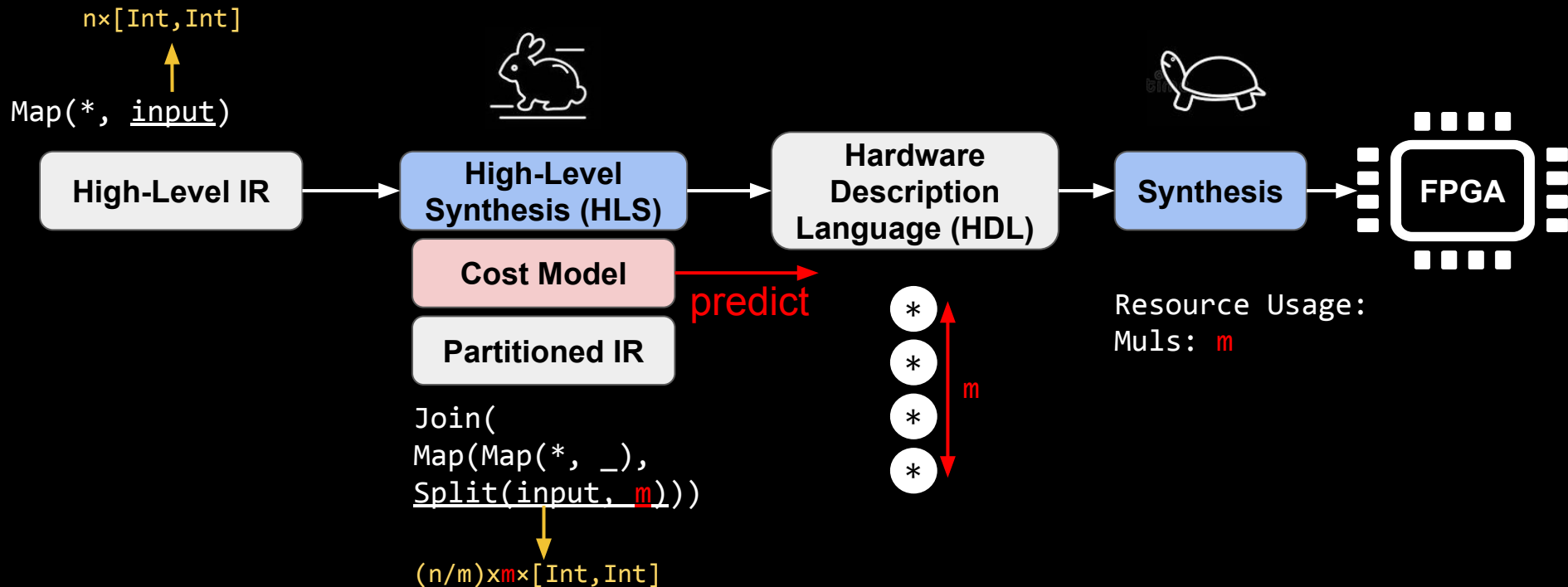


# Toward Automatic Hardware and Data Partitioning in HLS

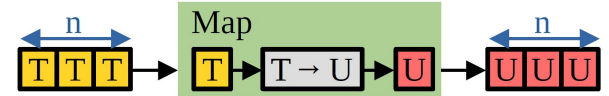


McGill

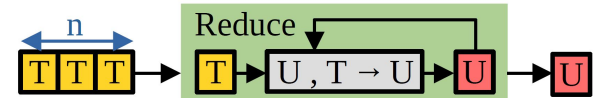
Tzung-Han Juang, Jonathan Van der Cruysse, and Christophe Dubach  
McGill University – CDP '24 – Tuesday November 12, 2024

# Functional Intermediate Representation (IR) [1]

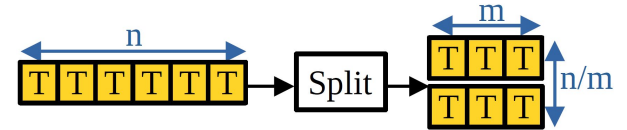
Map( $f: T \rightarrow U$ ,  $in: n \times T$ ):  $n \times U$



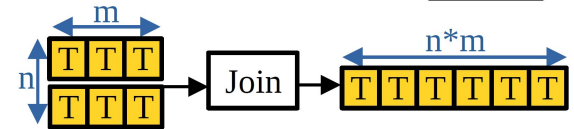
Reduce( $f: (U, T) \rightarrow U$ ,  $in: n \times T$ ):  $U$



Split( $in: n \times T$ ,  $size: m$ ):  $(n/m) \times m \times T$



Join( $in: n \times m \times T$ ):  $(nm) \times T$



DotProd( $in: n \times [Int, Int]$ ):  $Int = \text{Reduce}(+, \text{Map}(*, in))$

# Lower Functional IR

## Algo Level:

Map(\*3, in: n×Int): n×Int



## Arch Level:

MapStm(\*3, in: Stm[Int,n]): Stm[Int,n]

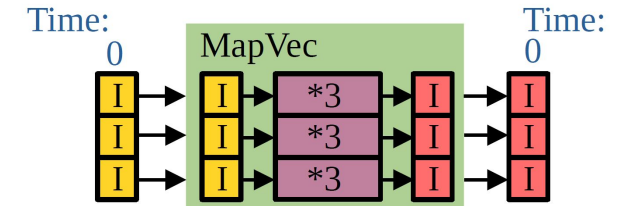
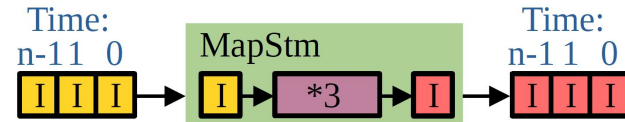
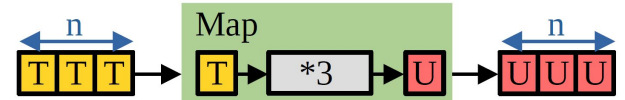


MapVec(\*3, in: Vec[Int,n]): Vec[Int,n]



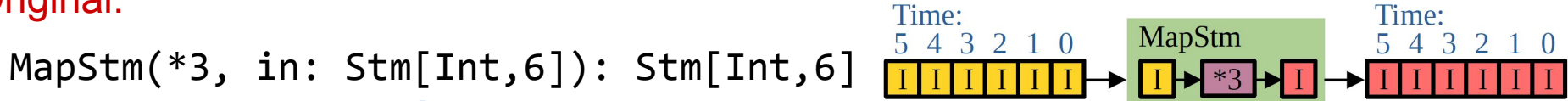
## HDL Level:

entity MapVec is port(input in std\_logic\_vector ...)

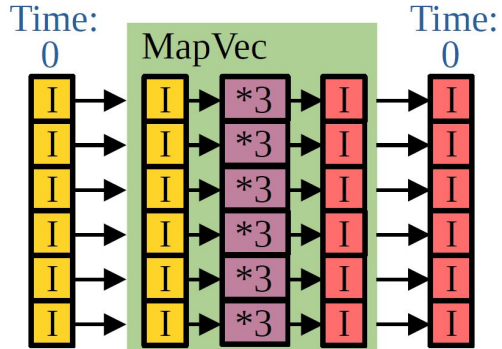


# Motivation: Determining Partition Size

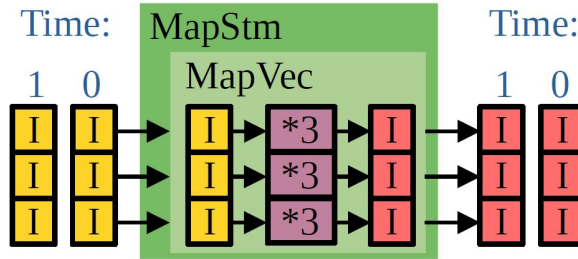
Original:



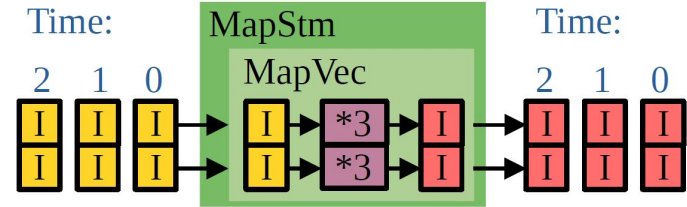
↓ After Optimizing



1 cycle  
6 resources



2 cycles  
3 resources



3 cycles  
2 resources

# Preserve Tuning Parameters for Partitioning

Algo Level:

```
Map(*3, in: 6×Int): 6×Int
```

↓ Express Partitioning

```
Join(Map(Map(*3, _), Split(in, tp)))
```

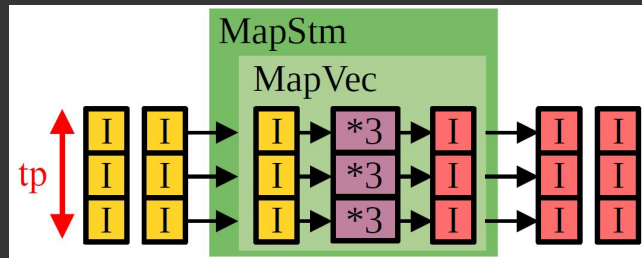
↓ Lowering

Arch Level:

```
JoinStm(MapStm(MapStm(*3, _), SplitStm(in, tp)))
```

↓ Optimizing

```
JoinStm(MapStm(Vec2Stm,  
  MapStm(MapVec(*3, _),  
    MapStm(Stm2Vec, SplitStm(in, tp))))))
```



# Insert Tuning Parameters

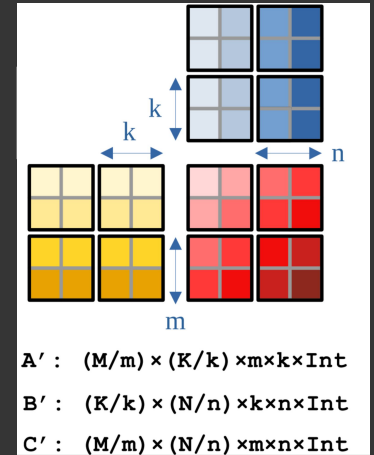
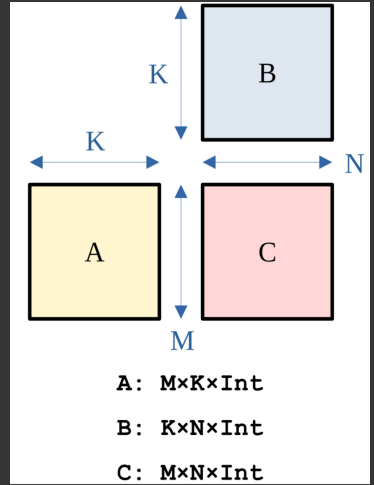
```
MatMul(A: MxKxInt, B: KxNxInt): MxNxInt =
  Map(a -> Map(b -> DotProd(a, b), Transpose(B)), A)
```



Add Tuning Parameters:  $m$ ,  $n$ ,  $k$

```
A' = Transpose(Split(A, (m, k))): (M/m)x(K/k)xmxkxInt
B' = Transpose(Split(B, (k, n))): (K/k)x(N/n)xkxnxInt
C' = Map(a' =>
  Map(b' =>
    Reduce(+,
      Map(tup =>
        MatMul(tup.0: mxkxInt, tup.1: kxnxInt)
        Zip(a', b')
      )), Transpose(B')
  ), A'): (M/m)x(N/n)xmxnxInt
```

Linked tuning parameter



# Partitioning with Single Input

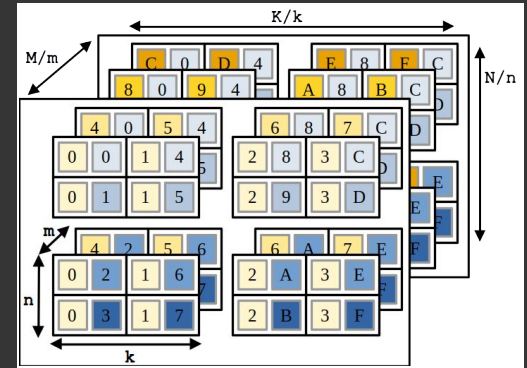
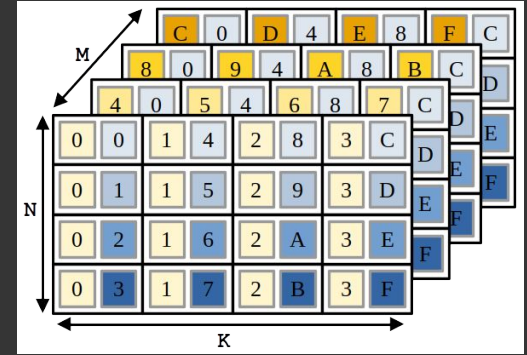
MatMul(AB:  $M \times N \times K \times [\text{Int}, \text{Int}]$ ):  $M \times N \times \text{Int}$

DCM(divideFunc, conquerFunc, mergeFunc, AB)

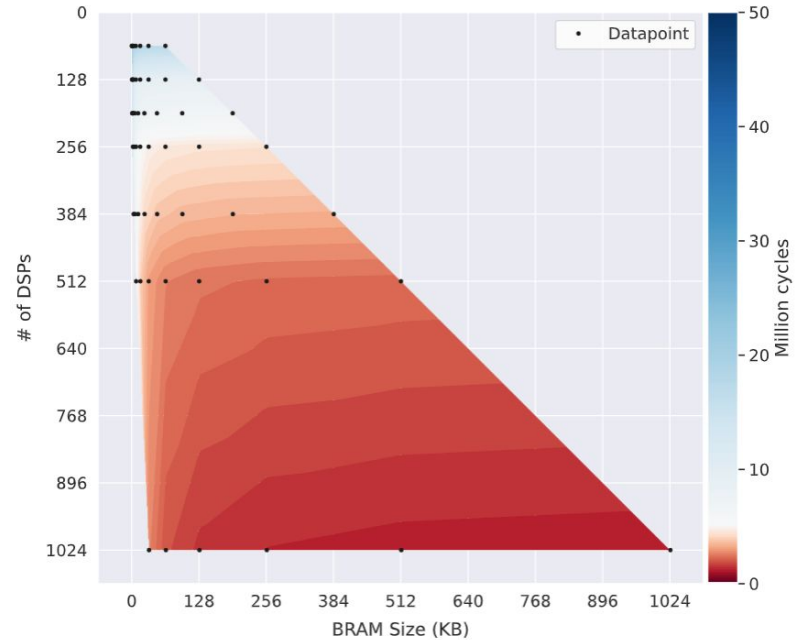
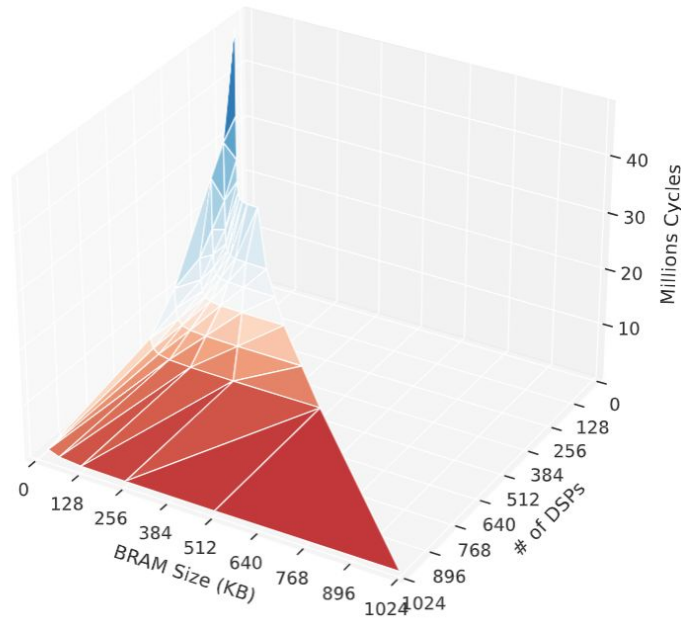
divided:  $(M/m) \times (N/n) \times (K/k) \times m \times n \times k \times [\text{Int}, \text{Int}]$   
= Transpose(Split(AB, (m, n, k)))

conquered:  $(M/m) \times (N/n) \times (K/k) \times m \times n \times \text{Int}$   
= Map(Map(Map(MatMul(\_), \_), \_), divided)

merged:  $(M/m) \times (N/n) \times m \times n \times \text{Int}$   
= Map(Map(Reduce(+, \_), \_), conquered)

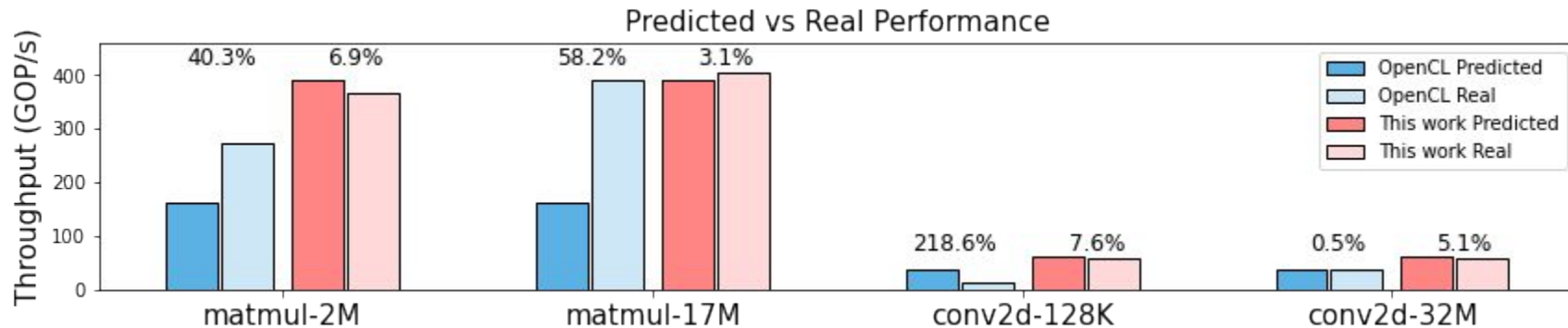


# Exploration for Matrix Multiplication

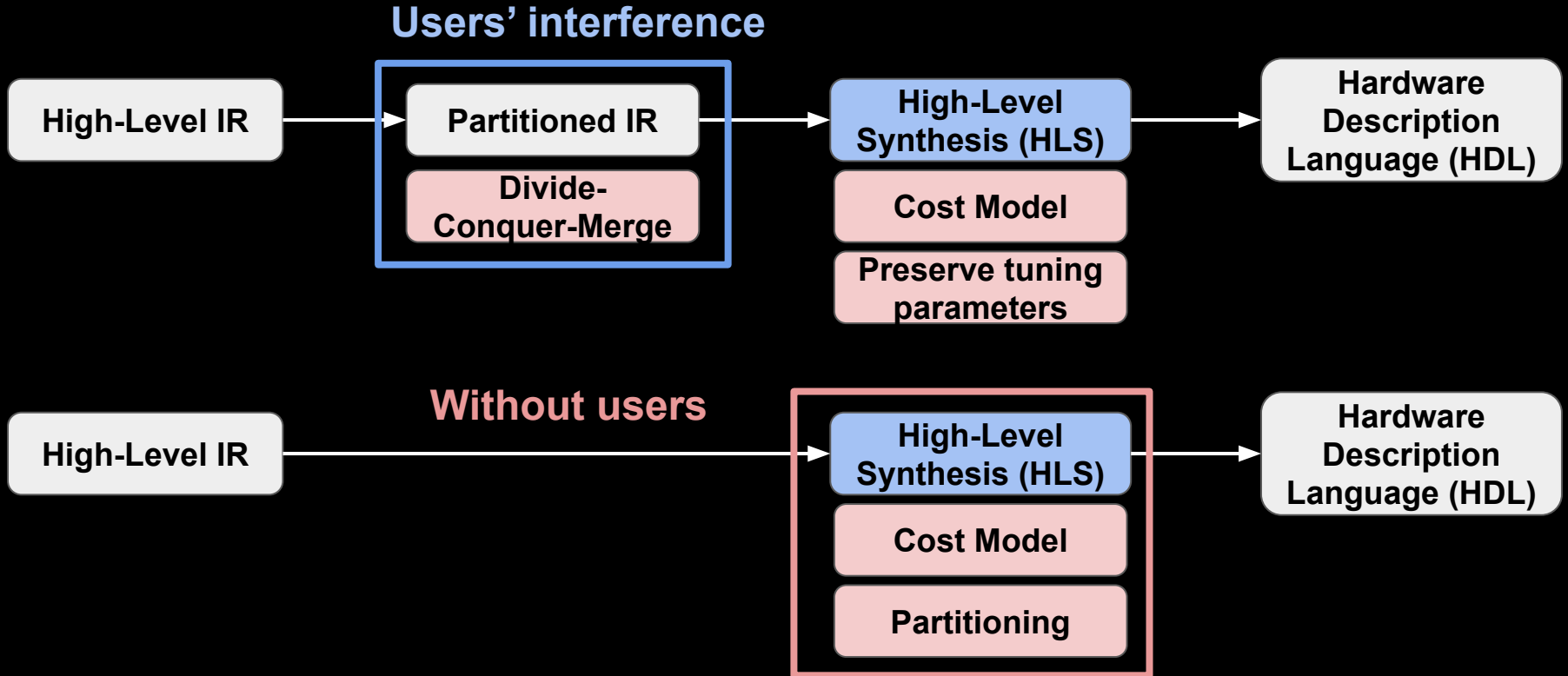




# Performance Prediction against Intel OpenCL

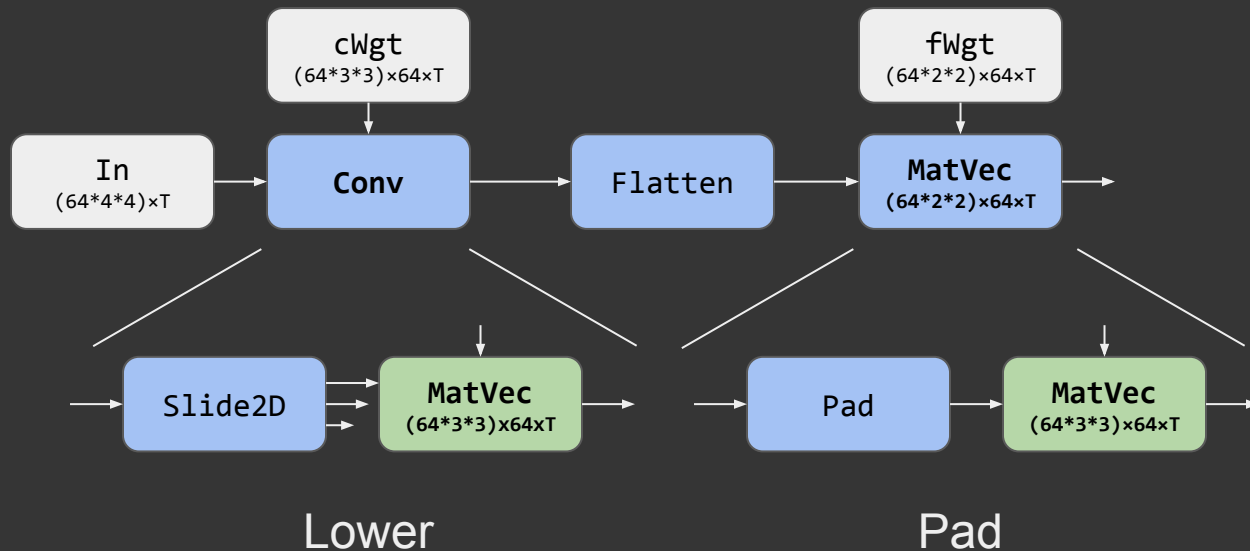


# From Semi-Automatic to Automatic Partitioning



# Automatic Partitioning with Resource Sharing

`MatVec(fWgt, Flatten(Conv(cWgt, In)))`

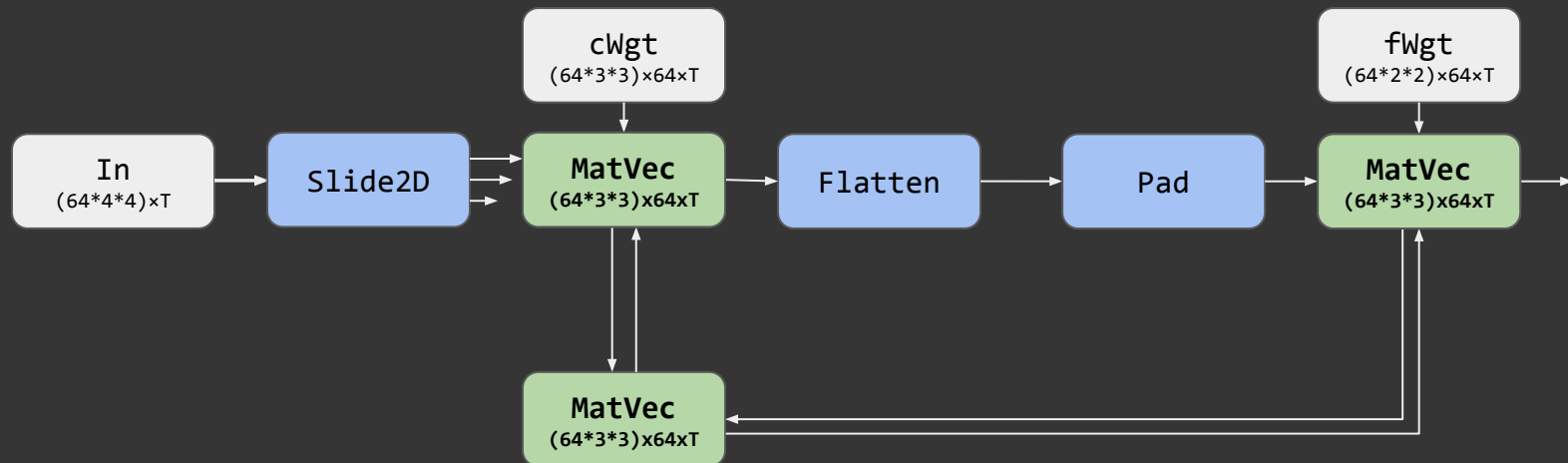


# Resource Sharing

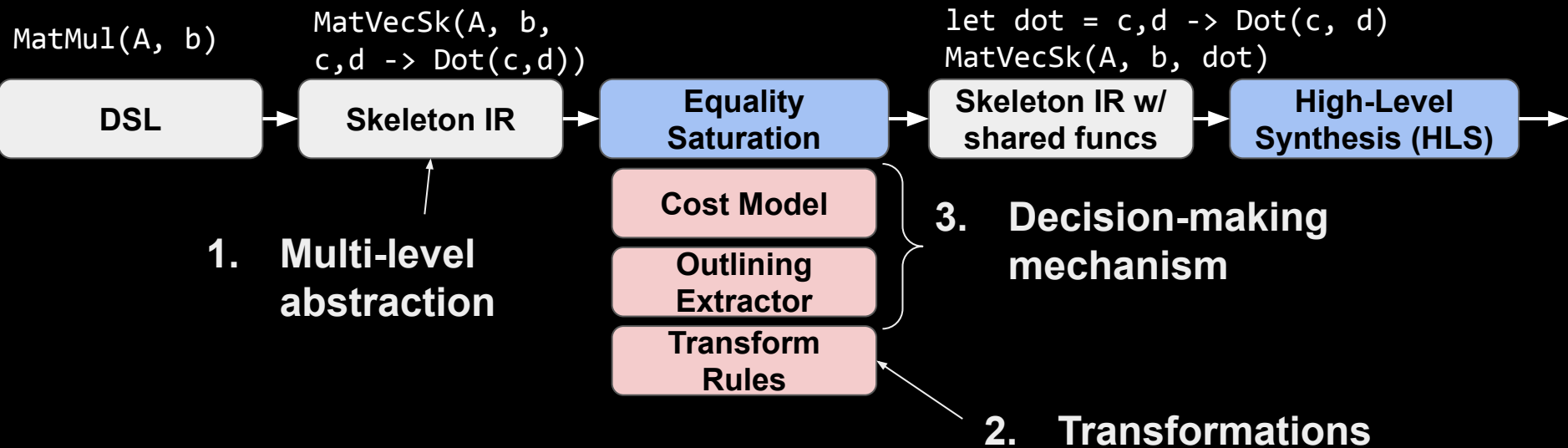
$\text{MatVec}(\text{fWgt}, \text{Flatten}(\text{Conv}(\text{cWgt}, \text{In})))$

## Ingredients:

1. Multi-level abstraction
2. Transformations
3. Decision-making mechanism



# Overview



# Ingredient #1: Skeleton IR

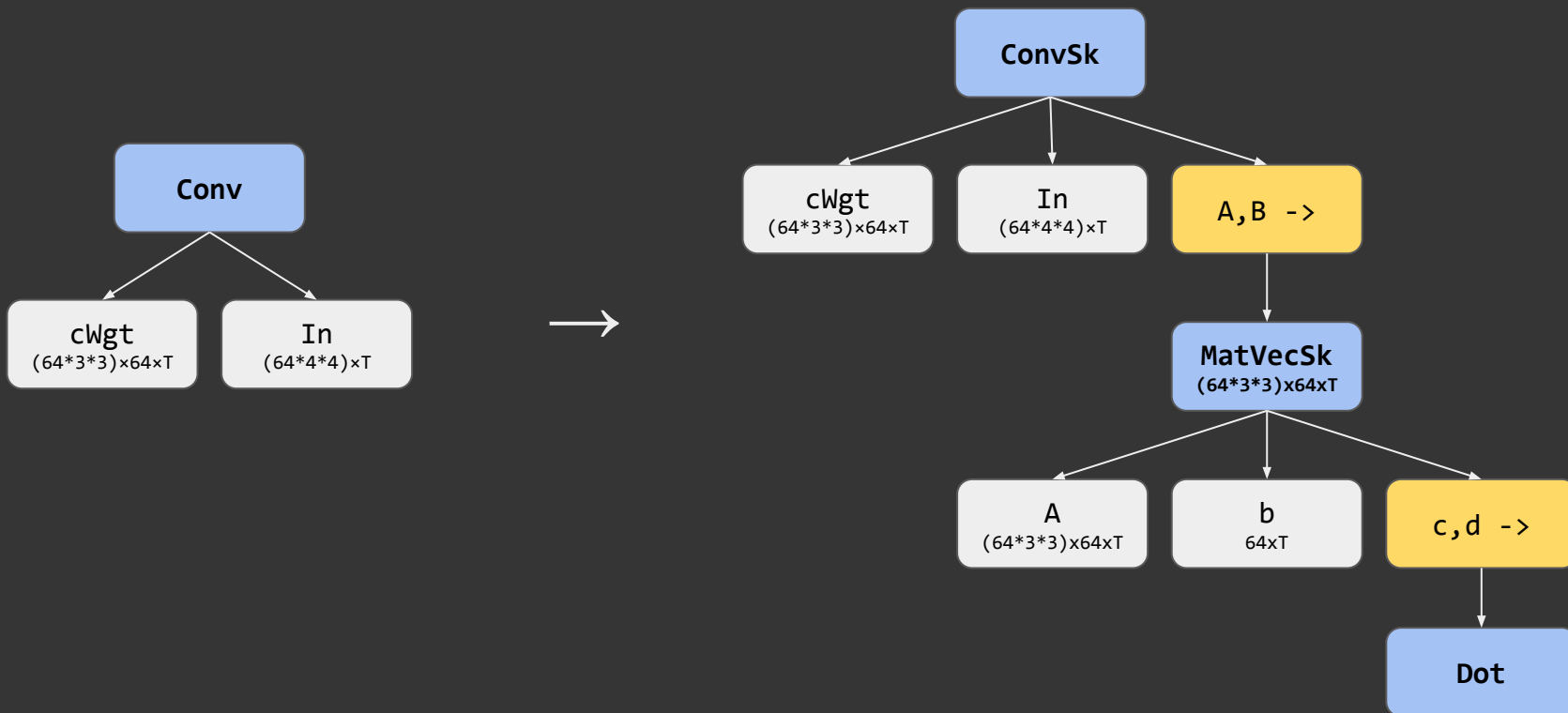
**DSL**

**Skeleton IR**

**MatVec(A, b) → MatVecSk(A, b, c, d -> Dot(c, d))**

**Conv(Wgt, In) → ConvSk(Wgt, In, A, b -> MatVecSk(A, b, ...))**

# Skeleton IR



# Ingredient #2: Transformations

## Conv

- Tile
    - Width/height
    - Output channels
    - Input channels
  - Pad width/height
- 

## MatVec

- Tile
  - Pad
- 

## Dot

- Change parallelism



# Transformations

**DSL**

**Skeleton IR**

**MatVec(A, b) → MatVecSk(A, b, c,d -> Dot(c, d))**

**Conv(Wgt, In) → ConvSk(Wgt, In, A,B -> MatVecSk(A, B, ...))**

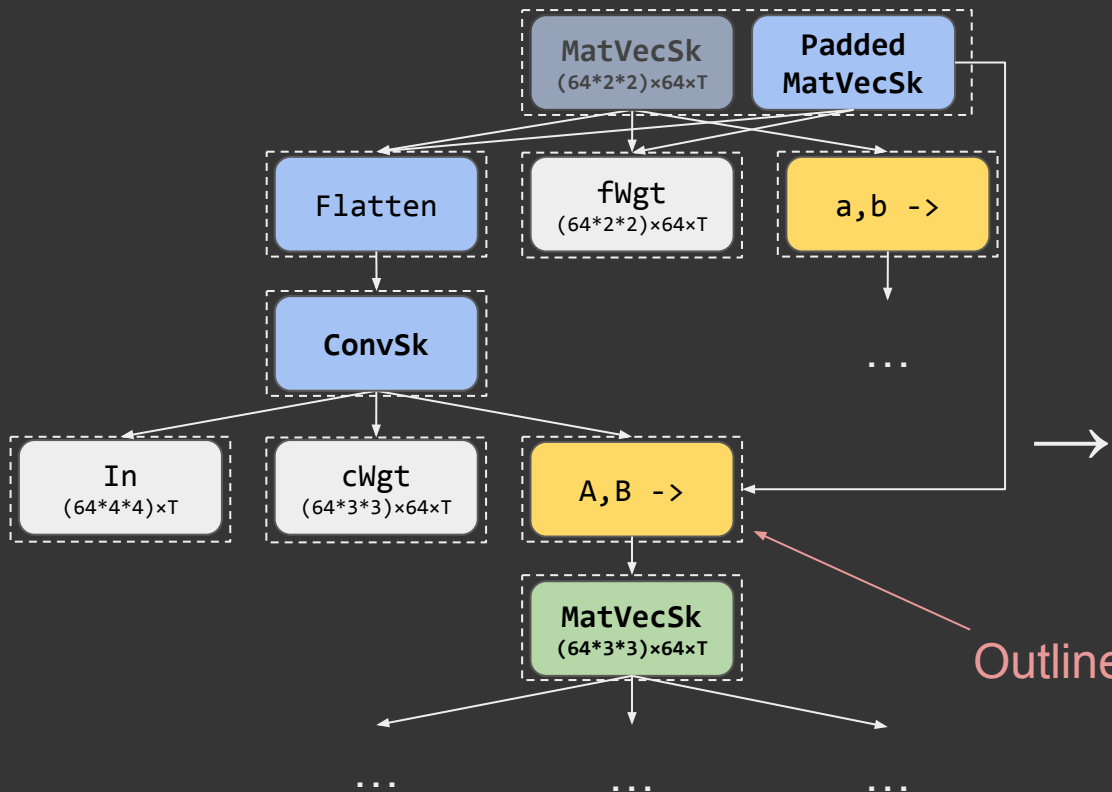
**TiledMatVecSk(A, B, C,D -> MatVecSk(C, D, ...))**

**PaddedMatVecSk(A, B, C,D -> MatVecSk(C, D, ...))**

# Ingredient #3: Equality Saturation

## Execution construction

- MatVec padding rule



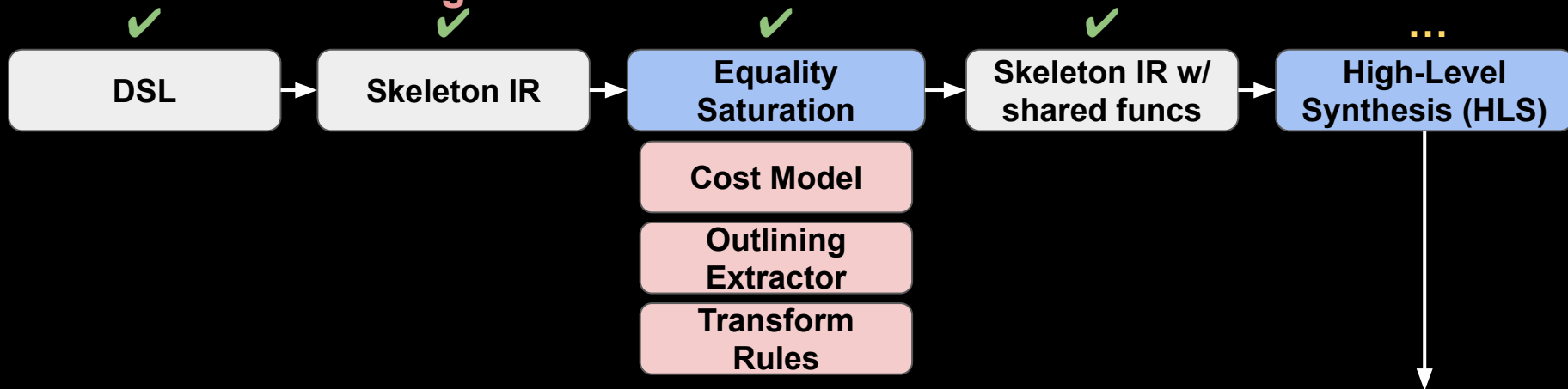
```
let mv = A, B ->
  MatVecSk(64*3*3)x64xT(...)
```

```
PaddedMatVecSk(
  fWgt,
  Flatten(ConvSk(
    cWgt, In, mv)),
  mv)
```

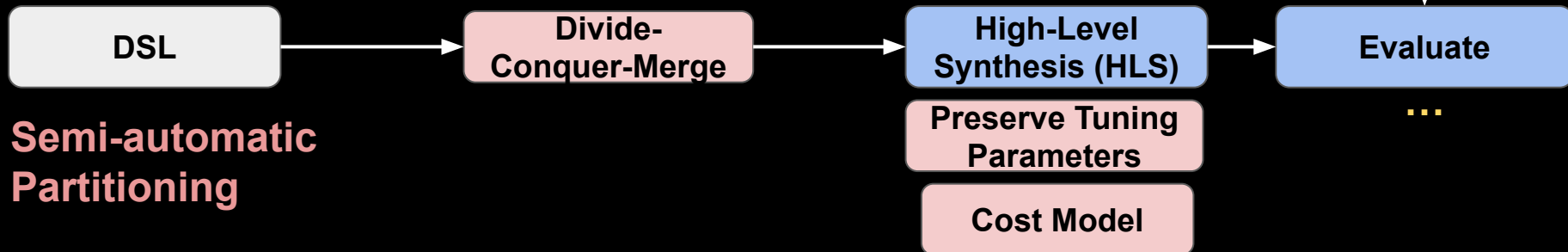
Outline

# Summary + Future Work

## Automatic Partitioning



## Semi-automatic Partitioning



# Toward Automatic Hardware and Data Partitioning in HLS

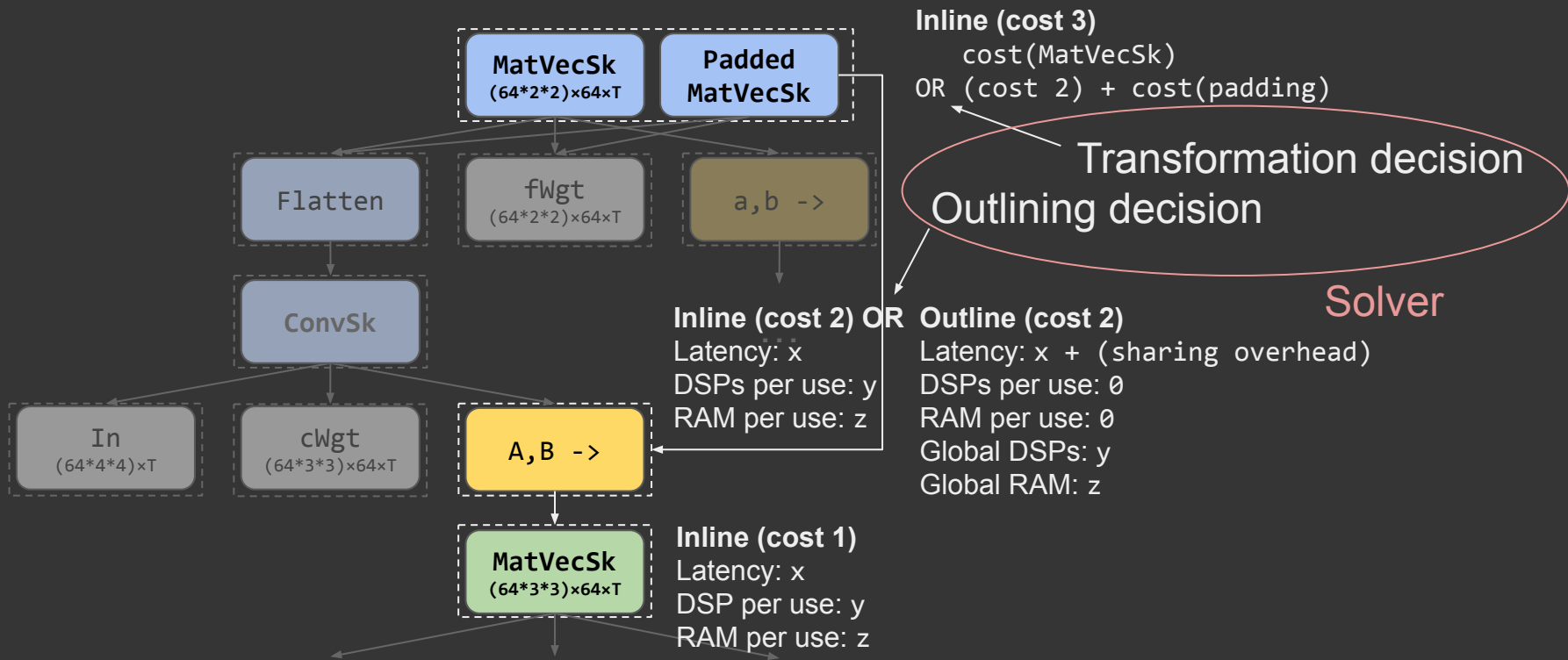


Tzung-Han Juang  
McGill University  
tzung-han.juang@mail.mcgill.ca



Jonathan Van der Cruysse  
McGill University  
jonathan.vandercruysse@mail.mcgill.ca

# Equality Saturation: Extraction and Outlining



# MatVec Skeleton

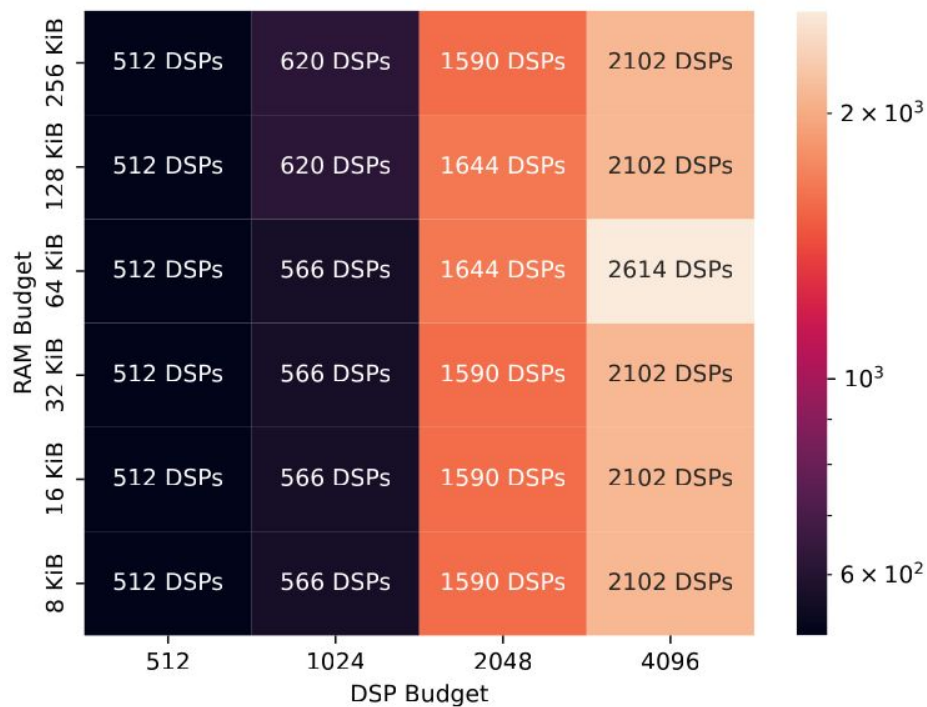
```
MatVec(A, b) = map(  
  (a, b') -> Dot(a, b'),  
  zip(A, repeat(B)))
```

```
MatVecSk(A, b, f) = map(  
  f,  
  zip(A, repeat(B)))
```

# Latency (in steps) per resource budget, first 4 layer of VGG



# DSP use per resource budget, first 4 layer of VGG





# RAM use per resource budget, first 4 layer of VGG

