

# T2S-Tensor: Productively Generating High-Performance Spatial Hardware for Tensor Computations

Nitish Srivastava<sup>1</sup>, Hongbo Rong<sup>2</sup>, Prithayan Barua<sup>3</sup>, Guanyu Feng<sup>4</sup>, Huanqi Cao<sup>3</sup>, Zhiru Zhang<sup>1</sup>, David Albonesi<sup>1</sup>, Vivek Sarkar<sup>3</sup>, Wenguang Chen<sup>3</sup>, Paul Petersen<sup>2</sup>, Geoff Lowney<sup>2</sup>, Adam Herr<sup>2</sup>, Christopher Hughes<sup>2</sup>, Timothy Mattson<sup>2</sup>, Pradeep Dubey<sup>2</sup>

<sup>1</sup>Cornell University

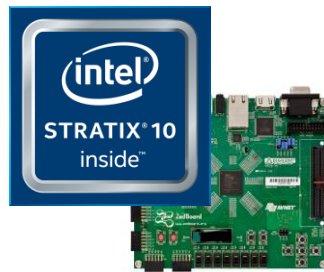
<sup>2</sup>Intel Parallel Computing Labs

<sup>3</sup>Georgia Institute of Technology

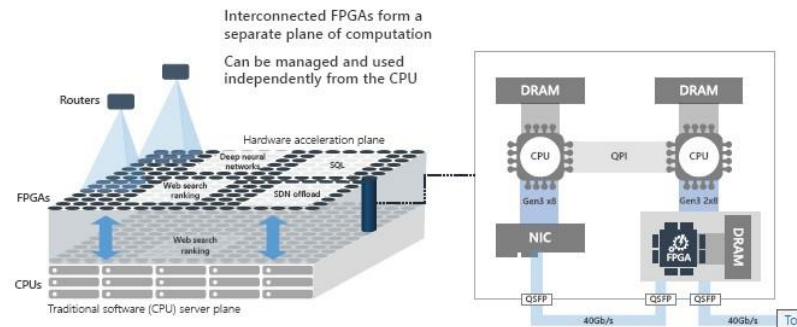
<sup>4</sup>Tsinghua University

# Tensor Computations

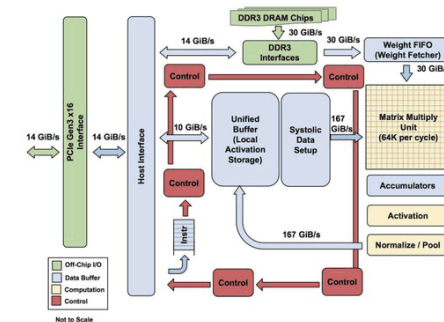
- ▶ Tensor computations are everywhere
  - Data analytics (e.g. graphs in social interactions)
  - Machine learning (e.g. CNN, recommender systems)
  - Scientific computing (e.g. finite element method)
- ▶ Increasing popularity of deploying tensor applications to spatial architectures



FPGAs



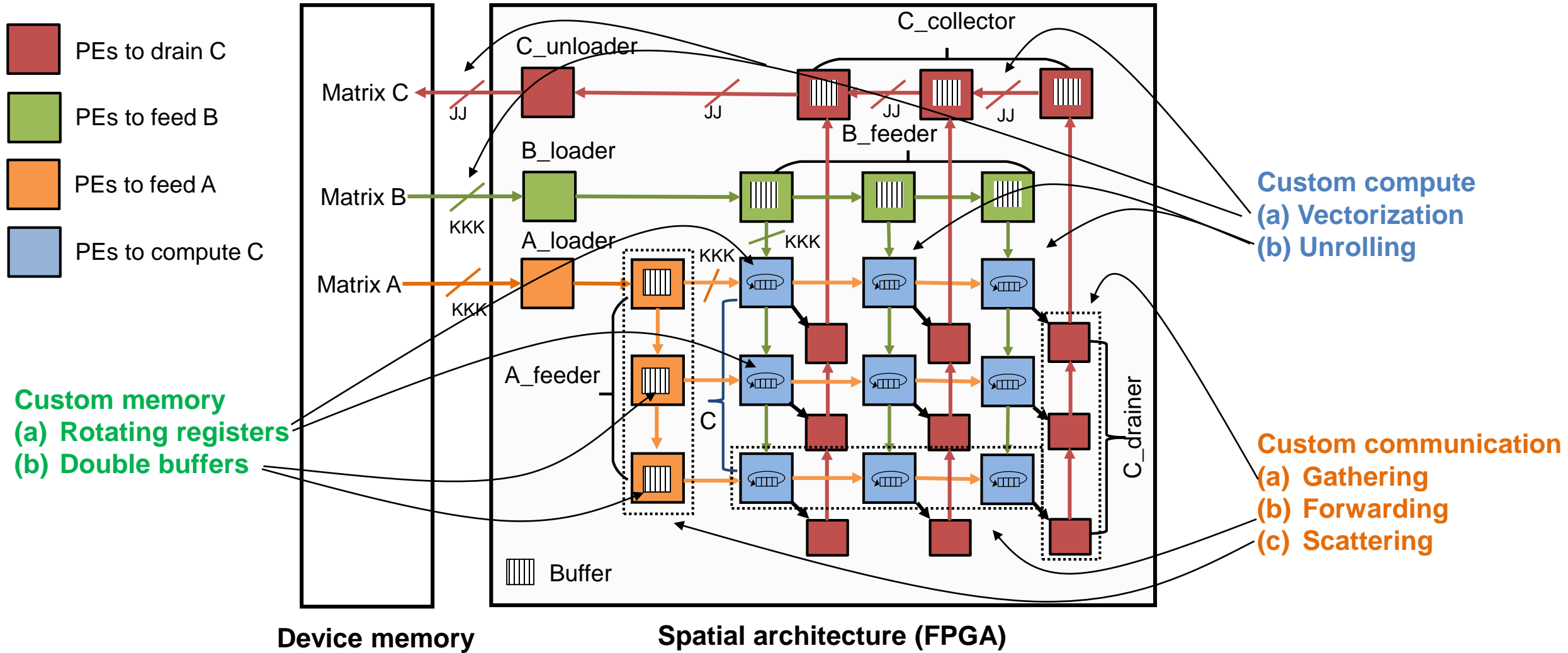
Microsoft Brainwave



ASICs (Google TPU)

- ▶ There exists a gap on how to map these applications on different spatial architectures

# Driving Example – Matrix Multiplication (GEMM)



High-performance GEMM design on FPGA



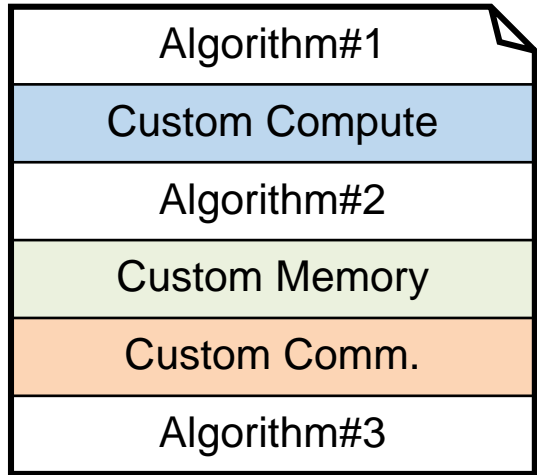






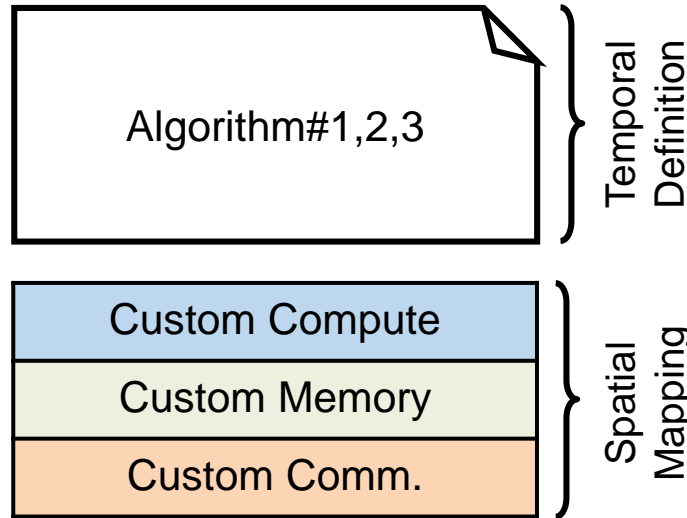
# Decoupling Hardware Customization and Algorithm

HLS C



Entangled algorithm specification & customization schemes [1,2,3]

T2S



Decoupled customization & clean abstraction

```

C(i, j) = 0;
C(i, j) += A(i, k) * B(k, j);
C.update().tile(k, j, i, kk, jj, ii, KK, JJ, II);
    .isolate_producer_chain(A, A_loader, A_feeder)
    .isolate_producer_chain(B, B_loader, B_feeder)
    .isolate_consumer_chain(C, C_drainer, C_unloader);
A_loader.unroll(ii).remove(jj).vload(kk);
A_feeder.buffer(ii, Buffer::Double).unroll(ii);
B_loader.unroll(jj).remove(ii).vload(kk);
B_feeder.buffer(k, Buffer::Double).unroll(jj);
C.update().unroll(jj, ii)
    .forward(A_feeder, {1, 0}) .forward(B_feeder, {0, 1});
C_drainer.unroll(jj, ii).gather(C, {1, 0})
C_unloader.buffer(ii).unroll(ii).vstore(jj);
    
```

Temporal to Spatial → T2S

T2S is an extension over Halide for spatial architectures

- [1] Intel HLS
- [2] Xilinx Vivado HLS
- [3] Canis, et al. FPGA'11



# Temporal Definition in T2S

Func C

$C(i, j) = 0$

$C(i, j) += A(i, k) * B(k, j)$

$C.\text{tile}(i,j,k,ii,jj,kk,ll,JJ,KK)$

Algorithm

C

for i, j, k

for ii, jj, kk

$i' = i * ll + ii$

$j' = j * JJ + jj$

$k' = k * KK + kk$

$C[i', j'] += A[i', k'] * B[k', j']$



C

# Spatial Mapping in T2S

```
Func C
C(i, j) = 0
C(i, j) += A(i, k) * B(k, j)
C.tile(i,j,k,ii,jj,kk,ll,JJ,KK)
C.isolate_producer(A, A_feeder)
```

Spatial  
Algorithm  
Mapping

**C**

```
for i, j, k
  for ii, jj, kk
    i',j',k' = ...
    C[i', j'] += A[i',k'] * B[k',j']
```



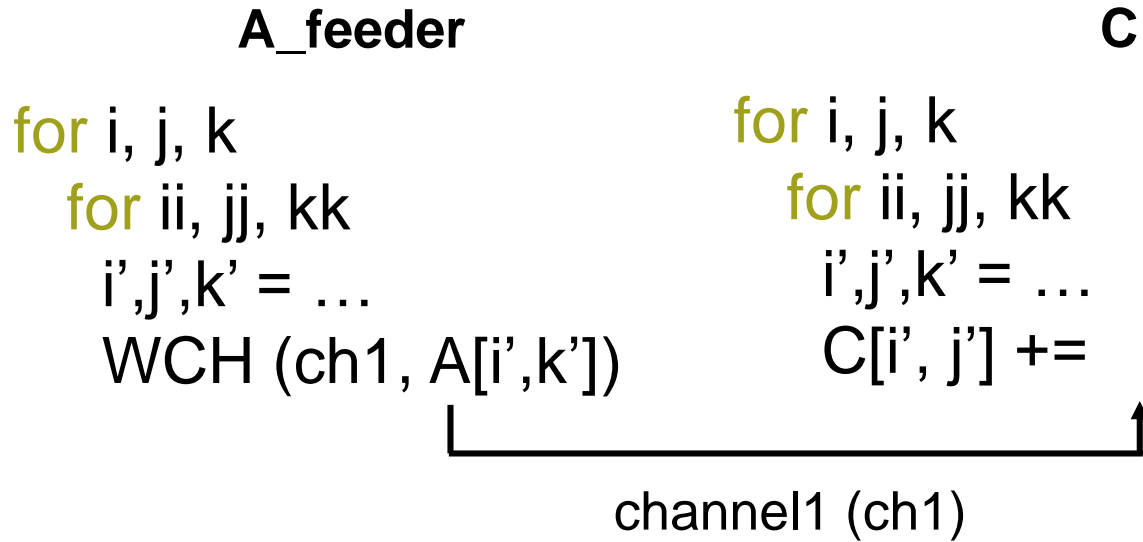
**C**

# Spatial Mapping in T2S

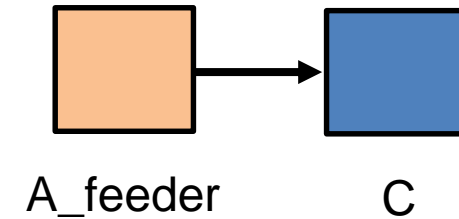
```

Func C
C(i, j) = 0
C(i, j) += A(i, k) * B(k, j)
C.tile(i,j,k,ii,jj,kk,II,JJ,KK)
C.isolate_producer(A, A_feeder)
  
```

Spatial Mapping Algorithm

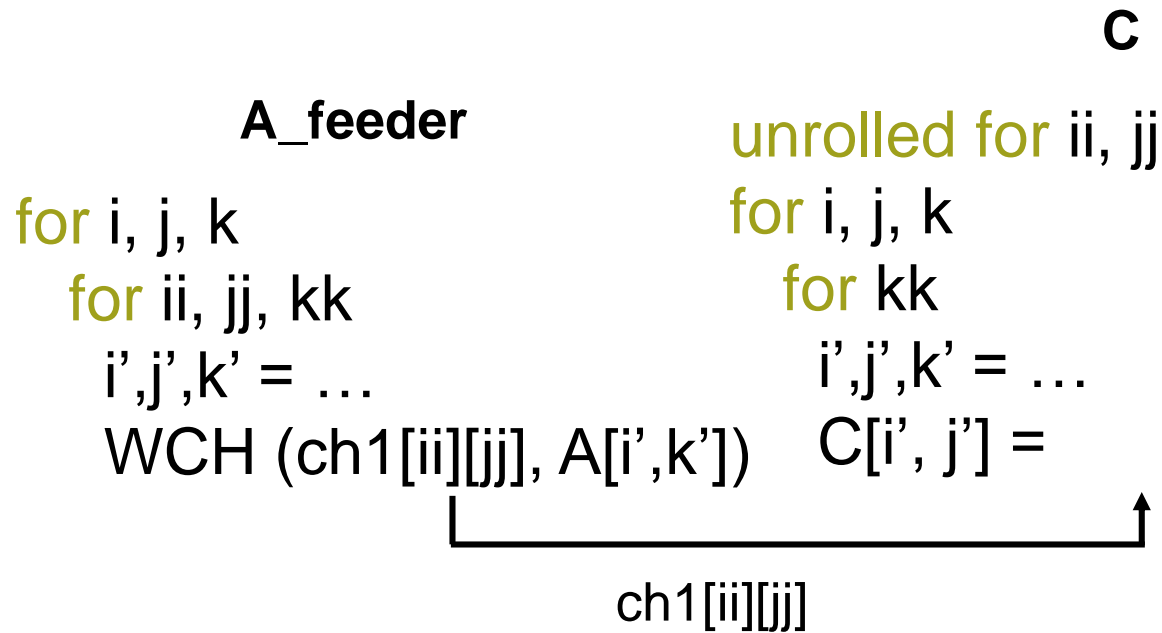


\* B[k',j']

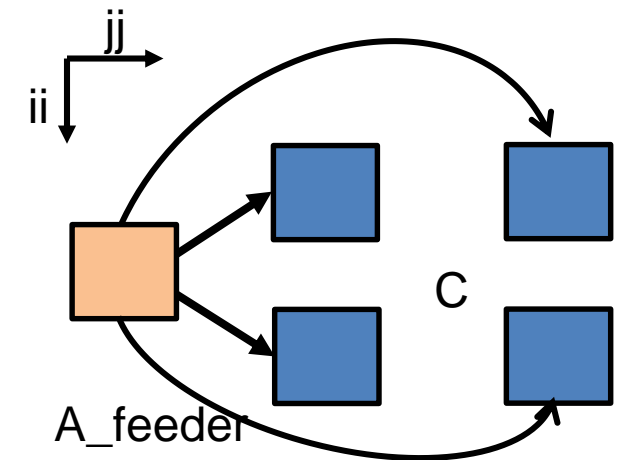


# Spatial Mapping in T2S

Func C $C(i, j) = 0$ $C(i, j) += A(i, k) * B(k, j)$ $C.tile(i, j, k, ii, jj, kk, II, JJ, KK)$	Algorithm
$C.isolate\_producer(A, A\_feeder)$ $C.unroll(ii, jj)$	Spatial Mapping

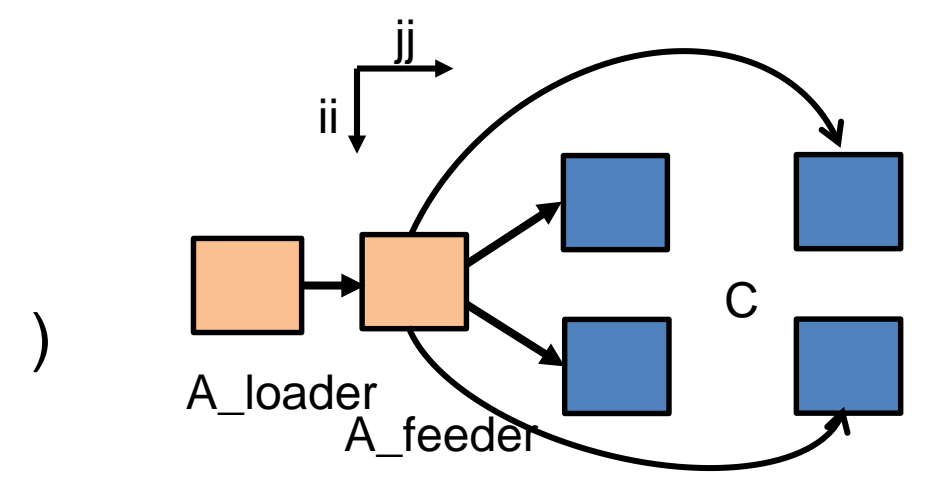
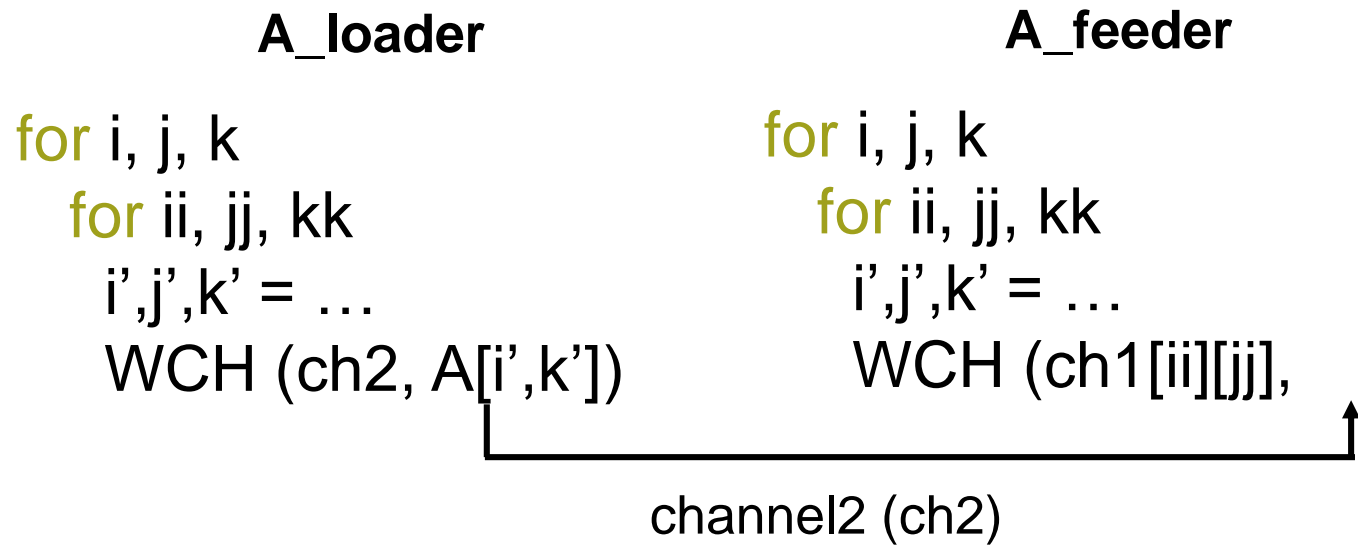


\* B[k',j']



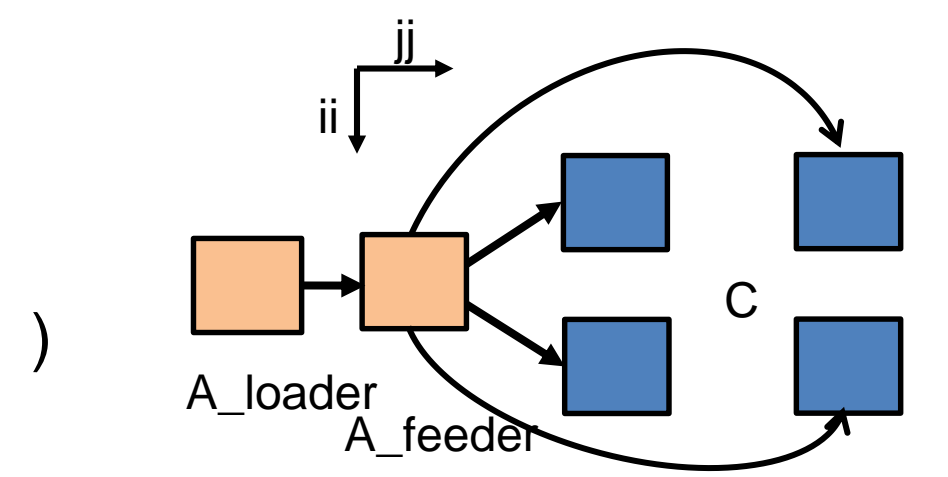
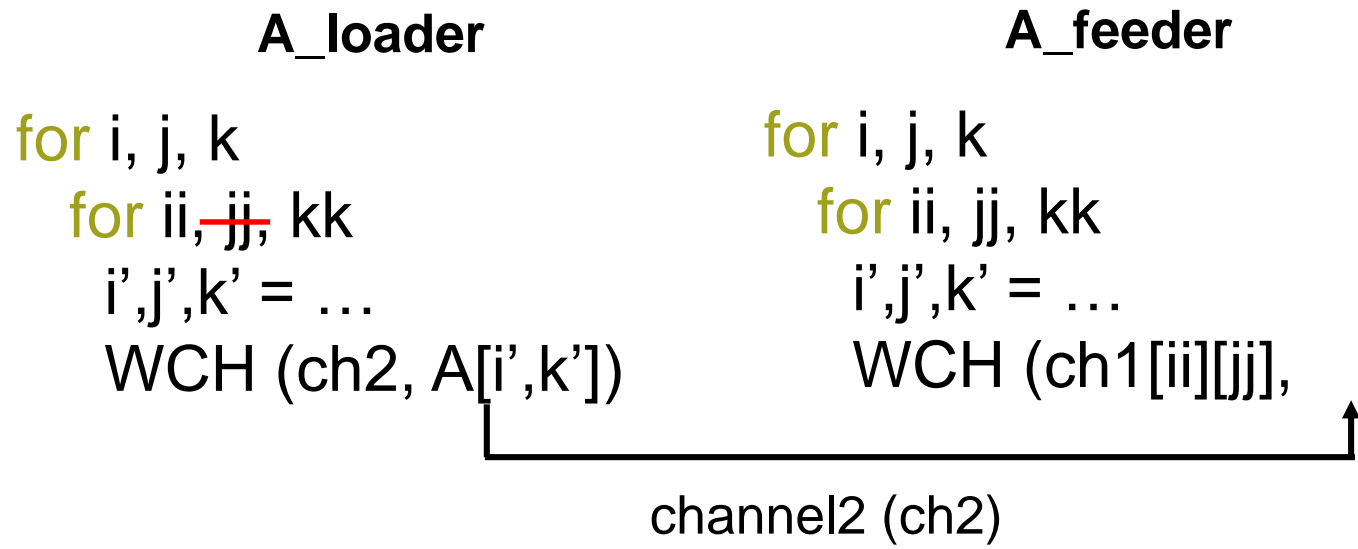
# Spatial Mapping in T2S

<pre> Func C C(i, j) = 0 C(i, j) += A(i, k) * B(k, j) C.tile(i,j,k,ii,jj,kk,II,JJ,KK) </pre>	Algorithm
<pre> C.isolate_producer(A, A_feeder) C.unroll(ii, jj) </pre>	Spatial Mapping
<pre> A_feeder.isolate_producer(A, A_loader) </pre>	



# Spatial Mapping in T2S

<pre> Func C C(i, j) = 0 C(i, j) += A(i, k) * B(k, j) C.tile(i,j,k,ii,jj,kk,II,JJ,KK) </pre>	Algorithm
<pre> C.isolate_producer(A, A_feeder) C.unroll(ii, jj) </pre>	Spatial Mapping
<pre> A_feeder.isolate_producer(A, A_loader) A_loader.remove(jj) A_feeder.buffer(ii, DOUBLE) </pre>	



# Spatial Mapping in T2S

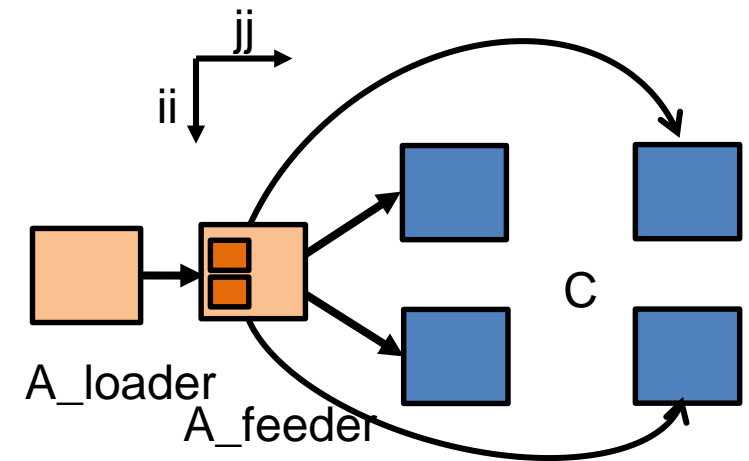
<pre> Func C C(i, j) = 0 C(i, j) += A(i, k) * B(k, j) C.tile(i,j,k,ii,jj,kk,II,JJ,KK)         </pre>	Algorithm
<pre> C.isolate_producer(A, A_feeder) C.unroll(ii, jj)         </pre>	Spatial Mapping
<pre> A_feeder.isolate_producer(A, A_loader) A_loader.remove(jj) A_feeder.buffer(ii, DOUBLE)         </pre>	

```

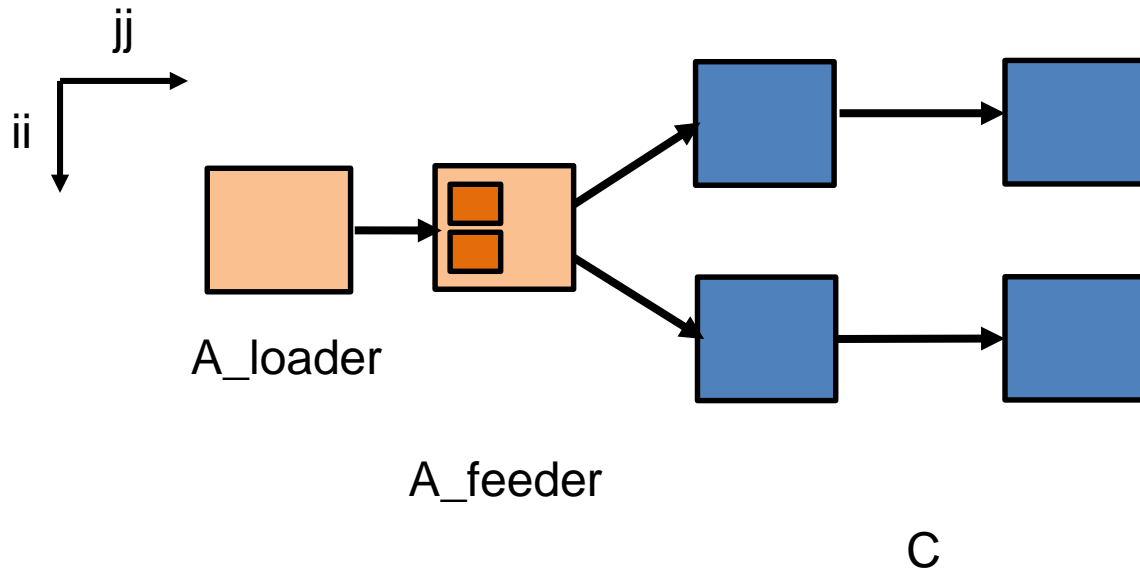
A_loader
for i, j, k
  for ii, jj, kk
    i',j',k' = ...
    WCH (ch2, A[i',k'])

A_feeder
for i, j, k, ii
  float buf [KK]
  for kk = 0 .. KK
    buf [kk] = RCH(.,.)

for jj, kk
  i',j',k' = ...
  WCH (ch1[ii][jj], buf[kk])
    
```



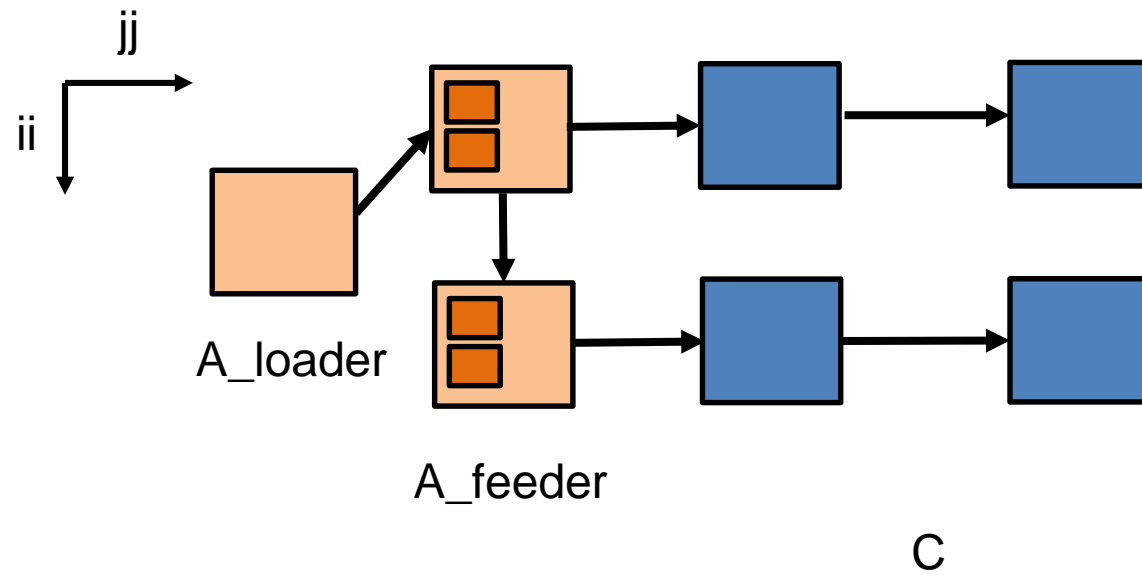
# Spatial Mapping in T2S



<pre> Func C C(i, j) = 0 C(i, j) += A(i, k) * B(k, j) C.tile(i,j,k,ii,jj,kk,ll,JJ,KK)         </pre>	Algorithm
<pre> C.isolate_producer(A, A_feeder) C.unroll(ii, jj)         </pre>	Spatial Mapping
<pre> A_feeder.isolate_producer(A, A_loader) A_loader.remove(jj) A_feeder.buffer(ii, DOUBLE)         </pre>	
<pre> C.forward(A_feeder, +jj)         </pre>	



# Spatial Mapping in T2S



```

Func C
C(i, j) = 0
C(i, j) += A(i, k) * B(k, j)
C.tile(i,j,k,ii,jj,kk,ll,JJ,KK)

```

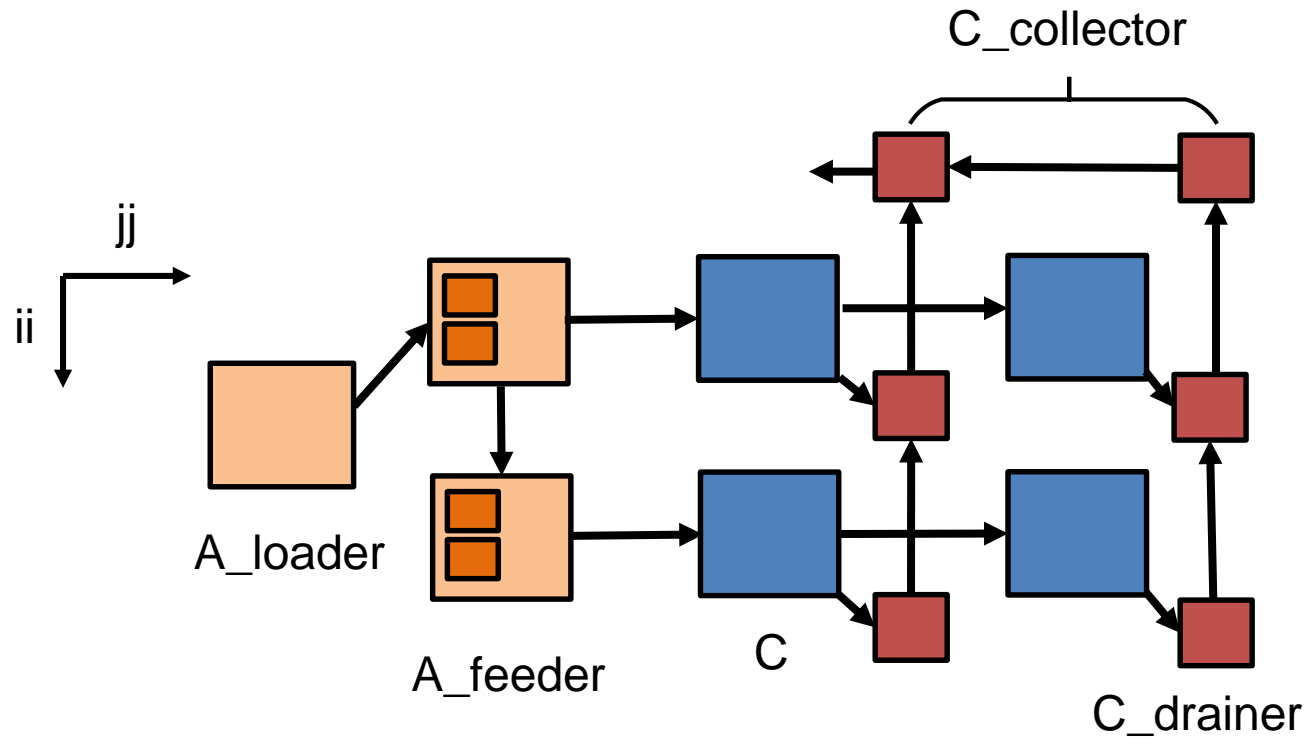
```

C.isolate_producer(A, A_feeder)
C.unroll(ii, jj)
A_feeder.isolate_producer(A, A_loader)
A_loader.remove(jj)
A_feeder.buffer(ii, DOUBLE)
C.forward(A_feeder, +jj)
A_feeder.unroll(ii).scatter(A, +ii)

```

Algorithm  
Spatial Mapping

# Spatial Mapping in T2S



Func C

$C(i, j) = 0$

$C(i, j) += A(i, k) * B(k, j)$

$C.tile(i, j, k, ii, jj, kk, ll, JJ, KK)$

$C.isolate\_producer(A, A\_feeder)$

$C.unroll(ii, jj)$

$A\_feeder.isolate\_producer(A, A\_loader)$

$A\_loader.remove(jj)$

$A\_feeder.buffer(ii, DOUBLE)$

$C.forward(A\_feeder, +jj)$

$A\_feeder.unroll(ii).scatter(A, +ii)$

$C.isolate\_consumer(C, C\_drainer)$

$C\_drainer.isolate\_consumer\_chain(C, C\_collector, C\_unloader)$

$C\_drainer.unroll(ii).unroll(jj).gather(C, -ii)$

$C\_collector.unroll(jj).gather(C, -jj)$

Algorithm

Spatial Mapping

# T2S Code

~20 LOC vs 750 lines of HLS code

```
C(j, i) = 0.0f;
C(j, i) += A(k, i) * B(j, k);
C.tile(j, i, jj, ii, JJ, II).tile(jj, ii, jjj, iii, JJJ, III);
C.update(0).tile(k, j, i, kk, jj, ii, KK, JJ, II).tile(kk, jj, ii, kkk, jjj, iii, KKK, JJJ, III);
C.update(0).isolate_producer_chain(A, A_serializer, A_loader, A_feeder)
    .isolate_producer_chain(B, B_serializer, B_loader, B_feeder)
    .isolate_consumer_chain(C, C_drainer, C_collector, C_unloader, C_deserializer);
A_serializer.sread().swrite();
B_serializer.sread().swrite();
C.update(0).vread({A,B});
C.update(0).unroll(ii)
    .unroll(jj)
C.update(0).forward(A_feeder, { 0, 1 })
    .forward(B_feeder, { 1, 0 });
A_serializer.remove(jjj, jj, j);
A_loader.remove(jjj, jj);
A_feeder.buffer(ii, true).unroll(ii);
B_serializer.remove(iii, ii, i);
B_loader.remove(iii, ii);
B_feeder.buffer(k, true).unroll(jj);
A_feeder.scatter(A, {1});
B_feeder.scatter(B, {1});
C_drainer.unroll(ii).unroll(jj).gather(C, jjj, { 1, 0 });
C_collector.unroll(jj).gather(C_drainer, jjj, { 1 });
```

# T2S Code

~20 LOC vs 750 lines of HLS code

```
C(j, i) = 0.0f;
```

```
C(j, i) += A(k, i) * B(j, k);
```

```
C.tile(j, i, jj, ii, JJ, II).tile(jj, ii, jjj, iii, JJJ, III);
```

```
C.update(0).tile(k, j, i, kk, jj, ii, KK, JJ, II).tile(kk, jj, ii, kkk, jjj, iii, KKK, JJJ, III);
```

```
C.update(0).isolate_producer_chain(A, A_serializer, A_loader, A_feeder)
```

```
    .isolate_producer_chain(B, B_serializer, B_loader, B_feeder)
```

```
    .isolate_consumer_chain(C, C_drainer, C_collector, C_unloader, C_deserializer);
```

```
A_serializer.sread().swrite();
```

```
B_serializer.sread().swrite();
```

```
C.update(0).vread({A,B});
```

```
C.update(0).unroll(ii)
```

```
    .unroll(jj)
```

```
C.update(0).forward(A_feeder, { 0, 1 })
```

```
    .forward(B_feeder, { 1, 0 });
```

```
A_serializer.remove(jjj, jj, j);
```

```
A_loader.remove(jjj, jj);
```

```
A_feeder.buffer(ii, true).unroll(ii);
```

```
B_serializer.remove(iii, ii, i);
```

```
B_loader.remove(iii, ii);
```

```
B_feeder.buffer(k, true).unroll(jj);
```

```
A_feeder.scatter(A, {1});
```

```
B_feeder.scatter(B, {1});
```

```
C_drainer.unroll(ii).unroll(jj).gather(C, jjj, { 1, 0 });
```

```
C_collector.unroll(jj).gather(C_drainer, jjj, { 1 });
```

**Custom compute  
(Loop Tiling)**

**Custom compute  
(Compute Partitioning)**

**Custom compute  
(Data Vectorization)**

**Custom compute  
(Loop Unrolling)**

# T2S Code

~20 LOC vs 750 lines of HLS code

```
C(j, i) = 0.0f;
```

```
C(j, i) += A(k, i) * B(j, k);
```

```
C.tile(j, i, jj, ii, JJ, II).tile(jjj, iii, JJJ, III);
```

**Custom compute  
(Loop Tiling)**

```
C.update(0).tile(k, j, i, kk, jj, ii, KK, JJ, II).tile(kkk, jjj, iii, KKK, JJJ, III);
```

```
C.update(0).isolate_producer_chain(A, A_serializer, A_loader, A_feeder)
```

```
    .isolate_producer_chain(B, B_serializer, B_loader, B_feeder)
```

```
    .isolate_consumer_chain(C, C_drainer, C_collector, C_unloader, C_deserializer);
```

**Custom compute  
(Compute Partitioning)**

```
A_serializer.sread().swrite();
```

```
B_serializer.sread().swrite();
```

```
C.update(0).vread({A,B});
```

**Custom compute  
(Data Vectorization)**

```
C.update(0).unroll(ii)
```

```
    .unroll(jj)
```

**Custom compute  
(Loop Unrolling)**

```
C.update(0).forward(A_feeder, { 0, 1 })
```

```
    .forward(B_feeder, { 1, 0 });
```

```
A_serializer.remove(jjj, jj, j);
```

```
A_loader.remove(jjj, jj);
```

```
A_feeder.buffer(ii, true).unroll(ii);
```

```
B_serializer.remove(iii, ii, i);
```

```
B_loader.remove(iii, ii);
```

```
B_feeder.buffer(k, true).unroll(jj);
```

**Custom memory  
(Buffer Insertion)**

```
A_feeder.scatter(A, {1});
```

```
B_feeder.scatter(B, {1});
```

```
C_drainer.unroll(ii).unroll(jj).gather(C, jjj, { 1, 0 });
```

```
C_collector.unroll(jj).gather(C_drainer, jjj, { 1 });
```

# T2S Code

~20 LOC vs 750 lines of HLS code

```
C(j, i) = 0.0f;
```

```
C(j, i) += A(k, i) * B(j, k);
```

```
C.tile(j, i, jj, ii, JJ, II).tile(jjj, ii, jjj, iii, JJJ, III);
```

```
C.update(0).tile(k, j, i, kk, jj, ii, KK, JJ, II).tile(kk, jj, ii, kkk, jjj, iii, KKK, JJJ, III);
```

```
C.update(0).isolate_producer_chain(A, A_serializer, A_loader, A_feeder)
```

```
    .isolate_producer_chain(B, B_serializer, B_loader, B_feeder)
```

```
    .isolate_consumer_chain(C, C_drainer, C_collector, C_unloader, C_deserializer);
```

```
A_serializer.sread().swrite();
```

```
B_serializer.sread().swrite();
```

```
C.update(0).vread({A,B});
```

```
C.update(0).unroll(ii)
```

```
    .unroll(jj)
```

```
C.update(0).forward(A_feeder, { 0, 1 })
```

```
    .forward(B_feeder, { 1, 0 });
```

```
A_serializer.remove(jjj, jj, j);
```

```
A_loader.remove(jjj, jj);
```

```
A_feeder.buffer(ii, true).unroll(ii);
```

```
B_serializer.remove(iii, ii, i);
```

```
B_loader.remove(iii, ii);
```

```
B_feeder.buffer(k, true).unroll(jj);
```

```
A_feeder.scatter(A, {1});
```

```
B_feeder.scatter(B, {1});
```

```
C_drainer.unroll(ii).unroll(jj).gather(C, jjj, { 1, 0 });
```

```
C_collector.unroll(jj).gather(C_drainer, jjj, { 1 });
```

**Custom compute  
(Loop Tiling)**

**Custom compute  
(Compute Partitioning)**

**Custom compute  
(Data Vectorization)**

**Custom compute  
(Loop Unrolling)**

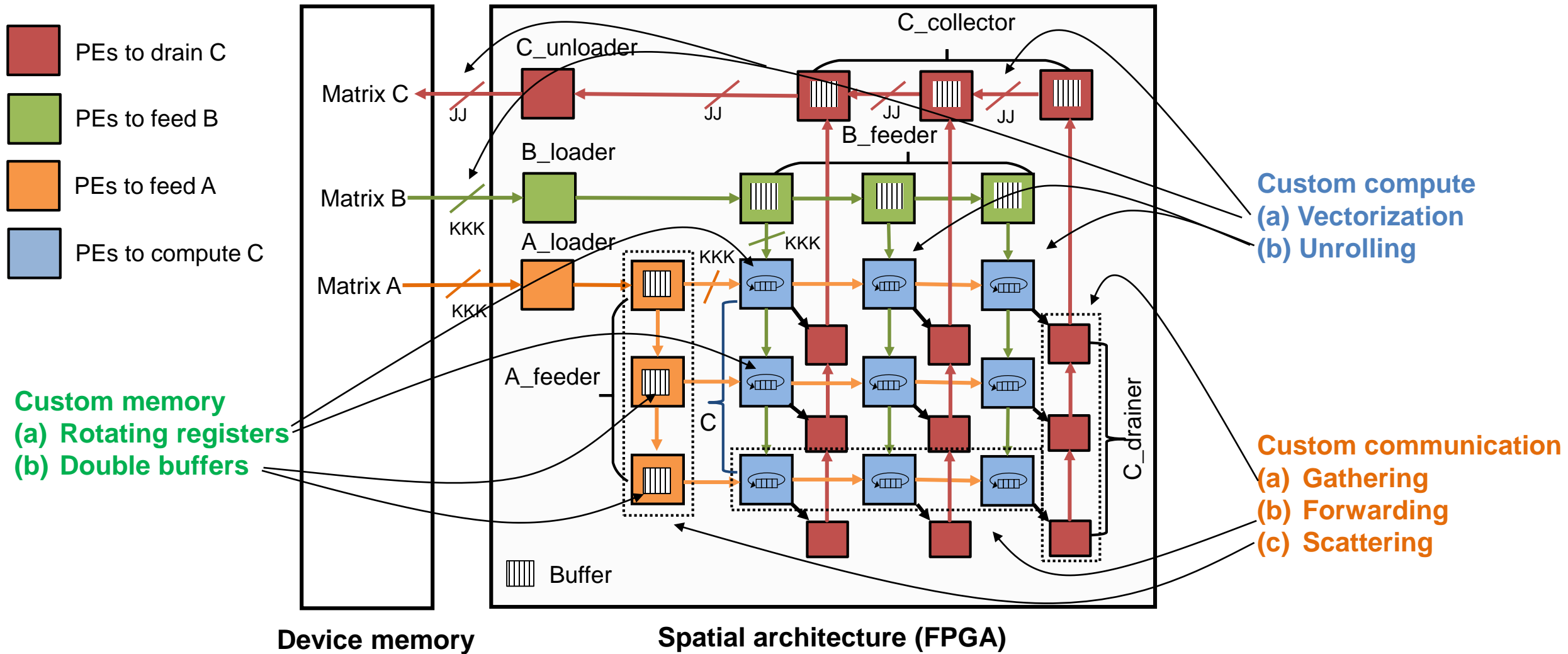
**Custom comm.  
(Data Forwarding)**

**Custom memory  
(Buffer Insertion)**

**Custom comm.  
(Data Scattering)**

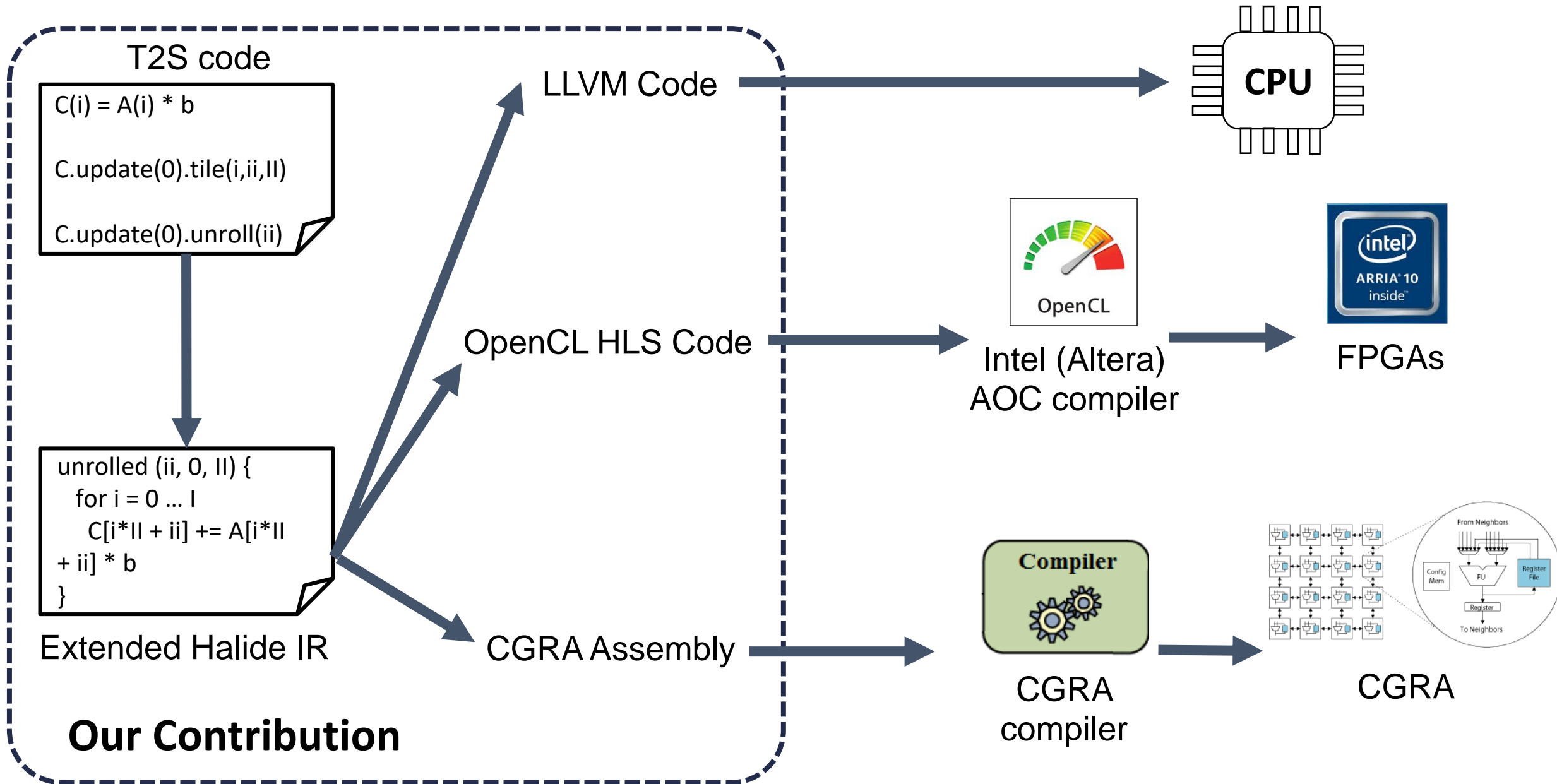
**Custom comm.  
(Data Gathering)**

# Driving Example – Matrix Multiplication (GEMM)



High-performance GEMM design on FPGA

# T2S Compilation Flow





# GEMM on Arria 10

- ▶ Baseline
  - Open-source NDRange-style OpenCL code, tuned on the specific FPGA
- ▶ Ninja
  - handwritten and manually optimized design from industry

	Baseline	T2S	Ninja
LOC	70	<b>20</b>	750
Systolic array size	--	10 x 8	10 x 8
Vector length	16 x float	16 x float	16 x float
# Logic elements	131 K (31%)	214 K (50%)	230 K (54%)
# DSPs	1,032 (68%)	<b>1,282 (84%)</b>	1,280 (84%)
# RAMs	1,534 (57%)	1,384 (51%)	1,069 (39%)
Frequency (MHz)	189	215	245
Throughput (GFLOPS)	311	<b>549</b>	626

1.8x speedup over the baseline with 3.5x less code  
82% performance of ninja implementation with 3% code

# Tensor Decomposition Kernels on FPGA & CGRA

## ► Tensor decomposition kernels

- **MTTKRP:**  $D(i, j) += A(i, k, l) * B(k, j) * C(l, j)$
- **TTM:**  $C(i, j, k) += A(i, j, l) * B(l, k)$
- **TTMc:**  $D(i, j, k) += A(i, l, m) * B(l, j) * C(m, k)$

Evaluation on CGRA

	LOC	Throughput wrt. Ninja GEMM	FMA Usage
GEMM	40	92 %	<b>100 %</b>
MTTKRP	32	99 %	<b>100 %</b>
TTM	47	104 %	<b>100 %</b>
TTMc	38	103 %	<b>95 %</b>

Evaluation on Arria-10 FPGA

Benchmark	LOC	Systolic Array Size	Logic Usage	DSP Usage	RAM Usage	Frequency (MHz)	Throughput (GFLOPS)
MTTKRP	28	8 x 9	53 %	81 %	56 %	204	<b>700</b>
TTM	30	8 x 11	64 %	93 %	88 %	201	<b>562</b>
TTMc	37	8 x 10	54 %	90 %	62 %	205	<b>738</b>

~100 % FMA usage and ~100 % throughput compared to ninja GEMM for CGRA  
 ~80-90 % DSP utilization and 560-740 GFLOPS for FPGA

# Conclusions

- ▶ T2S (Temporal to Spatial)
  - Provides a concise yet expressive programming abstraction that decouples spatial mapping from temporal definition
- ▶ Identifies a set of key compiler optimizations
  - essential for creating high-performance spatial hardware for tensor kernels
- ▶ Demonstrated high-performance (close to ninja) for
  - GEMM, MTTKRP, TTM and TTMc for both FPGA and CGRA

# **T2S-Tensor: Productively Generating High-Performance Spatial Hardware for Tensor Computations**

**Nitish Srivastava**, Hongbo Rong, Prithayan Barua, Guanyu Feng, Huanqi Cao, Zhiru Zhang, David Albonesi, Vivek Sarkar, Wenguang Chen, Paul Petersen, Geoff Lowney, Adam Herr, Christopher Hughes, Timothy Mattson, Pradeep Dubey

# **Question?**