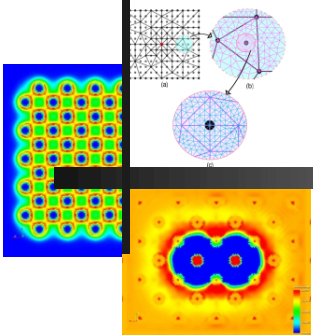# Fast, scalable and accurate finite-element based *ab initio* calculations using mixed precision computing

Vikram Gavini

*Department of Mechanical Engineering*
*Department of Materials Science and Engineering*
*University of Michigan, Ann Arbor*

*Collaborators: Sambit Das (U. Mich); Phani Motamarri (U. Mich);*
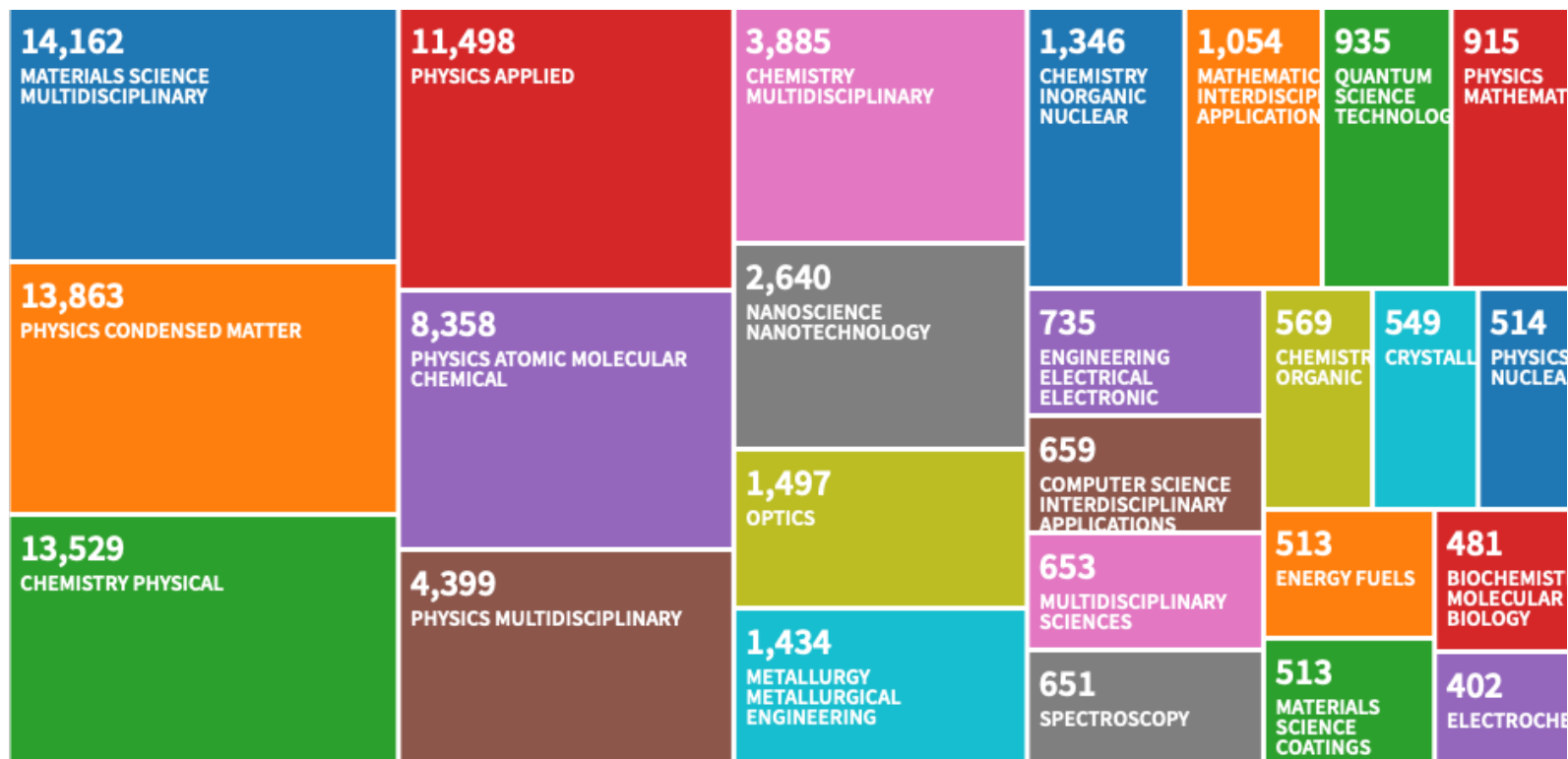*Bruno Turcksin (ORNL); Ying Wai Li (ORNL/LANL); Brent Leback (Nvidia)*

SMC 2019

# Impact of Density Functional Theory

## Citations to seminal work of Walter Kohn (1964,1965)



| | | | | | | |
|---|---|---|---|---|---|---|
| 14,162 MATERIALS SCIENCE MULTIDISCIPLINARY | 11,498 PHYSICS APPLIED | 3,885 CHEMISTRY MULTIDISCIPLINARY | 1,346 CHEMISTRY INORGANIC NUCLEAR | 1,054 MATHEMATIC INTERDISCIPL APPLICATION | 935 QUANTUM SCIENCE TECHNOLOG | 915 PHYSICS MATHEMAT |

14,162 MATERIALS SCIENCE MULTIDISCIPLINARY

13,863 PHYSICS CONDENSED MATTER

13,529 CHEMISTRY PHYSICAL

11,498 PHYSICS APPLIED

8,358 PHYSICS ATOMIC MOLECULAR CHEMICAL

4,399 PHYSICS MULTIDISCIPLINARY

3,885 CHEMISTRY MULTIDISCIPLINARY

2,640 NANOSCIENCE NANOTECHNOLOGY

1,497 OPTICS

1,434 METALLURGY METALLURGICAL ENGINEERING

1,346 CHEMISTRY INORGANIC NUCLEAR

735 ENGINEERING ELECTRICAL ELECTRONIC

659 COMPUTER SCIENCE INTERDISCIPLINARY APPLICATIONS

653 MULTIDISCIPLINARY SCIENCES

651 SPECTROSCOPY

1,054 MATHEMATIC INTERDISCIPL APPLICATION

935 QUANTUM SCIENCE TECHNOLOG

915 PHYSICS MATHEMAT

569 CHEMISTR ORGANIC

549 CRYSTALL

514 PHYSICS NUCLEA

513 ENERGY FUELS

481 BIOCHEMIST MOLECULAR BIOLOGY

513 MATERIALS SCIENCE COATINGS

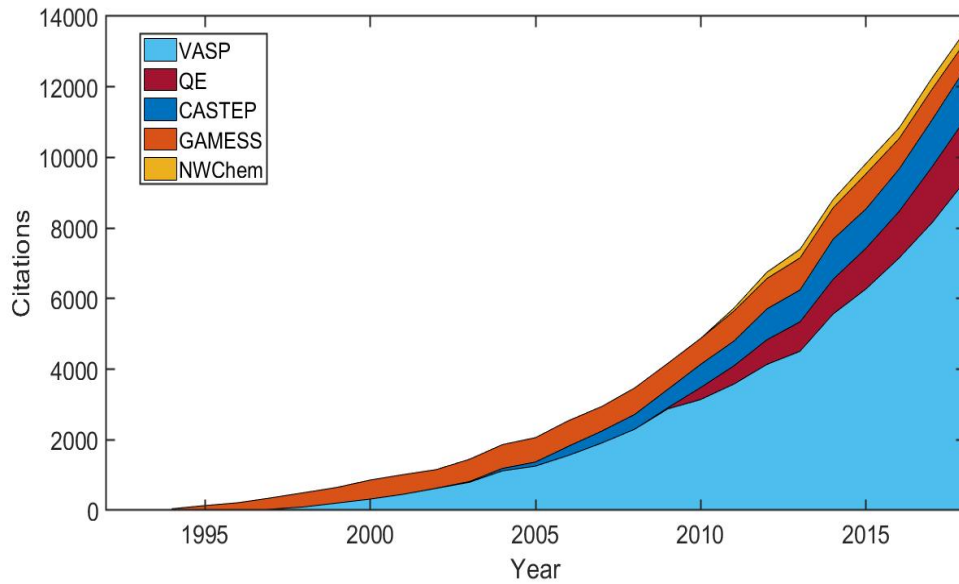402 ELECTROCHE

Data compiled from Web of Science

**12 of the 100 most-cited papers in scientific literature pertain to DFT!**
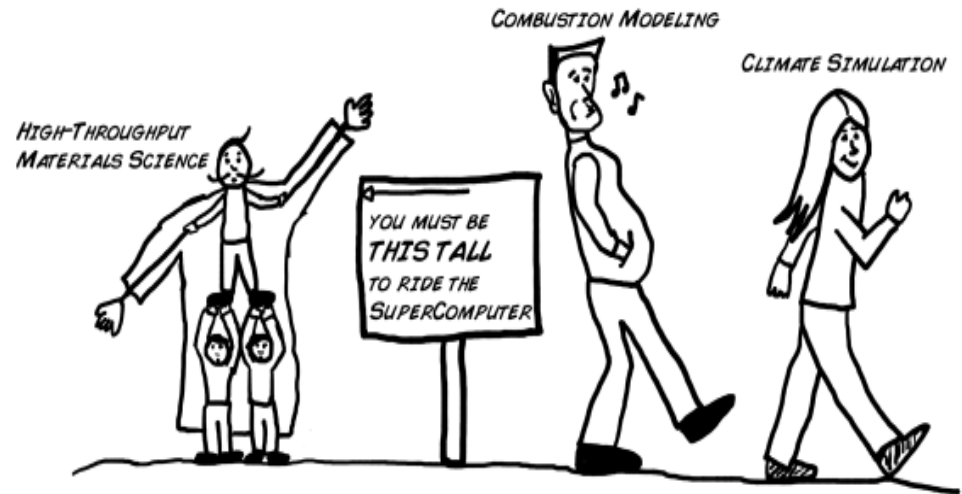
(Nature **514**, 550 (2014))

# DFT codes

~100 available DFT codes developed since 1980

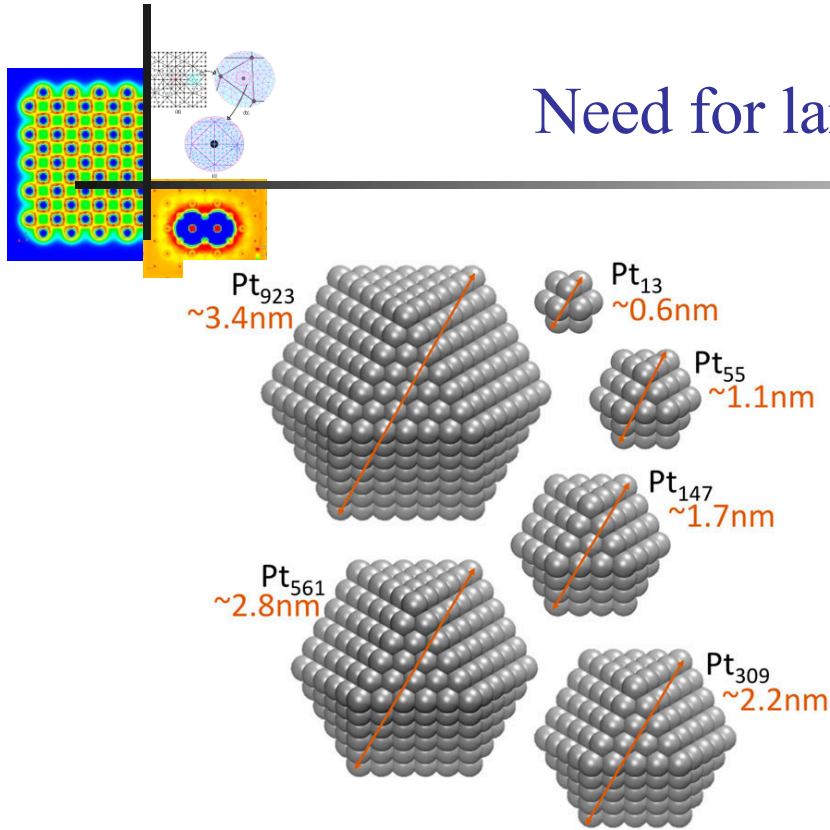Relationship to HPC



Data compiled from Web of Science
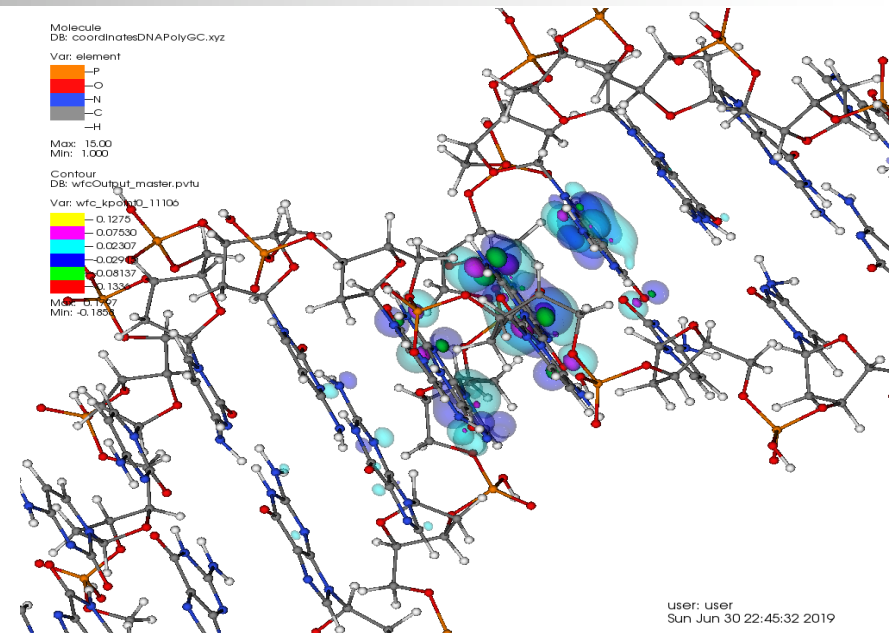
Courtesy: Anubhav Jain

## Key Issues

❖ Lack of good parallel scalability of existing DFT codes
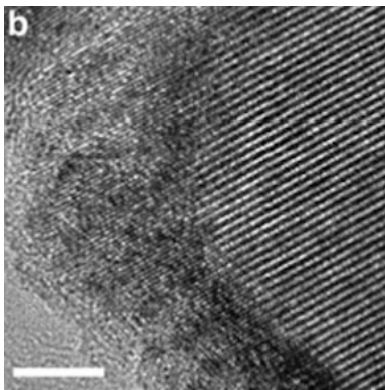❖ Computational complexity of DFT calculations ($O(N^3)$)
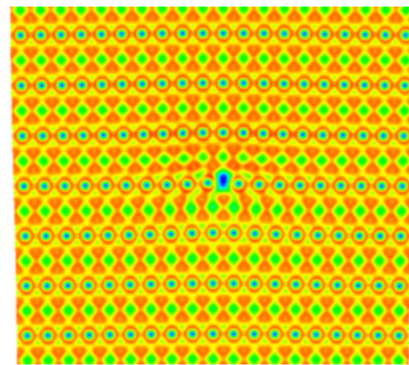
SMC 2019

# Need for large scale DFT calculations
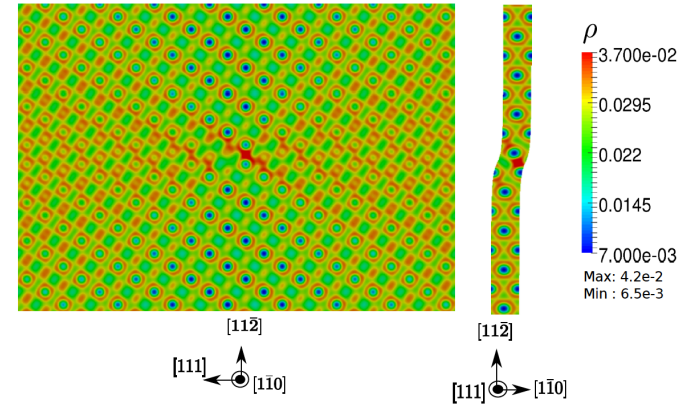


**Chemical properties of nanoparticles**

Pt$_{923}$ ~3.4nm
Pt$_{13}$ ~0.6nm
Pt$_{55}$ ~1.1nm
Pt$_{147}$ ~1.7nm
Pt$_{561}$ ~2.8nm
Pt$_{309}$ ~2.2nm

**Biological systems**

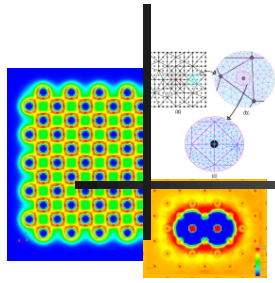Rocksalt phase formation during Litihiation of Magnetite
He et. al, Nature Comm, 2016

Edge dislocation:
Iyer et al. J. Mech, Phys. Solids (2015)

Screw dislocation:
Das & Gavini J. Mech, Phys. Solids (2017)

[11$\bar{2}$]
[111] [1$\bar{1}$0]

[11$\bar{2}$]
[111] [1$\bar{1}$0]

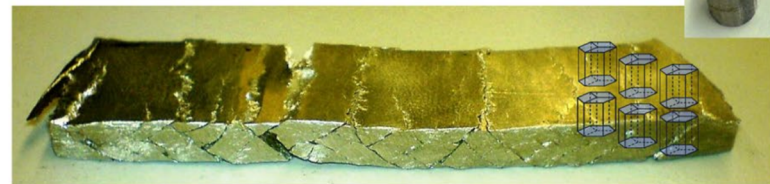**Defects in Materials**
SMC 2019

4

# Technological challenge of low ductility in Mg

➢ **Magnesium is the lightest structural material with high strength to weight ratio**
  - ❖ 75% lighter than Steel and 30% lighter than Aluminum

➢ **Every 10% reduction in the weight of a vehicle will result in 6-8% increase in fuel efficiency.**
  - ❖ Important implications to fuel efficiency and reducing carbon footprint

➢ **Low ductility key issue in the manufacturability of structural components. Main limitation in the adoptability of Mg and Mg alloys in automotive and aerospace sectors.** (T.M. Pollock, *Science* **328**, 986-987 (2010))

Courtesy: https://www.audi-technology-portal.de/en/body
Current state of art: Hybrid Steel and Aluminum construction

Brittle Mg (pure)

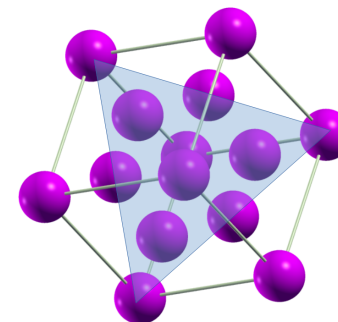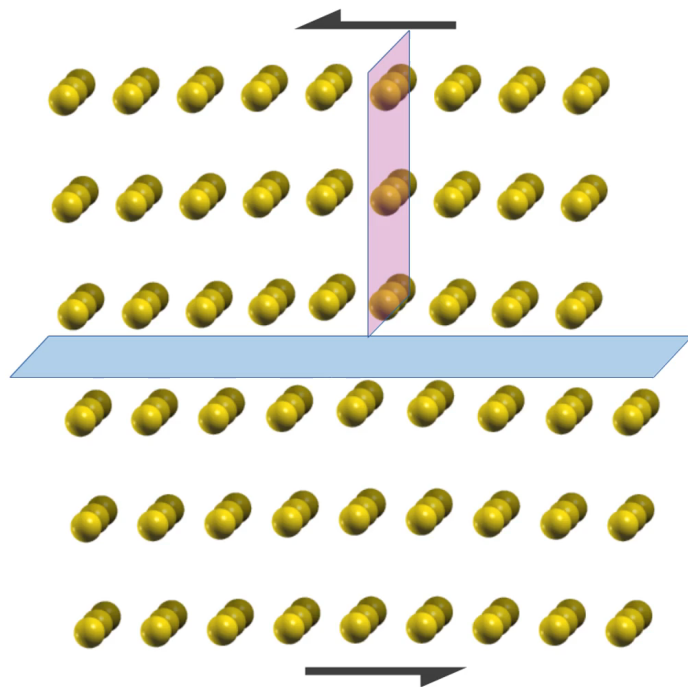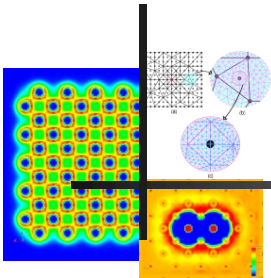10% CR
CR: cold rolled

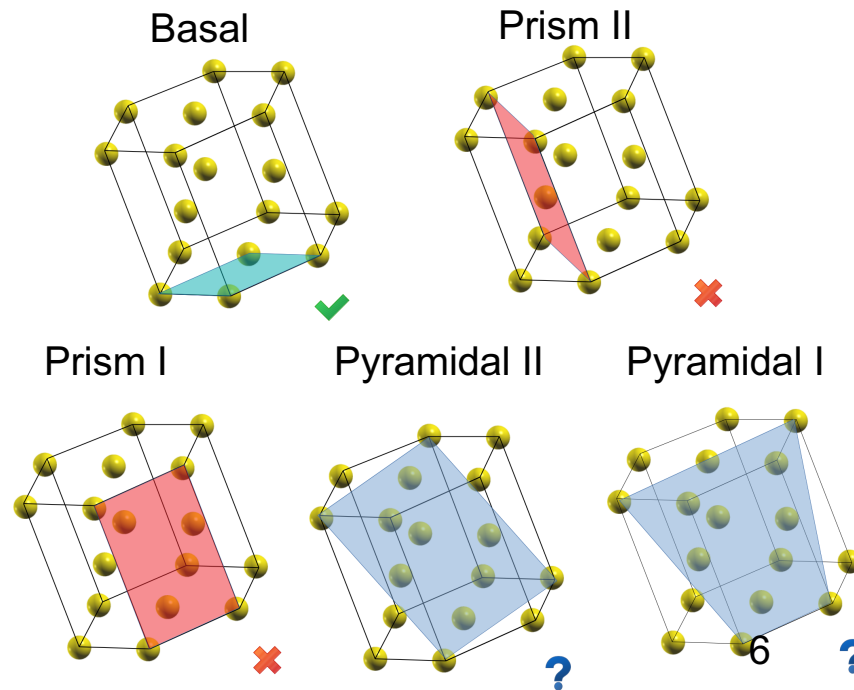S. Sandlöbes et al. Scientific Reports 7, 10458 (2017).
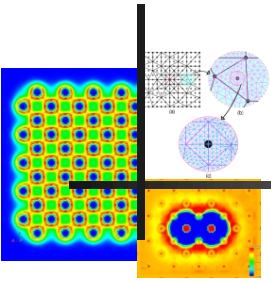
# Technological challenge of low ductility in Mg



4 slip planes in Face Centered Cubic
Crystals→ higher ductility



- ❖ Dislocations are energetically more favorable to reside on certain slip systems. (**Energetics**)

- ❖ Dislocation glide occurs after the applied shear stress is greater than the Perils barrier.

   (**Activation barrier**)

- ❖ More the number of slip systems where dislocation can glide easily higher is the ductility.

Basal    Prism II

Prism I    Pyramidal II    Pyramidal I

6

# Density Functional Theory

Kohn-Sham eigenvalue problem:

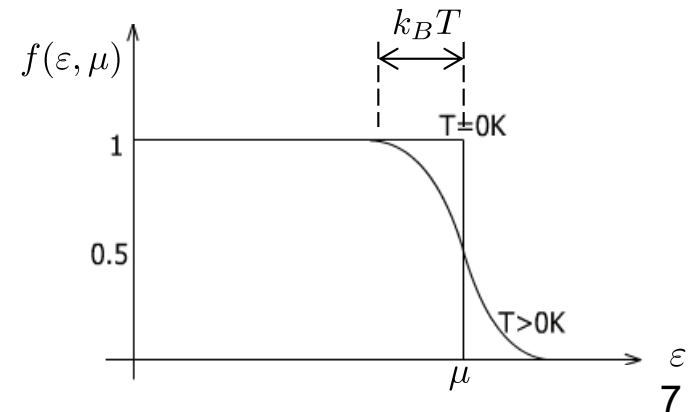$$\left(-\frac{1}{2}\nabla^2 + V_{eff}\right)\psi_i = \epsilon_i\psi_i$$

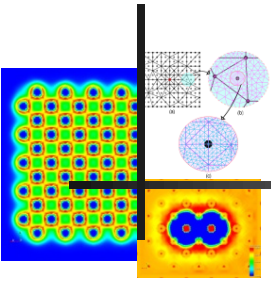Self consistent iteration
(Kohn-Sham map)

$$\rho = \sum_i f_i|\psi_i|^2, \qquad V_{eff}(\mathbf{r}) = V_H(\rho(\mathbf{r})) + V_{xc}(\rho(\mathbf{r})) + V_{ext}(\mathbf{R})$$

$$T_s(\Psi) = \frac{1}{2}\sum_i f_i \int |\nabla\psi_i(\mathbf{r})|^2 d\mathbf{r} \quad E_0(\Psi) = T_s(\Psi) + E_{xc}(\rho) + E_H(\rho) + E_{ext}(\rho) + E_{zz}$$

Orbital occupancy:

$$f_i = f(\varepsilon_i, \mu) = \frac{1}{1 + e^{\frac{\varepsilon_i - \mu}{k_B T}}} \qquad \sum_i f_i = N$$
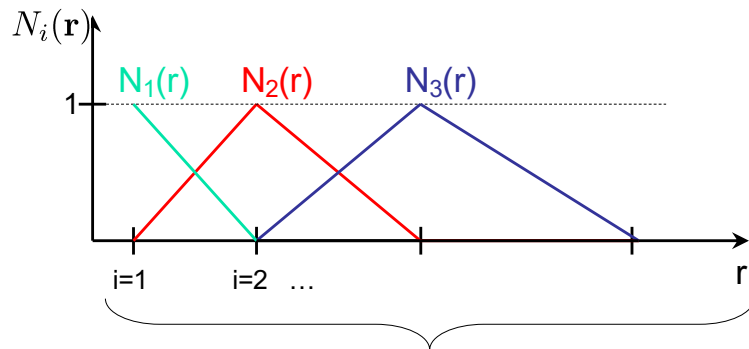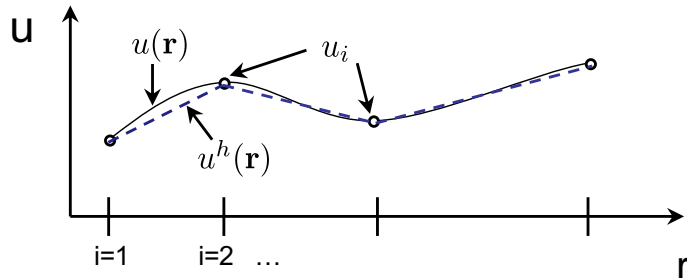
# DFT – Finite Element discretization



➢ Use finite-element basis for computing –

$$\psi_k^h(\mathbf{r}) = \sum_i \psi_{k_i} N_i(\mathbf{r}) \ \ k = 1, \ldots, N \,, \qquad \phi^h(\mathbf{r}) = \sum_i \phi_i N_i(\mathbf{r})$$

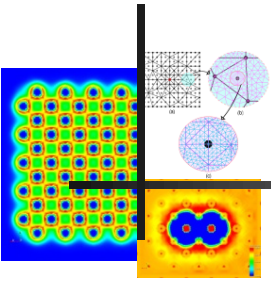$\psi_{k_i}, \ \phi_i \ldots$ – Nodal values
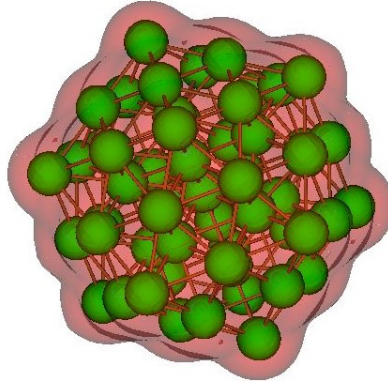$N_i(\mathbf{r})$ – Shape functions



## Features of FE basis

➢ Systematic convergence
  ❖ Element size
  ❖ Polynomial order

➢ Adaptive refinement

➢ Complex geometries and boundary conditions

➢ Potential for excellent parallel scalability

By changing the positioning of the nodes the spatial resolution of basis can be changed/adapted
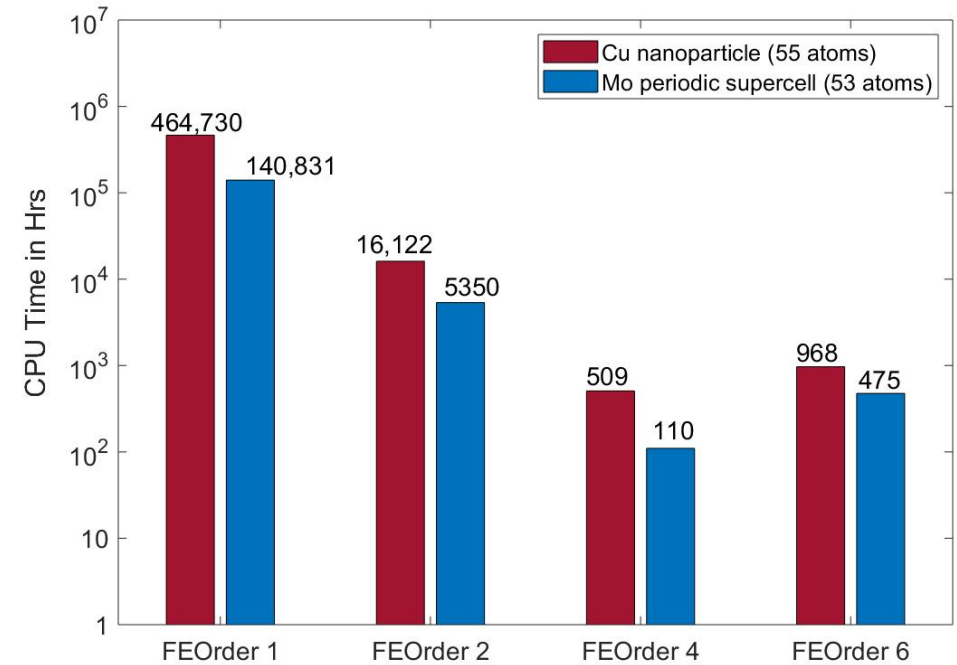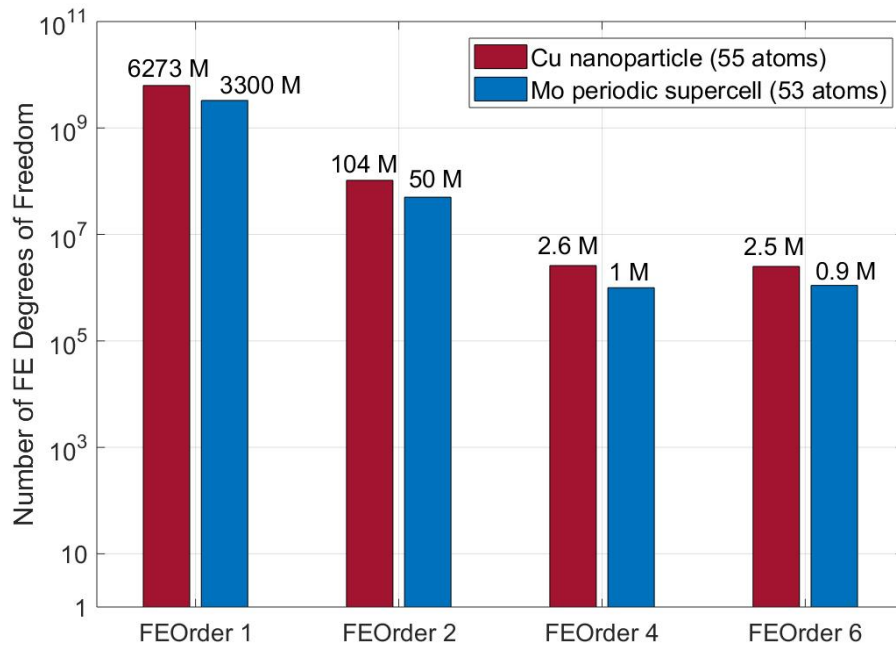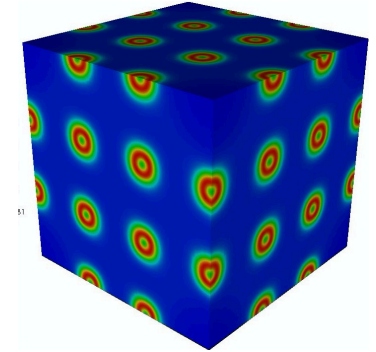
8

# Higher (polynomial) order FE basis

I. Cu nanoparticle 55 atoms

II. Mo periodic supercell w/ vacancy 53 atoms



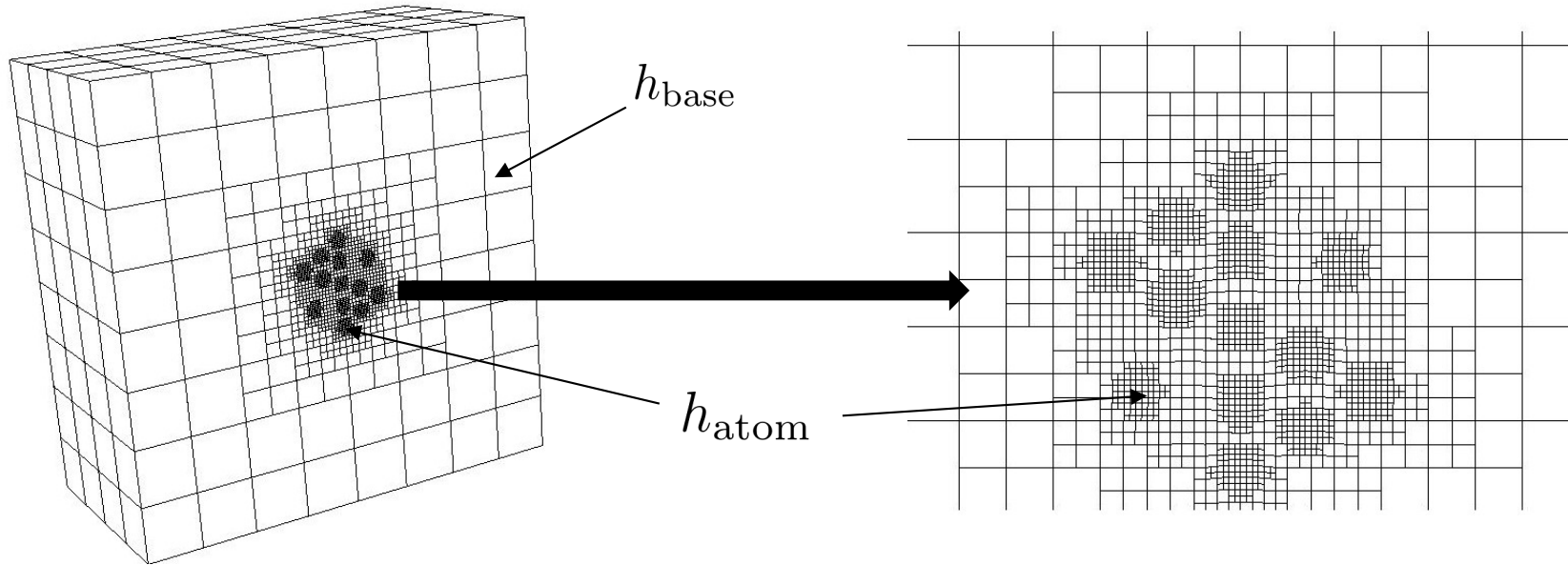~1000x advantage by using higher-order FE basis !

# Spatial adaptivity of the FE basis

➤ **Error Analysis:** $$|E - E^h| \le C \left( \sum_i |\bar{\psi}_i - \bar{\psi}_i^h|^2_{1,\Omega} \right) \le \mathcal{C} \sum_e h_e^{2k} \left[ \sum_i |\bar{\psi}_i|^2_{k+1,\Omega_e} \right]$$

➤ **Optimal FE mesh:** $$\min_h \int_\Omega \left\{ h^{2k}(\mathbf{x}) \left[ \sum_i |D^{k+1}\bar{\psi}_i(\mathbf{x})|^2 \right] \right\} d\mathbf{x} \qquad \text{subject to} \int_\Omega \frac{d\mathbf{x}}{h^3(\mathbf{x})} = N_E$$
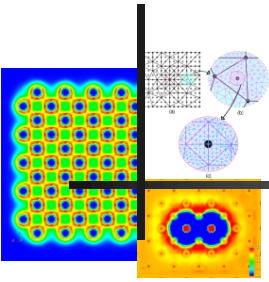


$h_{\text{base}}$

$h_{\text{atom}}$

| System Type pyr II dislocation | DoFs Uniform Mesh | DoFs for Adaptive Mesh |
|---|---|---|
| 1848 atom Mg | 347,206,614 | 55,112,161 |
| 6164 atom Mg | 892,047,315 | 179,034,231 |

# Eigen-space computation: Chebyshev acceleration

**Kohn-Sham eigenvalue problem:** $\widetilde{\mathbf{H}}\widetilde{\psi}_k = \epsilon_k \widetilde{\psi}_k$ for $k = 1,2,...N$ $(N \sim 1.1 N_e/2)$



$$\bar{\mathbf{H}} = c_1 \widetilde{\mathbf{H}} + c_2$$
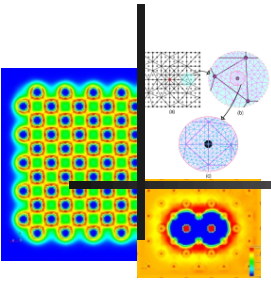
**Chebyshev Filtering:** $T_m(\bar{\mathbf{H}})\widetilde{\mathbf{\Psi}} = \widetilde{\mathbf{\Psi}}_F$

$$T_m(\bar{\mathbf{H}})\mathbf{X} = [2\bar{\mathbf{H}}T_{m-1}(\bar{\mathbf{H}}) - T_{m-2}(\bar{\mathbf{H}})]\mathbf{X}$$
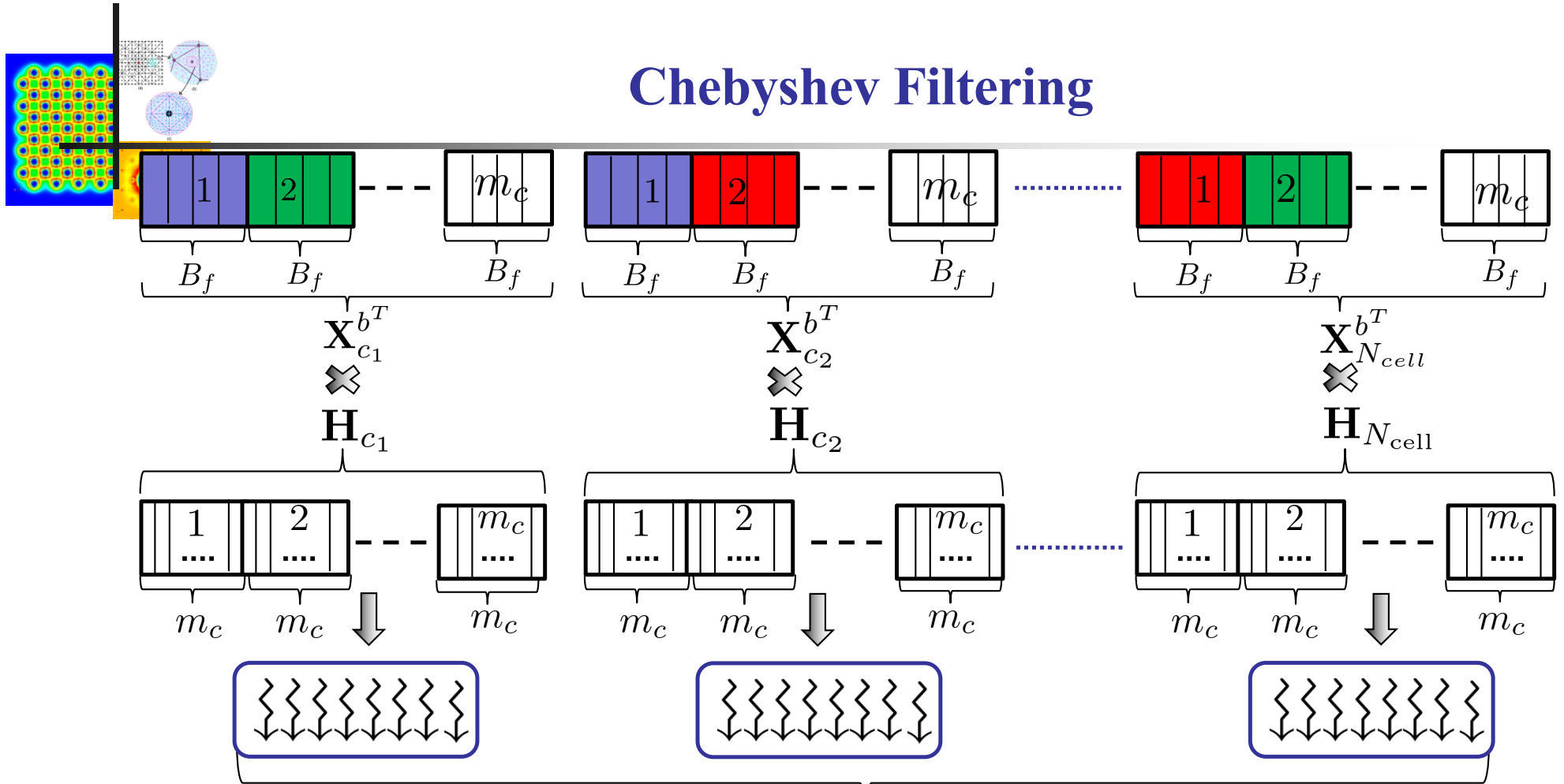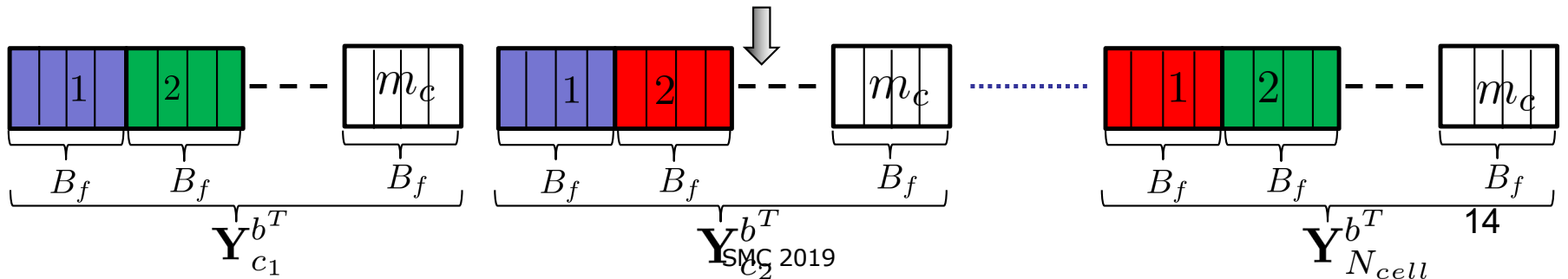
# Numerical algorithm

1. Start with initial guess for electron density $\rho_{in}^h(\mathbf{r}) = \rho_0(\mathbf{r})$ and the initial wavefunctions

2. Compute the discrete Hamiltonian $\bar{\mathbf{H}}$ using the input electron density $\rho_{in}^h$

3. **CF:** Chebyshev filtering: $\widetilde{\boldsymbol{\Psi}}_F = T_m(\bar{\mathbf{H}})\widetilde{\boldsymbol{\Psi}}$

4. **Orthonormalize** CF basis: $\widetilde{\boldsymbol{\Psi}}_F \rightarrow \widetilde{\boldsymbol{\Psi}}_F^o$

5. **Rayleigh-Ritz procedure**:
   - ❖ Compute projected Hamiltonian: $\hat{\mathbf{H}} = \widetilde{\boldsymbol{\Psi}}_F^{o\dagger}\widetilde{\mathbf{H}}\widetilde{\boldsymbol{\Psi}}_F^o$
   - ❖ Diagonalize $\hat{\mathbf{H}}$: $\hat{\mathbf{H}}\mathbf{Q} = \mathbf{Q}\mathbf{D}$
   - ❖ Subspace rotation: $\widetilde{\boldsymbol{\Psi}}^{\mathbf{R}} = \widetilde{\boldsymbol{\Psi}}_{\mathbf{F}}^o\mathbf{Q}$

6. Compute electron density $\rho_{\text{out}}^h(\mathbf{x}) = 2\sum_{i=1}^{N} f(\epsilon_i^h, \mu)|\psi_i^h(\mathbf{x})|^2$

7. If $||\rho_{\text{out}}^h(\mathbf{r}) - \rho_{in}^h(\mathbf{r})|| < tol$, EXIT; else, compute new $\rho_{in}^h$ using a mixing scheme and go to (2).
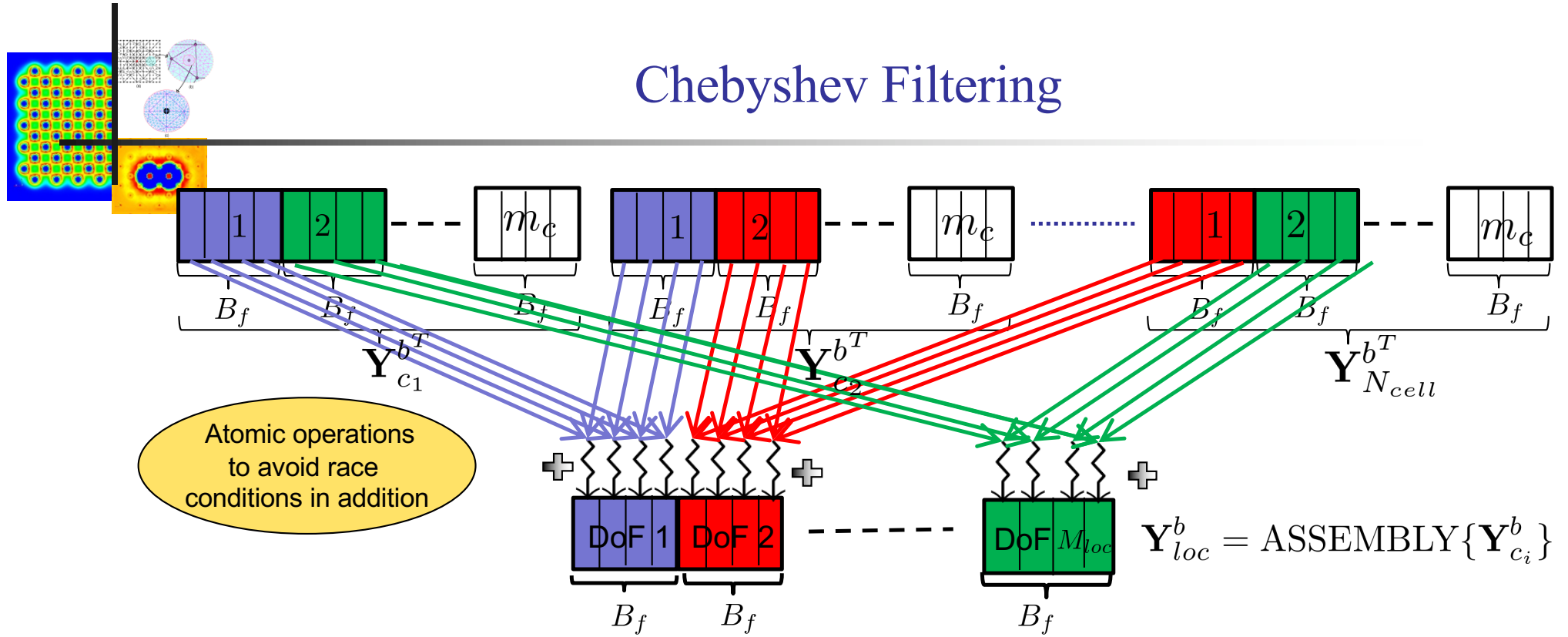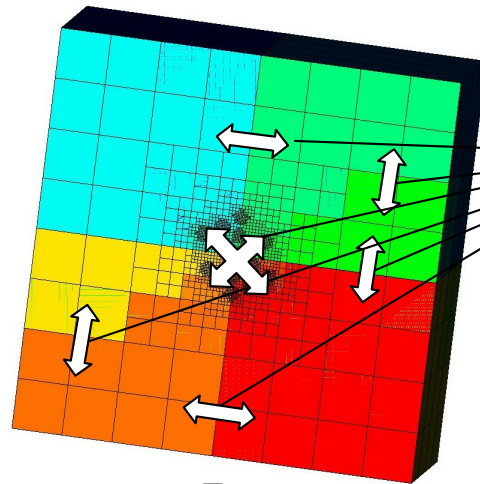
# Chebyshev Filtering

$$\mathbf{Y} = \mathbf{H}\,\mathbf{X}$$

$\mathbf{H} \rightarrow$ Sparse Matrix $(M \times M)$
$\mathbf{X} \rightarrow$ Dense Matrix $(M \times N)$
$\mathbf{Y} \rightarrow$ Dense Matrix $(M \times N)$

FE Cell $c_i$

$[\mathbf{H}_{c_i}]_{m_c \times m_c}\,[\mathbf{X}^b_{c_i}]_{m_c \times B_f}$

$N_{\text{cell}}$ : Number of FE cells



$N$

DoF 1    DoF 2    DoF $M_{loc}$    $\mathbf{X}_{loc}$

DoF 1  DoF 2    DoF$M_{loc}$    $\mathbf{X}^b_{loc}$

$B_f$

1  2  $m_c$  1  2  $m_c$  1  2  $m_c$

$B_f$  $B_f$  $B_f$  $B_f$  $B_f$  $B_f$  $B_f$  $B_f$

$\mathbf{X}^{b^T}_{c_1}$   $\mathbf{X}^{b^T}_{c_2}$   $\mathbf{X}^{b^T}_{N_{cell}}$

# Chebyshev Filtering

# Chebyshev Filtering



$$\mathbf{Y}_{loc}^{b} = \text{ASSEMBLY}\{\mathbf{Y}_{c_i}^{b}\}$$

Atomic operations to avoid race conditions in addition

$$\mathbf{Y}^{b} = \text{ASSEMBLY}\{\mathbf{Y}_{loc}^{b}\}$$

Assembly across processor boundaries: Communication in FP32

$$\mathbf{Y}^{b} = T_m(\mathbf{H})\mathbf{X}^{b} = [2\,\mathbf{H}\,T_{m-1}(\mathbf{H}) - T_{m-2}(\mathbf{H})]\mathbf{X}^{b}$$

Repeat for $b = 1 \cdots \dfrac{N}{B_f}$

SMC 2019

15

# Performance of Chebyshev filtering (Summit)

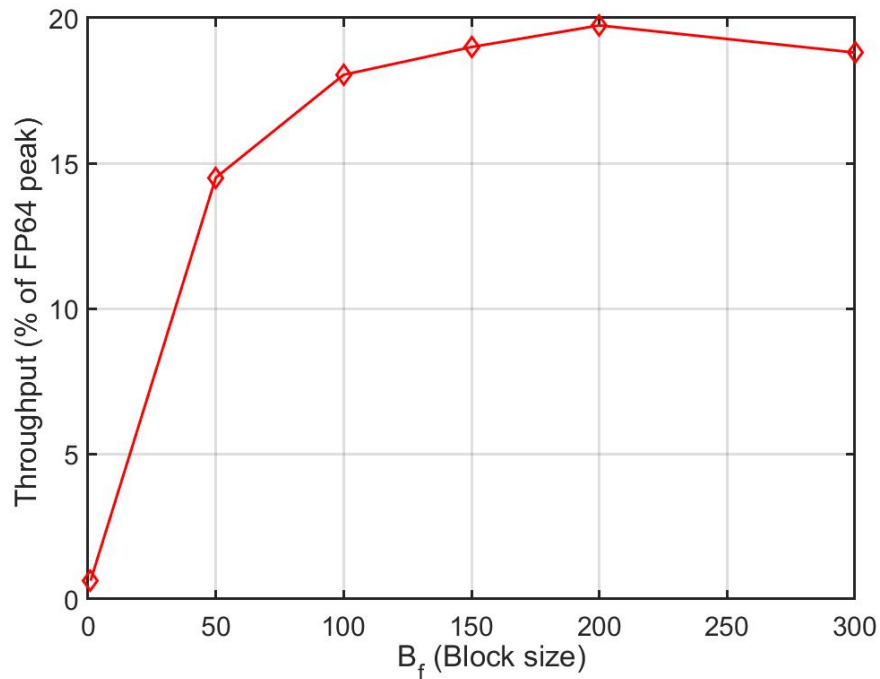**Case study**: Mg 3x3x3 supercell with a vacancy. (1070 electrons)



**Fig**: Chebyshev filtering throughput on 2 Summit nodes using 12 GPUs (3 MPI tasks per GPU) for various block sizes. FP64 peak of 2 Summit nodes is 87.6 TFLOPS
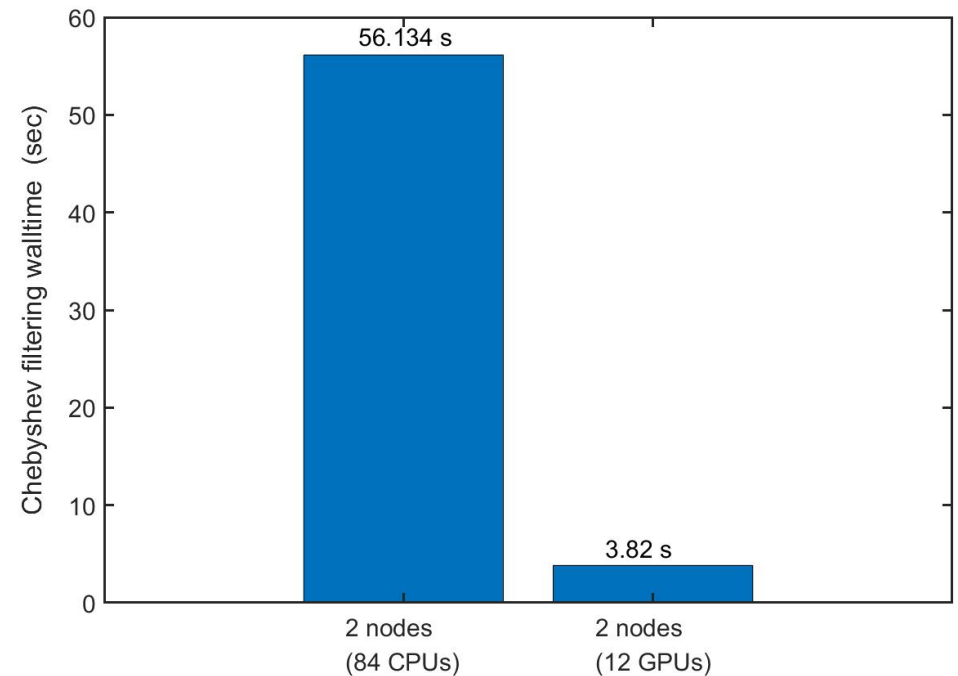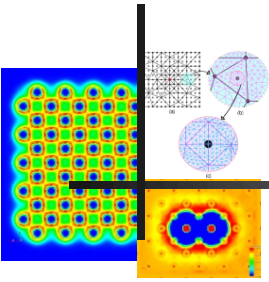


**Fig**: 14.7x GPU speed up for Chebyshev filtering. CPU run used 2 Summit nodes with 42 MPI tasks per node while GPU run used 2 Summit nodes with 12 GPUs (3 MPI tasks per GPU)
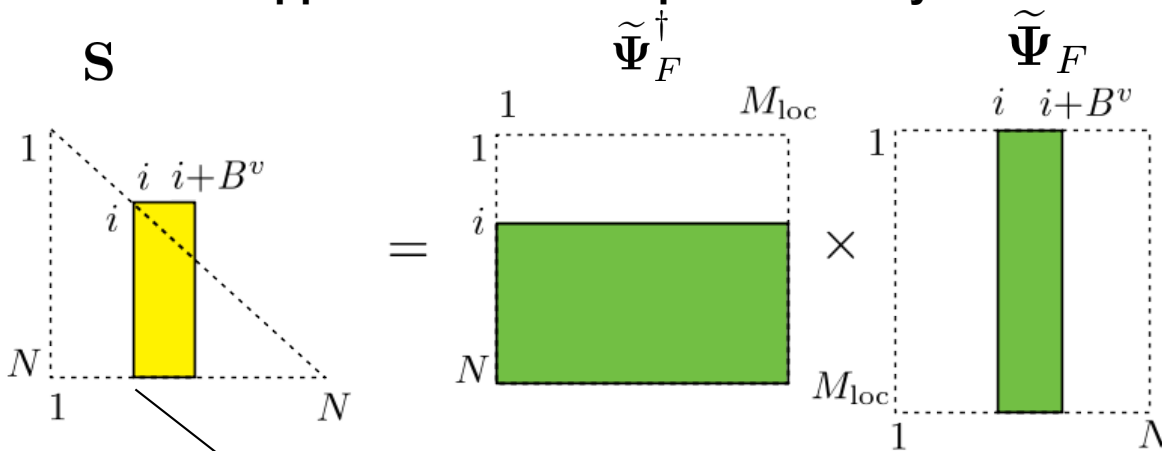
# Orthogonalization: Cholesky Gram-Schmidt

➤ Cholesky factorization of the overlap matrix: $\mathbf{S} = \widetilde{\mathbf{\Psi}}_F^{\dagger} \widetilde{\mathbf{\Psi}}_F = \mathbf{LL}^{\dagger}$. $\mathcal{O}(MN^2)$

➤ Orthonormal basis construction: $\widetilde{\mathbf{\Psi}}_F^{\mathrm{o}} = \widetilde{\mathbf{\Psi}}_F \mathbf{L}^{-1^{\dagger}}$. $\mathcal{O}(MN^2)$

**Blocked approach to reduce peak memory**



**Mixed precision computation for Chol-GS**

1. $\mathbf{S} = \mathrm{DP}\{\mathbf{S_d}\} + \mathrm{SP}\{\mathbf{S_{od}}\}$

2. $\mathbf{S} = \mathbf{LL}^{\dagger}$ in double precision.

3. Orthonormal basis construction:

$$\widetilde{\mathbf{\Psi}}_F^{\mathrm{o}} = \mathrm{DP}\left\{\widetilde{\mathbf{\Psi}}_F \mathbf{L_d}^{-1^{\dagger}}\right\} + \mathrm{SP}\left\{\widetilde{\mathbf{\Psi}}_F \mathbf{L}^{-1^{\dagger}}{}_{\mathbf{od}}\right\}$$
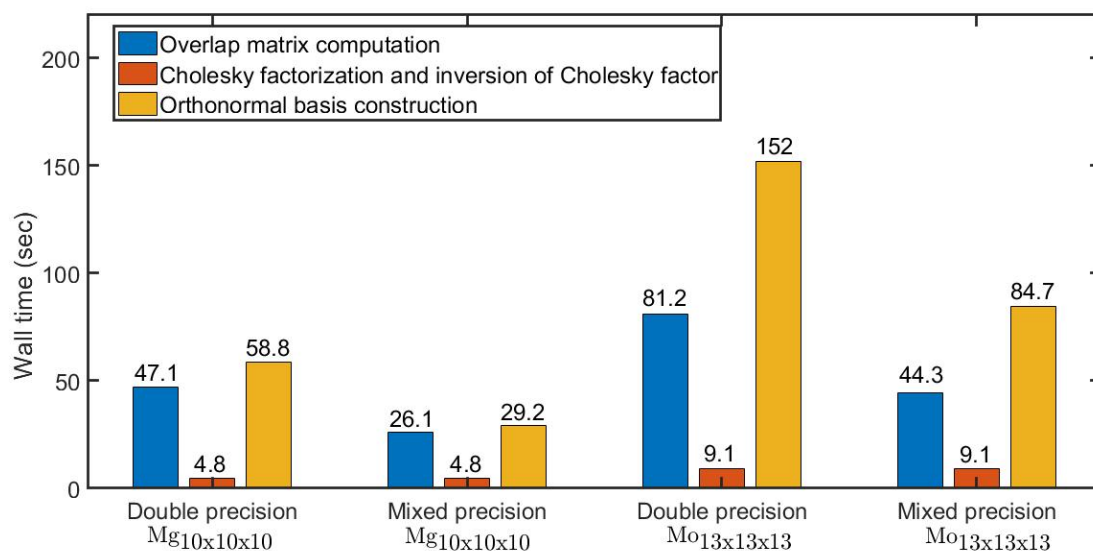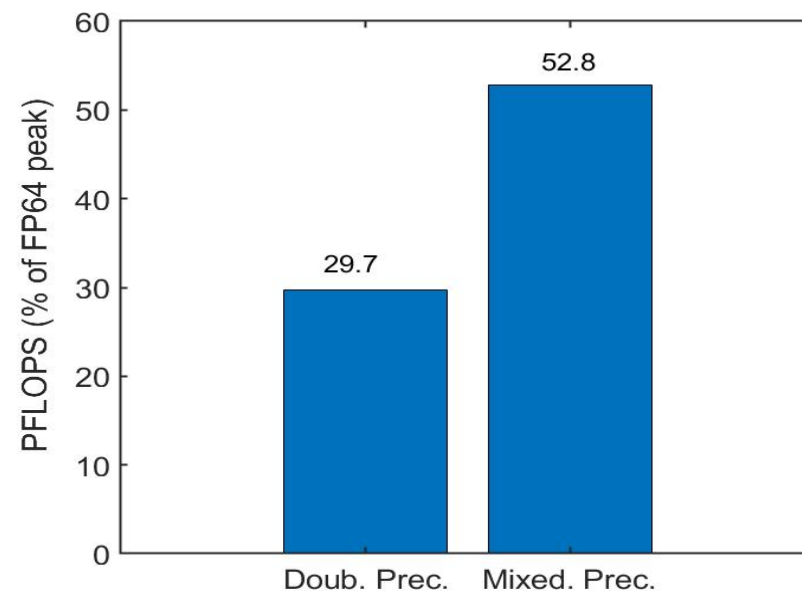
# Orthogonalization: Cholesky Gram-Schmidt

## NERSC Cori CPU cluster benchmark

Performance improvement in CholGS due to mixed precision algorithm. Case study: Mg10x10x10 (39,990 electrons) and Mo13x13x13 (61,502 electrons)

## Summit GPU cluster benchmark

Performance improvement in computation of **S** due to mixed precision algorithm. Case study: 61,640 electrons system using 1300 Summit nodes
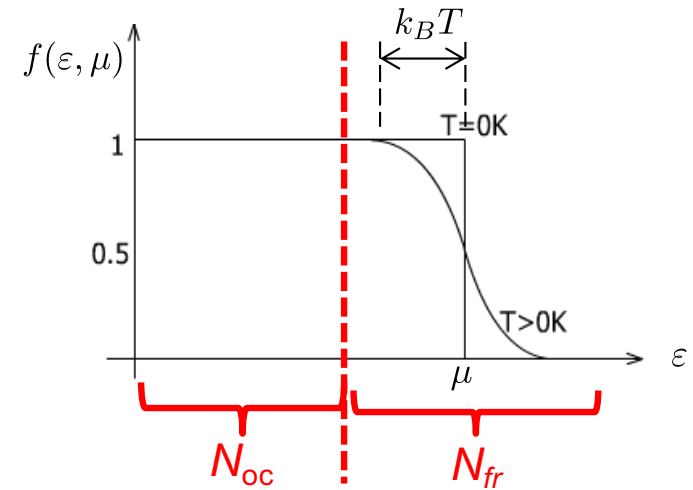
# Rayleigh-Ritz procedure

❖ Compute projected Hamiltonian:   $\hat{\mathbf{H}} = \widetilde{\boldsymbol{\Psi}}_F^{\mathrm{o}\dagger} \widetilde{\mathbf{H}} \widetilde{\boldsymbol{\Psi}}_F^{\mathrm{o}}.\ \mathcal{O}(MN^2)$

❖ Diagonalization of  $\hat{\mathbf{H}}$ :    $\hat{\mathbf{H}}\mathbf{Q} = \mathbf{Q}\mathbf{D}.\ \mathcal{O}(N^3)$

❖ Subspace rotation step:  $\widetilde{\boldsymbol{\Psi}}^{\mathbf{R}} = \widetilde{\boldsymbol{\Psi}}_{\mathbf{F}}^{\mathrm{o}}\mathbf{Q}.\ \mathcal{O}(MN^2)$

**Mixed precision computation for RR**

**1.**   Compute projected Hamiltonian:

$$\rho_{\mathrm{out}}^{h}(\mathbf{x}) = 2\sum_{i=1}^{N} f(\epsilon_i^h, \mu)|\psi_i^h(\mathbf{x})|^2$$

$$\widetilde{\boldsymbol{\Psi}}_F^{\mathrm{o}} = \left[ \widetilde{\boldsymbol{\Psi}}_{\mathrm{oc}}^{\mathrm{o}}\ \widetilde{\boldsymbol{\Psi}}_{\mathrm{fr}}^{\mathrm{o}} \right]$$



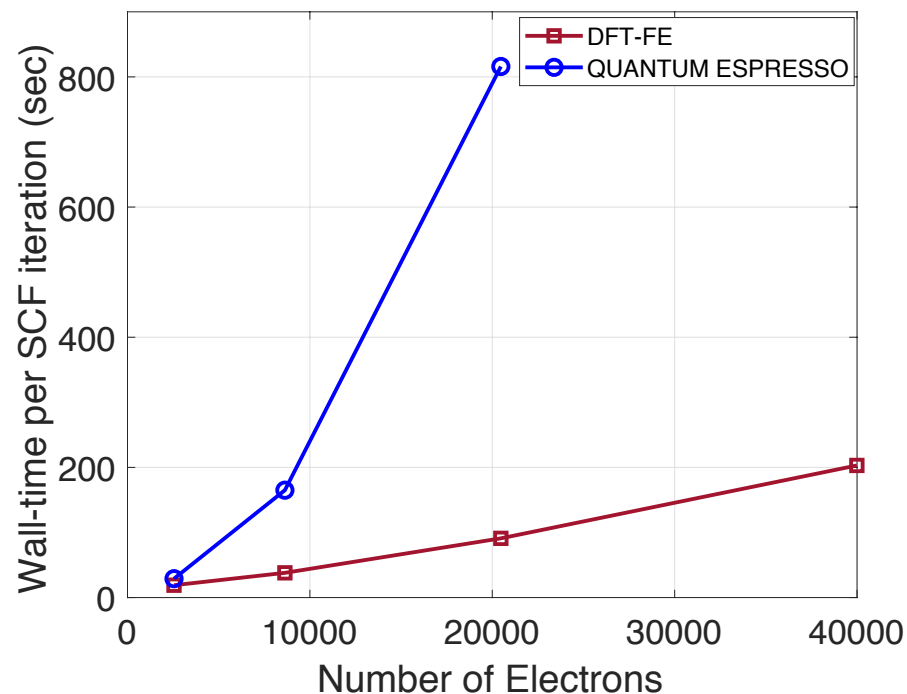$$\left[ \begin{array}{c|c} \hat{\mathbf{H}}_{\mathbf{oc-oc}} & \hat{\mathbf{H}}_{\mathbf{oc-fr}} \\ \hline \hat{\mathbf{H}}_{\mathbf{fr-oc}} & \hat{\mathbf{H}}_{\mathbf{fr-fr}} \end{array} \right] = \left[ \begin{array}{c|c} \mathrm{SP}\left\{ \widetilde{\boldsymbol{\Psi}}_{\mathrm{oc}}^{\mathrm{o}\dagger} \widetilde{\mathbf{H}} \widetilde{\boldsymbol{\Psi}}_{\mathrm{oc}}^{\mathrm{o}} \right\} & \mathrm{SP}\left\{ \widetilde{\boldsymbol{\Psi}}_{\mathrm{oc}}^{\mathrm{o}\dagger} \widetilde{\mathbf{H}} \widetilde{\boldsymbol{\Psi}}_{\mathrm{fr}}^{\mathrm{o}} \right\} \\ \hline \mathrm{SP}\left\{ \widetilde{\boldsymbol{\Psi}}_{\mathrm{fr}}^{\mathrm{o}\dagger} \widetilde{\mathbf{H}} \widetilde{\boldsymbol{\Psi}}_{\mathrm{oc}}^{\mathrm{o}} \right\} & \mathrm{DP}\left\{ \widetilde{\boldsymbol{\Psi}}_{\mathrm{fr}}^{\mathrm{o}\dagger} \widetilde{\mathbf{H}} \widetilde{\boldsymbol{\Psi}}_{\mathrm{fr}}^{\mathrm{o}} \right\} \end{array} \right]$$

# Rayleigh-Ritz procedure

2.  Diagonalization of $\hat{\mathbf{H}} : \hat{\mathbf{H}}\mathbf{Q} = \mathbf{Q}\mathbf{D}$ in double precision.

3.  Subspace rotation step: $\widetilde{\boldsymbol{\Psi}}^{\mathrm{R}} = \mathrm{DP}\left[\widetilde{\boldsymbol{\Psi}}_F^{\mathrm{o}}\mathbf{Q_d}\right] + \mathrm{SP}\left[\widetilde{\boldsymbol{\Psi}}_F^{\mathrm{o}}\mathbf{Q_{od}}\right]$

## Summit GPU cluster benchmark

Performance improvement in computation of $\hat{\mathbf{H}}$ due to
mixed precision algorithm. Case study: 61,640
electrons system using 1300 Summit nodes



SMC 2019

(Motamarri et al. Comput. Phys. Commun. (2019))

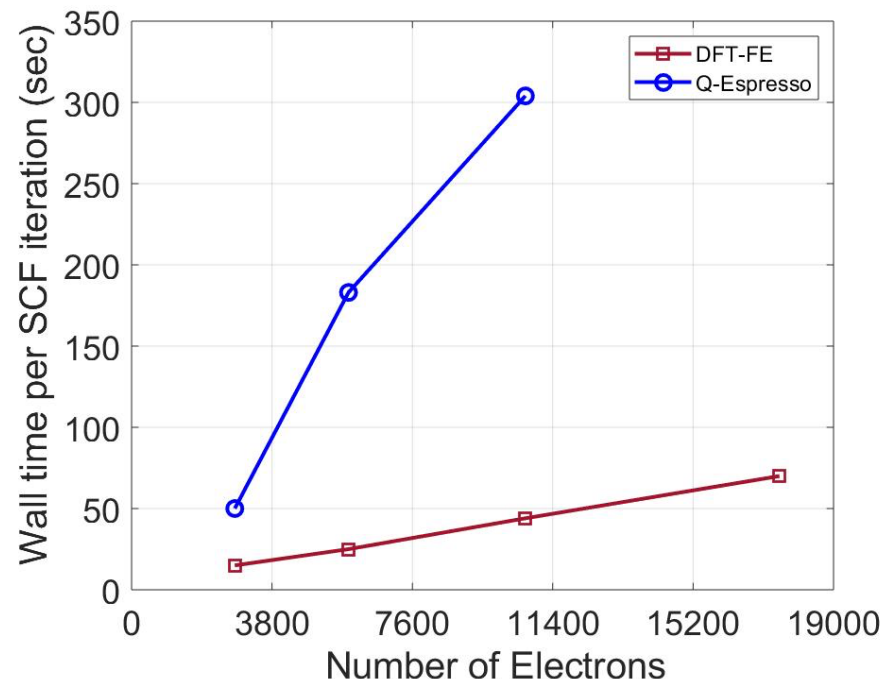Monovacancy in HCP Mg – periodic calculation ; ONCV pseudopotential
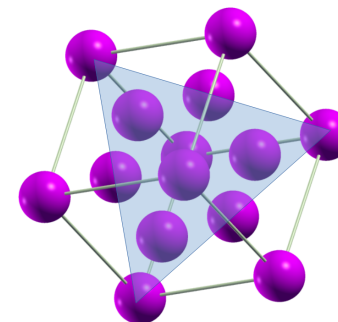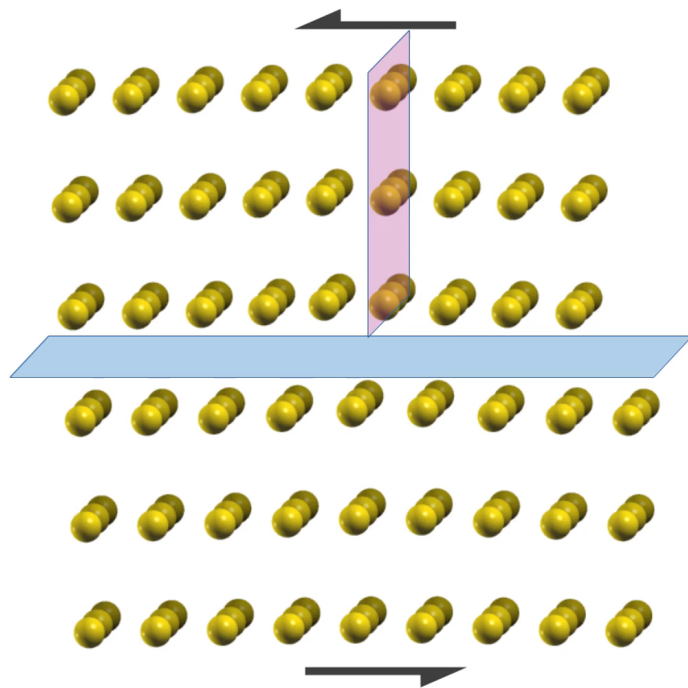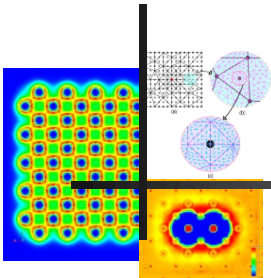Accuracy for all calculations <0.1mHa/atom (~2meV/atom)

Time per SCF in Node-Hrs for various system sizes
(NERSC Cori KNL)

| System size | Q-Espresso (Ecut: 45 Ha) | DFT-FE (h_min: 0.46, p=4) |
|---|---|---|
| 255 atoms ($N_e$ =2550) | 0.1 | 0.3 |
| 863 atoms ($N_e$ =8630) | 4.4 | 3.3 |
| 2047 atoms ($N_e$ =20470) | 123.5 | 21.6 |
| 3999 atoms ($N_e$ =39990) | - | 103.4 |



21

# Comparison with Quantum Espresso (Cori KNL)

Cu nanoparticles – non periodic calculation; ONCV pseudopotential

Accuracy for all calculations <0.1mHa/atom (~2meV/atom)

Time per SCF in Node-Hrs for various system sizes
(NERSC Cori KNL)

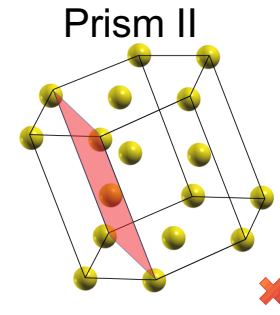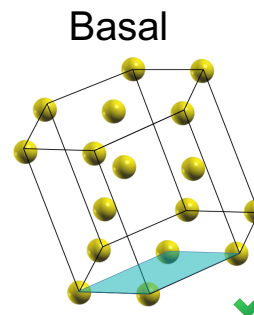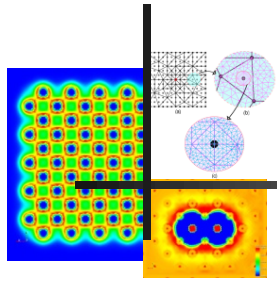| System size | Q-Espresso (Ecut: 50 Ha) | DFT-FE (h_min: 0.4; p=4) |
|---|---|---|
| 147 atoms ($N_e$ =2793) | 0.2 | 0.3 |
| 309 atoms ($N_e$ =5871) | 5.5 | 1.7 |
| 561 atoms ($N_e$ =10569) | 63.4 | 5.3 |
| 923 atoms ($N_e$ =17537) | - | 12.7 |

# Technological challenge of low ductility in Mg



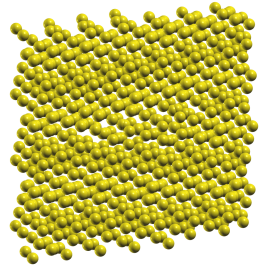12 slip systems in Face Centered Cubic Crystals→ higher ductility

- ❖ Dislocations are energetically more favorable to reside on certain slip systems. (**Energetics**)

- ❖ Dislocation glide occurs after the applied shear stress is greater than the Perils barrier.
    (**Activation barrier**)

- ❖ More the number of slip systems where dislocations can glide easily higher is the ductility.



Basal ✔    Prism II ✖

Prism I ✖    Pyramidal II ?    Pyramidal I ?
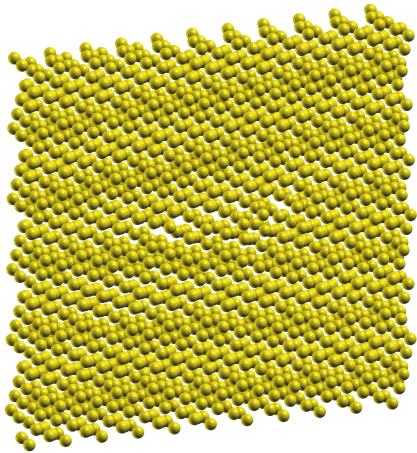
23

SMC 2019

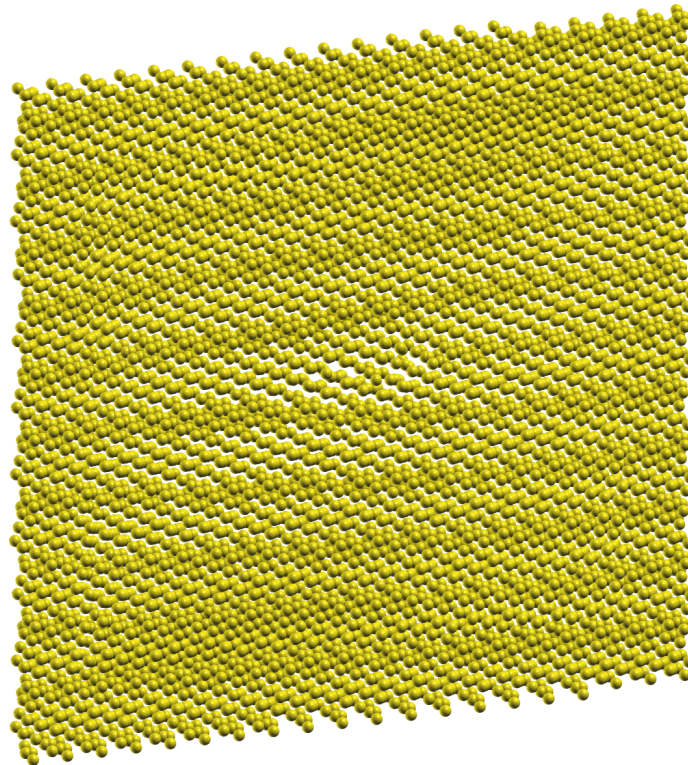# Mg Pyramidal dislocation systems

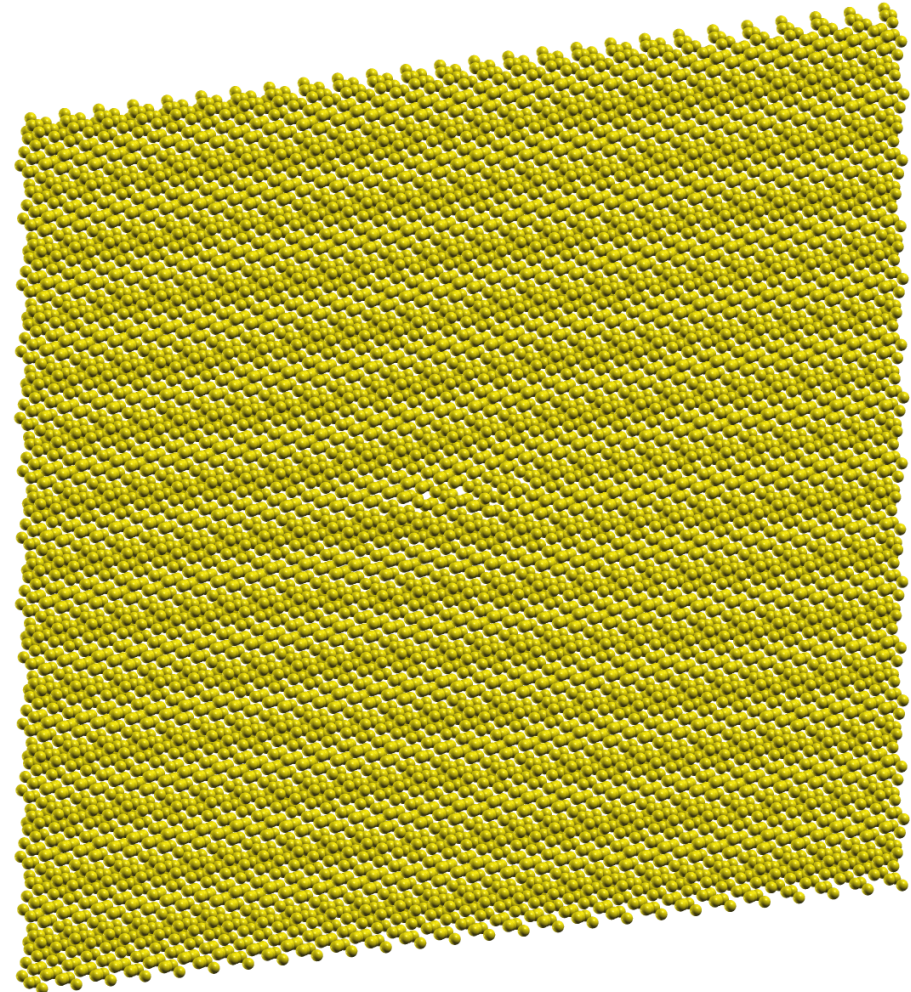Pyramidal I and II dislocation systems of various sizes
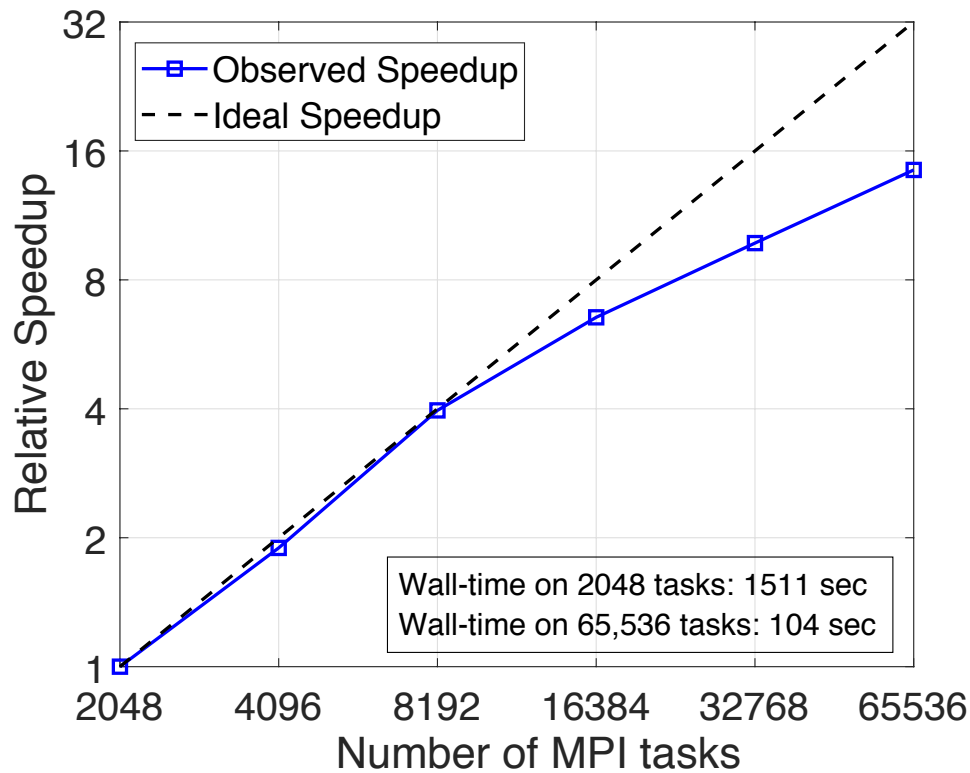
728 Mg atoms
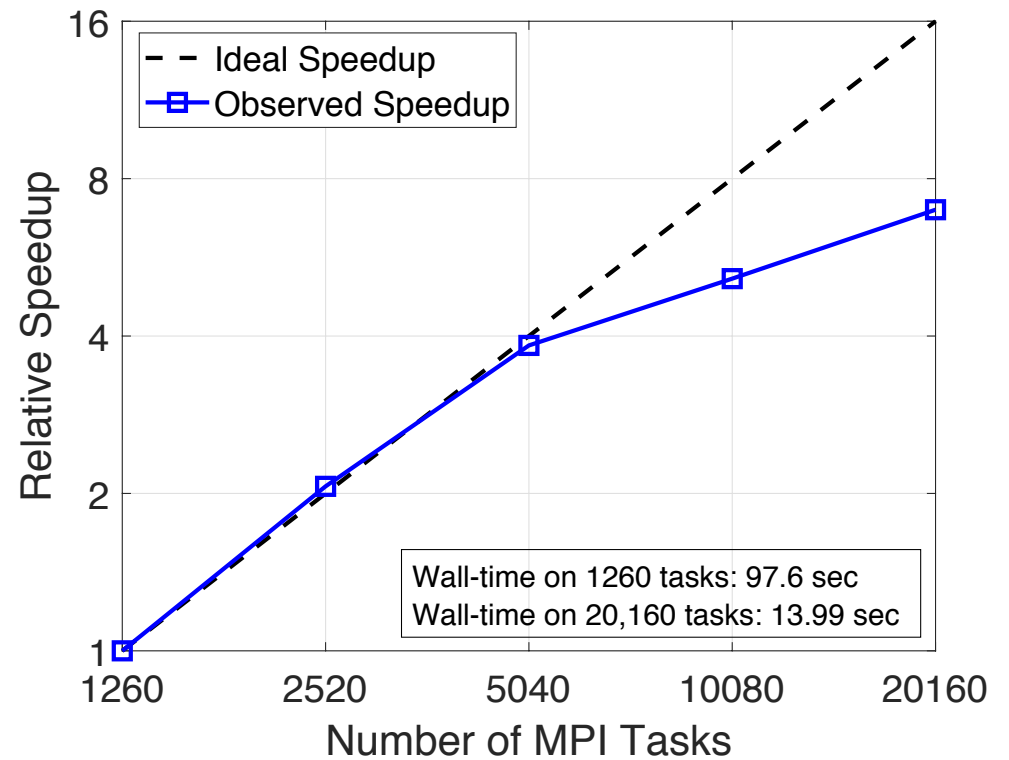
1848 Mg atoms

6164 Mg atoms

10,508 Mg atoms

# Performance Benchmarks – Strong Scaling/time to solution

## Mg pyr II screw dislocation – 1,848 atoms (18,480 e⁻); 55.11 million FE DoFs
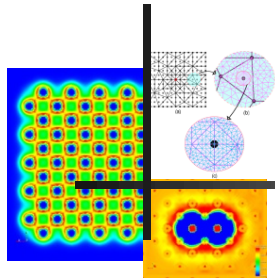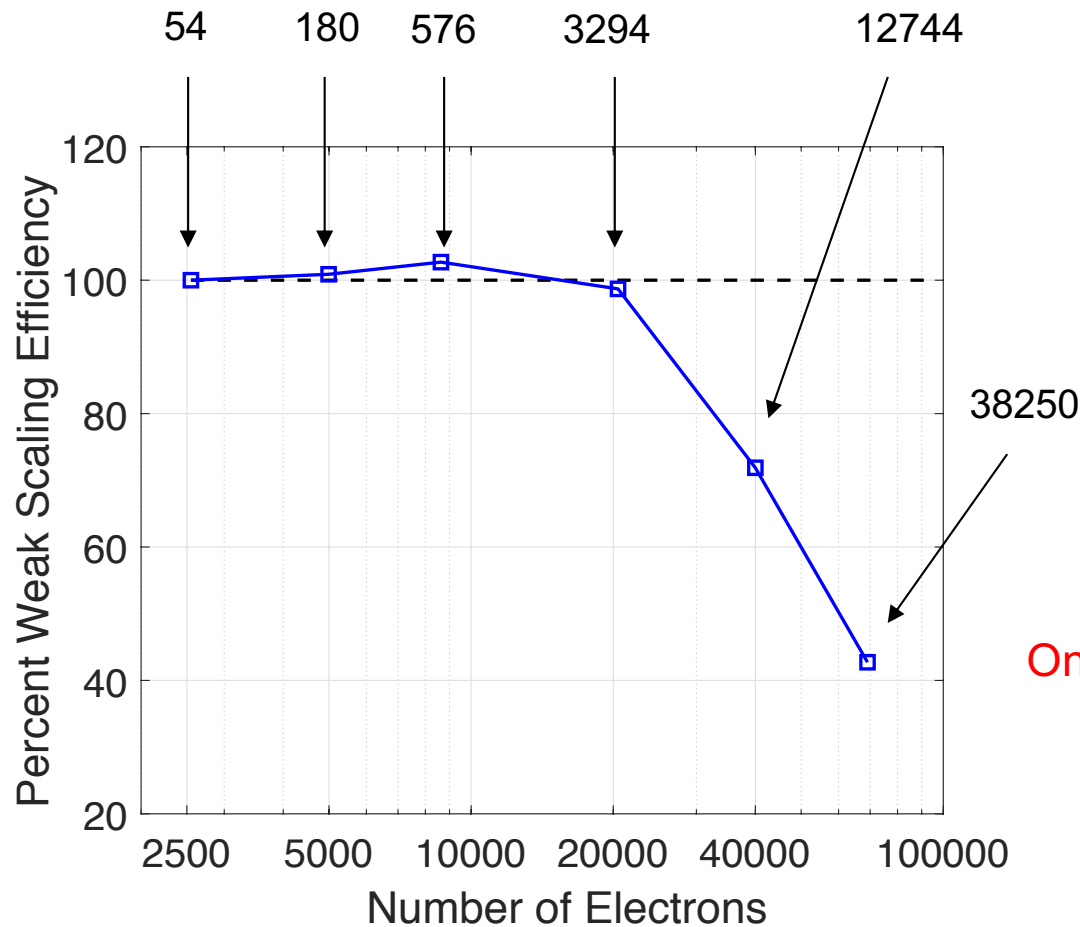
### Theta



Wall-time on 2048 tasks: 1511 sec
Wall-time on 65,536 tasks: 104 sec

### Summit GPUs



Wall-time on 1260 tasks: 97.6 sec
Wall-time on 20,160 tasks: 13.99 sec

3 MPI tasks per GPU via MPS

SMC 2019

# Performance Benchmarks – Weak Scaling (Summit)

Total MPI tasks (3 MPI tasks per GPU; via MPS )



**Computational Complexity**

Chebyshev filtering: $O(MN)$
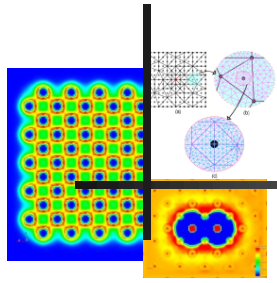
Orthonormalization: $O(MN^2)$

Rayleigh Ritz procedure: $O(MN^2)$

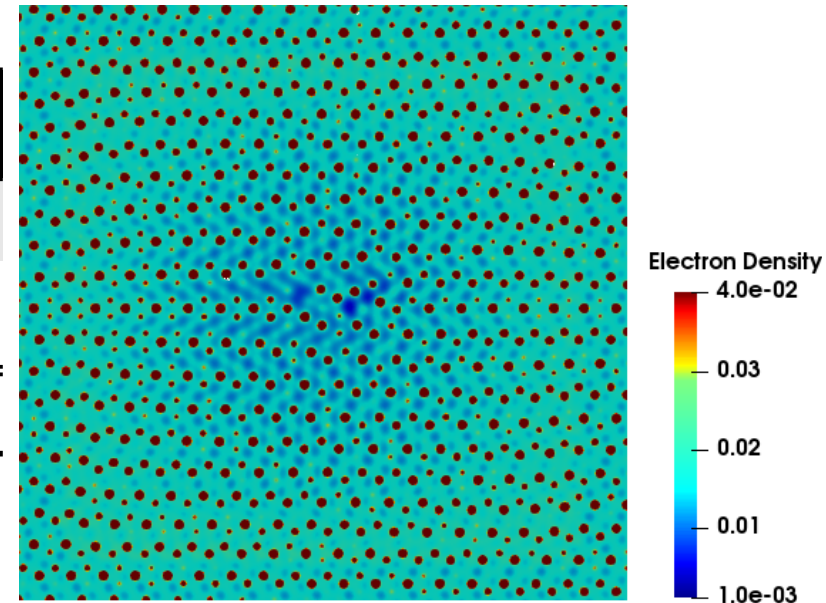<span style="color:red">Onset of cubic scaling significantly delayed !</span>

# Large-scale dislocation systems performance:
## Time-to-solution & Sustained Performance (Summit)

Mg Pyr II dislocation – 6,1640 atoms (61,640 e⁻); 1300 Summit nodes (FP64 peak: 56.65 PFLOPS)

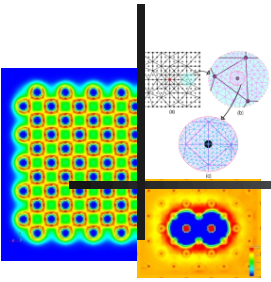| Procedure | Wall-time (sec) | FLOP count (PFOLPS) | PFLOPS (% of FP64 peak) |
|---|---|---|---|
| Initialization | 981 | - | - |
| Ground-state | 7377 | 123174 | 16.7 (29.5%) |
| **Total** | **8358** | **123174** | **14.7 (26.0%)** |



Electron Density

Mg Pyr II dislocation – 10,508 atoms (105,080 e⁻) ; 3800 Summit nodes (FP64 peak: 165.58 PFLOPS)

| Step | Wall-time (sec) | FLOP count (PFLOP) | PFLOPS (% of FP64 peak) |
|---|---|---|---|
| Single SCF | 142.7 | 6563.7 | 46.0 (27.8%) |

# Concluding remarks
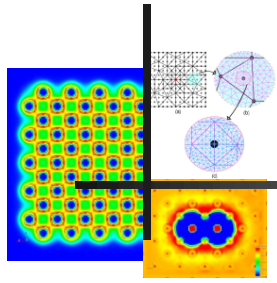
➢ Computational framework
   ❖ Higher-order FE basis
   ❖ Spatial adaptivity
   ❖ Spectral finite-elements w/ GLL quadratures

➢ Algorithms
   ❖ Chebyshev filtering
   ❖ Mixed precision ideas in Orthogonalization and Rayleigh Ritz

➢ Parallel implementation
   ❖ Cell level matrix-matrix operations in Chebyshev filtering with single precision communication
   ❖ Optimizations to reduce peak memory foot print in Orthogonalization and Rayleigh Ritz steps

➢ Fast and accurate large-scale DFT calculations
   ❖ Significant outperformance of some widely used plane-wave codes in both computational efficiency and minimum time-to-solution
   ❖ ~20x speedup using GPUs on a node-to-node comparison
   ❖ Sustained performance of 46 PFOLPS in DFT

28

# THANK YOU!