream=dfft%bstreams(j/dfft%subbatchsize+1) ) )
  ENDDO
  !
  CALL fft_scatter_many_columns_to_planes_store_omp( dfft, aux(j*dfft_nnr+1:), nx3, &
                                   dfft_nnr, f(j*dfft_nnr+1:), &
                                   sticks, dfft%nr3p, isgn, currsize, &
                                   j/dfft%subbatchsize+1 )
  !
  iadc = hipStreamAddCallback(dfft%bstreams(j/dfft%subbatchsize+1),ptr_callback ,c_loc(event), 0)
  !
  !$omp end task
  !
ENDDO
!
```
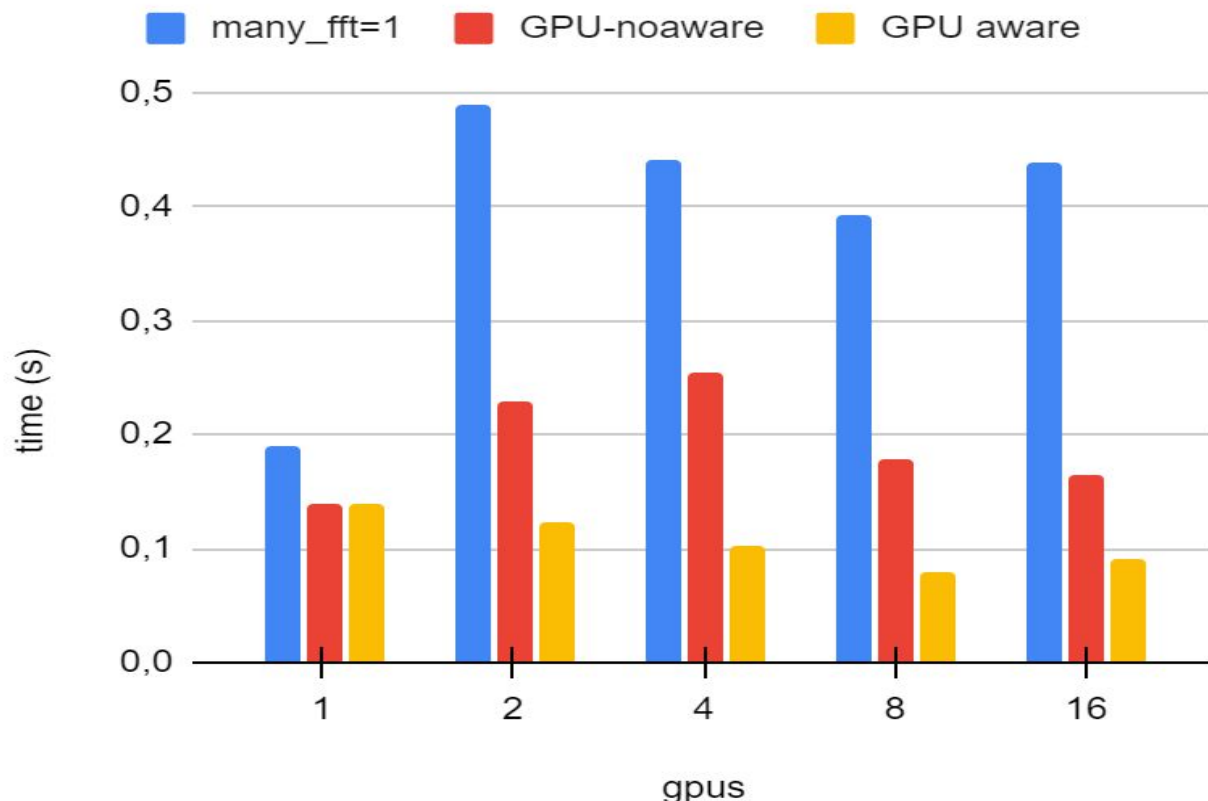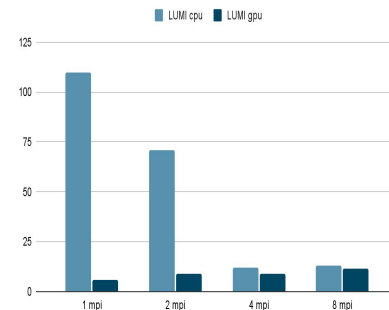
# Batched FFTs - OMP tasks + dep.



- Starting point: **oneMKL** does not get explicit streams as input;

- Simplest scheme given by **n tasks** associated to **n subbatches;**

- **still in progress**
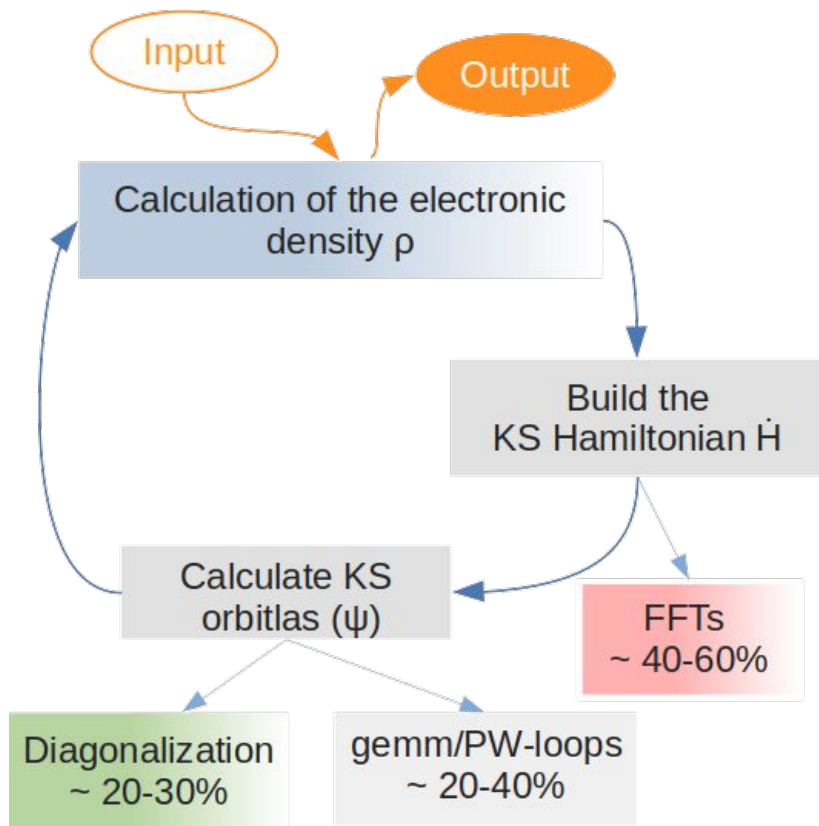
# Batched FFTs - performance



vloc_psi/call

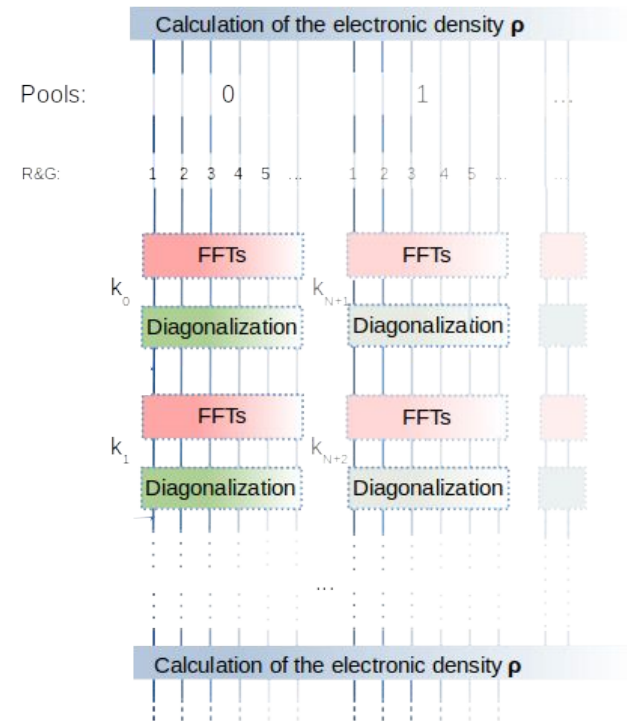- Gold surface;
- 112 atoms;
- ~1600 electrons.
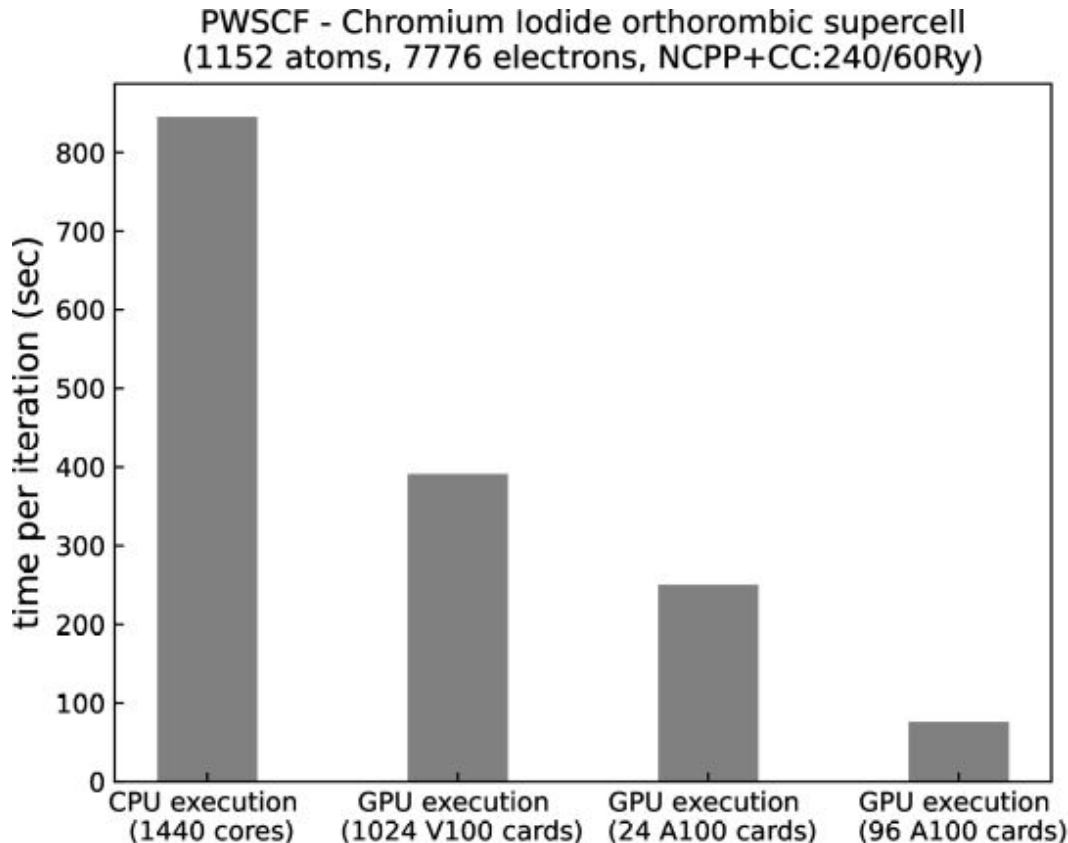
- vloc_psi only.

# QE codes and performance results

Input

Output

Calculation of the electronic density ρ

Build the KS Hamiltonian Ḣ

Calculate KS orbitlas (ψ)

FFTs
~ 40-60%

Diagonalization
~ 20-30%

gemm/PW-loops
~ 20-40%

MPI + OpenMP

Calculation of the electronic density ρ

Pools: 0 1 ...

R&G: 1 2 3 4 5 ... 1 2 3 4 5 ...

FFTs    FFTs

$k_0$    $k_{N+1}$

Diagonalization    Diagonalization

FFTs    FFTs

$k_1$    $k_{N+2}$

Diagonalization    Diagonalization

...

Calculation of the electronic density ρ

# PWscf: performance



PWSCF - Chromium Iodide orthorombic supercell
(1152 atoms, 7776 electrons, NCPP+CC:240/60Ry)

- ◆ 'only' **16 GB** for **V100** (so no pools) but still 1024 GPUs better than optimized run on 1440 cores (Marconi m100);

- ◆ **80 GB** for **A100** (Selene): 24 GPUs (no pools) better than 1024 V100;

- ◆ **~3x** speed-up with **96 A100** (3 pools) **vs 24 A100** (no pools).

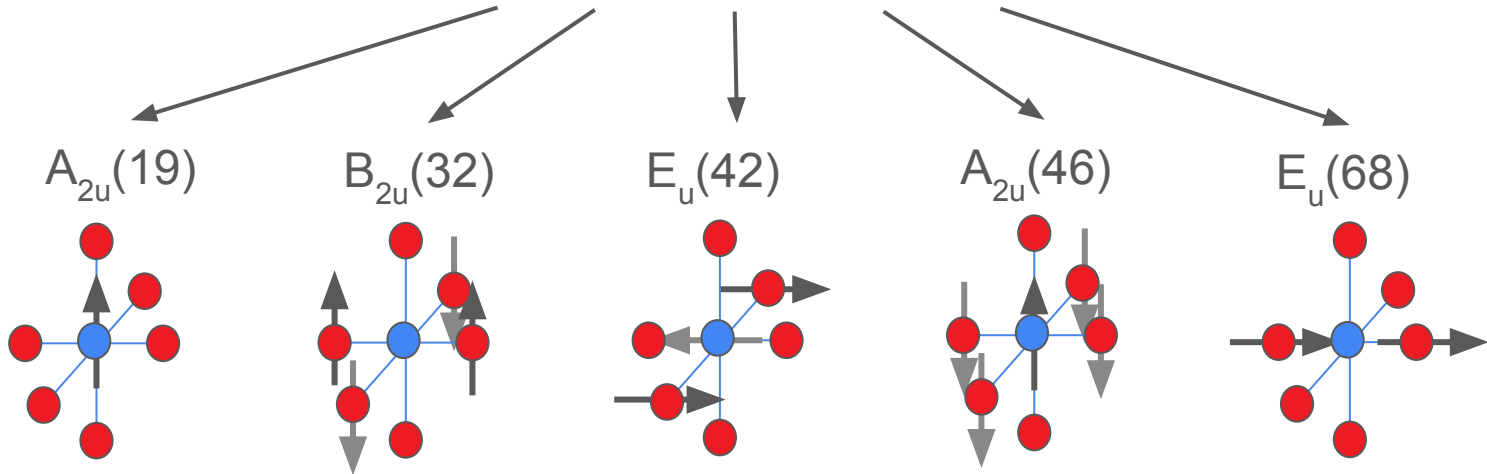The **PHonon** code works for a rather wide variety of systems and methods:

✓ **Insulators** (also polar insulators, with LO-TO splitting)

✓ **Metals**

✓ **Magnetic systems** at the scalar relativistic collinear level

✓ Spin-orbit coupling (fully relativistic approach)

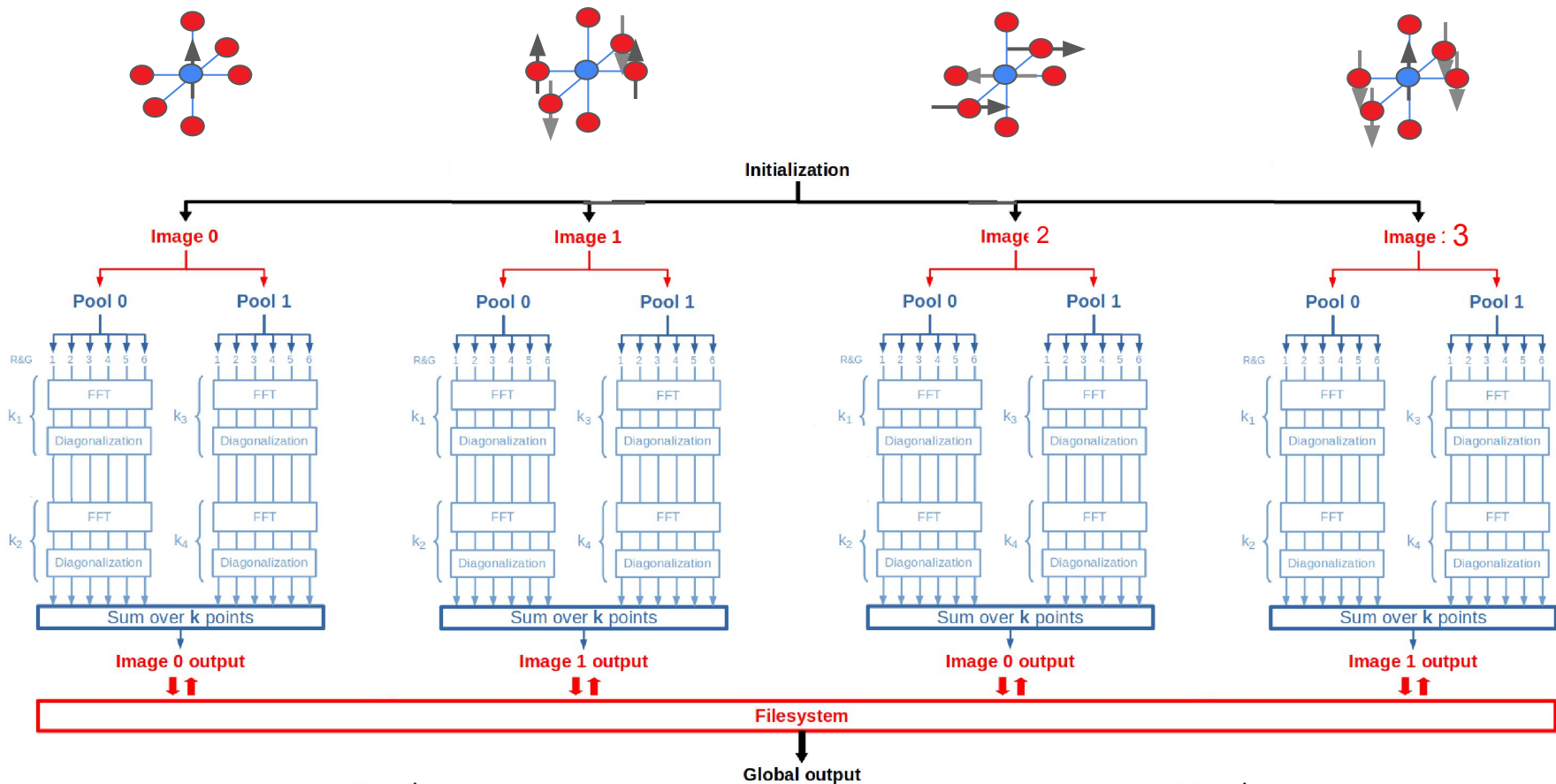✓ Electric field calculations: Born effective charges, **dielectric tensor**

# PH: phonon

*Interatomic Force Constants (IFC)*

$$\tilde{D}_{s\alpha,s'\beta}(\mathbf{q}) = \frac{1}{\sqrt{M_s M_{s'}}} \sum_{\mathbf{R},\mathbf{R'}} \boxed{\frac{\partial^2 E_{tot}}{\partial \mathbf{u}_{s\alpha}(\mathbf{R}) \partial \mathbf{u}_{s'\beta}(\mathbf{R'})}} e^{i\mathbf{q}(\mathbf{R'}-\mathbf{R})}$$
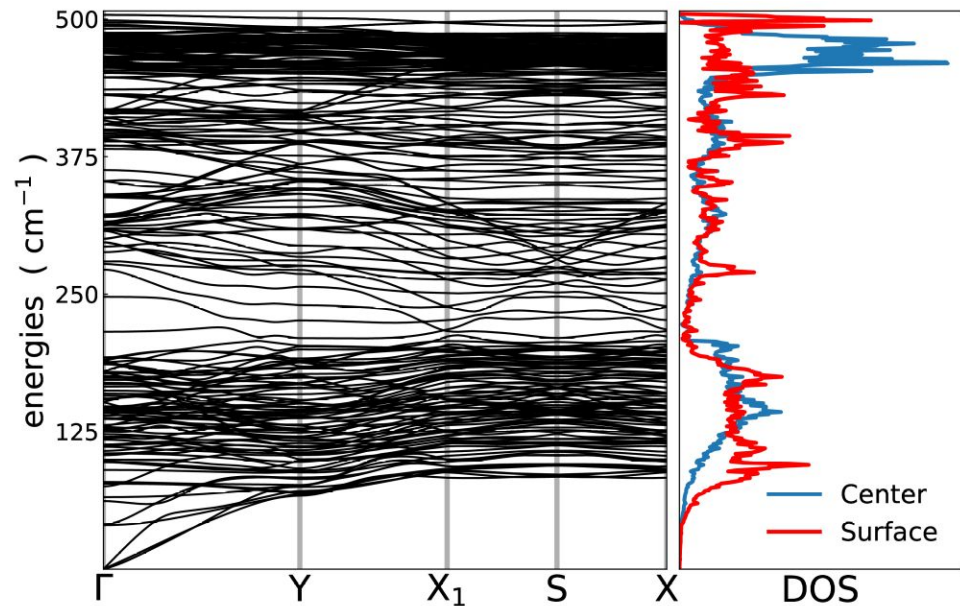
Irreducible vibration modes



$A_{2u}(19)$    $B_{2u}(32)$    $E_u(42)$    $A_{2u}(46)$    $E_u(68)$

# PH: parallel scheme



Initialization

Image 0      Image 1      Image : 2      Image : 3

Pool 0   Pool 1    Pool 0   Pool 1    Pool 0   Pool 1    Pool 0   Pool 1

R&G   1 2 3 4 5 6   1 2 3 4 5 6

$k_1$   FFT    $k_3$   FFT

Diagonalization    Diagonalization

$k_2$   FFT    $k_4$   FFT

Diagonalization    Diagonalization

Sum over **k** points    Sum over **k** points    Sum over **k** points    Sum over **k** points

Image 0 output    Image 1 output    Image 0 output    Image 1 output

Filesystem

Global output

# PH: performance

## Silicon 100: phonon dispersion

## Silicon 100: phonon dispersion



PHonon: best CPU vs best GPU Si(100)-32L

Time-Dependent Density-Functional Perturbation Theory (TDDFpT):

✓ optical absorption spectroscopy;

✓ Electron energy loss spectroscopy (EELS);

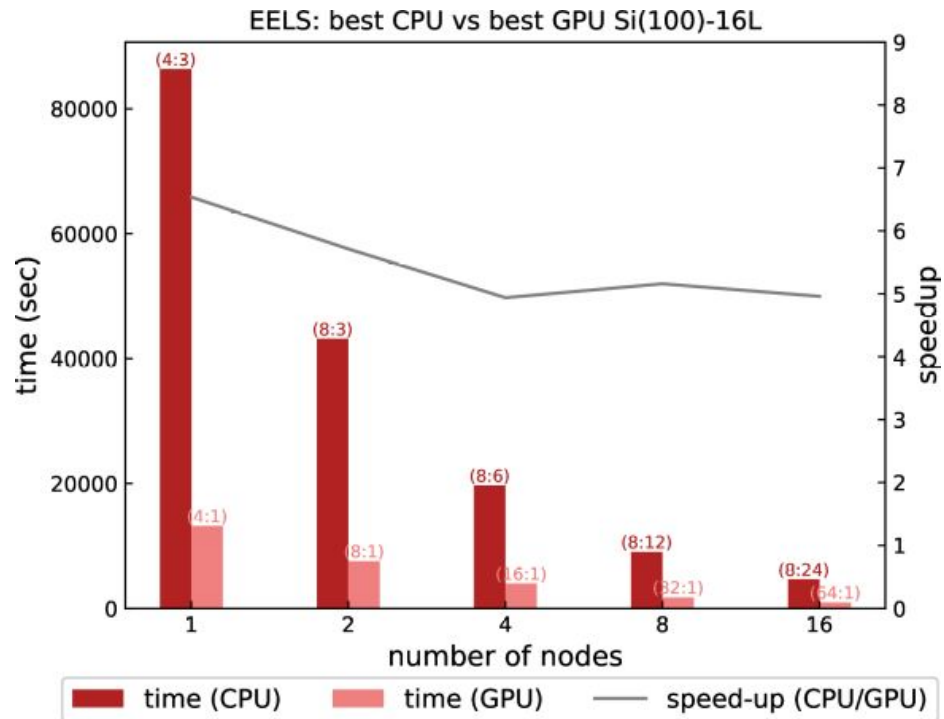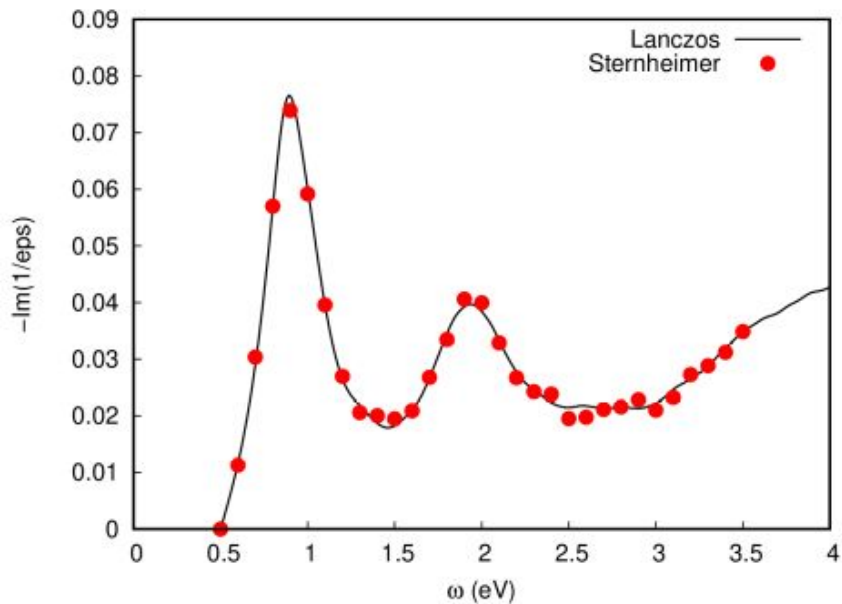✓ Inelastic X-ray scattering (IXS);

✓ Inelastic neutron scattering (INS);

$e^-, n^\circ, h\nu$

$\Psi_i$

$e^-, n^\circ, h\nu$

$\Psi_f$

Dynamical susceptibilities

$$\varphi_{\text{ext}}(t) \longrightarrow A(t) \approx A^\circ + A'(t)$$

$$A'(t) = \int \mathrm{d}t' \chi(t - t') \varphi_{\text{ext}}(t')$$

Time-Dependent Density-Functional Perturbation Theory (TDDFpT):

✓ optical absorption spectroscopy;

✓ Electron energy loss spectroscopy (EELS);

✓ Inelastic X-ray scattering (IXS);

✓ Inelastic neutron scattering (INS);

Silicon 100: spectrum of electron energy loss

# EELS: performance

Silicon 100: spectrum of electron energy loss





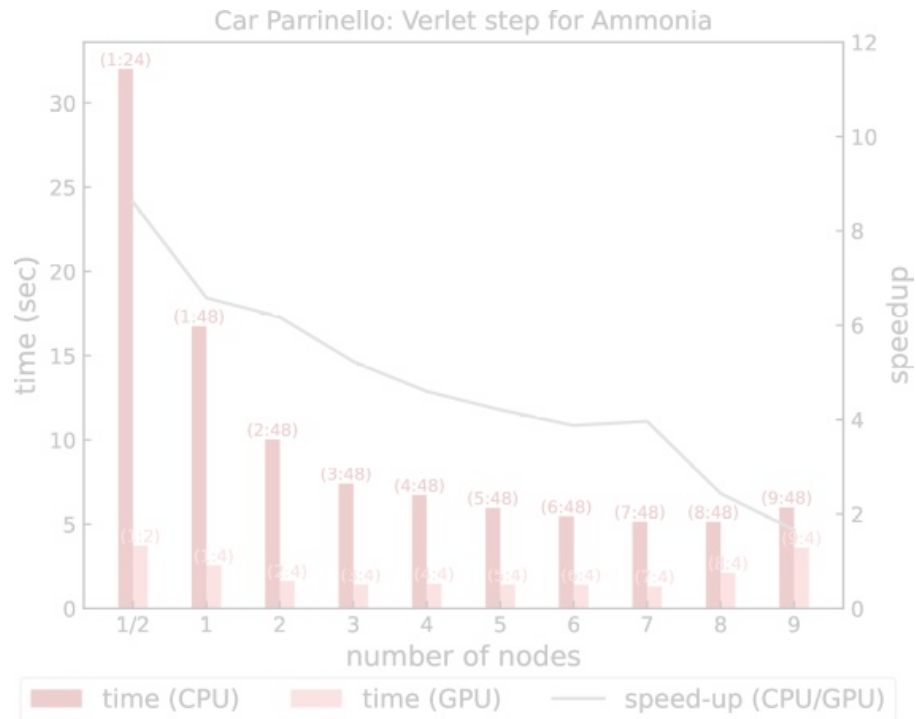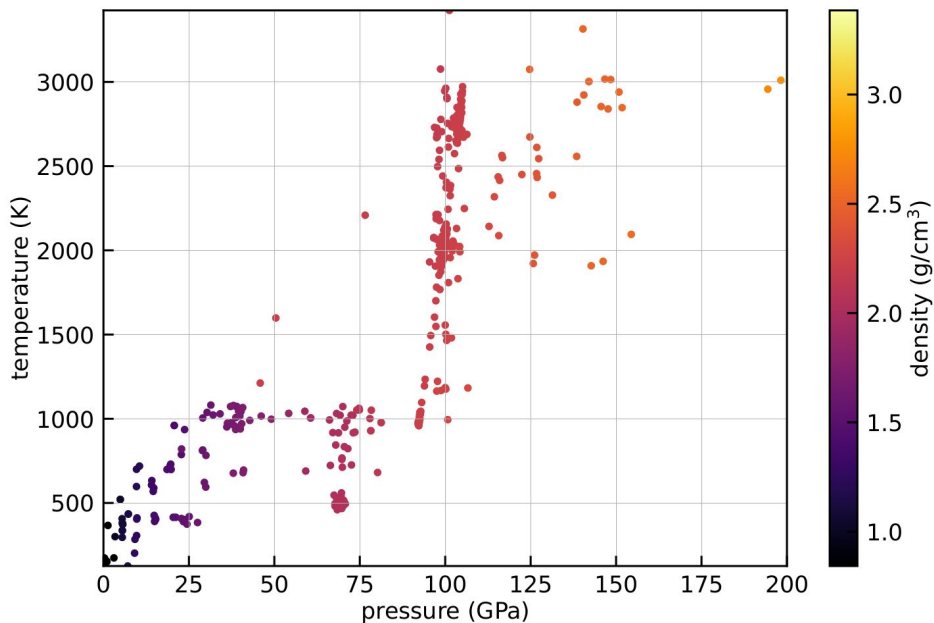EELS: best CPU vs best GPU Si(100)-16L

**Ab-initio** molecular dynamics (**MD**):

- classical molecular dynamics + **QM** electronic structure;

- combines MD with **DFT**;

- accounts for formation or break of **bonds**;

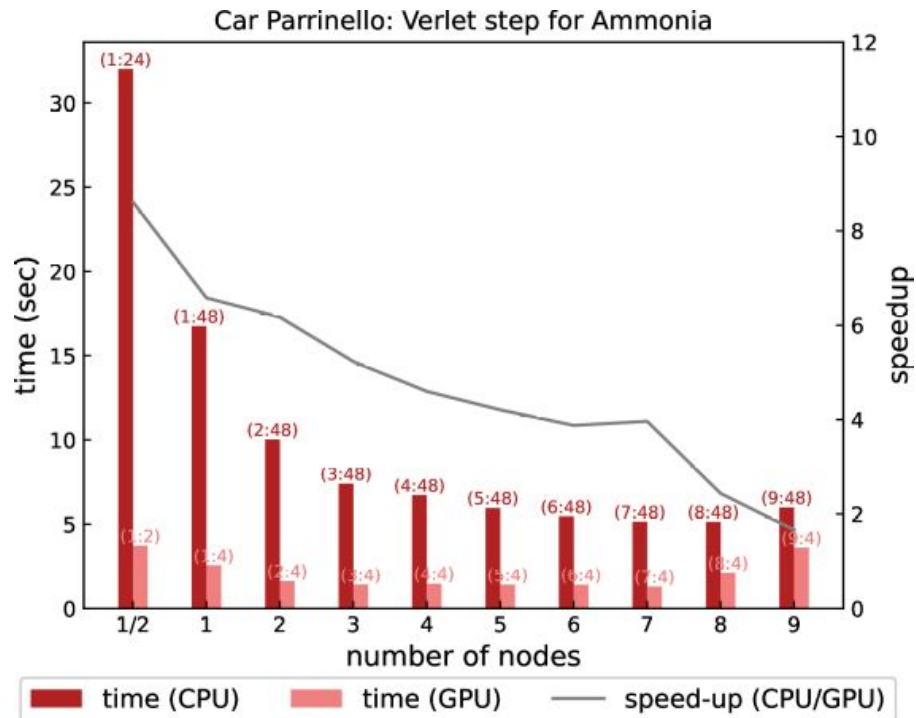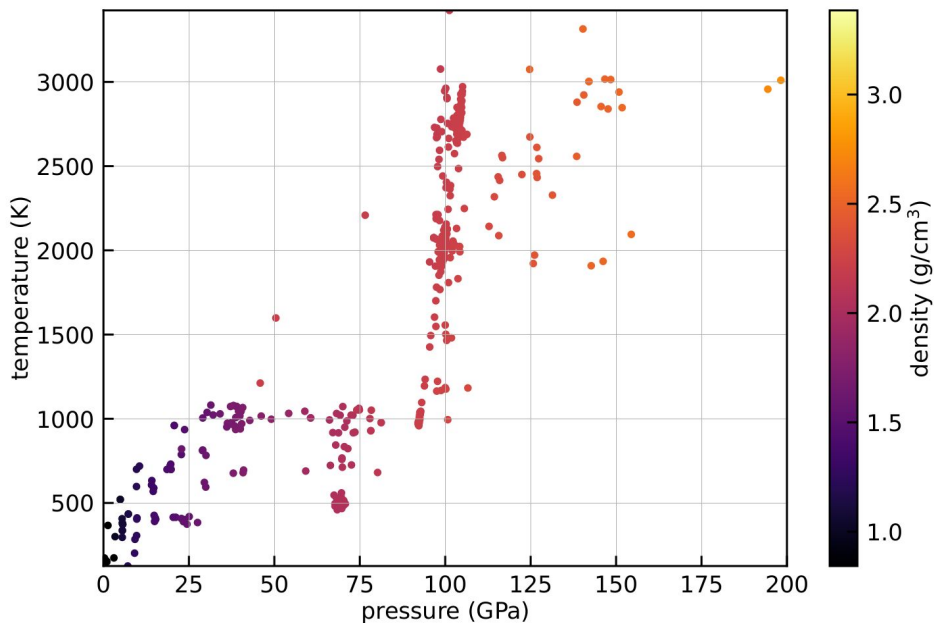- accounts for complex bindings, e.g. **transition metal ions**

Superionic Ammonia: Eq. of state
144 Nitrogen + 432 Hydrogen (1152 e)

# CP: performance

Superionic Ammonia: Eq. of state
144 Nitrogen + 432 Hydrogen (1152 e)



Car Parrinello: Verlet step for Ammonia

# Conclusions

# Summary

- **Modularity** of QE;

- **directive** based porting;

- **CPU and GPU** low level routines **at the same time**;

- **multiple standards** with multiple backends;

- full porting on **Nvidia**[®] side (still transitioning to full openACC);

- ongoing porting on **AMD**[®]**/Intel**[®] side (advanced status on PWscf).

# Outlook

- Full port of **PWscf on AMD® and Intel®** by this year likely;

- **merge** Nvidia® standard branch with AMD®/Intel® one;

- port of **codes other than PW** on AMD®/Intel®;

- **extensive benchmarks** on the main the euroHPC machines;

- incorporation of **devXlib**;

- new features….

# Acknowledgments

# Thank you