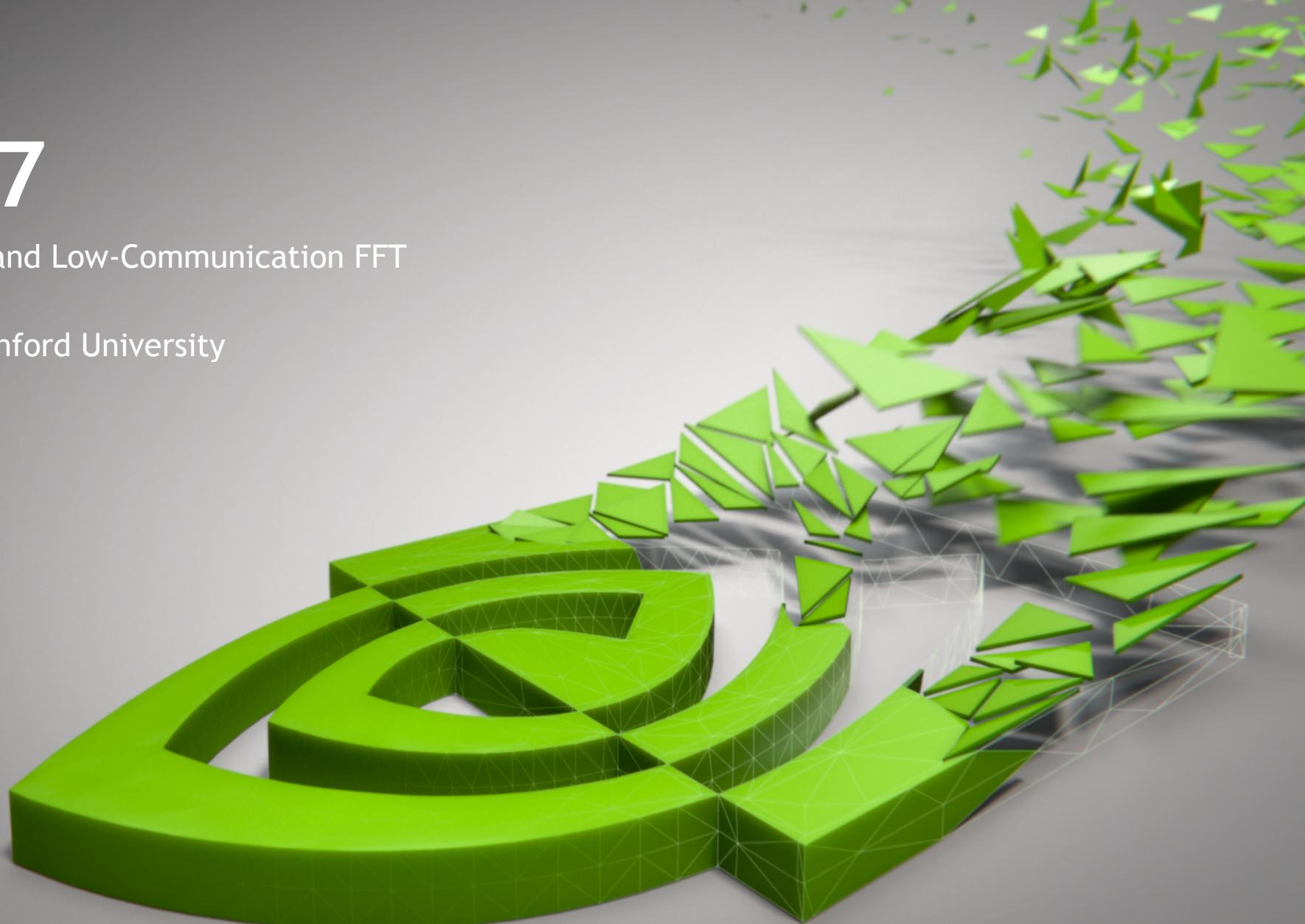


PPAM 17

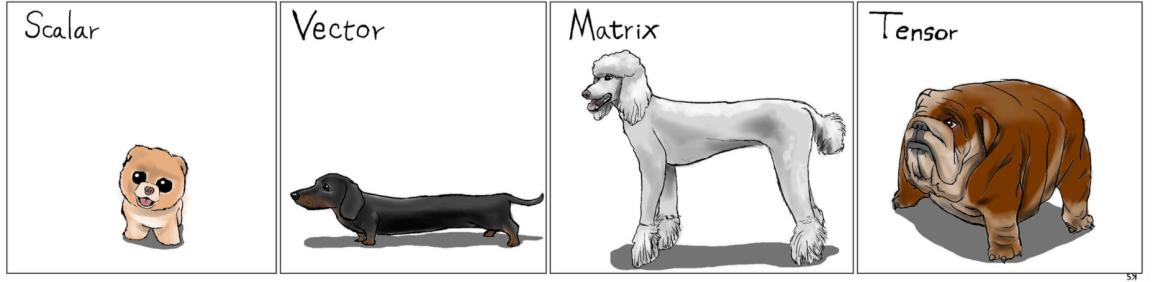
Tensor Contractions and Low-Communication FFT

Cris Cecka

NVIDIA Research, Stanford University



TENSOR CONTRACTIONS



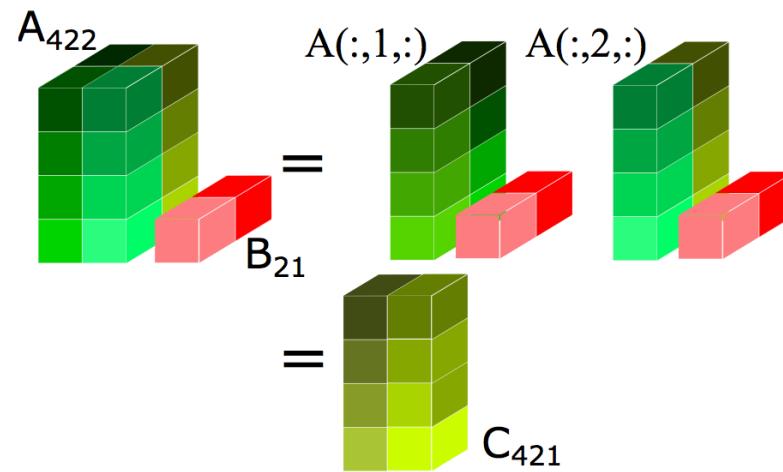
Modern data is inherently multi-dimensional.

NVIDIA especially interested due to ML.

Portability, BLAS spec, efficiency.

TENSOR CONTRACTIONS

$$C_{\mathcal{C}} = A_{\mathcal{A}} B_{\mathcal{B}}$$



$$\text{e.g. } C_{mnp} = A_{mnk} B_{kp}$$

- Core primitive of multilinear algebra
- BLAS level 3 – unbounded compute intensity.

What do we have?

Tensor computation libraries

- ① Arbitrary/restricted tensor operation of any order and dimension
 - ① Tensortoolbox (Matlab)
 - ② FTensor (C++)
 - ③ Cyclops (C++)
 - ④ BTAS (C++)
 - ⑤ All the Python...

Efficient computing frame

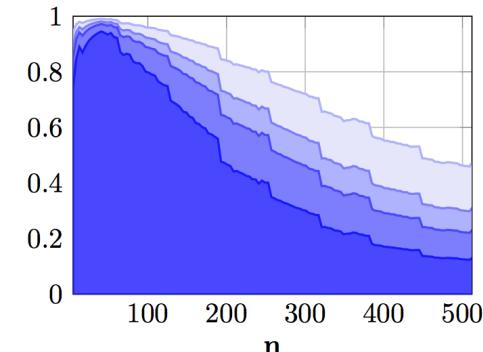
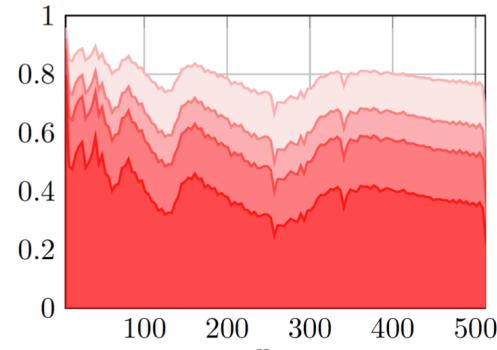
- ① Static analysis solutions
 - ① PPCG [ISL] (polyhedral)
 - ② TCE (DSL)
- ② Parallel and distributed primitives
 - ① BLAS, cuBLAS
 - ② BLIS, BLASX, cuBLASXT

TENSOR LIBRARIES

Explicit permutation dominates.

Consider $C_{mnp} = A_{km} B_{pkn}$.

- ① $A_{km} \rightarrow A_{mk}$
- ② $B_{pkn} \rightarrow B_{kpn}$
- ③ $C_{mnp} \rightarrow C_{mpn}$
- ④
$$C_{m(pn)} = A_{mk} B_{k(pn)}$$
- ⑤ $C_{mpn} \rightarrow C_{mnp}$



(Top) CPU. (Bottom) GPU. The fraction of time spent in copies/transpositions. Lines are shown with 1, 2, 3, and 6 transpositions.

EXISTING PRIMITIVES

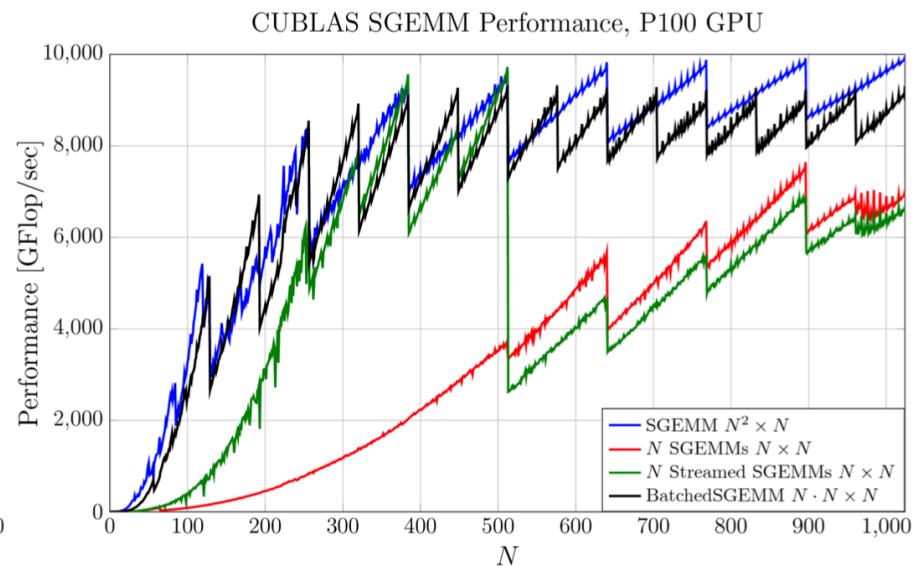
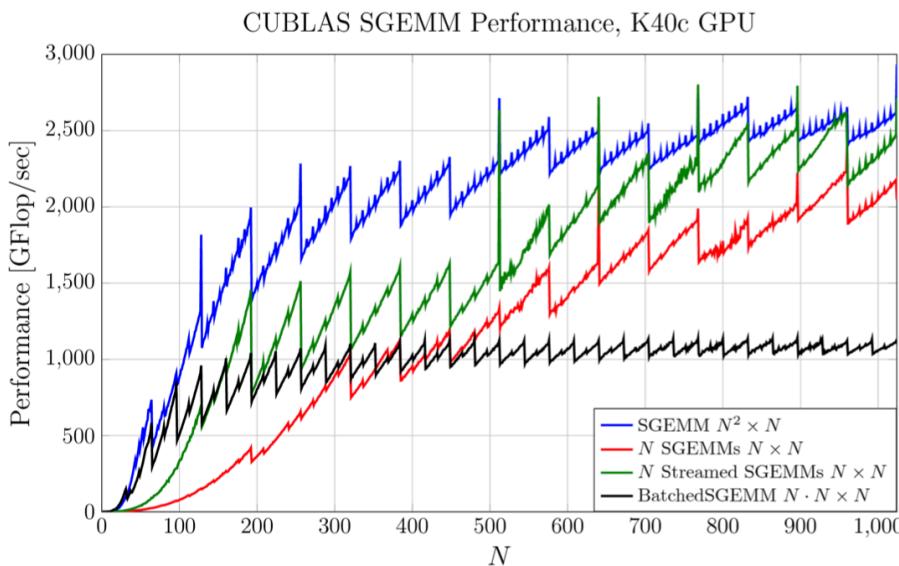
- GEMM
 - Suboptimal for many small matrices
- Pointer-to-pointer BatchedGEMM
 - Available since MKL 13.1b and cuBLAS 4.0

$$C[p] = \alpha \operatorname{op}(A[p]) \operatorname{op}(B[p]) + \beta C[p]$$

```
cublas<T>gemmBatched(cublasHandle_t handle,
                        cublasOperation_t transA, cublasOperation_t transB,
                        int M, int N, int K,
                        const T* alpha,
                        const T** A, int ldA,
                        const T** B, int ldB,
                        const T* beta,
                        T** C, int ldC,
                        int batchCount)
```

EXISTING PRIMITIVES

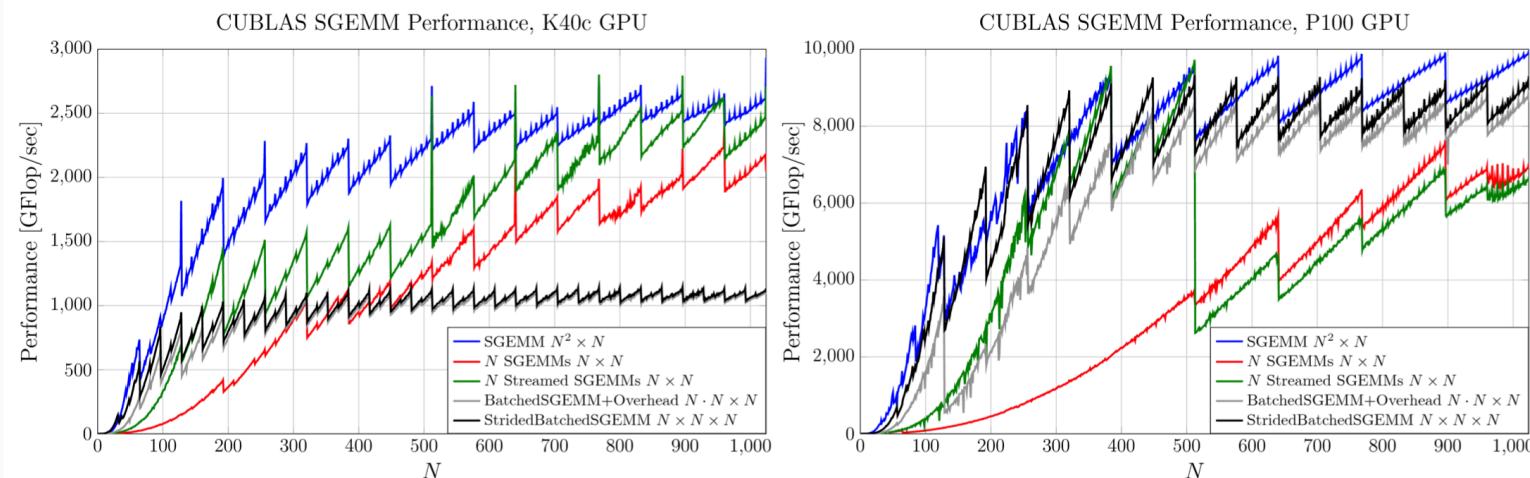
Pointer-to-Pointer BatchedGEMM



EXISTING PRIMITIVES

Pointer-to-Pointer BatchedGEMM

Except actually...



Solution: StridedBatchedGEMM

STRIDED BATCHED GEMM

```
cublas<T>gemmStridedBatched(cublasHandle_t handle,
                               cublasOperation_t transA, cublasOperation_t transB,
                               int M, int N, int K,
                               const T* alpha,
                               const T* A, int ldA1, int strideA,
                               const T* B, int ldb1, int strideB,
                               const T* beta,
                               T* C, int ldc1, int strideC,
                               int batchCount)
```

- Common use case for pointer-to-pointer batched GEMM
- No pointer-to-pointer data structure or overhead
- Performance on par with pure GEMM (P100 and beyond)

CONTRACTIONS

$$C_{mnp} = A_{**} \times B_{***}$$

$$C_{mnp} \equiv C[m + n \cdot \text{IdC1} + p \cdot \text{IdC2}]$$

 : Single GEMM
(Provided compact layout)

Case	Contraction	Kernel1	Kernel2	Case	Contraction	Kernel1	Kernel2
1.1	$A_{mk}B_{knp}$	$C_{m(np)} = A_{mk}B_{k(np)}$	$C_{mn[p]} = A_{mk}B_{kn[p]}$	4.1	$A_{kn}B_{kmp}$	$C_{mn[p]} = B_{km[p]}^\top A_{kn}$	
1.2	$A_{mk}B_{kp\bar{n}}$	$C_{mn[p]} = A_{mk}B_{k[p]n}$	$C_{m[n]p} = A_{mk}B_{kp[n]}$	4.2	$A_{kn}B_{kpm}$	$C_{mn[p]} = B_{k[p]m}^\top A_{kn}$	
1.3	$A_{mk}B_{nk\bar{p}}$	$C_{mn[p]} = A_{mk}B_{nk[p]}^\top$		4.3	$A_{kn}B_{mkp}$	$C_{mn[p]} = B_{mk[p]}A_{kn}$	
1.4	$A_{mk}B_{pk\bar{n}}$	$C_{m[n]p} = A_{mk}B_{pk[n]}^\top$		4.4	$A_{kn}B_{pk\bar{m}}$		
1.5	$A_{mk}B_{n\bar{p}k}$	$C_{m(np)} = A_{mk}B_{(np)k}^\top$	$C_{mn[p]} = A_{mk}B_{n[p]k}^\top$	4.5	$A_{kn}B_{mpk}$	$C_{mn[p]} = B_{m[p]k}A_{kn}$	
1.6	$A_{mk}B_{pn\bar{k}}$	$C_{m[n]p} = A_{mk}B_{p[n]k}^\top$		4.6	$A_{kn}B_{pm\bar{k}}$		
2.1	$A_{km}B_{knp}$	$C_{m(np)} = A_{km}^\top B_{k(np)}$	$C_{mn[p]} = A_{km}^\top B_{kn[p]}$	5.1	$A_{pk}B_{kmn}$	$C_{(mn)p} = B_{k(mn)}^\top A_{pk}^\top$	$C_{m[n]p} = B_{km[n]}^\top A_{pk}^\top$
2.2	$A_{km}B_{kp\bar{n}}$	$C_{mn[p]} = A_{km}^\top B_{k[p]n}$	$C_{m[n]p} = A_{km}^\top B_{kp[n]}$	5.2	$A_{pk}B_{knm}$	$C_{m[n]p} = B_{k[n]m}^\top A_{pk}^\top$	
2.3	$A_{km}B_{nk\bar{p}}$	$C_{mn[p]} = A_{km}^\top B_{nk[p]}^\top$		5.3	$A_{pk}B_{mkn}$	$C_{m[n]p} = B_{mk[n]}A_{pk}^\top$	
2.4	$A_{km}B_{pk\bar{n}}$	$C_{m[n]p} = A_{km}^\top B_{pk[n]}^\top$		5.4	$A_{pk}B_{nkm}$		
2.5	$A_{km}B_{n\bar{p}k}$	$C_{m(np)} = A_{km}^\top B_{(np)k}^\top$	$C_{mn[p]} = A_{km}^\top B_{n[p]k}^\top$	5.5	$A_{pk}B_{mnk}$	$C_{(mn)p} = B_{(mn)k}A_{pk}^\top$	$C_{m[n]p} = B_{m[n]k}A_{pk}^\top$
2.6	$A_{km}B_{pn\bar{k}}$	$C_{m[n]p} = A_{km}^\top B_{p[n]k}^\top$		5.6	$A_{pk}B_{nmk}$		
3.1	$A_{nk}B_{kmp}$	$C_{mn[p]} = B_{km[p]}^\top A_{nk}^\top$		6.1	$A_{kp}B_{kmn}$	$C_{(mn)p} = B_{k(mn)}^\top A_{kp}^\top$	$C_{m[n]p} = B_{km[n]}^\top A_{kp}^\top$
3.2	$A_{nk}B_{kpm}$	$C_{mn[p]} = B_{k[p]m}^\top A_{nk}^\top$		6.2	$A_{kp}B_{knm}$	$C_{m[n]p} = B_{k[n]m}^\top A_{kp}^\top$	
3.3	$A_{nk}B_{mkp}$	$C_{mn[p]} = B_{mk[p]}A_{nk}^\top$		6.3	$A_{kp}B_{mkn}$	$C_{m[n]p} = B_{mk[n]}A_{kp}^\top$	
3.4	$A_{nk}B_{pk\bar{m}}$			6.4	$A_{kp}B_{nkm}$		
3.5	$A_{nk}B_{mpk}$	$C_{mn[p]} = B_{m[p]k}A_{nk}^\top$		6.5	$A_{kp}B_{mnk}$	$C_{(mn)p} = B_{(mn)k}A_{kp}^\top$	$C_{m[n]p} = B_{m[n]k}A_{kp}^\top$
3.6	$A_{nk}B_{pm\bar{k}}$			6.6	$A_{kp}B_{nmk}$		

CONTRACTIONS

$$C_{mnp} = A_{**} \times B_{***}$$

$$C_{mnp} \equiv C[m + n \cdot \text{ldC1} + p \cdot \text{ldC2}]$$

 : Single SB-GEMM
(Any layout)

Case	Contraction	Kernel1	Kernel2	Case	Contraction	Kernel1	Kernel2
1.1	$A_{mk}B_{knp}$	$C_{m(np)} = A_{mk}B_{k(np)}$	$C_{mn[p]} = A_{mk}B_{kn[p]}$	4.1	$A_{kn}B_{kmp}$	$C_{mn[p]} = B_{km[p]}^\top A_{kn}$	
1.2	$A_{mk}B_{kpn}$	$C_{mn[p]} = A_{mk}B_{k[p]n}$	$C_{m[n]p} = A_{mk}B_{kp[n]}$	4.2	$A_{kn}B_{kpm}$	$C_{mn[p]} = B_{k[p]m}^\top A_{kn}$	
1.3	$A_{mk}B_{nkp}$	$C_{mn[p]} = A_{mk}B_{nk[p]}^\top$		4.3	$A_{kn}B_{mkp}$	$C_{mn[p]} = B_{mk[p]}A_{kn}$	
1.4	$A_{mk}B_{pkn}$	$C_{m[n]p} = A_{mk}B_{pk[n]}^\top$		4.4	$A_{kn}B_{pkm}$		
1.5	$A_{mk}B_{npk}$	$C_{m(np)} = A_{mk}B_{(np)k}^\top$	$C_{mn[p]} = A_{mk}B_{n[p]k}^\top$	4.5	$A_{kn}B_{mpk}$	$C_{mn[p]} = B_{m[p]k}A_{kn}$	
1.6	$A_{mk}B_{pnk}$	$C_{m[n]p} = A_{mk}B_{p[n]k}^\top$		4.6	$A_{kn}B_{pmk}$		
2.1	$A_{km}B_{knp}$	$C_{m(np)} = A_{km}^\top B_{k(np)}$	$C_{mn[p]} = A_{km}^\top B_{kn[p]}$	5.1	$A_{pk}B_{kmn}$	$C_{(mn)p} = B_{k(mn)}^\top A_{pk}^\top$	$C_{m[n]p} = B_{km[n]}^\top A_{pk}^\top$
2.2	$A_{km}B_{kpn}$	$C_{mn[p]} = A_{km}^\top B_{k[p]n}$	$C_{m[n]p} = A_{km}^\top B_{kp[n]}$	5.2	$A_{pk}B_{knm}$	$C_{m[n]p} = B_{k[n]m}^\top A_{pk}^\top$	
2.3	$A_{km}B_{nkp}$	$C_{mn[p]} = A_{km}^\top B_{nk[p]}^\top$		5.3	$A_{pk}B_{mkn}$	$C_{m[n]p} = B_{mk[n]}A_{pk}^\top$	
2.4	$A_{km}B_{pkn}$	$C_{m[n]p} = A_{km}^\top B_{pk[n]}^\top$		5.4	$A_{pk}B_{nkm}$		
2.5	$A_{km}B_{npk}$	$C_{m(np)} = A_{km}^\top B_{(np)k}^\top$	$C_{mn[p]} = A_{km}^\top B_{n[p]k}^\top$	5.5	$A_{pk}B_{mnk}$	$C_{(mn)p} = B_{(mn)k}A_{pk}^\top$	$C_{m[n]p} = B_{m[n]k}A_{pk}^\top$
2.6	$A_{km}B_{pnk}$	$C_{m[n]p} = A_{km}^\top B_{p[n]k}^\top$		5.6	$A_{pk}B_{nmk}$		
3.1	$A_{nk}B_{kmp}$	$C_{mn[p]} = B_{km[p]}^\top A_{nk}^\top$		6.1	$A_{kp}B_{kmn}$	$C_{(mn)p} = B_{k(mn)}^\top A_{kp}$	$C_{m[n]p} = B_{km[n]}^\top A_{kp}$
3.2	$A_{nk}B_{kpm}$	$C_{mn[p]} = B_{k[p]m}^\top A_{nk}^\top$		6.2	$A_{kp}B_{knm}$	$C_{m[n]p} = B_{k[n]m}^\top A_{kp}$	
3.3	$A_{nk}B_{mkp}$	$C_{mn[p]} = B_{mk[p]}A_{nk}^\top$		6.3	$A_{kp}B_{mkn}$	$C_{m[n]p} = B_{mk[n]}A_{kp}$	
3.4	$A_{nk}B_{pkm}$			6.4	$A_{kp}B_{nkm}$		
3.5	$A_{nk}B_{mpk}$	$C_{mn[p]} = B_{m[p]k}A_{nk}^\top$		6.5	$A_{kp}B_{mnk}$	$C_{(mn)p} = B_{(mn)k}A_{kp}$	$C_{m[n]p} = B_{m[n]k}A_{kp}$
3.6	$A_{nk}B_{pmk}$			6.6	$A_{kp}B_{nmk}$		

FAST FOURIER TRANSFORM

$$[\mathbf{F}_N]_{ij} = e^{-2\pi\imath ij/N}$$

Ubiquitous primitive in computational sciences

Arguably one of the oldest and most well-studied

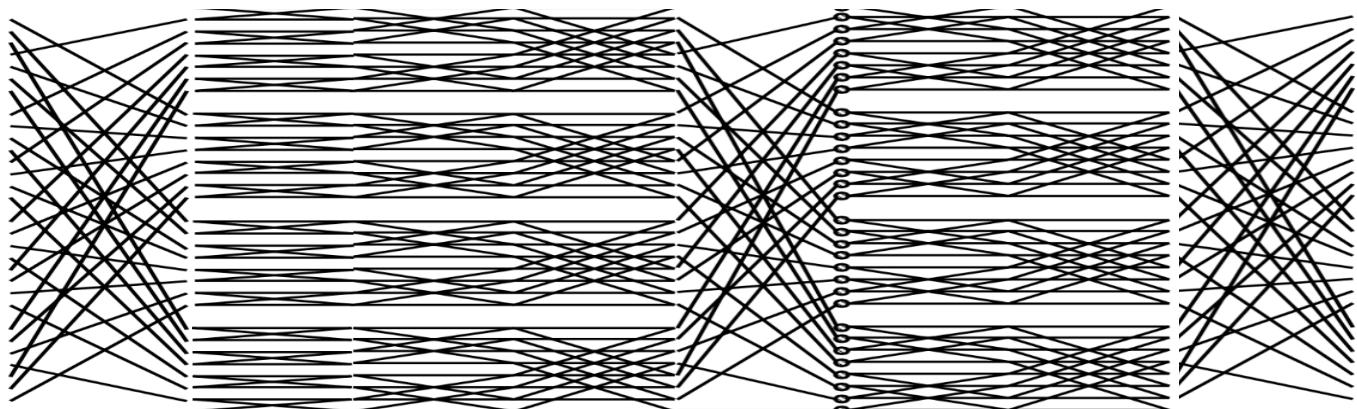
Performance on modern architectures...

SPLIT-RADIX FFT

Algorithm

$$\mathbf{F}_N = \boldsymbol{\Pi}_{M,P} (\mathbf{I}_M \otimes \mathbf{F}_P) \boldsymbol{\Pi}_{P,M} \mathbf{T}_{M,P} (\mathbf{I}_P \otimes \mathbf{F}_M) \boldsymbol{\Pi}_{M,P}$$

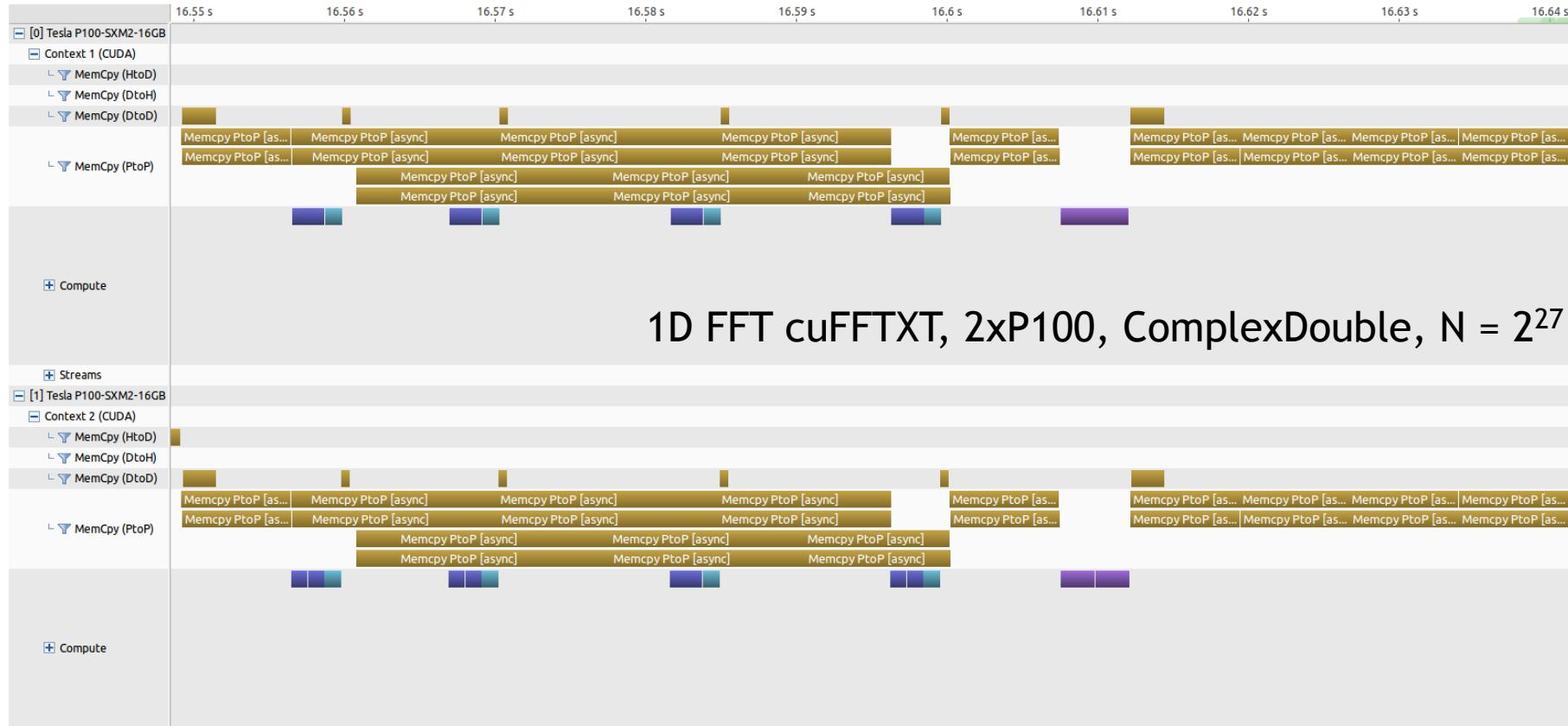
1. Transpose $P \rightarrow M$.
2. P local FFTs of size M .
3. Twiddle factors \mathbf{T} .
4. Transpose $M \rightarrow P$.
5. M local FFTs of size P .
6. Transpose $P \rightarrow M$.



SPLIT-RADIX FFT

Profile

GPU 1



GPU 2

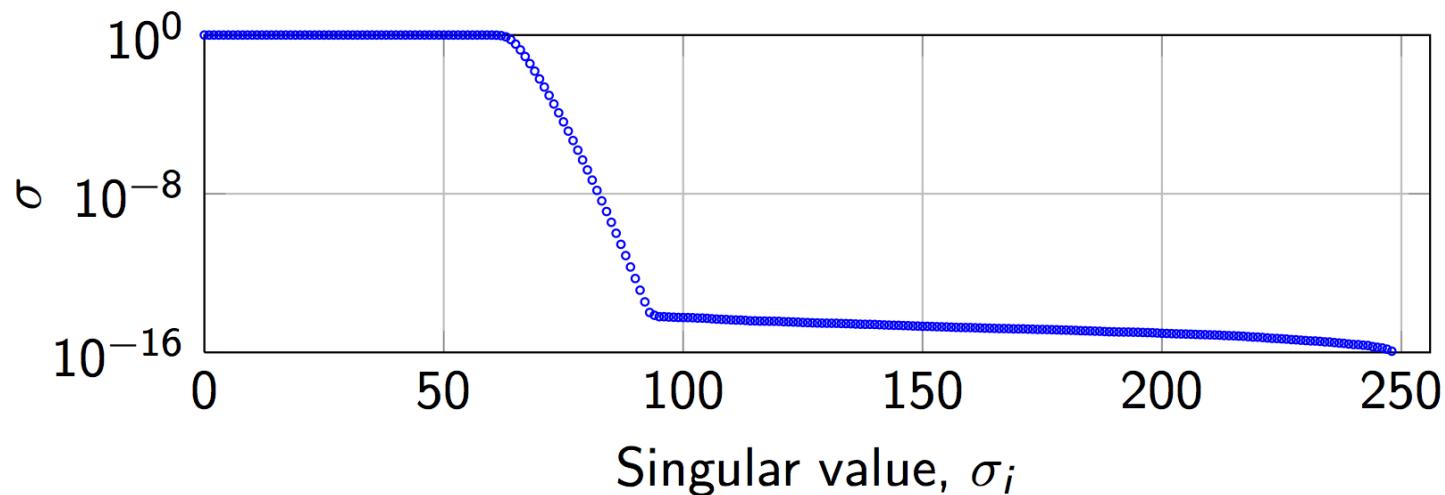
FMM-FFT

Edelman et al. 1999

$$F_N = \begin{pmatrix} F_{N|P} & \cdots & D^{P-1}F_{N|P} \\ F_{N|P}D & \cdots & D^{P-1}F_{N|P}D \\ \vdots & & \vdots \\ F_{N|P}D^{P-1} & \cdots & D^{P-1}F_{N|P}D^{P-1} \end{pmatrix}$$

Low-rank structure

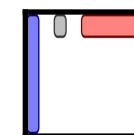
$$\sigma_i(\mathcal{F}_{1024}[0:256, 0:256])$$



STRUCTURED DENSE MATRICES AND FMM

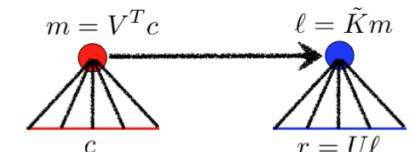
- SVD:

$$A = U D V^*$$



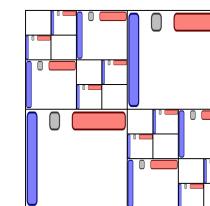
- Low-Rank:

$$K = U \tilde{K}_{r \times r} V^*$$



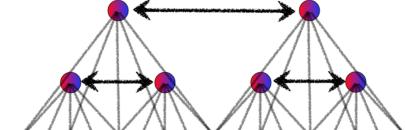
- Hierarchically LR:

$$K_{\mathcal{I}\mathcal{J}} = U_{\mathcal{I}} \tilde{K}_{\mathcal{I}\mathcal{J}} V_{\mathcal{J}}^*$$

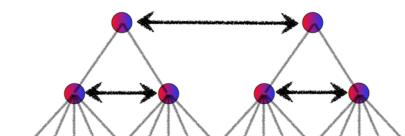


- H-Semi-Separable:

$$K_{\mathcal{I}\mathcal{J}} = U_{\mathcal{I}} \tilde{U}_{\tilde{\mathcal{I}}} \tilde{K}_{\tilde{\mathcal{I}}\tilde{\mathcal{J}}} \tilde{V}_{\tilde{\mathcal{J}}}^* V_{\mathcal{J}}^*$$



- H²-Matrix/FMM



FMM-FFT

Edelman et al. 1999 Algorithm

$$\mathbf{F}_N = (\mathbf{I}_P \otimes \mathbf{F}_M) \boldsymbol{\Pi}_{M,P} (\mathbf{I}_M \otimes \mathbf{F}_P) \mathbf{M}_{M,P}$$
$$\mathbf{M}_{M,P} = \boldsymbol{\Pi}_{P,M} \operatorname{diag}(\mathbf{I}_M, \mathbf{C}_1, \mathbf{C}_2, \dots, \mathbf{C}_p) \boldsymbol{\Pi}_{M,P}$$

Block diagonal matrix with P blocks of M x M.

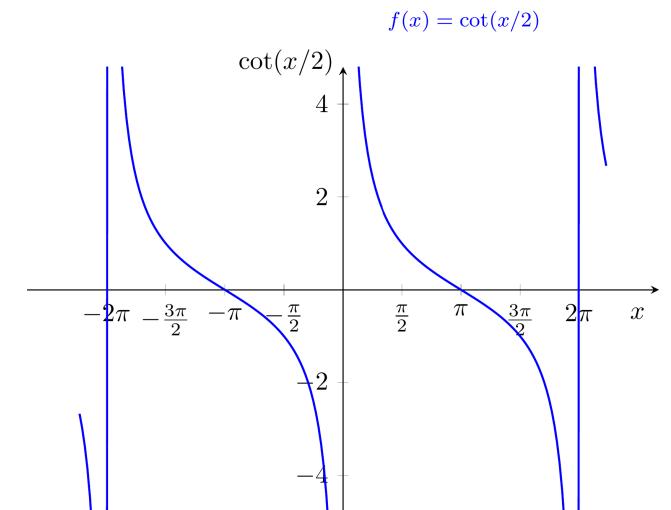
Direct application would destroy asymptotic complexity of the FFT and not effect communication.

However, each C is a *dense kernel matrix*.

COTANGENT FMM

$$[\mathbf{C}_p]_{mn} = \rho_p \left[\cot\left(\frac{\pi}{M}\left(n - m + \frac{p}{P}\right)\right) + i \right]$$

- One dimensional
- Uniform – source/target are the integers
- Periodic
- Distributed
- Size M-by-M
- P of them!
 - Interleaved



FMM-FFT

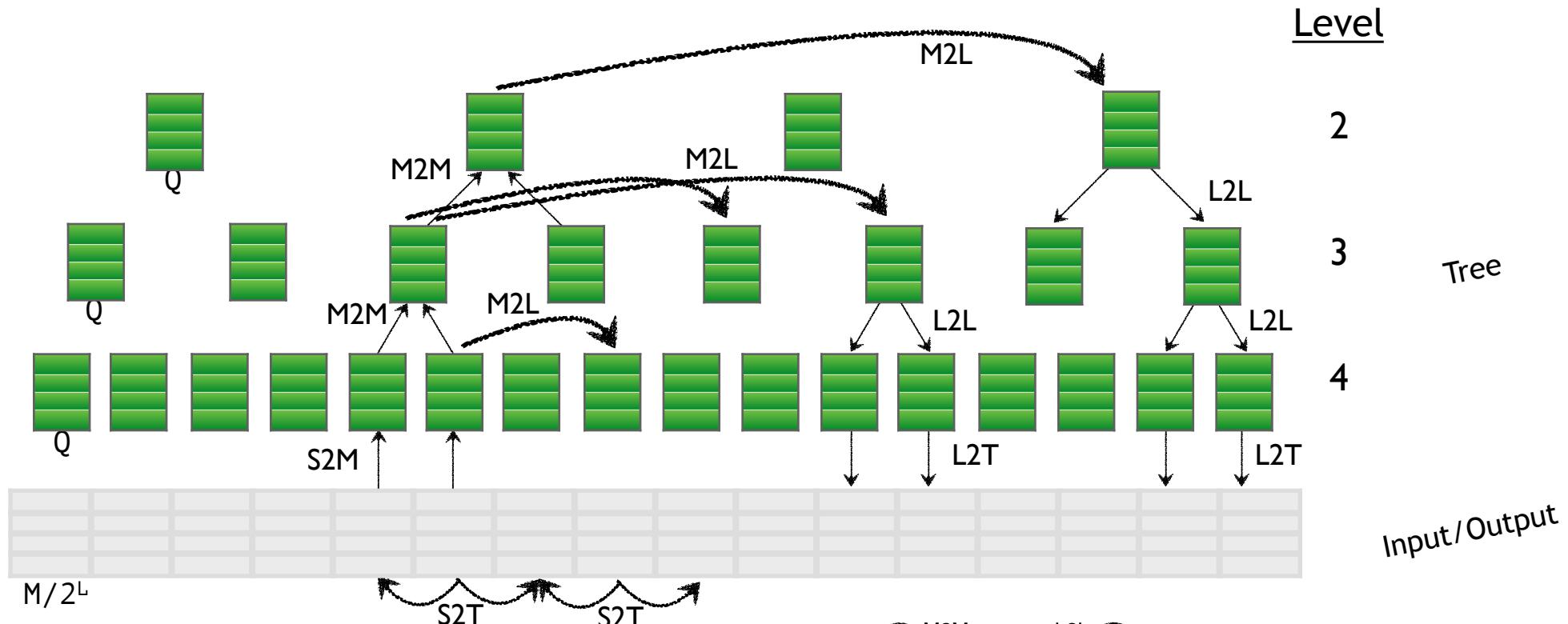
Edelman et al. 1999 Algorithm

$$\mathbf{F}_N = (\mathbf{I}_P \otimes \mathbf{F}_M) \boldsymbol{\Pi}_{M,P} (\mathbf{I}_M \otimes \mathbf{F}_P) \mathbf{M}_{M,P}$$
$$\mathbf{M}_{M,P} = \boldsymbol{\Pi}_{P,M} \operatorname{diag}(\mathbf{I}_M, \mathbf{C}_1, \mathbf{C}_2, \dots, \mathbf{C}_p) \boldsymbol{\Pi}_{M,P}$$

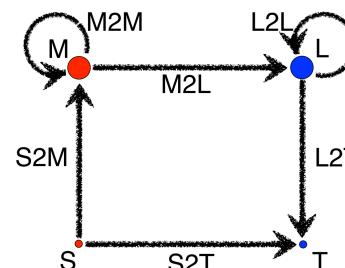
1. Compute $P - 1$ distributed FMMs of size $M \times M$.
 2. M local FFTs of size P .
 3. Transpose $P \rightarrow M$.
 4. P local FFTs of size M .
- } 2D $M \times P$ FFT

FMM OPERATORS

- S: “Source”
- T: “Target”
- M: “Multipole”
- L: “Local”



Each operator is an (implicit) matrix.



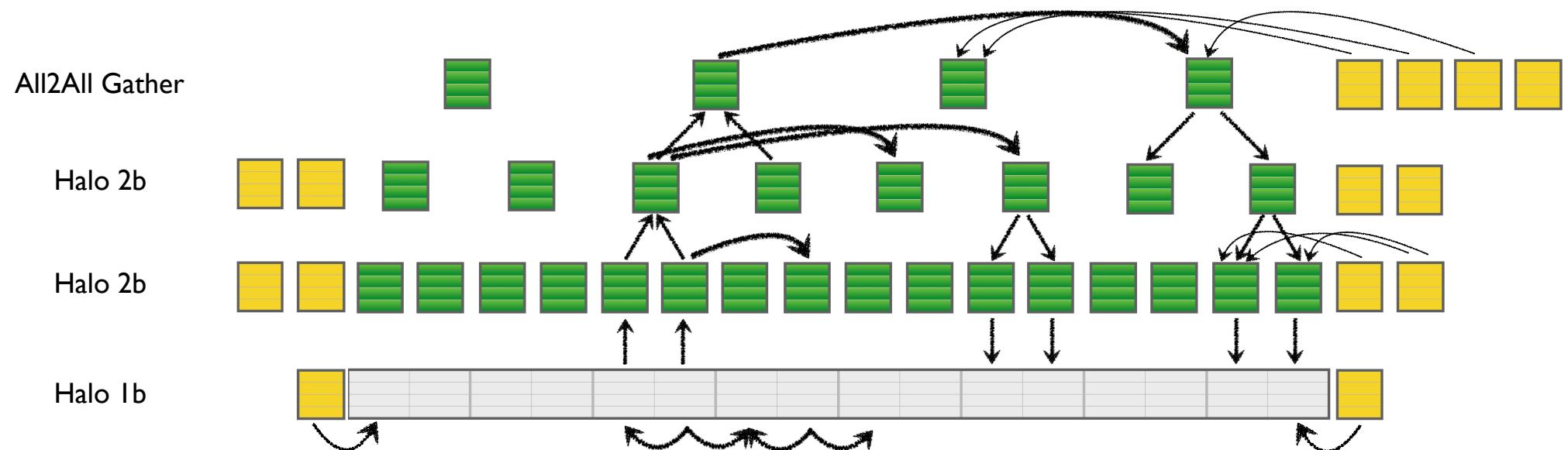
PARAMETERS OF THE FMM-FFT

- FFT $N = M P$

P FMMs to perform of size $M \times M$

- FMM
 - Rank Q
 - Leaf box size M_L
 - Base level B
 - Leaf level $L = \log_2(M/M_L)$

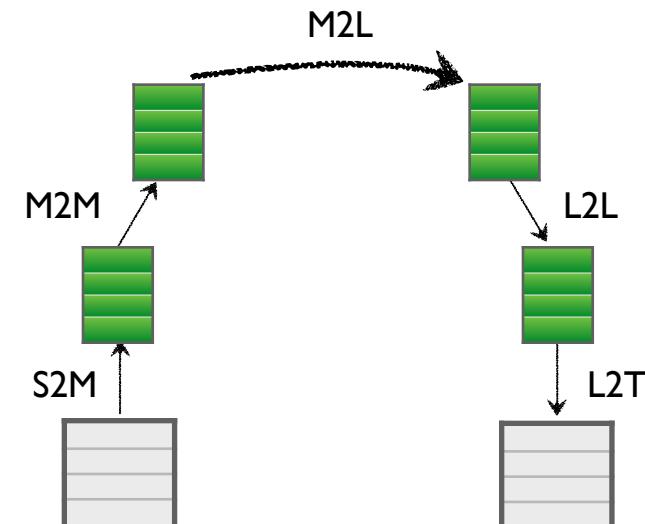
DISTRIBUTED FMM



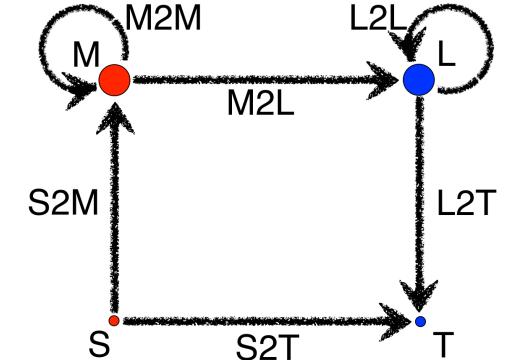
INTERPOLATIVE FMM

Chebyshev interpolations

- Use Chebyshev interpolations as the intermediate data
 - Same operators across all boxes
 - Same operators across all levels
 - Almost same operators across all FMMs
- Express as tensors and efficiently compute.

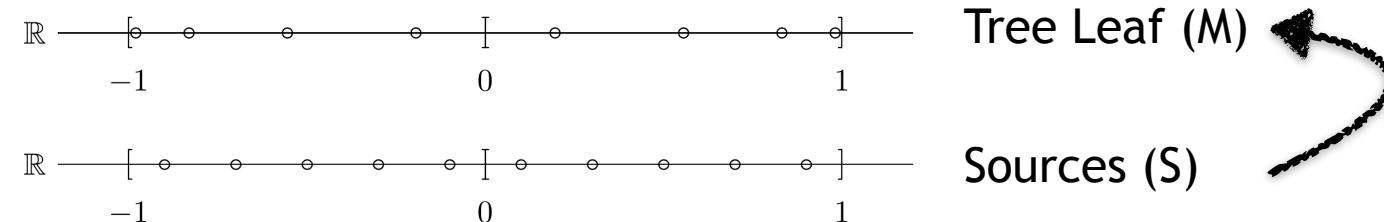


S2M/L2T



$$\mathcal{M}_{(p-1)qb}^L = S2M_{qm} \mathcal{S}_{pmb}$$

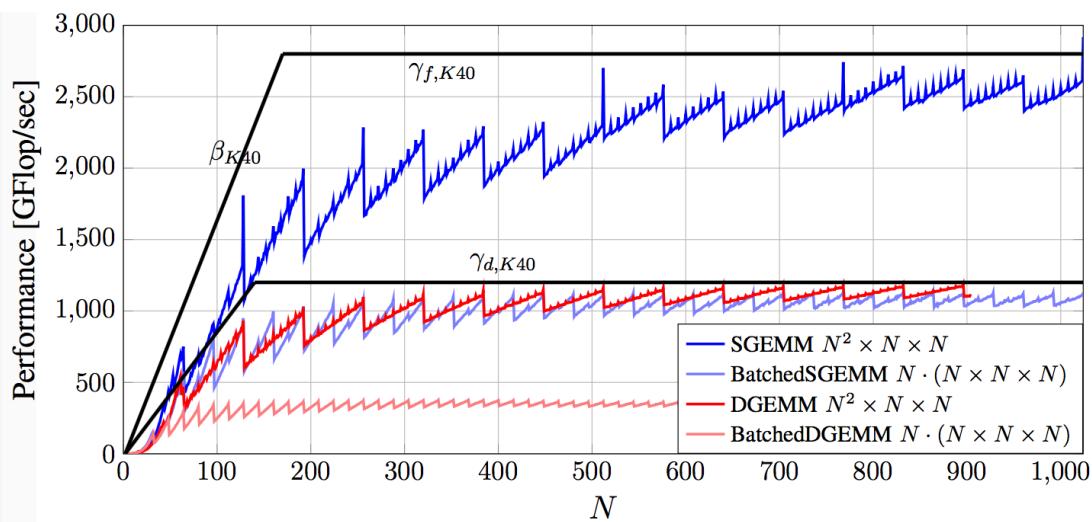
$$\mathcal{T}_{pmb} = L2T_{mq} \mathcal{L}_{(p-1)qb} + \mathcal{T}_{pmb}$$



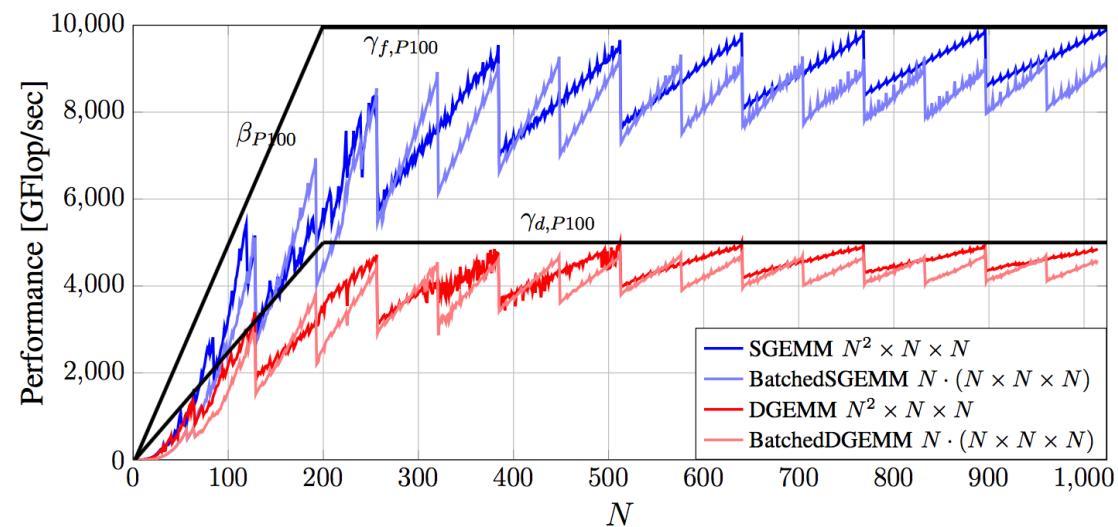
Computed with single BatchedGEMM

BATCHED MATRIX-MATRIX MULTIPLY

cublas<T>gemmStridedBatched
in cuBLAS 8.0



(a) K40c GPU

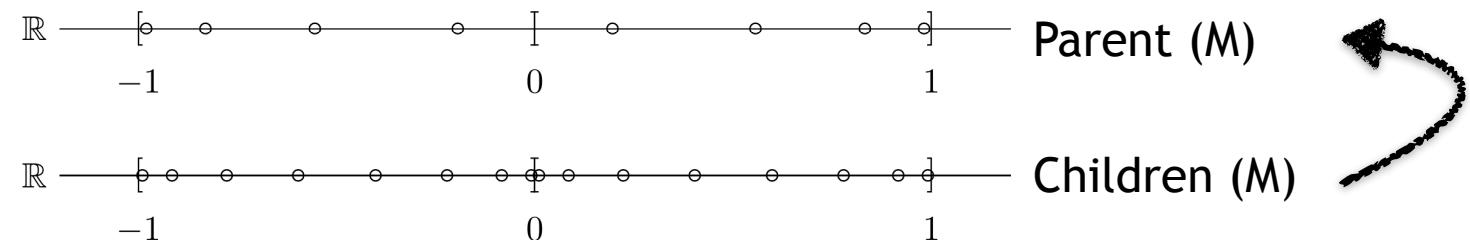
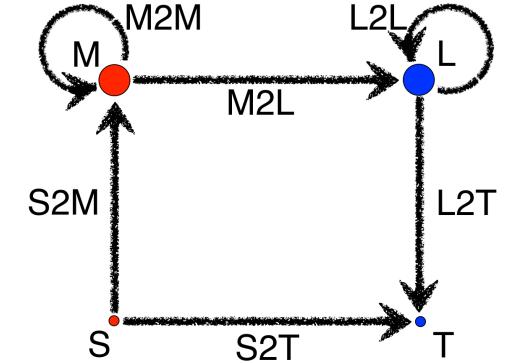


(b) P100 GPU

M2M/L2L

$$\mathcal{M}_{pqb}^{\ell} = M2M_{qk} \mathcal{M}_{pk(2b)}^{\ell+1}$$

$$\mathcal{L}_{pq(2b)}^{\ell+1} = L2L_{qk} \mathcal{L}_{pkb}^{\ell} + \mathcal{L}_{pq(2b)}^{\ell+1}$$



Computed with single BatchedGEMM

S2T/M2L

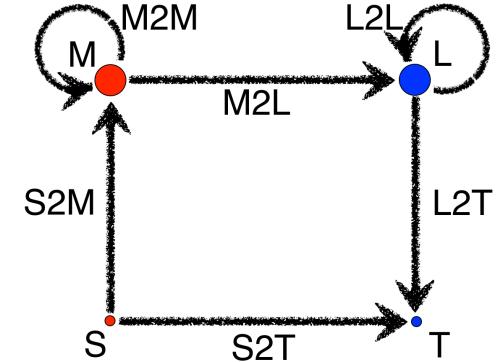
$$\mathcal{T}_{pib} = S2T_{p(j-i)} \mathcal{S}_{pj b}$$

$$S2T_{pk} = \begin{cases} \cot\left(\frac{\pi}{N}(p + Pk)\right) & p > 0 \\ \delta_{k0} & p = 0 \end{cases}$$

$$\mathcal{L}_{pib}^\ell = M2L_{pijs}^\ell \mathcal{M}_{pj(b+s)}^\ell$$

$$M2L_{pijs}^\ell = \cot\left(\frac{\pi}{2^\ell}\left(\frac{z_j}{2} - \frac{z_i}{2} + s\right) + \frac{\pi}{N}(p+1)\right)$$

- Also Level-3 Linear Algebra computations, but no BLAS primitives.
 - CUSTOM KERNELS



INTERPOLATIVE FMM

Operator

$$S2T_{ijs}^{(p)} = \cot\left(\frac{\pi}{M}(j-i) + \frac{\pi}{2^L}s + \frac{\pi}{N}p\right)$$

$$S2M_{ij} = \ell_i(s_j) = \begin{cases} 1 & s_j = z_i \\ \frac{\lambda_i \ell(s_j)}{s_j - z_i} & \text{else} \end{cases}$$

$$M2M_{ij}^{\pm} = \ell_i\left(\frac{z_j \pm 1}{2}\right)$$

$$M2L_{ijs}^{(p)} = \cot\left(\frac{\pi}{2^L}\left(\frac{z_j}{2} - \frac{z_i}{2}\right) + \frac{\pi}{2^L}s + \frac{\pi}{N}p\right)$$

$$L2L_{ij}^{\pm} = \ell_j\left(\frac{z_i \pm 1}{2}\right) = M2M_{ji}^{\pm}$$

$$L2T_{ij} = \ell_j(t_i) = S2M_{ji}$$

Storage

$$P(4M_L - I)$$

$$QM_L$$

$$2Q^2$$

$$4(L-B)PQ^2$$

$$2Q^2$$

$$QM_L$$

Compute

$$3P2^L M_L^2$$

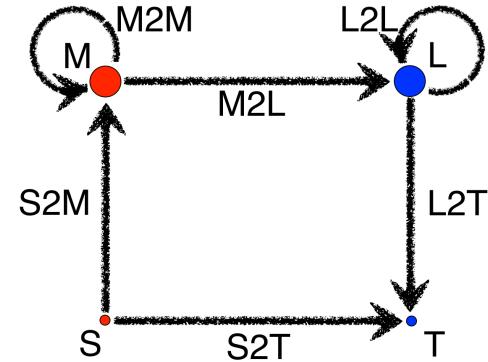
$$2PMQ$$

$$4(2^L - 2^B)PQ^2$$

$$3(2^L - 2^B)PQ^2$$

$$4(2^L - 2^B)PQ^2$$

$$2PMQ$$



Compute on the fly!

Total Precomputation ($S2M + M2M$):
 $Q^2 M_L + Q^3 \sim 20k$ Flops

Total Overhead ($S2M + M2M$):
 $QM_L + Q^2 \sim 12k$ Bytes

ALGORITHM

Algorithm 1 The FMM-FFT via tensor contractions.

```

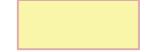
1:  $\mathcal{M}_{(p-1)qb}^L = S2M_{qm} \mathcal{S}_{pmb}$  // S2M
2: SendRecv  $\mathcal{S}$  halo. // COMM  $\mathcal{S}$ 
3:  $\mathcal{T}_{pib} = S2T_{p(j-i)} \mathcal{S}_{pj}$  // S2T
4: for  $\ell = L - 1, \dots, B$  do
5:    $\mathcal{M}_{pq}^\ell = M2M_{qk} \mathcal{M}_{pk(2b)}^{\ell+1}$  // M2M
6:   for  $\ell = L, \dots, B - 1$  do
7:     SendRecv  $\mathcal{M}^\ell$  halo. // COMM  $\mathcal{M}^\ell$ 
8:      $\mathcal{L}_{pib}^\ell = M2L_{pijs}^\ell \mathcal{M}_{pj(b+s)}^\ell$  // M2L  $\ell$ 
9:   All-to-all gather  $\mathcal{M}^B$ . // COMM  $\mathcal{M}^B$ 
10:   $\mathcal{L}_{pib}^B = M2L_{pijs} \mathcal{M}_{pj(b+s)}^B$  // M2L  $B$ 
11:   $r_p = 1_{qb} \mathcal{M}_{(p-1)qb}^B$  // REDUCE
12:  for  $\ell = B, \dots, L - 1$  do
13:     $\mathcal{L}_{pq(2b)}^{\ell+1} = L2L_{qk} \mathcal{L}_{pkb}^\ell + \mathcal{L}_{pq(2b)}^{\ell+1}$  // L2L
14:     $\mathcal{T}_{pmb} = L2T_{mq} \mathcal{L}_{(p-1)qb}^L + \mathcal{T}_{pmb}$  // L2T
15:     $\mathcal{T}_{pmb} = \rho_p (\mathcal{T}_{pmb} + \imath r_p)$  // POST
16:     $\mathbf{F}_{M,P} \mathcal{T}$  // 2D FFT

```

 : Primitive – vendor optimized

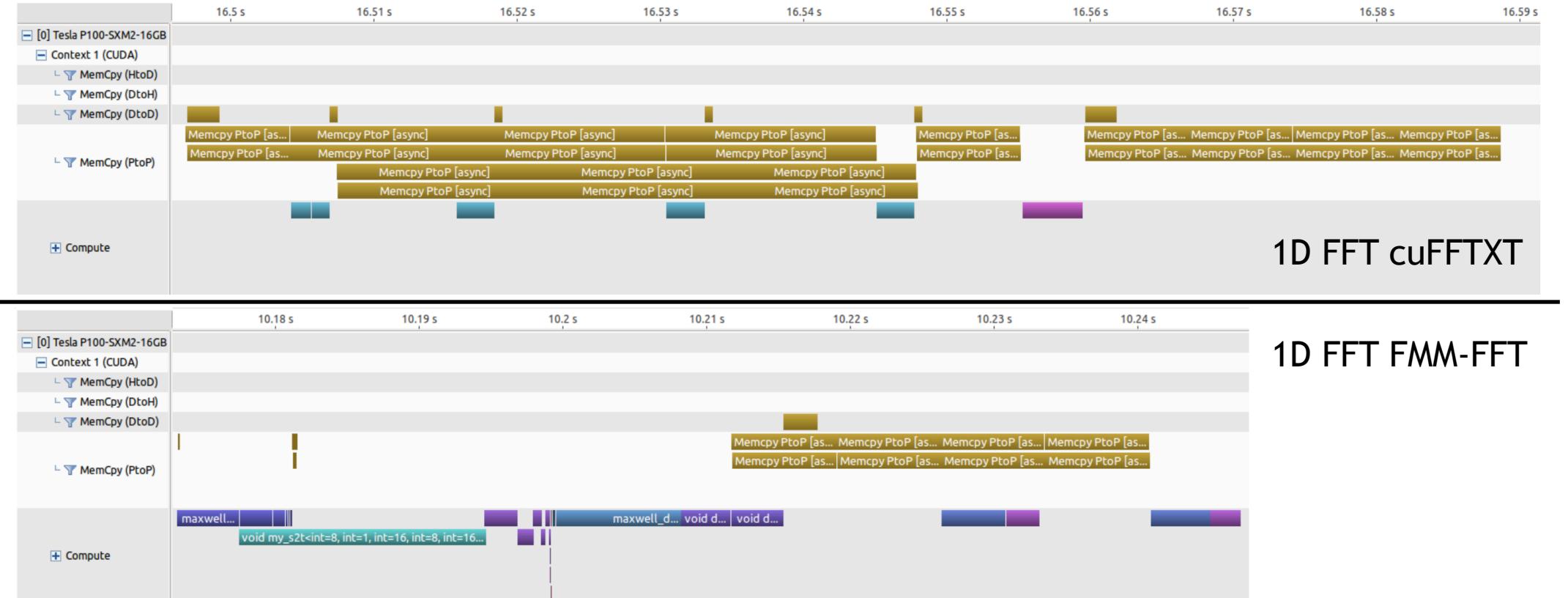
cublas<T>gemmStridedBatched

 : Custom kernel

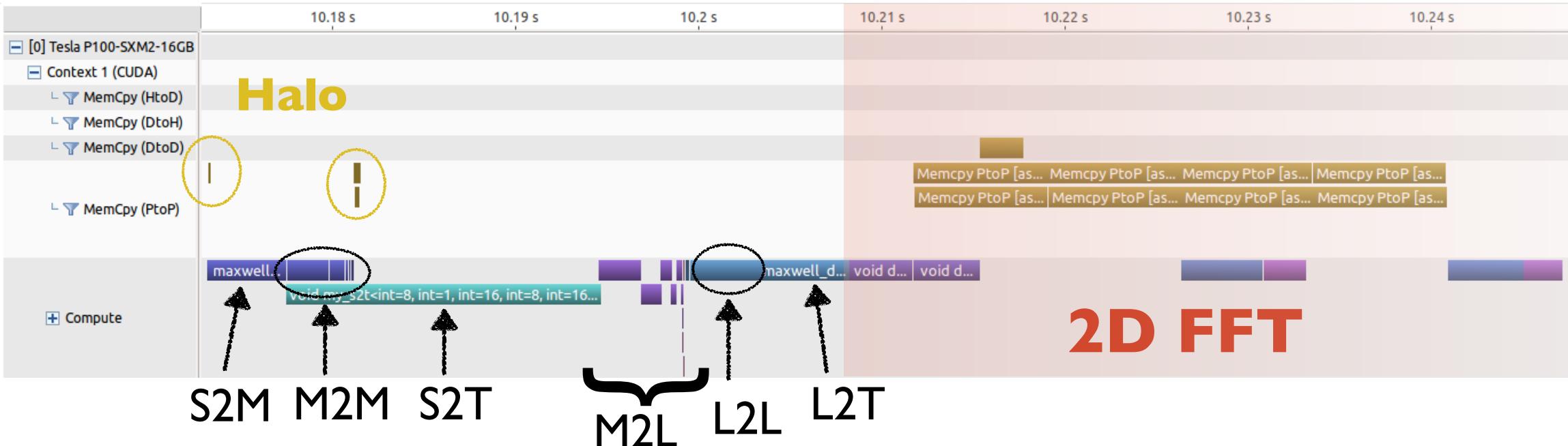
 : FMM Communication

COMPARISON

Profile



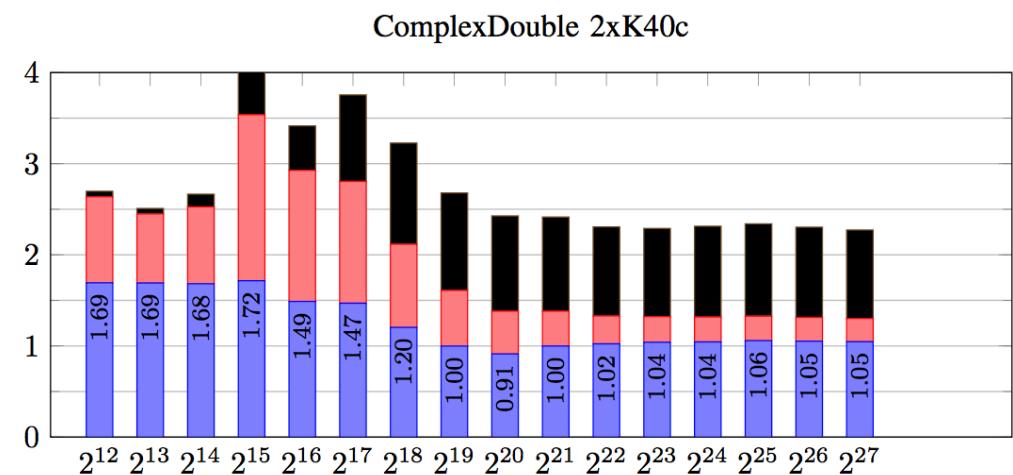
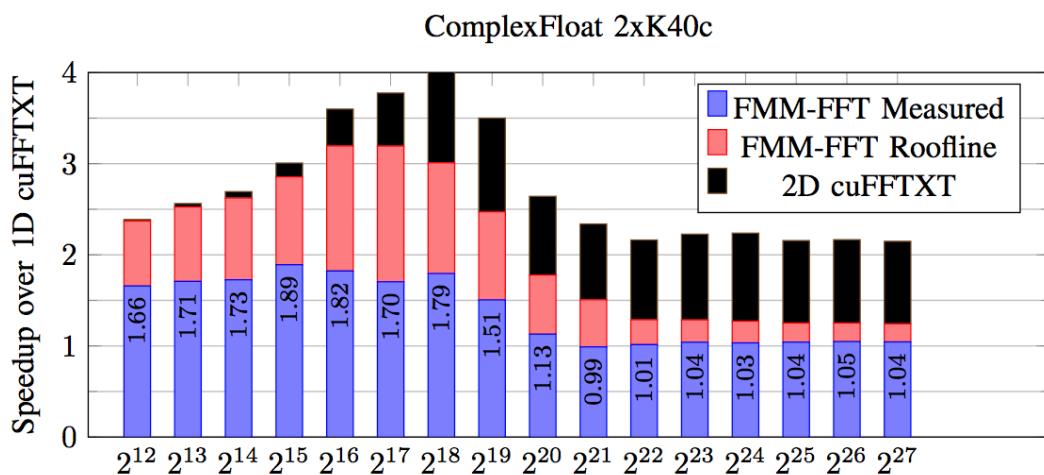
FMM-FFT PROFILE



Emph: 255 FMMs of size 524k x 524k in 32ms

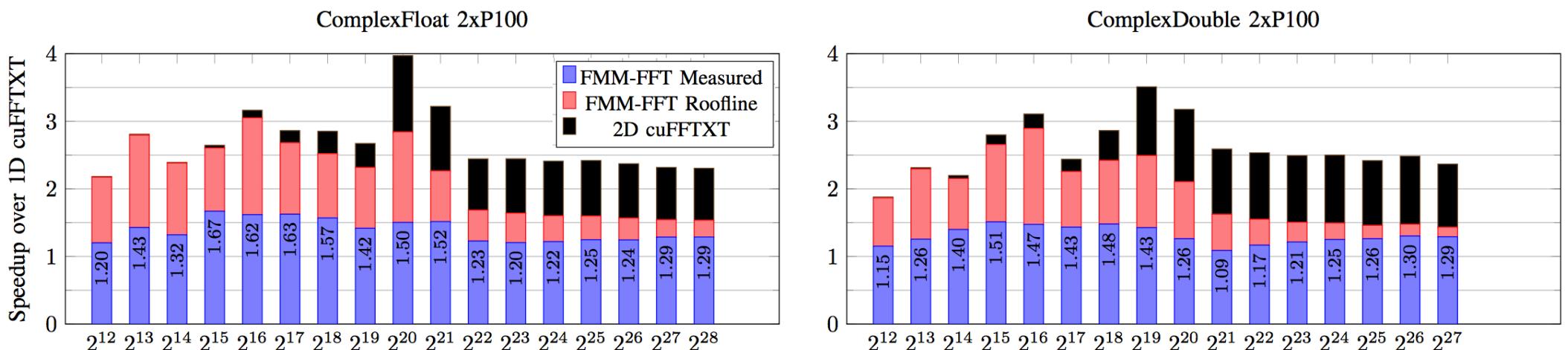
2xK40c, PCIe

FMM-FFT Performance



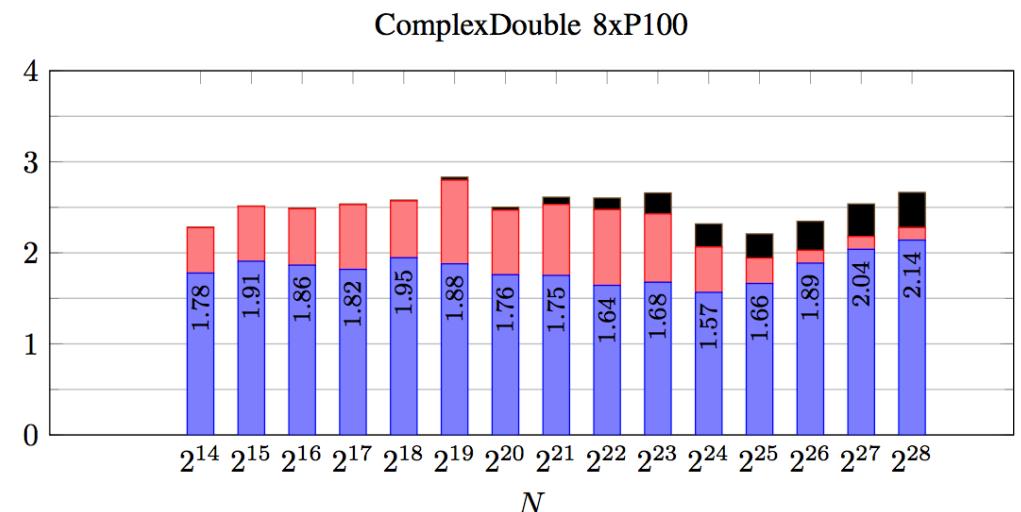
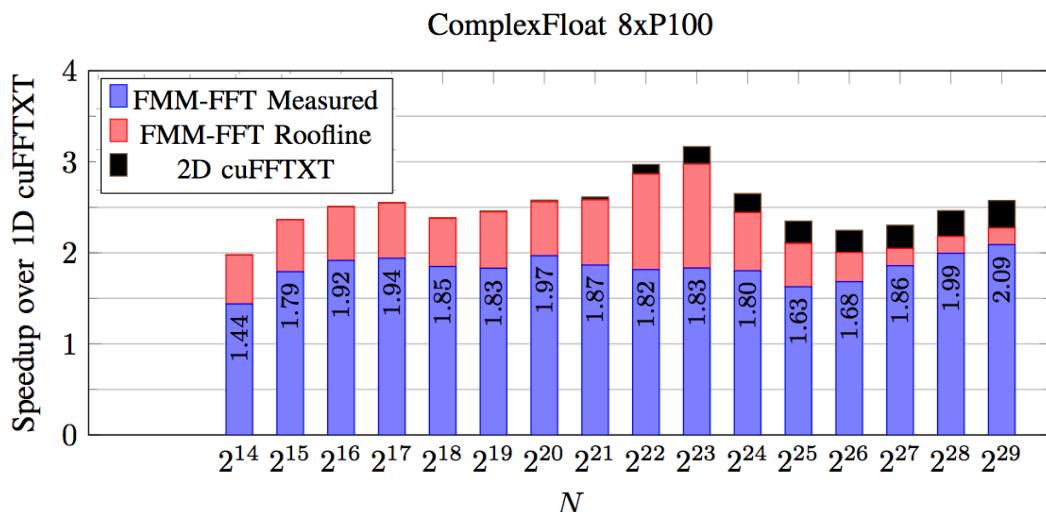
2xP100, NVLINK

FMM-FFT Performance

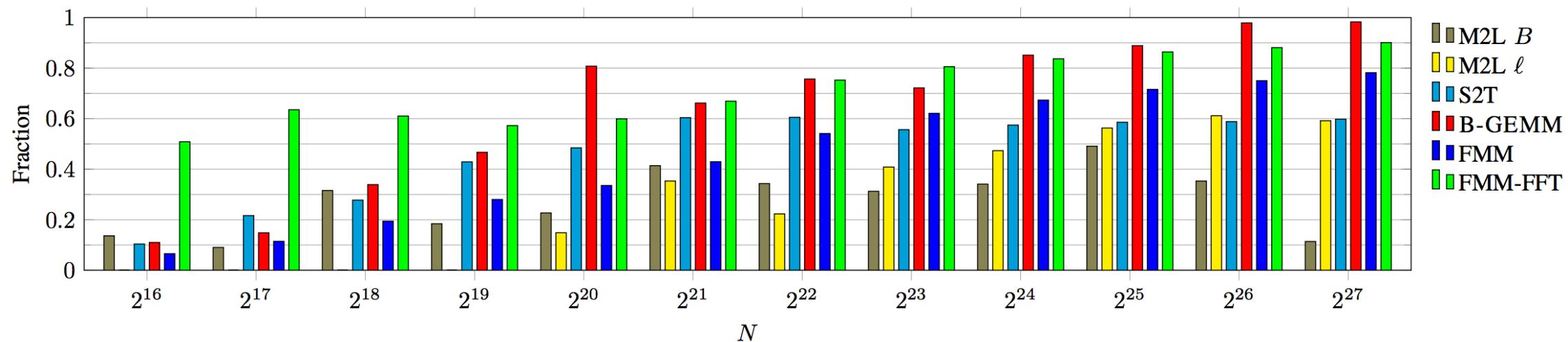


8xP100, NVLINK

FMM-FFT Performance



EFFICIENCY

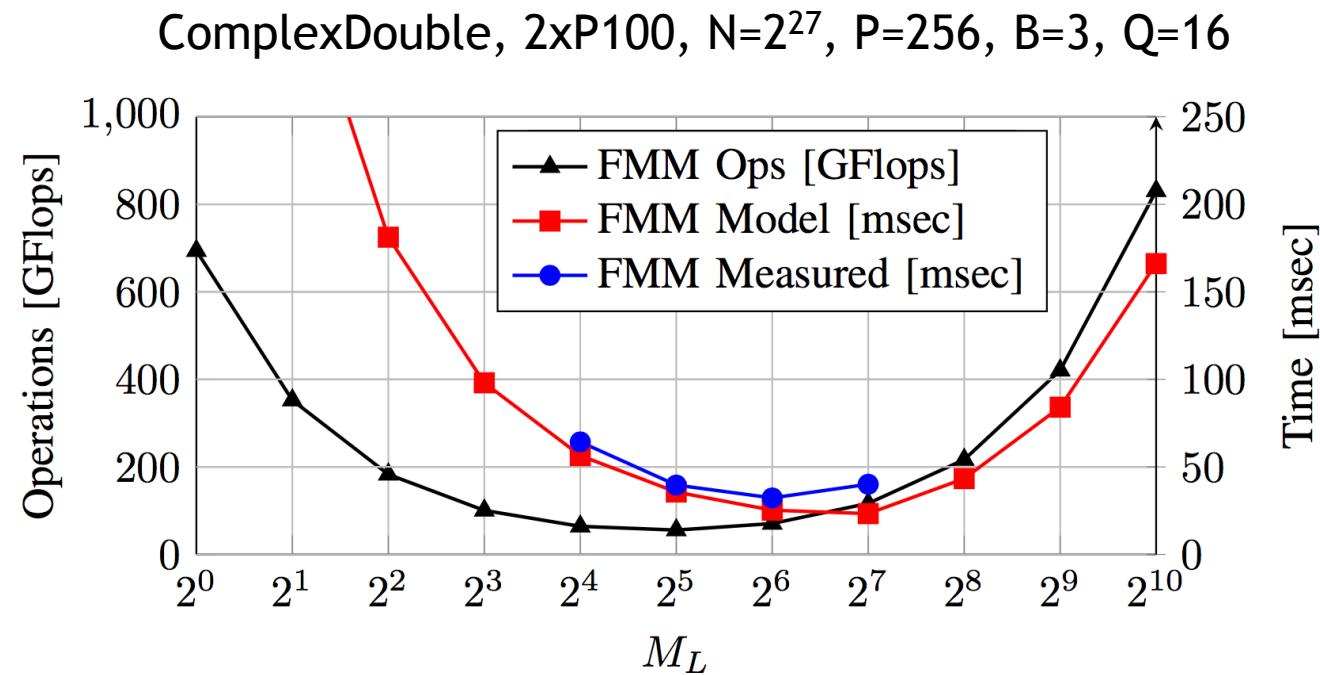


- >95% BatchedGEMM
- 60% S2T/M2L
- >90% FMM-FFT

PARAMETER DEPENDENCE – M_L

Points per leaf box per FMM

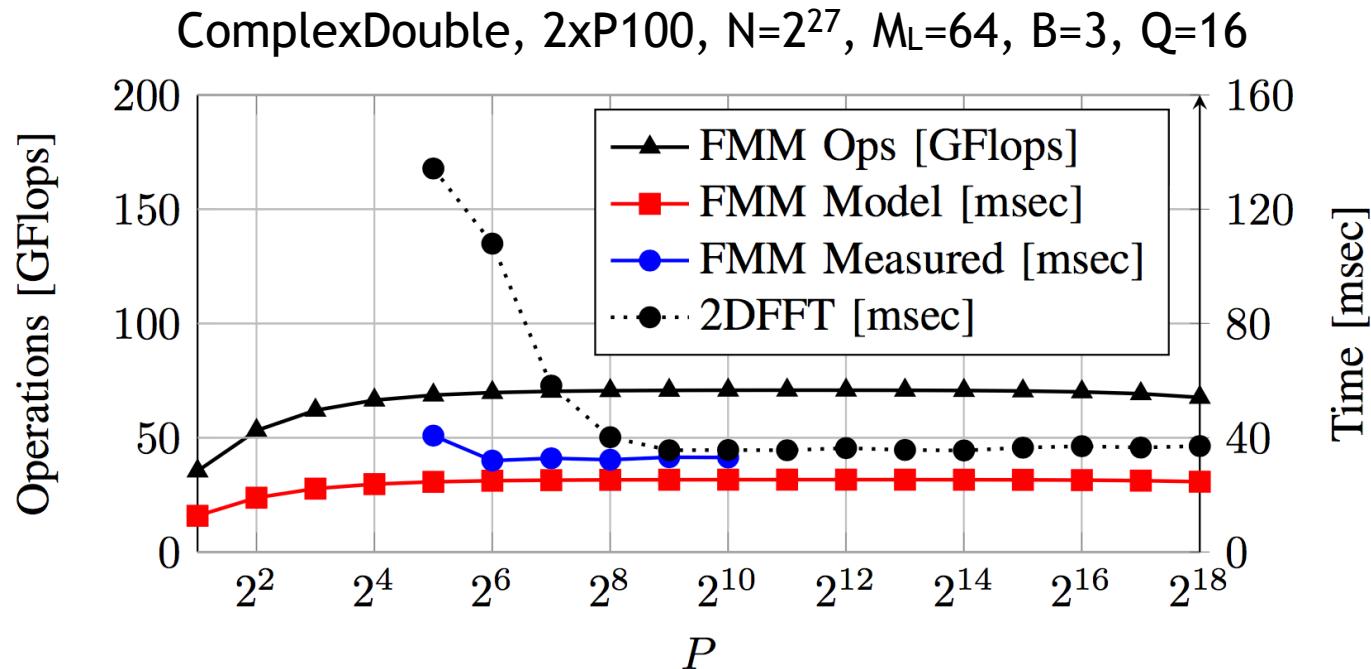
- Trade #levels for S2T comp
- Counting flops not enough
 - Want intensity increase
- Tune performance for $M_L=64$



PARAMETER DEPENDENCE – P

Number of FMMs

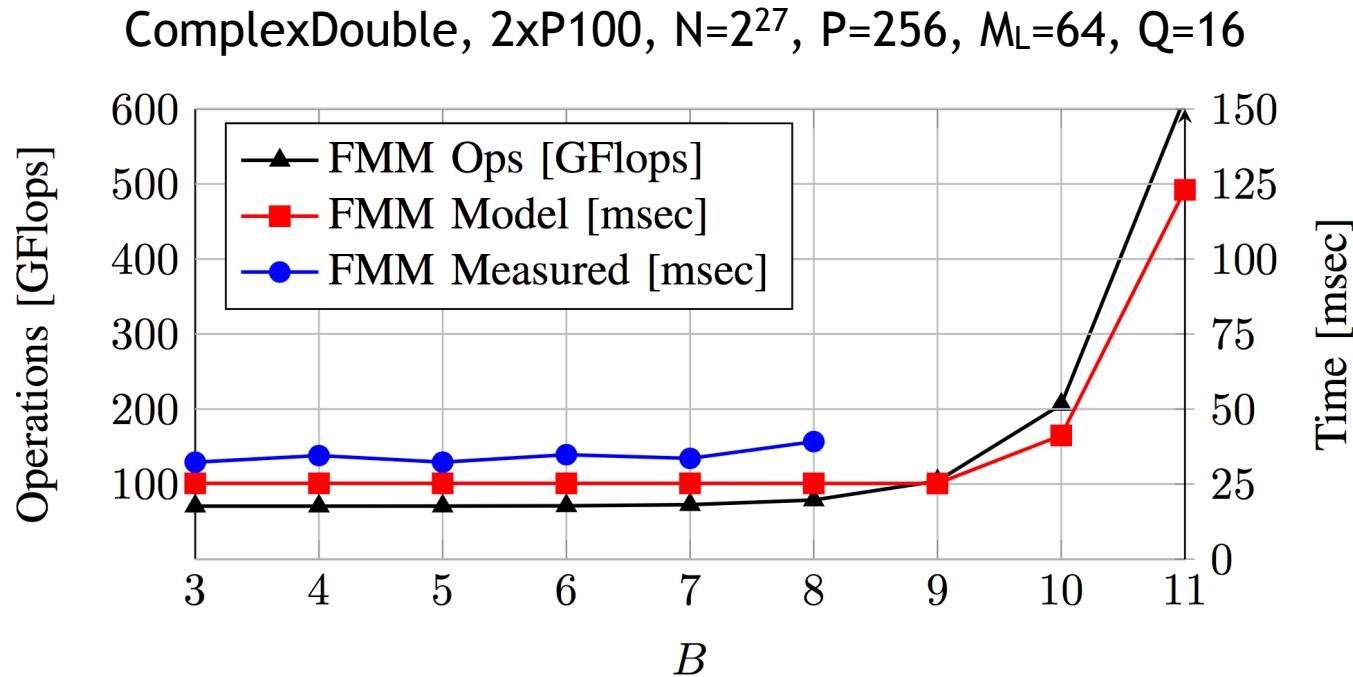
- Flops/Intensity approx constant
 - Trade #levels for #FMMs
- Large P good
 - Fill up B-GEMM
 - More square 2D FFT



PARAMETER DEPENDENCE – B

Base level of FMM trees

- Not very significant
- Scale to 128 GPUs w/o complications

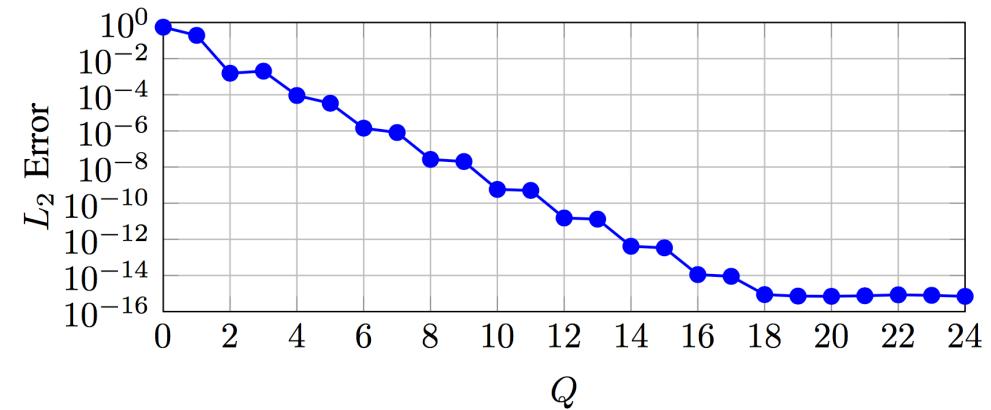
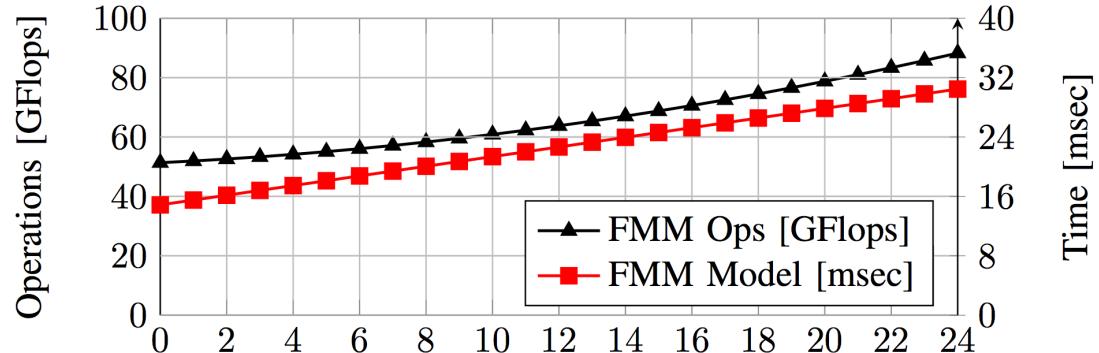


PARAMETER DEPENDENCE – Q

Quadrature order of FMM interpolation

- Accuracy tuning with performance dependence

ComplexDouble, 2xP100, N=2²⁷, P=256, M_L=64, B=3



FUTURE

- Integration into cuFFT
- Application to 2D/3D FFTs?
 - Convolutions
- NUFFT, Sparse FFT nearly out-of-the-box
- Volta predictions and measurements
 - Mixed precision (e.g. FP16 far-field) to use Tensor Core?
- Persistent Matrix Batched GEMM (cuBLAS optimization)
 - Staged Persistent Matrix Batched GEMM (cooperative groups, RNNs)

CONCLUSION

- FMM-FFT trades 2/3 communication in 1D FFT for P FMMs
 - Viable on highest comp:comm architecture available
- Detailed implementation that relies heavily on existing primitives
 - Primitives >95% efficient
 - Two custom dense kernels >60% efficient
 - Entire FMM-FFT >90% efficient
- Tunable accuracy-performance tradeoff
- Compute model accurately predicts performance

Cris Cecka, “Low-communication FMM-accelerated FFTs on GPUs” *International Conference for High Performance Computing, Networking, Storage, and Analysis (SC)*, Denver, CO, November 12-17, 2017

Yang Shi, U.N. Niranjan, Animashree Anandkumar, and **Cris Cecka**.
“Tensor contractions with extended BLAS kernels on CPU and GPU” In *2016 IEEE 23rd International Conference on High Performance Computing (HiPC)*, pages 193–202, Dec 2016