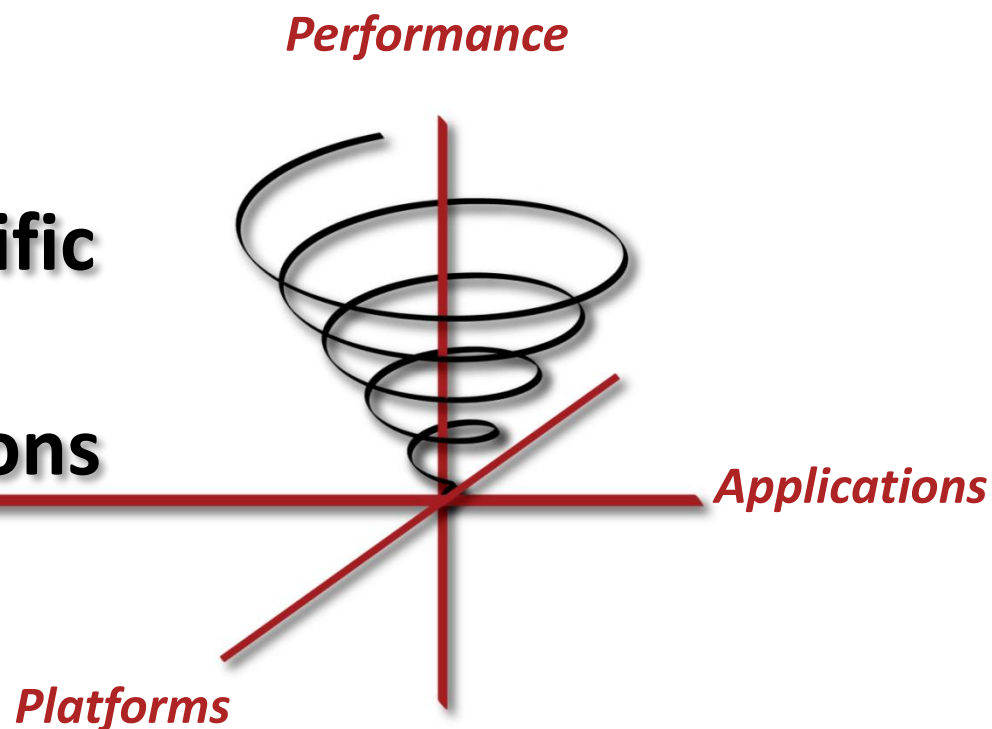


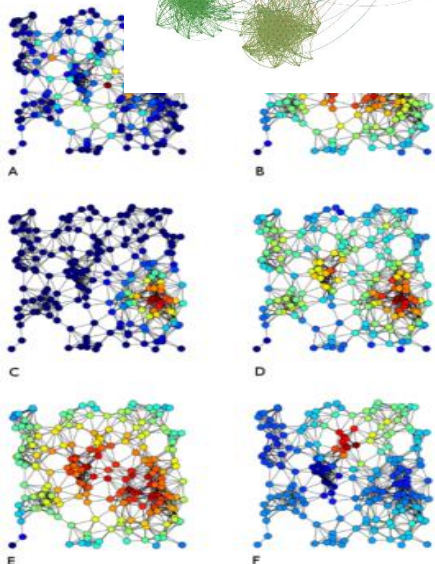
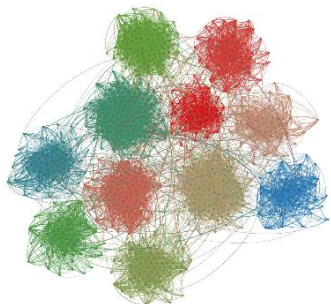
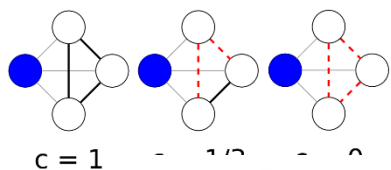
An Algorithmic Specific Code Generator for GEMM-Like Operations

Richard Michael Veras

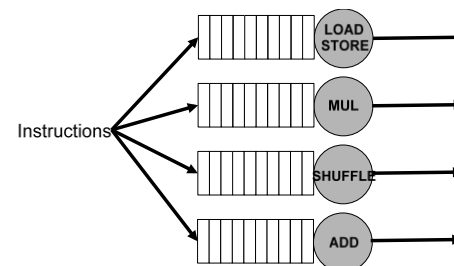
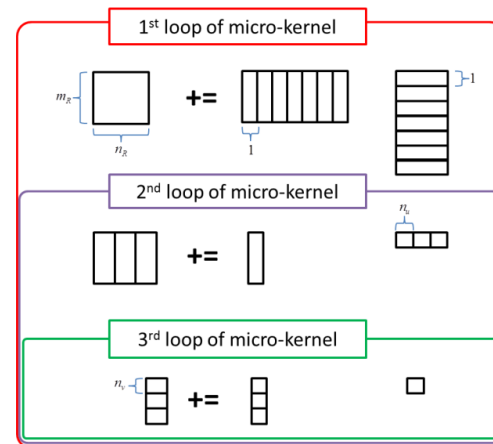


Want Automatic High Performance

GEMM-Like Operations

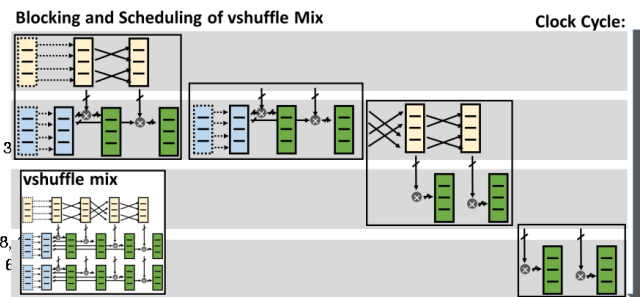
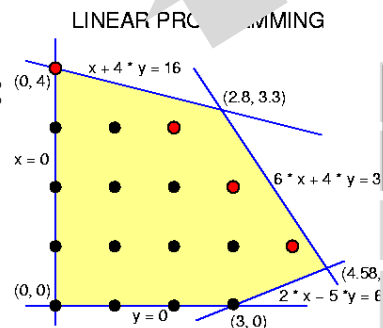


Model Driven approach for generating DGEMM:



High Perf.

Compiler Techniques:



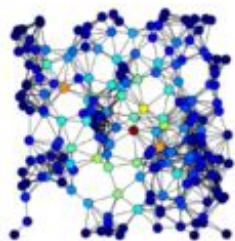
Function to maximize: $f(x, y) = 6 * x + 5 * y$
 Optimum LP solution $(x, y) = (2.4, 3.4)$
 Pareto optima: $(0, 4), (2, 3), (3, 2), (4, 1)$
 Optimum ILP solution $(x, y) = (4, 1)$

GEMM Like Operations

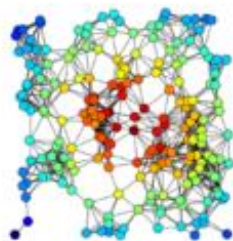
Centrality: $\sim A^T A$

[betweenness]

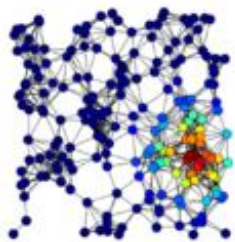
(Z, +, MIN, 0, 1)



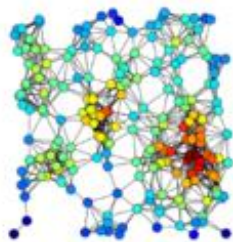
A



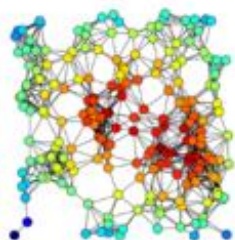
B



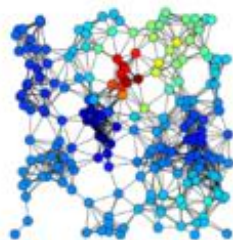
C



D



E

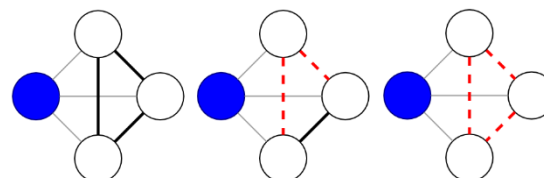


F

Clustering: $\sim AA^T A$

(triangle)

(Z, \wedge , +, 1, 0)



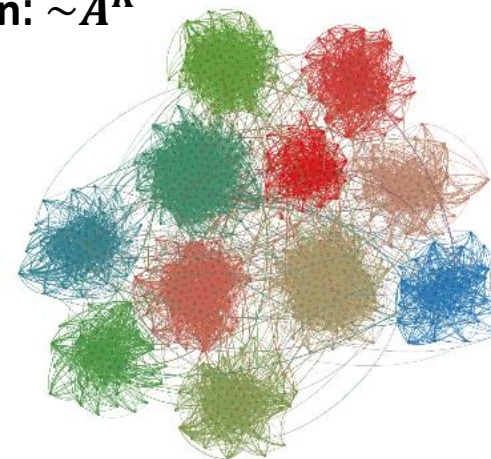
$$c = 1$$

$$c = 1/3$$

$$c = 0$$

Community

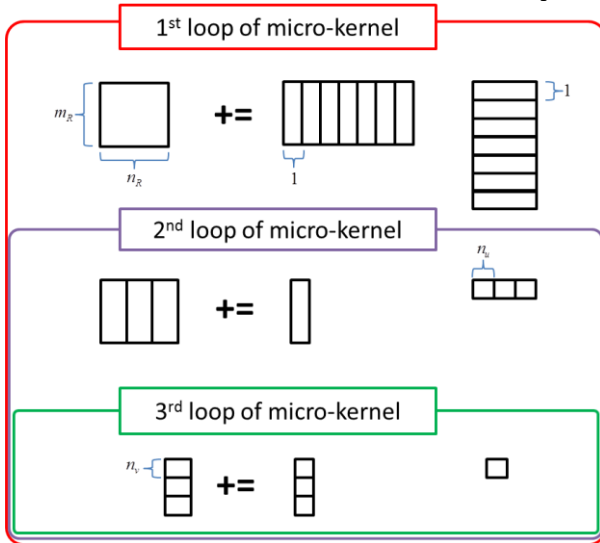
Detection: $\sim A^K$



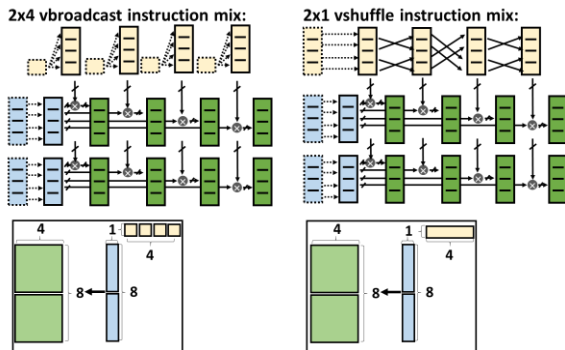
Check out GraphBLAS

High Performance Micro-Kernels

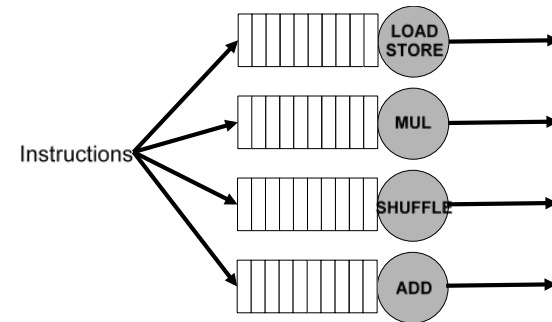
Cast micro kernel as outer product:



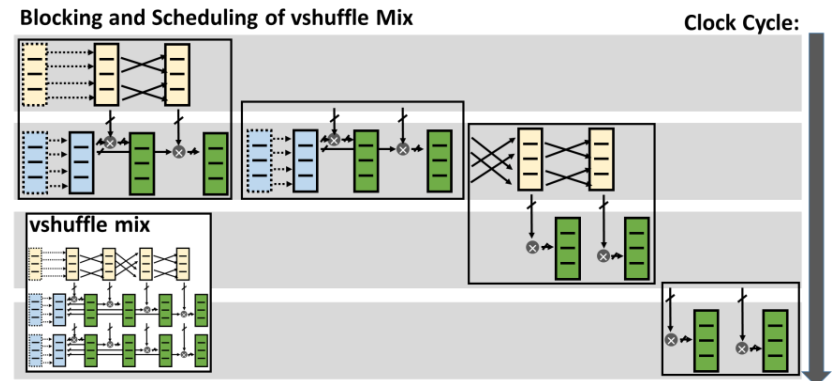
Enumerate all possible tilings given ISA



Use Models to Select from Design Space:



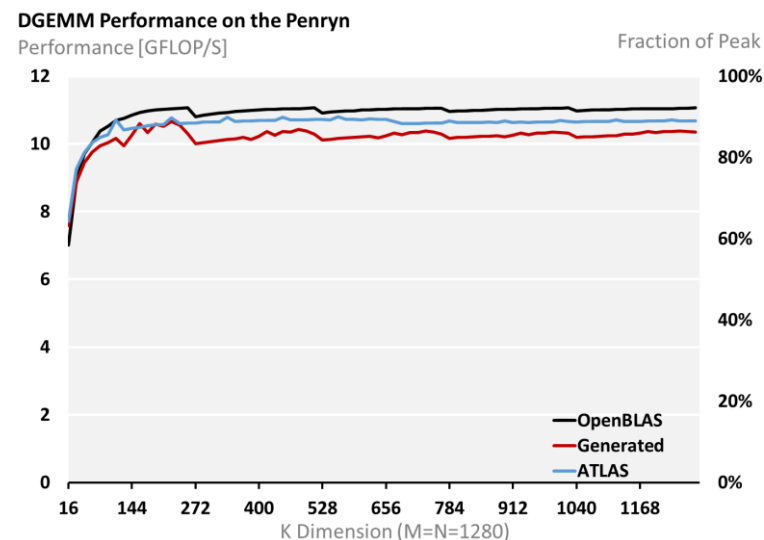
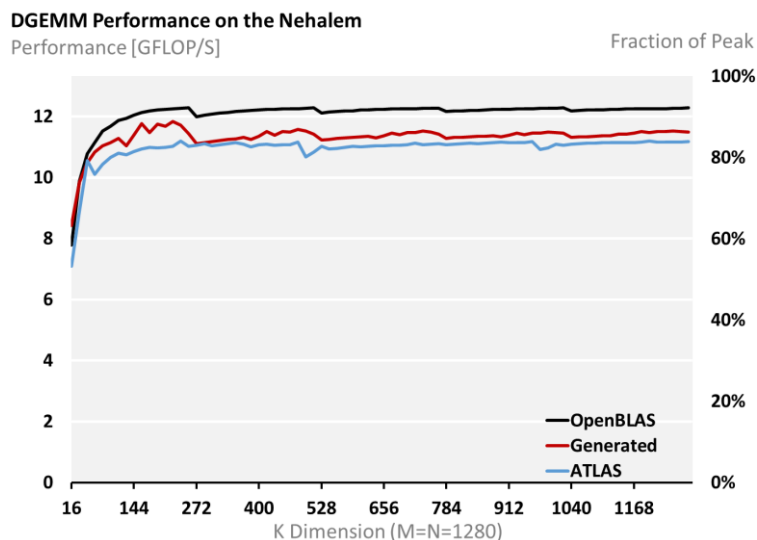
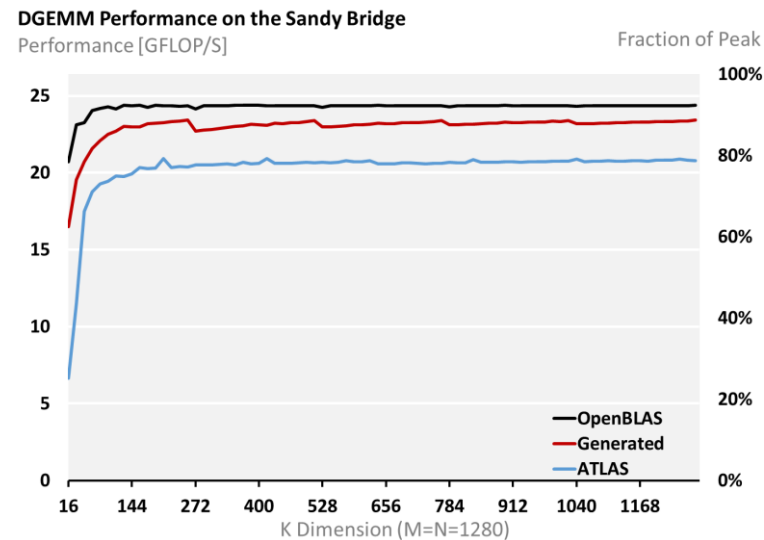
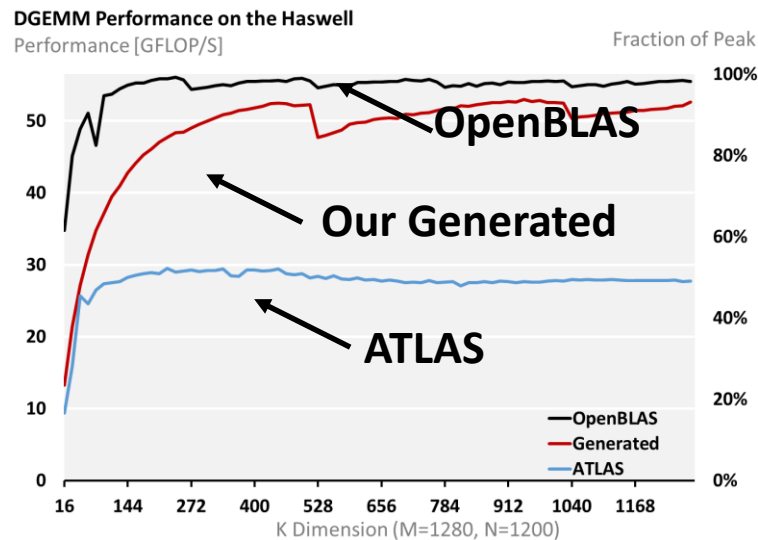
Aggressively Schedule and Optimize:



Veras, R., Smith T., Low T.M., Franchetti, F. van de Geijn, R. [CGO 2017 Submitted]

Richard Veras (rveras@cmu.edu)

High Performance Micro-Kernels

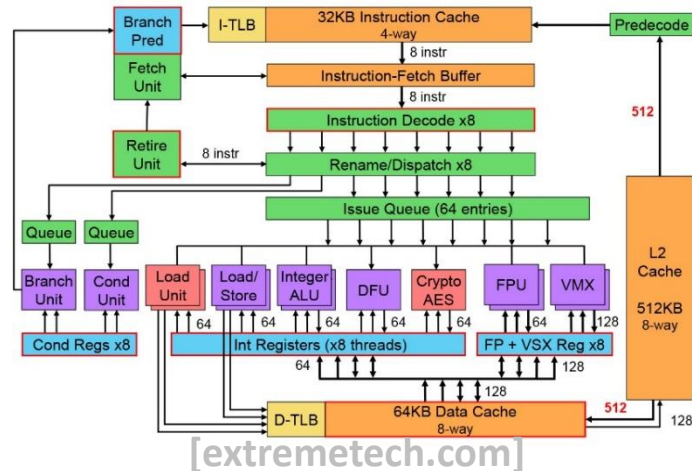


Veras, R., Smith T., Low T.M., Franchetti, F. van de Geijn, R. [CGO 2017 Submitted]

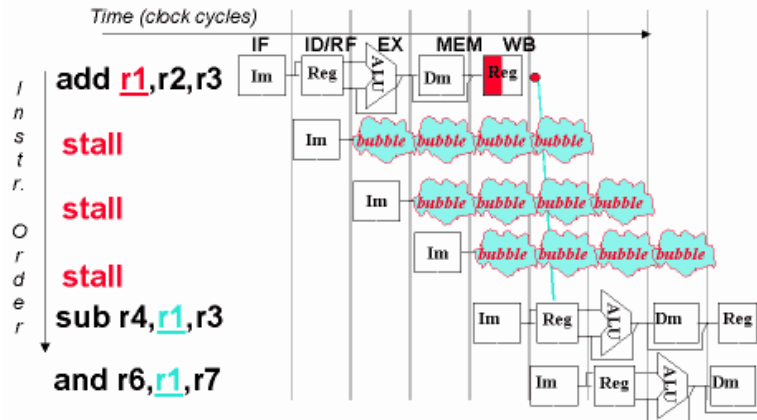
Richard Veras (rveras@cmu.edu)

Automating with Compiler Techniques

Subtle Semiring Changes Impact Performance:

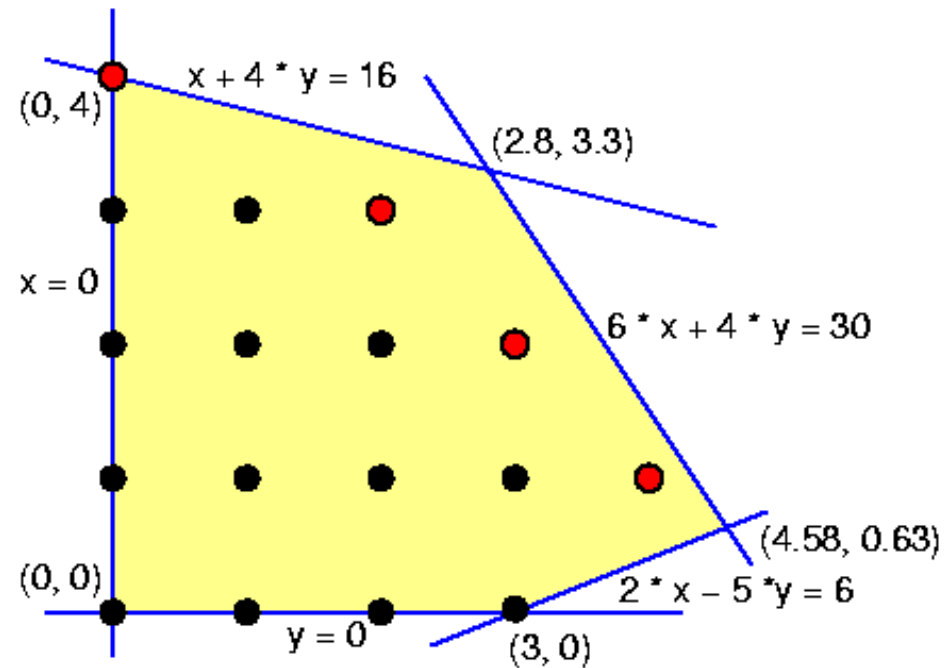


Minimizing Stalls



Cast as ILP Problem

LINEAR PROGRAMMING



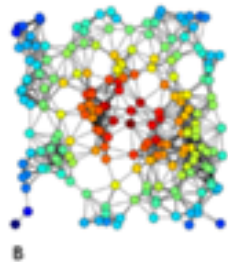
Function to maximize: $f(x, y) = 6 * x + 5 * y$
 Optimum LP solution $(x, y) = (2.4, 3.4)$
 Pareto optima: $(0, 4), (2, 3), (3, 2), (4, 1)$
 Optimum ILP solution $(x, y) = (4, 1)$

[cs.duke.edu]

A Generator for GEMM-Like Kernels

Input:

Betweenness Centrality
(Floyd-Warshall)



GEMM Like:

$$C \leftarrow AB + C$$

$$C_{acc} \leftarrow 0 \quad \text{Initialize}$$

$$C_{acc} \leftarrow \sum_p ab^T \quad \text{Compute}$$

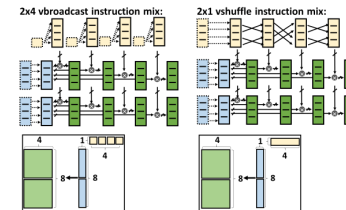
$$C \leftarrow f(C_{acc}) \quad \text{Accumulate}$$

Semiring

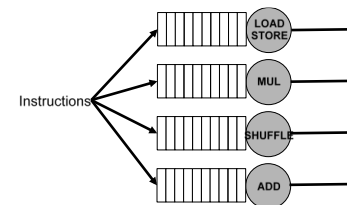
$$(Z, +, MIN, 0, 1)$$

Output:

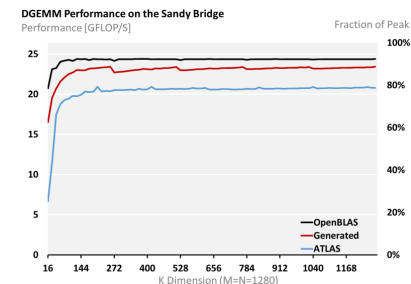
Kernel Algorithm from
our design space:



Kernel Sustains High
Throughput:

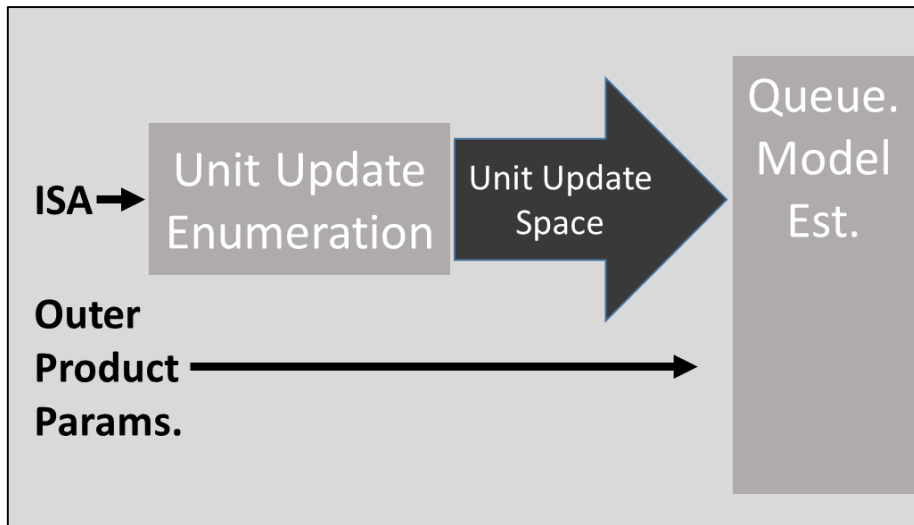


Tuned to the Target
Architecture

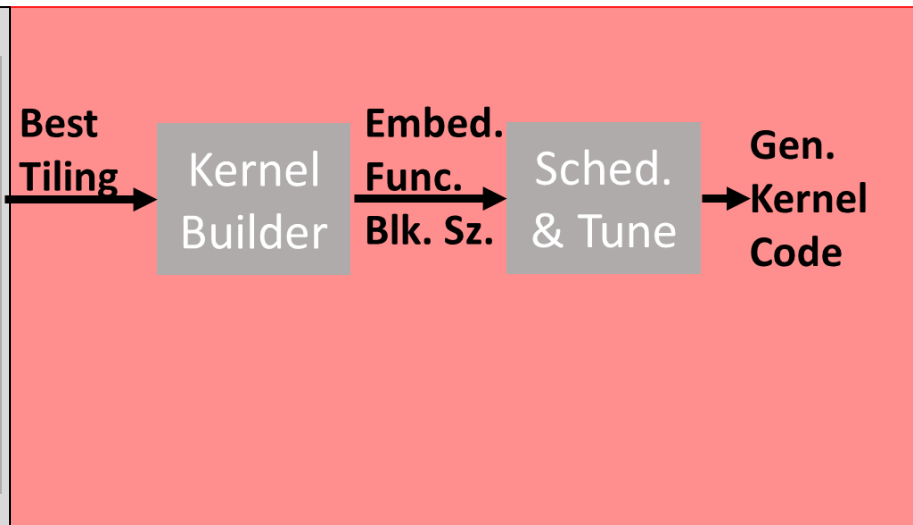


Our GEMM Generator Pipeline

Find Efficient Instructions Mix (Algo)



Turn Mix into Efficient Code (Implementation)



User defines
block size, ISA
and semiring

Enumerate
algorithm space

Top Candidate
selected

Template
Created for
Candidate

Template
transformed into
optimized and
scheduled code

From Math to Tiling

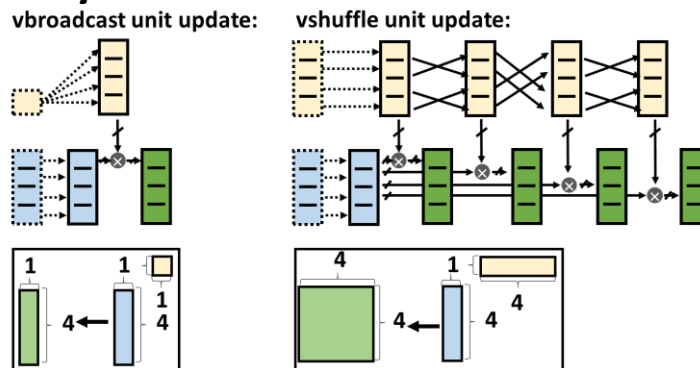
$$C \leftarrow AB + C$$

$$C_{acc} \leftarrow 0 \quad \text{Initialize}$$

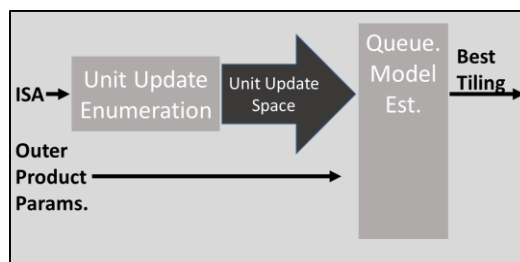
$$C_{acc} \leftarrow \sum_p ab^T \quad \text{Compute}$$

$$C \leftarrow f(C_{acc}) \quad \text{Accumulate}$$

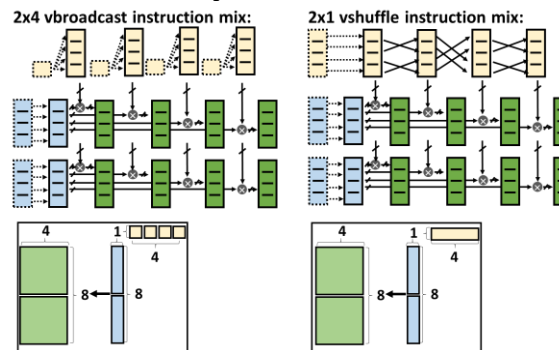
Identify Small Outer Products from ISA



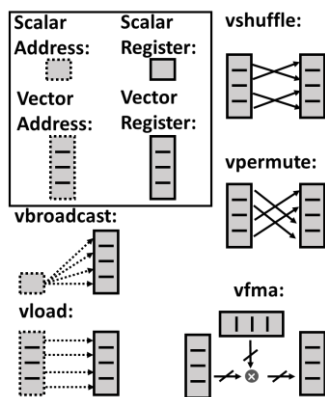
A High Throughput Mix



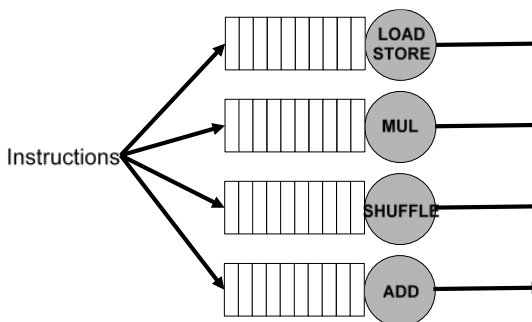
Enumerate Space of Outer Products



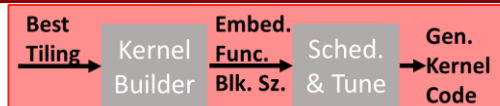
Start with ISA



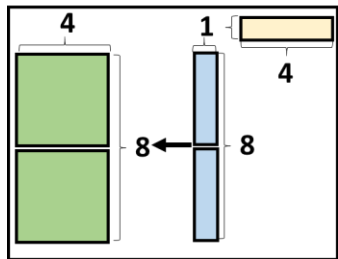
Select Best Mix with Queueing Model



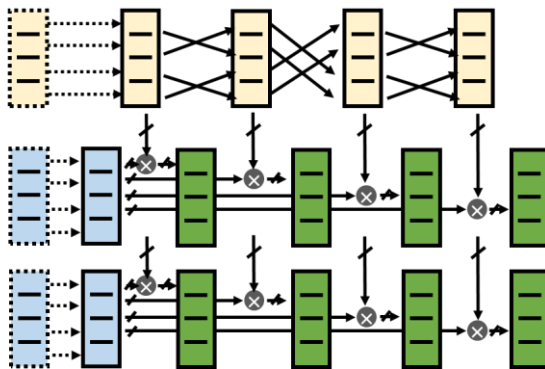
From Tiling to Template



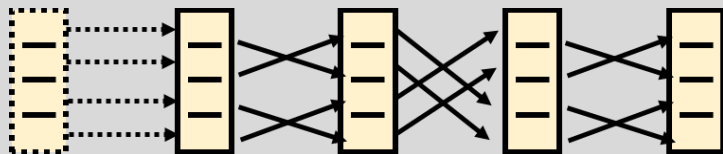
Selected Outer Product:



2x1 vshuffle instruction mix:



Embedding Function (get_b):



```
def get_b_element( var array b_reg[][],  
                  ptr B, ii, jj )
```

```
  opts = {  
    0: assign( b_reg[jj], vload(B,jj)),  
    1: assign( b_reg[jj], shuffle(b_reg[jj-1]) ),  
    2: assign( b_reg[jj], permute(b_reg[jj-1]) ),  
    3: assign( b_reg[jj], shuffle(b_reg[jj-1]) ) }  
  if ii mod v = 0  
    return opts[jj]
```

$$C_{acc} \leftarrow 0$$

```
for( i = 0; i < m_r; i++ )  
  for( j = 0; j < n_r; j++ )  
    init(c_reg, ii, jj )
```

$$C_{acc} \leftarrow \sum_p ab^T$$

```
for( pp = 0; pp < k_b; pp++ )  
  /* perform the outer products */  
  for( i = 0; i < m_r; i+=m_s )  
    for( j = 0; j < n_r; j+=n_s )  
      for( ii = i; ii < i+m_s; ii++ )  
        get_a_elem(a_reg, ii, j );  
        for( jj = j; jj < j+n_s; jj++ )  
          get_b_elem(b_reg, ii, jj );  
          apply(c_reg, a_reg, b_reg, ii, jj, pp);
```

$$C \leftarrow f(C_{acc})$$

```
for( i = 0; i < m_r; i++ )  
  for( j = 0; j < n_r; j++ )  
    accumulate(C, c_reg, ii, jj );
```

Floyd-Warshall Embedded in GEMM

DGEMM

Floyd-Warshall

Semiring:

$(R, *, +, 1, 0)$

$(Z, +, MIN, 0, \infty)$

Initialize:

$C_{acc} \leftarrow 0$

```
assign( c_reg[ii,jj], 0)
```

```
assign( c_reg[ii,jj],  
(ii==jj)? 0 : INFINITY)
```

Compute:

$C_{acc} \leftarrow \sum_p ab^T$

```
assign( c_reg[ii,jj],  
a_reg[ii,pp]*  
b_reg[pp,jj],  
+c_reg[ii,jj]))
```

```
assign( c_reg[ii,jj],  
MIN(a_reg[ii,pp]+  
b_reg[pp,jj],  
c_reg[ii,jj]))
```

Accumulate:

$C \leftarrow f(C_{acc})$

```
assign( C[(ii,jj)],  
c_reg[ii,jj]+  
C[(ii,jj)])
```

```
assign( C[(ii,jj)],  
MIN(c_reg[ii,jj],  
C[(ii,jj)]) )
```

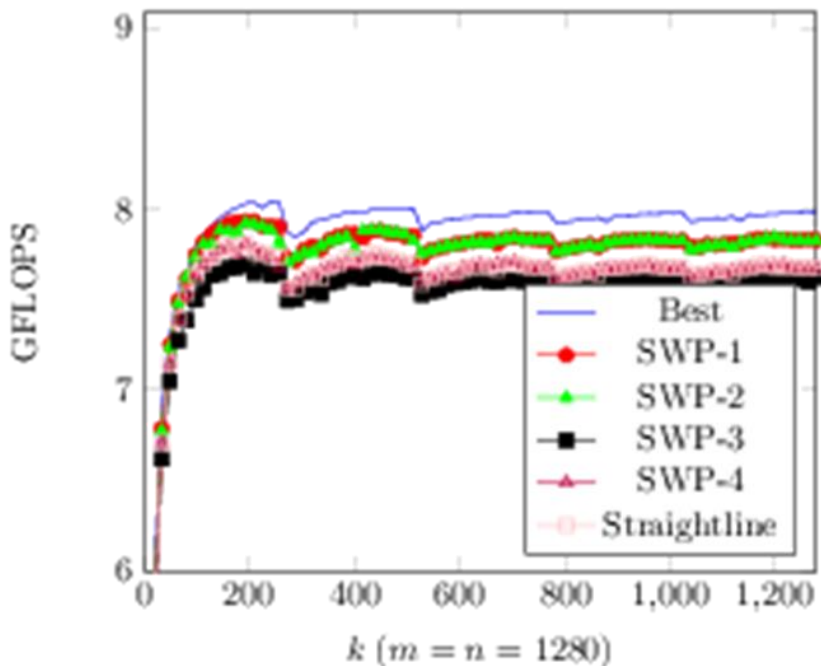
Scheduling the Problem

We have built the kernel code,
Now we need to schedule:



Static Scheduling still matters
on OOO Processors:

Various Schedules for Nehalem (gcc-macro)



Pipeline for Scheduling:
Built Kernel +
Embedding Func.

Express as
Decision Vars

Formulate
Constraints
over Vars

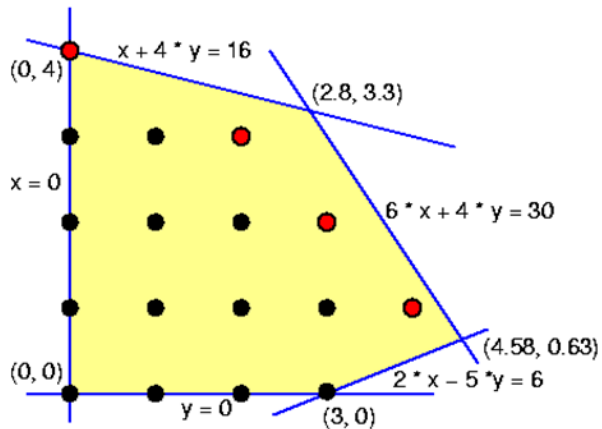
Minimize with
ILP Solver

Scheduled Kernel Code

OASIC approach for ILP Scheduling

Representing Design Space as
Polytope:

LINEAR PROGRAMMING



Decision Variable: **instruction n**
is executed on **functional unit k**
at **time step t**

x_{nt}^k ← Functional unit
← cycle
← Instruction label

Instruction label

Expressing Constraints in
terms of X:

Every instruction **n** is executed **once**

$$\sum_k \sum_t x_{nt}^k = 1$$

At any timestep t, functional unit
k is used no more than it can

$$\sum_k x_{nt}^k \leq R_k$$

If t_m depends on t_n , then t_m will
not execute until l cycles after t_n

$$\sum_k \sum_t x_{nt_m}^k + \sum_k \sum_t x_{nt_m}^k \leq 1$$

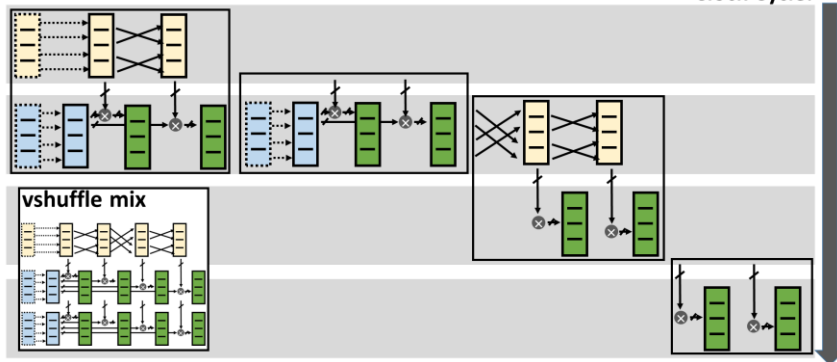
But wait, there's more!

Emitting The Code

Code is now Scheduled

Blocking and Scheduling of vshuffle Mix

Clock Cycle:



Need Custom ANSI C compliant SIMD wrappers to schedule in Compiler:

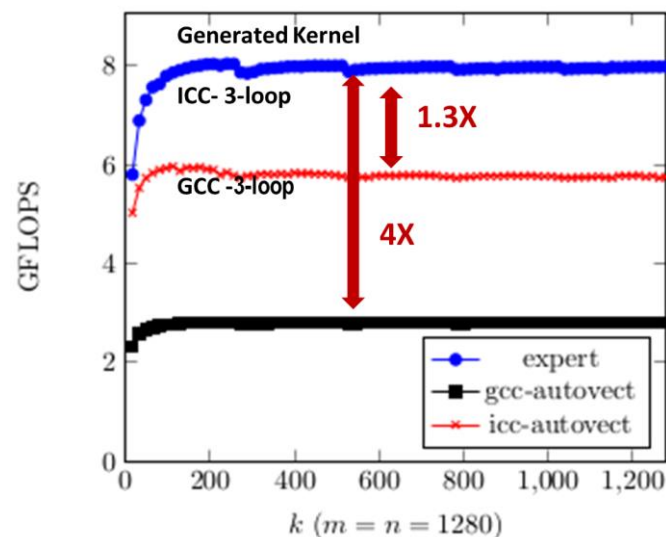
```
#define VADD(srca, srcb, dest)
asm volatile(
    "vaddpd %[vsrca], %[vsrcb],
        %[vdest]"
    : [vdest] "=x" (dest)
    : [vsrca] "x" (srca),
      [vsrcb] "x" (srcb);
```

Code is emitted:

```
for( pp = 0; pp < k_b; pp+=KUNR )
/* STEADY STATE CODE */
VLOAD_IA(GET_A_ADDR(0), GET_A_REG(0))
VLOAD_IA(GET_A_ADDR(1), GET_A_REG(1))
VLOAD_IA(GET_B_ADDR(0), GET_B_REG(0))
VSHUFFLE_IA(GET_B_REG(0), GET_B_REG(1))
VFMA(GET_A_REG(0),
    GET_B_REG(0), GET_C_REG(0, 0))
VFMA(GET_A_REG(0), GET_B_REG(1), GET_C_REG(0, 1))
VPERM2F128_IA(0x01, GET_B_REG(1), GET_B_REG(2))
VSHUFFLE_IA(0x05, GET_B_REG(2), GET_B_REG(3))
VFMA(GET_A_REG(1), GET_B_REG(0), GET_C_REG(0, 0))
VFMA(GET_A_REG(1), GET_B_REG(1), GET_C_REG(0, 1))
```

Is this necessary?

Compiler Autovect versus Expert Nehalem



Putting it All Together

Have Operation that we can express like GEMM:

$$C \leftarrow AB + C$$

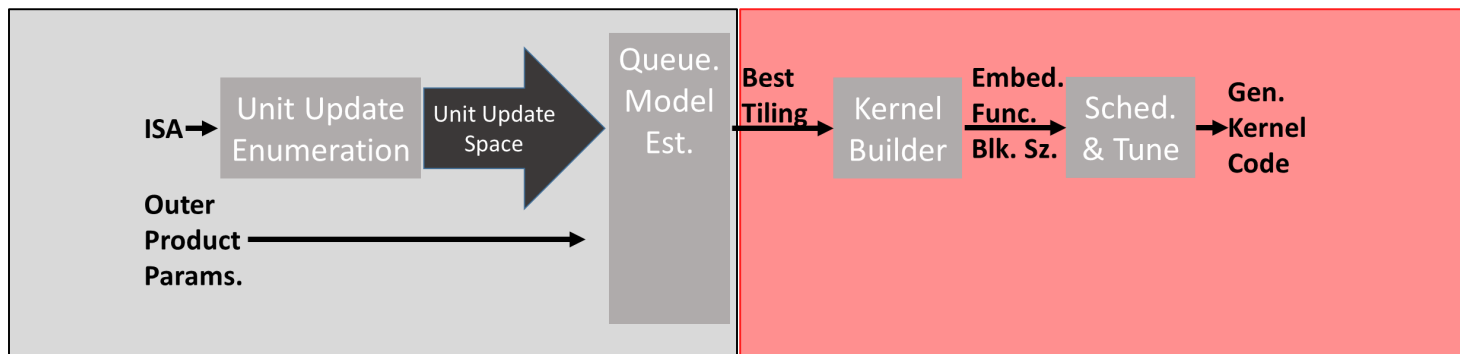
$$C_{acc} \leftarrow 0 \quad \text{Initialize}$$

$$C_{acc} \leftarrow \sum_p ab^T \quad \text{Compute}$$

$$C \leftarrow f(C_{acc}) \quad \text{Accumulate}$$

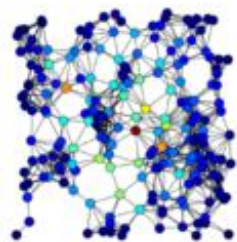
	DGEMM	Floyd-Warshall
Semi-ring:	$(Double, *, +, 1, 0)$	$(Z, +, MIN, 0, 1)$
Initialize: $C_{acc} \leftarrow 0$	<code>assign(c_reg[ii,jj], 0)</code>	<code>assign(c_reg[ii,jj], (ii==jj)? 0 : INFINITY)</code>
Compute: $C_{acc} \leftarrow \sum_p ab^T$	<code>assign(c_reg[ii,jj], a_reg[ii,pp]* b_reg[pp,jj], +c_reg[ii,jj])</code>	<code>assign(c_reg[ii,jj], MIN(a_reg[ii,pp]+ b_reg[pp,jj], c_reg[ii,jj]))</code>
Accumulate: $C \leftarrow f(C_{acc})$	<code>assign(C[(ii,jj)], c_reg[ii,jj]+ C[(ii,jj)])</code>	<code>assign(C[(ii,jj)], MIN(c_reg[ii,jj], C[(ii,jj)]))</code>

Run through our Pipeline:

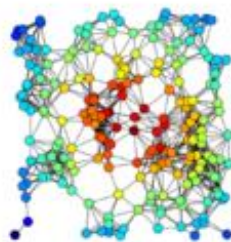


Moving Forward

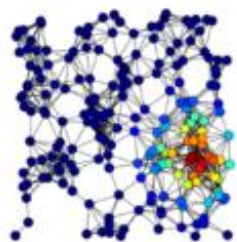
Centrality: $\sim A^T A$
(betweenness)



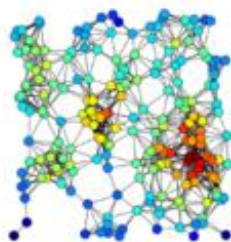
A



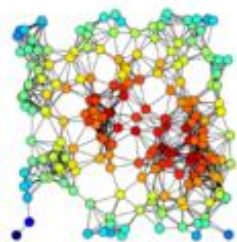
B



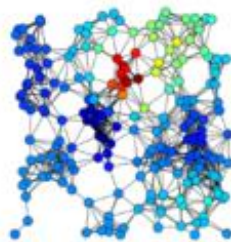
C



D

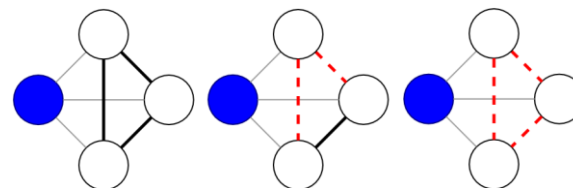


E



F

Clustering: $\sim AA^T A$
(triangle)

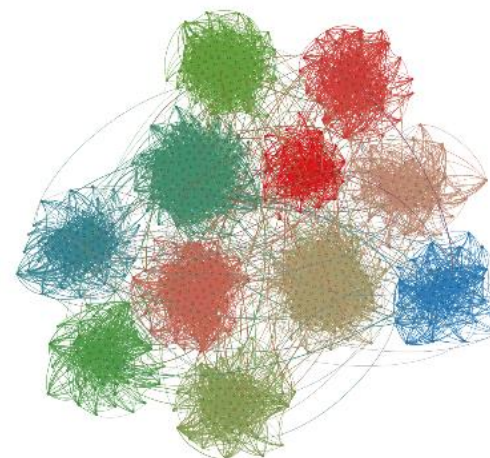


$$c = 1$$

$$c = 1/3$$

$$c = 0$$

Community
Detection: $\sim A^K$



Summary

- There exists a **large class** of **GEMM-like Operations**
- Obtaining DGEMM level **performance** for each of these operations **requires automation**
- We have a **systematic approach** for automatically **generating DGEMM**
- We are extending it by allowing the user to **define a semi-ring** with an initialize and accumulate function