



# Optimizing FPGA-based Accelerator Design for Deep Convolutional Neural Network

**Chen Zhang**<sup>1</sup>, Peng Li<sup>3</sup>, Guangyu Sun<sup>1,2</sup>, Yijin Guan<sup>1</sup>, Bingjun Xiao<sup>3</sup>, Jason Cong<sup>1,2,3</sup>

<sup>1</sup>Peking University

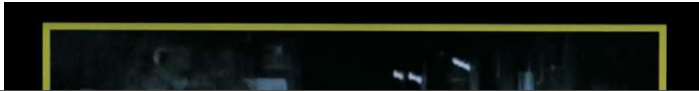
<sup>2</sup>PKU/UCLA Joint Research Institution

<sup>3</sup>University of California, Los Angeles



# Convolutional Neural Network

## Automotive Safety



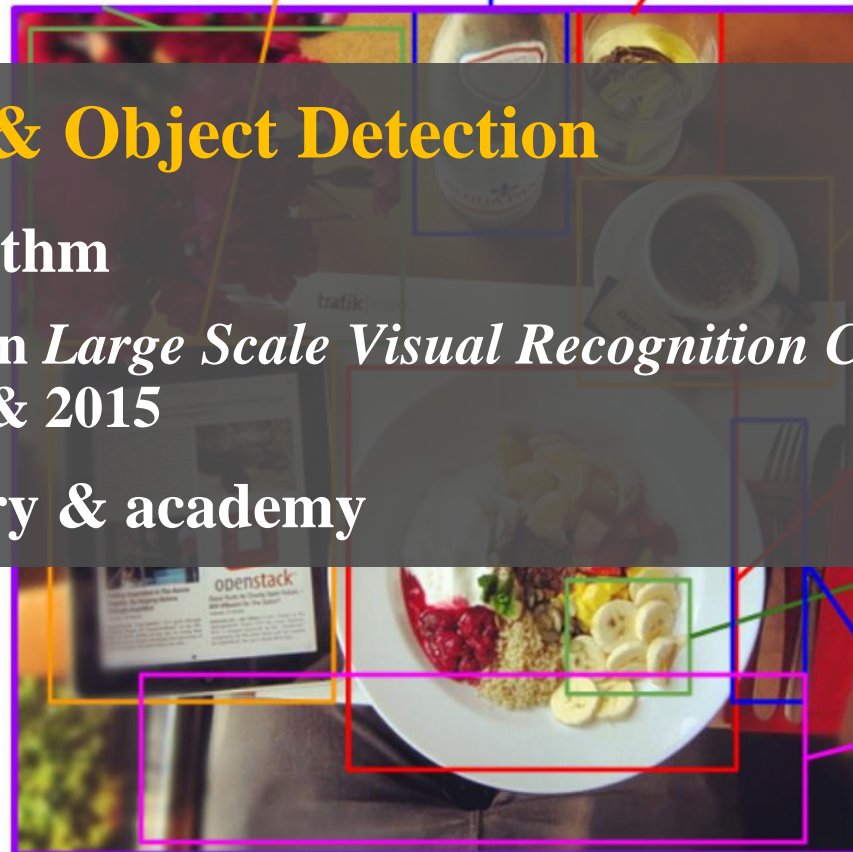
## Image Classification & Object Detection

- State-of-the-art Algorithm
  - Highest Correctness in *Large Scale Visual Recognition Challenge* 2012 & 2013 & 2014 & 2015
- Widely used in industry & academy



## Image Descriptions @ Stanford

bouquet of red flowers    tablet    bottle of water    glass of water with ice and lemon



cup of coffee

dining table with breakfast items

plate of fruit

banana

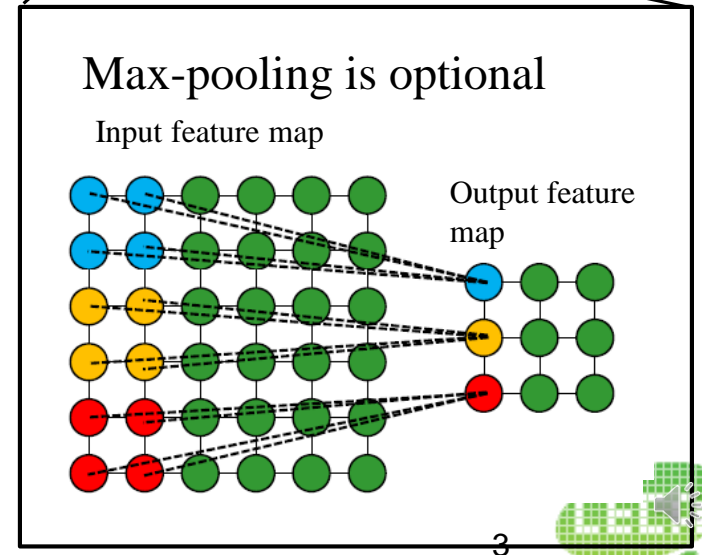
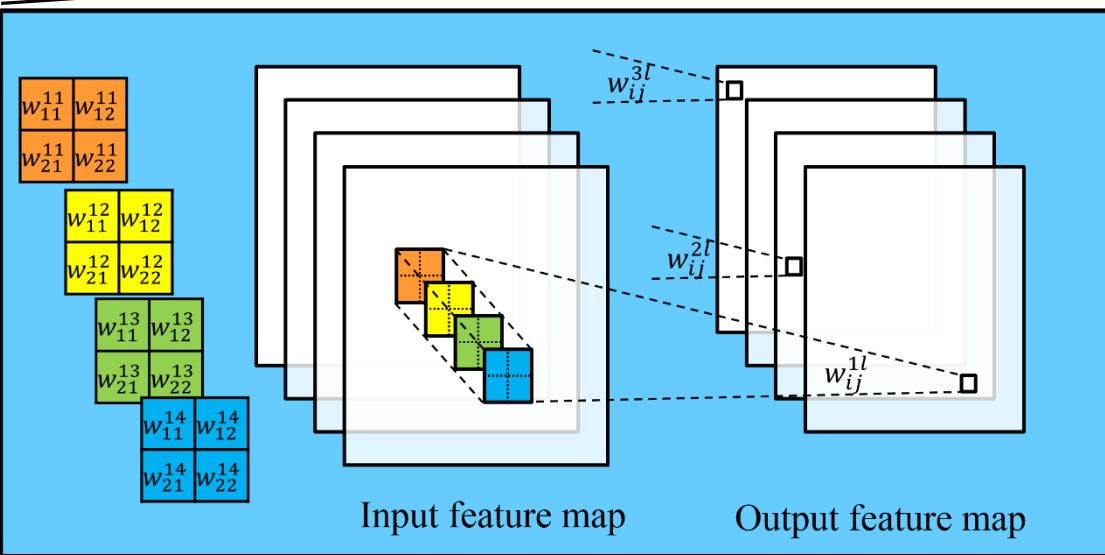
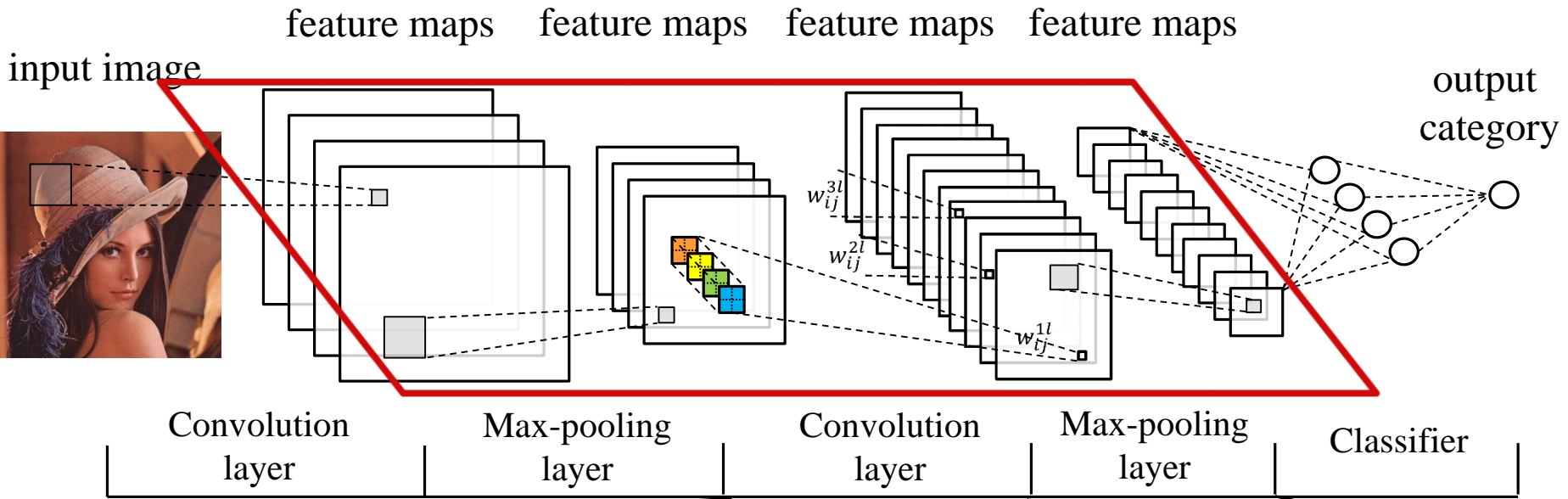
slices

fork

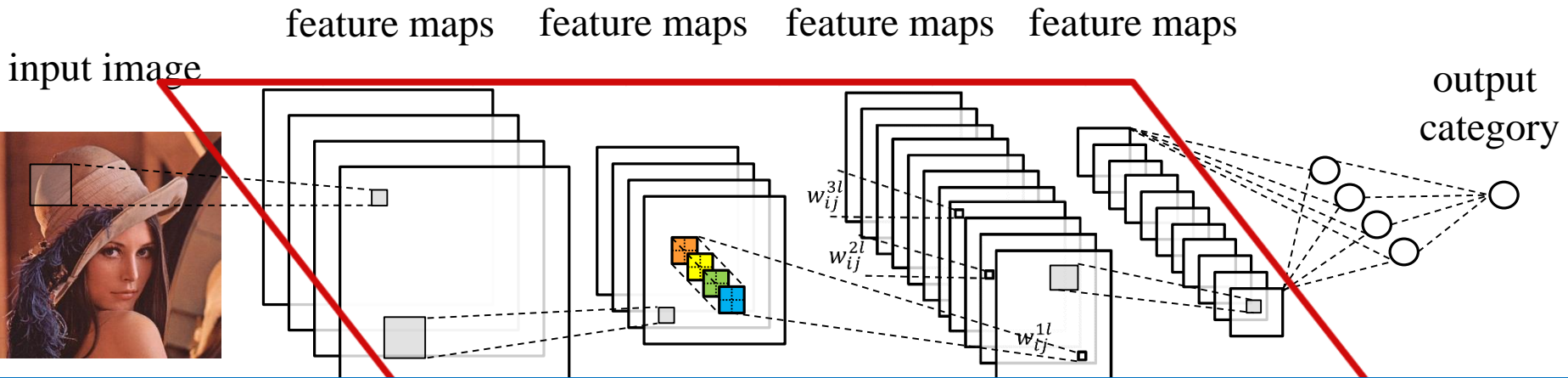
a person sitting at a table



# Convolutional Neural Network

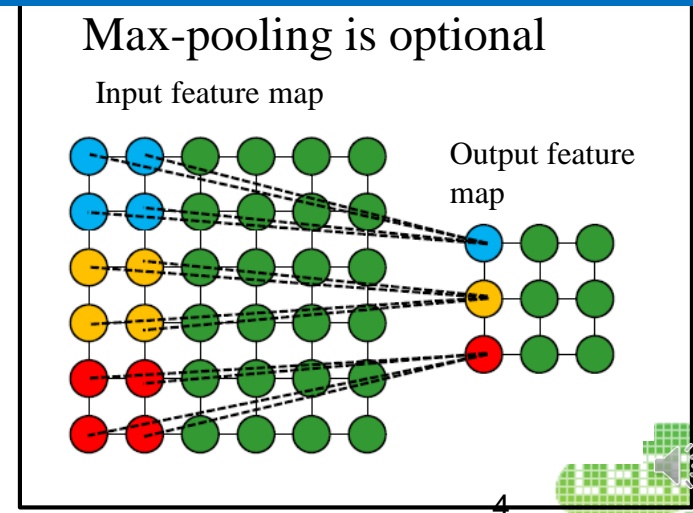
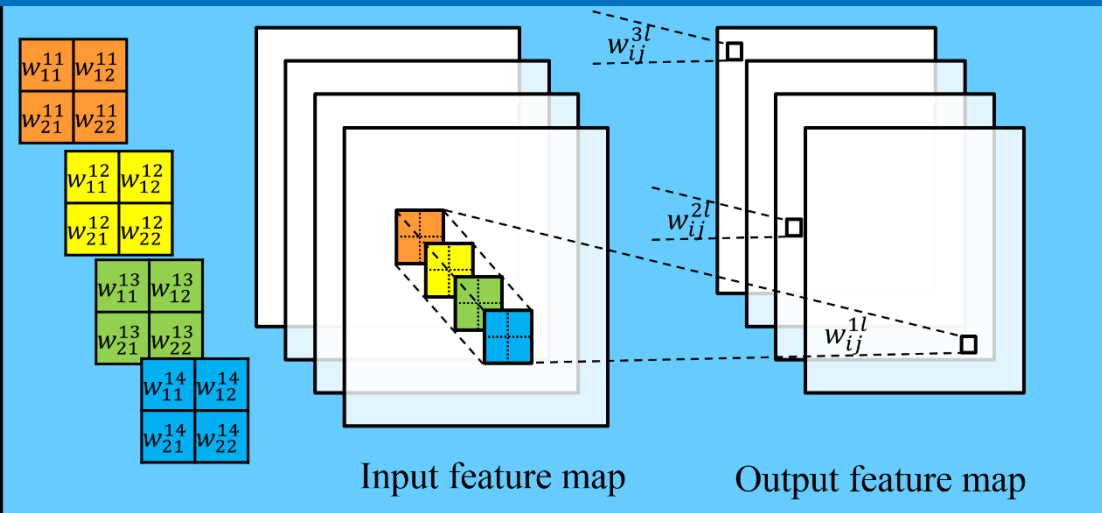


# Convolutional Neural Network



**Convolutional layers account for over 90% computation**

- [1] A. Krizhevsky, et al. Imagenet classification with deep convolutional neural networks. NIPS 2012.
- [2] J. Cong and B. Xiao. Minimizing computation in convolutional neural networks. ICANN 2014



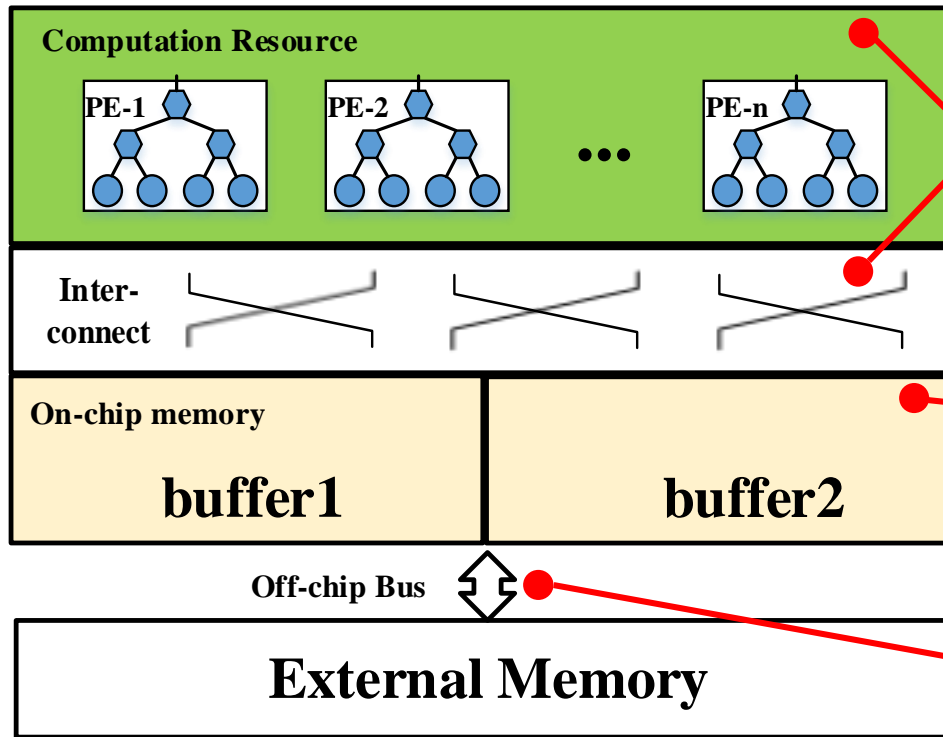
# Related Work

---

1. **CNP: An FPGA-based processor for convolutional networks. FPL 2009;**
2. **A massively parallel coprocessor for convolutional neural networks. ASAP 2009;**
3. **A programmable parallel accelerator for learning and classification. PACT 2010;**
4. **A Dynamically Configurable Coprocessor for Convolutional Neural Networks. ISCA 2010;**
5. **Design and Implementation of an FPGA-based Real-Time Face Recognition System. (short paper) FCCM 2011;**
6. **A massively parallel digital learning processor. NIPS 2008;**
7. **NeuFlow: A Runtime Reconfigurable Dataflow Processor for Vision. CVPRW 2011;**
8. **Accelerating deep neural networks on mobile processor with embedded programmable logic. (Poster) NIPS 2013;**
9. **Memory-Centric Accelerator Design for Convolutional Neural Networks. ICCD 2013.**



# Hardware Computing on FPGA



Challenge: Resrc. Util.

Solution: Unroll & Pipeline



Challenge: BRAM Limits

Solution: Loop Tiling



Challenge: BW Limits

Solution: Data reuse



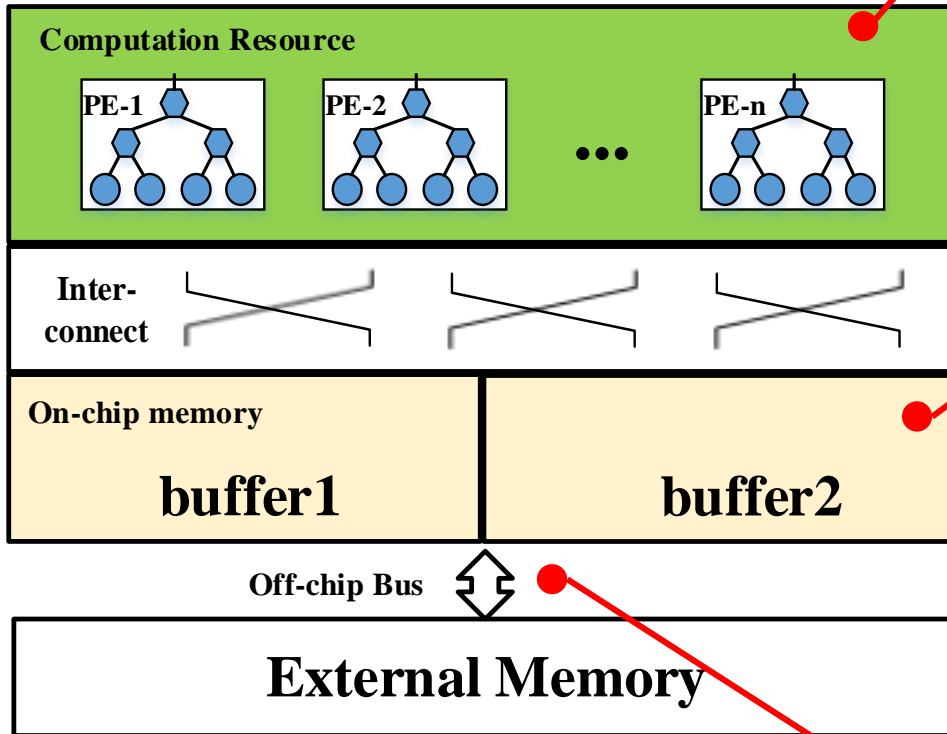
Challenge: Co-optimization,  
-- match 'computation' & 'communication'

Solution: Roofline Model

Main  
Contribution



# FPGA design in Roofline Model



Computational Perf.

$$= \frac{\text{total number of operations}}{\text{execution cycles}}$$

GFLOP/S

Computation To  
Communication Ratio

$$= \frac{\text{Total number of operations}}{\text{Amount of external data access}}$$

FLOP / (Mem. Byte Access)

Bandwidth

Gbyte/S



# FPGA design in Roofline Model

Computational Performance

GFLOP/S

Design

$$\frac{\text{Gbyte}}{S} = \frac{\text{GFLOP/s}}{\text{FLOP}/(\text{Mem. Byte Access})}$$

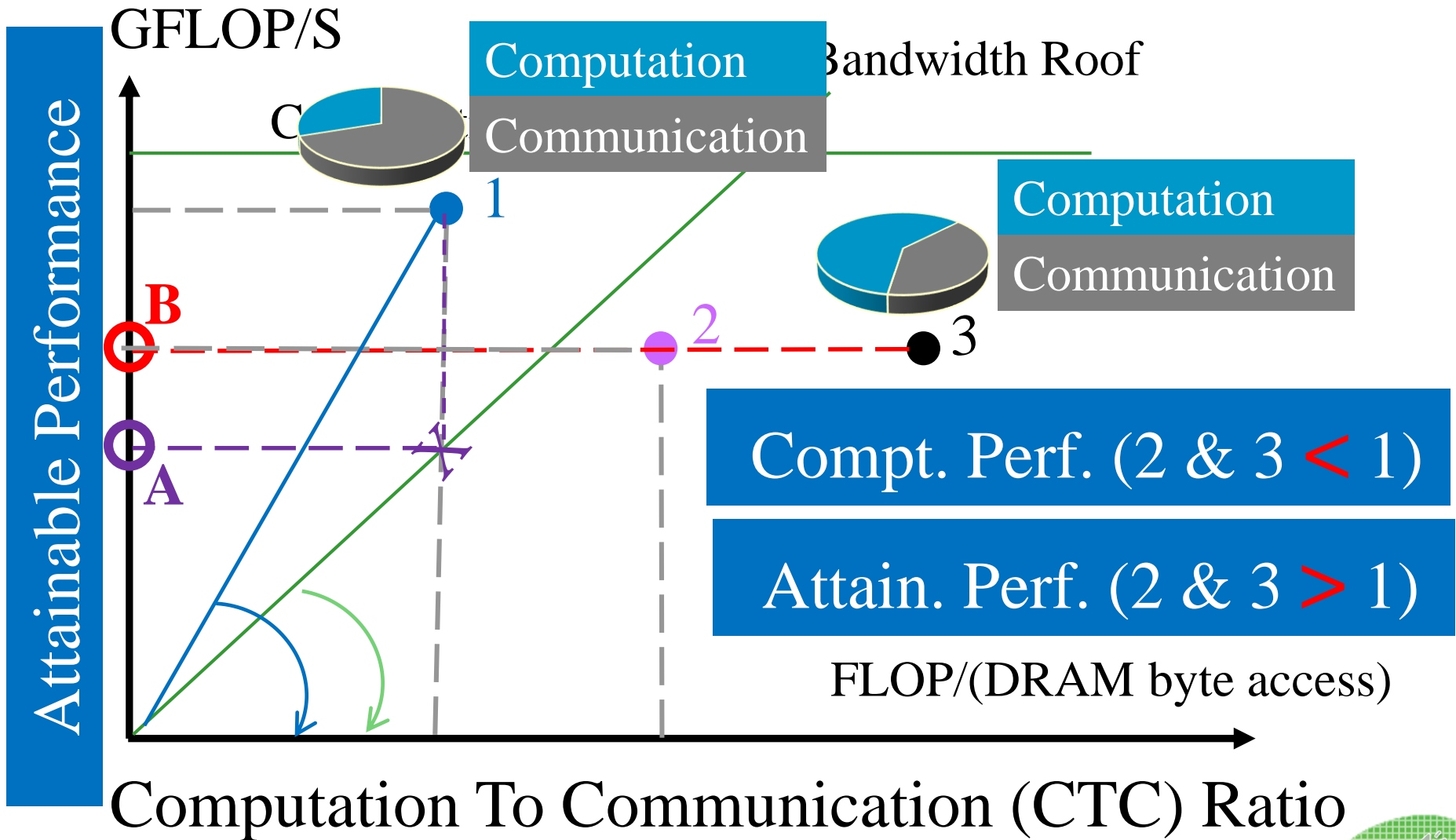
BandWidth

FLOP/(Memory byte access)

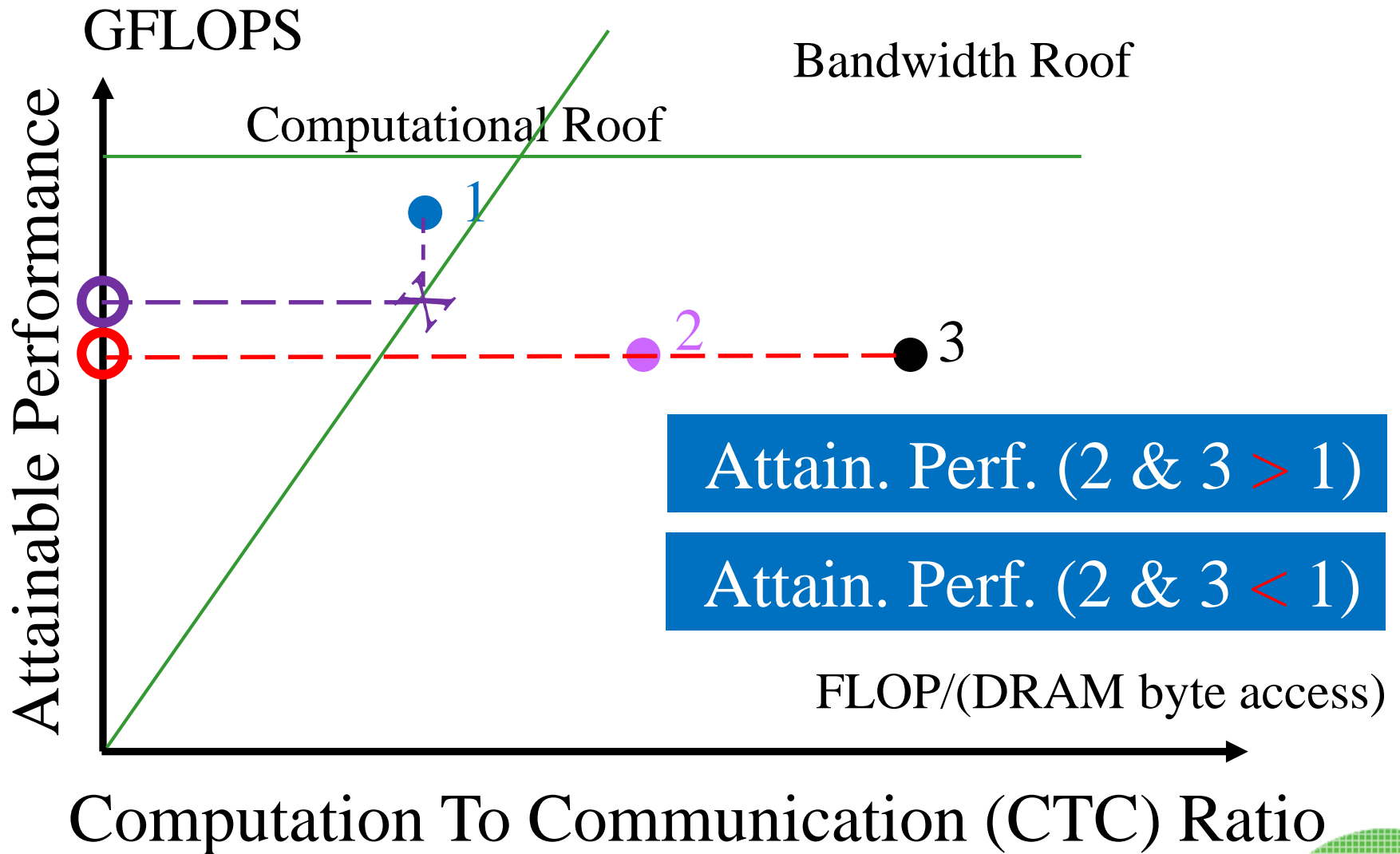
Computation To Communication (CTC) Ratio



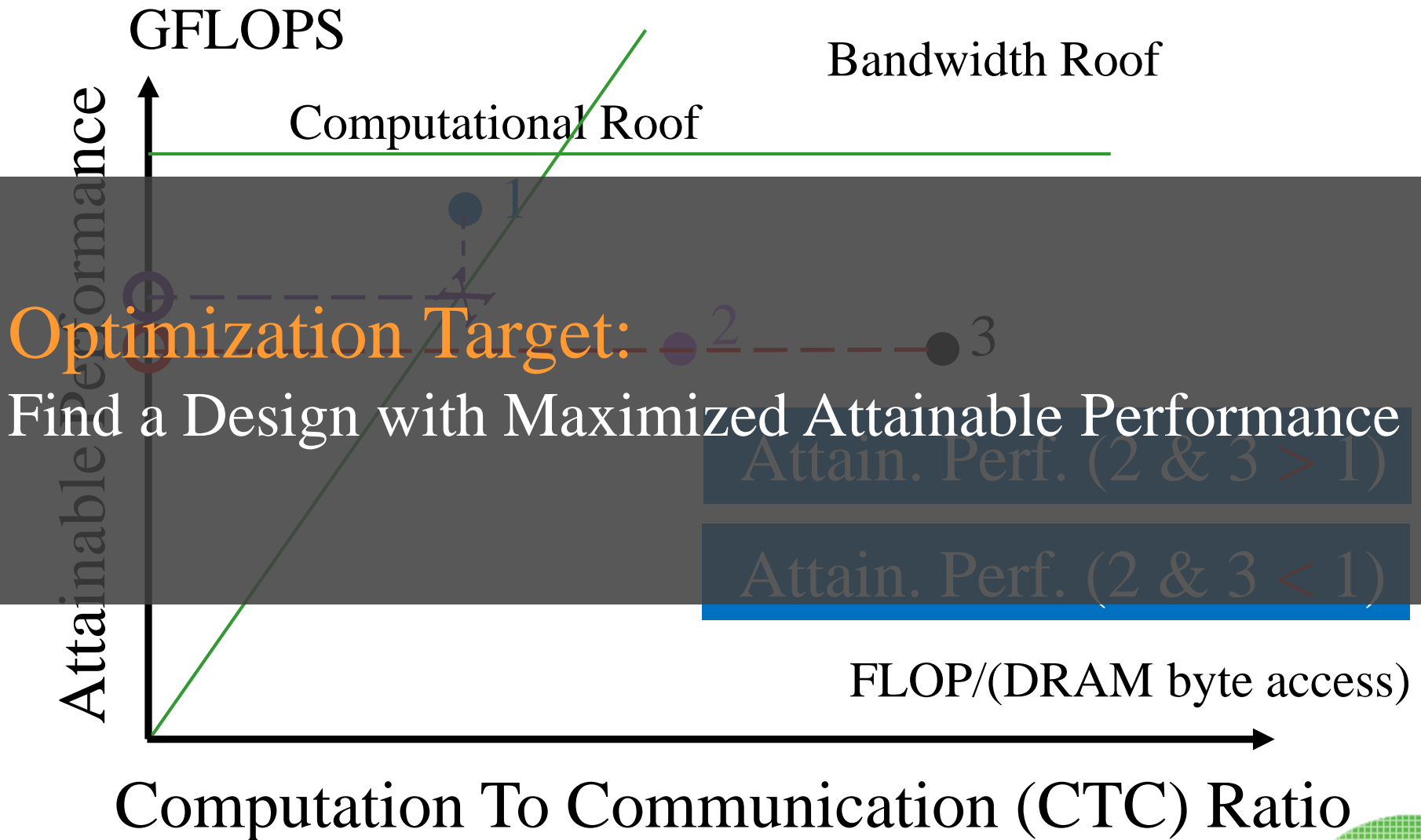
# Attainable Performance



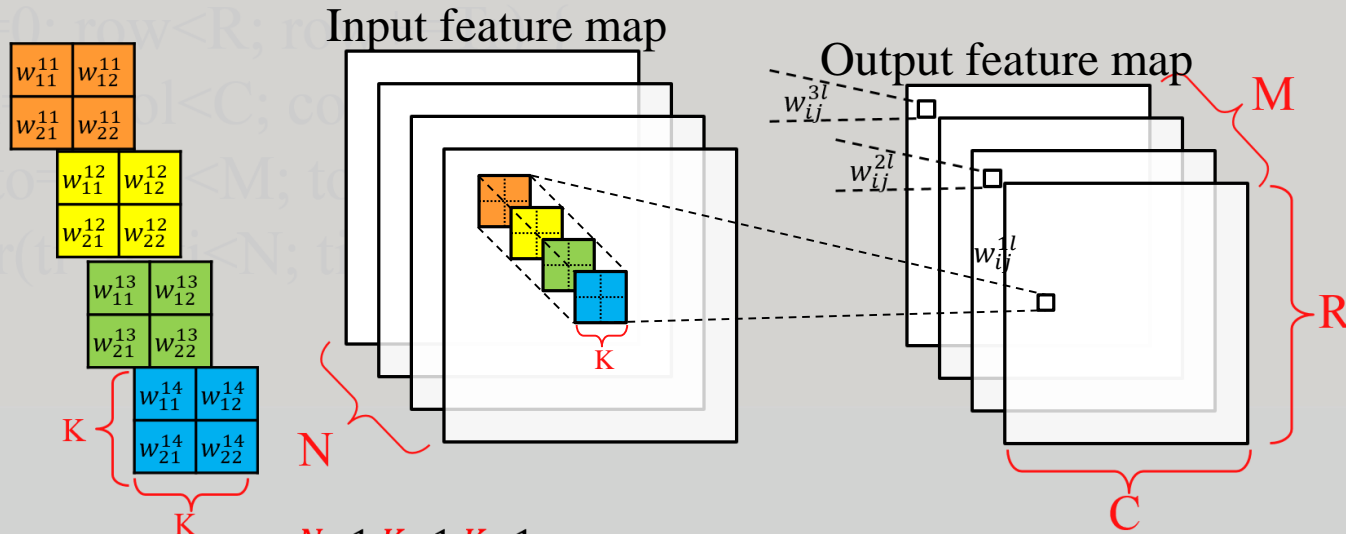
# Attainable Performance



# Attainable Performance



# Loop tiling



$$Y[m][r][c] = \sum_{n=0}^{N-1} \sum_{i=0}^{K-1} \sum_{j=0}^{K-1} W[m][n][i][j] \times X[n][S \times r + i][S \times c + j]$$

```

1 for(row=0; row<R; row++) {
2   for(col=0; col<C; col++) {
3     for(to=0; to<M; to++) {
4       for(ti=0; ti<N; ti++) {
5         for(i=0; i<K; i++) {
6           for(j=0; j<K; j++) {
               output_fm[to][row][col] +=
               weights[to][ti][i][j]*input_fm[ti][S*row+i][S*col+j];
           }
         }
       }
     }
   }
 }

```

**R, C, M, N, K, S** are all configuration parameters of the convolutional layer



# Loop tiling

```
1 for(row=0; row<R; row+=Tr) { (Tile loop)
2   for(col=0; col<C; col+=Tc) { (Tile loop)
3     for(to=0; to<M; to+=Tm) { (Tile loop)
4       for(ti=0; ti<N; ti+=Tn) { (Tile loop)
```

## Off-chip Data Transfer: Memory Access Optimization

### On-chip Data: Computation Optimization

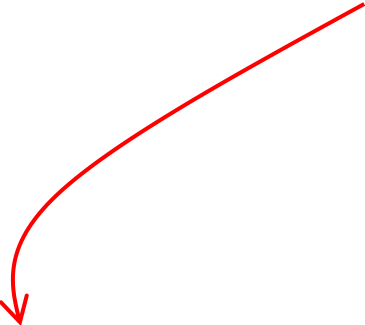
```
5   for(trr=row; trr<min(row+Tr, R); trr++) { (Point loop)
6     for(tcc=col; tcc<min(tcc+Tc, C); tcc++) { (Point loop)
7       for(too=to; too<min(to+Tn, M); too++) { (Point loop)
8         for(tii=ti; tii<(ti+Tn, N); tii++) { (Point loop)
9           for(i=0; i<K; i++) { (Point loop)
10            for(j=0; j<K; j++) { (Point loop)
                output_fm[to][row][col] +=
                weights[to][ti][i][j]*input_fm[ti][S*row+i][S*col+j];
            }
          }
        }
      }
    }
  }
}
```



# Computation Optimization

(Tile loops)

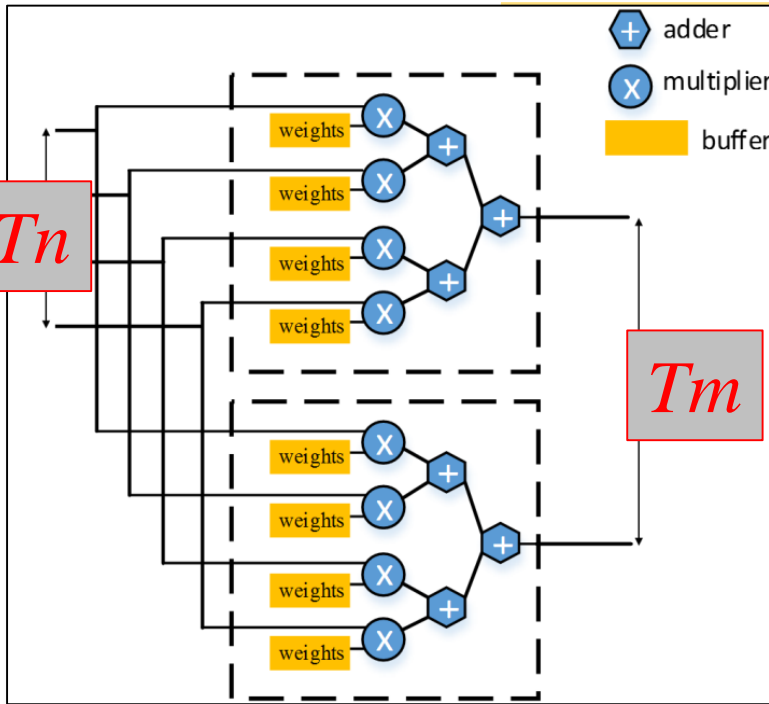
```
5 for(trr=row; trr<min(row+Tr, R); trr++) (Point loop)
6   for(tcc=col; tcc<min(tcc+Tc, C); tcc++) (Point loop)
7     for(too=to; too<min(to+Tn, M); too++) (Point loop)
8       for(tii=ti; tii<(ti+Tn, N); tii++) (Point loop)
9         for(i=0; i<K; i++) {
10          for(j=0; j<K; j++) {
              output_fm[to][row][col] +=
              weights[to][ti][i][j]*input_fm[ti][S*row+i][S*col+j];
          }}}}}}
```



```
5 for(i=0; i<K; i++) {
6   for(j=0; j<K; j++) {
7     for(trr=row; trr<min(row+Tr, R); trr++) {
8       for(tcc=col; tcc<min(tcc+Tc, C); tcc++) {
#pragma HLS pipeline
9         for(too=to; too<min(to+Tm, M); too++) {
#pragma HLS UNROLL
10          for(tii=ti; tii<(ti+Tn, N); tii++) {
#pragma HLS UNROLL
              output_fm[to][row][col] +=
              weights[to][ti][i][j]*input_fm[ti][S*row+i][S*col+j];
          }}}}}}
```



# Computation Optimization



(Tile loops)

```

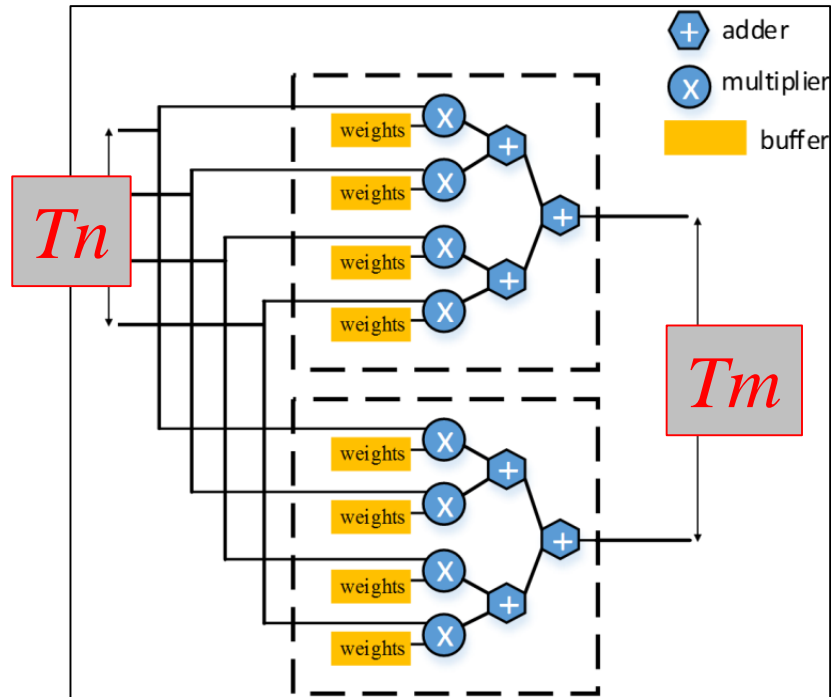
for(i=0; i<K; i++) {
  for(j=0; j<K; j++) {
    for(trr=row; trr<min(row+Tr, R); trr++) {
      for(tcc=col; tcc<min(tcc+Tc, C); tcc++) {
        HLS pipeline
        for(too=to; too<min(to+Tm, M); too++) {
          HLS UNROLL
          for(tii=ti; tii<(ti+Tn, N); tii++) {
            HLS UNROLL
            output_fm[to][row][col] +=
            weights[to][tj][i][j]*input_fm[ti][S*row+i][S*col+j];
          }}}}}
    }
  }
}
    
```

$$\text{Computational Performance} = \frac{\text{total number of operations}}{\text{execution cycles}}$$

$$\begin{aligned} \text{execution cycles} &= \frac{R}{T_r} \times \frac{C}{T_c} \times \left\lceil \frac{M}{T_m} \right\rceil \times \left\lceil \frac{N}{T_n} \right\rceil \times (K \times K \times T_c \times T_r + P) \\ &\approx \left\lceil \frac{M}{T_m} \right\rceil \times \left\lceil \frac{N}{T_n} \right\rceil \times R \times C \times K \times K \end{aligned}$$



# Computational Performance



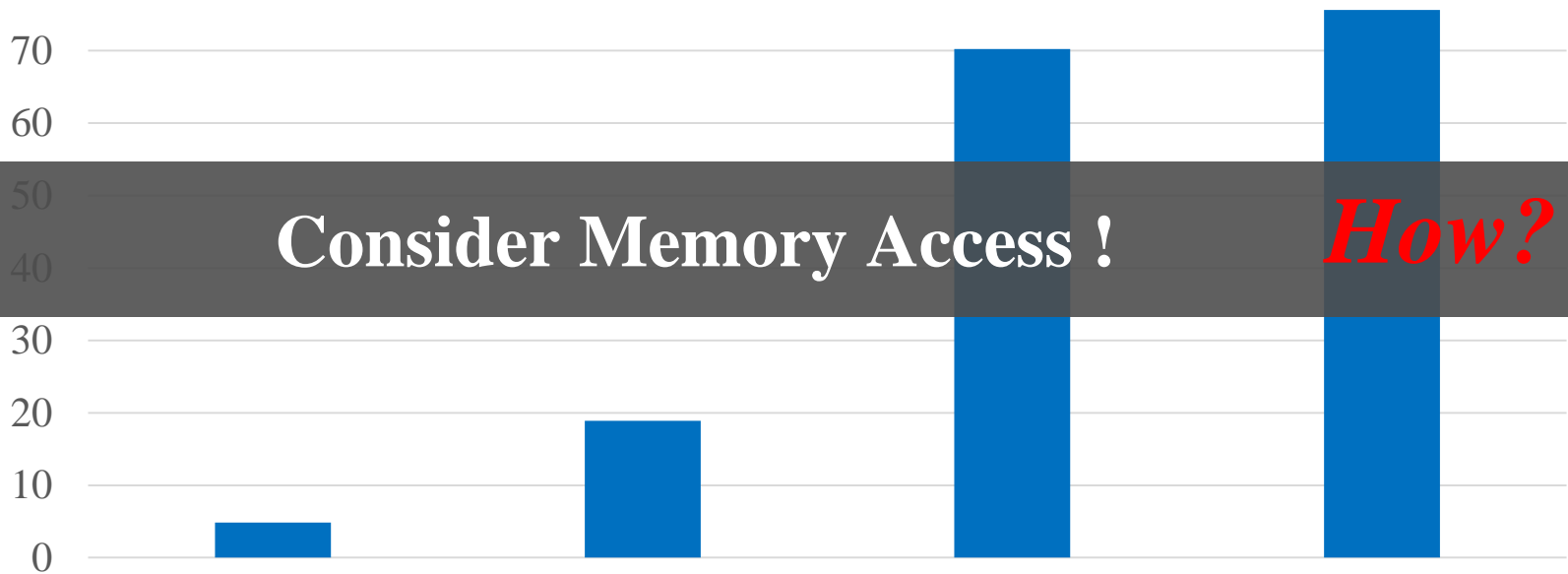
	Design 1	Design 2	Design 3	Design 4
Output ( $T_m$ )	5	10	20	30
Input ( $T_n$ )	5	10	20	15
DSPs	125	500	2000	2250





# Computational Performance

$$\text{Computational Performance} = \frac{\text{total number of operations}}{\text{execution cycles}}$$



Consider Memory Access !

How?

	Design 1	Design 2	Design 3	Design 4
Output ( $T_m$ )	5	10	20	30
Input ( $T_n$ )	5	10	20	15
DSPs	125	500	2000	2250



# Memory Access Optimization

```
1 for(row=0; row<R; row+=Tr) {  
2   for(col=0; col<C; col+=Tc) {  
3     for(to=0; to<M; to+=Tm) {  
4       for(ti=0; ti<N; ti+=Tn) {
```

```
    load output feature map
```

```
    S: foo(output_fm(to, row, col));
```

```
    store output feature map
```

```
  }
```

```
}}
```

```
5   for(trr=row; trr<min(row+Tr, R); trr++)  
6     for(tcc=col; tcc<min(tcc+Tc, C); tcc++)  
7       for(too=to; too<min(to+Tn, M); too++)  
8         for(tii=ti; tii<(ti+Tn, N); tii++)  
9           for(i=0; i<K; i++) {  
10            for(j=0; j<K; j++) {  
                output_fm[to][row][col] +=  
                    weights[to][ti][i][j]*  
                    input_fm[ti][S*row+i][S*col+j];  
            }  
        }  
    }  
}
```



# Memory Access Optimization

```

1 for(row=0; row<R; row+=Tr) {
2   for(col=0; col<C; col+=Tc) {
3     for(to=0; to<M; to+=Tm) {
4       for(ti=0; ti<N; ti+=Tn) {
         load output feature map
         S: foo(output_fm(to, row, col));
         store output feature map
       }
     }
  }
}

```

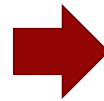
```

1 for(row=0; row<R; row+=Tr) {
2   for(col=0; col<C; col+=Tc) {
3     for(to=0; to<M; to+=Tm) {
4       for(ti=0; ti<N; ti+=Tn) {
         load output feature map
         S: foo(output_fm(to, row, col));
       }
         store output feature map
     }
  }
}

```

Before local memory promotion

$$\begin{aligned}
 & \text{'output\_fm' memory access} \\
 & = 2 \times \frac{R}{Tr} \times \frac{C}{Tc} \times \frac{M}{Tm} \times \frac{N}{Tn} \times Size_{array}
 \end{aligned}$$



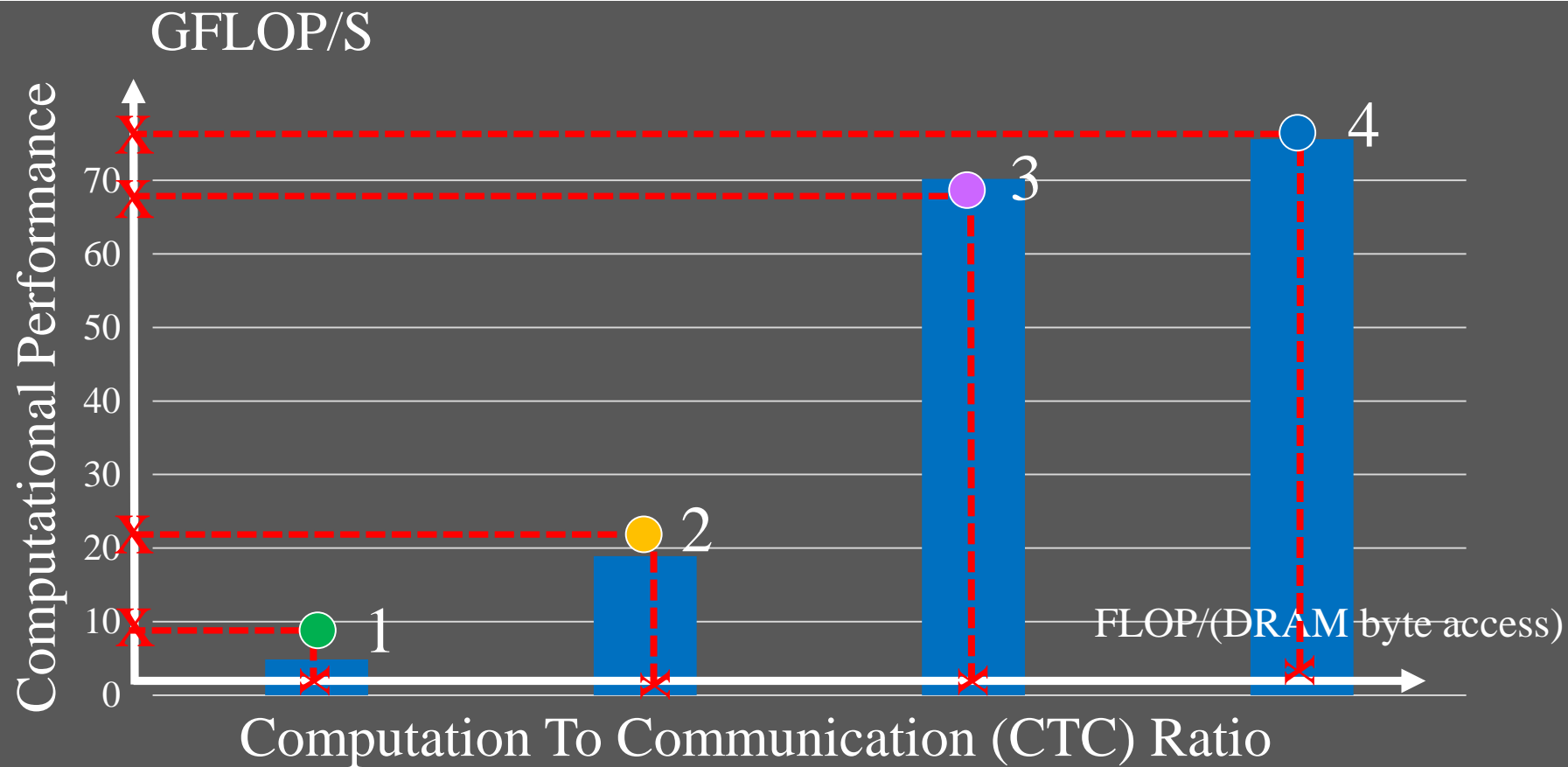
After local memory promotion

$$\begin{aligned}
 & \text{'output\_fm' memory access} \\
 & = 2 \times \frac{\bar{R}}{Tr} \times \frac{C}{Tc} \times \frac{M}{Tm} \times Size_{array}
 \end{aligned}$$

$$\text{Computation to Communication Ratio} = \frac{\text{Total number of operations}}{\text{Total amount of external data access}}$$



# Computation To Communication Ratio



	Design 1	Design 2	Design 3	Design 4
Output ( $T_m$ )	5	10	20	30
Input ( $T_n$ )	5	10	20	15



# Design Space

## Computation Engine:

Constraints for CNN configurations:

$$T_m \in (\text{Integer}, 1 < T_m < M) \quad N=128$$

$$T_n \in (\text{Integer}, 1 < T_n < N) \quad N=192$$

Constraints for FPGA resource:

$$T_n \times T_m \in (\text{Integer}, 1 < T_m \times T_n < \# \text{ of PE})$$

# of PE = 450

Legal Solutions  
of  $T_m$  &  $T_n$ :

2097

## Communication:

# of memory access methods

Legal Solutions:

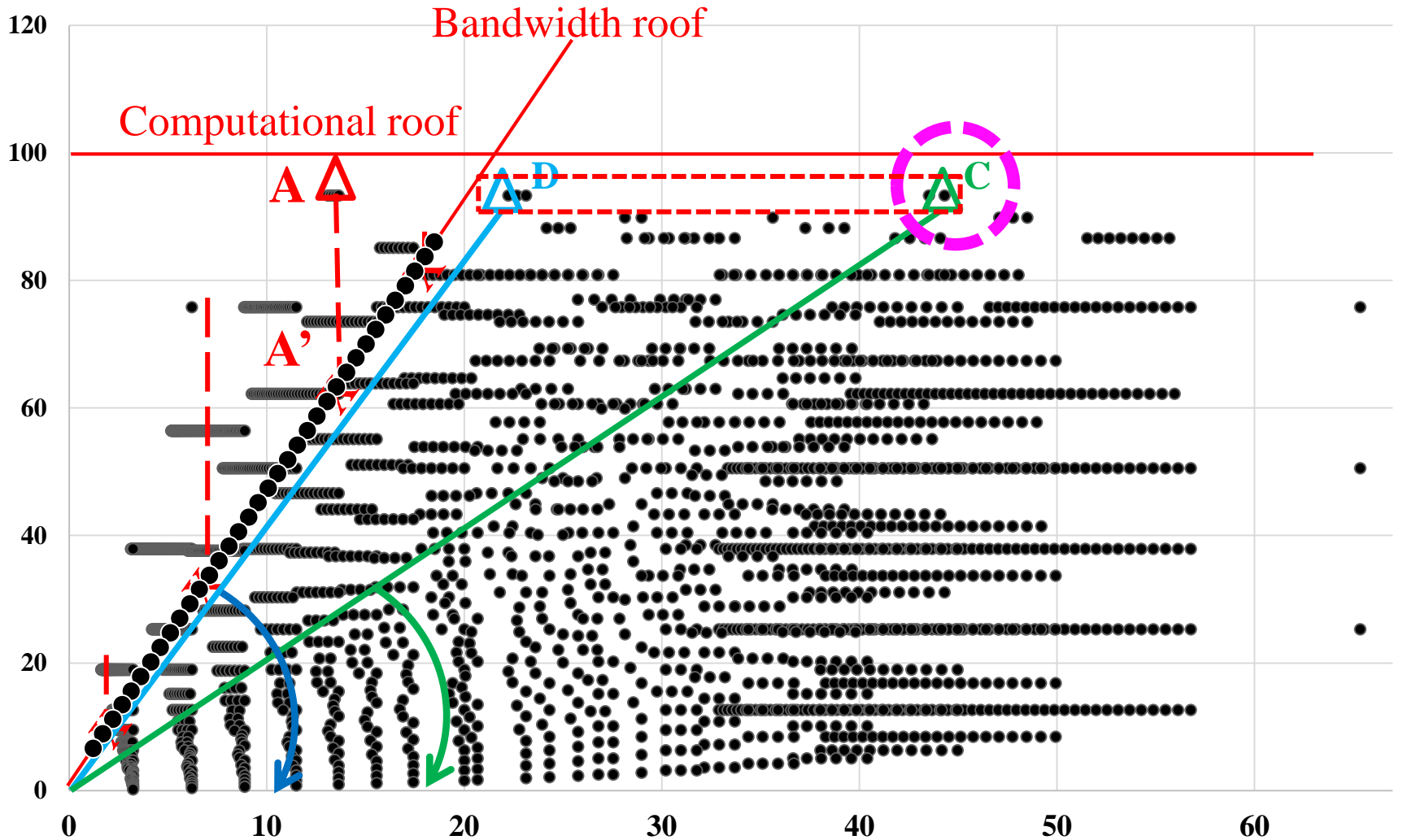
3

Total Legal Solutions:

6291

# Design Space Exploration

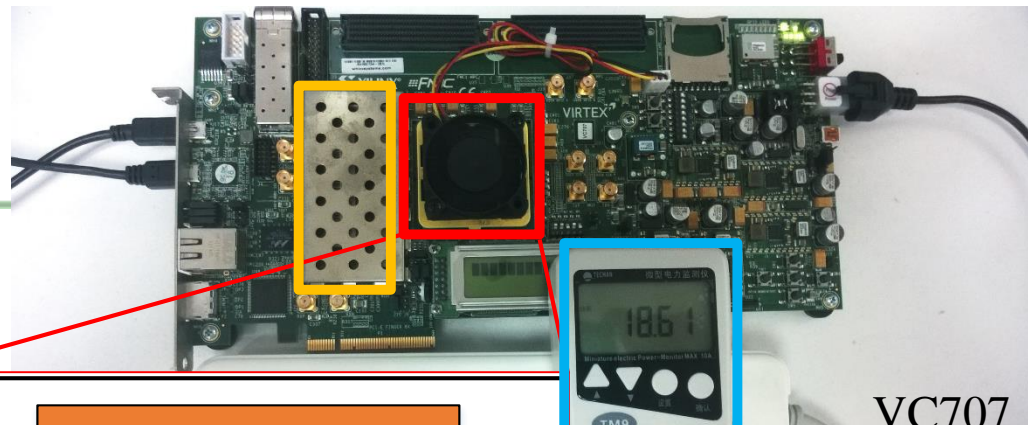
Attainable performance (GFLOPS)



CTC Ratio

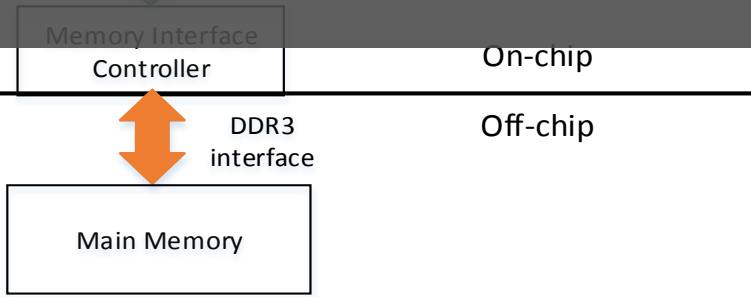
Ratio of  
(FLOP)/(DRAM byte access)

# System Overview

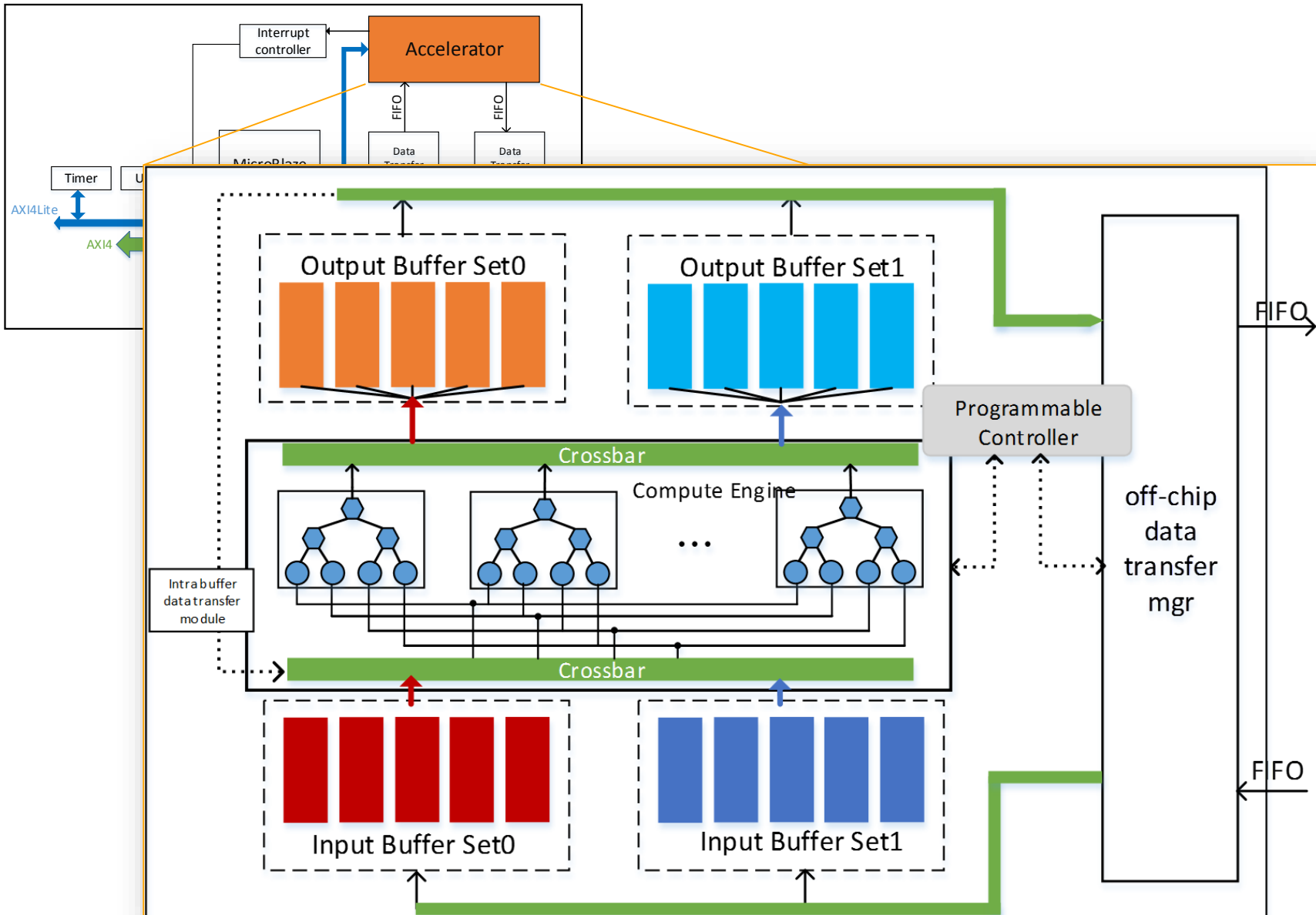


VC707

	Estimated	On-board run	Difference
Layer 1	7.32 ms	7.67 ms	~5%
Layer 2	5.11 ms	5.35 ms	~5%
Layer 3	3.42 ms	3.79 ms	~9%
Layer 4	2.59 ms	2.88 ms	~10%
Layer 5	1.73 ms	1.93 ms	~10%
<b>Overall</b>	<b>20.17 ms</b>	<b>21.61 ms</b>	<b>~7%</b>



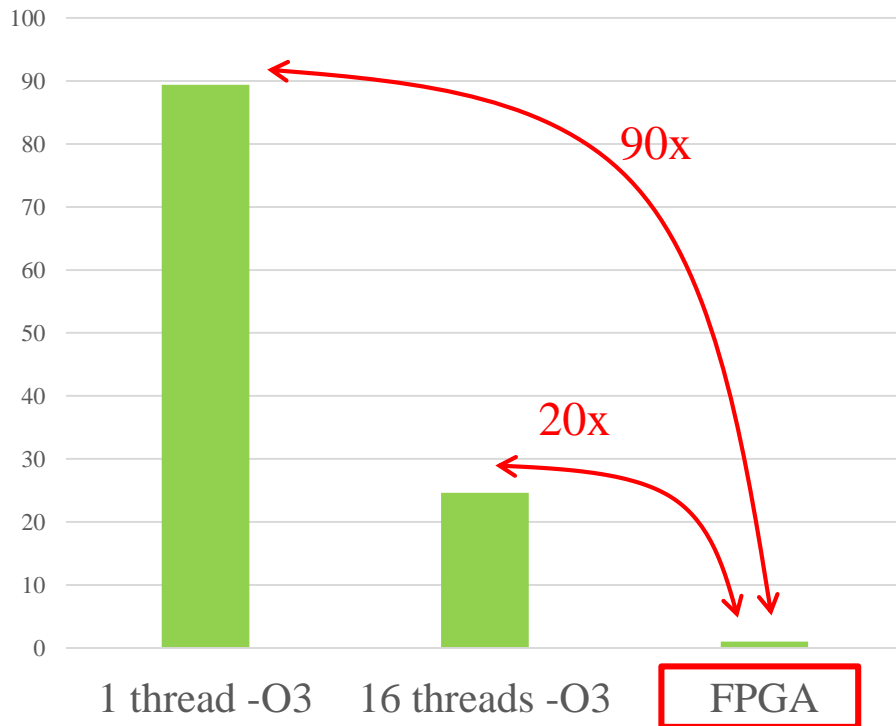
# Accelerator



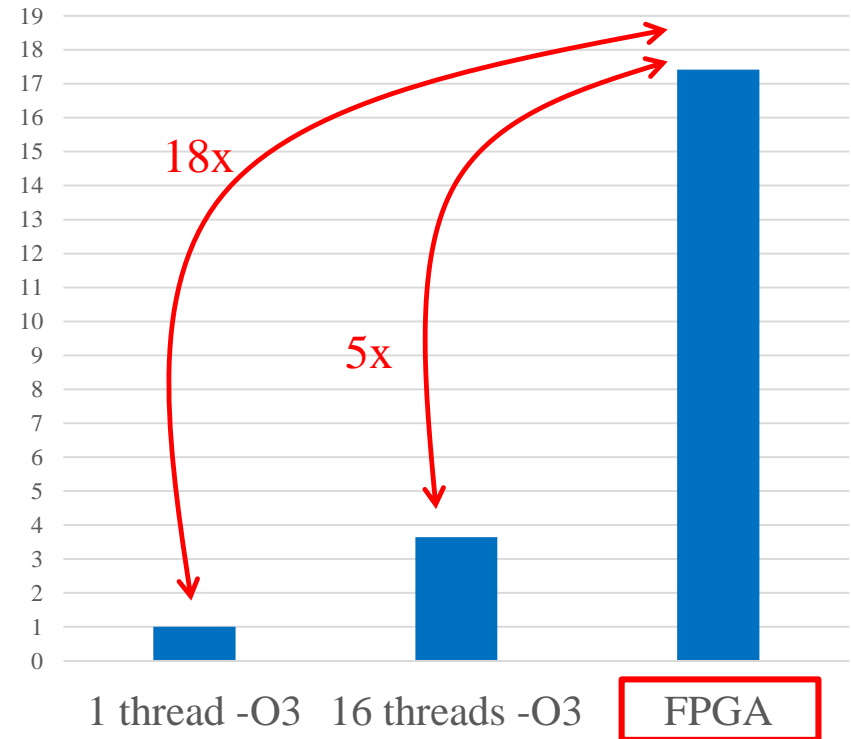


# Experimental Results: vs. CPU

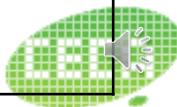
## Energy



## Speedup



CPU	Xeon E5-2430 (32nm)	16 cores	2.2 GHz	gcc 4.7.2 -O3 OpenMP 3.0
FPGA	Virtex7-485t (28nm)	448 PEs	100MHz	Vivado 2013.4 Vivado HLS 2013.4

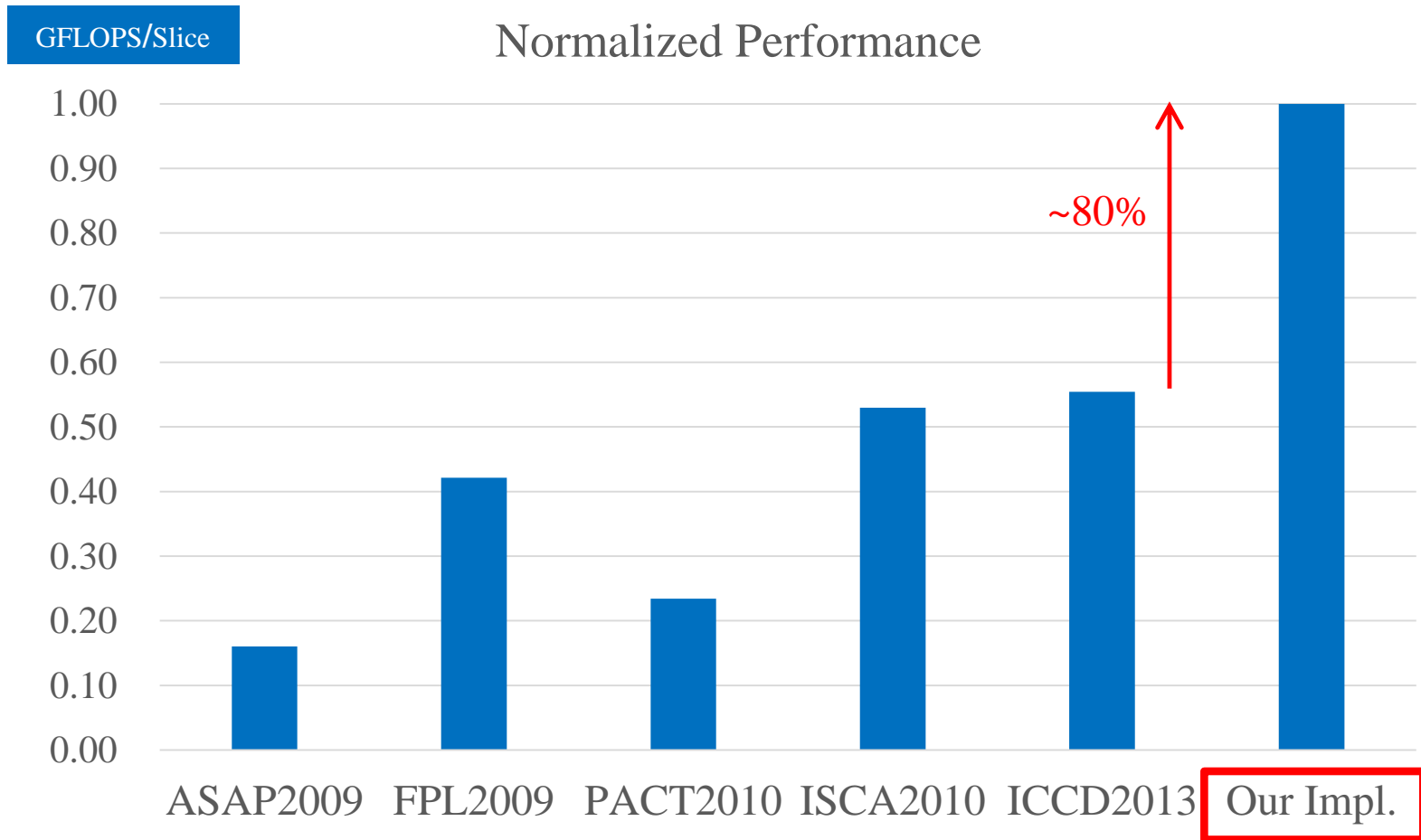


# *Experimental Results: vs. Other FPGAs*

FPL2009	ASAP2009	PACT2010	ISCA2010	ICCD2013	Our Impl.
Virtex 4	Virtex 5	Virtex 5	Virtex 5	Virtex 6	Virtex 7
125MHz	115MHz	125MHz	200MHz	150MHz	100MHz
5.25 GOPS	6.74 GOPS	7 GOPS	16 GOPS	17 GOPS	61.6 GOPS



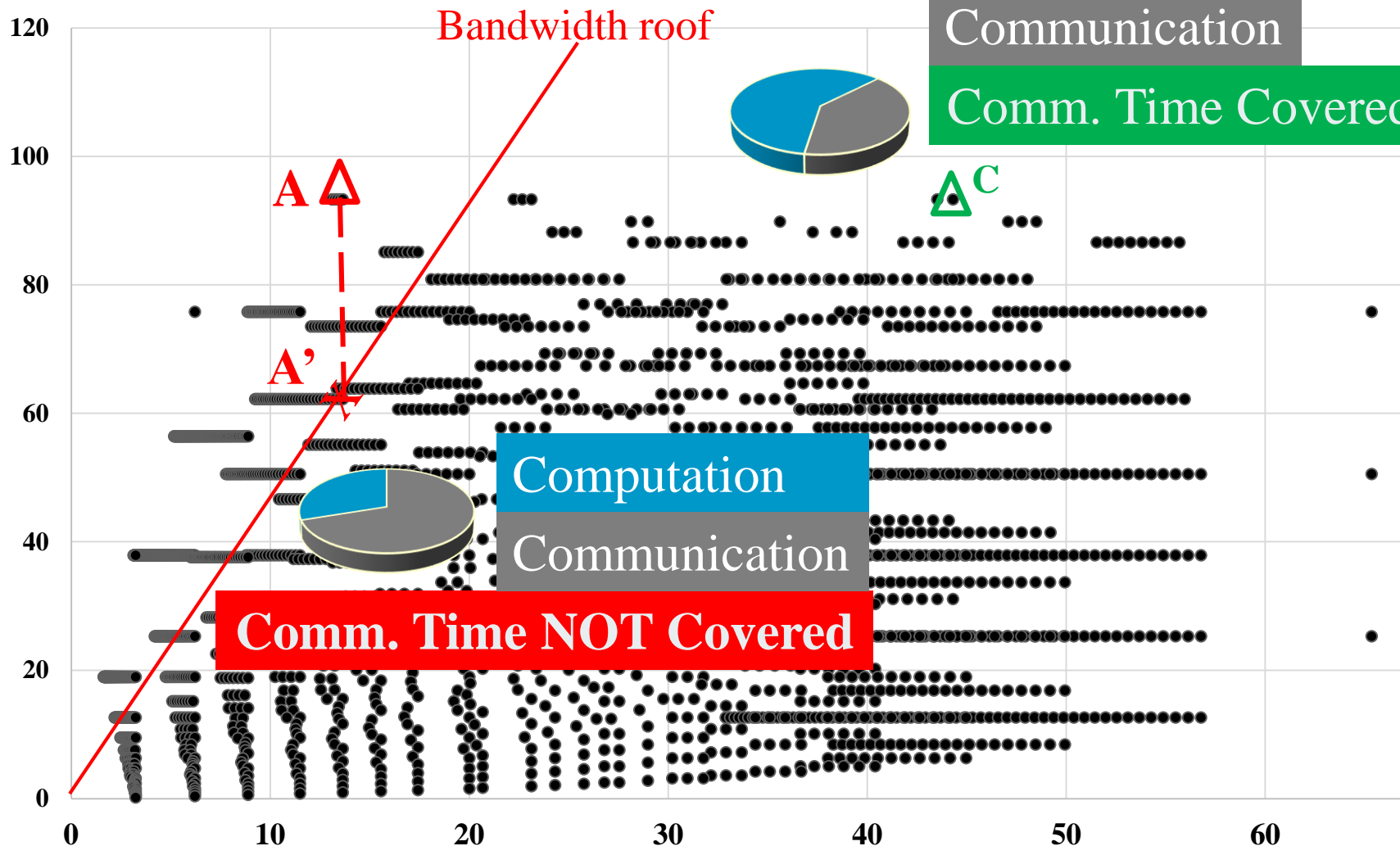
# Experimental Results: vs. Other FPGAs



# Experimental Results

Attainable performance (GFLOPS)

Computation  
Communication  
Comm. Time Covered



CTC Ratio

Ratio of  
(FLOP)/(DRAM byte  
access)

# Conclusions

□ An accelerator for convolutional neural network

## ➤ Contribution:

- Accurate Analytical model for **computation & communication**
- Find the best solution with **roofline model**

## ➤ Result:

- On-board run implementation
- **~3.5x** better performance over other FPGA implementations
- **~80%** performance/area improvement



---

***Thank You***

Q & A

