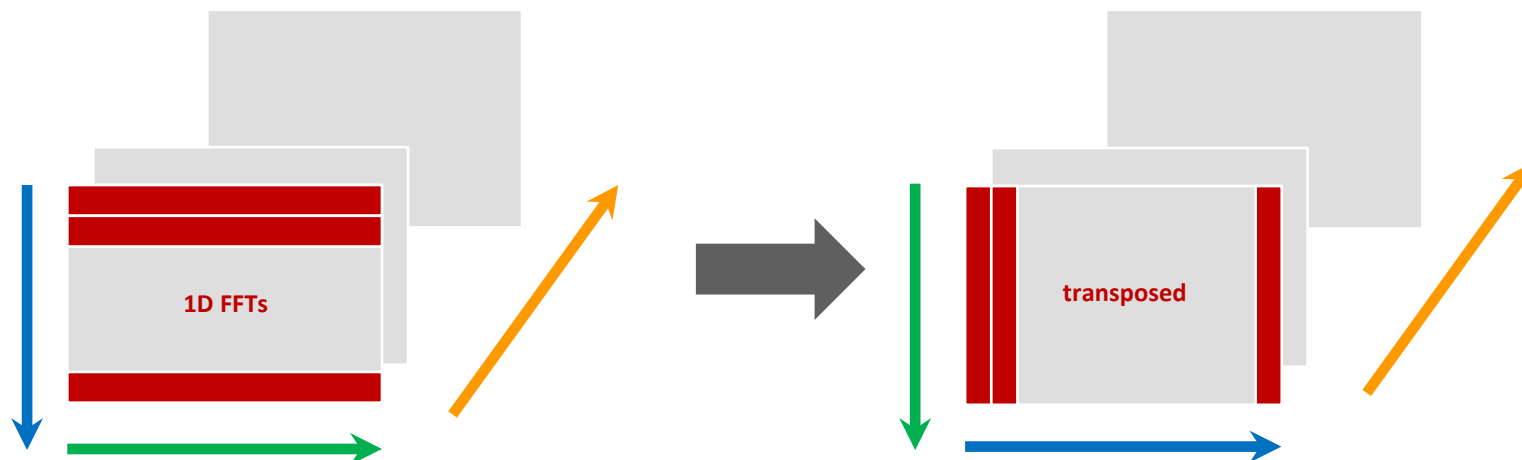# HPC Libraries as DSL

**Franz Franchetti**
Carnegie Mellon University

**In collaboration with**
Tze Meng Low, Qi Guo, and Berkin Akin

**Supported by the DARPA PERFECT Program**

# How (Not To) Program a Computational Kernel



```
for (chan = 0; chan < N_CHAN; ++chan)
{
  for (range = 0; range < N_RANGE; ++range)
  {
    /* Apply a 1D FFT to convert to Doppler space */
    fft(
      (float *)&datacube_pulse_major_padded[chan][range][0],
      N_DOP,
      N_DOP_LOG2,
      FFT_FORWARD);
  }
}
/* Transpose from doppler major to range major*/
corner_turn_to_range_major_order
(doppler_datacube, datacube_pulse_major_padded);
```
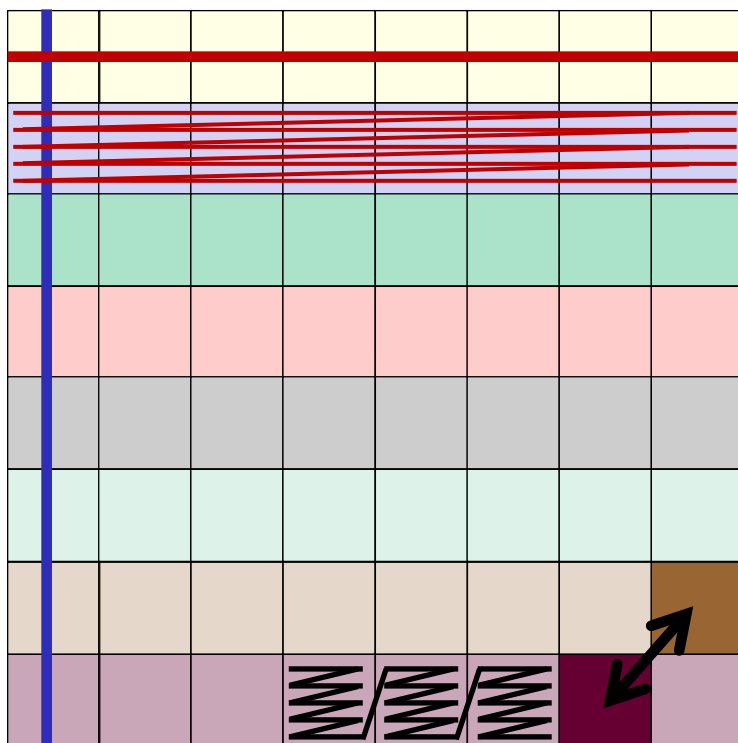
**7.5x**

```
/* Replace with fftw guru interface*/
plan = fftwf_plan_guru_dft(
         1, dims,
         2, howmany_dims,
         datacube_pulse_major_padded,
         doppler_datacube,
         FFTW_FORWARD,
         FFTW_WISDOM_ONLY);

fftwf_execute(plan);
fftwf_destroy_plan(plan);
```
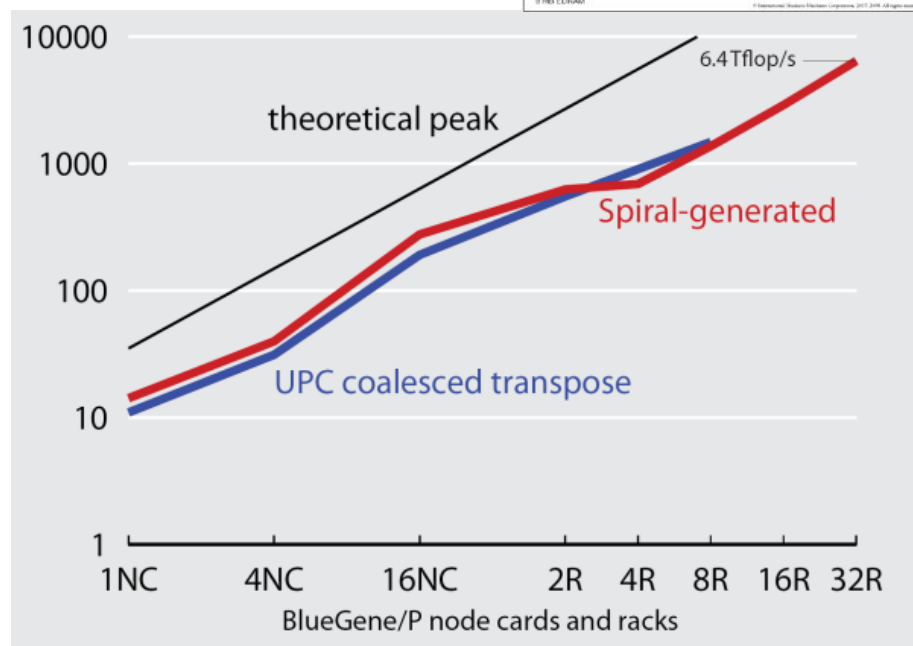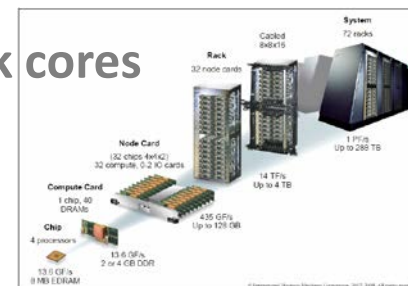
T. M. Low, Q. Guo, F. Franchetti: **Optimizing Space Time Adaptive Processing Through Accelerating Memory-Bounded Operations.** IEEE High Performance Extreme Computing Conference (HPEC), 2015.

Electrical & Computer
**ENGINEERING**

# But: Libraries Impose Incompatible Constraints

**FFTW: contiguous rows/columns**

**128k cores**



**MPI: contiguous sub-blocks**

**30% gain from MPI-aware FFTW**

# Using Libraries—Isn't that best HPC practice?

**In theory, yes. But**

- **Sometimes looks like overkill**
  Add a whole library for a one-page function?!

- **People often do not use them**
  dependencies, don't know them,
  "not invented here"

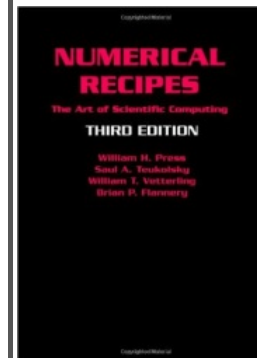- **Need to combine many libraries**
  MPI, FFTW, LAPACK, Boost, STL,…

- **Uneven performance and coverage across functions**
  not whole FFTW interface well-supported by MKL, ESSL,…

**1 page vs. 30 MB C code**

# Idea 1: Standard HPC Libraries as Kernel DSL

## DSL definition and syntax

- **API of standard HPC libraries**
  FFTW, LAPACK, BLAS, SparseBLAS, GraphBLAS

- **Libraries and language extensions for parallelism**
  OpenMP, MPI, OpenACC

- **Subset of C/C++**
  Single threaded, side-effect free, only standard C library calls

## DSL semantics and user knowledge

- **C semantics + library semantics**
  Program can be executed

- **Use OpenMP/OpenACC annotations to communicate meta-information**
  ```
  #pragma omp for private(…) shared(…)
  ```

# Idea 2: Interpret Program as Specification

- **Make a *subset* of "C+OpenMP+MKL" a DSL with DSL compiler**
  Combines code synthesis, telescoping languages, and HPC best practices

- **Treat DSL program as *specification*, not as program**
  computational kernels only, thus moderate code complexity

- **Library calls are a DSL instructions, C fragments are "call-backs"**
  overcomes the usual problem of code between library calls

- **Establish side effect-freeness and parallelization opportunities**
  use OpenMP/OpenACC to express independent loops, variable visibility,…

- **Enables whole program optimization in domain specific compiler**
  data layout transformations, kernel merging, target novel accelerators,…

# Result: Performance Portability

- **Same source code, good performance across architectures**
  Intel Haswell, Intel Xeon PHI and Near Memory Accelerator

- **Parallel cross-kernel optimized library-based code is fast**
  40x speed-up over C baseline (PNNL TAV STAP benchmark)

- **Library-based code is pre-requisite for *Near Memory Accelerator***
  130x performance and 8,000x power efficiency gain on accelerator over C base line

**STAP Performance on Haswell**

Speedup [times]



Library-free C code

Obvious BLAS/ FFTW calls

DSL compiler gain

Full cross-kernel optimization

45 40 35 30 25 20 15 10 5 0

Small    Medium    Large

# Outline

- **Example: STAP**

- **Cross-call/cross library optimization with Spiral**

- **Library-based hardware acceleration**

- **Summary**

T. M. Low, Q. Guo, F. Franchetti: **Optimizing Space Time Adaptive Processing Through Accelerating Memory-Bounded Operations.** IEEE High Performance Extreme Computing Conference (HPEC), 2015.

# Space Time Adaptive Processing (STAP)

**dense data cube**



Input: 32 MB
Output: 128 MB
Millions of Math Ops

**SWAP requirement: high performance, low power**

# Main Computational Stages in STAP

| | |
|---|---|
| **Doppler Transforms** | **O(10k)**    FFTs + Data Layout Transforms |
| ↓ | |
| **Estimate Covariance Matrices** | **O(1M)**    Outer-products |
| ↓ | |
| **Compute Adaptive Weights** | **O(1M)**    Matrix Factorization Triangular Solve (vector) |
| ↓ | |
| **Normalize & Apply Adaptive Weights** | **O(10M)** Inner-products **O(100k)** Vector Scaling |

**Many memory-bounded operations**

Based on PNNL's Third-order Doppler STAP implementation (DARPA PERFECT STAP Benchmark)

# Background: FFTW 3

- **Latest version of FFTW**
  supports threading, SIMD, MPI

- **De-facto standard FFT library for HPC**
  distributed with Linux, installed everywhere

- **Autotuning + program generation**
  small hand-written core

- **Interface widely supported**
  Intel MKL, IBM ESSL, AMD ACML,…

```
#include <fftw3.h>
...
{
    fftw_complex *in, *out;
    fftw_plan p;
    ...
    in = (fftw_complex*) fftw_malloc(sizeof(fftw_complex) * N);
    out = (fftw_complex*) fftw_malloc(sizeof(fftw_complex) * N);
    p = fftw_plan_dft_1d(N, in, out, FFTW_FORWARD, FFTW_ESTIMATE);
    ...
    fftw_execute(p); /* repeat as needed */
    ...
    fftw_destroy_plan(p);
    fftw_free(in); fftw_free(out);}
```



Matteo Frigo and Steven G. Johnson: **The Design and Implementation of FFTW3.** Proceedings of the IEEE 93 (2), 216–231 (2005). Invited paper, Special Issue on Program Generation, Optimization, and Platform Adaptation.

# Doppler Transform: FFT + Corner Turn



1D FFTs

transposed

```
for (chan = 0; chan < N_CHAN; ++chan)
{
  for (range = 0; range < N_RANGE; ++range)
  {
    /* Apply a 1D FFT to convert to Doppler space */
    fft(
     (float *)&datacube_pulse_major_padded[chan][range][0],
     N_DOP,
     N_DOP_LOG2,
     FFT_FORWARD);
  }
}
/* Transpose from doppler major to range major*/
corner_turn_to_range_major_order
(doppler_datacube, datacube_pulse_major_padded);
```

```
/* Replace with fftw guru interface*/
plan = fftwf_plan_guru_dft(
            1, dims,
            2, howmany_dims,
            datacube_pulse_major_padded,
            doppler_datacube,
            FFTW_FORWARD,
            FFTW_WISDOM_ONLY);

fftwf_execute(plan);
fftwf_destroy_plan(plan);
```

**Combine loop of FFTs plus subsequent corner turn into one FFTW call**

# Use FFTW Guru Interface for Data Copy

**Zero pad plus batch transpose**



```
const fftwf_iodim howmany_dims[3] =
{{n:N_RANGE, is:1, os:N_DOP},
 {n:N_PULSES, is:N_RANGE, os:1},
 {n:N_CHAN,
  is:N_RANGE*N_PULSES,
  os:N_RANGE*N_DOP}};

plan =
fftwf_plan_guru_dft(0, NULL,
          3, howmany_dims,
          datacube,
          datacube_pulse_major_padded,
          FFTW_FORWARD,
          FFTW_ESTIMATE);

fftwf_execute(plan);
```

**Batch transpose**



```
const fftwf_iodim howmany_dims[3] =
{{n:N_DOP, is:1, os:N_RANGE},
 {n:N_RANGE, is:N_DOP, os:1},
 {n:N_CHAN,
  is:N_RANGE*N_DOP,
  os:N_RANGE*N_DOP}};

plan =
fftwf_plan_guru_dft(0, NULL,
          3, howmany_dims,
          doppler_major,
          range_major,
          FFTW_FORWARD,
          FFTW_ESTIMATE);

fftwf_execute(plan);
```

**FFTW Rank-0 FFT abstracts copy, transpose, gather, scatter**

# Borrowing FFTW's Guru Interface for BLAS

- ## FFTW guru interface
  strong support for batch FFTs and n-dimensional data cubes

```
typedef struct{
  int n; // size of the dimension
  int is; // stride for input
  int os; // stride for output
} fftw_iodim;
```

```
/* FFTW Guru Interface */
fftw_plan fftw_plan_guru_dft(
    int rank, const fftw_iodim *dims,
    int howmany_rank, const fftw_iodim *howmany_dims,
    fftw_complex *in, fftw_complex *out,
    int sign, unsigned flags);
```

- ## Standard BLAS interface
  Fortran 77 style interface

```
void cblas_cdotc_sub (const int N, const void * x, const int incx,
const void * y, const int incy, void * dotc);
```

- ## Generalization: BLAS guru interface
  borrow FFTW's data and batch representation for BLAS operations

```
typedef struct{
  int n;
  int i0s; // stride for input0
  int i1s; // stride for input1
} blas_indim;
```

```
/* BLAS Guru Configuration Interface */
blas_conf blas_conf_guru(
// (rank, dims[rank]) describes the basic BLAS operation size
    int rank,
    blas_indim *dims,
// (howmany_rank, howmany_dims[rank]) describes the "vector" size
    int howmany_rank,
    blas_indim *howmany_dims);
```

**Use BLAS guru interface in backend but do not expose to end user**

# Use OpenMP to Avoid Changing BLAS Interface

- **Absolutely cannot change BLAS**
  set-in-stone standard since 1979

- **Supported by everybody**
  MKL, ESSL, ACML, CUBLAS,…

- **Idea: Use OpenMP to express batch**
  mark loops as independent

- **Use compiler to extract batch BLAS descriptor**
  automatically translate BLAS + OpenMP to BLAS Guru Interface

```
#pragma omp for
for (sv = 0; sv < N_STEERING; ++sv) {
#pragma omp for
  for (block = 0; block < N_DOP*N_BLOCKS; ++block) {
    cblas_cdotc_sub(TDOF*N_CHAN,
      (float*)adaptive_weights[block][sv],
      1,
      (float*)steering_vectors[sv],
      1,
      (float*)&accums[block][sv]);
  }
}
```
                    **STAP rewritten with BLAS + OpenMP**

compiler →

```
blas_guru(
    1,
    {{n:TDOF*N_CHAN, i0s:1, i1s:1}}
    2,
    {{n:N_DOP*N_BLOCKS,
      i0s:N_STEERING*TDOF*N_CHAN,
      i1s:1},
     {n:N_STEERING,
      i0s:TDOF*N_CHAN,
      i1s:TDOF*N_CHAN}});
```

# After Transformation: STAP = 4 Library Calls

| | |
|---|---|
| **Doppler Transforms** | **FFTW call** |
| ↓ | |
| **Estimate Covariance Matrices** | **BLAS3 call** |
| ↓ | |
| **Compute Adaptive Weights** | **BLAS3 call** |
| ↓ | |
| **Normalize & Apply Adaptive Weights** | **BLAS2 call** |

**Result: STAP is expressed using four library calls + OpenMP directives**

# Outline

- **Example: STAP**

- **Cross-call/cross library optimization with Spiral**

- **Library-based hardware acceleration**

- **Summary**

# DSL Compiler for Global Transformations

- **Parse library calls and OpenMP**
  convert to Spiral's operator language (OL)

- **Spiral performs global optimizations**
  kernel merging, data layout transformations,…

- **Output CPU or accelerator code**
  run with Intel MKLK or our own accelerator

- **Utilize BLAS Guru interface**
  our accelerator implements BLAS Guru



M. Püschel, J. Moura, J. Johnson, D. Padua, M. Veloso, B. Singer, J. Xiong, F. Franchetti, A. Gacic, Y. Voronenko, K. Chen, R. W. Johnson, N. Rizzolo,
**"SPIRAL: Code Generation for DSP Transforms,"** Proceedings of the IEEE *Special Issue on Program Generation, Optimization, and Adaptation*, 2005

# What is Spiral?



*Traditionally*

*Spiral Approach*

*Spiral*

High performance library
optimized for given platform

*Comparable
performance*

High performance library
optimized for given platform

# Platform-Aware Formal Program Synthesis

**Model:** common abstraction
= spaces of matching formulas



Architectural parameter:
Vector length,
#processors, …

**optimization**

Kernel:
problem size,
algorithm choice

# Operators

## Definition

- **Operator: Multiple complex vectors → multiple complex vectors**
- **Higher-dimensional data is linearized**
- **Operators are potentially nonlinear**

$$\mathsf{M} : \begin{cases} \mathbb{C}^{n_0} \times \cdots \times \mathbb{C}^{n_{k-1}} \longrightarrow \mathbb{C}^{N_0} \times \cdots \times \mathbb{C}^{N_{\ell-1}} \\ (\mathbf{x}_0, \mathbf{x}_1, \ldots, \mathbf{x}_{k-1}) \mapsto \mathsf{M}(\mathbf{x}_0, \mathbf{x}_1, \ldots, \mathbf{x}_{k-1}) \end{cases}$$

## Example: Matrix-matrix-multiplication (MMM)

$$\mathrm{MMM}_{m,k,n} : \mathbb{R}^{mk} \times \mathbb{R}^{kn} \to \mathbb{R}^{mn}$$
$$(\mathbf{A}, \mathbf{B}) \mapsto \mathbf{A}\mathbf{B}$$

$k$

$m$ A

$n$

$k$ B

$\mathrm{MMM}_{m,k,n}$

$n$

C $m$

**Key to capture FFTs, numerical linear algebra, message passing in one framework**

# Operator Language

| name | definition |
|------|------------|
| *Linear, arity (1,1)* | |
| identity | $I_n : \mathbb{C}^n \to \mathbb{C}^n;\ \mathbf{x} \mapsto \mathbf{x}$ |
| vector flip | $J_n : \mathbb{C}^n \to \mathbb{C}^n;\ (x_i) \mapsto (x_{n-i})$ |
| transposition of an $m \times n$ matrix | $L_m^{mn} : \mathbb{C}^{mn} \to \mathbb{C}^{mn};\ \mathbf{A} \mapsto \mathbf{A}^T$ |
| matrix $M \in \mathbb{C}^{m \times n}$ | $\mathrm{M} : \mathbb{C}^n \to \mathbb{C}^m; \mathbf{x} \mapsto M\mathbf{x}$ |
| *Multilinear, arity (2,1)* | |
| Point-wise product | $\mathrm{P}_n : \mathbb{C}^n \times \mathbb{C}^n \to \mathbb{C}^n;\ ((x_i),(y_i)) \mapsto (x_i y_i)$ |
| Scalar product | $\mathrm{S}_n : \mathbb{C}^n \times \mathbb{C}^n \to \mathbb{C};\ ((x_i),(y_i)) \mapsto \Sigma(x_i y_i)$ |
| Kronecker product | $\mathrm{K}_{m \times n} : \mathbb{C}^m \times \mathbb{C}^n \to \mathbb{C}^{mn};\ ((x_i),\mathbf{y})) \mapsto (x_i \mathbf{y})$ |
| *Others* | |
| Fork | $\mathrm{Fork}_n : \mathbb{C}^n \to \mathbb{C}^n \times \mathbb{C}^n;\ \mathbf{x} \mapsto (\mathbf{x},\mathbf{x})$ |
| Split | $\mathrm{Split}_n : \mathbb{C}^n \to \mathbb{C}^{n/2} \times \mathbb{C}^{n/2};\ \mathbf{x} \mapsto (\mathbf{x}^U, \mathbf{x}^L)$ |
| Concatenate | $\oplus_n : \mathbb{C}^n \times \mathbb{C}^m \to \mathbb{C}^{n+m};\ (\mathbf{x},\mathbf{y}) \mapsto \mathbf{x} \oplus \mathbf{y}$ |
| Duplication | $\mathrm{dup}_n^m : \mathbb{C}^n \to \mathbb{C}^{nm};\ (\mathbf{x} \mapsto \mathbf{x} \otimes I_m$ |
| Min | $\min_n : \mathbb{C}^n \times \mathbb{C}^n \to \mathbb{C}^n;\ (\mathbf{x},\mathbf{y}) \mapsto (\min(x_i, y_i))$ |
| Max | $\max_n : \mathbb{C}^n \times \mathbb{C}^n \to \mathbb{C}^n;\ (\mathbf{x},\mathbf{y}) \mapsto (\max(x_i, y_i))$ |

# Some Application Domains in OL

## Linear Transforms

$$\mathbf{DFT}_n \rightarrow (\mathbf{DFT}_k \otimes I_m)\, \mathsf{T}_m^n (I_k \otimes \mathbf{DFT}_m)\, \mathsf{L}_k^n, \quad n = km$$

$$\mathbf{DFT}_n \rightarrow P_n(\mathbf{DFT}_k \otimes \mathbf{DFT}_m) Q_n, \quad n = km,\; \gcd(k,m) = 1$$

$$\mathbf{DFT}_p \rightarrow R_p^T (I_1 \oplus \mathbf{DFT}_{p-1}) D_p (I_1 \oplus \mathbf{DFT}_{p-1}) R_p, \quad p \text{ prime}$$

$$\mathbf{DCT\text{-}3}_n \rightarrow (I_m \oplus J_m)\, \mathsf{L}_m^n (\mathbf{DCT\text{-}3}_m(1/4) \oplus \mathbf{DCT\text{-}3}_m(3/4))$$

$$\cdot (\mathsf{F}_2 \otimes I_m) \begin{bmatrix} I_m & 0 \oplus -J_{m-1} \\ & \frac{1}{\sqrt{2}}(I_1 \oplus 2 I_m) \end{bmatrix}, \quad n = 2m$$
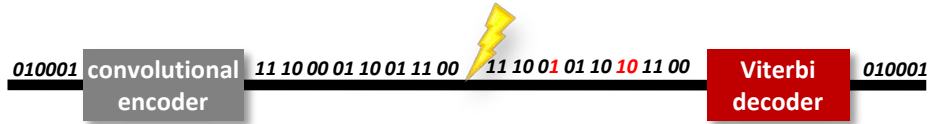
$$\mathbf{DCT\text{-}4}_n \rightarrow S_n \mathbf{DCT\text{-}2}_n \operatorname{diag}_{0 \le k < n}(1/(2\cos((2k+1)\pi/4n)))$$

$$\mathbf{IMDCT}_{2m} \rightarrow (J_m \oplus I_m \oplus I_m \oplus J_m)\left(\left(\begin{bmatrix} 1 \\ -1 \end{bmatrix} \otimes I_m\right) \oplus \left(\begin{bmatrix} -1 \\ -1 \end{bmatrix} \otimes I_m\right)\right) J_{2m}\, \mathbf{DCT\text{-}4}_{2m}$$

$$\mathbf{WHT}_{2^k} \rightarrow \prod_{i=1}^{t} (I_{2^{k_1+\cdots+k_{i-1}}} \otimes \mathbf{WHT}_{2^{k_i}} \otimes I_{2^{k_{i+1}+\cdots+k_t}}), \quad k = k_1 + \cdots + k_t$$

$$\mathbf{DFT}_2 \rightarrow \mathsf{F}_2$$

$$\mathbf{DCT\text{-}2}_2 \rightarrow \operatorname{diag}(1, 1/\sqrt{2})\, \mathsf{F}_2$$

$$\mathbf{DCT\text{-}4}_2 \rightarrow J_2\, \mathsf{R}_{13\pi/8}$$

## Software Defined Radio



$$F_{K,F} \rightarrow \prod_{i=1}^{F} \left( (I_{2^{K-2}} \otimes_j B_{F-i,j}) L_{2^{K-2}}^{2^{K-1}} \right)$$

$$\underline{\mathbf{F}}_{K,F}\; \nu \rightarrow \prod_{i=1}^{F} \left( \left( I_{2^{K-2}/\nu} \otimes_{j_1} \vec{L}_\nu^{2\nu} \vec{B}_{F-i,j_1}^\nu \right) (L_{2^{K-2}/\nu}^{2^{K-1}/\nu} \overline{\otimes} I_\nu) \right)$$

$$B_{i,j} : \begin{cases} \pi_U = \min_{d_U}(\pi_A + \beta_{A \to U}, \pi_B + \beta_{B \to U}) \\ \pi_V = \min_{d_V}(\pi_A + \beta_{A \to V}, \pi_B + \beta_{B \to V}) \end{cases}$$

## Matrix-Matrix Multiplication



$$\mathrm{MMM}_{1,1,1} \rightarrow (\cdot)_1$$

$$\mathrm{MMM}_{m,n,k} \rightarrow (\otimes)_{m/m_b \times 1} \otimes \mathrm{MMM}_{m_b,n,k}$$

$$\mathrm{MMM}_{m,n,k} \rightarrow \mathrm{MMM}_{m,nb,k} \otimes (\otimes)_{1 \times n/nb}$$

$$\mathrm{MMM}_{m,n,k} \rightarrow ((\Sigma_{k/k_b} \circ (\cdot)_{k/k_b}) \otimes \mathrm{MMM}_{m,n,k_b}) \circ$$
$$((L_{k/k_b}^{mk/k_b} \otimes I_{k_b}) \times I_{kn})$$

$$\mathrm{MMM}_{m,n,k} \rightarrow (L_m^{mn/n_b} \otimes I_{n_b}) \circ$$
$$((\otimes)_{1 \times n/n_b} \otimes \mathrm{MMM}_{m,n_b,k}) \circ$$
$$(I_{km} \times (L_{n/n_b}^{kn/n_b} \otimes I_{n_b}))$$

## Synthetic Aperture Radar (SAR)



$$\mathrm{SAR}_{k \times m \to n \times n} \rightarrow \mathrm{DFT}_{n \times n} \circ \mathrm{Interp}_{k \times m \to n \times n}$$

$$\mathrm{DFT}_{n \times n} \rightarrow (\mathrm{DFT}_n \otimes I_n) \circ (I_n \otimes \mathrm{DFT}_n)$$

$$\mathrm{Interp}_{k \times m \to n \times n} \rightarrow (\mathrm{Interp}_{k \to n} \otimes_i I_n) \circ (I_k \otimes_i \mathrm{Interp}_{m \to n})$$

$$\mathrm{Interp}_{r \to s} \rightarrow \left( \bigoplus_{i=0}^{n-2} \mathrm{InterpSeg}_k \right) \oplus \mathrm{InterpSegPruned}_{k,\ell}$$

$$\mathrm{InterpSeg}_k \rightarrow \mathsf{G}_f^{u \cdot n \to k} \circ \mathrm{iPrunedDFT}_{n \to u \cdot n} \circ \left(\frac{1}{n}\right) \circ \mathrm{DFT}_n$$

# Formal Approach for all Types of Parallelism

- **Multithreading (Multicore)**

  $I_p \otimes_\parallel A_{\mu n}, \quad \mathbf{L}_m^{mn} \bar{\bar{\otimes}} I_\mu$

- **Vector SIMD (SSE, VMX/Altivec,…)**

  $A \hat{\otimes} I_\nu \quad \underbrace{\mathbf{L}_2^{2\nu}}_{\text{isa}}, \quad \underbrace{\mathbf{L}_\nu^{2\nu}}_{\text{isa}}, \quad \underbrace{\mathbf{L}_\nu^{\nu^2}}_{\text{isa}}$

- **Message Passing (Clusters, MPP)**

  $I_p \otimes_\parallel A_n, \quad \underbrace{\mathbf{L}_p^{p^2} \bar{\otimes} I_{n/p^2}}_{\text{all-to-all}}$

- **Streaming/multibuffering (Cell)**

  $I_n \otimes_2 A_{\mu n}, \quad \mathbf{L}_m^{mn} \bar{\bar{\otimes}} I_\mu$

- **Graphics Processors (GPUs)**

  $\prod_{i=0}^{n-1} A_i, \quad A_n \hat{\otimes} I_w, \quad P_n \otimes Q_w$

- **Gate-level parallelism (FPGA)**

  $\prod_{i=0}^{n-1}{}^{\text{ir}} A, \quad I_s \tilde{\otimes} A, \quad \underbrace{\mathbf{L}_n^m}_{\text{bram}}$

- **HW/SW partitioning (CPU + FPGA)**

  $\underbrace{A_1}_{\text{fpga}}, \quad \underbrace{A_2}_{\text{fpga}}, \quad \underbrace{A_3}_{\text{fpga}}, \quad \underbrace{A_4}_{\text{fpga}}$

# Autotuning in Constraint Solution Space

**Intel Core i7 (2ⁿᵈ Gen)**

**DFT₂₅₆**

**OL specification**

**Expansion + backtracking**

**OL (dataflow) expression**

### Base cases

$$I_4 \otimes_{\parallel} A_{16n}$$
$$L_m^{mn} \bar{\otimes} I_{16}$$
$$A \bar{\otimes} I_4$$
$$\underbrace{L_2^8}_{\text{SSE}}, \underbrace{L_4^8}_{\text{SSE}}, \underbrace{L_4^{16}}_{\text{SSE}}, \underbrace{L_2^4 \otimes I_2}_{\text{SSE}}$$

### Transformation rules

$$\frac{AB}{\mathrm{smp}(p,\mu)} \to \frac{A}{\mathrm{smp}(p,\mu)} \frac{B}{\mathrm{smp}(p,\mu)}$$

$$\frac{L_m^{mn}}{\mathrm{smp}(p,\mu)} \to \begin{cases} \frac{\left(I_p \otimes L_{m/p}^{mn/p}\right)}{\mathrm{smp}(p,\mu)} \frac{\left(L_p^{pn} \otimes I_{m/p}\right)}{\mathrm{smp}(p,\mu)} \\ \frac{\left(L_m^{pm} \otimes I_{n/p}\right)}{\mathrm{smp}(p,\mu)} \frac{\left(I_p \otimes L_m^{mn/p}\right)}{\mathrm{smp}(p,\mu)} \end{cases}$$

$$\frac{I_m \otimes A_n}{\mathrm{smp}(p,\mu)} \to I_p \otimes_{\parallel} \left(I_{m/p} \otimes A_n\right)$$

$$\dots$$

### Breakdown rules

$$\mathrm{DFT}_n \to \left(\mathrm{DFT}_k \otimes I_m\right) T_m^n \cdot \left(I_k \otimes \mathrm{DFT}_m\right) L_k^n$$

$$\mathrm{DFT}_n \to P_n(\mathrm{DFT}_k \otimes \mathrm{DFT}_m)Q_n$$

$$\mathrm{DFT}_p \to R_p^T(I_1 \oplus \mathrm{DFT}_{p-1})D_p \cdot (I_1 \oplus \mathrm{DFT}_{p-1})R_p$$

$$\mathrm{DFT}_2 \to \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$$

**Recursive descent**

**Σ-OL (loop) expression**

**Confluent term rewriting**

**Optimized Σ-OL expression**

**Recursive descent**

**Abstract code**

**Confluent term rewriting**

**Optimized abstract code**

**Recursive descent**

**C code**

$$\left(\left(L_m^{mp} \otimes I_{n/p\mu}\right) \bar{\otimes} I_\mu\right)\left(I_p \otimes_{\parallel} (\mathrm{DFT}_m \otimes I_{n/p})\right)\left(\left(L_p^{mp} \otimes I_{n/p\mu}\right) \bar{\otimes} I_\mu\right)\left(\bigoplus_{i=0}^{p-1}{}_{\parallel} T_n^{mn,i}\right)\left(I_p \otimes_{\parallel}(I_{m/p} \otimes \mathrm{DFT}_n)\right)\left(I_p \otimes_{\parallel} L_{m/p}^{mn/p}\right)\left(\left(L_p^{pn} \otimes I_{m/p\mu}\right) \bar{\otimes} I_\mu\right)$$

Electrical *&* Computer
**ENGINEERING**

# Translating an OL Expression Into Code

**Constraint Solver Input:**

$$\underbrace{\text{DFT}_8}_{\text{double}}$$

**Output =**

**OL Expression:** $(\text{DFT}_2 \otimes I_4)\, T_4^8 \left( I_2 \otimes \left( (\text{DFT}_2 \otimes I_2)\, T_2^4\, (I_2 \otimes \text{DFT}_2)\, L_2^4 \right) \right) L_2^8$

**Σ-OL:**

$$\sum_{j=0}^{3} \left( S_j\, \text{DFT}_2\, G_j \right) \sum_{k=0}^{1} \left( \sum_{l=0}^{1} \left( S_{k,l}\, \text{diag}\!\left( t_{k,l} \right) \text{DFT}_2\, G_l \right) \sum_{m=0}^{1} \left( S_m\, \text{diag}\!\left( t_m \right) \text{DFT}_2\, G_{k,m} \right) \right)$$

**C Code:**

```
void sub(double *y, double *x) {
  double f0, f1, f2, f3, f4, f7, f8, f10, f11;
  f0 = x[0] - x[3];
  f1 = x[0] + x[3];
  f2 = x[1] - x[2];
  f3 = x[1] + x[2];
  f4 = f1 - f3;
  y[0] = f1 + f3;
  y[2] = 0.7071067811865476 * f4;
  f7 = 0.9238795325112867 * f0;
  f8 = 0.3826834323650898 * f2;
  y[1] = f7 + f8;
  f10 = 0.3826834323650898 * f0;
  f11 = (-0.9238795325112867) * f2;
  y[3] = f10 + f11;
}
```

**OL specification**

Expansion + backtracking

**OL (dataflow) expression**

Recursive descent

**Σ-OL (loop) expression**

Confluent term rewriting

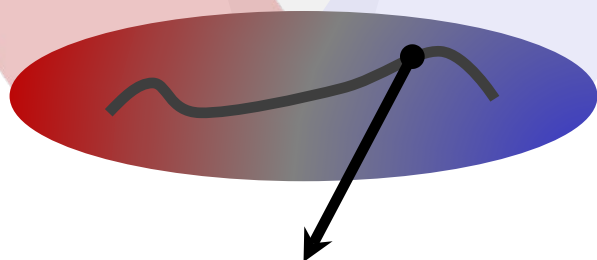**Optimized Σ-OL expression**

Recursive descent

**Abstract code**

Confluent term rewriting

**Optimized abstract code**

Recursive descent

**C code**

# STAP in SPIRAL's OL specification

$$(I_D \otimes I_B \otimes (I_V \otimes I_S \otimes (Dot \circ Extract))) \circ (I_V \otimes \frac{1}{\| Dot^{CF} \|_2})) \circ$$ **Inner Products**

$$(I_D \otimes I_B \otimes I_V \otimes (Trsv_{bwd}^{CF} \circ Trsv_{fwd}^{CF})) \circ$$

**Linear System Solver**

$$(I_D \otimes I_B \otimes Chol^{CF}) \circ$$

$$(I_B \otimes I_D \otimes (\frac{1}{S} \times (I_S \otimes (Her^{CF} \circ Extract)))) \circ$$ **Covariance Estimate**

$$(I_C \otimes L_R^{RD}) \circ$$

**FFT + Corner Turn**

$$(I_C \otimes I_R \otimes DFT_D)$$

# STAP: Necessary Optimization

$$(I_D \otimes I_B \otimes (I_V \otimes I_S \otimes (Dot \circ Extract)) \circ (I_V \otimes \frac{1}{\| Dot^{CF} \|_2})) \circ$$

$$(I_D \otimes I_B \otimes I_V \otimes (Trsv_{bwd}^{CF} \circ Trsv_{fwd}^{CF})) \circ$$

$$(I_D \otimes I_B \otimes Chol^{CF}) \circ$$

$$(I_B \otimes I_D \otimes (\frac{1}{S} \times (I_S \otimes (Her^{CF} \circ Extract)))) \circ$$

$$(I_C \otimes L_R^{RD}) \circ$$

$$(I_C \otimes I_R \otimes DFT_D)$$

**BLAS2 operations to be optimized into BLAS3 operations**

# Problem: BLAS/LAPACK interface

■ **Merged Operation**



**for i = 0 to 3**
**cblas_cher(…)**

Equivalent operation

**cblas_cherk(…)**

■ **Batched Operation**



**for i = 0 to 3**
**cblas_cdotu_sub(…)**

No Equivalent operation

**?**

**blas_guru(…)**

# STAP: SPIRAL Optimization

**BLAS 2 operation**

$$(I_B \otimes I_D \otimes (\frac{1}{S} \times (I_S \otimes (Her^{CF} \circ Extract))))$$

**Data stored column-major instead of row-major as required**

**After Spiral's transformations**

$$(I_{DB} \otimes (\frac{1}{S} \times (Herk^{CF,S} \circ I_S \otimes (L_{CF}^{SCF} \circ Extract))))$$

**Loop interchange + Loop coalescing**

**BLAS2 to BLAS3 call**

**Loop fission**

**Transposition from column to row major**

# Interpreting the OL Specification

**After SPIRAL's transformations**

$$(I_{DB} \otimes (I_V \otimes (I_S \otimes Dot \circ I_S \otimes Extract)) \circ (I_V \otimes \frac{1}{\| Dot^{CF} \|_2})) \circ$$ **Batched BLAS1**

$$(I_{DB} \otimes (Trsm^{CF}_{bwd,V} \circ Trsm^{CF}_{fwd,V} \circ (I_V \otimes L^{CF}_V))) \circ$$

$$(I_{DB} \otimes Chol^{CF}) \circ$$

**Loop fission**

$$(I_{DB} \otimes (\frac{1}{S} \times (Herk^{CF,S} \circ I_S \otimes (L^{SCF}_{CF} \circ Extract)))) \circ$$

$$(I_C \otimes (L^{RD}_R \circ (I_R \otimes DFT_D)))$$

**Transposition from column to row major**

**Loop interchange +** **BLAS2 to BLAS3 call**
**Loop coalescing**

# OL vs. FFTW Guru Inteface

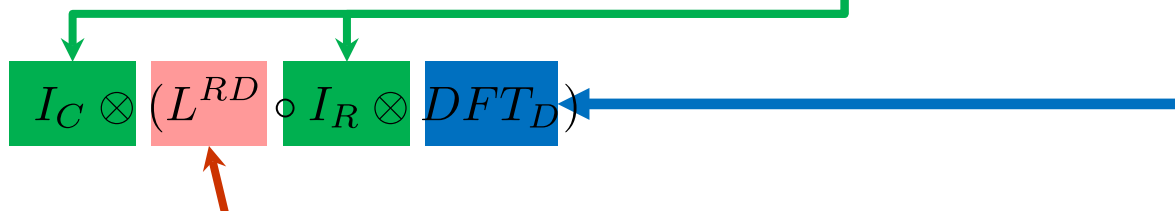- **FFTW guru-interface**
  Batched descriptor for FFT computation

```
fftwf_plan_guru_dft(
    1, dims,
    2, howmany_dims,
    datacube_pulse_major_padded,
    doppler_datacube,
    FFTW_FORWARD,FFTW_WISDOM_ONLY);
```
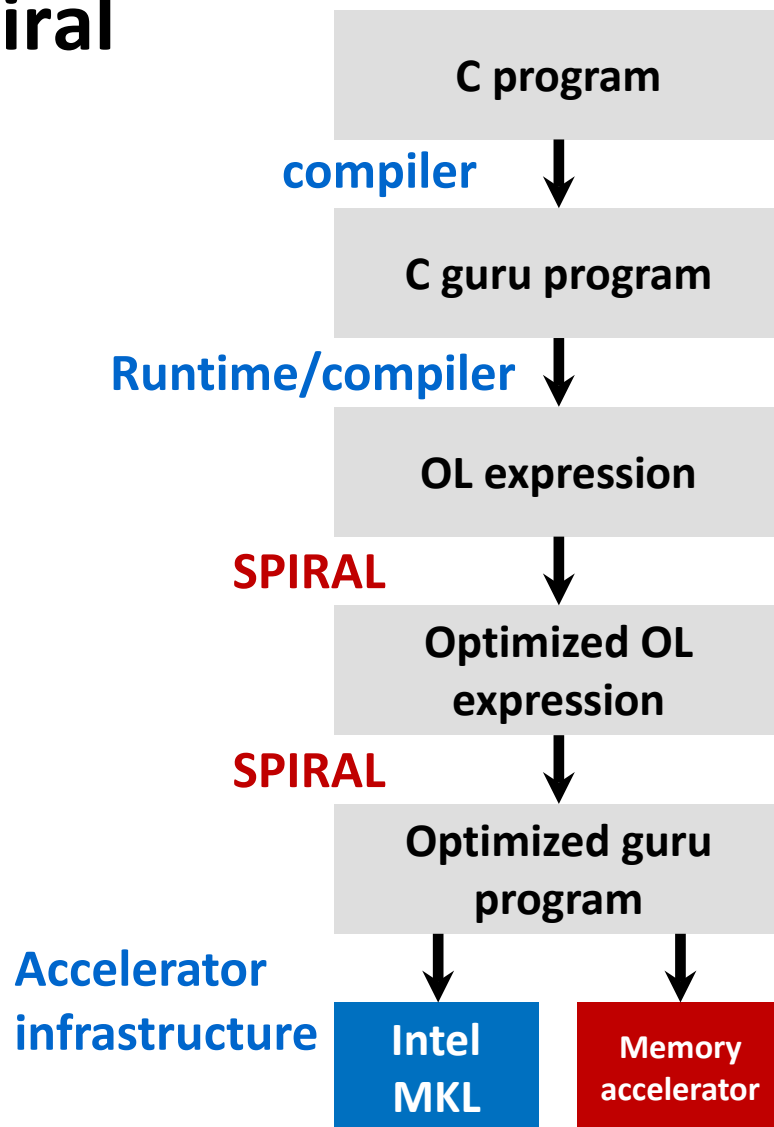
**Compute descriptor**

**Batch descriptor**

- **OL representation of FFTW guru interface**

$$I_C \otimes (L^{RD} \circ I_R \otimes DFT_D)$$

**Implicit Reshape in Compute Descriptor**

# Optimizing STAP with Spiral

- **Input: C program with BLAS and FFTW calls**

- **Spiral Input: STAP C program only using FFTW and BLAS guru calls**

- **This is equivalent to a OL formula in Spiral**

- **Enables full-program data layout optimization and kernel merging**

- **Final program: efficient on CPU, can target novel accelerators**

C program

**compiler** ↓

C guru program

**Runtime/compiler** ↓

OL expression

**SPIRAL** ↓

Optimized OL expression

**SPIRAL** ↓

Optimized guru program

↓ ↓

**Accelerator infrastructure**

Intel MKL   Memory accelerator

**This paves the way for our 3DIC memory side accelerator (discussed next)**

# Outline

- **Example: STAP**

- **Cross-call/cross library optimization with Spiral**

- **Library-based hardware acceleration**

- **Summary**

Q. Guo, T.-M. Low, N. Alachiotis, B. Akin, L. Pileggi, J. C. Hoe, and F. Franchetti: **Enabling Portable Energy Efficiency with Memory Accelerated Library.** IEEE/ACM MICRO-48, 2015, *to appear*.

# Overview: Memory-Side Accelerators

- **Accelerator behind DRAM interface**
  - No off-DIMM data traffic and SERDES
  - Huge problem sizes possible
  - 3D stacking is enabling technology
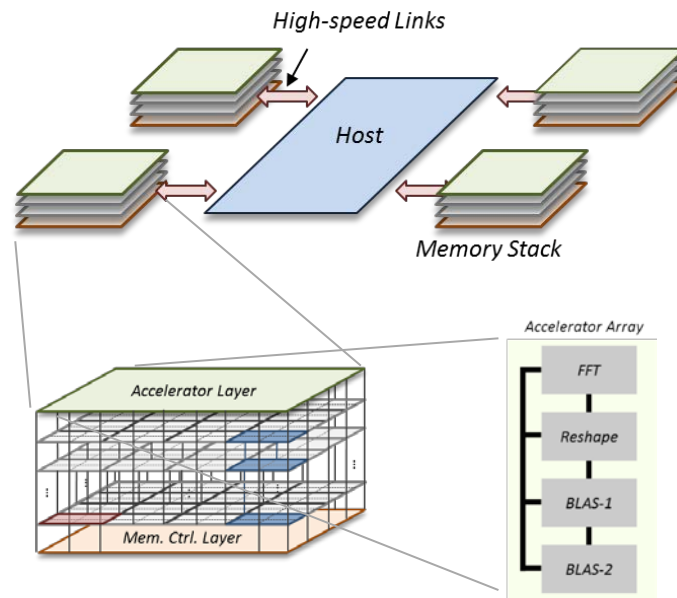
- **Configurable array of accelerators**
  - Domain specific, highly configurable
  - Cover DoD-relevant kernels

- **System CPU**
  - Standard: Multicore CPU+GPU

- **Software stack**
  - User level: standard numeric libraries BLAS1/2, FFTW, SpMV/SpGEMM,…
  - OS, driver, compiler, runtime system support necessary for drop-in replacement

# Accelerator Hardware Architecture
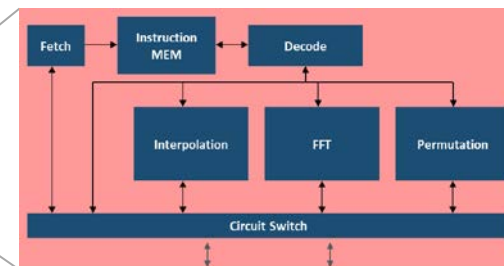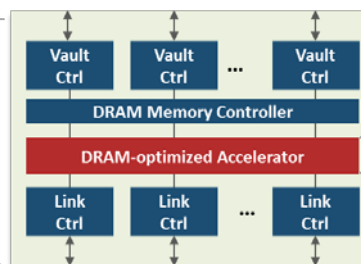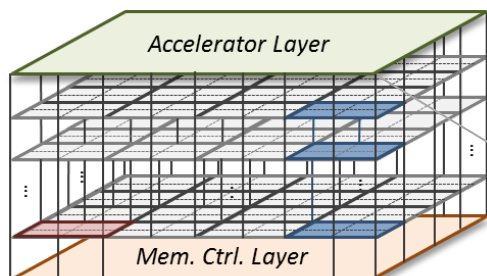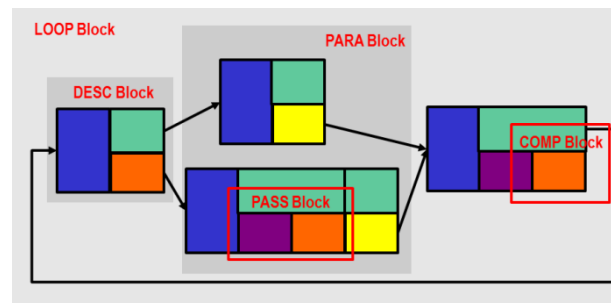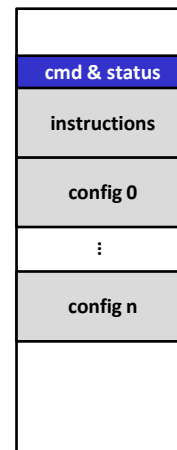
- **Accelerator tiles**
  - One accelerator tile per 3D-DRAM vault
  - Can be linked to one big accelerator
  - Configurable/programmable:
    units, interconnect, iterators

- **Programming model**
  - Memory mapped device
  - Physical memory addressing
  - Data and command segment in DRAM
  - Controlled via read/writes to
    command addresses

```
LOOP
 numLoop = 10;
 DESC
  descID = 0;
  numPass = 1;
  PASS
   passID = 0;
   numComp = 1;
   COMP
    compID = 0;
    Inst = FFT;
    Conf = "fft.conf"
   ENDCOMP
  ENDPASS
 ENDDESC
 ...
ENDLOOP
```

# Accelerator Software Architecture

- **Accelerator software abstraction**
  - memory mapped device
  - command and data address space
  - kernel/virtual memory configuration
- **Low level user API**
  - device driver interface
  - low level C configuration library
- **Application level API**
  - transparent device-aware malloc/free
  - standard low intensity math libraries: BLAS-1, BLAS-2, FFTW, sparse mat/vec, copy/reshape, corner turn
  - "guru" math library interfaces
  - OpenACC directives + proto compiler

```c
// TAV STAP example
//
// kernel reimplemented with BLAS and OpenACC
//

float *adaptive_weights;
float *steering_vectors;

#pragma declare \
    device_resident(adaptive_weights);
#pragma declare \
    device_resident(steering_vectors);

adaptive_weights = XMALLOC(sizeof(complex)*
    num_adaptive_weight_elements);
steering_vectors = XMALLOC(sizeof(complex)*
    num_steering_vector_elements);

...
#pragma acc data copyin(adaptive_weights, \
    steering_vectors), copyout(accums)
#pragma acc kernels loop
for (sv = 0; sv < N_STEERING; ++sv)
{
  for (block = 0; block < N_DOP*N_BLOCKS; ++block)
  {
    accum.re = accum.im = 0.0f;
    cblas_cdotc_sub(TDOF*N_CHAN,
        (float*)adaptive_weights[block][sv],
        1,
        (float*)steering_vectors[sv],
        1,
        (float*)&accums[block][sv]);
  }
}
...
```

# Hardware in the Loop Simulation

- **Accelerator Simulation**
  - **Timing:** Synopsis DesignWare
  - **Power:** DRAMSim2, Cacti, Cacti-3DD, McPAT

- **Full System Evaluation**
  - Run code on real system (Haswell, Xeon Phi) or in simulator (SimpleScalar,…)
  - Normal DRAM access for CPU, but trap accelerator command memory space, invoke simulator

```
#include "fftw.h"
#include "papi_util.h"
...
{
  fftw_complex *in, *out;
  fftw_plan p;
  ...
  in = (fftw_complex*)
    fftw_malloc(sizeof(fftw_complex) * N * N);

  out = (fftw_complex*)
    fftw_malloc(sizeof(fftw_complex) * N * N);

  p = fftw_plan_dft_2d(N, N, in, out,
    FFTW_FORWARD, FFTW_MEASURE);

  ///// Measurement Code /////
  papi_begin();
  ///////////////////////////

  fftw_execute(p);

  ///// Measurement Code /////
  papi_end();
  ///////////////////////////

  fftw_destroy_plan(p);
  fftw_free(in); fftw_free(out);
}
```
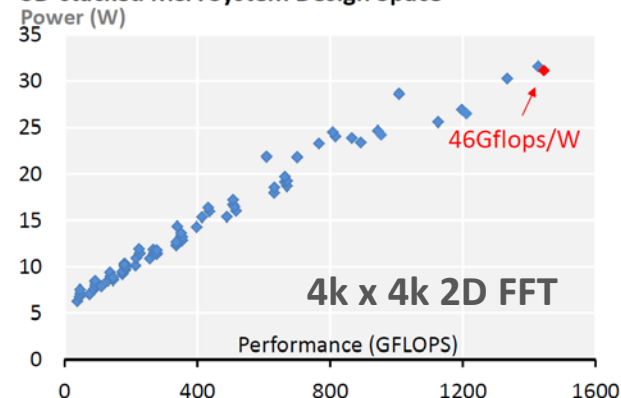


3D-stacked MSA System Design Space

46Gflops/W

4k x 4k 2D FFT

Electrical & Computer
ENGINEERING

# Accelerating the TAV STAP Benchmark

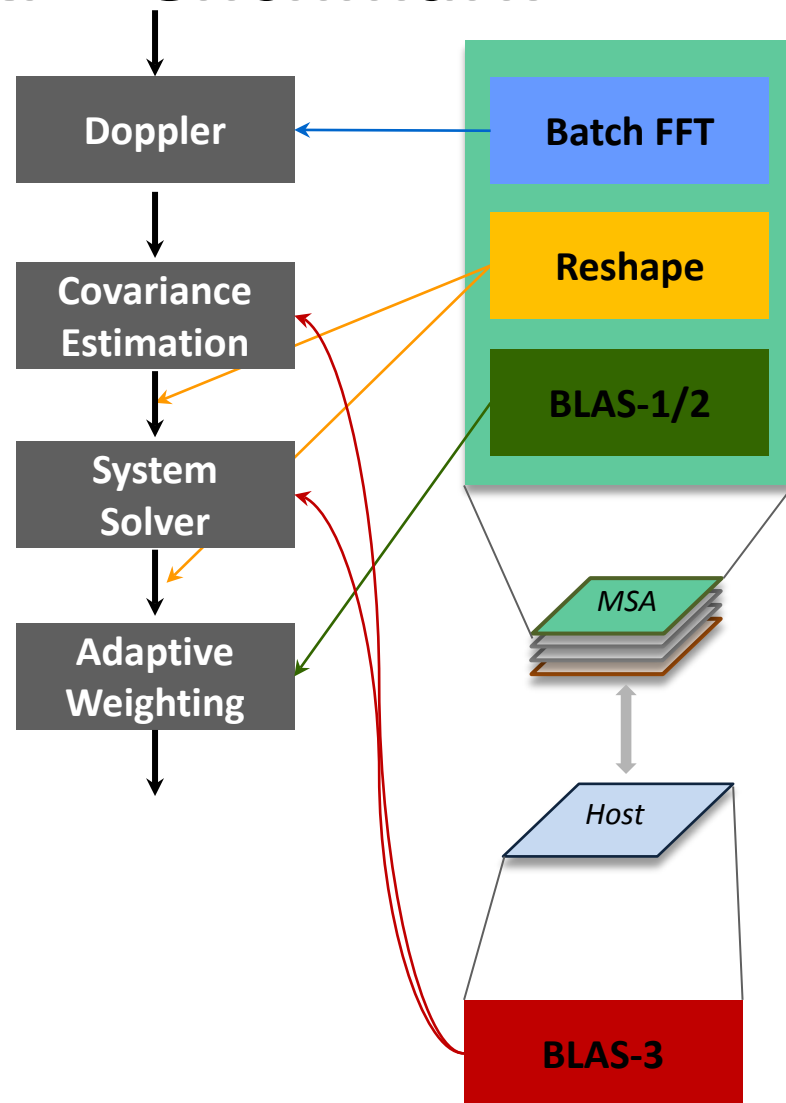■ **STAP structure: 4 phases**

  ▪ **Doppler:** batch FFT

  ▪ **Covariance estimation:** outer product

  ▪ **System solver:** Cholesky or QR

  ▪ **Adaptive weighting:** inner product

■ **STAP memory-side acceleration**

  ▪ **On CPU:** BLAS-3/Cholesky

  ▪ **On accelerator:** FFT, BLAS-1, BLAS-2

■ **Software interface**

  ▪ Drop-in replacement for Intel MKL:
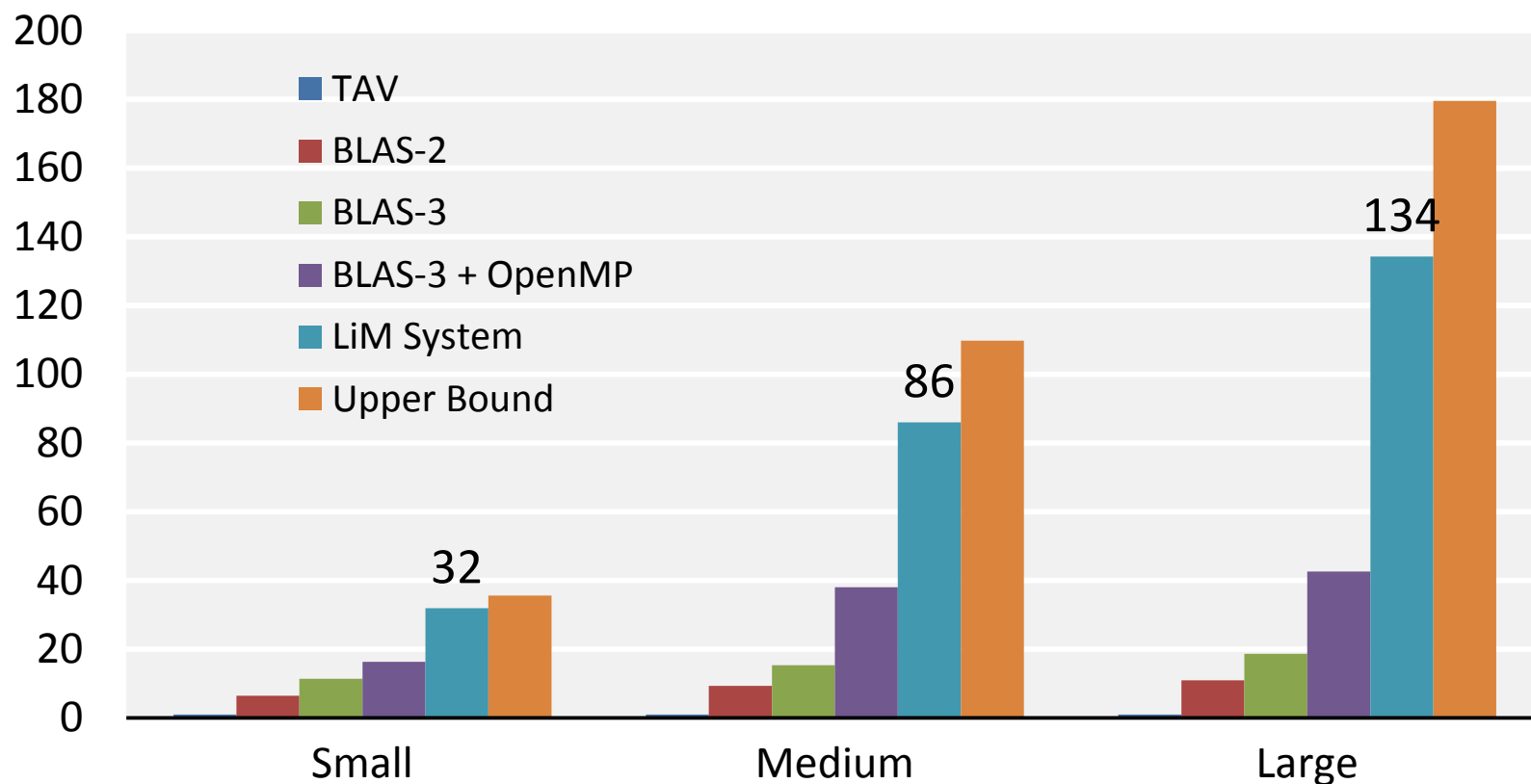    FFTW, BLAS-1/2/3

  ▪ OpenACC pragmas for batch BLAS-1/2

*No DMA/copy-in/copy-out necessary for accelerator*

Doppler

Covariance
Estimation

System
Solver

Adaptive
Weighting

Batch FFT

Reshape

BLAS-1/2

MSA

Host

BLAS-3

# STAP Performance Results

**Performance Gains over Baseline (TAV)**

Speedup (Times)



Legend:
- TAV
- BLAS-2
- BLAS-3
- BLAS-3 + OpenMP
- LiM System
- Upper Bound

Small: 32

Medium: 86

Large: 134

# STAP Power Efficiency Results

**GFLOPS/W Gain over Baseline (TAV)**

**EDP Reduction (Times)**



Legend: TAV, BLAS-2, BLAS-3, BLAS-3 + OpenMP, LiM System, Upper Bound

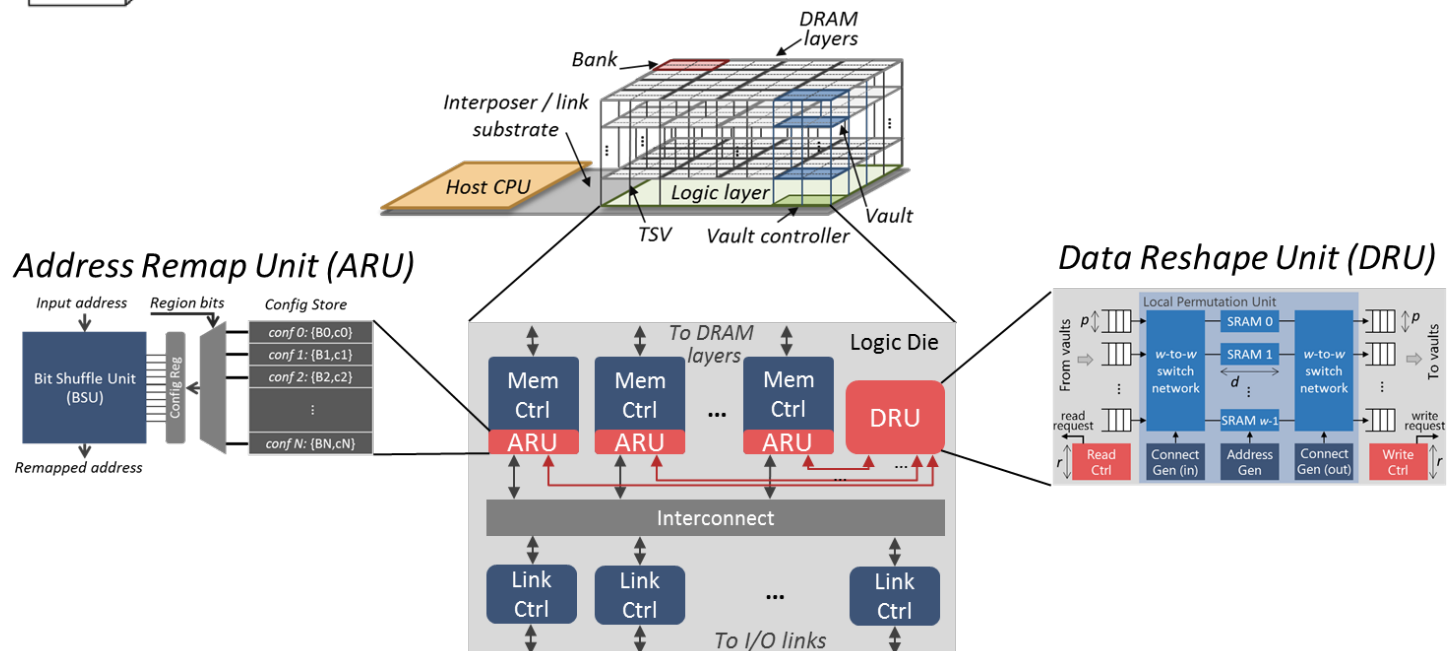Small: 595
Medium: 4,548
Large: 7,993

# Memory Layout Accelerator



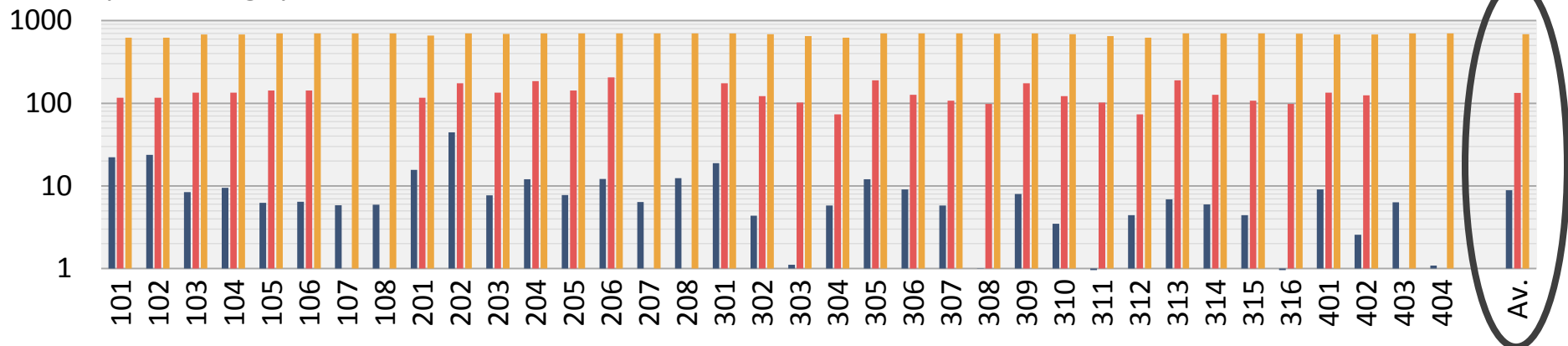**Intel MKL (Math Kernel Library) functions:**

```
void   mkl_simatcopy (const char ordering, char trans, size_t rows,size_t cols,
                          const float alpha, float * AB, size_t lda, size_t ldb);
void   mkl_somatcopy (char ordering, char trans, size_t rows, size_t cols,
                  const float alpha, float * A, size_t lda, float * B, size_t ldb);
void   pstrmr2d (char *uplo , char *diag , MKL_INT *m , MKL_INT *n , float *a,
                  , …, MKL_INT *jb , MKL_INT *descb , MKL_INT*ictxt ); // ScaLAPACK
vsPackI ( n, a, inca, y );
vsUnpackI ( n, a, y, incy );
void   cblas_scopy (const MKL_INT n, const float *x, const MKL_INT incx, …);
void   cblas_sswap (const MKL_INT n, float *x, const MKL_INT incx, float *y,…);
pslared2d (n, ia, ja, desc, byrow, byall, work, lwork)
pslaswp (direc, rowcol, n, a, ia, ja, desca, k1, k2, ipiv)
...
```
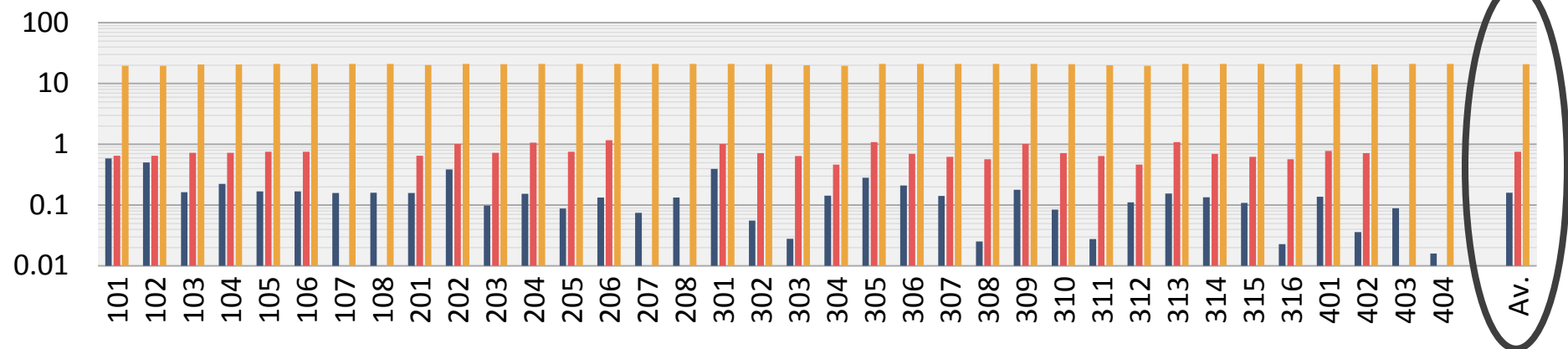


*Address Remap Unit (ARU)*

*Data Reshape Unit (DRU)*

B. Akin, F. Franchetti, J. C. Hoe: **Data Reorganization in Memory Using 3D-stacked DRAM,**
42nd International Symposium on Computer Architecture (ISCA), 2015.

Electrical & Computer
ENGINEERING

# Accelerating MKL Reorganization Routines

■ CPU (i7-4770K)    ■ GPU (GTX 780)    ■ DRU (MH)
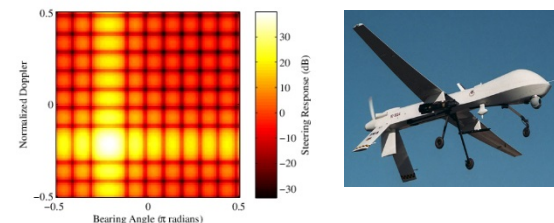


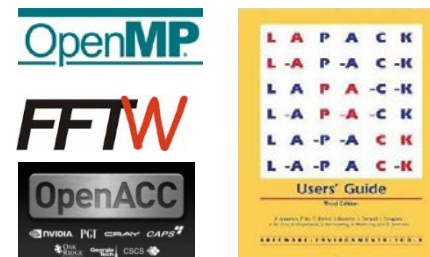Reshape Throughput [GB/s]



Energy Efficiency [GB/J]

# Outline

- **Example: STAP**

- **Cross-call/cross library optimization with Spiral**

- **Library-based hardware acceleration**

- **Summary**

# Summary

- **Target: dense scientific and HPEC kernels**
  **math heavy, regular, good library coverage**

- **HPC Library + OpenMP + C as DSL**
  **View program as specification**

- **Enormous efficiency gains are possible**
  **Novel accelerators + Spiral optimization**