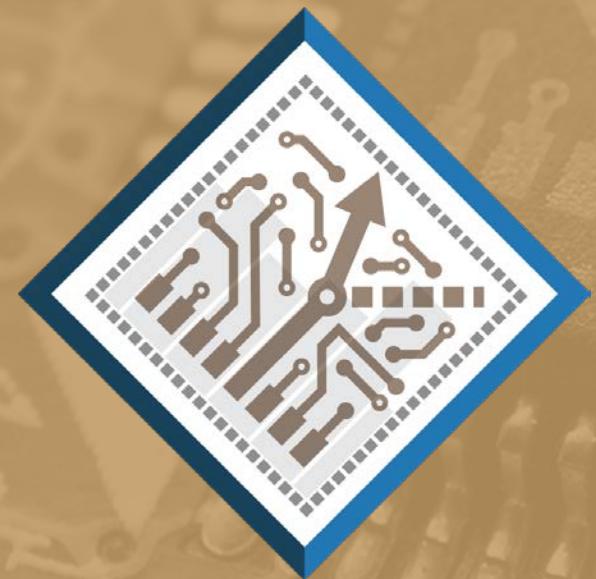




SAMAN AMARASINGHE

**PROFESSOR & ASSOC. DEPT. HEAD
DEPT. OF EECS, MIT**



DOMAIN SPECIFIC LANGUAGES for DOMAIN SPECIFIC ARCHITECTURES

This material is based on research sponsored by the Defense Advanced Research Projects Agency (DARPA) under agreement HR0011-18-3-0007; the Application Driving Architectures (ADA) Research Center, a JUMP Center co-sponsored by SRC and DARPA; the Toyota Research Institute; the U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing Research under Award Numbers DE-SC0008923 and DE-SC0018121; and the National Science Foundation under Grant No. CCF-1533753;

DISCLAIMER: The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the U.S. Government.

"A GOLDEN AGE FOR COMPUTER ARCHITECTURE" TURING AWARD LECTURE BY HENNESSY AND PATTERSON, 2018

- SW-
 - Modularity
 - Efficiency
- HW-
 - Off-chip
 - Jitter
- Comp-
 - D

Why Tailor

- Achieve characteristics
 - Not one-size-fits-all
 - Different domains
- Requirements processing
- Examples
 - Neural networks
 - GPUs
 - Programs

DSAs require architectural support:

- Hard to structure
- Need domain knowledge
- Domains
 - IEEE numbers
 - 32-64 bits
- Domain-specific

- ## Research Opportunities
- General-purpose applications:
 - Make Python run like C with compiler + HW
 - Deja vu: make HLLs fast on RISC
 - Domain-specific applications (bigger opportunity?):
 - What are the right DSLs for important applications?
 - Codesign of new DSLs and DSAs
 - Advanced compilation techniques for optimizing the matching:
 - New territory: not extraction of high level structure from C/Fortran but matching/optimization
 - Challenge: not to compromise DSLs with short-term ISA-specific or microarchitectural-specific compromises

36

From the Turing Award Lecture by Hennessy and Patterson

WHAT IS A DOMAIN SPECIFIC LANGUAGE?

- A Domain Specific Language capture expert knowledge about an application domain
- DSLs provide a high-level programming model for domain scientists
 - Capture the programmer intent at a higher level of abstraction
- DSLs provide the compiler more opportunities for optimization
 - Can encode domain expert knowledge of specific optimizations
 - Better view of the intended computation without heroic analysis
 - Less low-level decisions by the programmer that has to be undone
 - Able to generate high-performance code for multiple hardware platforms



HALIDE

A DSL for Image Processing

DISTRIBUTION STATEMENT A. Approved for public release

Image from the Smithsonian
National Museum of Natural History

A SIMPLE EXAMPLE: 3X3 BLUR

```
void box_filter_3x3(const Image &in, Image &blur) {  
    Image blurx(in.width(), in.height()); // allocate blurx array  
  
    for (int y = 0; y < in.height(); y++)  
        for (int x = 0; x < in.height(); x++)  
            blurx(x, y) = (in(x-1, y) + in(x, y) + in(x+1, y))/3;  
  
    for (int x = 0; x < in.width(); x++)  
        for (int y = 0; y < in.height(); y++)  
            blur(y, x) = (blurx(x, y-1) + blurx(x, y) + blurx(x, y+1))/3;  
}
```

HAND-OPTIMIZED C++

```
void box_filter_3x3(const Image &in, Image &blury) {
    __m128i one_third = _mm_set1_epi16(21846);
    #pragma omp parallel for
    for (int yTile = 0; yTile < in.height(); yTile += 32) {
        __m128i a, b, c, sum, avg;
        __m128i blurx[(256/8)*(32+2)]; // allocate tile blurx array
        for (int xTile = 0; xTile < in.width(); xTile += 256) {
            __m128i *blurxPtr = blurx;
            for (int y = -1; y < 32+1; y++) {
                const uint16_t *inPtr = &(in[yTile+y][xTile]);
                for (int x = 0; x < 256; x += 8) {
                    a = _mm_loadu_si128((__m128i*)(inPtr-1));
                    b = _mm_loadu_si128((__m128i*)(inPtr+1));
                    c = _mm_load_si128((__m128i*)(inPtr));
                    sum = _mm_add_epi16(_mm_add_epi16(a, b), c);
                    avg = _mm_mulhi_epi16(sum, one_third);
                    _mm_store_si128(blurxPtr++, avg);
                    inPtr += 8;
                }
                blurxPtr = blurx;
                for (int y = 0; y < 32; y++) {
                    __m128i *outPtr = ((__m128i *)(&(blury[yTile+y][xTile])));
                    for (int x = 0; x < 256; x += 8) {
                        a = _mm_load_si128(blurxPtr+(2*256)/8);
                        b = _mm_load_si128(blurxPtr+256/8);
                        c = _mm_load_si128(blurxPtr++);
                        sum = _mm_add_epi16(_mm_add_epi16(a, b), c);
                        avg = _mm_mulhi_epi16(sum, one_third);
                        _mm_store_si128(outPtr++, avg);
                    }
                }
            }
        }
    }
}
```

9.9 → 0.9 ms/megapixel

11x faster
(quad core x86)

Tiled, fused
Vectorized
Multithreaded
Redundant computation
Near roof-line optimum

HALIDE DECOUPLES ALGORITHM FROM SCHEDULE

- Algorithm: *what* is computed
 - The algorithm defines pipelines as pure functions
 - Pipeline stages are functions from coordinates to values
 - Execution order and storage are unspecified

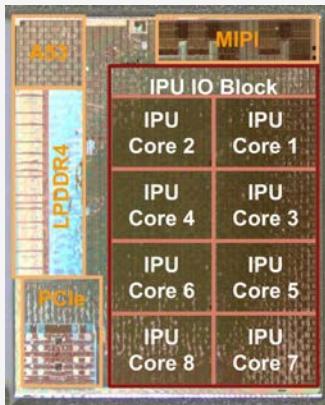
```
blurx(x, y) = (in(x-1, y) + in(x, y) + in(x+1, y))/3;  
blury(x, y) = (blurx(x, y-1) + blurx(x, y) + blurx(x, y+1))/3;
```

HALIDE DECOUPLES ALGORITHM FROM SCHEDULE

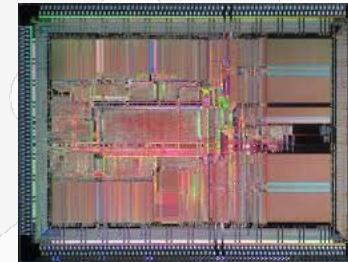
- Algorithm: *what* is computed
 - The algorithm defines pipelines as pure functions
 - Pipeline stages are functions from coordinates to values
 - Execution order and storage are unspecified
- Schedule: *where and when* it's computed
 - Architecture Specific
 - Single, unified model for all schedules
 - Simple enough to search, expose to user
 - Powerful enough to beat expert-tuned code

DSAS FOR IMAGE PROCESSING

Google Pixel Visual Core



Qualcomm Snapdragon Hexagon 682 DSP



- Halide is the DSL for programming these DSAs
- Programmers able to find the best hardware utilization using the scheduling language
- Thus, the DSAs are more general purpose and easy to use

A row of three soft-shell tacos filled with meat, cheese, lettuce, and salsa, arranged horizontally across the background.

TACO

A DSL for Sparse Tensor Algebra

TENSORS ARE EVERYWHERE

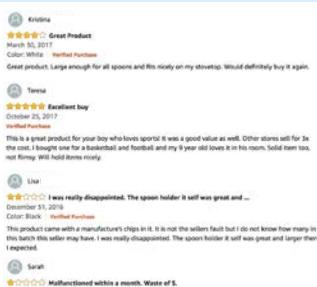
Data Analytics



Movies

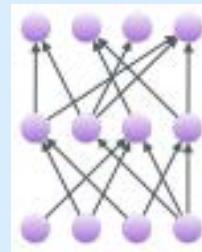


Social Networks

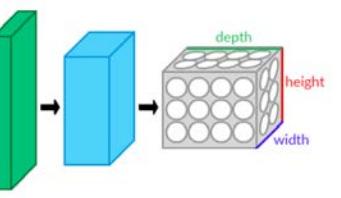


Product Reviews

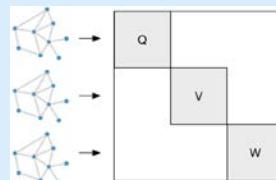
Machine Learning



Sparse Networks



Sparse Convolutional Networks

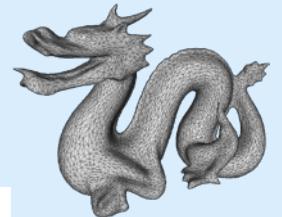


Graph Convolutional Network

Science and Engineering



Robotics

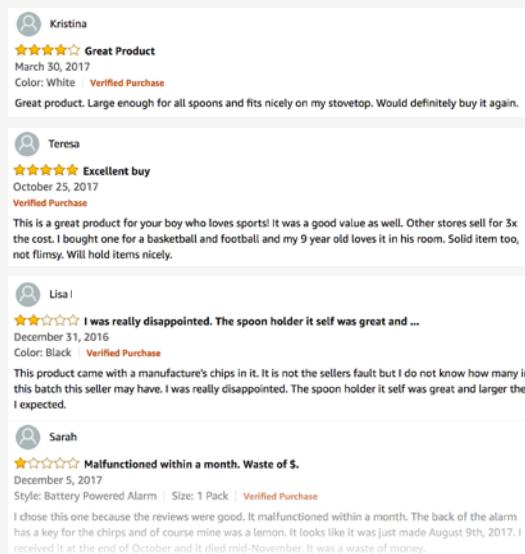


Simulations

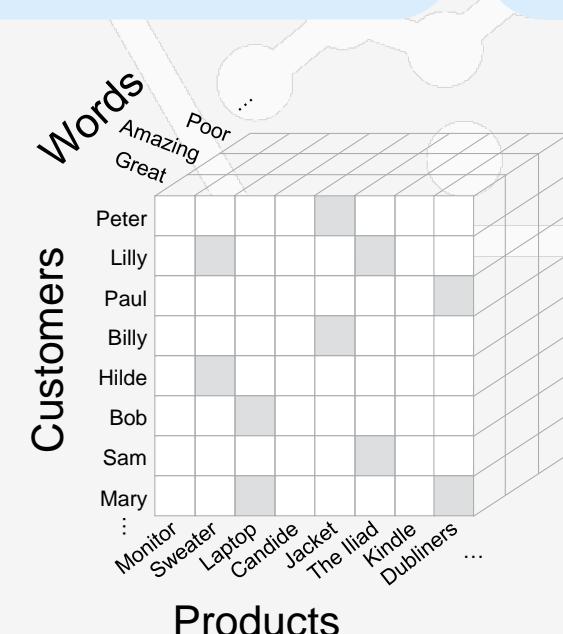


Computational Biology

Amazon Product Reviews



DISTRIBUTION STATEMENT A. Approved for public release



Products

Extremely sparse
Dense storage: 107 Exabytes
Sparse storage: 13 Gigabytes

Images from the
Wikimedia and
computational Imaging

COMPLEXITY OF SPARSE CODE

$$A_{ijk} = B_{ijk} + C_{ijk}$$

↑
CSF COO

```
int iB = 0;
int C0_pos = C0_pos[0];
while (C0_pos < C0_pos[1]) {
    int iC = C0_crd[C0_pos];
    int C0_end = C0_pos + 1;
    if (iC == iB)
        while ((C0_end < C0_pos[1]) && (C0_crd[C0_end] == iB))
            C0_end++;
    }
    if (iC == iB) {
        int B1_pos = B1_pos[iB];
        int C1_pos = C0_pos;
        while ((B1_pos < B1_pos[iB + 1]) && (C1_pos < C0_end)) {
            int jB = B1_crd[B1_pos];
            int iC = C1_crd[C1_pos];
            int jC = min(jB, iC);
            int A1_pos = (iB * A1_size) + jC;
            int C1_end = C1_pos + 1;
            if (jC == iB)
                while ((C1_end < C0_end) && (C1_crd[C1_end] == jC))
                    C1_end++;
            }
            if ((jB == jC) && (jC == iB)) {
                int B2_pos = B2_pos[B1_pos];
                int C2_pos = C1_pos;
                while ((B2_pos < B2_pos[B1_pos + 1]) && (C2_pos < C1_end)) {
                    int kB = B2_crd[B2_pos];
                    int kC = C2_crd[C2_pos];
                    int k = min(kB, kC);
                    int A2_pos = (A1_pos * A2_size) + k;
                    if ((kB == k) && (kC == k)) {
                        A[A2_pos] = B[B2_pos] + C[C2_pos];
                    } else if (kB == k) {
                        A[A2_pos] = B[B2_pos];
                    } else {
                        A[A2_pos] = C[C2_pos];
                    }
                    if (kB == k) B2_pos++;
                    if (kC == k) C2_pos++;
                }
                while (B2_pos < B2_pos[B1_pos + 1]) {
                    int kB0 = B2_crd[B2_pos];
                    int A2_pos0 = (A1_pos * A2_size) + kB0;
                    A[A2_pos0] = B[B2_pos];
                    B2_pos++;
                }
                while (C2_pos < C1_end) {
                    int kC0 = C2_crd[C2_pos];
                    int A2_pos1 = (A1_pos * A2_size) + kC0;
                    A[A2_pos1] = C[C2_pos];
                    C2_pos++;
                }
            } else if (jB == iB) {
                for (int B2_pos0 = B2_pos[B1_pos]; B2_pos0 < B2_pos[B1_pos + 1]; B2_pos0++) {
                    int kB1 = B2_crd[B2_pos0];
                    int A2_pos2 = (A1_pos * A2_size) + kB1;
                    A[A2_pos2] = B[B2_pos0];
                }
            } else {
                for (int C2_pos0 = C1_pos; C2_pos0 < C1_end; C2_pos0++) {
                    int kC1 = C2_crd[C2_pos0];
                    int A2_pos3 = (A1_pos * A2_size) + kC1;
                    A[A2_pos3] = C[C2_pos0];
                }
            }
        if (jB == iB) B1_pos++;
        if (jC == iB) C1_pos = C1_end;
    }
}
```

```
while (B1_pos < B1_pos[iB + 1]) {
    int jB0 = B1_crd[B1_pos];
    int A1_pos0 = (iB * A1_size) + jB0;
    for (int B2_pos1 = B2_pos[B1_pos]; B2_pos1 < B2_pos[B1_pos + 1]; B2_pos1++) {
        int kB2 = B2_crd[B2_pos1];
        int A2_pos4 = (A1_pos0 * A2_size) + kB2;
        A[A2_pos4] = B[B2_pos1];
    }
    B1_pos++;
}
while (C1_pos < C0_end) {
    int jC0 = C1_crd[C1_pos];
    int A1_pos1 = (iB * A1_size) + jC0;
    int C1_end0 = C1_pos + 1;
    while ((C1_end0 < C0_end) && (C1_crd[C1_end0] == jC0)) {
        C1_end0++;
    }
    for (int C2_pos1 = C1_pos; C2_pos1 < C1_end0; C2_pos1++) {
        int kC2 = C2_crd[C2_pos1];
        int A2_pos5 = (A1_pos1 * A2_size) + kC2;
        A[A2_pos5] = C[C2_pos1];
    }
    C1_pos = C1_end0;
} else {
    for (int B1_pos0 = B1_pos[iB]; B1_pos0 < B1_pos[iB + 1]; B1_pos0++) {
        int jB1 = B1_crd[B1_pos0];
        int A1_pos2 = (iB * A1_size) + jB1;
        for (int B2_pos2 = B2_pos[B1_pos0]; B2_pos2 < B2_pos[B1_pos0 + 1]; B2_pos2++) {
            int kB3 = B2_crd[B2_pos2];
            int A2_pos6 = (A1_pos2 * A2_size) + kB3;
            A[A2_pos6] = B[B2_pos2];
        }
    }
    if (iC == iB) C0_pos = C0_end;
    iB++;
}
while (iB < B0_size) {
    for (int B1_pos1 = B1_pos[iB]; B1_pos1 < B1_pos[iB + 1]; B1_pos1++) {
        int kB2 = B1_crd[B1_pos1];
        int A1_pos3 = (iB * A1_size) + kB2;
        for (int B2_pos3 = B2_pos[B1_pos1]; B2_pos3 < B2_pos[B1_pos1 + 1]; B2_pos3++) {
            int kB4 = B2_crd[B2_pos3];
            int A2_pos7 = (A1_pos3 * A2_size) + kB4;
            A[A2_pos7] = B[B2_pos3];
        }
    }
    iB++;
}
```

TOO MANY TENSOR KERNELS FOR A FIXED-FUNCTION LIBRARY

CSpars

$$a = Bc + a$$

Eigen (SpMV)

$$a = Bc$$

$$a = Bc + b \quad A = B + C \quad a = \alpha Bc + \beta a$$

PETSc

$$a = B^T c \quad A = \alpha B \quad a = B(c + d)$$

$$a = B^T c + d \quad A = B + C + D \quad A = BC$$

$$A = B \odot C \quad a = b \odot c \quad A = 0 \quad A = B \odot (CD)$$

$$A = BCd \quad A = B^T \quad a = B^T Bc$$

$$a = b + c \quad A = B \quad K = A^T CA$$

$$A_{ij} = \sum_{kl} B_{ikl} C_{lj} D_{kj} \quad A_{kj} = \sum_{il} B_{ikl} C_{lj} D_{ij}$$

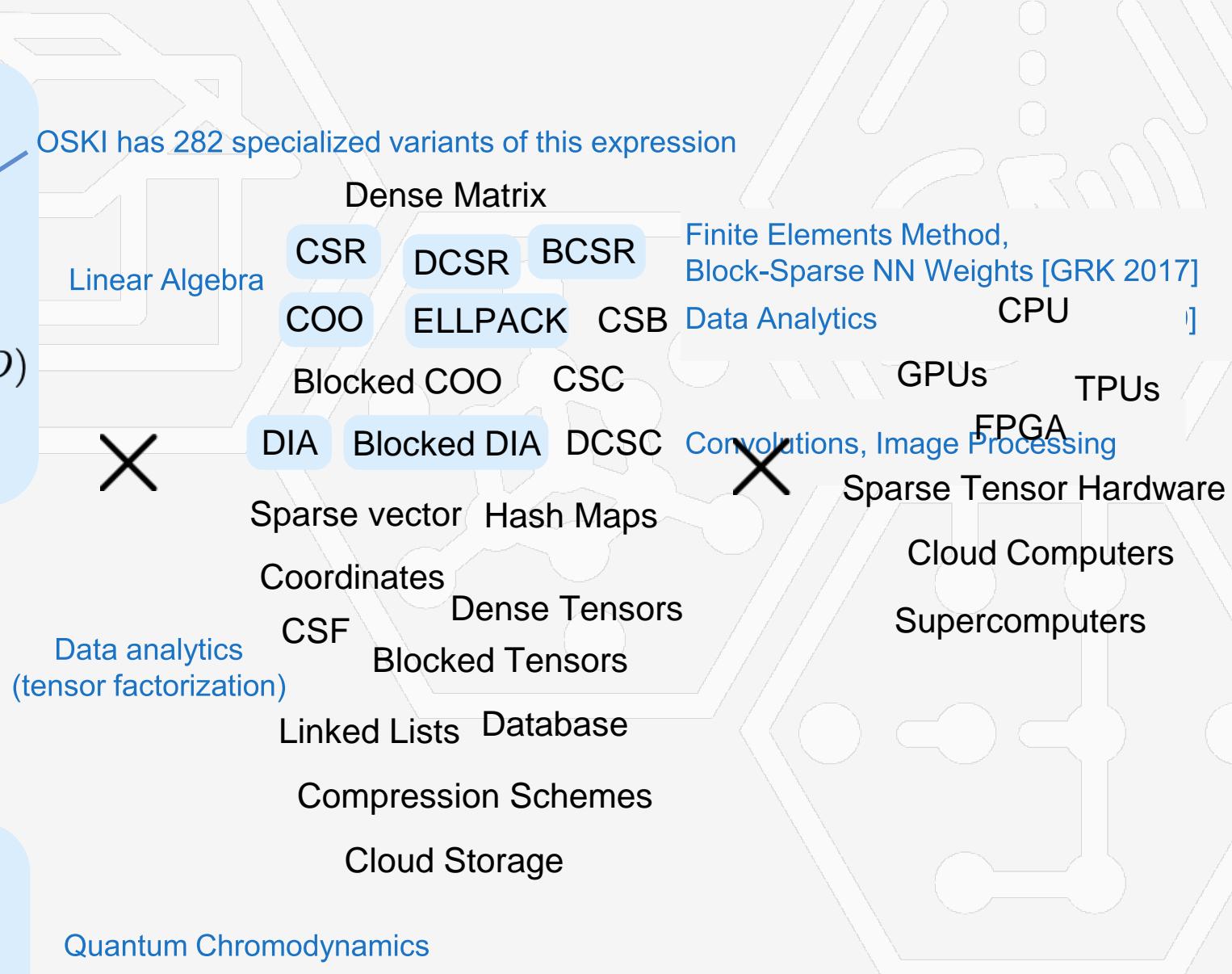
$$A_{lj} = \sum_{ik} B_{ikl} C_{ij} D_{kj} \quad A_{ij} = \sum_k B_{ijk} c_k$$

$$A_{ijk} = \sum_l B_{ikl} C_{lj} \quad A_{ik} = \sum_j B_{ijk} c_j$$

$$A_{jk} = \sum_i B_{ijk} c_i \quad A_{ijl} = \sum_k B_{ikl} C_{kj}$$

$$C = \sum_{ijkl} M_{ij} P_{jk} \overline{M_{lk}} \overline{P_{il}} \quad \tau = \sum_i z_i (\sum_j z_j \theta_{ij}) (\sum_k z_k \theta_{ik})$$

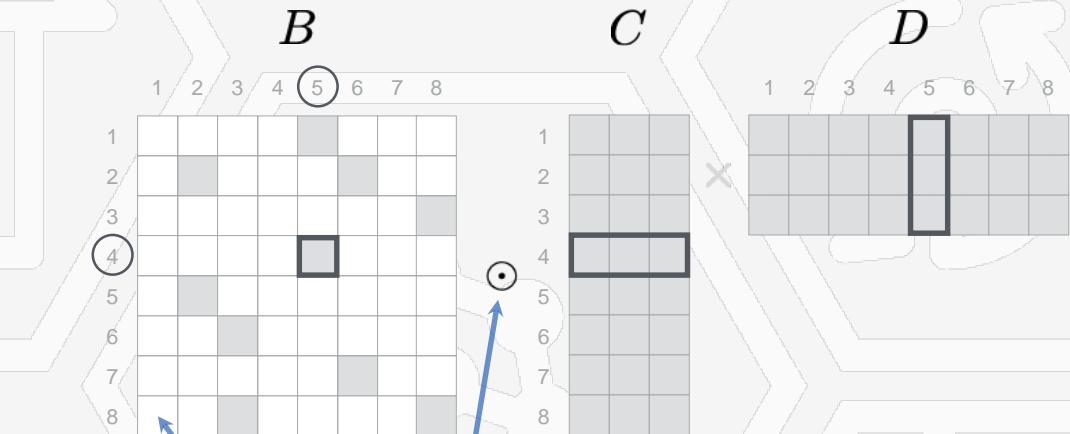
$$a = \sum_{ijklmno} M_{ij} P_{jk} M_{kl} P_{lm} \overline{M_{nm}} P_{no} \overline{M_{po}} \overline{P_{ip}}$$



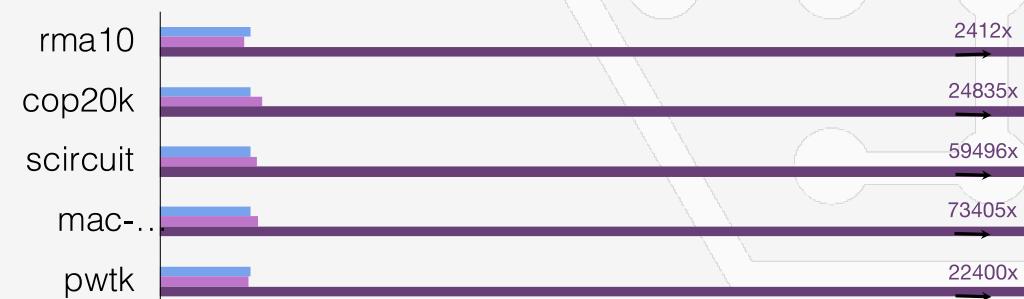
MANY NON-BINARY EXPRESSIONS MUST BE COMPUTED IN A SINGLE KERNEL

$$\begin{aligned}
 & a = Bc + a & a = Bc \\
 & a = Bc + b & A = B + C & a = \alpha Bc + \beta a \\
 & a = B^T c & A = \alpha B & a = B(c + d) \\
 & a = B^T c + d & A = B + C + D & A = BC \\
 & A = B \odot C & a = b \odot c & A = 0 \\
 & A = BCd & A = B^T & a = B^T Bc \\
 & a = b + c & A = B & K = A^T CA \\
 & A_{ij} = \sum_{kl} B_{ikl} C_{lj} D_{kj} & A_{kj} = \sum_{il} B_{ikl} C_{lj} D_{ij} \\
 & A_{lj} = \sum_{ik} B_{ikl} C_{ij} D_{kj} & A_{ij} = \sum_k B_{ijk} c_k \\
 & A_{ijk} = \sum_l B_{ikl} C_{lj} & A_{ik} = \sum_j B_{ijk} c_j \\
 & A_{jk} = \sum_i B_{ijk} c_i & A_{ijl} = \sum_k B_{ikl} C_{kj} \\
 & C = \sum_{ijkl} M_{ij} P_{jk} \overline{M}_{lk} \overline{P}_{il} & \tau = \sum_i z_i (\sum_j z_j \theta_{ij}) (\sum_k z_k \theta_{ik}) \\
 & a = \sum_{ijklmno} M_{ij} P_{jk} M_{kl} P_{lm} \overline{M}_{nm} \overline{P}_{no} \overline{M}_{po} \overline{P}_{ip}
 \end{aligned}$$

Sampled Dense-Dense Matrix Multiplication (SDDMM)



We will generate fused operations



MANY NON-BINARY EXPRESSIONS MUST BE COMPUTED IN A SINGLE KERNEL

$a = Bc$

$$a = Bc + a$$

$$a = Bc + b \quad A = B + C \quad a = \alpha Bc + \beta a$$

$$a = B^T c \quad A = \alpha B \quad a = B(c + d)$$

$$a = B^T c + d \quad A = B + C + D \quad A = BC$$

$$A = B \odot C \quad a = b \odot c \quad A = 0 \quad A = B \odot (CD)$$

$$A = BCd \quad A = B^T \quad a = B^T Bc$$

$$a = b + c \quad A = B \quad K = A^T CA$$

$$A_{ij} = \sum_{kl} B_{ikl} C_{lj} D_{kj} \quad A_{kj} = \sum_{il} B_{ikl} C_{lj} D_{ij}$$

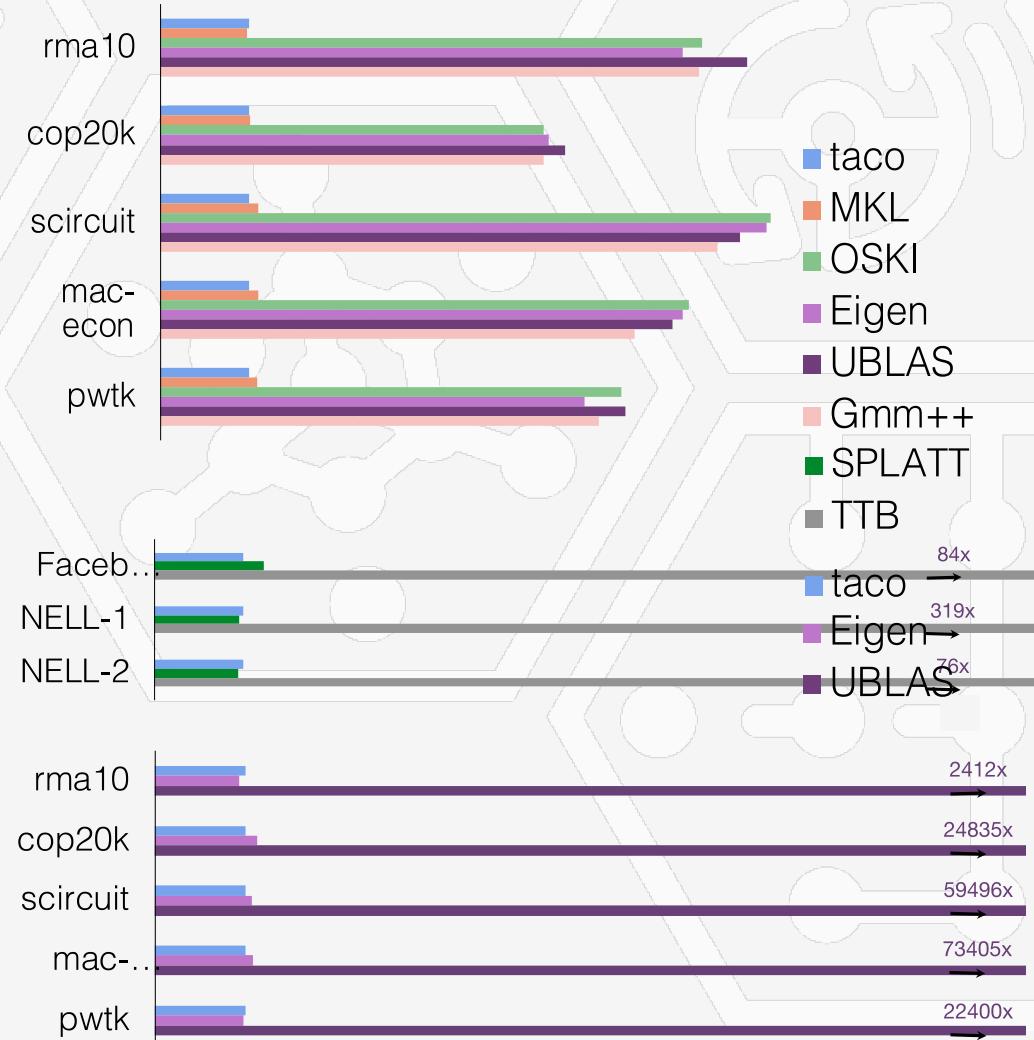
$$A_{lj} = \sum_{ik} B_{ikl} C_{ij} D_{kj} \quad A_{ij} = \sum_k B_{ijk} c_k$$

$$A_{ijk} = \sum_l B_{ikl} C_{lj} \quad A_{ik} = \sum_j B_{ijk} c_j$$

$$A_{jk} = \sum_i B_{ijk} c_i \quad A_{ijl} = \sum_k B_{ikl} C_{kj}$$

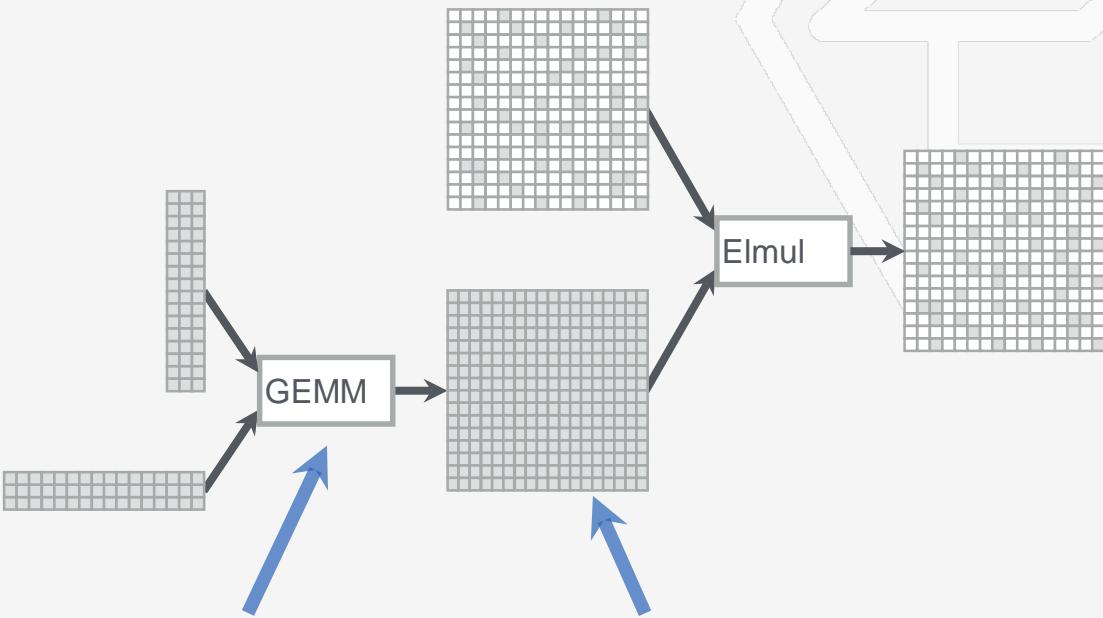
$$C = \sum_{ijkl} M_{ij} P_{jk} \overline{M}_{lk} \overline{P}_{il} \quad \tau = \sum_i z_i (\sum_j z_j \theta_{ij}) (\sum_k z_k \theta_{ik})$$

$$a = \sum_{ijklmno} M_{ij} P_{jk} M_{kl} P_{lm} \overline{M}_{nm} P_{no} \overline{M}_{po} \overline{P}_{ip}$$



PERFORMANCE COST OF LIBRARY ABSTRACTION IS HIGH

Fixed-function Libraries



Hand-Optimized Code

```

int iB = 0;
int C0_pos = C0_pos_arr[0];
while (C0_pos < C0_pos_arr[1]) {
    int iC = C0_idx_arr[C0_pos];
    int C0_end = C0_pos + 1;
    if (iC == iB)
        while ((C0_end < C0_pos_arr[1]) && (C0_idx_arr[C0_end] == iB)) {
            C0_end++;
        }
    if (iC == iB) {
        int B1_pos = B1_pos_arr[iB];
        int C1_pos = C0_pos;
        while ((B1_pos < B1_pos_arr[iB + 1]) && (C1_pos < C0_end)) {
            int jB = B1_idx_arr[B1_pos];
            int jC = C1_idx_arr[C1_pos];
            int j = min(jB, jC);
            int A1_pos = (iB * A1_size) + j;
            int C1_end = C1_pos + 1;
            if (jC == j)
                while ((C1_end < C0_end) && (C1_idx_arr[C1_end] == j)) {
                    C1_end++;
                }
            if ((jB == j) && (jC == j)) {
                int B2_pos = B2_pos_arr[B1_pos];
                int C2_pos = C1_pos;
                while ((B2_pos < B2_pos_arr[B1_pos + 1]) && (C2_pos < C1_end)) {
                    int kB = B2_idx_arr[B2_pos];
                    int kC = C2_idx_arr[C2_pos];
                    int k = min(kB, kC);
                    int A2_pos = (A1_pos * A2_size) + k;
                    if ((kB == k) && (kC == k)) {
                        A_val_arr[A2_pos] = B_val_arr[B2_pos] + C_val_arr[C2_pos];
                    } else if (kB == k) {
                        A_val_arr[A2_pos] = B_val_arr[B2_pos];
                    } else {
                        A_val_arr[A2_pos] = C_val_arr[C2_pos];
                    }
                    if (kB == k) B2_pos++;
                    if (kC == k) C2_pos++;
                }
            while (B2_pos < B2_pos_arr[B1_pos + 1]) {
                int kB0 = B2_idx_arr[B2_pos];
                int A2_pos0 = (A1_pos * A2_size) + kB0;
                A_val_arr[A2_pos0] = B_val_arr[B2_pos];
                B2_pos++;
            }
            while (C2_pos < C1_end) {
                int kC0 = C2_idx_arr[C2_pos];
                int A2_pos1 = (A1_pos * A2_size) + kC0;
                A_val_arr[A2_pos1] = C_val_arr[C2_pos];
                C2_pos++;
            }
            } else if (jB == j) {
                for (int B2_pos0 = B2_pos_arr[B1_pos];
                     B2_pos0 < B2_pos_arr[B1_pos + 1]; B2_pos0++) {
                    int kB1 = B2_idx_arr[B2_pos0];
                    int A2_pos2 = (A1_pos * A2_size) + kB1;
                    A_val_arr[A2_pos2] = B_val_arr[B2_pos0];
                }
            } else {
                for (int C2_pos0 = C1_pos;
                     C2_pos0 < C1_end; C2_pos0++) {
                    int kC1 = C2_idx_arr[C2_pos0];
                    int A2_pos3 = (A1_pos * A2_size) + kC1;
                    A_val_arr[A2_pos3] = C_val_arr[C2_pos0];
                }
            }
            if (jB == j) B1_pos++;
            if (jC == j) C1_pos = C1_end;
        }
    }
}

```

```

while (B1_pos < B1_pos_arr[iB + 1]) {
    int jB0 = B1_idx_arr[B1_pos];
    int A1_pos0 = (iB * A1_size) + jB0;
    for (int B2_pos1 = B2_pos_arr[B1_pos];
         B2_pos1 < B2_pos_arr[B1_pos + 1]; B2_pos1++) {
        int kB2 = B2_idx_arr[B2_pos1];
        int A2_pos4 = (A1_pos0 * A2_size) + kB2;
        A_val_arr[A2_pos4] = B_val_arr[B2_pos1];
    }
    B1_pos++;
}
while (C1_pos < C0_end) {
    int jC0 = C1_idx_arr[C1_pos];
    int A1_pos1 = (iB * A1_size) + jC0;
    int C1_end0 = C1_pos + 1;
    while ((C1_end0 < C0_end) && (C1_idx_arr[C1_end0] == jC0)) {
        C1_end0++;
    }
    for (int C2_pos1 = C1_pos;
         C2_pos1 < C1_end0; C2_pos1++) {
        int kC2 = C2_idx_arr[C2_pos1];
        int A2_pos5 = (A1_pos1 * A2_size) + kC2;
        A_val_arr[A2_pos5] = C_val_arr[C2_pos1];
    }
    C1_pos = C1_end0;
} else {
    for (int B1_pos0 = B1_pos_arr[iB];
         B1_pos0 < B1_pos_arr[iB + 1]; B1_pos0++) {
        int jB1 = B1_idx_arr[B1_pos0];
        int A1_pos2 = (iB * A1_size) + jB1;
        for (int B2_pos2 = B2_pos_arr[B1_pos0];
             B2_pos2 < B2_pos_arr[B1_pos0 + 1]; B2_pos2++) {
            int kB3 = B2_idx_arr[B2_pos2];
            int A2_pos6 = (A1_pos2 * A2_size) + kB3;
            A_val_arr[A2_pos6] = B_val_arr[B2_pos2];
        }
    }
    if (iC == iB) C0_pos = C0_end;
    iB++;
}
while (iB < B0_size) {
    for (int B1_pos1 = B1_pos_arr[iB];
         B1_pos1 < B1_pos_arr[iB + 1]; B1_pos1++) {
        int jB2 = B1_idx_arr[B1_pos1];
        int A1_pos3 = (iB * A1_size) + jB2;
        for (int B2_pos3 = B2_pos_arr[B1_pos1];
             B2_pos3 < B2_pos_arr[B1_pos1 + 1]; B2_pos3++) {
            int kB4 = B2_idx_arr[B2_pos3];
            int A2_pos7 = (A1_pos3 * A2_size) + kB4;
            A_val_arr[A2_pos7] = B_val_arr[B2_pos3];
        }
    }
    iB++;
}

```

$$A_{ijk} = B_{ijk} + C_{ijk}$$

ABSTRACTION WITHOUT FRICTION

Language and Compiler Abstractions

Domain
Specific
Languages



Compiler

Generated Code

```
int iB = 0;
int C0_pos = C0_pos_arr[0];
while (C0_pos < C0_pos_arr[1]) {
    int iC = C0_idx_arr[C0_pos];
    int C0_end = C0_pos + 1;
    if (iC == iB)
        while ((C0_end < C0_pos_arr[1]) && (C0_idx_arr[C0_end] == iB)) {
            C0_end++;
        }
    if (iC == iB) {
        int B1_pos = B1_pos_arr[iB];
        int C1_pos = C0_pos;
        while ((B1_pos < B1_pos_arr[iB + 1]) && (C1_pos < C0_end)) {
            int jB = B1_idx_arr[B1_pos];
            int jC = C1_idx_arr[C1_pos];
            int j = min(jB, jC);
            int A1_pos = (iB * A1_size) + j;
            int C1_end = C1_pos + 1;
            if (jC == j)
                while ((C1_end < C0_end) && (C1_idx_arr[C1_end] == j)) {
                    C1_end++;
                }
            if ((jB == j) && (jC == j)) {
                int B2_pos = B2_pos_arr[B1_pos];
                int C2_pos = C1_pos;
                while ((B2_pos < B2_pos_arr[B1_pos + 1]) && (C2_pos < C1_end)) {
                    int kB = B2_idx_arr[B2_pos];
                    int kC = C2_idx_arr[C2_pos];
                    int k = min(kB, kC);
                    int A2_pos = (A1_pos * A2_size) + k;
                    if ((kB == k) && (kC == k)) {
                        A_val_arr[A2_pos] = B_val_arr[B2_pos] + C_val_arr[C2_pos];
                    } else if (kB == k) {
                        A_val_arr[A2_pos] = B_val_arr[B2_pos];
                    } else {
                        A_val_arr[A2_pos] = C_val_arr[C2_pos];
                    }
                    if (kB == k) B2_pos++;
                    if (kC == k) C2_pos++;
                }
                while (B2_pos < B2_pos_arr[B1_pos + 1]) {
                    int kB0 = B2_idx_arr[B2_pos];
                    int A2_pos0 = (A1_pos * A2_size) + kB0;
                    A_val_arr[A2_pos0] = B_val_arr[B2_pos];
                    B2_pos++;
                }
                while (C2_pos < C1_end) {
                    int kC0 = C2_idx_arr[C2_pos];
                    int A2_pos1 = (A1_pos * A2_size) + kC0;
                    A_val_arr[A2_pos1] = C_val_arr[C2_pos];
                    C2_pos++;
                }
            } else if (jB == j) {
                for (int B2_pos0 = B2_pos_arr[B1_pos];
                     B2_pos0 < B2_pos_arr[B1_pos + 1]; B2_pos0++) {
                    int kB1 = B2_idx_arr[B2_pos0];
                    int A2_pos2 = (A1_pos * A2_size) + kB1;
                    A_val_arr[A2_pos2] = B_val_arr[B2_pos0];
                }
            } else {
                for (int C2_pos0 = C1_pos;
                     C2_pos0 < C1_end; C2_pos0++) {
                    int kC1 = C2_idx_arr[C2_pos0];
                    int A2_pos3 = (A1_pos * A2_size) + kC1;
                    A_val_arr[A2_pos3] = C_val_arr[C2_pos0];
                }
            }
            if (jB == j) B1_pos++;
            if (jC == j) C1_pos = C1_end;
        }
    }
}
```

```
while (B1_pos < B1_pos_arr[iB + 1]) {
    int jB0 = B1_idx_arr[B1_pos];
    int A1_pos0 = (iB * A1_size) + jB0;
    for (int B2_pos1 = B2_pos_arr[B1_pos];
         B2_pos1 < B2_pos_arr[B1_pos + 1]; B2_pos1++) {
        int kB2 = B2_idx_arr[B2_pos1];
        int A2_pos4 = (A1_pos0 * A2_size) + kB2;
        A_val_arr[A2_pos4] = B_val_arr[B2_pos1];
    }
    B1_pos++;
}
while (C1_pos < C0_end) {
    int jC0 = C1_idx_arr[C1_pos];
    int A1_pos1 = (iB * A1_size) + jC0;
    int C1_end0 = C1_pos + 1;
    while ((C1_end0 < C0_end) && (C1_idx_arr[C1_end0] == jC0)) {
        C1_end0++;
    }
    for (int C2_pos1 = C1_pos;
         C2_pos1 < C1_end0; C2_pos1++) {
        int kC2 = C2_idx_arr[C2_pos1];
        int A2_pos5 = (A1_pos1 * A2_size) + kC2;
        A_val_arr[A2_pos5] = C_val_arr[C2_pos1];
    }
    C1_pos = C1_end0;
} else {
    for (int B1_pos0 = B1_pos_arr[iB];
         B1_pos0 < B1_pos_arr[iB + 1]; B1_pos0++) {
        int jB1 = B1_idx_arr[B1_pos0];
        int A1_pos2 = (iB * A1_size) + jB1;
        for (int B2_pos2 = B2_pos_arr[B1_pos0];
             B2_pos2 < B2_pos_arr[B1_pos0 + 1]; B2_pos2++) {
            int kB3 = B2_idx_arr[B2_pos2];
            int A2_pos6 = (A1_pos2 * A2_size) + kB3;
            A_val_arr[A2_pos6] = B_val_arr[B2_pos2];
        }
    }
    if (iC == iB) C0_pos = C0_end;
    iB++;
}
while (iB < B0_size) {
    for (int B1_pos1 = B1_pos_arr[iB];
         B1_pos1 < B1_pos_arr[iB + 1]; B1_pos1++) {
        int jB2 = B1_idx_arr[B1_pos1];
        int A1_pos3 = (iB * A1_size) + jB2;
        for (int B2_pos3 = B2_pos_arr[B1_pos1];
             B2_pos3 < B2_pos_arr[B1_pos1 + 1]; B2_pos3++) {
            int kB4 = B2_idx_arr[B2_pos3];
            int A2_pos7 = (A1_pos3 * A2_size) + kB4;
            A_val_arr[A2_pos7] = B_val_arr[B2_pos3];
        }
    }
    iB++;
}
```

$$A_{ijk} = B_{ijk} + C_{ijk}$$

THE TENSOR ALGEBRA COMPILER (TACO)

Expressions

$$\begin{aligned} A &= Bc + a & a &= Bc \\ A &= B \odot C & A &= B + C & a &= \alpha Bc + \beta a \\ A &= BCd & A &= \alpha B & A &= 0 & A &= BC \\ A_{ij} &= \sum_{kl} B_{ikl}C_{lj}D_{kj} & A &= B^T & a &= B^T Bc \\ A_{ijk} &= \sum_l B_{ikl}C_{lj} & A_{ik} &= \sum_j B_{ijk}c_j & A_{kj} &= \sum_{il} B_{ikl}C_{lj}D_{ij} \\ C &= \sum_{ijkl} M_{ij}P_{jk}\overline{M_{lk}}\overline{P_{il}} & A_{ij} &= (\sum_k B_{ijk}C_{ijk}) + D_{ij} \\ a &= \sum_{ijklmnp} M_{ij}P_{jk}M_{kl}P_{lm}\overline{M_{nm}}\overline{P_{no}}\overline{M_{po}}\overline{P_{ip}} & \tau &= \sum_i z_i(\sum_j z_j\theta_{ij})(\sum_k z_k\theta_{ik}) \end{aligned}$$

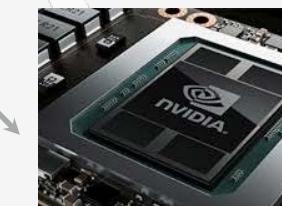
Formats

Dense Matrix	CSR	BCSR
COO	DCSR	ELLPACK
DIA	Blocked COO	CSB
Blocked DIA	DCSC	
Sparse vector	Hash Maps	
CSF	Dense Tensors	
	Blocked Tensors	

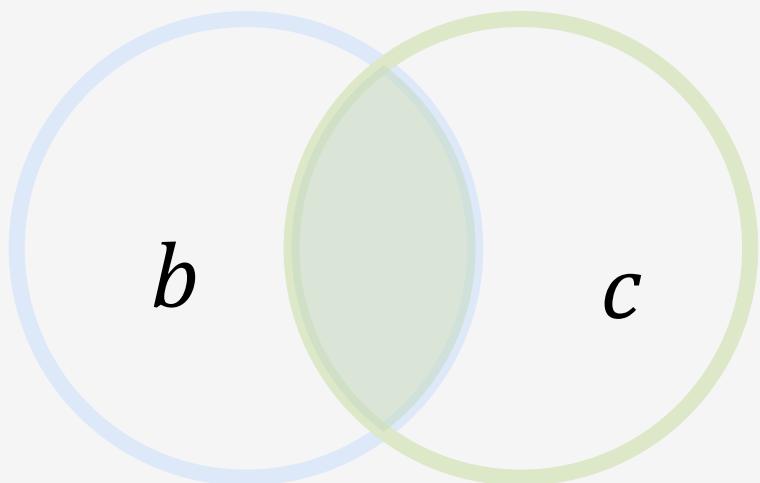
Tensor Algebra Compiler
(taco)



DARPA SDH
Nvidia Symphony DSA



AVOID WASTED WORK AND ITERATIONS

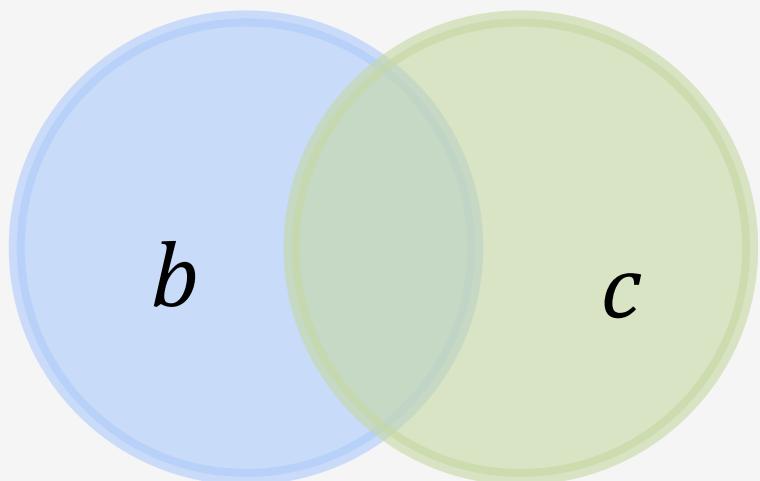


$$a_i = b_i c_i$$

Multiplication requires intersection

```
int pb1 = b1_pos[0];
int pc1 = c1_pos[0];
while (pb1 < b1_pos[1] && pc1 < c1_pos[1]) {
    int ib = b1_crd[pb1];
    int ic = c1_crd[pc1];
    int i = min(ib, ic);
    if (ib == i && ic == i) {
        a[i] = b[pb1] * c[pc1];
    }
    if (ib == i) pb1++;
    if (ic == i) pc1++;
}
```

AVOID WASTED WORK AND ITERATIONS

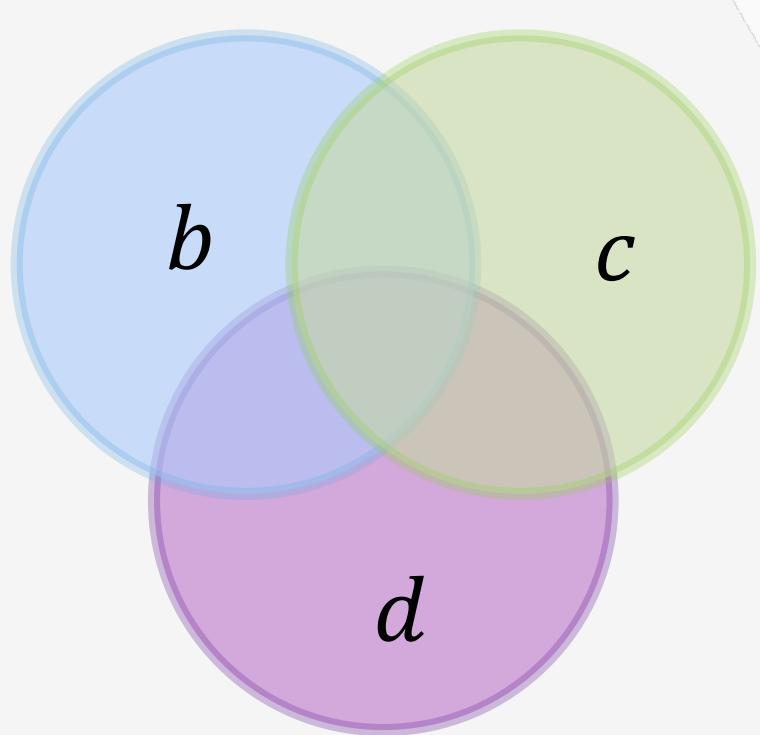


$$a_i = b_i + c_i$$

Addition requires union

```
int pb1 = b1_pos[0];
int pc1 = c1_pos[0];
while (pb1 < b1_pos[1] && pc1 < c1_pos[1]) {
    int ib = b1_crd[pb1];
    int ic = c1_crd[pc1];
    int i = min(ib, ic);
    if (ib == i && ic == i) {
        a[i] = b[pb1] + c[pc1];
    }
    else if (ib == i) {
        a[i] = b[pb1];
    }
    else {
        a[i] = c[pc1];
    }
    if (ib == i) pb1++;
    if (ic == i) pc1++;
}
while (pb1 < b1_pos[1]) {
    int i = b1_crd[pb1];
    a[i] = b[pb1++];
}
while (pc1 < c1_pos[1]) {
    int i = c1_crd[pc1];
    a[i] = c[pc1++];
}
```

AVOID WASTED WORK AND ITERATIONS



$$a_i = b_i + c_i + d_i$$

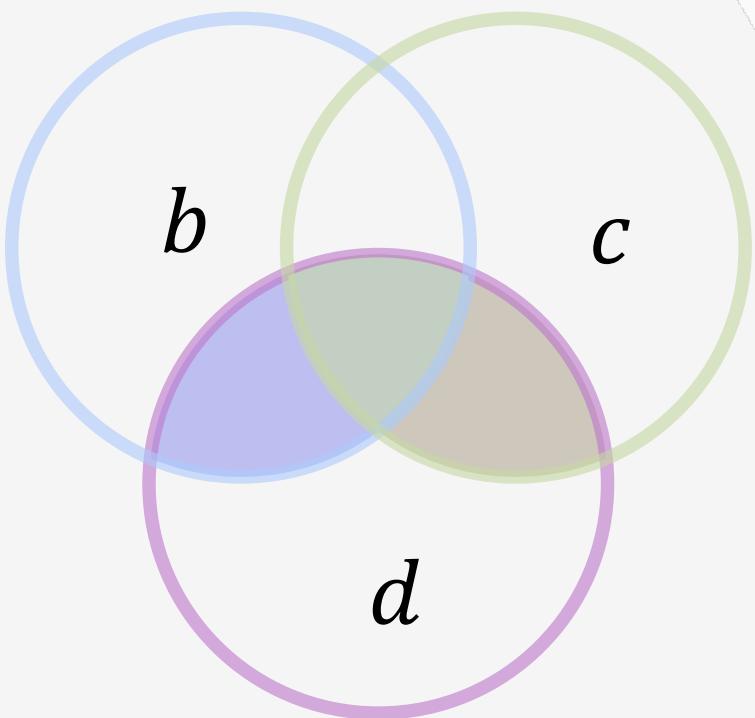
```
int pc1 = b1_pos[0];
int pc1 = c1_pos[0];
int id = 0;
while (pb1 < b1_pos[1] && pc1 < c1_pos[1]) {
    int ib = b1_crdl[pb1];
    int pd1 = id;
    int pat = id;
    if (ib == id && ic == id) {
        apat1 = bp1 + cp1 + dp1;
    }
    else if (ib == id) {
        apat1 = bp1 + dp1;
    }
    else if (ic == id) {
        apat1 = cp1 + dp1;
    }
    else {
        apat1 = dp1;
    }
    if (ib == id) pb1++;
    if (ic == id) pc1++;
    id++;
}

while (pc1 < c1_pos[1] && pd1 < d1_pos[1]) {
    int ic = c1_crdl[pc1];
    int id = d1_crdl[pd1];
    int i = min(ic, id);
    if (ic == i && id == i) {
        a[i] = cp1 + dp1;
    }
    else if (ic == i) {
        a[i] = cp1;
    }
    else {
        a[i] = dp1;
    }
    if (id == i) pd1++;
    if (ic == i) pc1++;
}

while (pb1 < b1_pos[1] && pc1 < c1_pos[1]) {
    int ib = b1_crdl[pb1];
    int ic = c1_crdl[pc1];
    int i = min(ib, ic);
    if (ib == i && ic == i) {
        a[i] = bp1 + cp1;
    }
    else if (ib == i) {
        a[i] = bp1;
    }
    else {
        a[i] = cp1;
    }
    if (ib == i) pb1++;
    if (ic == i) pc1++;
}

while (id < d1_dimension) {
    int pd1 = id;
    int pat = id;
    apat1 = dp1;
    id++;
}
```

AVOID WASTED WORK AND ITERATIONS



$$a_i = (b_i + c_i)d_i$$

```
int pb1 = b1_pos[0];
int pc1 = c1_pos[0];
int pd1 = d1_pos[0];
while (pb1 < b1_pos[1] && pc1 < c1_pos[1] && pd1 < d1_pos[1]) {
    int ib = b1_crd[pb1];
    int ic = c1_crd[pc1];
    int id = d1_crd[pd1];
    int i = min(ib, ic, id);
    if (ib == i && ic == i && id == i) {
        a[i] = (b[pb1] + c[pc1]) * d[pd1];
    }
}
```

```
else if (ib == i && id == i) {
    a[i] = b[pb1] * d[pd1];
}
else if (ic == i && id == i) {
    a[i] = c[pc1] * d[pd1];
}
```

```
if (ib == i) pb1++;
if (ic == i) pc1++;
if (id == i) pd1++;
}
```

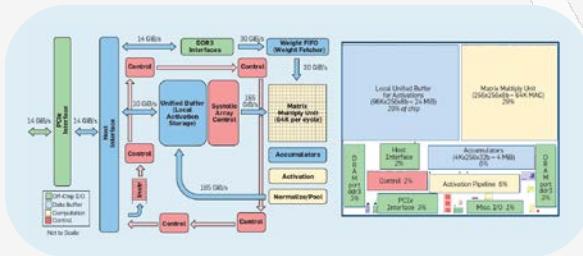
```
while (pc1 < c1_pos[1] && pd1 < d1_pos[1]) {
    int ic = c1_crd[pc1];
    int id = d1_crd[pd1];
    int i = min(ic, id);
    if (ic == i && id == i) {
        a[i] = c[pc1] * d[pd1];
    }
    if (ic == i) pc1++;
    if (id == i) pd1++;
}
```

```
while (pb1 < b1_pos[1] && pd1 < d1_pos[1]) {
    int ib = b1_crd[pb1];
    int id = d1_crd[pd1];
    int i = min(ib, id);
    if (ib == i && id == i) {
        a[i] = b[pb1] * d[pd1];
    }
    if (ib == i) pb1++;
    if (id == i) pd1++;
}
```

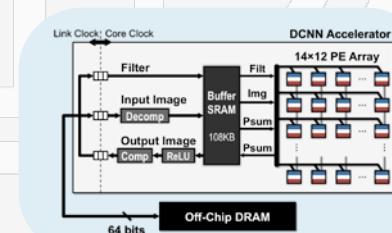
DSAS FOR ALL SPARSE TENSOR ALGEBRA

- Fixed-function Domain-Specific Architectures

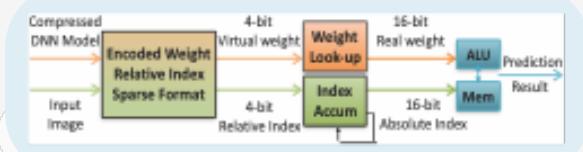
Google TPU



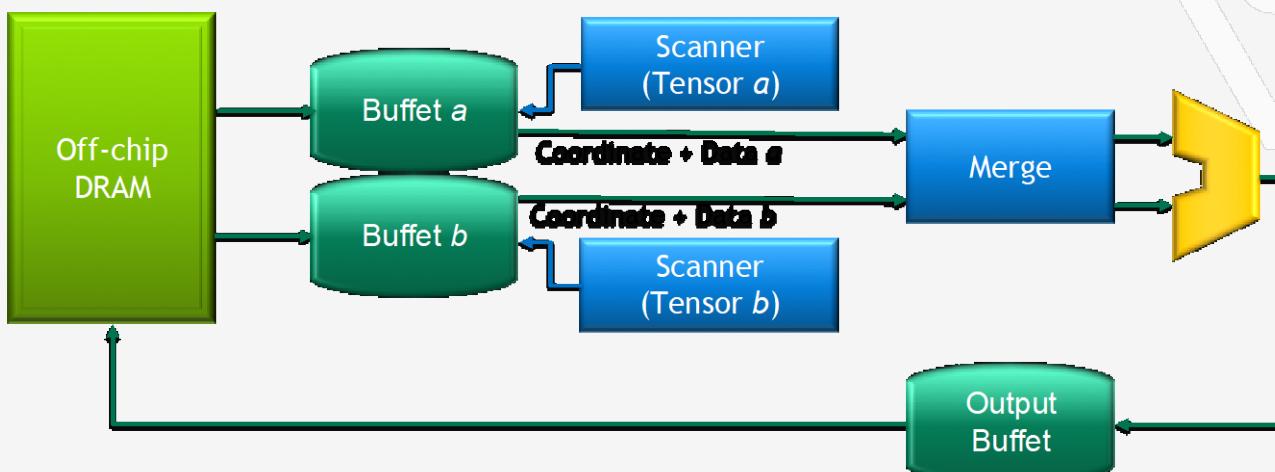
MIT Eyeriss



Stanford EIE

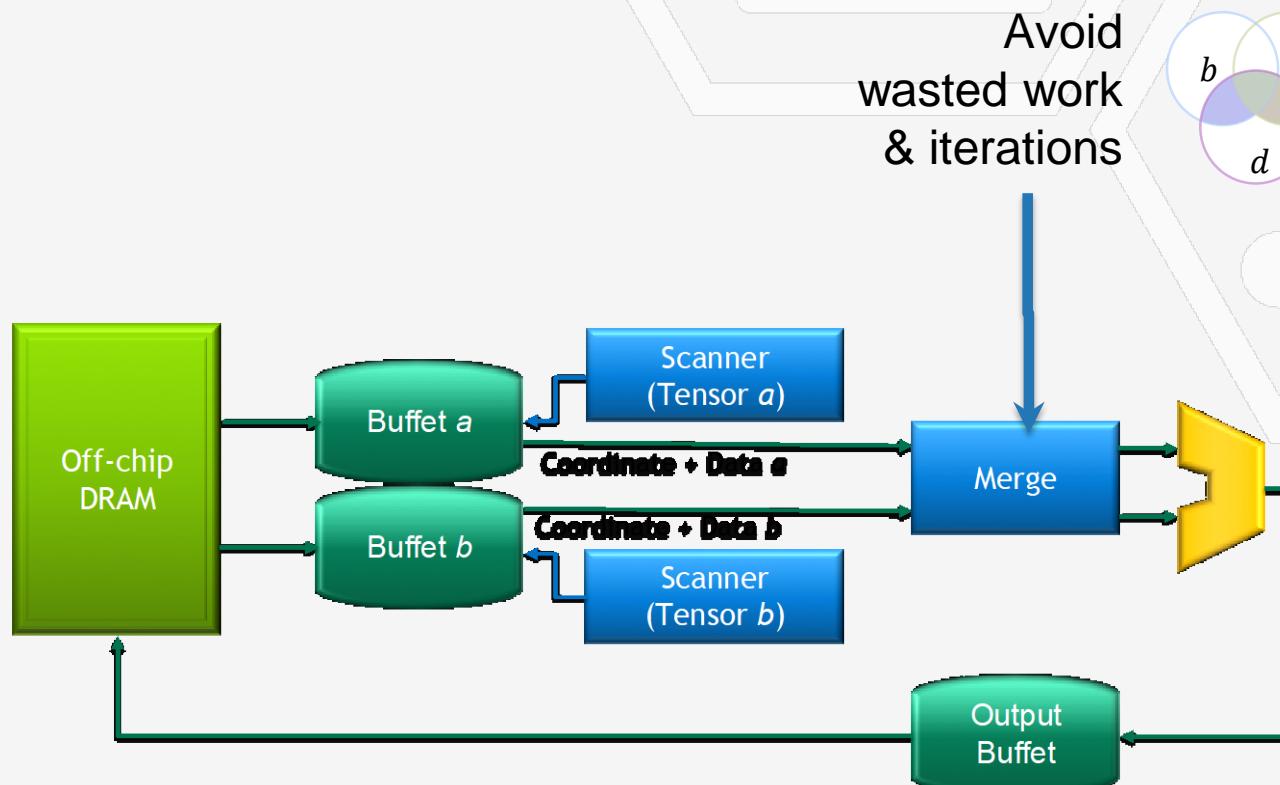


- Domain-Specific Architecture for all Sparse Tensor Algebra



Nvidia Symphony Architecture (DARPA SDH)

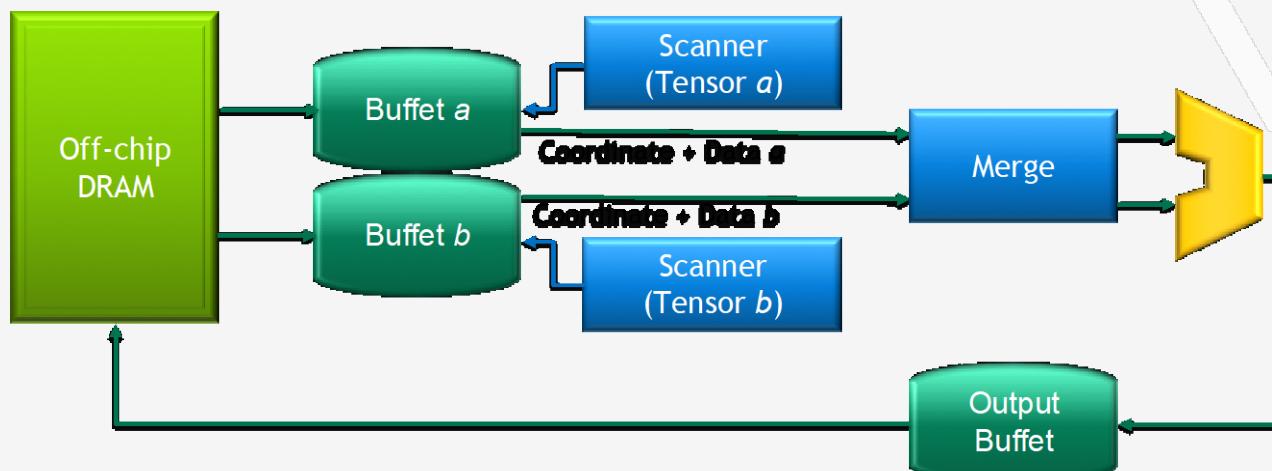
DSAS FOR ALL SPARSE TENSOR ALGEBRA



DSAS FOR ALL SPARSE TENSOR ALGEBRA

- TACO generates code for the Symphony Architecture
- Buffets, Scanners & Merge units provide acceleration
- Imperative C code for more complex outer loops

With TACO, Symphony can accelerate all tensor algebra expressions in any data format



$y = Ax + z$

union case #1

```
for (int py = 0; py < y1_dimension; py++) {  
    y[py] = 0.0;  
}  
  
int pA1 = A1_pos[0];  
int A1_end = A1_pos[1];  
int pz1 = z1_pos[0];  
int z1_end = z1_pos[1];  
while (pA1 < A1_end & pz1 < z1_end) {  
    int IA = A1_crd[pA1];  
    int iz = z1_crd[pz1];  
    if (IA == i & iz == j) {  
        double tj = 0.0;  
        int pA2 = A2_pos[pA1];  
        int pA3 = A2_pos[pz1];  
        y[i] = tj + z[zp1];  
    } else if (IA == i) {  
        do while (iz - n >= 0) {  
            int pA2 = A2_pos[pA1];  
            int pA3 = A2_pos[pz1];  
            px0 += (int)(jx0 == j0);  
            y[i] = tj0;  
        }  
    } else {  
        y[i] = z[zp1];  
    }  
    pA1 += (int)(IA == i);  
    pz1 += (int)(iz == j);  
}
```

union case #2

intersection in hardware

```
A2, x |  
A2, x |  
∅
```

union case #3

intersection in hardware

```
A2, x |  
A2, x |  
∅
```

intersection in hardware

```
A2, x |  
A2, x |  
∅
```

$y = A_2 x + z$

union case #2

```
while (pA1 < A1_end) {  
    int iA0 = A1_crd[pA1];  
    double tj1 = 0.0;  
    int pA21 = A2_pos[pA1];  
    int pA31 = A2_pos[pA1];  
    y[iA0] = tj1;  
    pA1++;  
}  
  
while (pz1 < z1_end) {  
    int iz0 = z1_crd[pz1];  
    y[iz0] = z[zp1];  
    pz1++;  
}
```

intersection in hardware

```
A2, x |  
A2, x |  
∅
```

CHALLENGES TO A HARMONIOUS FUTURE WITH DSLS

- Proliferation of DSLs
 - Each DSL is Costly to Build and Maintain
 - Common tools for building DSL
 - Delite from Stanford
 - Common backend from DSLs to DSAs
 - MIT Tiramisu compiler middle-end, similar to LLVM compiler backend for unicores
 - Interoperability between DSLs
 - Many programs will require multiple DSLs to implement
 - Ex: Multi-physics
 - A Unified Programming Interface: Abstraction without Friction
 - Python glue with common C data structures (NumPy)
- Quality of Life of Programmers
 - Debuggers, package managers, IDE support, Foreign Function Interfaces
 - Documentation, libraries, examples and a community

GraphIt
Graph
Analytics

TACO
Tensor
Algebra

Condensa
Neural Networks

Simlt
Physical
Simulation

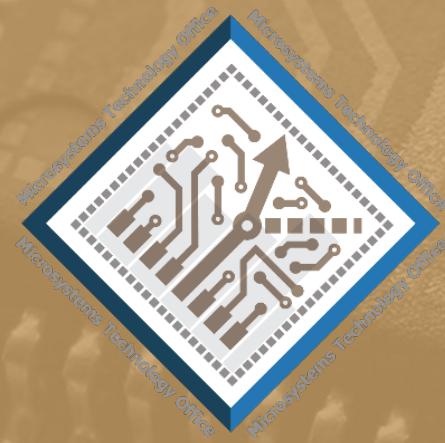
A LOT MORE PERFORMANCE LEFT TO MINE

- Current software practices are not performance efficient
 - Moore's Law era mindset of excess resources
- Case of Matrix Multiply

Implementation	Time (s)	GFLOPS	Speedup
1 Python	25,552.48	0.005	1

CONCLUSION

- Domain Specific Architectures are the future of performance
- Domain Specific Languages are essential to DSAs
 - Make DSAs generally available for the entire domain
 - Provide high performance
- Symphony Project is building multiple DSLs
 - TACO: Sparse Tensor Algebra
 - GraphIt: Graph Analytics
 - Condensa: Neural Networks
- Looking for other DSL opportunities



ERI ELECTRONICS RESURGENCE INITIATIVE

S U M M I T

2019 | Detroit, MI | July 15 - 17