# Coulomb and Landau Gauge Fixing in GPUs using CUDA and MILC

## Nuno Cardoso

NCSA, University of Illinois



Lattice 2014

- MILC and QUDA library does not have gauge fixing on GPUs
- When using MILC + QUDA: lot of computational time can be waisted on CPU gauge fixing
- Main goal was to create a library to be called from MILC code with support for single and multi-GPUs to perform gauge fixing

• On the lattice, the Coulomb/Landau gauge is defined by maximising the functional

$$F_{U}[g] = \frac{1}{4N_{c}V} \sum_{x} \sum_{\mu} \operatorname{Re} \left[ \operatorname{Tr} \left( g(x)U_{\mu}(x)g^{\dagger}(x+\hat{\mu}) \right) \right]$$

with  $N_c$  the dimension of the gauge group and V the lattice volume.

• On the gauge fixing process, the quality of the gauge fixing is measured by

$$\theta = \frac{1}{N_c V} \sum_{x} \operatorname{Tr} \left[ \Delta(x) \Delta^{\dagger}(x) \right]$$

where

$$\Delta(x) = \sum_{
u} \left[ U_{
u}(x - a \hat{
u}) - U_{
u}(x) - ext{h.c.} - ext{trace} 
ight]$$

is the lattice version of  $\partial_{\mu}A_{\mu} = 0$ .

- Two well known methods to fix the gauge:
  - Relaxation algorithm: overrelaxation
  - Steepest descent method with FFTs

2/15

### Overrelaxation

 $\bullet$  the relaxation algorithm aims to optimize the value of  $F_U[g]$  locally, i.e., searching the maximum of

$$f^{g}(x) = \operatorname{Re}\operatorname{Tr}\left[g(x)K(x)\right]$$

for all x, where

$$\mathcal{K}(x) = \sum_{\mu} \left( U_{\mu}(x)g^{\dagger}(x+\hat{\mu}) + U_{\mu}^{\dagger}(x-\hat{\mu})g^{\dagger}(x-\hat{\mu}) \right)$$

the local solution is then given by

$$g(x) = rac{K^{\dagger}(x)}{\sqrt{\det K^{\dagger}(x)}}$$

in the case of the gauge group SU(2).

- for N > 2 one iteratively operates in the (N(N-1)/2) SU(2) subgroups.
- the overrelaxation algorithm replaces the gauge transformation g(x) by  $g^{\omega}(x)$  with  $\omega \in [1, 2[$  in each step of the iteration.

< 🗆 🕨

UV6	Overrelaxation		
	calculate $F_{g}[U]$ and $\theta$ (optional)		
2:	while $\theta \ge \epsilon$ do		
	<b>for</b> site parity = even, odd <b>do</b>		
4:	for all x with same parity do		
	for all SU(2) subgroups do		
6:	local optimization, find $g(x) \in SU(2)$		
	which is function of $U_\mu(x)$ and $U_\mu(x-\hat\mu)$		
8:	for all $\mu$ do		
	apply $g(x)$ to $U_\mu(x)$ and $U_\mu(x-\hat{\mu})$		
10:	end for		
	end for		
12:	end for		
	end for		
14:	calculate $F_{g}[U]$ and $ heta$		
	end while		

 ${ \blacksquare } { \blacksquare } { > }$ 



## Coulomb and Landau Gauge Fixing in GPUs using CUDA and MILC $_{\mbox{Overrelaxation}}$

## Overrelaxation implementation using multi-GPUs

CUDA + MILC + MPI

## 2D Lattice Example:

- Gauge array split by even and odd lattice sites
- Pre calculate interior lattice array index sites and lattice array index at the boundaries for each dimension partitioned
- Exchange face links
- Measure gauge fixing quality
- Main loop to perform gauge fixing,
  - Loop over all lattice sites by even and odd links separately
    - update top/bottom and right/left face links
    - pack top (same parity) and ghost (opposite parity) face links
    - exchange top face (same parity) and ghost (opposite parity) links
    - exchange ghost links
    - update interior links
    - unpack received links
  - reunitarize gauge at X steps
  - measure gauge fixing quality

< 🗆 🕨

## Coulomb and Landau Gauge Fixing in GPUs using CUDA and MILC Overrelaxation

## Overrelaxation implementation using multi-GPUs

CUDA + MILC + MPI

2D Lattice Example:

- Gauge array split by even and odd lattice sites
- Pre calculate interior lattice array index sites and lattice array index at the boundaries for each dimension partitioned
- Exchange face links
- Measure gauge fixing quality
- Main loop to perform gauge fixing,
  - Loop over all lattice sites by even and odd links separately
    - update top/bottom and right/left face links
    - pack top (same parity) and ghost (opposite parity) face links
    - exchange top face (same parity) and ghost (opposite parity) links
    - exchange ghost links
    - update interior links
    - unpack received links
  - reunitarize gauge at X steps
  - measure gauge fixing quality



< 🗆 🕨

## Coulomb and Landau Gauge Fixing in GPUs using CUDA and MILC Overrelaxation

## Overrelaxation implementation using multi-GPUs

CUDA + MILC + MPI

2D Lattice Example:

- Gauge array split by even and odd lattice sites
- Pre calculate interior lattice array index sites and lattice array index at the boundaries for each dimension partitioned
- Exchange face links
- Measure gauge fixing quality
- Main loop to perform gauge fixing,
  - Loop over all lattice sites by even and odd links separately
    - update top/bottom and right/left face links
    - pack top (same parity) and ghost (opposite parity) face links
    - exchange top face (same parity) and ghost (opposite parity) links
    - exchange ghost links
    - update interior links
    - unpack received links
  - reunitarize gauge at X steps
  - measure gauge fixing quality



#### Steepest descent method with FFTs

• The naive steepest descent method chooses at each step of the iterative procedure

$$g(x) = \exp\left[rac{lpha}{2}\left(\sum_{
u} \Delta_{-
u}\left[U_{
u}(x) - U_{
u}^{\dagger}(x)\right] - \text{trace}
ight)
ight]$$

- However, when it is applied to larger lattices, this method faces the problem of critical slowing down.
- This problem can be attenuated by Fourier acceleration, Davies et. al, doi:10.1103/PhysRevD.37.1581, 1987.
- At each iteration one chooses

$$g(x) = \exp\left[\hat{F}^{-1}\frac{\alpha}{2}\frac{p_{\max}^2 a^2}{p^2 a^2}\hat{F}\left(\sum_{\nu} \Delta_{-\nu}\left[U_{\nu}(x) - U_{\nu}^{\dagger}(x)\right] - \text{trace}\right)\right]$$

with

$$\Delta_{-\nu}\left(U_{\mu}(x)\right) = U_{\mu}(x - a\hat{\nu}) - U_{\mu}(x)$$

 $p^2$  are the eigenvalues of  $\left(-\partial^2\right),~a$  is the lattice spacing and  $\hat{F}$  represents a fast Fourier transform (FFT).

- For the parameter  $\alpha$ , the optimal value is 0.08.
- For numerical purposes, it is enough to expand to first order the exponential, followed by a reunitarization.

Steepest descent method with FFTs		
	calculate $\Delta(x)$ , $F_g[U]$ and $ heta$	
2:	while $\theta \ge \epsilon$ do	
	for all elements of $\Delta(x)$ matrix do	
4:	apply FFT forward	
	apply $p_{\max}^2/p^2$	
6:	apply FFT backward	
	normalize	
8:	end for	
	for all $\times$ do	
10:	obtain $g(x)$ from $\Delta(x)$ and reunitarize	
	end for	
12:	for all x do	
	for all $\mu$ do	
14:	$U_\mu(x)  o g(x) U_\mu(x) g^\dagger(x+\hat{\mu})$	
	end for	
16:	end for	
	calculate $\Delta(x)$ , $F_{g}[U]$ and $ heta$	
18:	end while	

< 🗆 🕨



•  $2 \times 2D$  FFTs is 7-8% faster than 3D+1D FFTs.



7/15

- Performance tests done in Blue Waters
  - GPUs Tesla K20X with ECC code enabled by default
  - one GPU per node
- Performance results only for SU(3)
- Overrelaxation:
  - single and multi-GPU
  - Coulomb and Landau
  - with/without using AtomicAdd()
  - with/without using texture memory
  - full SU(3) matrix and 12 real parameters memory storage
- Steepest descent method with FFTs:
  - only single GPU
  - Coulomb and Landau
  - with/without using texture memory
  - full SU(3) matrix and 12 real parameters memory storage

< 🗆 🕨

## Results: Landau gauge fixing with Overrelaxation



#### Single GPU performance

AtomicAdd means using CUDA atomicAdd function

SP/DP means single/double precision

Tex means using Texture memory

12/18 parameters are the number of SU(3) real parameters used to store gauge array in memory

Also supports 8 real parameters, no performance results done vet.

Gauge Fixing in GPUs using CUDA and MILC



GPU Speedup over CPU MILC code in Single node

The performance impact in converting MILC format to/from GPU format as well as the copies depends on the number of the gauge fixing steps. In this example the number of gage fixing iterations was around 170.

Gauge Fixing in GPUs using CUDA and MILC

## Results: Landau gauge fixing with Overrelaxation

### Weak Scaling

- Fixed node local lattice volume to 32<sup>4</sup>
- Lattice volume:  $32^3 \times Nt$  with  $Nt = 32 \times (Number of nodes)$



#### Double precision

11/15

## Results: Landau gauge fixing with Overrelaxation

### Weak Scaling

- Fixed node local lattice volume to 32<sup>4</sup>
- Lattice volume:  $32^3 \times Nt$  with  $Nt = 32 \times (Number of nodes)$



N. Cardoso

Gauge Fixing in GPUs using CUDA and MILC

11/15



#### Single GPU performance

N. Cardoso

Gauge Fixing in GPUs using CUDA and MILC

12/15

< □ →

i.

Multi-GPU performance, 32<sup>4</sup> lattice volume



Adapted previous code to support SU(N) with N>2, added Coulomb gauge fixing and even/odd lattice site order

N. Cardoso, P. J. Silva, P. Bicudo, O. Oliveira, Landau Gauge Fixing on GPUs, Comput.Phys.Commun. 184 (2013) 124-129. arXiv:1206.0675, doi:10.1016/j.cpc.2012.09.007

13/15

## Conclusion

- Coulomb and Landau gauge fixing implemented on GPUs
- Big speedup using GPUs to perform gauge fixing
- Overrelaxation implemented in single and multi-GPUs
- Steepest descent algorithm with Fourier acceleration converges faster than overrelaxation algorithm, in general However, steepest descent algorithm with Fourier acceleration requires  $1.5 \times$  more memory.
- Steepest descent algorithm with Fourier acceleration not implemented in multi-GPU
  - needs multi-GPU 4D FFTs
  - needs much more data exchange between GPUs
  - needs to perform 2k FFTs per step, with k the number of  $\Delta(x)$  elements
  - one solution to avoid using FFTs is to use a multigrid implementation of the Fourier acceleration method Cucchieri et. al, Phys.Rev.D57, arXiv:hep-lat/9711047, 1998. not implemented here!
- this library supports:
  - SU(N) with N>2, maximum N supported is limited by hardware resources
  - SU(3) parameterization with 8 and 12 real numbers

14/15

## Thanks

15/15

 $\leftarrow \Box \rightarrow$ 

## Results

### Landau gauge fixing: GPU performance, single node



Using MILC stop criterium  $10^{-7}$ .

In this example, the steepest descent method with FFTs only needed around 17 iterations to converged and the overrelaxation around 140-180 iterations.

• Speed up of the GPU Landau gauge fixing with the steepest descent method with FFTs over GPU Landau gauge fixing with overrelaxation. Right figure includes the time to copy and reorder from MILC format to GPU format, the left figure does not include this time.

Gauge fixing quality results example MILC input file for su3 ora: prompt 0 nx 32 ny 32 nz 32 nt 32 (...) trajecs 5 traj\_between\_meas 1 beta 6.0 steps\_per\_trajectory 4 qhb\_steps 4 (...)

< □ ►



### Landau gauge fixing quality results

Number of iterations to achieve a gauge fixing quality of  $\theta < 10^{-15}$ . In these results,  $\theta$  parameter was used as stop criterium instead of the default MILC criterium.



### Coulomb gauge fixing quality results

Number of iterations to achieve a gauge fixing quality of  $\theta < 10^{-15}$ . In these results,  $\theta$  parameter was used as stop criterium instead of the default MILC criterium.