## **% TRIUMF**

## GPU accelerated *ab initio* nuclear structure calculations on ORNL Summit

TRIUMF-Helmholtz Workshop on Scientific Computing DESY, Hamburg, Germany, 16-17 September, 2019

Petr Navratil TRIUMF

Collaborator: Kostas Kravvaris (LLNL)



Discovery, accelerated

1

2019-09-16

### Outline

- Ab initio calculations in nuclear physics
- Implementation of large-basis sparse matrix diagonalization algorithm
- GPU acceleration of matrix vector multiplication on ORNL Summit and LLNL Lassen
- Example of an application to a nuclear reaction relevant for astrophysics

### First principles or ab initio nuclear theory



First principles or *ab initio* nuclear theory – what we do at present



Ab initio

•

- ♦ Degrees of freedom: Nucleons
- ♦ All nucleons are active
- ♦ Exact Pauli principle
- Realistic inter-nucleon interactions
  Accurate description of NN (and 3N) data

4

♦ Controllable approximations

### Conceptually simplest *ab initio* method: No-Core Shell Model (NCSM)

- Basis expansion method
  - Harmonic oscillator (HO) basis truncated in a particular way (N<sub>max</sub>)
  - Why HO basis?
    - Lowest filled HO shells match magic numbers of light nuclei (2, 8, 20 – <sup>4</sup>He, <sup>16</sup>O, <sup>40</sup>Ca)
    - Equivalent description in relative-coordinate and Slater determinant basis
- Short- and medium range correlations
- Bound-states, narrow resonances

(A) 
$$\Psi^{A} = \sum_{N=0}^{N_{\text{max}}} \sum_{i} c_{Ni} \Phi_{Ni}^{HO}(\vec{\eta}_{1}, \vec{\eta}_{2}, ..., \vec{\eta}_{A-1})$$

(A) 
$$\Psi_{SD}^{A} = \sum_{N=0}^{N_{max}} \sum_{j} c_{Nj}^{SD} \Phi_{SDNj}^{HO}(\vec{r}_{1}, \vec{r}_{2}, ..., \vec{r}_{A}) = \Psi^{A} \varphi_{000}(\vec{R}_{CM})$$





	Progress in Particle and Nuclear Physics 69 (2013) 131-181	
	Contents lists available at SciVerse ScienceDirect	The second secon
	Progress in Particle and Nuclear Physics	Parameter -
ELSEVIER	journal homepage: www.elsevier.com/locate/ppnp	And the second s

Review *Ab initio* no core shell model Bruce R. Barrett<sup>a</sup>, Petr Navrátil<sup>b</sup>, James P. Vary<sup>c,\*</sup>

#### Multi-nucleon states in the Slater Determinant basis

Many-body HO Slater determinants

 $\left\langle \vec{r}_1 \vec{\sigma}_1 \vec{\tau}_1, \vec{r}_2 \vec{\sigma}_2 \vec{\tau}_2, \cdots, \vec{r}_A \vec{\sigma}_A \vec{\tau}_A \middle| a_l^{\dagger} \cdots a_j^{\dagger} a_i^{\dagger} \middle| 0 \right\rangle$  $= \frac{1}{\sqrt{A!}} \begin{vmatrix} \varphi_i(\vec{r}_1) & \varphi_i(\vec{r}_2) & \dots & \varphi_i(\vec{r}_A) \\ \varphi_j(\vec{r}_1) & \varphi_j(\vec{r}_2) & \varphi_j(\vec{r}_A) \\ \vdots & \ddots & \vdots \\ \varphi_l(\vec{r}_1) & \varphi_l(\vec{r}_2) & \dots & \varphi_l(\vec{r}_A) \end{vmatrix}$  $\varphi_{nljm_j\frac{1}{2}m_t}(\vec{r},\vec{\sigma},\vec{\tau})$  $= R_{nl}(r) \Big[ Y_l(\hat{r}) \otimes \chi^{S}_{\frac{1}{2}}(\vec{\sigma}) \Big]_{m_i}^j \chi^{T}_{\frac{1}{2}m_l}(\vec{\tau})$ 

M-scheme basis space dimension 10 10 4He 10 6Li 8Be 10 10B 104 12C 16Ō 10 19F 23Na 27Al  $10^{\circ}$ 10 100 2 8 10 12 14 6 Δ Nmax

10<sup>10</sup>

10

- Antisymmetrization is trivial
- Good M, M<sub>T</sub> and parity quantum numbers, but not J and T
  - Huge number of basis states

#### **Basis states: occupation representation**

- How are Slater determinants (SD) actually represented in a computer program?
  - We are dealing with fermions, so a single-particle state is either occupied or empty, which in computer language translates to either 1's or 0's

7

- A very useful approach is a bit representation known as M-scheme
  - HO single-particle states will have good *j*, *m<sub>j</sub>*

- A single integer represents a complicated Slater determinant
- While the many-body SD states will have good M, they do not have good J. States of good J must be projected and will be a combination of Slater determinants. Same for T and  $M_T$
- Eigenstates obtained by the diagonalization of the Hamiltonian in SD basis will have a good J and (approximately) good T

#### Getting the eigenvalues and eigenfunctions: Code NCSD (No-Core shell model in Slater Determinants)

- M-scheme basis:
  - Total M,  $M_T = (Z-N)/2$  and parity conserved
- Setup Hamiltonian matrix  $\langle \Phi_i | H | \Phi_j \rangle$
- Diagonalize using Lanczos algorithm
  - Bring matrix to tri-diagonal form ( $v_1$ ,  $v_2$  ... orthonormal, H Hermitian)

 $H\mathbf{v}_{1} = \alpha_{1}\mathbf{v}_{1} + \beta_{1}\mathbf{v}_{2}$   $H\mathbf{v}_{2} = \beta_{1}\mathbf{v}_{1} + \alpha_{2}\mathbf{v}_{2} + \beta_{2}\mathbf{v}_{3}$   $H\mathbf{v}_{3} = \beta_{2}\mathbf{v}_{2} + \alpha_{3}\mathbf{v}_{3} + \beta_{3}\mathbf{v}_{4}$   $H\mathbf{v}_{4} = \beta_{3}\mathbf{v}_{3} + \alpha_{4}\mathbf{v}_{4} + \beta_{4}\mathbf{v}_{5}$ 

- n<sup>th</sup> iteration computes 2n<sup>th</sup> moment
- Eigenvalues converge to extreme (largest in magnitude) values
- $\sim 150-200$  iterations needed for 10 eigenvalues (even for  $10^9$  states)
- Eigenvectors obtained with simplified form of Gaussian elimination

tridiagonal matrix Pick a random vector:  $|v_1\rangle$ Compute a new vector:  $|w'\rangle = H|v_1\rangle$ Store the overlap:  $\alpha_1 = \langle v_1 | w' \rangle$ Orthogonalize the two:  $|w\rangle = |w'\rangle - \langle v_1 |w'\rangle |v_1\rangle$ Loop for as many as we can: For j = 1, ..., mStore the norm of the new vector:  $\beta_j = \sqrt{\langle w | w \rangle}$ Normalize the new vector:  $|v_i\rangle = |w_{i-1}\rangle/\beta_i$ Compute new vector:  $|w'_i\rangle = H|v_j\rangle$ Store Overlap:  $\alpha_j = \langle w'_j | v_j \rangle$ Obtain new vector:  $|w_j\rangle = |w'_j\rangle - \alpha_j |v_j\rangle - \beta_j |v_{j-1}\rangle$ 

**Diagonalize only small** 

#### **NCSD Main Loop Structure**

NCSD – Fortran 90 with MPI, OpenMP and CUDA GPU accelerations, capable of running ~12,000 MPI tasks Lanczosh diagonalization, bit operations, hashing, partial or full storing of non-zero Hamiltonian matrix elements in memory 10 NN+3N interaction input, can calculate  $N_{max}$  basis sequence (0(1),2(3),... $N_{max}$ ) in a single run





- We construct the matrix elements in RAM as it is generally more readily available than GPU memory.
- Matrix elements are transferred to GPU and perform the matrix-vector multiplications on the GPU.
- Due to mixed float/double precision for matrix elements/vector elements, no CUBLAS intrinsic function is used.

- Matrix is currently saved in Compressed Sparse Row format. Will change to Compressed Sparse Block in the future for better GPU efficiency.
- Current code does SpMV multiplication; by using SpMM for a block of vectors, GPU usage will further increase.

#### SUMMIT - Oak Ridge National Laboratory's 200 petaflop supercomputer

- Specifications
  - Processor: IBM POWER9<sup>™</sup> (2/node)
  - GPUs: 27,648 NVIDIA Volta V100s (6/node)
  - Nodes: 4,608
  - Node Performance: 42TF
  - Memory/node: 512GB DDR4 + 96GB HBM2
  - NV Memory/node: 1600GB
  - Total System Memory: >10PB DDR4 + HBM + Non-volatile
  - Interconnect Topology: Mellanox EDR 100G InfiniBand, Non-blocking Fat Tree
  - Peak Power Consumption: 13MW
- INCITE award "Nuclear Structure and Reactions"



### Summit Node (2) IBM Power9 + (6) NVIDIA Volta V100

Anatomy of a Summit Node:

- Fat Node allows for more freedom in resource management.
- 16GB of memory/GPU ٠ means they can be used for storage of not only the vectors, but also the matrix elements.
- 512GB of RAM/node ٠ allow for the entire Hamiltonian Matrix to be stored.
- The 42 hyper-threaded • cores accelerate the construction of the matrix.



NVLink2

1 (900 GB/s)



- Speedup decreases at high node count as the non-GPU accelerated part of the code dominates.
- Memory footprint also decreased by a factor of 2 from more efficient storage of matrix elements.
- SUMMIT has enough nodes so that the entire matrix fits on GPU memory.



### LASSEN machine at LLNL

- 4 GPUs/node
- 256GB RAM
- 756 Nodes total

Need different strategy for large calculations!

- Use multiple MPI tasks/GPU.
- One task fully utilizes the GPU and the rest work in pure CPU mode.
- Allows for better distribution of the matrix between RAM/GPU memory.
- Implement MPI window object for sharing memory within a single node, further reducing memory footprint.



5

Example of an application of the GPU accelerated code: p+<sup>11</sup>C scattering and <sup>11</sup>C(p,γ)<sup>12</sup>N capture

<sup>11</sup>C(p,γ)<sup>12</sup>N capture relevant in hot *p*-*p* chain: Link between pp chain and the CNO cycle - bypass of slow triple alpha capture <sup>4</sup>He(αα,γ)<sup>12</sup>C



 ${}^{3}He(\alpha,\gamma)^{7}Be(\alpha,\gamma)^{11}C(p,\gamma)^{12}N(p,\gamma)^{13}O(\beta^{+},\nu)^{13}N(p,\gamma)^{14}O$  ${}^{3}He(\alpha,\gamma)^{7}Be(\alpha,\gamma)^{11}C(p,\gamma)^{12}N(\beta^{+},\nu)^{12}C(p,\gamma)^{13}N(p,\gamma)^{14}O$  ${}^{11}C(\beta^{+},\nu)^{11}B(p,\alpha)^{8}Be({}^{4}\text{He},{}^{4}\text{He})$ 

# Example of an application of the GPU accelerated code: p+<sup>11</sup>C scattering and <sup>11</sup>C(p,γ)<sup>12</sup>N capture

- NCSM with continuum calculations of <sup>11</sup>C(p,p) with chiral NN+3N under way
  - <sup>11</sup>C and <sup>12</sup>N NCSM eigenstates calculated with NCSD on Summit using GPU acceleration
  - 1024 nodes, 6144 MPI tasks with 1 GPU/task and 7 OpenMP threads/task
  - Largest matrix dimensions: 131 million for <sup>11</sup>C and 167 million for <sup>12</sup>N (~7 hours to get 9 eigenstates)



- NCSM with continuum calculations of <sup>11</sup>C(p,p) with chiral NN+3N under way
  - <sup>11</sup>C and <sup>12</sup>N NCSM eigenstates calculated with NCSD on Summit using GPU acceleration
  - 1024 nodes, 6144 MPI tasks with 1 GPU/task and 7 OpenMP threads/task
  - Largest matrix dimensions: 131 million for <sup>11</sup>C and 167 million for <sup>12</sup>N (~7 hours to get 9 eigenstates)



NCSMC calculations to be validated by measured cross sections and applied to calculate the  ${}^{11}C(p,\gamma){}^{12}N$  capture

#### Conclusions

• *Ab initio* calculations of nuclear structure and reactions becoming feasible beyond the lightest nuclei

• Make connections between the low-energy QCD, many-body systems, and nuclear astrophysics

#### Large-scale computational problems

- Codes running on massively parallel supercomputers
- Taking advantage of GPU acceleration with a significant speedup demonstrated on ORNL Summit and LLNL Lassen
- Enables nuclear physics applications such as calculations of nuclear reactions relevant for astrophysics

# **∂**TRIUMF

## Thank you! Merci! Danke!

