

WILSON INVERTER USING CUDA

MICHAEL CLARK

CENTRE FOR COMPUTATIONAL SCIENCE
BOSTON UNIVERSITY

WITH
R. BABICH, K. BARROS,
R. BROWER AND C. REBBI



GPU HARDWARE

GTX 280 / 285

FP32: 0.933/1.06 Tflops, FP64: 78/89 Gflops

Memory 1GB, Bandwidth 141 / 159 GBs⁻¹

230 Watts, \$224 / \$250



Tesla S1070

FP32: 4 Tflops, FP64: 320 Gflops

Memory 16GB, Bandwidth 408 GBs⁻¹

900 Watts, \$8000

Tesla C1060

FP32: 0.933 Tflop, FP64: 78 Gflops

Memory 4GB, Bandwidth 102 GBs⁻¹

230 Watts, \$1000





GPU PROGRAMMING

- Need to completely rethink algorithms/problems accordingly
 - Data parallelism
 - Flops/Bandwidth ratio huge : minimize memory traffic at all costs
 - Reorder data for memory coalescing
 - Limited register size requires each thread must be *light weight*
 - Use the shared memory to share data between threads
 - Fit problems within memory constraints
 - Single precision strong, double precision weak

GPU WILSON-DIRAC MATRIX-VECTOR

$$D = \delta_{x,y} - \kappa \sum_{\mu} \left(P^{-\mu} U_{x,y}^{\mu} \delta_{x+\mu,y} + P^{+\mu} U_{y,x}^{\mu\dagger} \delta_{x-\mu,y} \right)$$

- Assign a single spacetime point to each thread
 - Lightweight threads
- (Spin-projected) dslash operator has
 - 1320 flops per site
 - 1440 bytes per site (32 bit)
- Extremely bandwidth limited
- Need to reduce memory bandwidth at all costs
 - Exploit all symmetries of Wilson-Dirac operator

FIELD ORDERING

- Can order data optimally for any given hardware
 - Reorder field elements for coalescing (internal d.o.f. outside spacetime)
 - Usually store field of 24 vectors
 - Instead 6 fields of float4s
- Boundary conditions destroy coalescing -> use texture cache

γ BASIS

- Usually use chiral basis
- γ_4 projector is $P^{\pm 4} = \begin{pmatrix} 1 & 0 & \pm 1 & 0 \\ 0 & 1 & 0 & \pm 1 \\ \pm 1 & 0 & 1 & 0 \\ 0 & \pm 1 & 0 & 1 \end{pmatrix}$
- Instead use non-relativistic basis

$$P^{+4} = \begin{pmatrix} 2 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \quad P^{-4} = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 2 \end{pmatrix}$$

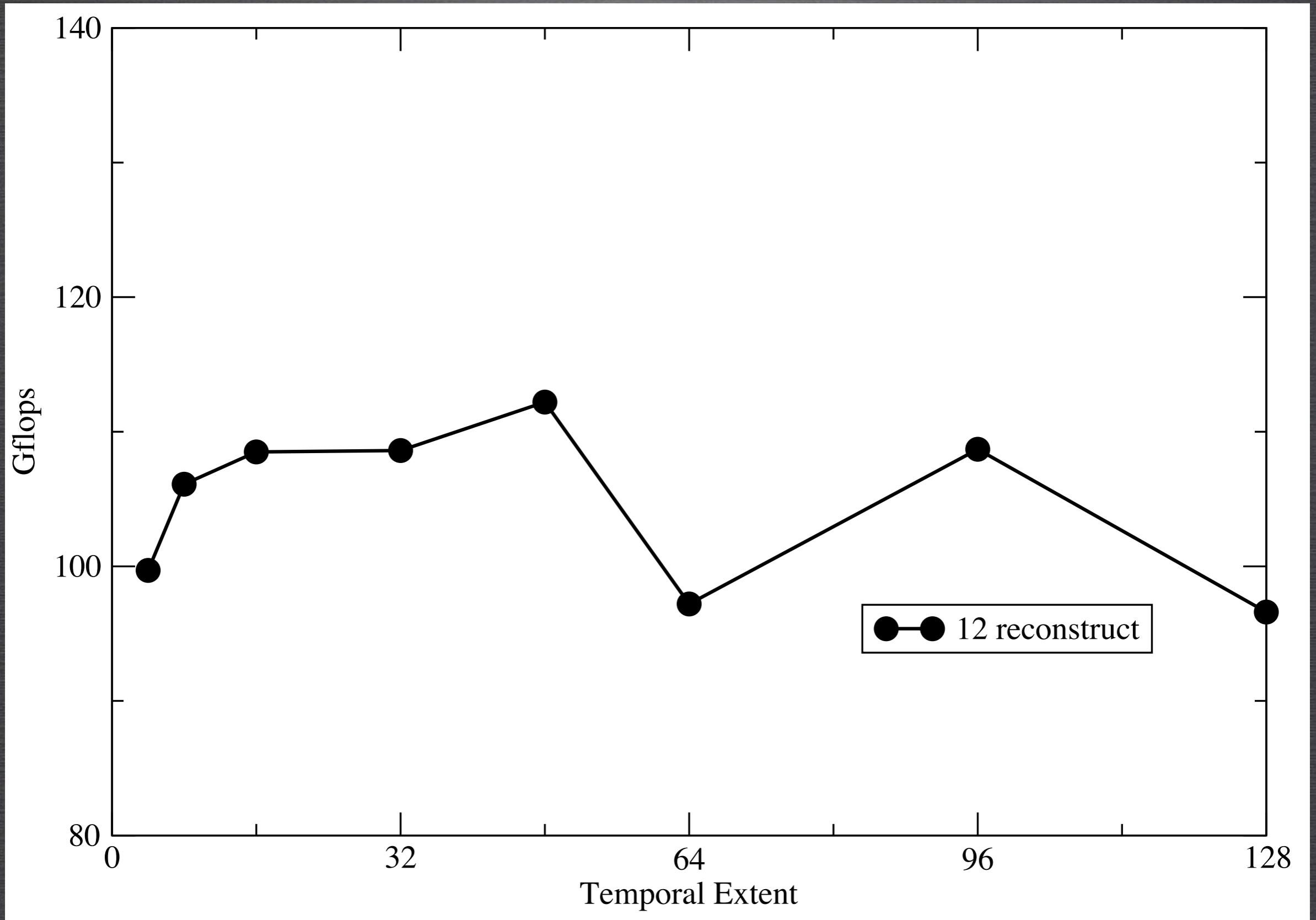
- Only need to load half-spinors in T direction
- Bytes per site **1440** \rightarrow **1248**

SU(3) REPRESENTATION

- SU(3) matrices are all unitary complex matrices with $\det = 1$
 - 18 real numbers, but only 8 free parameters (generators)
 - 12 number parameterization

$$\begin{pmatrix} a_1 & a_2 & a_3 \\ b_1 & b_2 & b_3 \\ c_1 & c_2 & c_3 \end{pmatrix} \rightarrow \begin{pmatrix} a_1 & a_2 & a_3 \\ b_1 & b_2 & b_3 \end{pmatrix} \quad \mathbf{c} = (\mathbf{a} \times \mathbf{b})^*$$

- Reconstruct full matrix on the fly
- 1152 Bytes per site (c.f. 1440)
- Additional 384 flops per site (c.f. 1368)



WILSON MATRIX-VECTOR PERFORMANCE

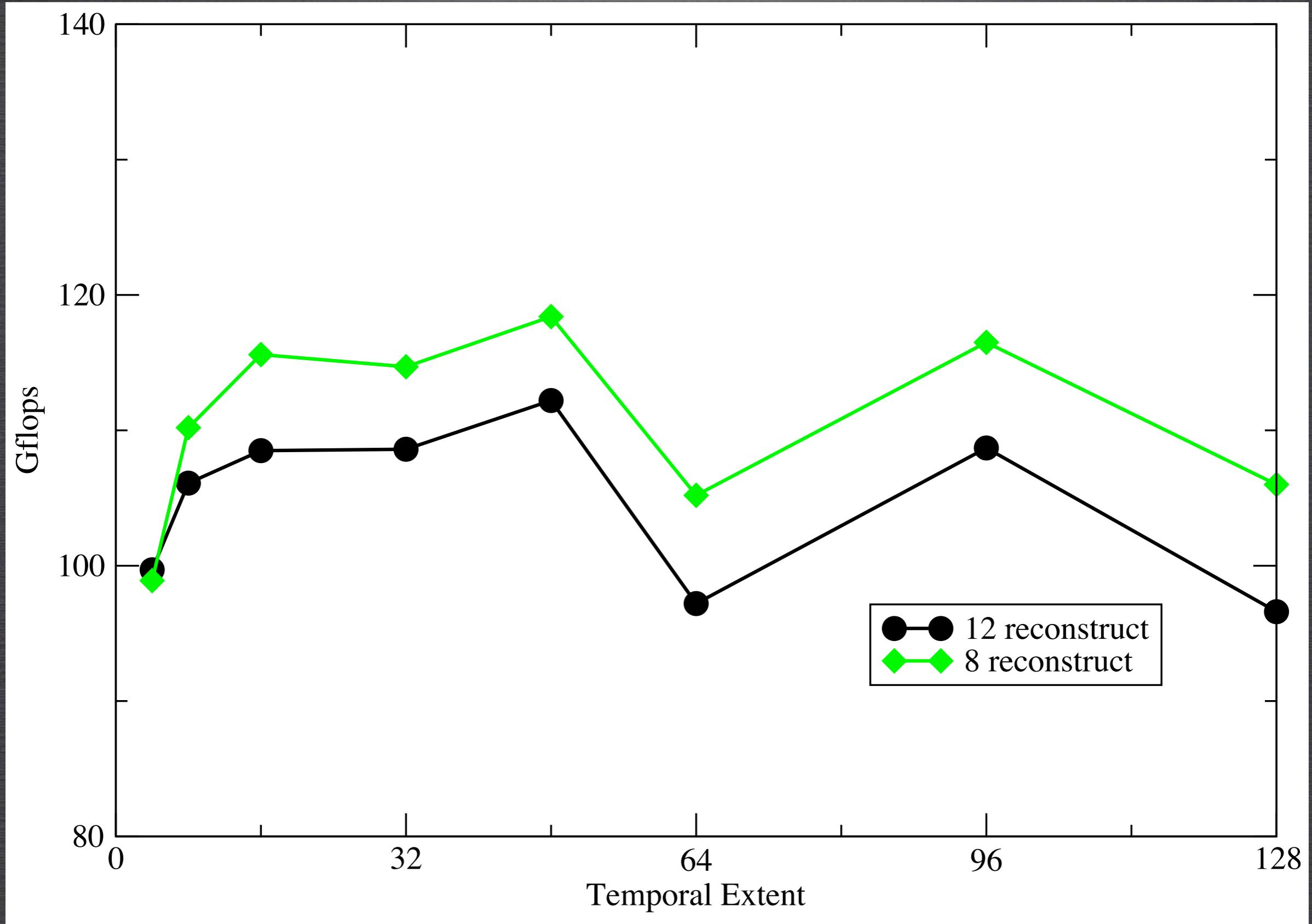
SINGLE PRECISION ($V=24^3 \times T$)

SU(3) REPRESENTATION

- Minimal 8 number parameterization

$$\begin{pmatrix} a_1 & a_2 & a_3 \\ b_1 & b_2 & b_3 \\ c_1 & c_2 & c_3 \end{pmatrix} \longrightarrow \begin{pmatrix} \arg(a_1) & \arg(c_1) & \operatorname{Re}(a_2) & \operatorname{Im}(a_2) \\ \operatorname{Re}(a_3) & \operatorname{Im}(a_3) & \operatorname{Re}(b_1) & \operatorname{Im}(b_1) \end{pmatrix}$$

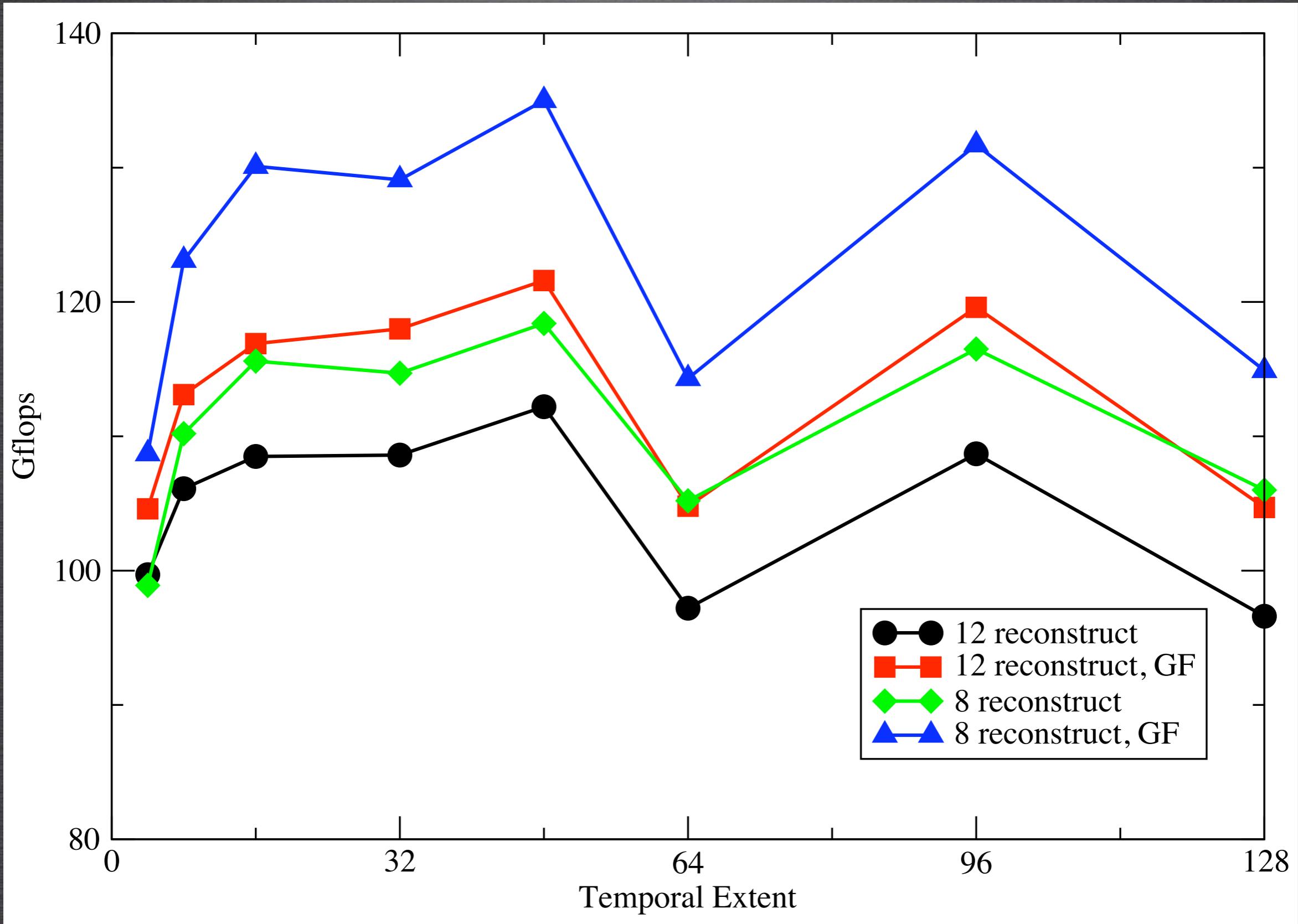
- Obtain a_1 and c_1 from normality
- Reconstruct b_2, b_3, c_2, c_3 from SU(2) rotation
- **1024** Bytes per site (c.f. 1440)
- Additional **856** flops per site
 - Including 2 sqrt, 4 trigonometric, 2 divide



WILSON MATRIX-VECTOR PERFORMANCE
SINGLE PRECISION ($V=24^3 \times T$)

MORE TRICKS

- Impose local gauge transformation -> temporal gauge
 - SU(3) field = unit matrix in temporal direction
 - Must calculate this transformation (done once only)
- 960 Bytes per site (c.f. 1440)
- In total: 33% bandwidth reduction



WILSON MATRIX-VECTOR PERFORMANCE

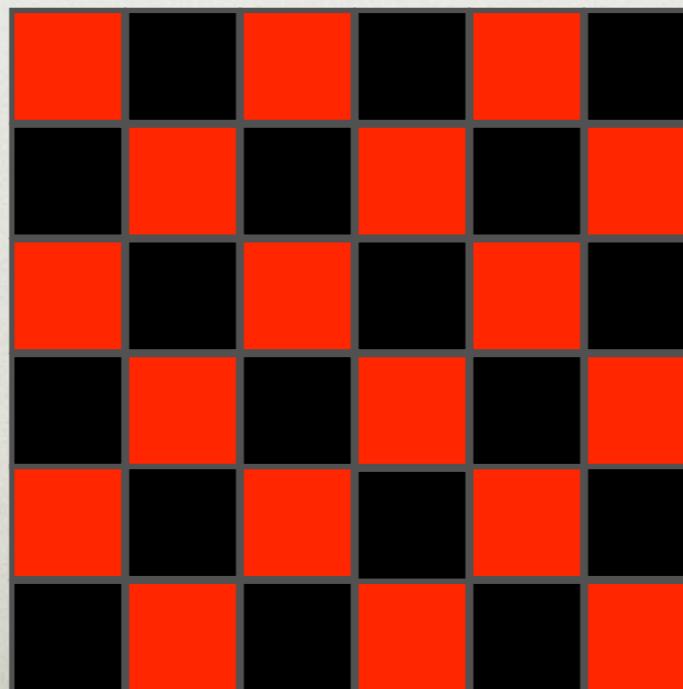
SINGLE PRECISION ($V=24^3 \times T$)

SPINOR SHARING

- Stencil methods on GPUs have large gains when vector components are shared between threads
 - Load spinor components into shared memory
 - Share spinors between threads
- Can we reuse some of the spinor loads?
 - Even-odd updating complicates things
 - 24 components per site -> Limited by shared memory

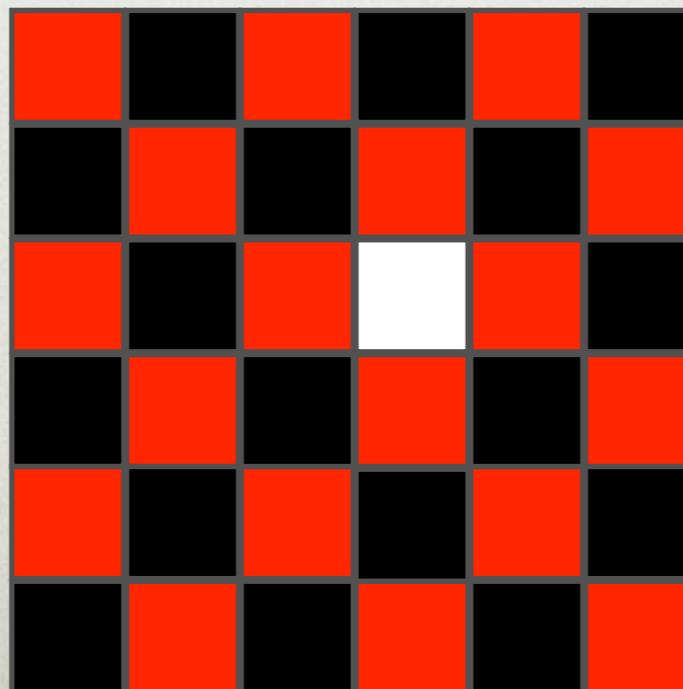
SPINOR SHARING

- Stencil methods on GPUs have large gains when vector components are shared between threads
 - Load spinor components into shared memory
 - Share spinors between threads
- Can we reuse some of the spinor loads?
 - Even-odd updating complicates things
 - 24 components per site -> Limited by shared memory



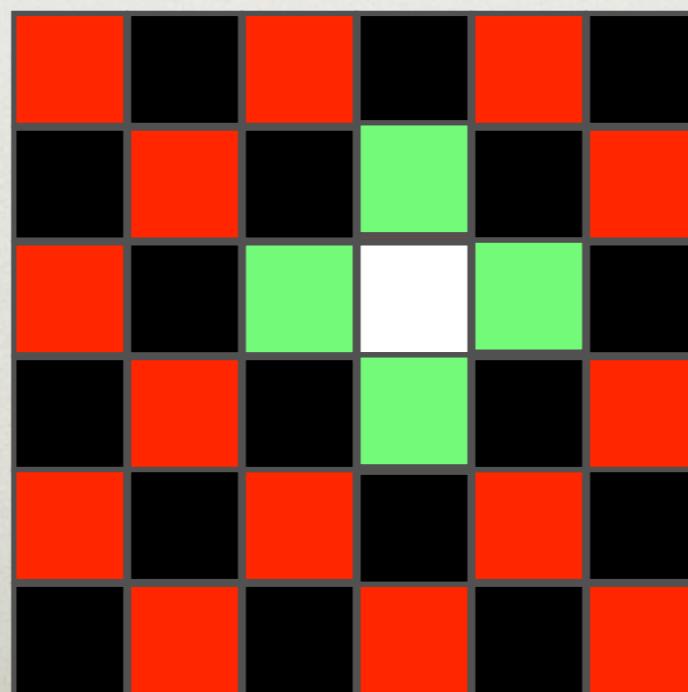
SPINOR SHARING

- Stencil methods on GPUs have large gains when vector components are shared between threads
 - Load spinor components into shared memory
 - Share spinors between threads
- Can we reuse some of the spinor loads?
 - Even-odd updating complicates things
 - 24 components per site -> Limited by shared memory



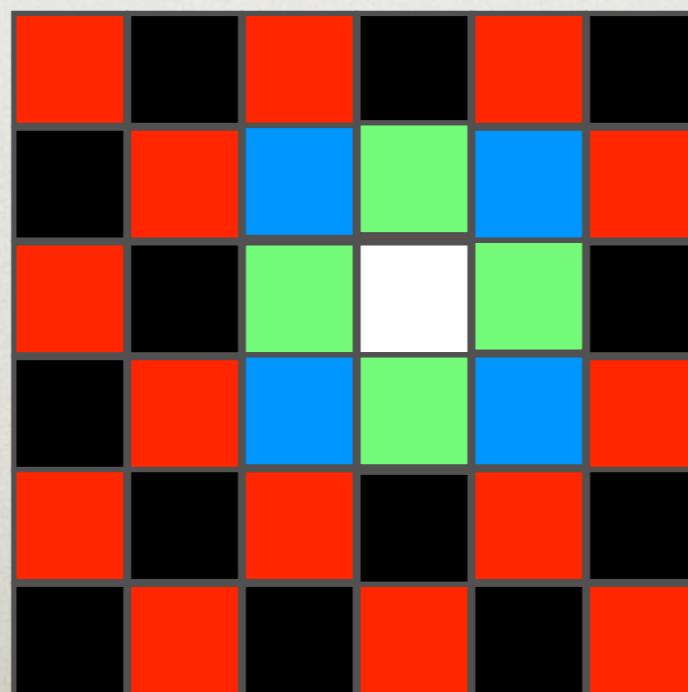
SPINOR SHARING

- Stencil methods on GPUs have large gains when vector components are shared between threads
 - Load spinor components into shared memory
 - Share spinors between threads
- Can we reuse some of the spinor loads?
 - Even-odd updating complicates things
 - 24 components per site -> Limited by shared memory



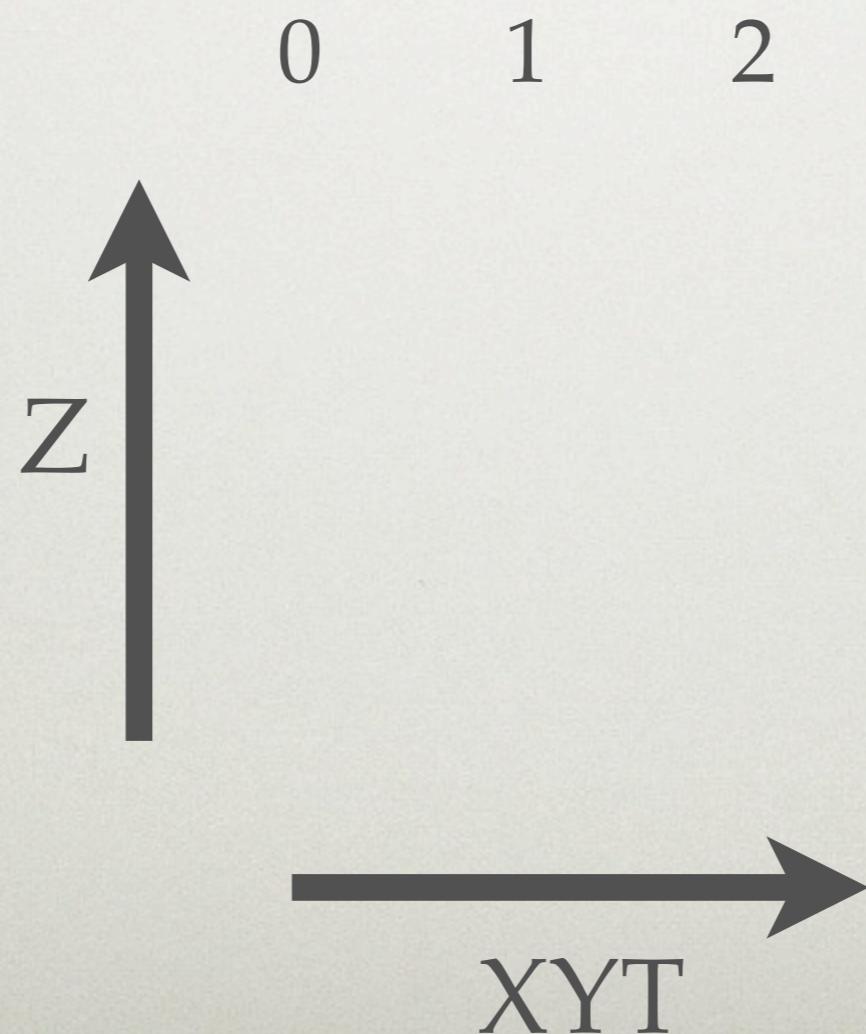
SPINOR SHARING

- Stencil methods on GPUs have large gains when vector components are shared between threads
 - Load spinor components into shared memory
 - Share spinors between threads
- Can we reuse some of the spinor loads?
 - Even-odd updating complicates things
 - 24 components per site -> Limited by shared memory



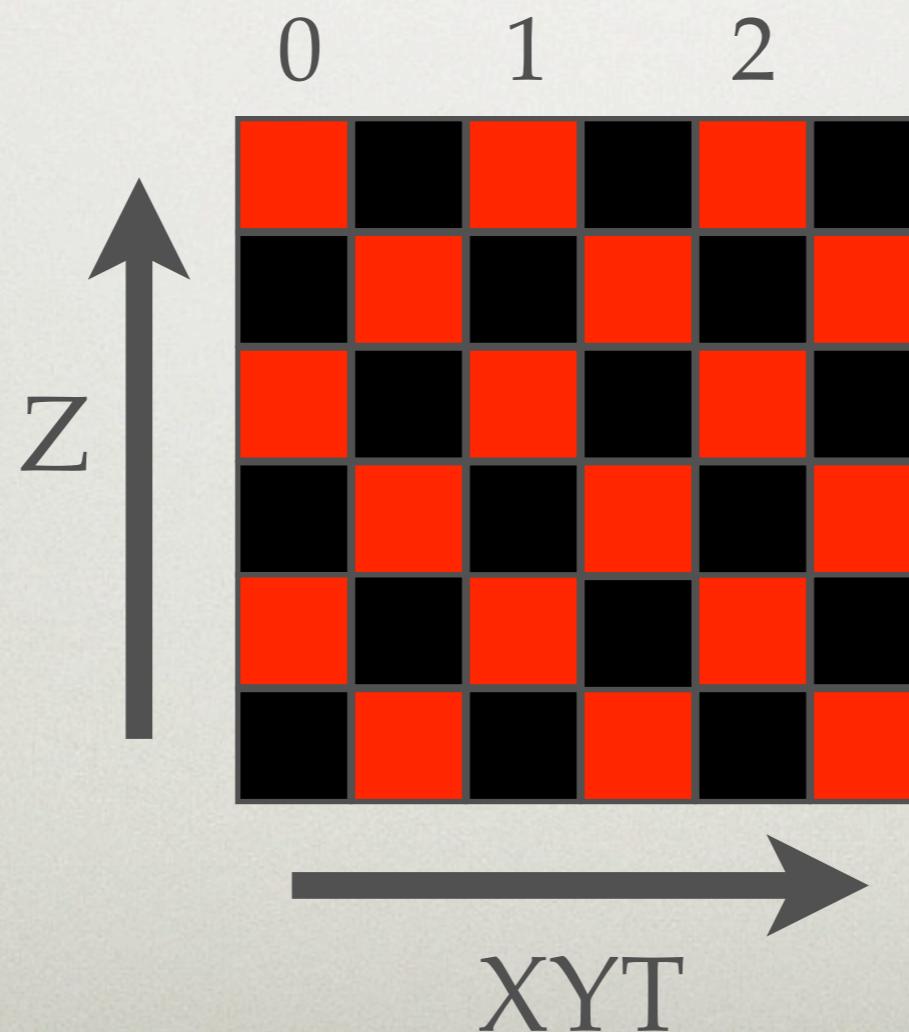
Z STREAMING

- Easy to reuse loads in one dimension
- Use XYT threads, each thread streams down Z dimension



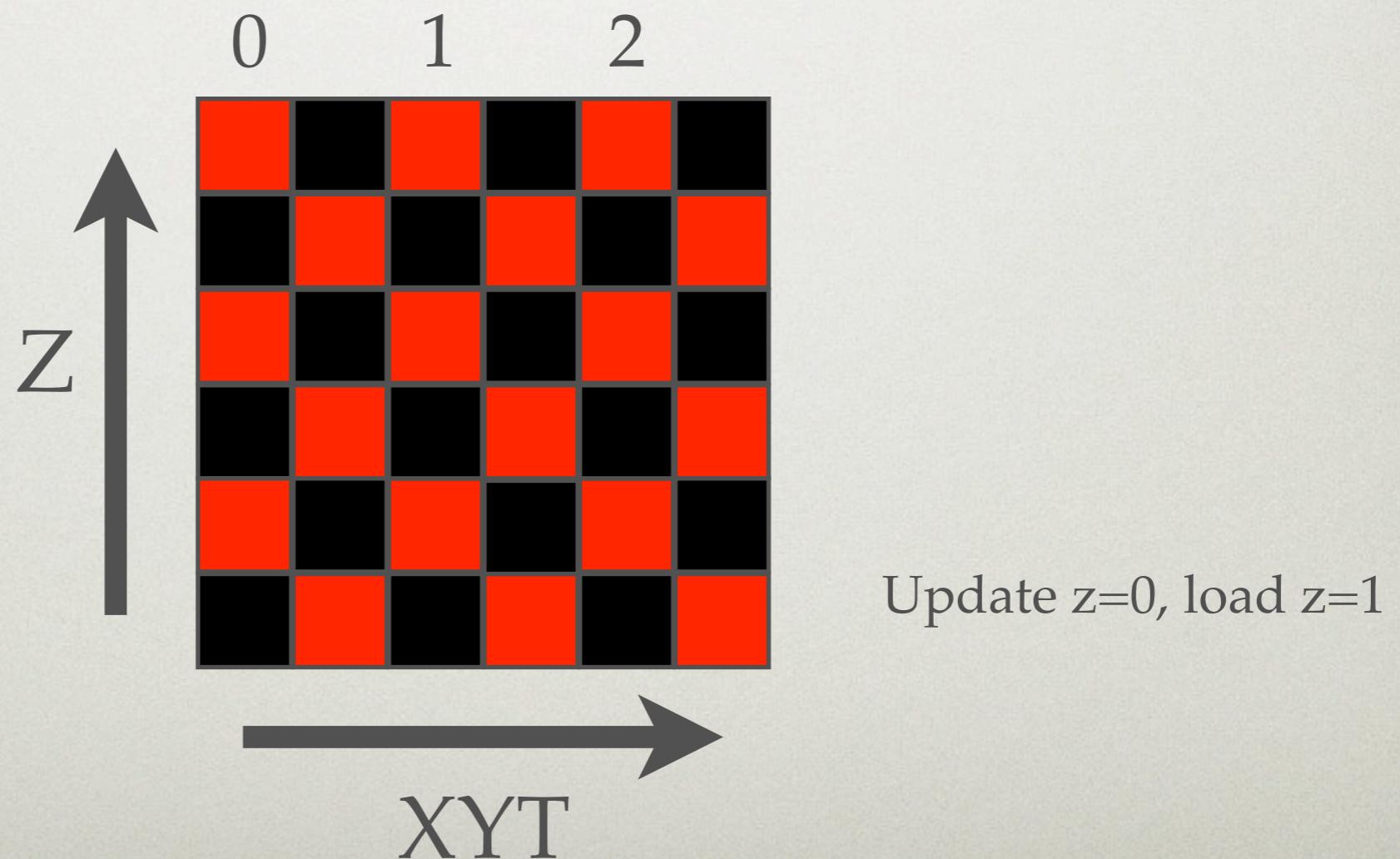
Z STREAMING

- Easy to reuse loads in one dimension
- Use XYT threads, each thread streams down Z dimension



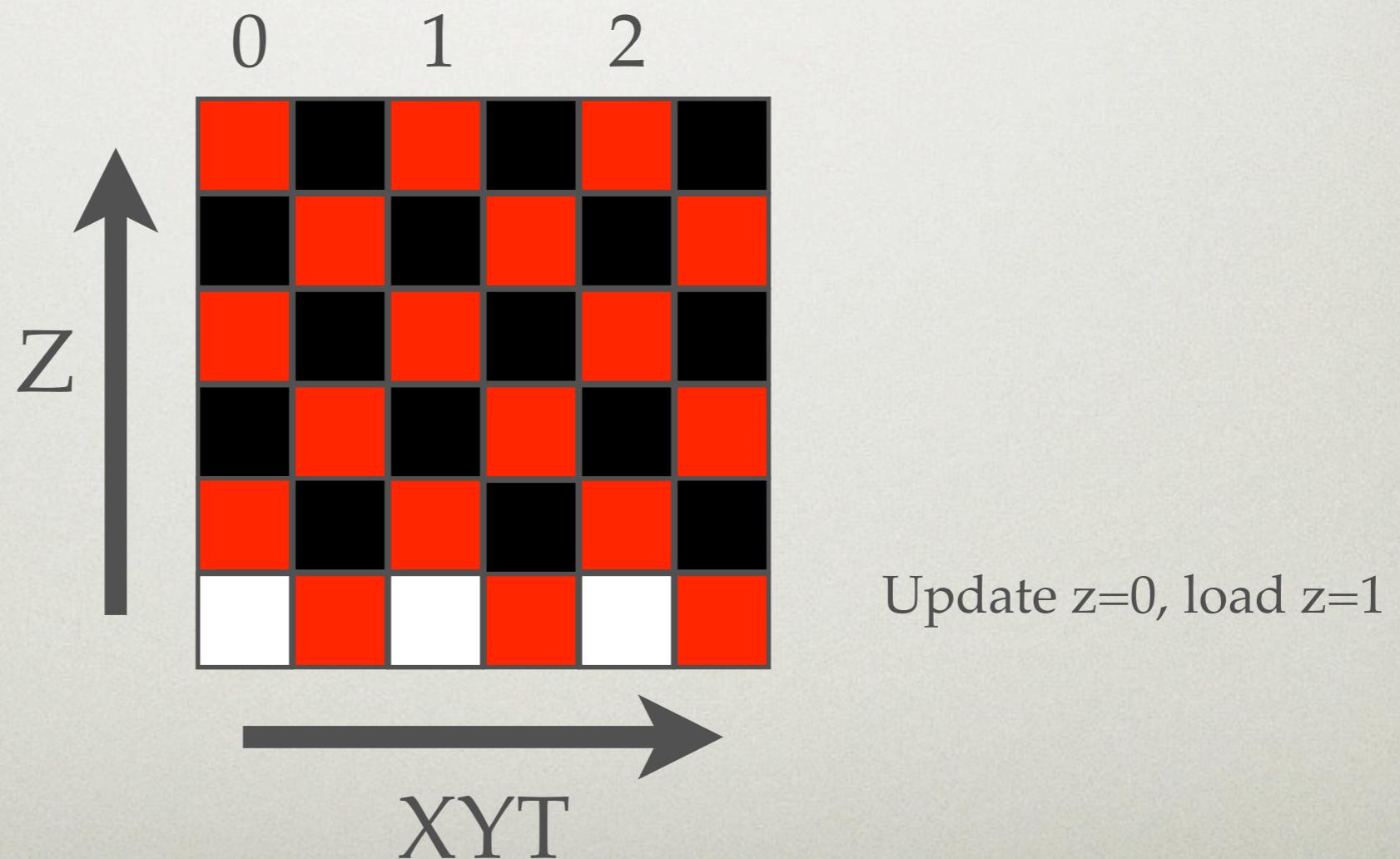
Z STREAMING

- Easy to reuse loads in one dimension
- Use XYT threads, each thread streams down Z dimension



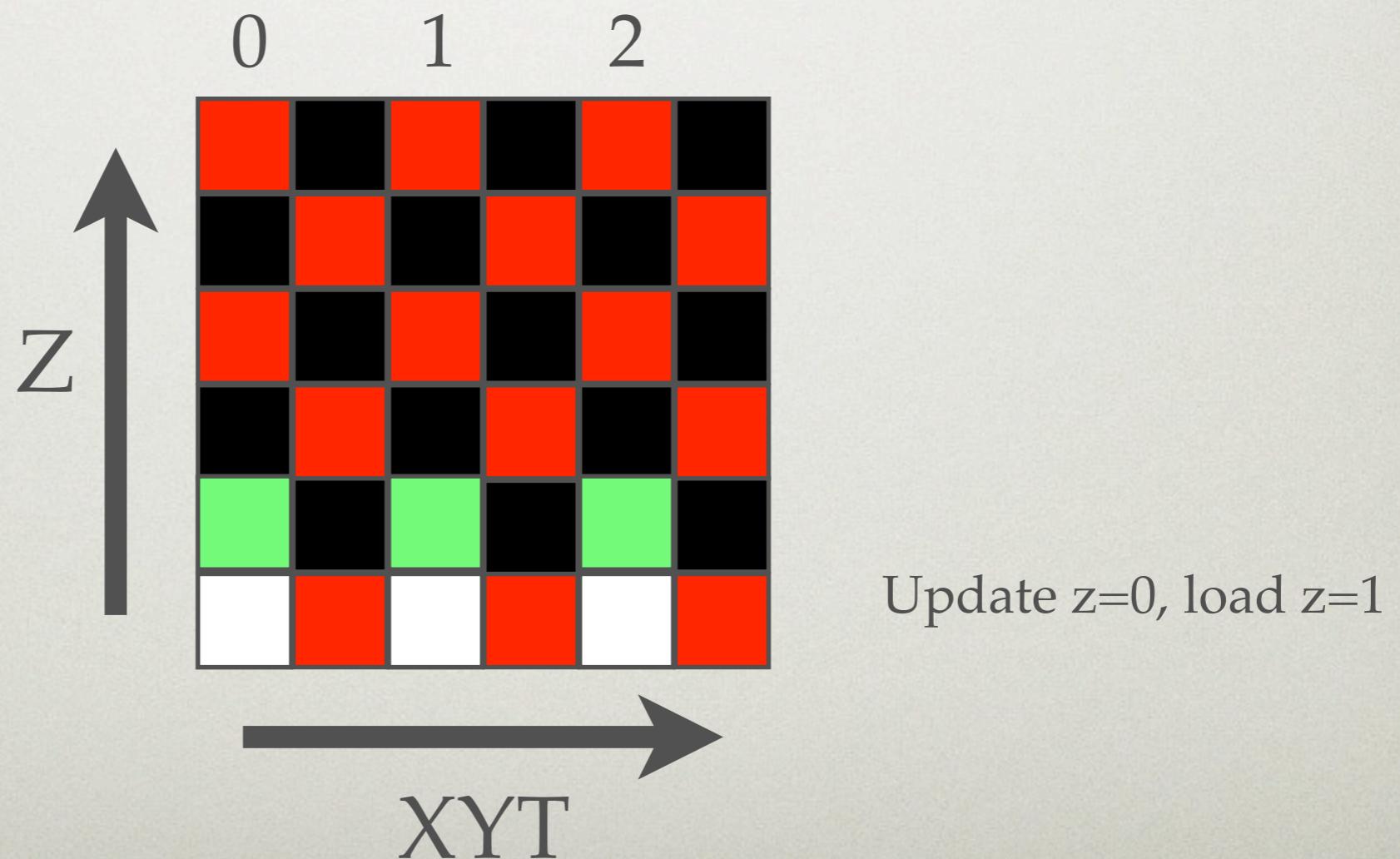
Z STREAMING

- Easy to reuse loads in one dimension
- Use XYT threads, each thread streams down Z dimension



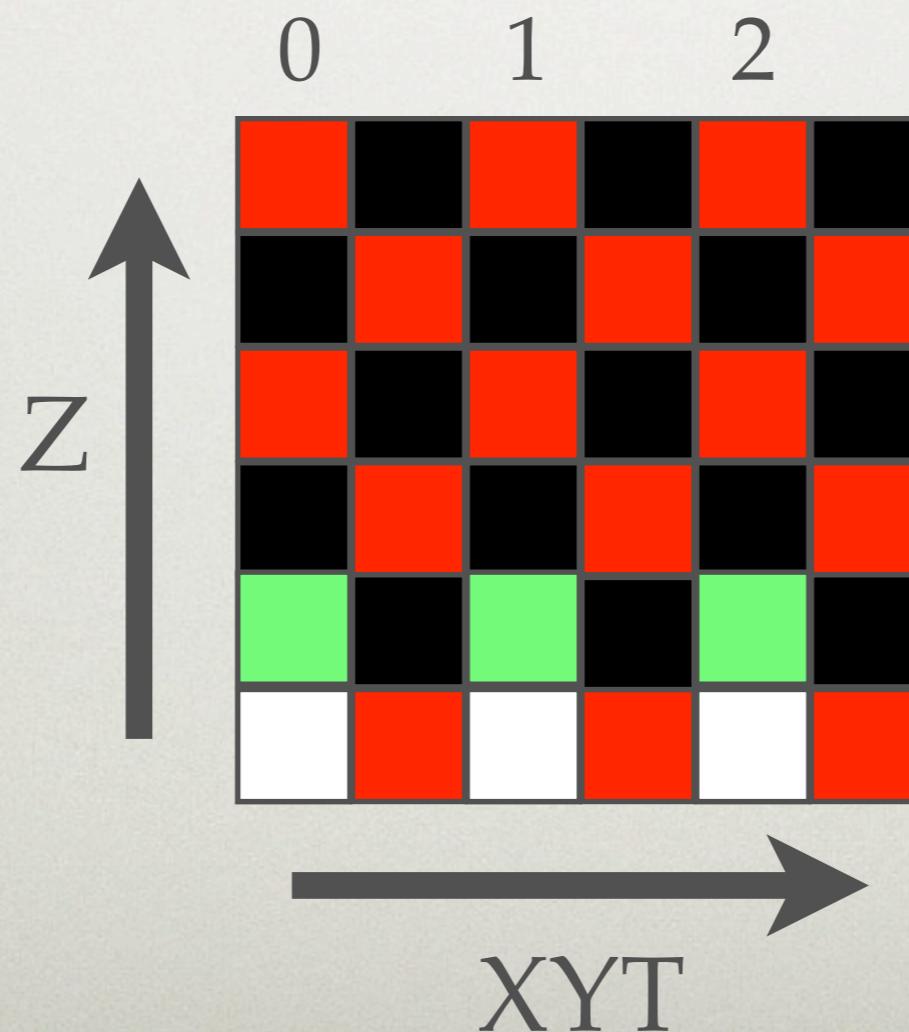
Z STREAMING

- Easy to reuse loads in one dimension
- Use XYT threads, each thread streams down Z dimension



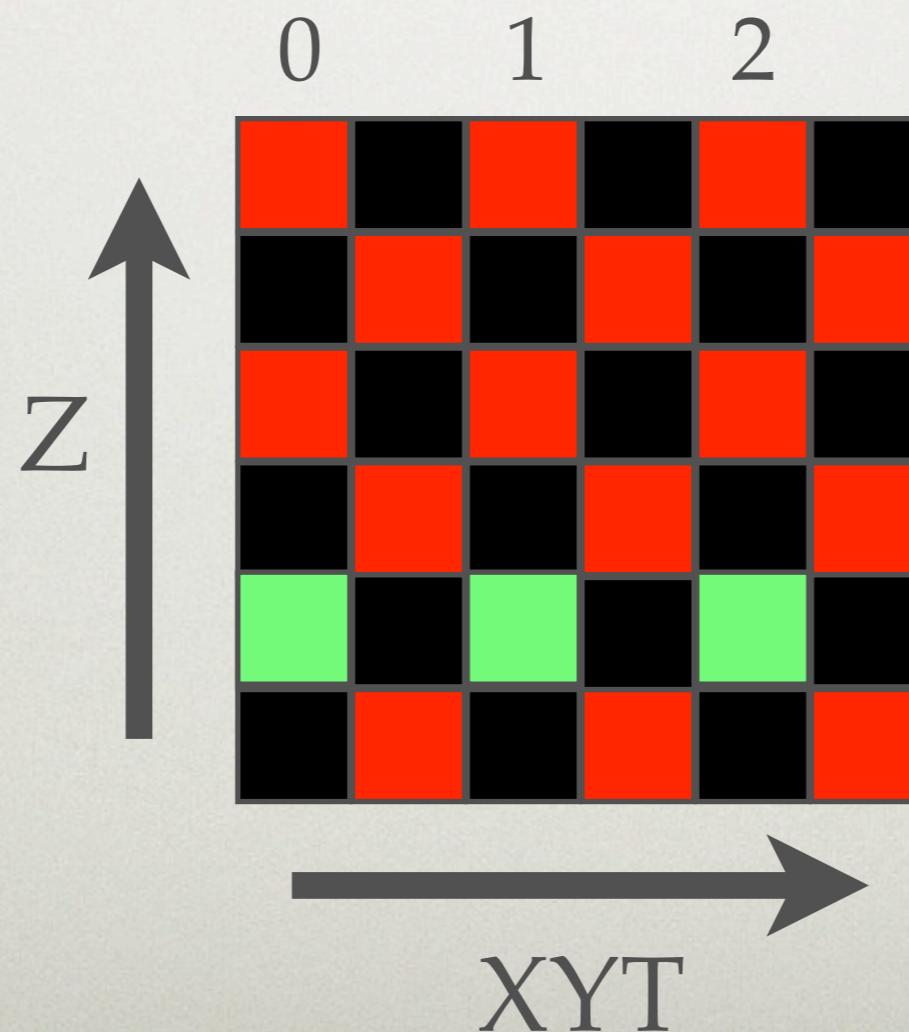
Z STREAMING

- Easy to reuse loads in one dimension
 - Use XYT threads, each thread streams down Z dimension



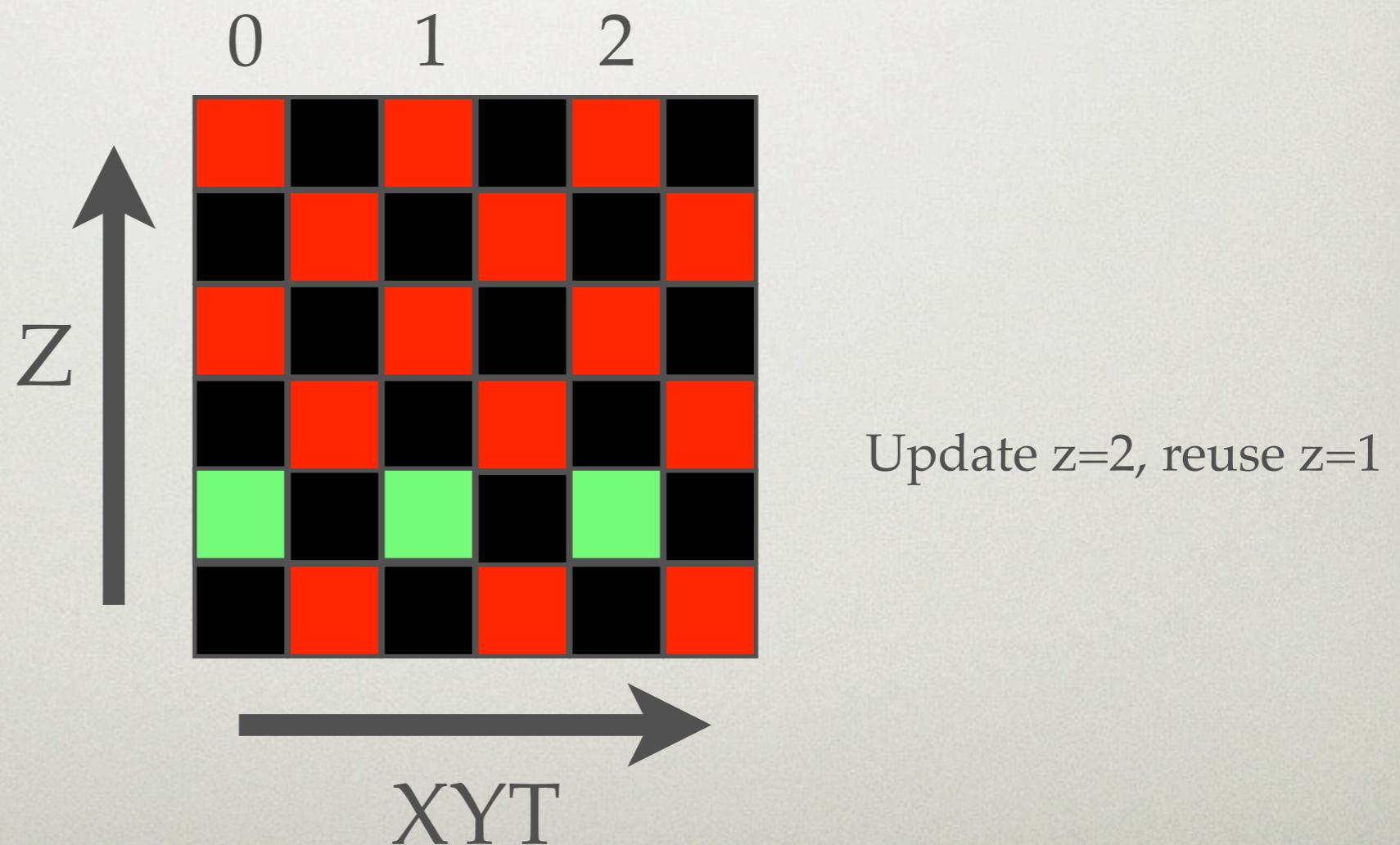
Z STREAMING

- Easy to reuse loads in one dimension
- Use XYT threads, each thread streams down Z dimension



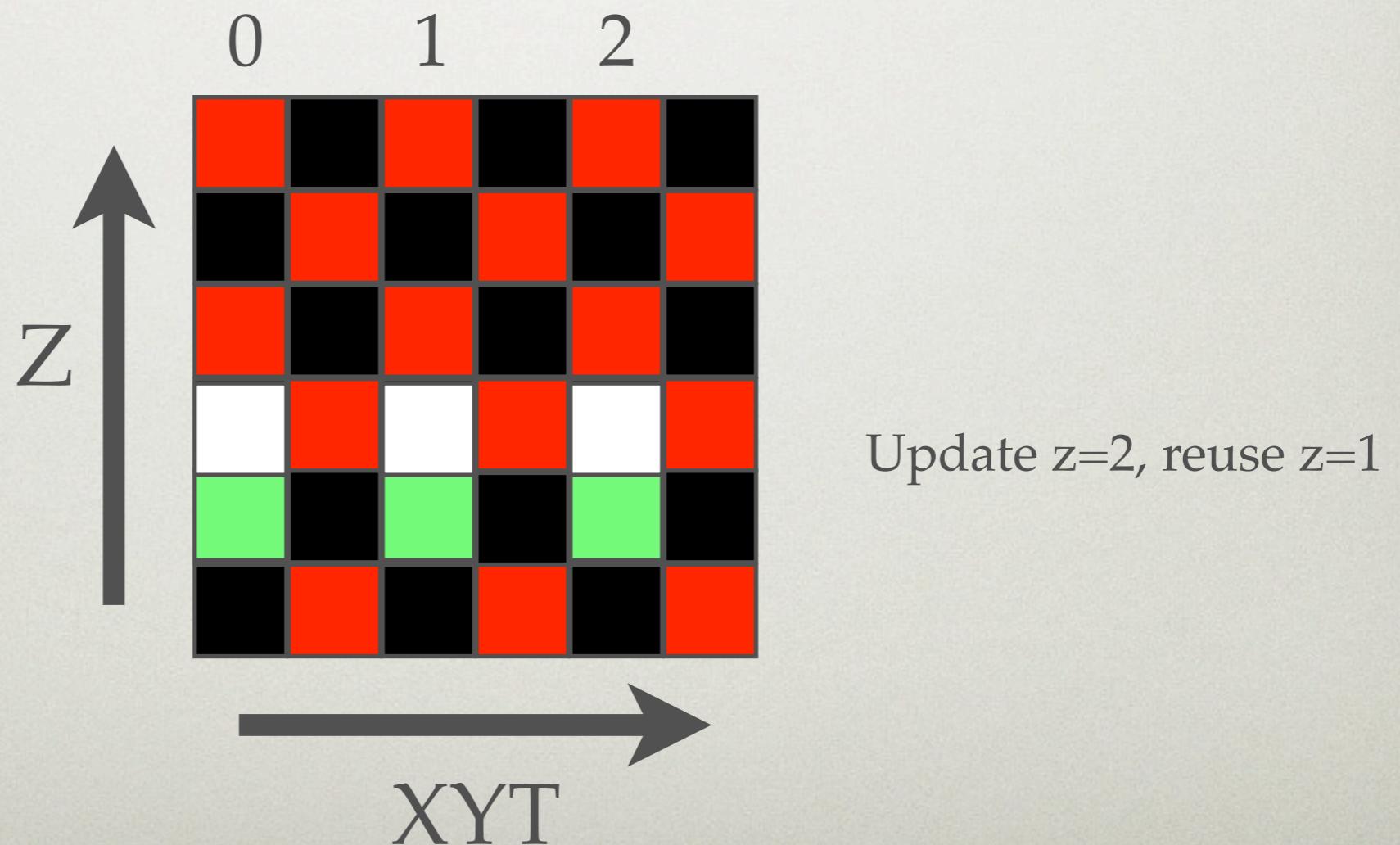
Z STREAMING

- Easy to reuse loads in one dimension
- Use XYT threads, each thread streams down Z dimension



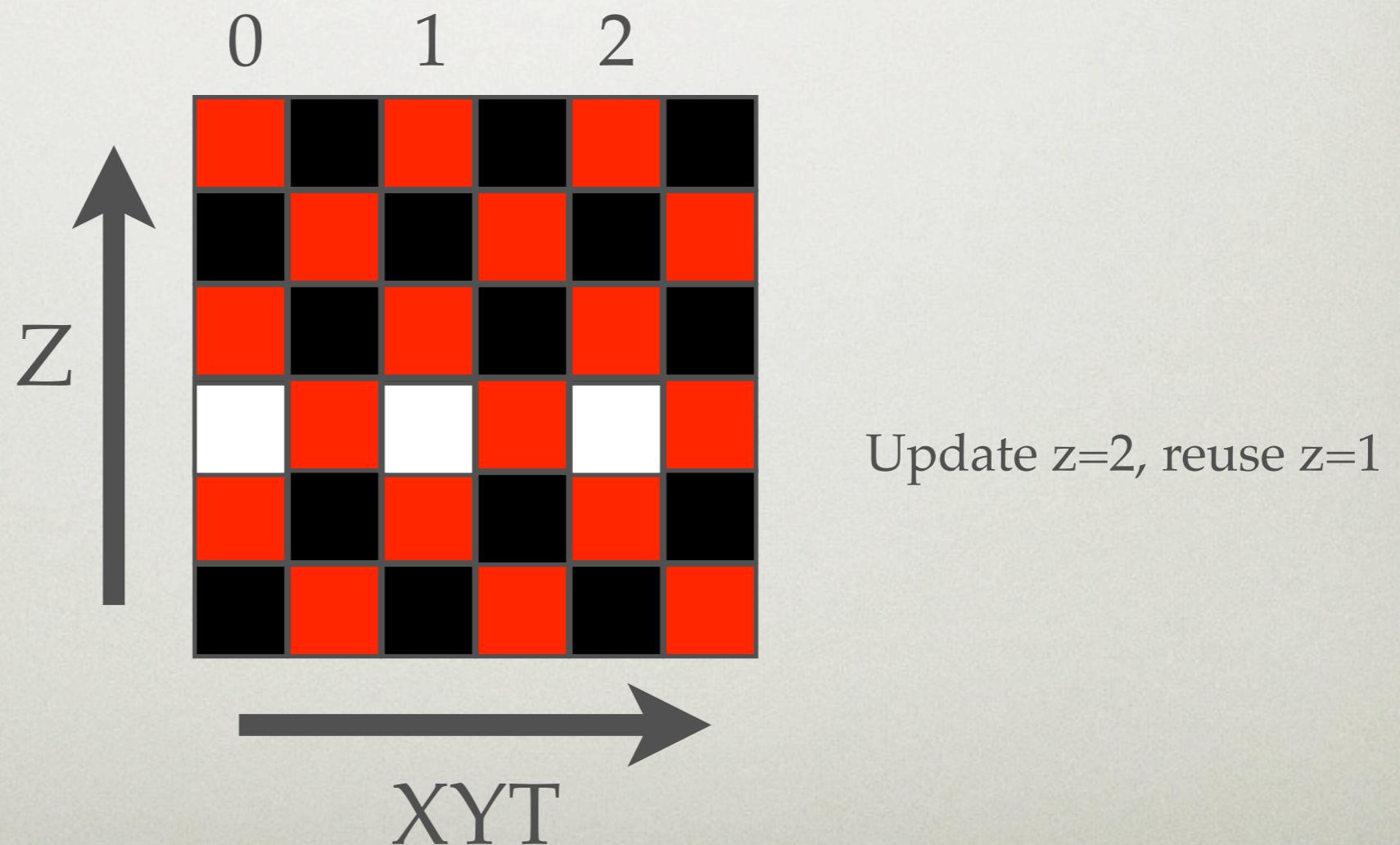
Z STREAMING

- Easy to reuse loads in one dimension
- Use XYT threads, each thread streams down Z dimension



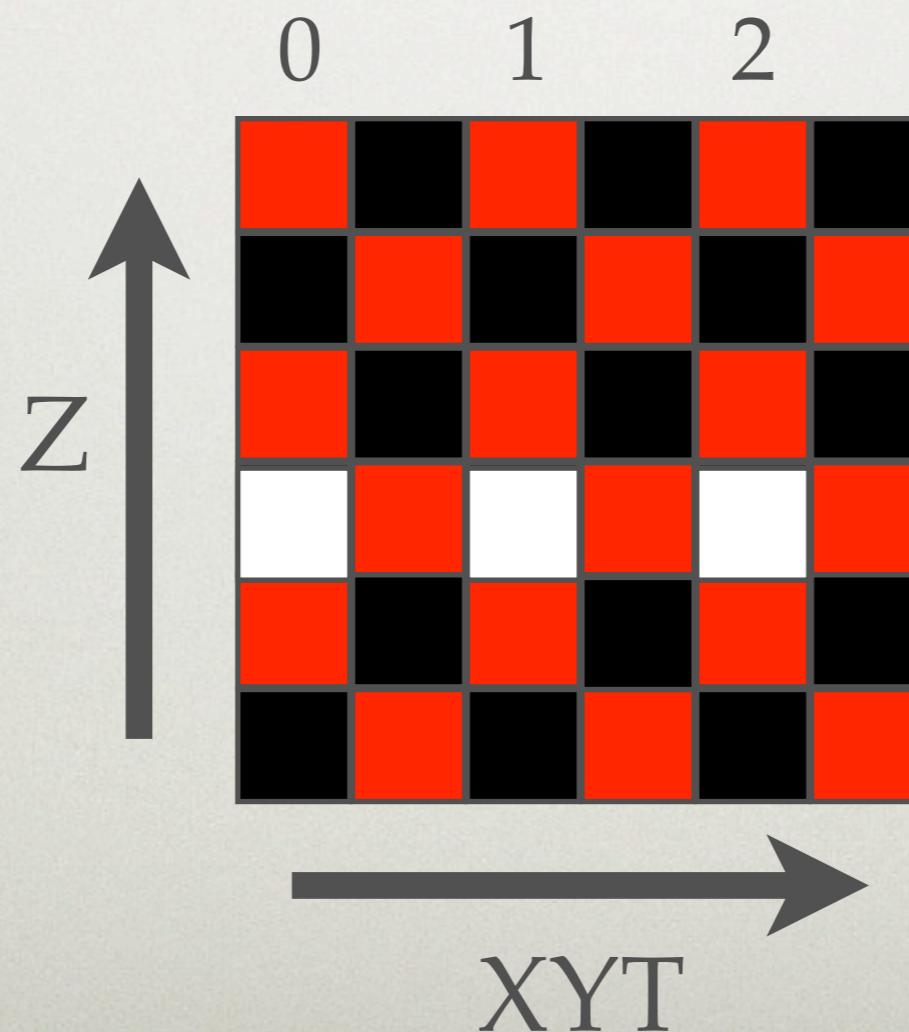
Z STREAMING

- Easy to reuse loads in one dimension
- Use XYT threads, each thread streams down Z dimension



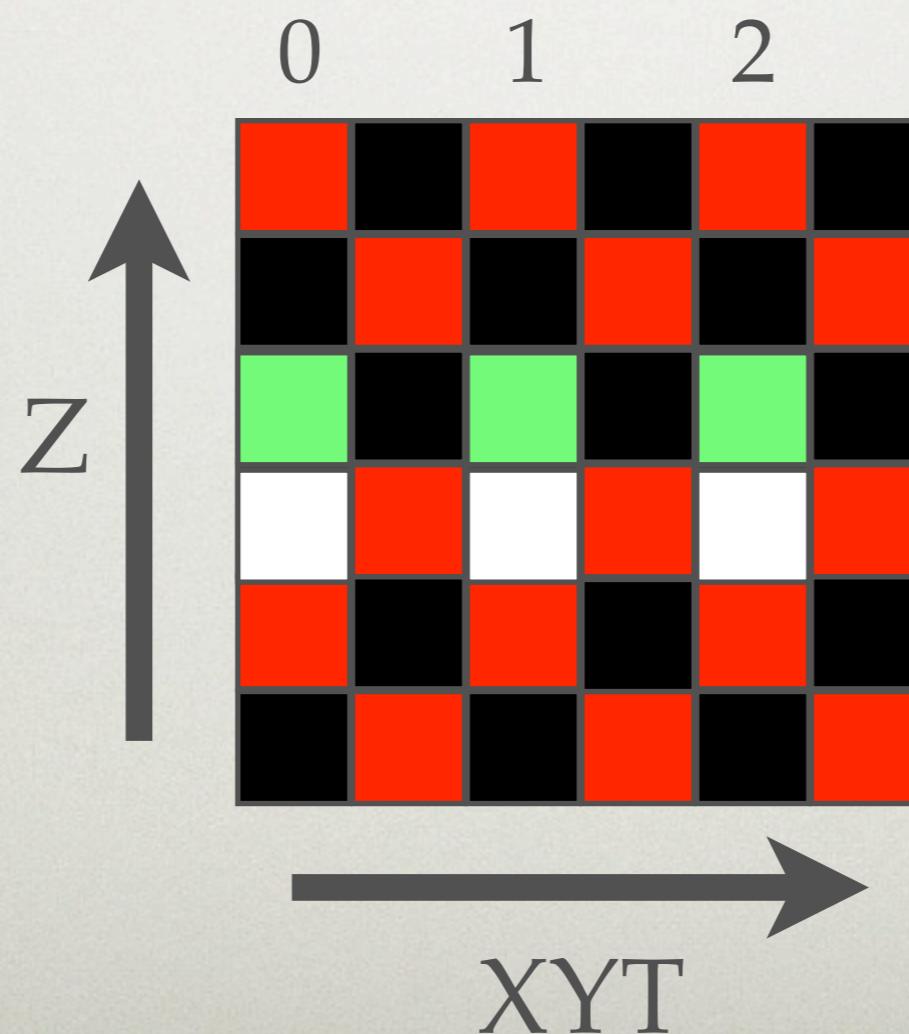
Z STREAMING

- Easy to reuse loads in one dimension
- Use XYT threads, each thread streams down Z dimension



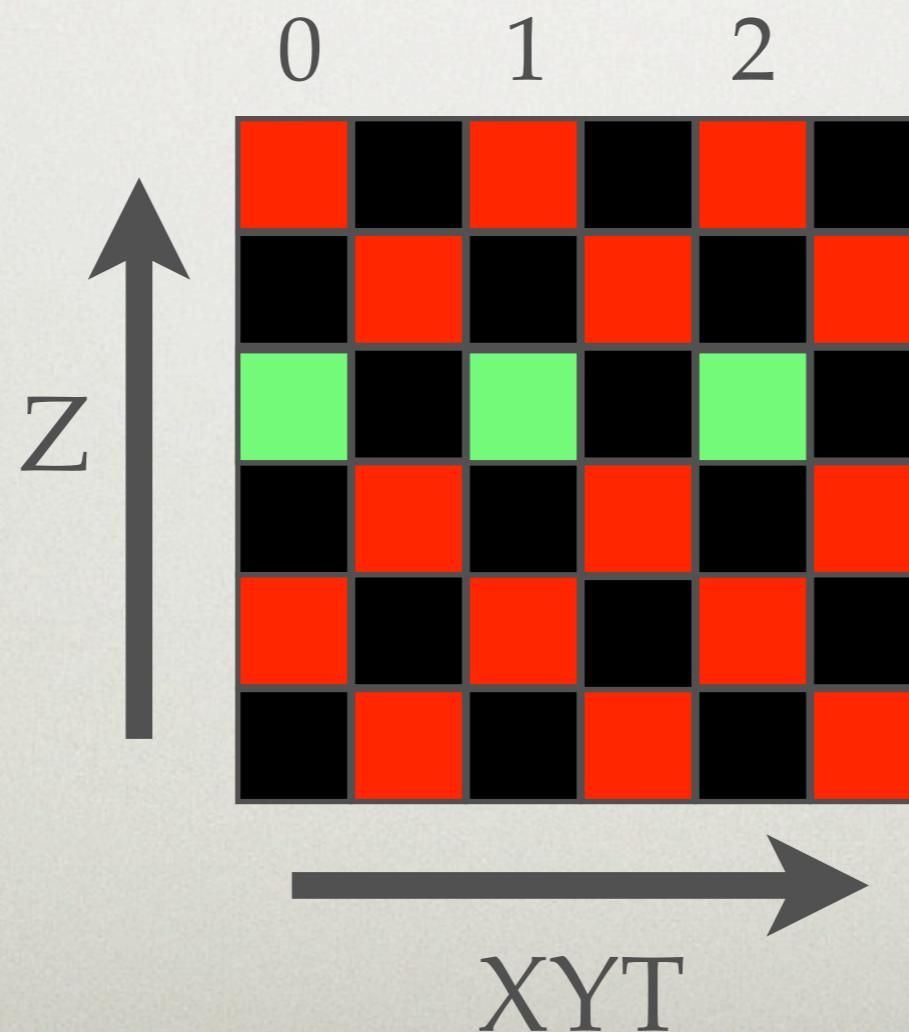
Z STREAMING

- Easy to reuse loads in one dimension
- Use XYT threads, each thread streams down Z dimension



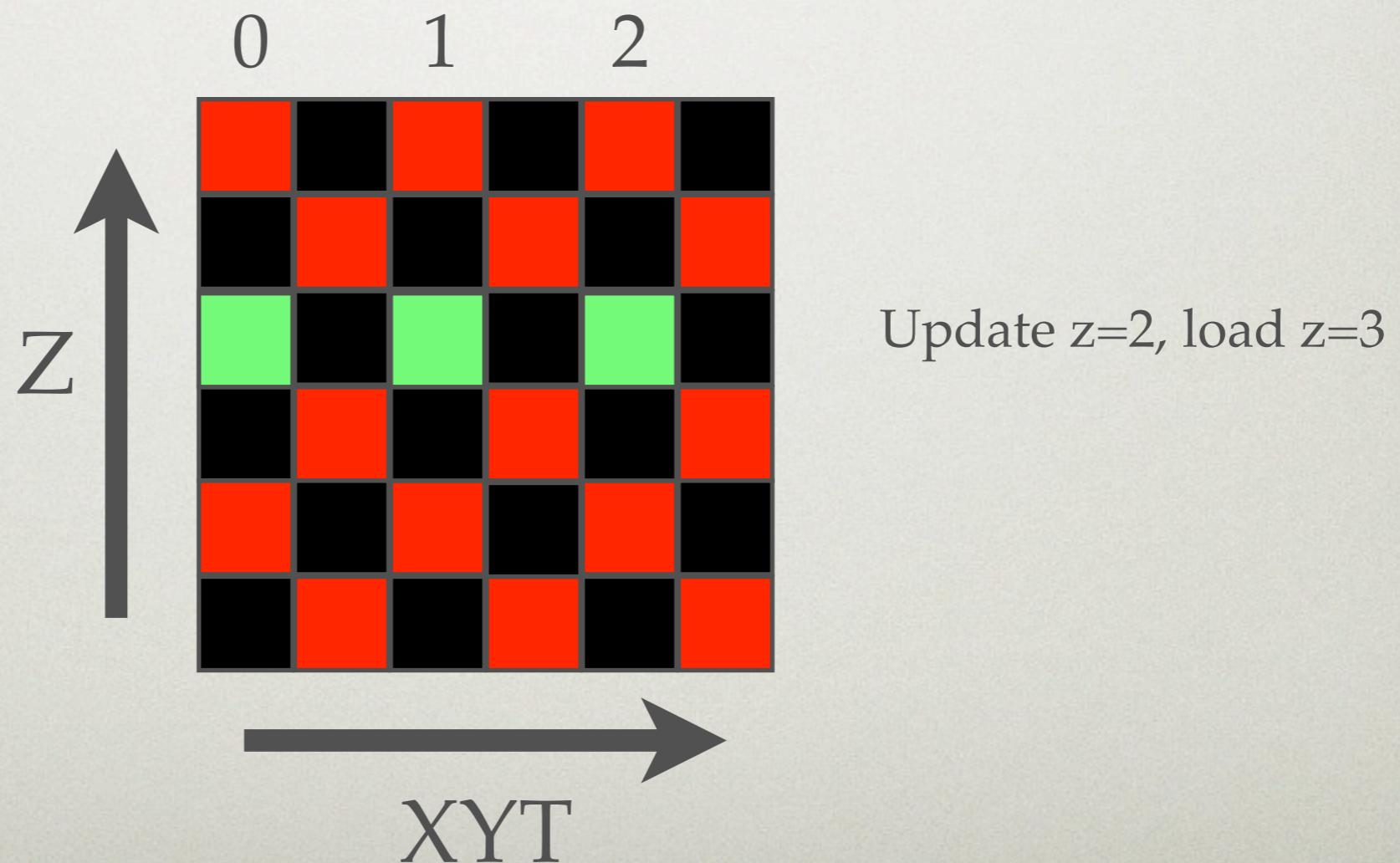
Z STREAMING

- Easy to reuse loads in one dimension
- Use XYT threads, each thread streams down Z dimension



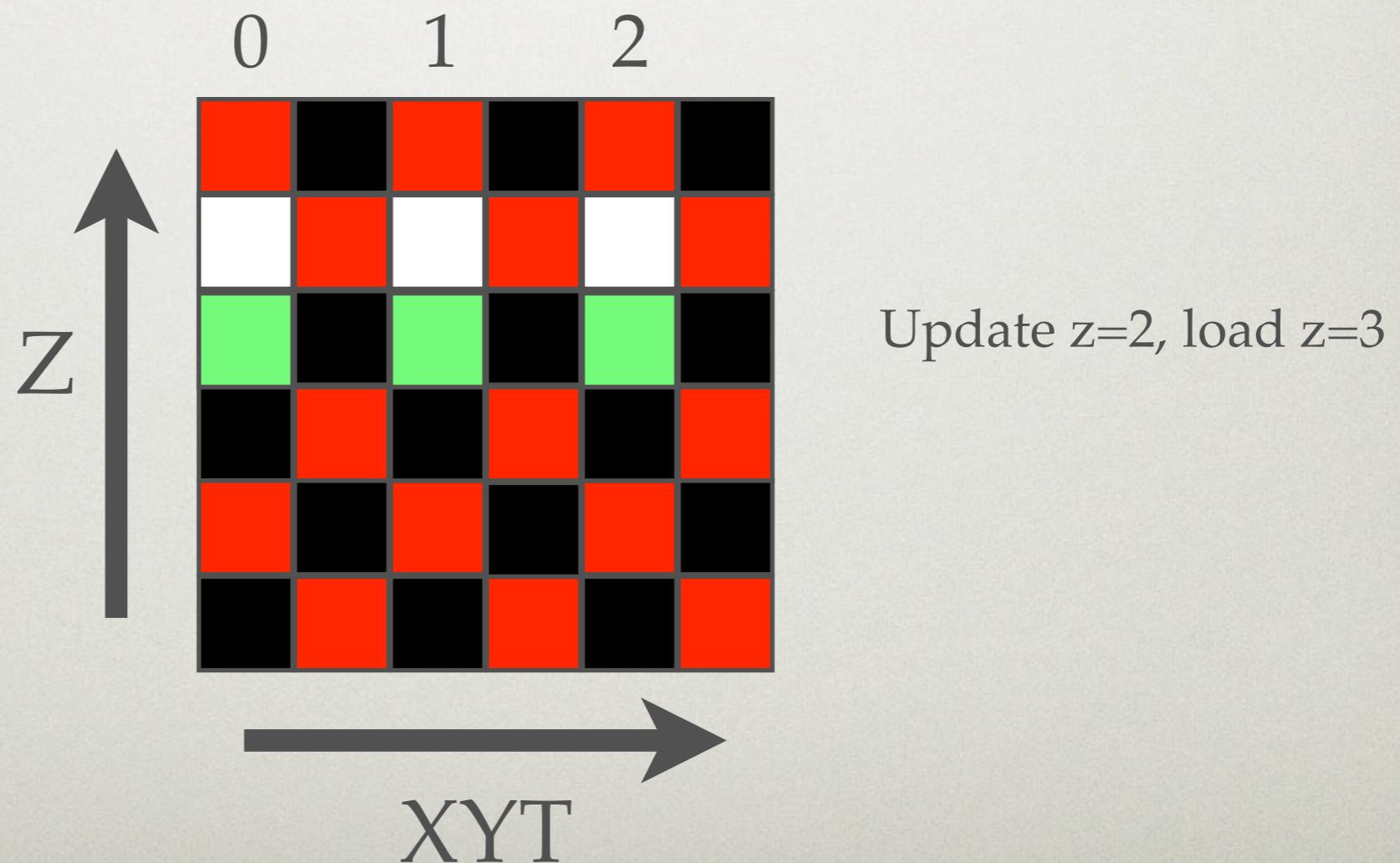
Z STREAMING

- Easy to reuse loads in one dimension
- Use XYT threads, each thread streams down Z dimension



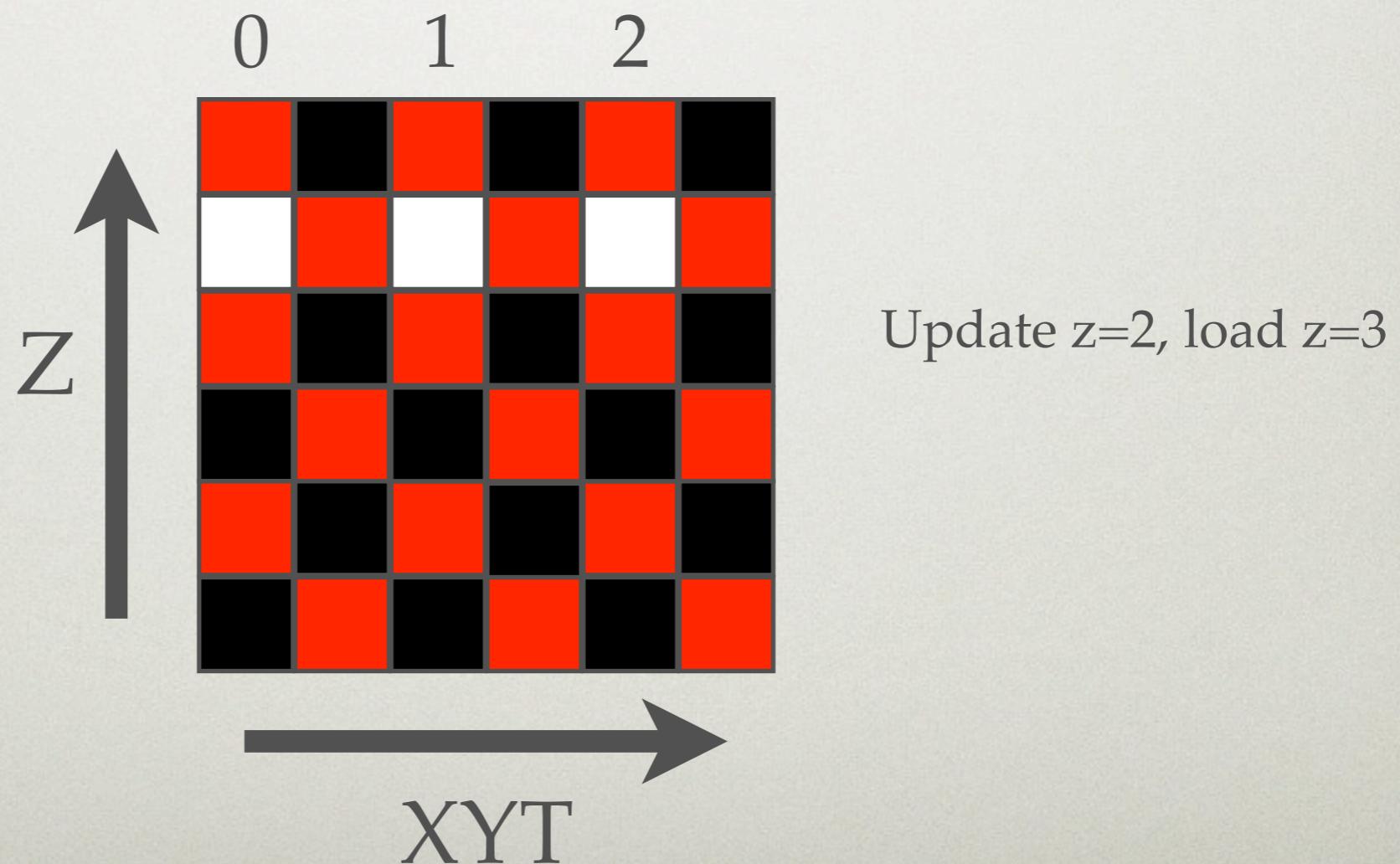
Z STREAMING

- Easy to reuse loads in one dimension
- Use XYT threads, each thread streams down Z dimension



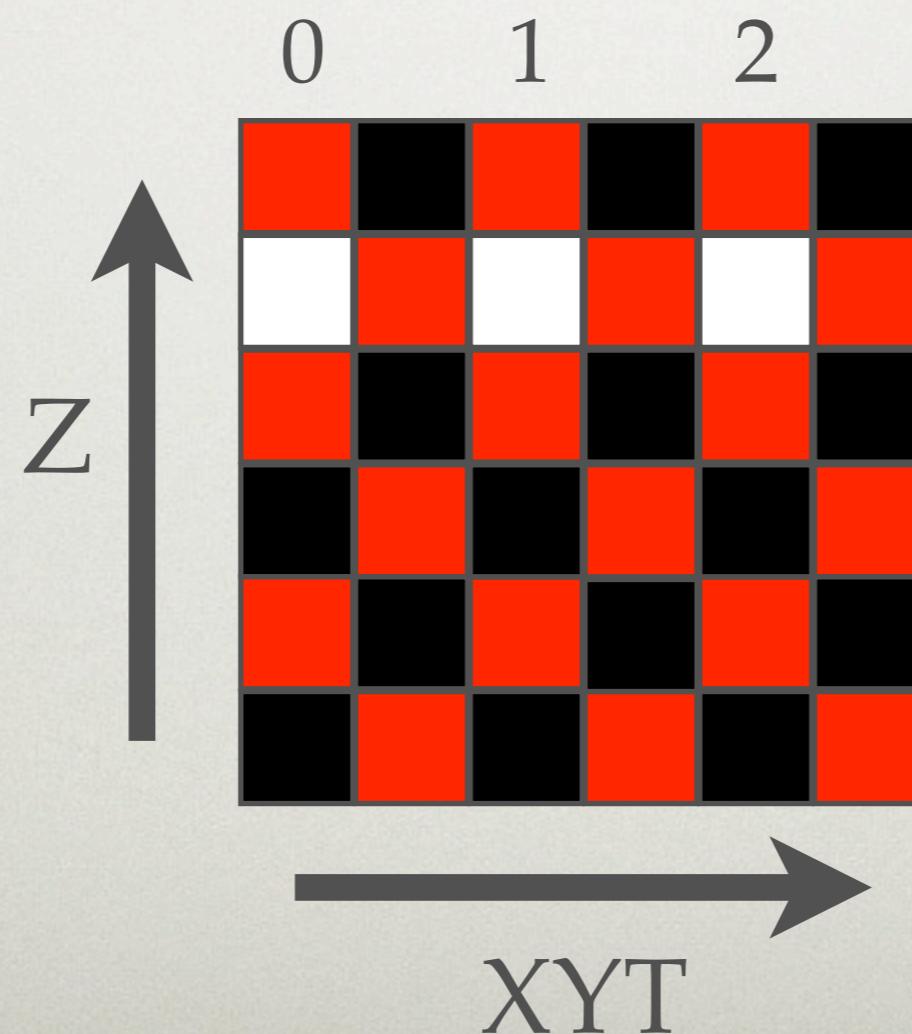
Z STREAMING

- Easy to reuse loads in one dimension
- Use XYT threads, each thread streams down Z dimension



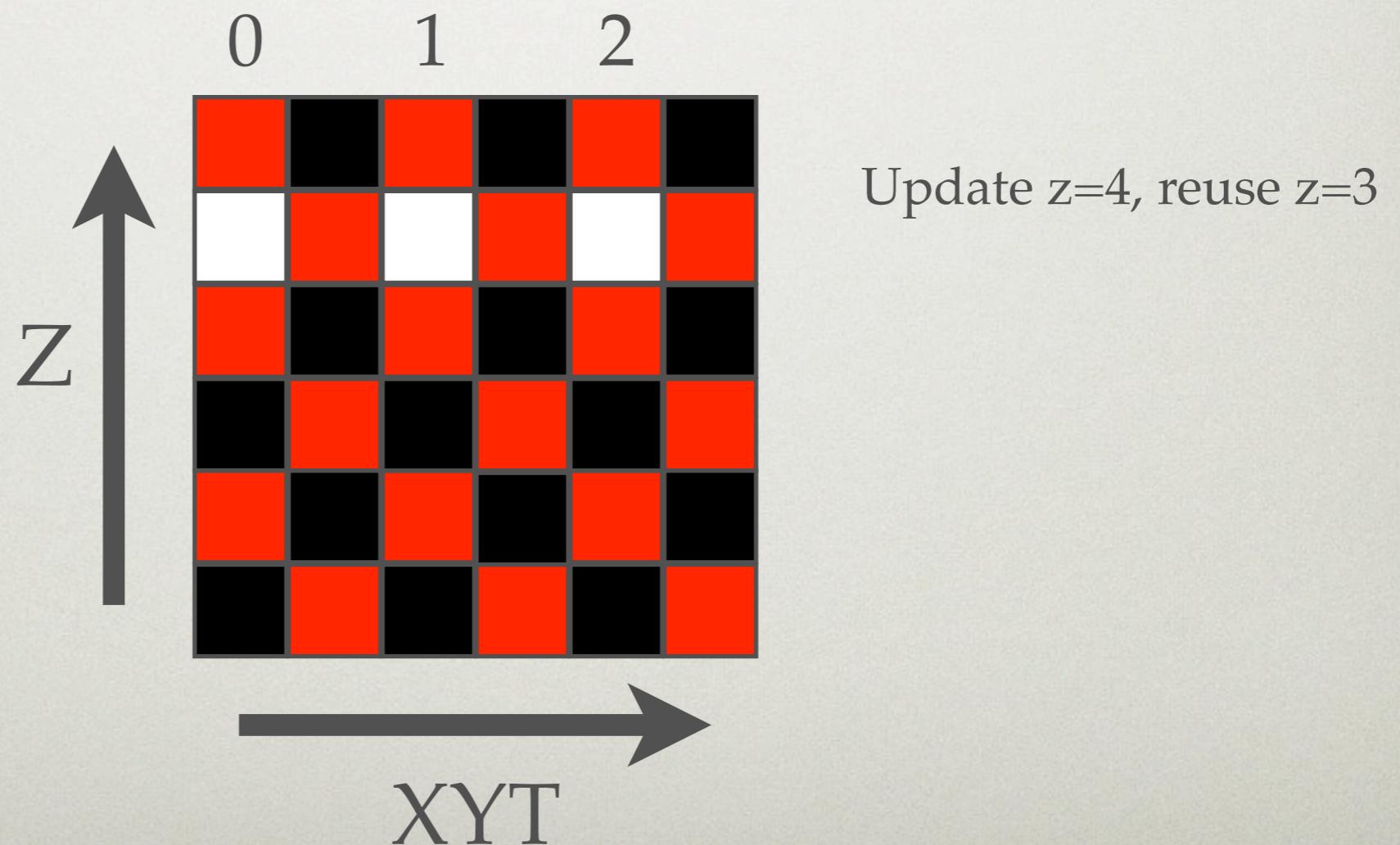
Z STREAMING

- Easy to reuse loads in one dimension
- Use XYT threads, each thread streams down Z dimension



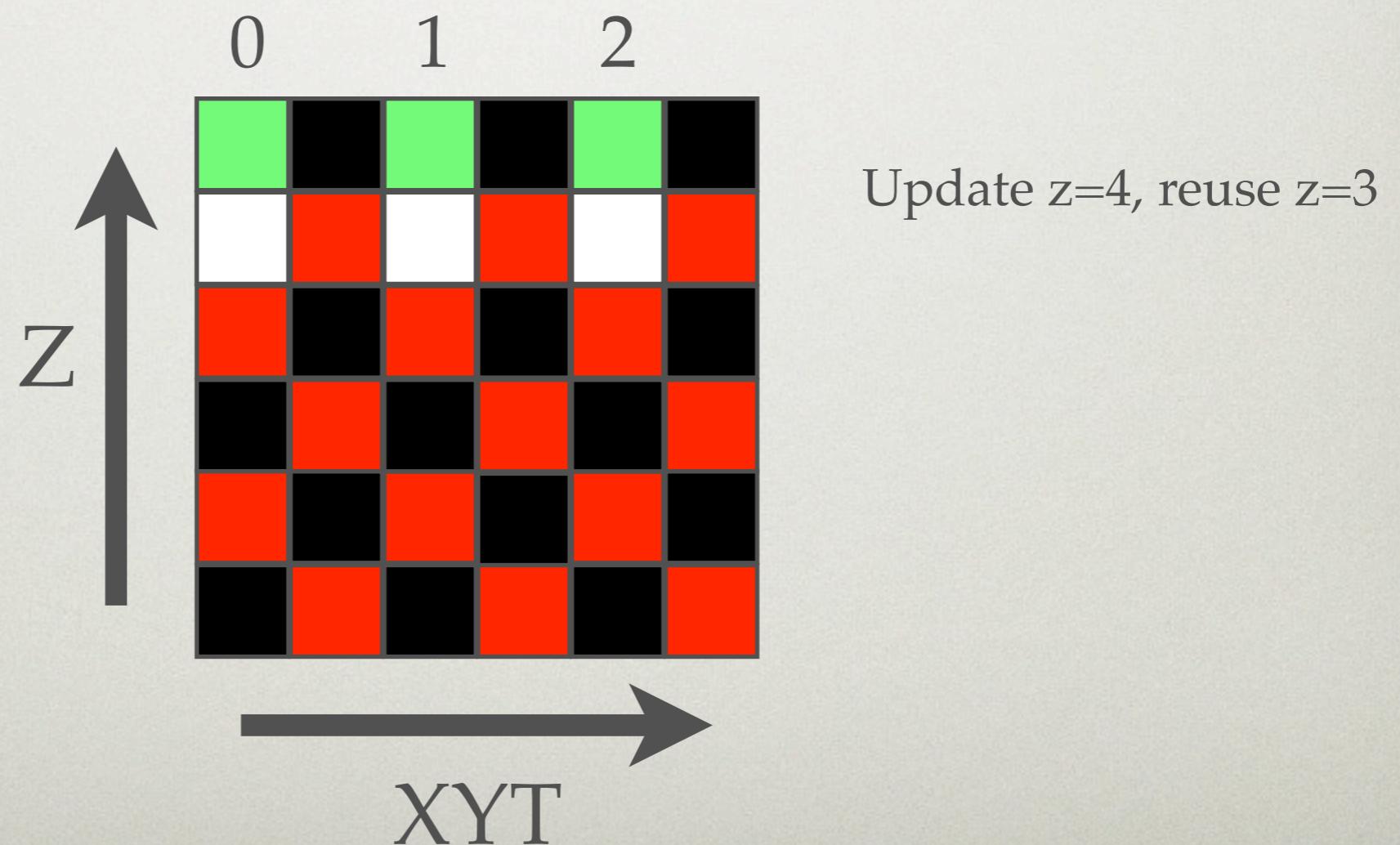
Z STREAMING

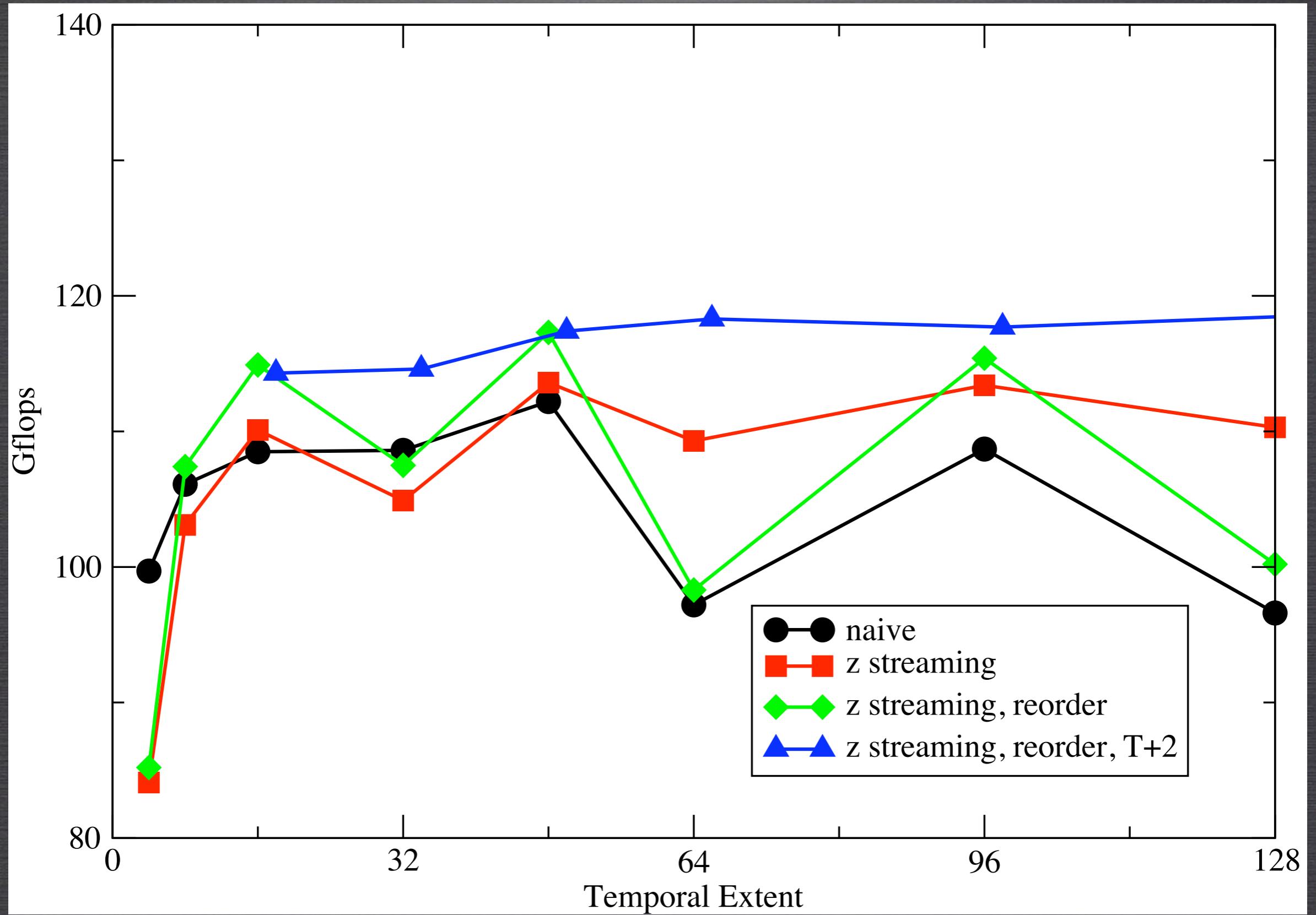
- Easy to reuse loads in one dimension
- Use XYT threads, each thread streams down Z dimension



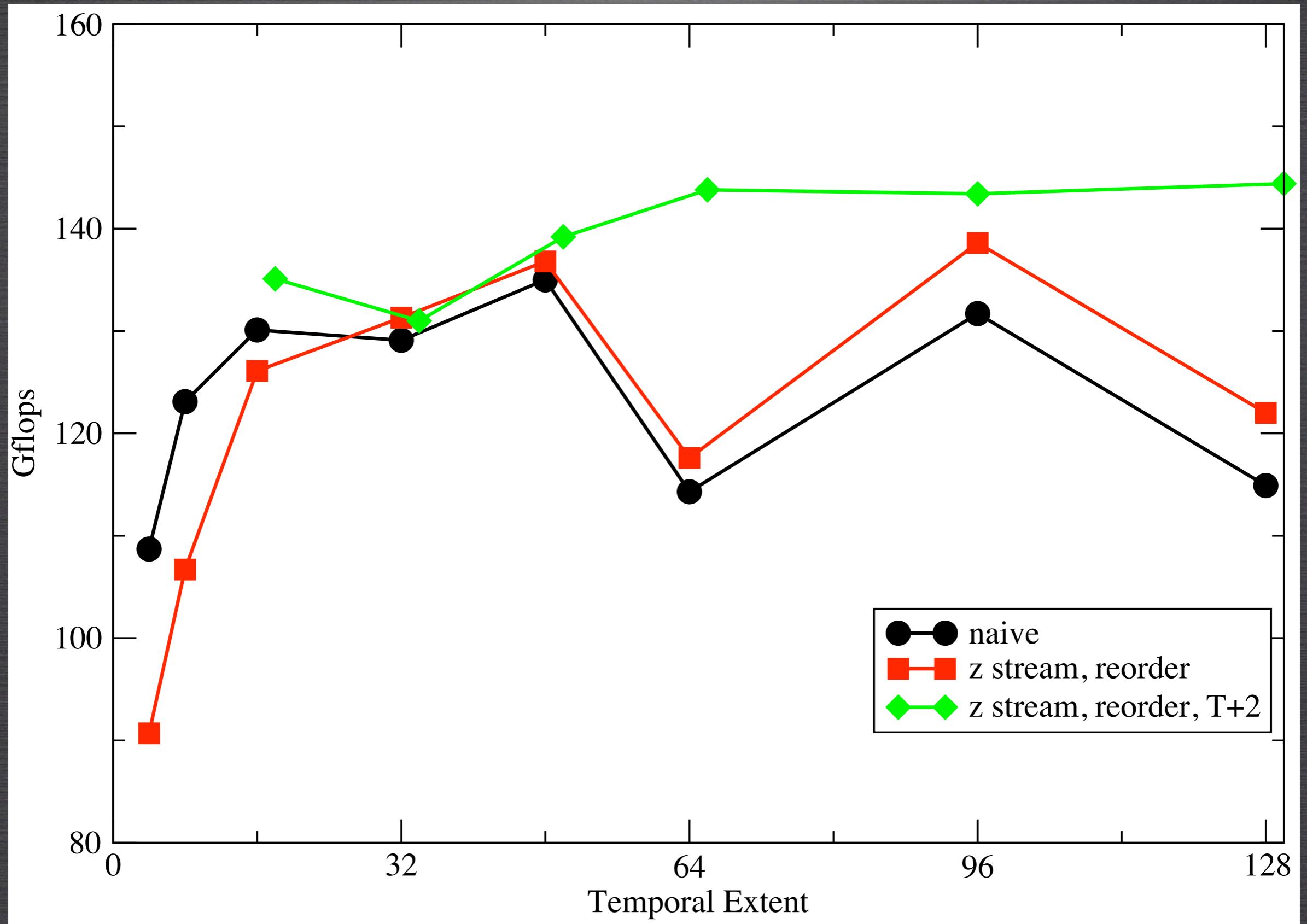
Z STREAMING

- Easy to reuse loads in one dimension
- Use XYT threads, each thread streams down Z dimension





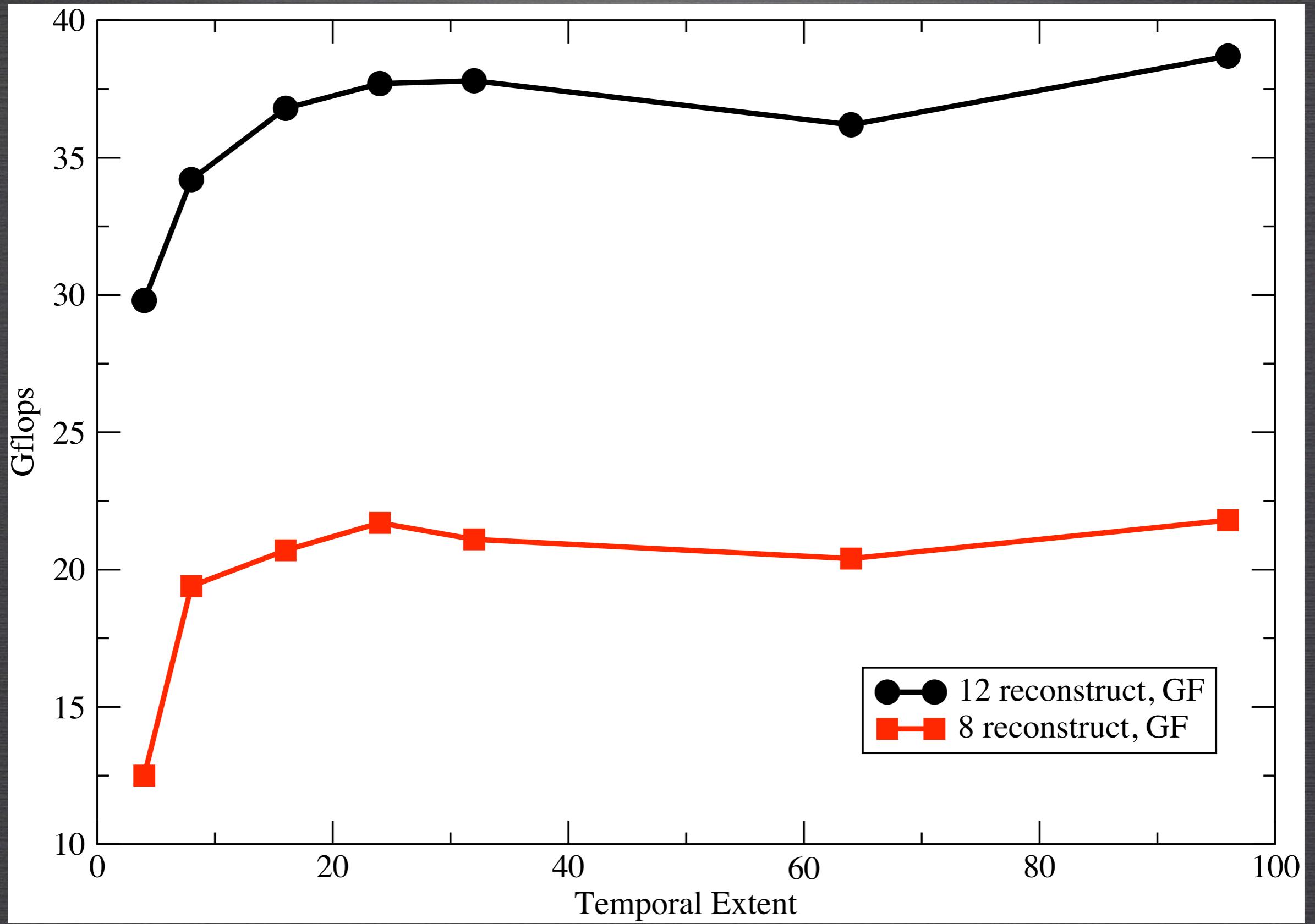
WILSON MATRIX-VECTOR PERFORMANCE
SINGLE PRECISION Z STREAMING(12 RECONSTRUCT, $V=24^3 \times T$)



WILSON MATRIX-VECTOR PERFORMANCE
SINGLE PRECISION Z STREAMING (8 RECONSTRUCT GF, V=24³XT)

DOUBLE PRECISION

- Double precision peak ~ 78 Gflops
 - Flop / Bandwidth ratio much more forgiving
 - Not all tricks are useful anymore...
- Half the amount of shared memory and registers
 - Simpler is better
 - Performance penalty only a factor ~3 vs. single



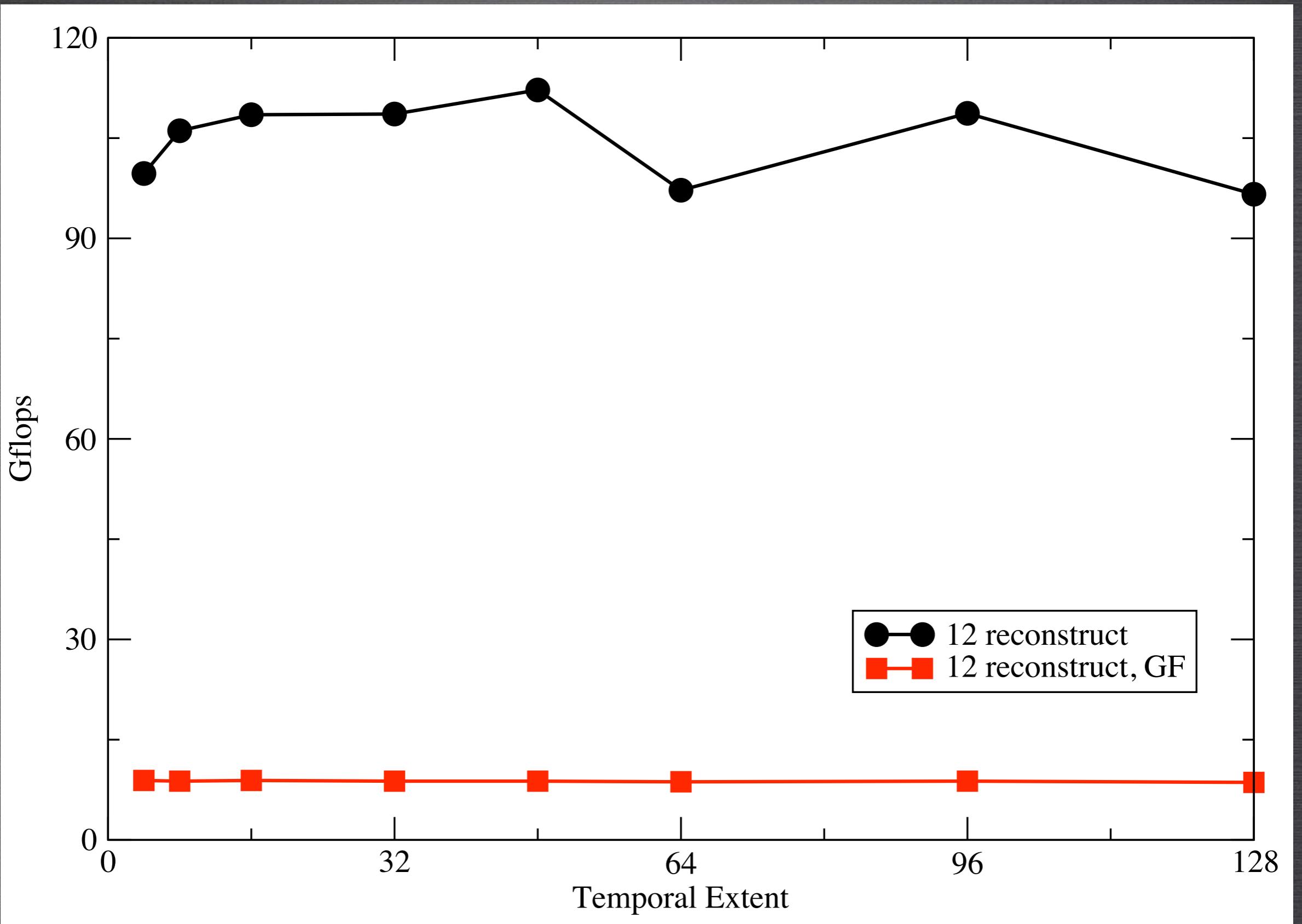
WILSON MATRIX-VECTOR PERFORMANCE

DOUBLE PRECISION ($V = 24^3 \times T$)

KRYLOV SOLVER IMPLEMENTATION

- Use GPU as an accelerator
- Replace matrix-vector in solver with GPU matrix-vector?
 - Transfer \mathbf{p}_{k+1} to GPU
 - Perform computation
 - Transfer $A\mathbf{p}_{k+1}$ back to CPU
- Problem
 - Transfers and matrix-vector are $O(N)$ operations

```
while (|rk| > ε) {  
    βk = (rk, rk) / (rk-1, rk-1)  
    pk+1 = rk - βkpk  
  
    α = (rk, rk) / (pk+1, A pk+1)  
    rk+1 = rk - α A pk+1  
    xk+1 = xk + α pk+1  
    k = k+1  
}
```

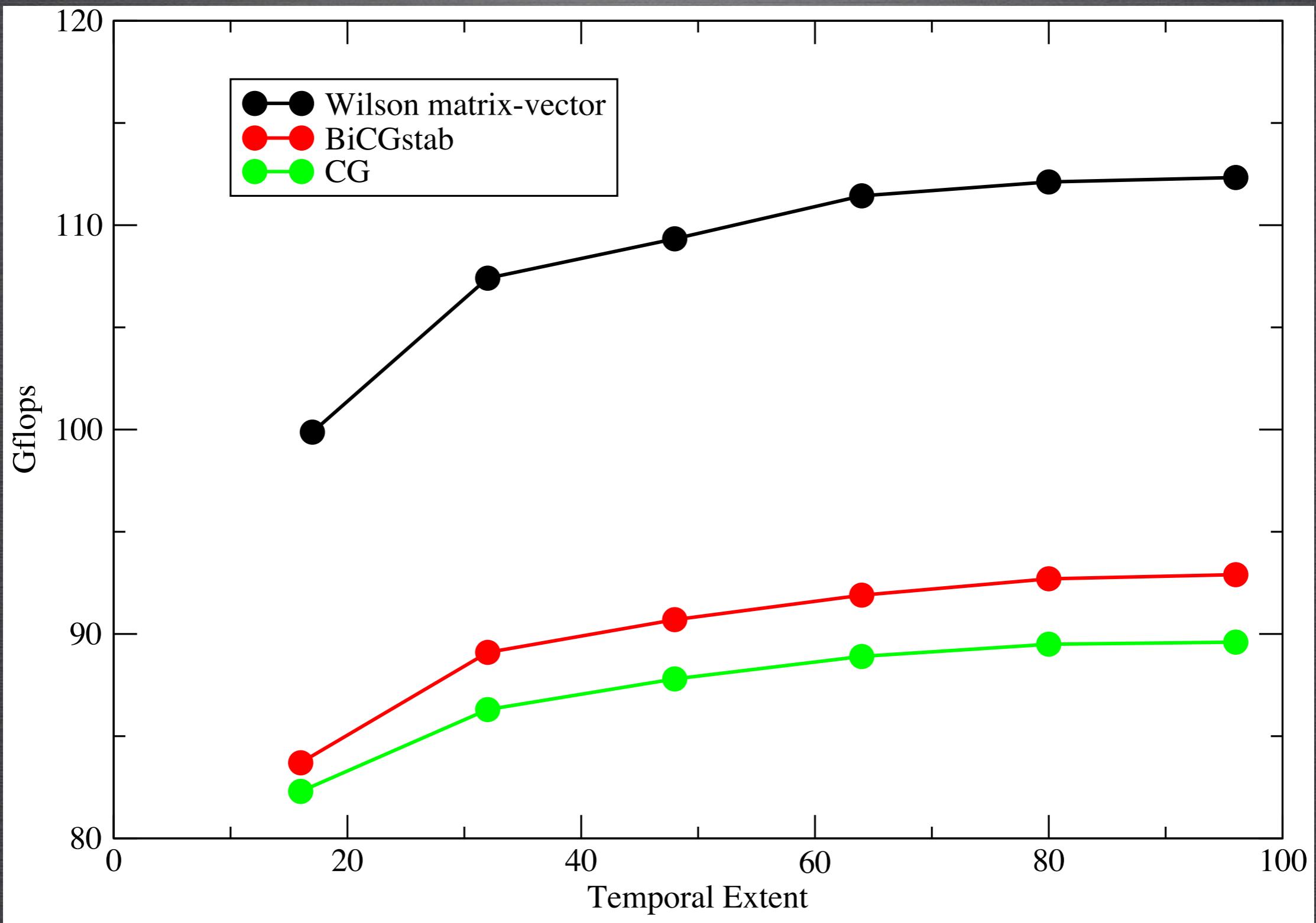


WILSON MATRIX-VECTOR PERFORMANCE

SINGLE PRECISION ACCELERATOR ($V = 24^3 \times T$)

KRYLOV SOLVER IMPLEMENTATION

- Complete solver must be on GPU
 - Transfer \mathbf{b} to GPU
 - Solve $\mathbf{Ax} = \mathbf{b}$
 - Transfer \mathbf{x} to CPU
- Besides matrix-vector, require BLAS level 1 type operations
 - AXPY operations: $\mathbf{b} += \mathbf{ax}$
 - NORM operations: $c = (\mathbf{b}, \mathbf{b})$
- CUBLAS library available
- Better to coalesce operations to minimize bandwidth
 - e.g., AXPY_NORM



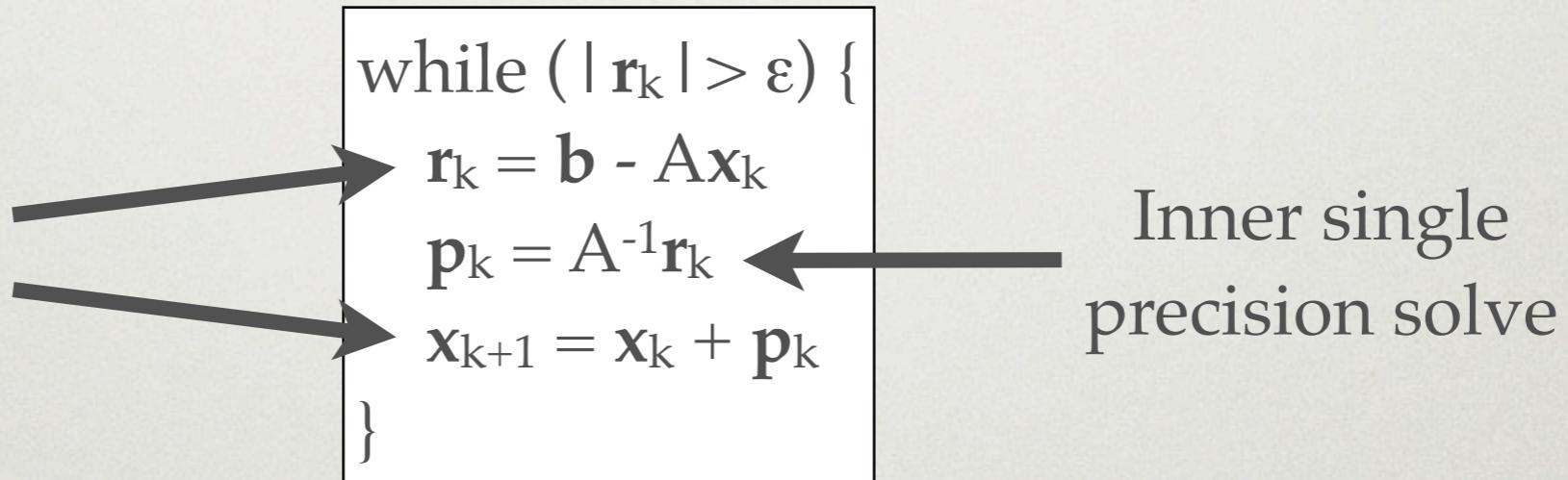
WILSON INVERTER PERFORMANCE

SINGLE PRECISION (12 RECONSTRUCT, V=24³XT)

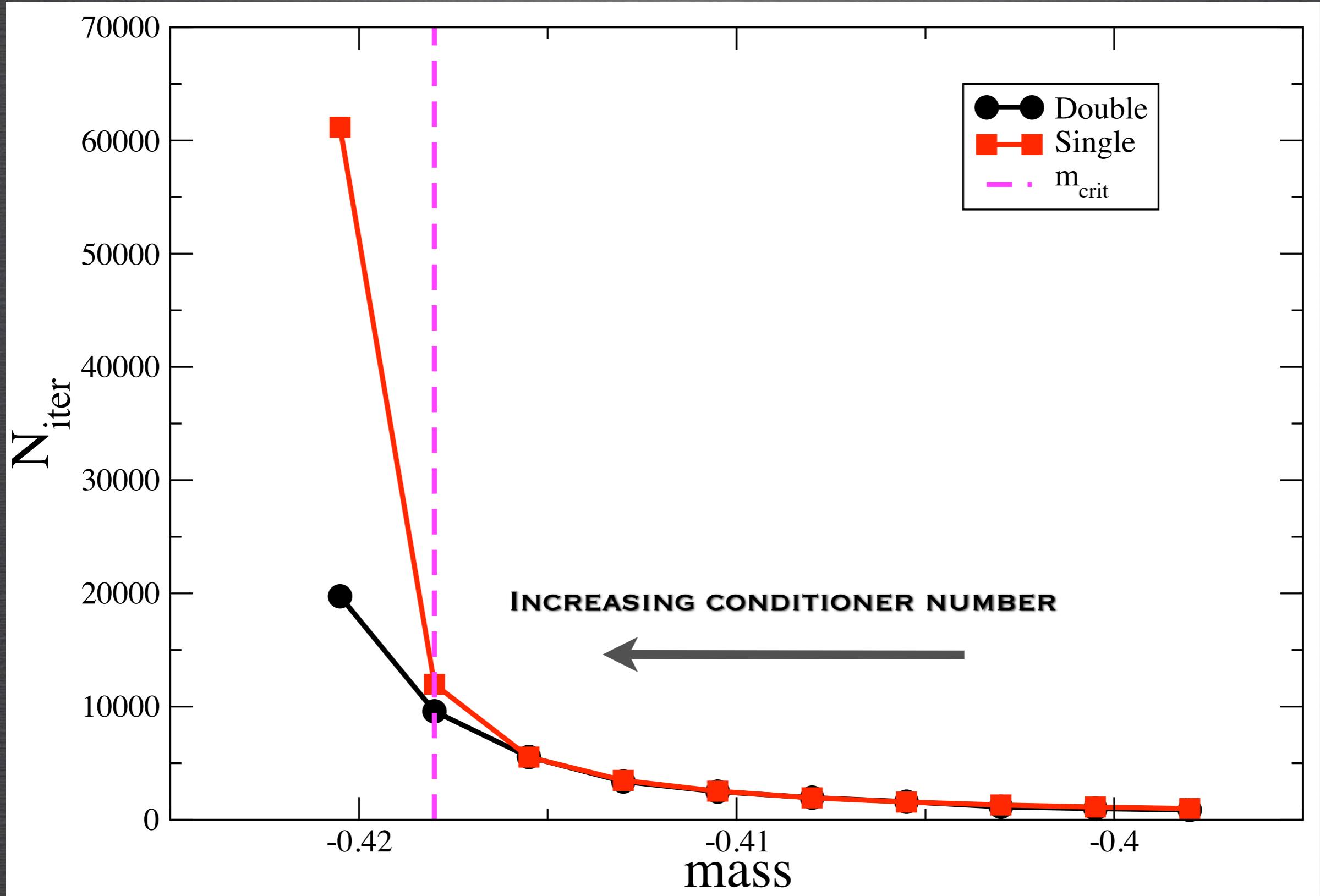
MULTI-PRECISION SOLVERS

- Require solver tolerance beyond limit of single precision
- Use **defect-correction**

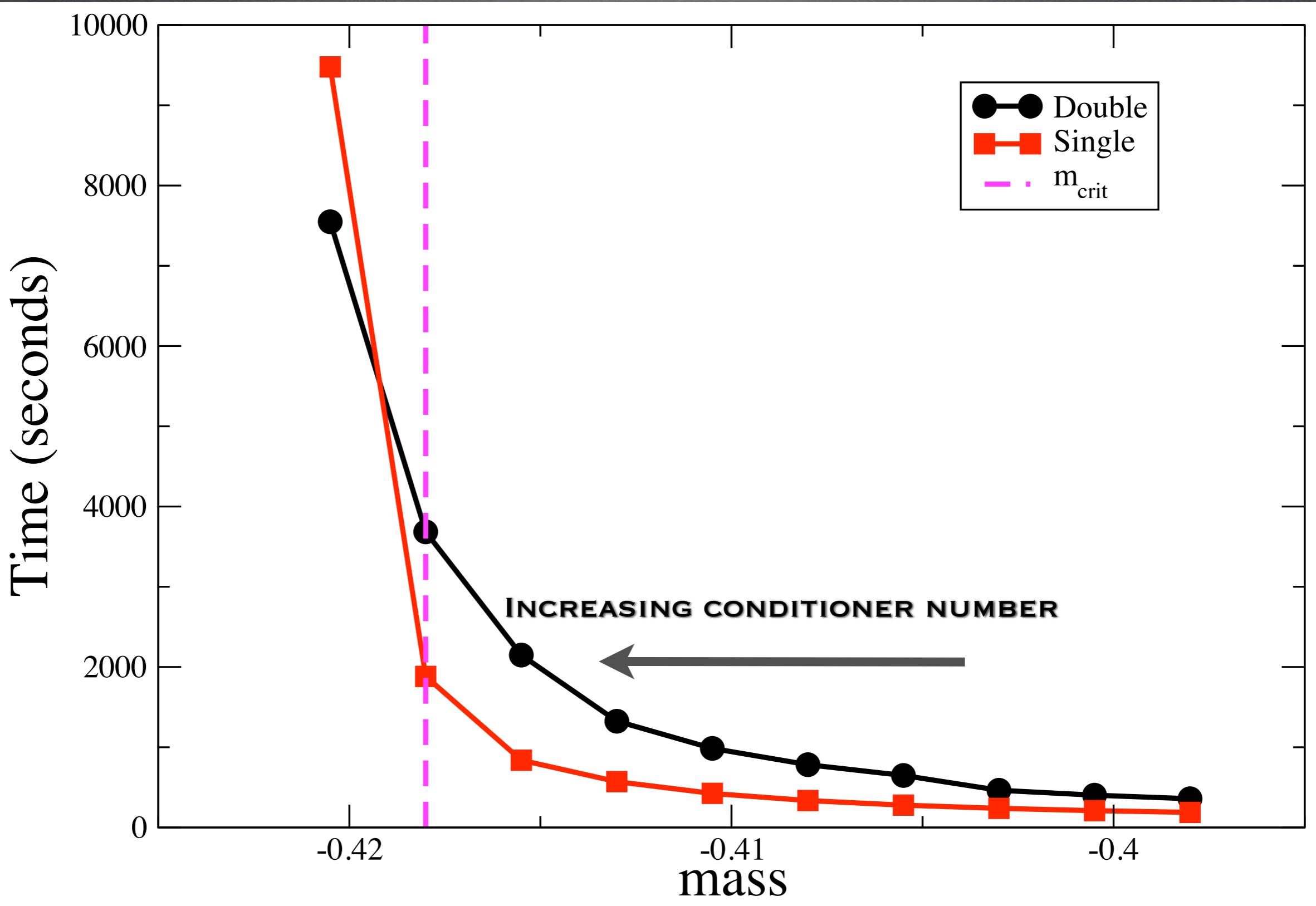
Double precision
mat-vec and
accumulate



- Double precision done can be done on CPU or GPU
 - Can always check GPU gets correct answer
 - Disadvantage is each new single precision solve is a restart



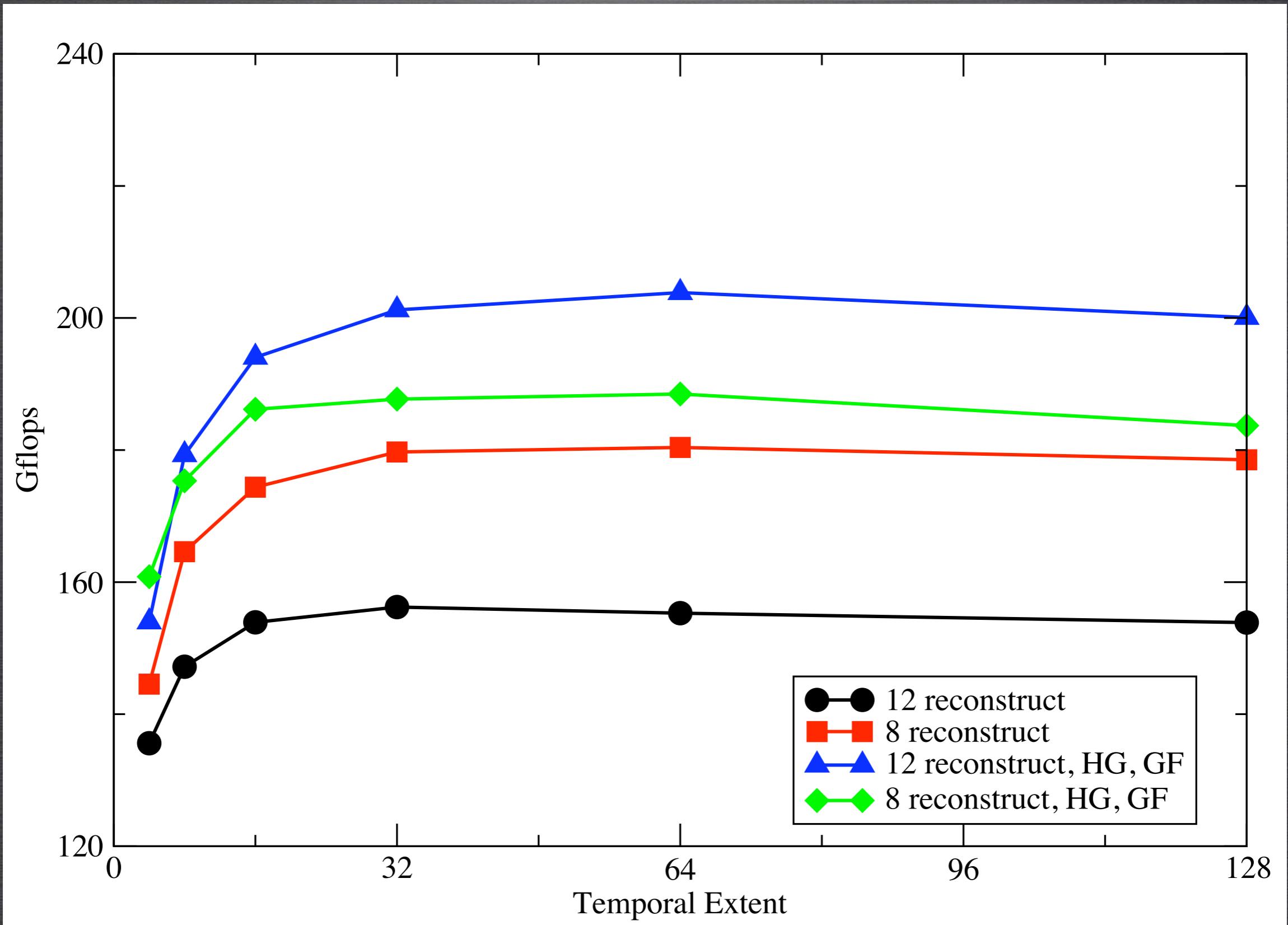
WILSON INVERTER ITERATIONS
($\varepsilon=10^{-8}$, $V=32^3 \times 96$)



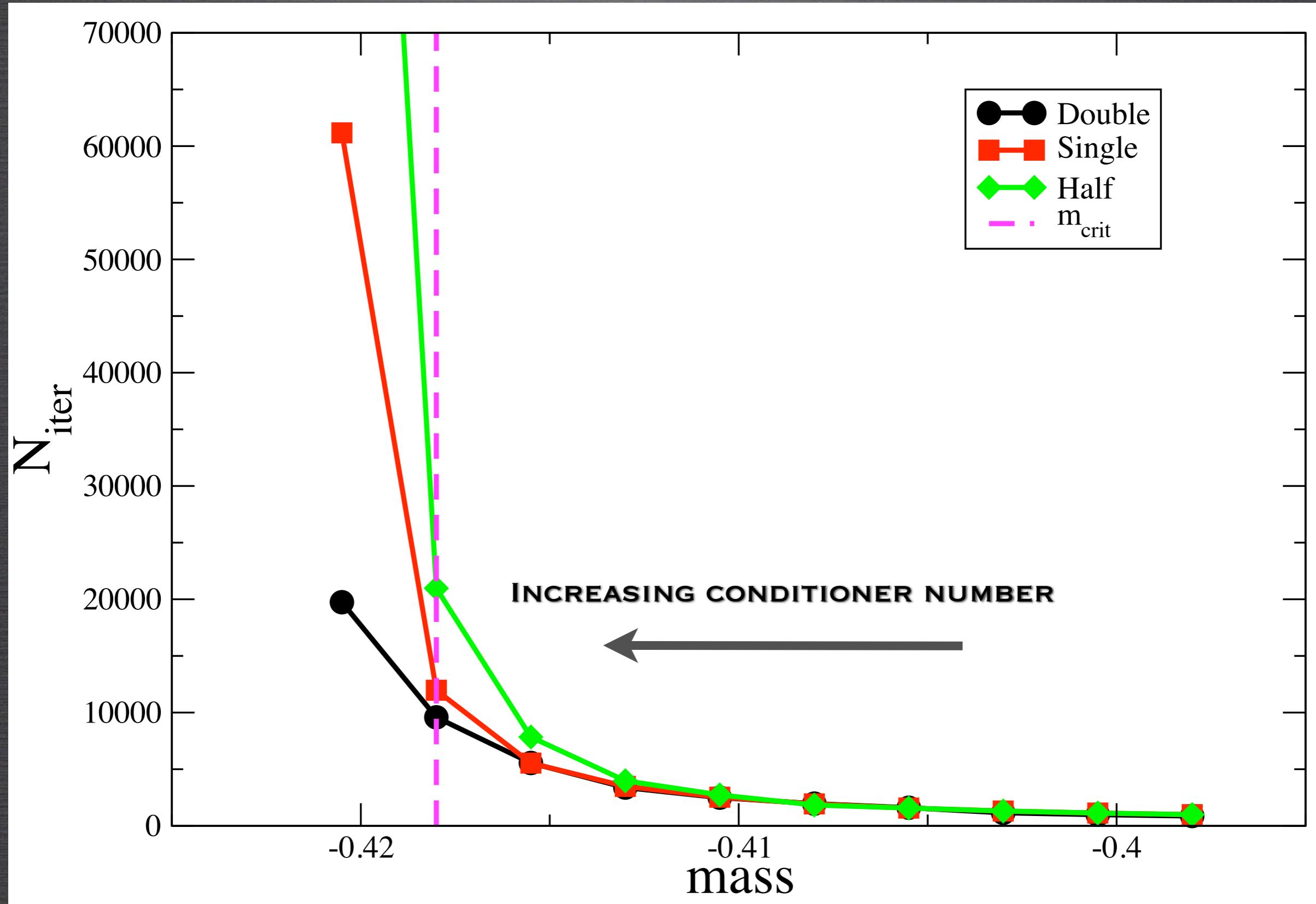
WILSON INVERTER TIME TO SOLUTION
($\varepsilon=10^{-8}$, $V=32^3 \times 96$)

HALF PRECISION

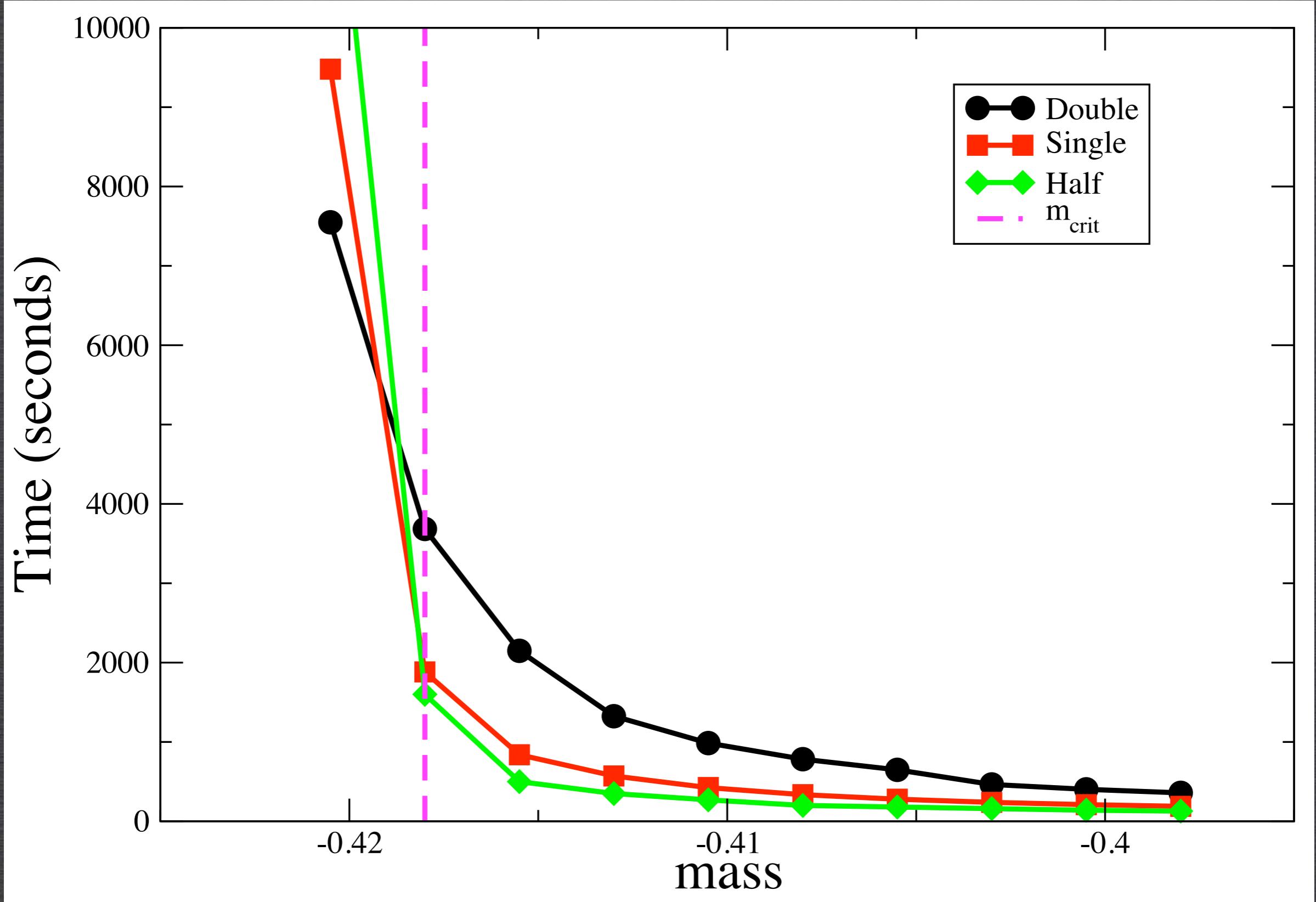
- Single-precision can be used to find double-precision result
- GPU kernel is still bandwidth bound
 - Use half precision for inner solve?
- ~~No support for FP16 in CUDA~~ (new in CUDA 2.3)
- Does support texture reading 16-bit integer -> float [-1,1]
 - Internal calculation done in 32-bit
- Natural format for SU(3) fields
- Vector field has no exponent restrictions
- Dirac operator off-diagonals are O(1) mixing of site-wise components
 - Scale all components at a site by maximum element
 - 24 floats per site -> 24 half-integers + 1 float



WILSON MATRIX-VECTOR PERFORMANCE HALF PRECISION



WILSON INVERTER ITERATIONS
($\varepsilon = 10^{-8}$, $V = 32^3 \times 96$)

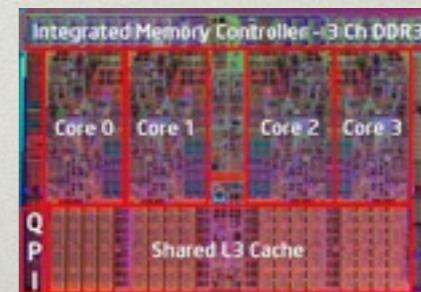
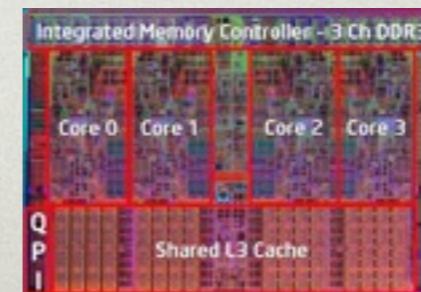
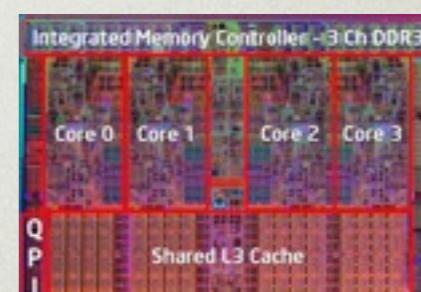
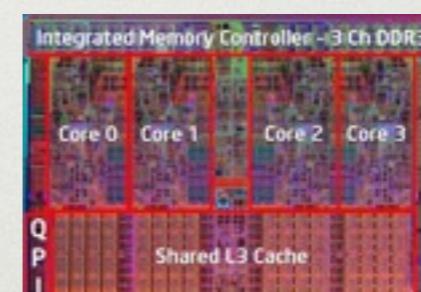
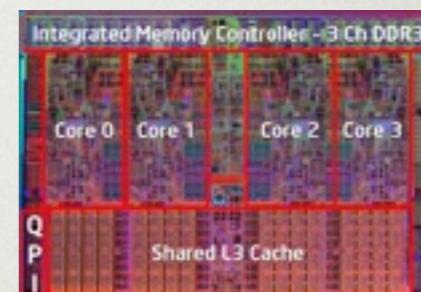
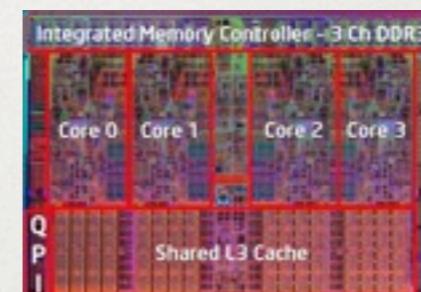
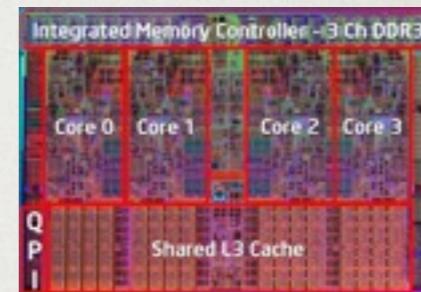
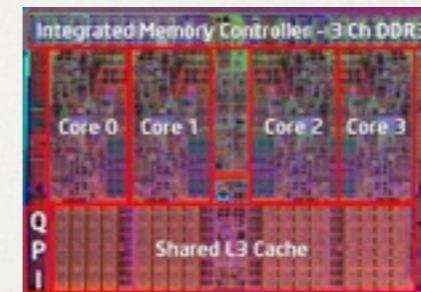


WILSON INVERTER TIME TO SOLUTION
($\varepsilon=10^{-8}$, $V=32^3 \times 96$)

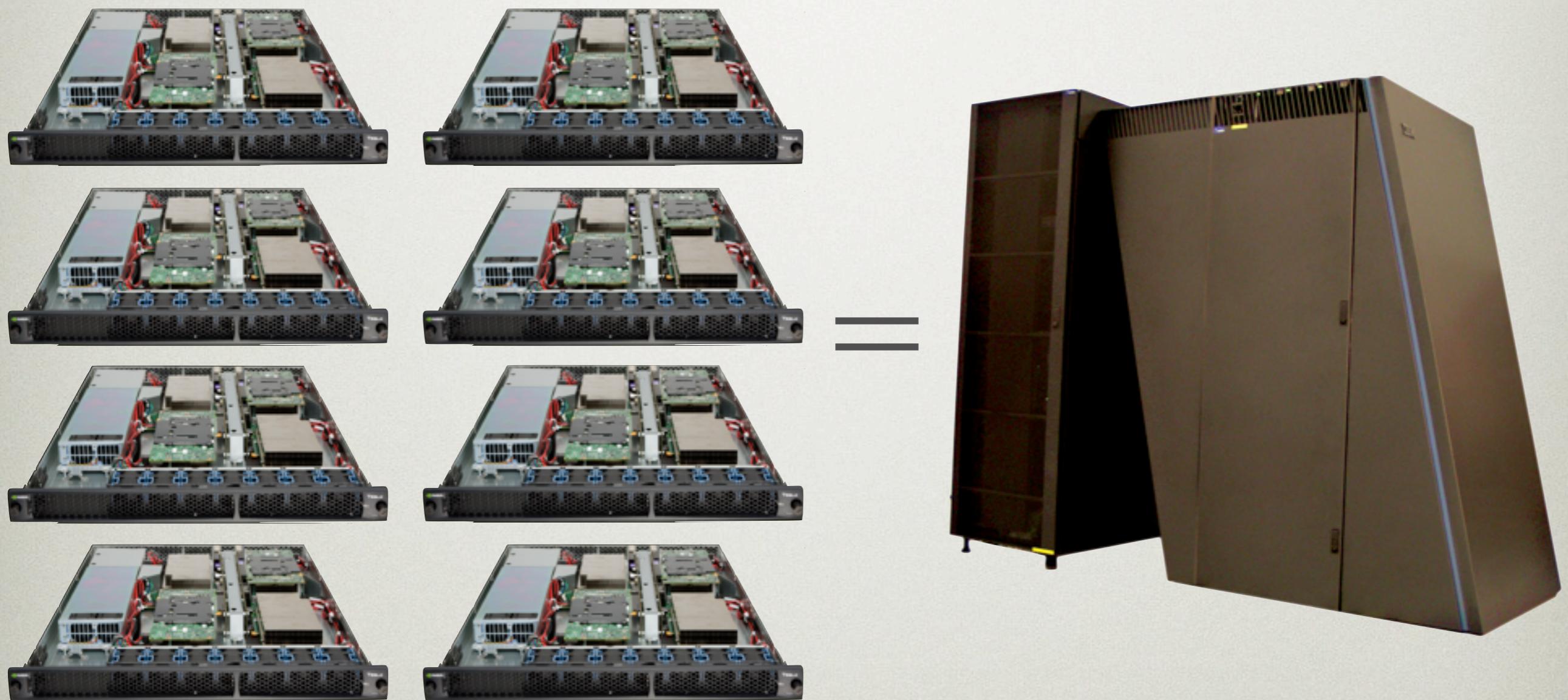
PERFORMANCE PER MFLOP



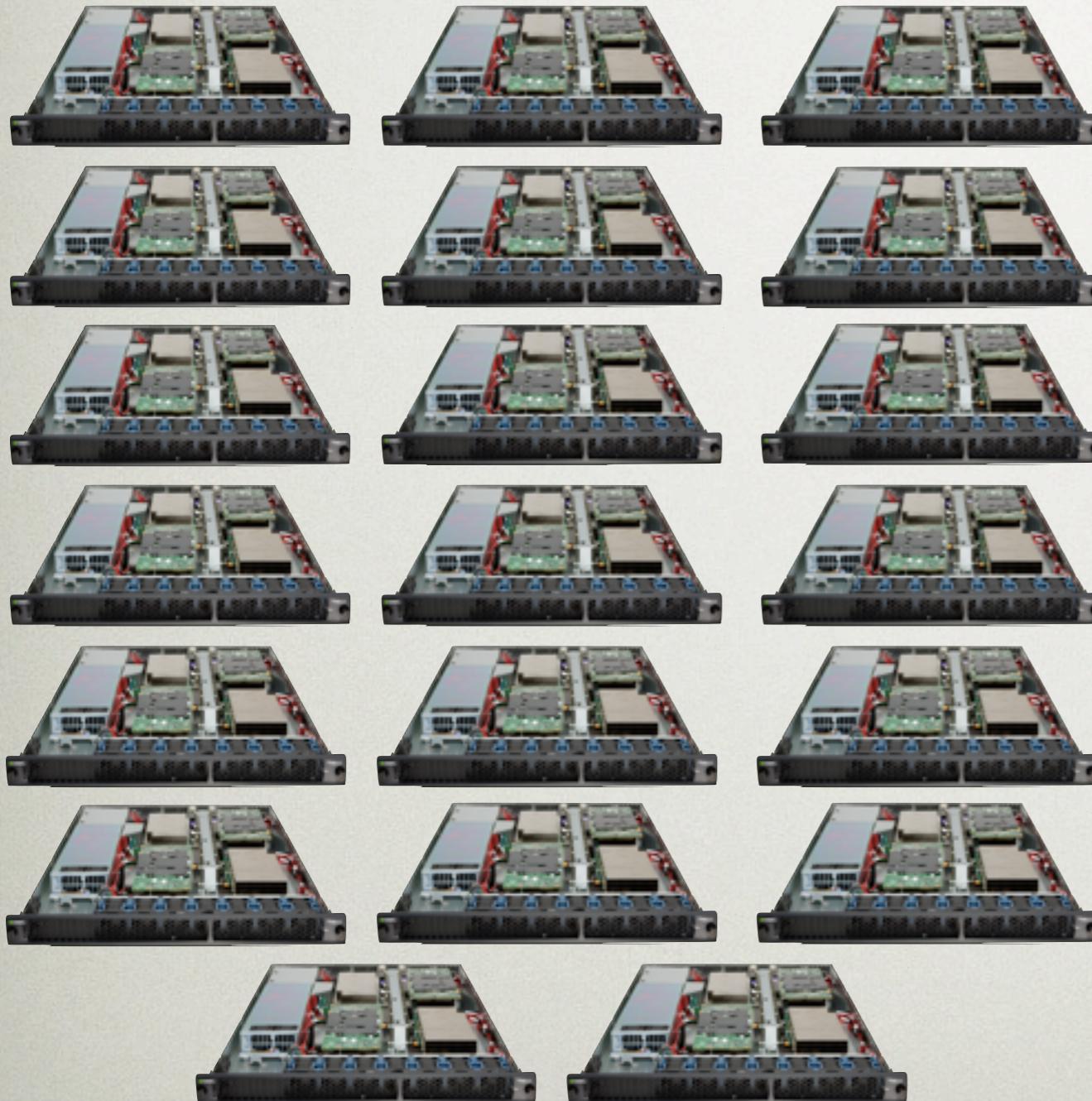
=



PERFORMANCE PER MFLOP



PERFORMANCE PER WATT



=



PERFORMANCE PER \$



=

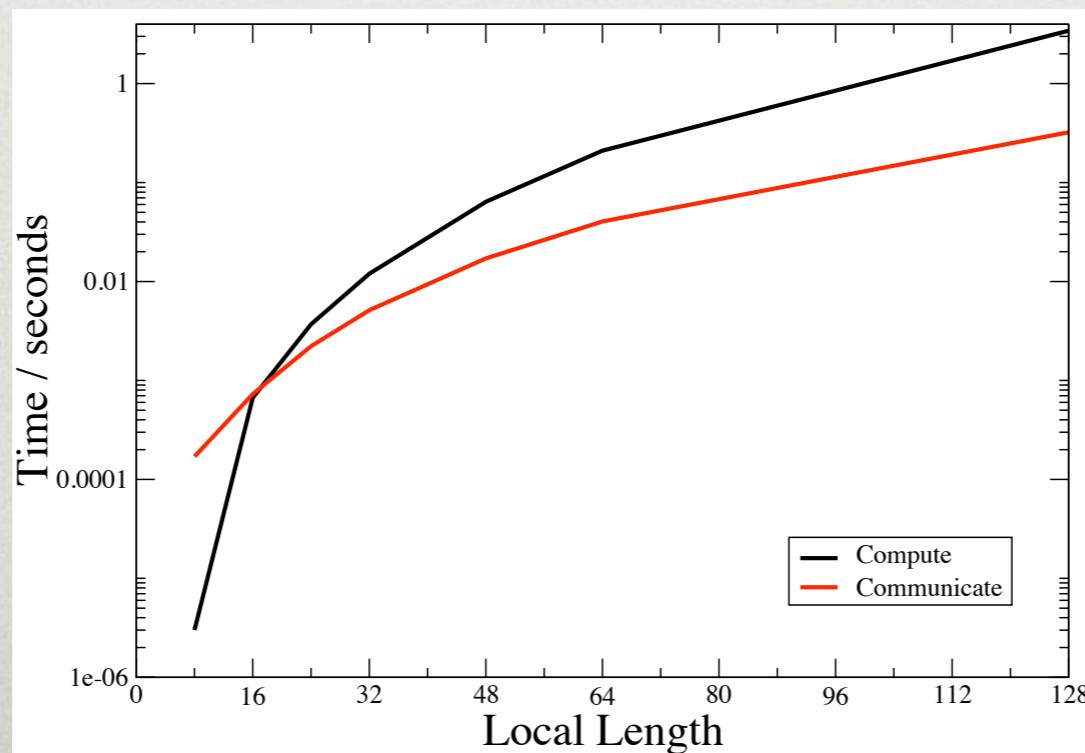


MULTI-GPU

- Want to implement multi-GPU kernels
 - Overlap computation and communication
 - Limited by latency and inter-GPU bandwidth
- Multi-GPUs within a single PC or clusters with GPU in each?
- No-one's got hard numbers yet
- Model **conservative** performance scaling
 - Wilson dslash kernel @ 100 Gflops
 - PCIe bandwidth 5GBs^{-1} , latency = $5\ \mu\text{s}$
 - 8 way half-spinor communication

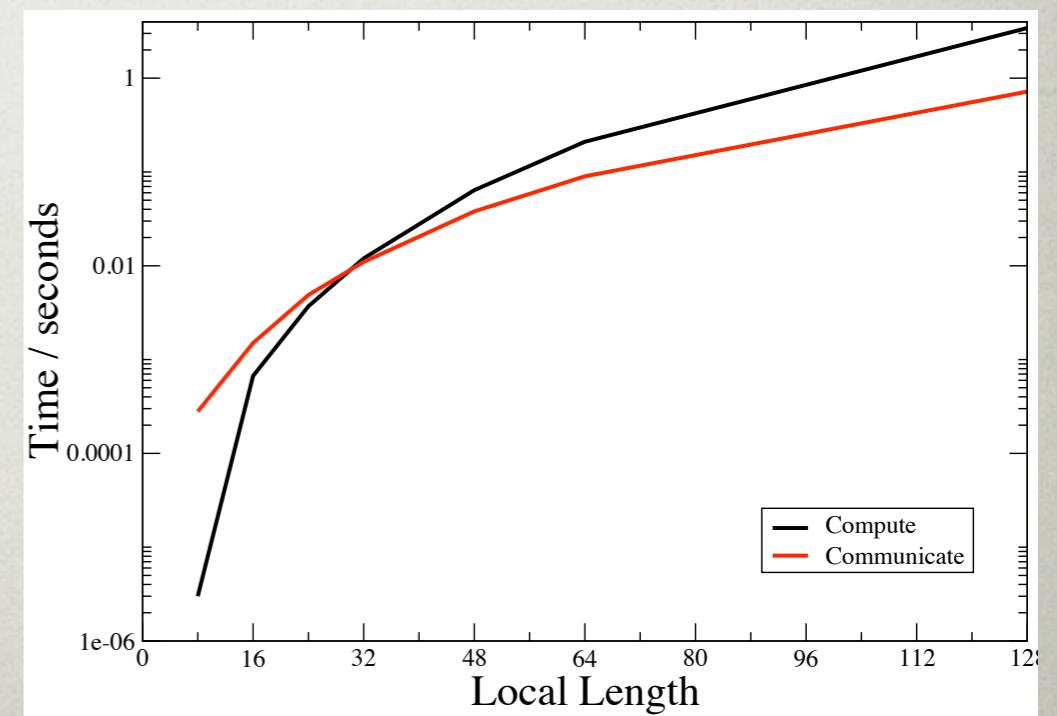
POSSIBLE MULTI-GPU CONFIGURATIONS

- Multi-GPU within a box
 - Cheap and easy to set up
 - Should provide excellent scaling
 - $O(1)$ Tflop under your desk!



POSSIBLE MULTI-GPU CONFIGURATIONS

- True GPU cluster
 - Many nodes each with GPUs
 - E.g., NCSA
 - 2x dual-core 2.4 GHz AMD Opterons, 8 GB of memory
 - 1xTesla S1070
 - 2GB/sec Infiniband connection
 - Good scaling very challenging
 - Need large local volume
 - Better algorithms (DD solver?)
 - Custom hardware?



ISN'T IT A LOT OF WORK TO WRITE **FAST** GPU CODE?

- Many optimization parameters to consider
 - Loop unrolling
 - Depth of parallelism
 - Order of execution
- Also algorithm/physics parameters
 - Matrix transpose, SU(3) parameterization, gamma matrix basis
- Remove grunt work by writing code to write code
 - Python code generator creates CUDA code
 - Optimization parameters are just that, parameters
- Quickly find maximum in performance space





CONCLUSIONS

- Fantastic algorithmic performance obtained on today GPUs
 - Flops per Watt, Flops per \$ -> new metric Time to solution per \$
 - Some work required to get best performance
 - Algorithm design critical component
 - Bandwidth constraints force rethink of problem
 - New algorithms
 - Many groups harnessing the power of GPUs
 - Next step is multi-GPU
 - Case yet to be proven
 - **Prediction:** next year lots to report

QUDA LIBRARY

- QCD-GPU library
 - Wilson and Wilson-clover fermions
 - Inverter interface (SciDac level 3)
 - Dslash interface (accelerator mode)
- Hooks into QDP / CPS++ etc.
- Future...
 - Domain wall, twisted mass
 - Fermion and Gauge forces
 - Multi-GPU
 - Adaptive multigrid