



復旦大學
FUDAN UNIVERSITY

Towards Efficient Data Parallelism on Spatial CGRA via Constraint Satisfaction and Graph Coloring

Authors: Yuan Dai, Xuchen Gao*, et al.

Presenter: Xuchen Gao

Institution: Fudan University, China

Email: xcgao22@m.fudan.edu.cn

Contents

- Introduction
- Problem Formulation
- Proposed Methodology
- Experiment
- Conclusion

Contents

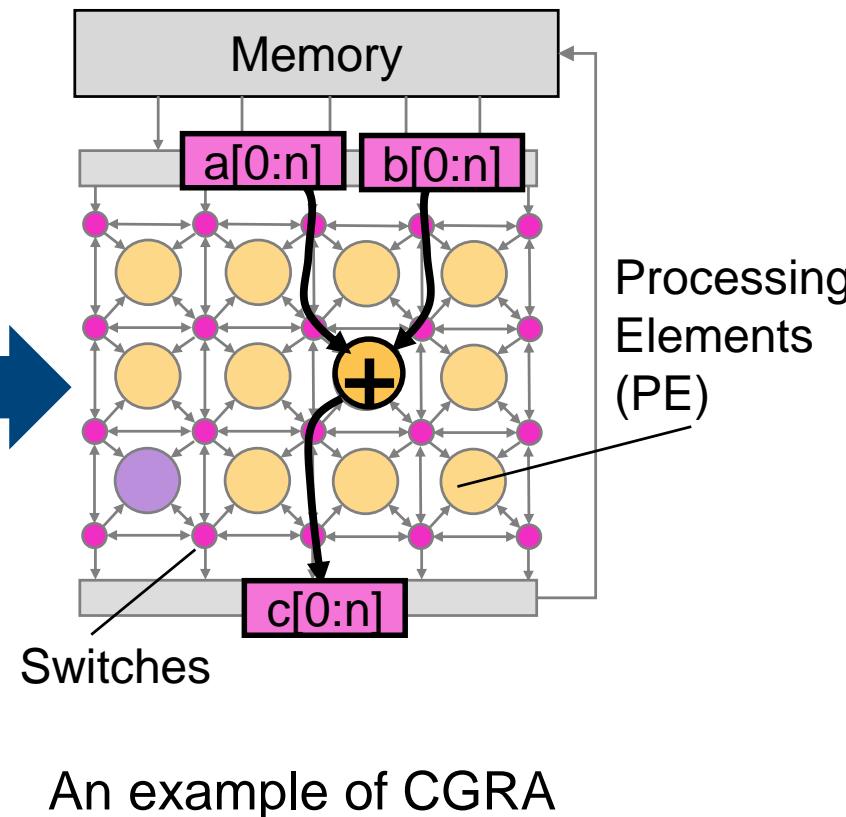
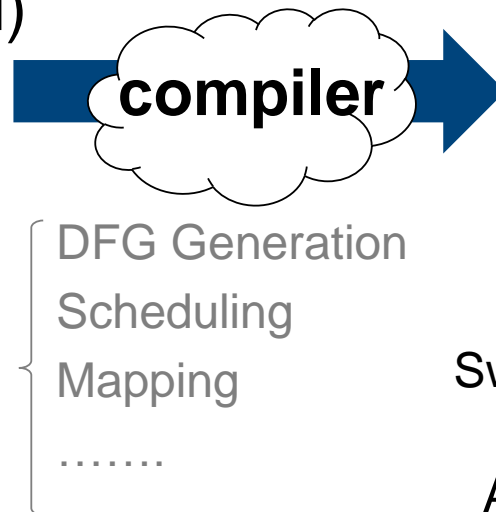
- Introduction
- Problem Formulation
- Proposed Methodology
- Experiment
- Conclusion

Background

- The Workflow of Coarse-Grained Reconfigurable Architecture (CGRA)

```
for (int i = 0; i < n; ++i)  
  c[i] = a[i] + b[i];
```

application code

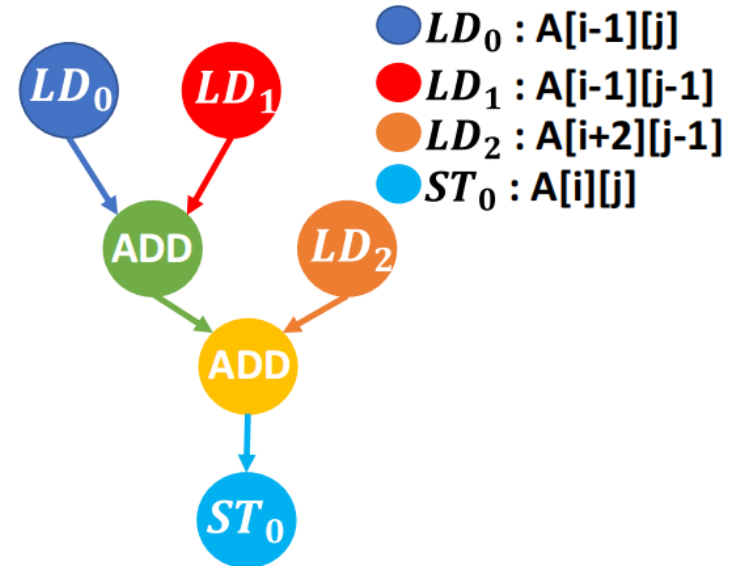


Motivation

■ Motivation Example

```
#define N 32
int A[N][N];
void example ( ) {
    for( int i = 1; i < N ; i ++ )
        for(int j = 1; j < N ; j ++ )
            A[i][j] = A[i-1][j] + A[i-1][j-1]
                + A[i+2][j-1];
}
```

Motivating example code

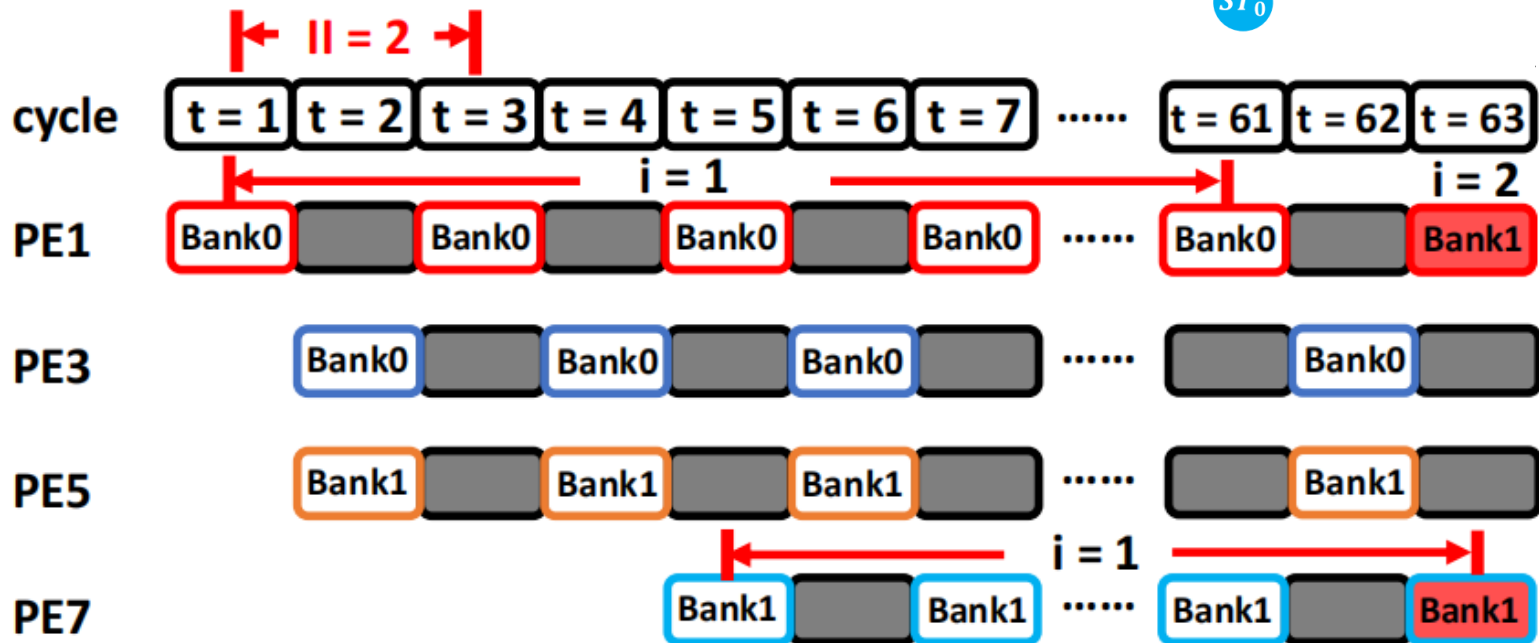
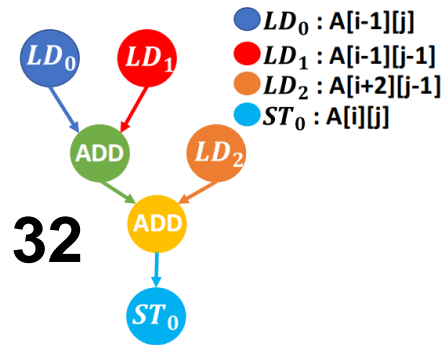


DFG of example code

Motivation

■ Motivation Example

- Invalid scheme with $N = 2$ and $B = 32$



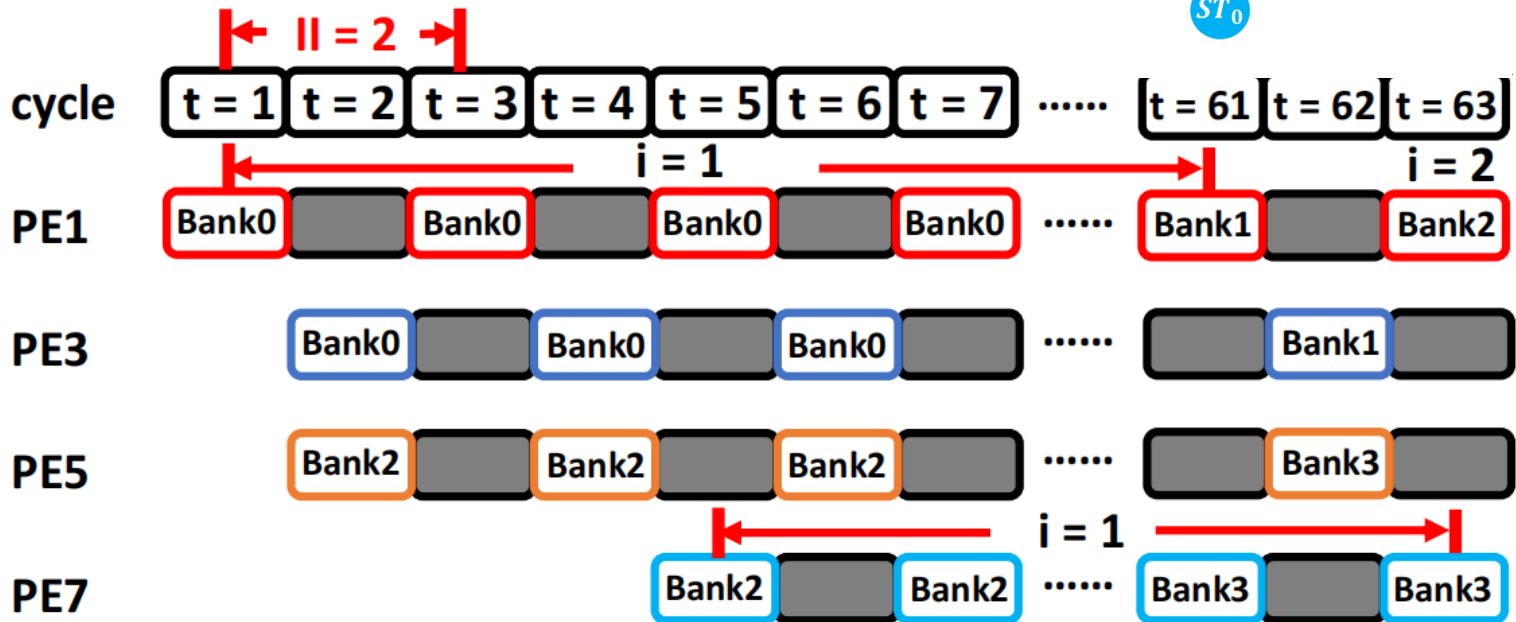
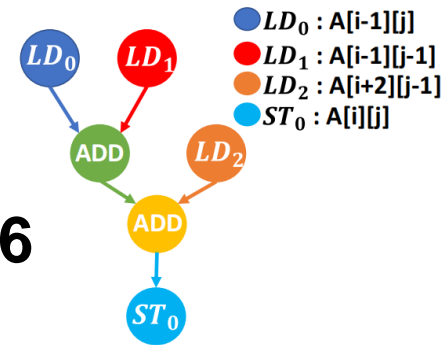
Bank0	Bank1
$A[0][0 \sim 31]$	$A[1][0 \sim 31]$
$A[2][0 \sim 31]$	$A[3][0 \sim 31]$
.....

Bank index = $\left\lfloor \frac{x}{32} \right\rfloor \% 2$
 x : the address of data

Motivation

■ Motivation Example

- Valid scheme with $N = 4$ and $B = 16$



Bank0	Bank1	Bank2	Bank3
$A[0][0 \sim 15]$	$A[0][16 \sim 31]$	$A[1][0 \sim 15]$	$A[1][16 \sim 31]$
$A[2][0 \sim 15]$	$A[2][16 \sim 31]$	$A[3][0 \sim 15]$	$A[3][16 \sim 31]$
.....

Bank index = $\left\lfloor \frac{x}{16} \right\rfloor \% 4$
 x : the address of data

Contents

- Introduction
- **Problem Formulation**
- Proposed Methodology
- Experiment
- Conclusion

Definitions

- Affine Address Access

$$\phi(\vec{i}) = [AS_0 \ AS_1 \dots \ AS_{l-1}] \begin{bmatrix} i_0 \\ i_1 \\ \dots \\ i_{l-1} \end{bmatrix} + BA$$

- Access Pattern

$$P_A = \{\phi_1(\vec{i}_1), \dots, \phi_m(\vec{i}_m)\}$$

- Control Step Access Pattern

$$P_A^n = \{\phi_1(\vec{i}_1), \dots, \phi_q(\vec{i}_q)\}$$

q ($q \leq m$) memory accesses at n^{th} cycle

Definitions

- Memory Partition

$$f(\phi) = \lfloor \frac{\phi}{B} \rfloor \% N; \quad g(\phi) = \lfloor \frac{\phi}{N \cdot B} \rfloor \cdot B + \phi \% B$$

$f(\phi)$: bank function

$g(\phi)$: offset function

- Iteration Distance

$$iter_d^{\alpha\beta} = \frac{data_delay}{Mapped_II}$$

Target

- ① Access II

$$AII = MAX(Con_i), 1 \leq i \leq N$$

- ② Iteration Count Vector

$$\forall \vec{i}_\alpha, \vec{i}_\beta \in M, \forall \phi_\alpha(\vec{i}_\alpha), \phi_\beta(\vec{i}_\beta) \in P_A^n, \alpha \neq \beta$$

$$(\vec{i}_\alpha - \vec{i}_\beta) \cdot (1, ic_0, \dots, \prod_{t=0}^{l-2} ic_t) = iter_d^{\alpha\beta}, f(\phi_\alpha(\vec{i}_\alpha)) \neq f(\phi_\beta(\vec{i}_\beta))$$

- Given an l -level loop in the iteration domain M with m affine memory accesses on the same array

$$P_A = \{\phi_1(\vec{i}_1), \dots, \phi_m(\vec{i}_m)\}$$

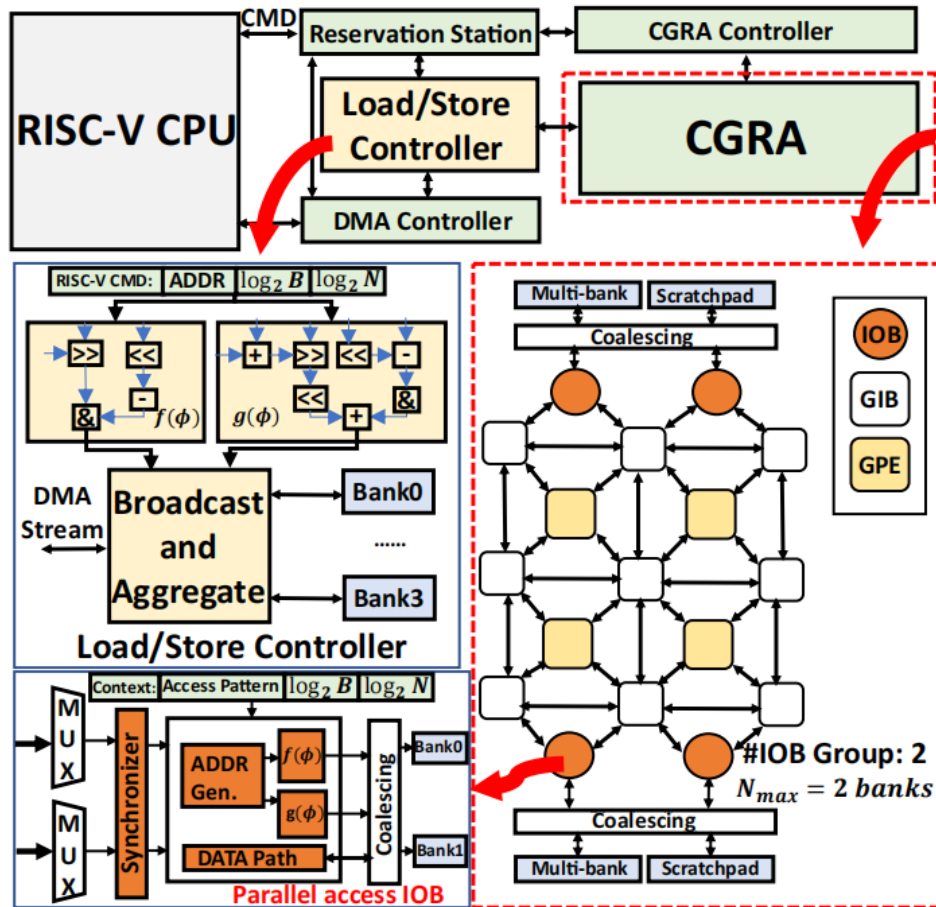
find the N and B such that:

✂ Minimize AII in ① & Hold the equation ②

Contents

- Introduction
- Problem Formulation
- **Proposed Methodology**
- Experiment
- Conclusion

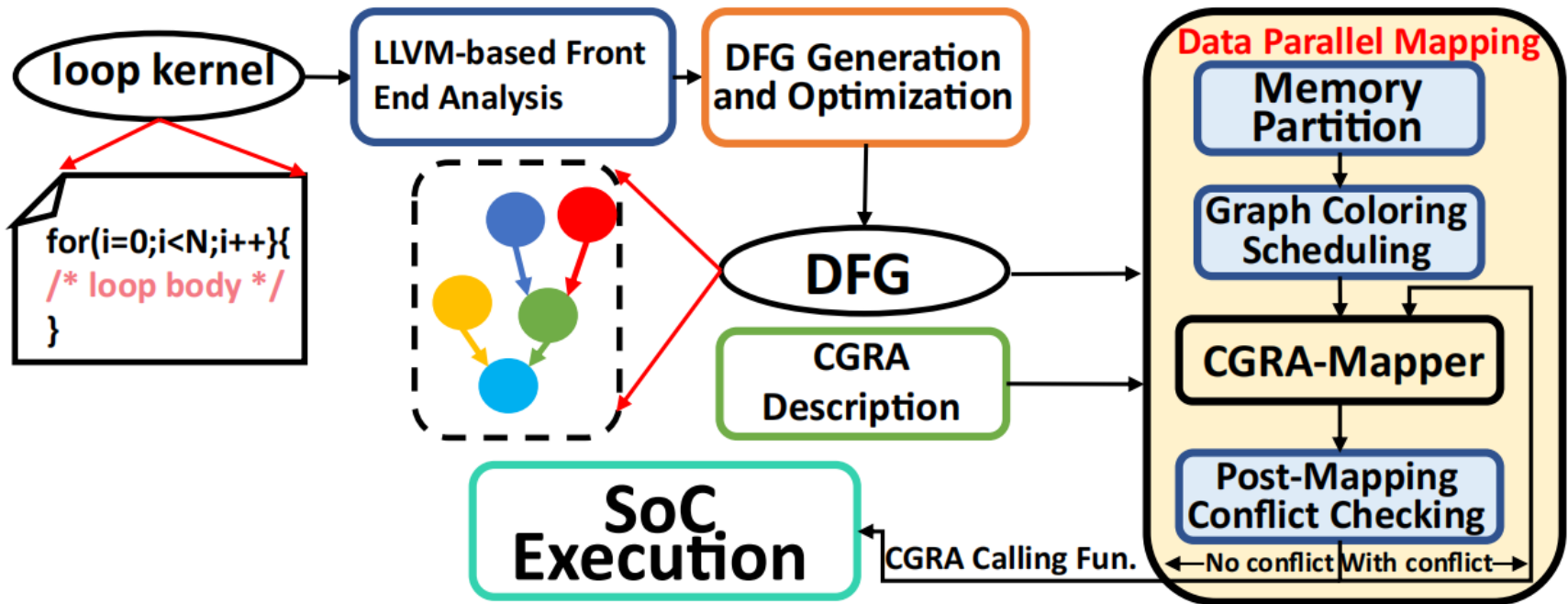
End-to-End Framework: Hardware



- Enhanced IOB to access different banks during execution
- Specific mechanism for loading/storing data from/to external memory

The SoC supports parallel memory access

End-to-End Framework: Software



Software flow with the proposed memory partition algorithm

End-to-End Framework: Software

```

1  #define PB_NI 32
2  #define PB_NJ 32
3  #define PB_NK 32
4  int C[_PB_NI][_PB_NI];
5  int A[_PB_NI][_PB_NI];
6  int B[_PB_NI][_PB_NI];
7  //kernel_23
8  void kernel13() {
9      #ifdef CGRA_COMPILER
10         loop_begin(); ①: added pragma for front-end tool
11     #endif
12     int alpha = 5;
13     for (int i = 0; i < _PB_NI; i++)
14         for (int j = 0; j < _PB_NJ; j++)
15             {
16                 for (int k = 0; k < _PB_NK; k = k + 4)
17                     C[i][j] = alpha * A[i][k] * B[k][j] + A[i][k+1] * B[k+1][j] + A[i][k+2] * B[k+2][j] + A[i][k+3] * B[k+3][j];
18             }
19
20     #ifdef CGRA_COMPILER
21         loop_end(); ①
22     #endif
23 }
24

```

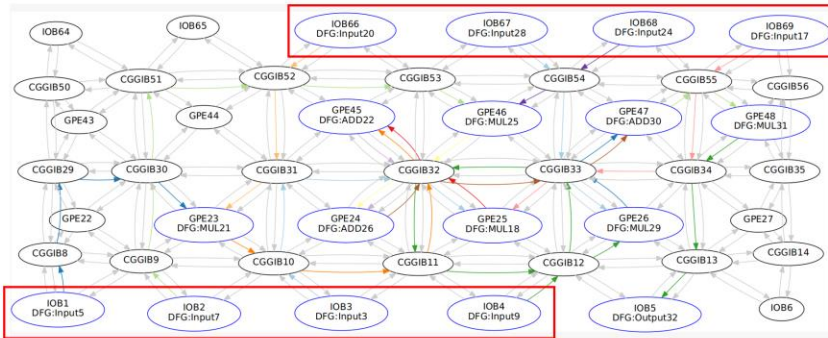
Gemm Source Code

```
Succeed to map DFG to ADG!<<<<<<
Final II = 1 !<<<<<
***** Memory partition and schedule results *****
Array Name: A; N: 4; B: 1
>> Step 0<<
Input3 Input5 Input7 Input9

Array Name: B; N: 4; B: 32
>> Step 0<<
Input17 Input20 Input24 Input28

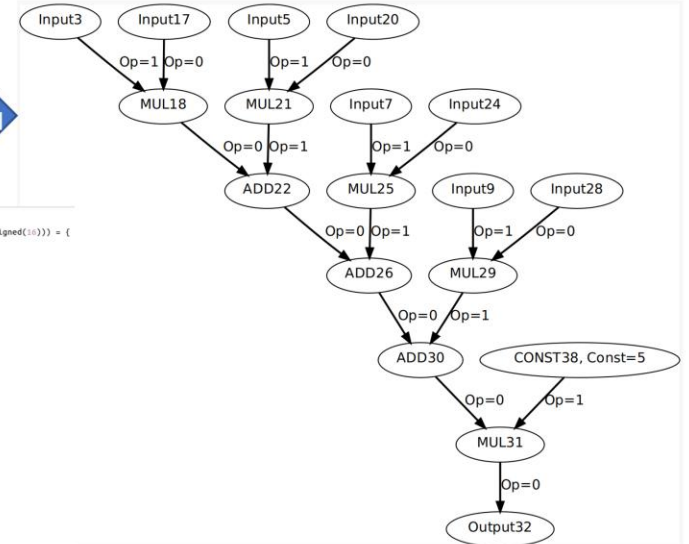
PE usage: 0.666667
IO usage: 0.75
DFG latency: 10
```

Mapper Report



Mapping Result

LLVM-based front-end tool



Generated DFG

```
*Verdi* : Begin traversing the scopes, layer (0).
*Verdi* : Enable +all dumping.
*Verdi* : End of traversing.
*Verdi* : fsdbDumpon - All FSDB files at 0 ps.
Initialization finished!
CPU add finished!
It takes 30611 cycles for CPU to finish the task.
It takes 2432 cycles for CGRA to finish the task(1).
CGRA comput finished!
Checking results!
Done!
*Verdi* : fsdbDumpoff - All FSDB files at 862.145.500 ps.
```

Simulation Result

CGRA Calling Function

CSP formulation

- Constraint Satisfaction Problem(CSP)
 - **Iteration Constraint(pre-mapping):** Before mapping, all the accesses in the P_A are regarded in the same iteration

$$\vec{i}_\alpha = \vec{i}_\beta = \vec{i}; \vec{i} \in M$$

- **Bank Constraint:** $\phi_\alpha(\vec{i}_\alpha)$ and $\phi_\beta(\vec{i}_\beta)$ access the same bank equation:

$$f(\phi_\alpha(\vec{i}_\alpha)) = f(\phi_\beta(\vec{i}_\beta)) \Leftrightarrow \lfloor \frac{\phi_\alpha(\vec{i}_\alpha)}{B} \rfloor \% N = \lfloor \frac{\phi_\beta(\vec{i}_\beta)}{B} \rfloor \% N$$

$$\Leftrightarrow \exists k \in \mathbb{Z}, s.t. \lfloor \frac{\phi_\alpha(\vec{i}_\alpha)}{B} \rfloor = \lfloor \frac{\phi_\beta(\vec{i}_\beta)}{B} \rfloor + kN$$



$$\exists k \in \mathbb{Z}, s.t. -B < \boxed{\phi_\alpha(\vec{i}_\alpha) - \phi_\beta(\vec{i}_\beta) - kNB} < B$$



CSP formulation

- Constraint Satisfaction Problem(CSP)

$$-B < \left[(AS_0^\alpha - AS_0^\beta) \dots (AS_{l-1}^\alpha - AS_{l-1}^\beta) \right] \begin{bmatrix} i_0 \\ i_1 \\ \dots \\ i_{l-1} \end{bmatrix} + (BA^\alpha - BA^\beta) - kNB < B$$

※**CSP problem**: the conflict exists when one solution is found under such the constraints

Graph Coloring-based Scheduling

- When failing to find a conflict-free solution
 - schedule the conflict accesses into different control steps
 - minimizing the All

$$P_A = \{ \phi(LD_0), \phi(LD_1), \phi(LD_2), \phi(ST_0) \}$$

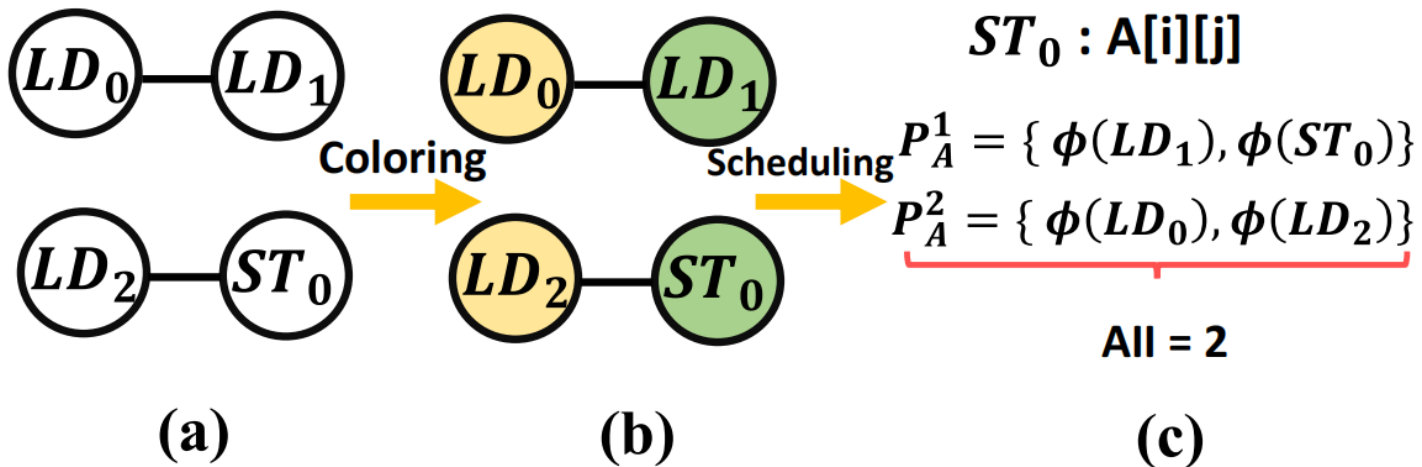
Partition parameters: $N = 4, B = 16$

$LD_0 : A[i-1][j]$

$LD_1 : A[i-1][j-1]$

$LD_2 : A[i+2][j-1]$

$ST_0 : A[i][j]$



Example of graph coloring and scheduling process

CSP-based Post-mapping Checking

- Iteration Constraint: now $iter_d^{\alpha\beta} \neq 0$

$$(\vec{i}_\alpha - \vec{i}_\beta) \cdot (1, ic_0, \dots, \prod_{t=0}^{l-2} ic_t) = iter_d^{\alpha\beta}$$

It can be transformed into a CSP problem, where the conflict exists when one solution is found

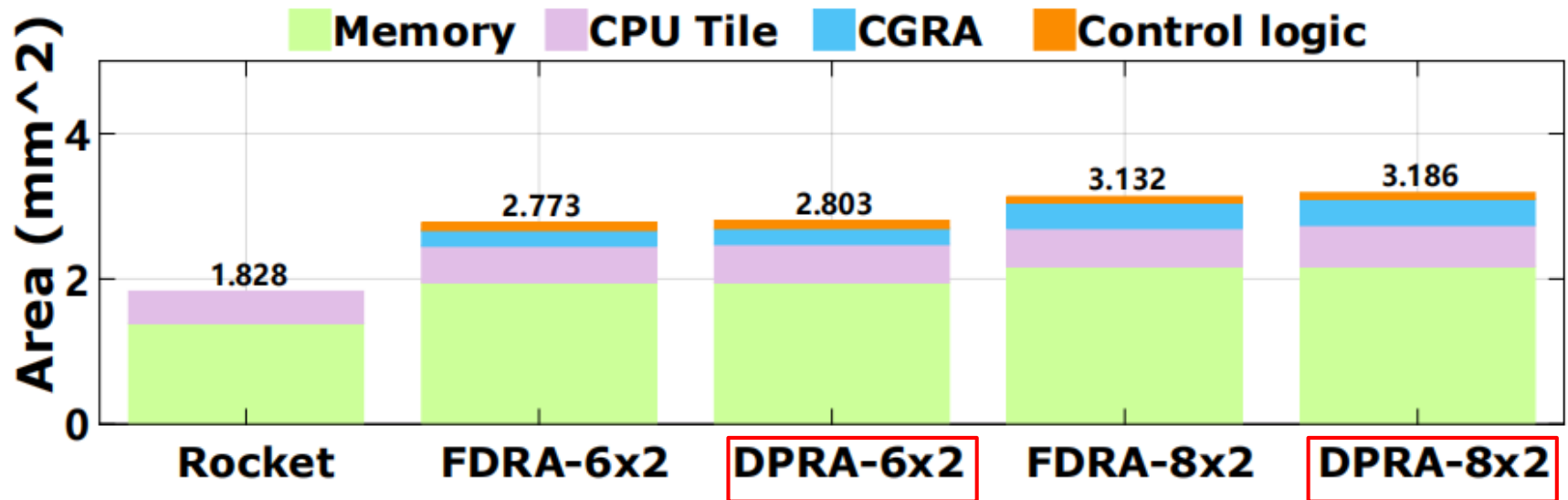
- CSP Check

$$\exists k \in \mathbb{Z}, \text{ s.t. } -B < \phi_\alpha(\vec{i}_\alpha) - \phi_\beta(\vec{i}_\beta) - kNB < B$$

Contents

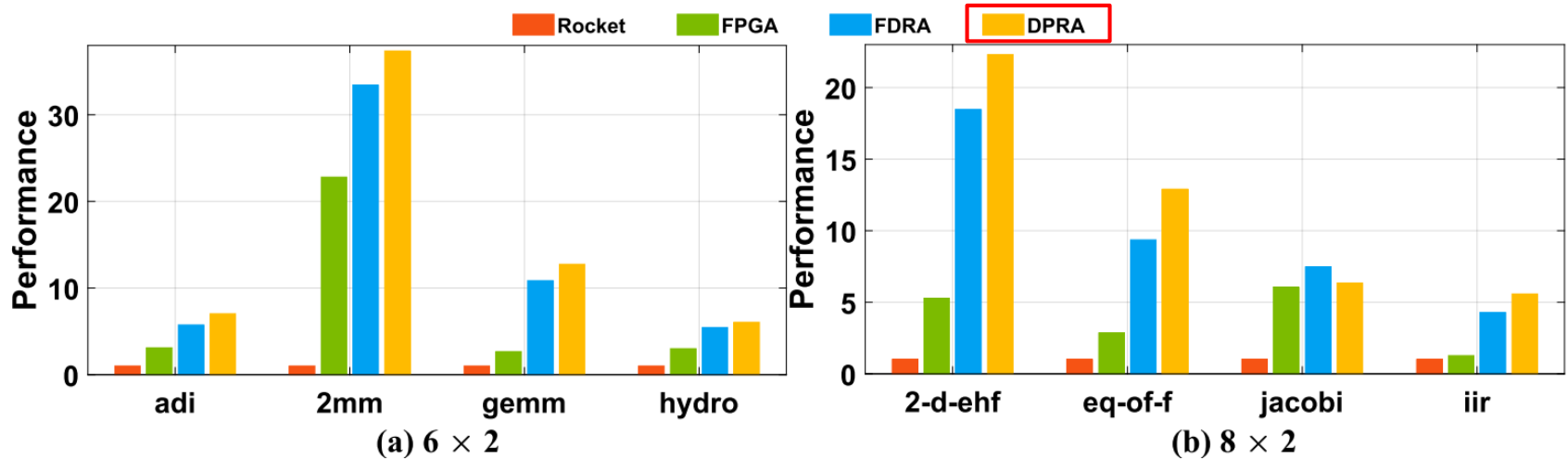
- Introduction
- Problem Formulation
- Proposed Methodology
- **Experiment**
- Conclusion

System Level Evaluation

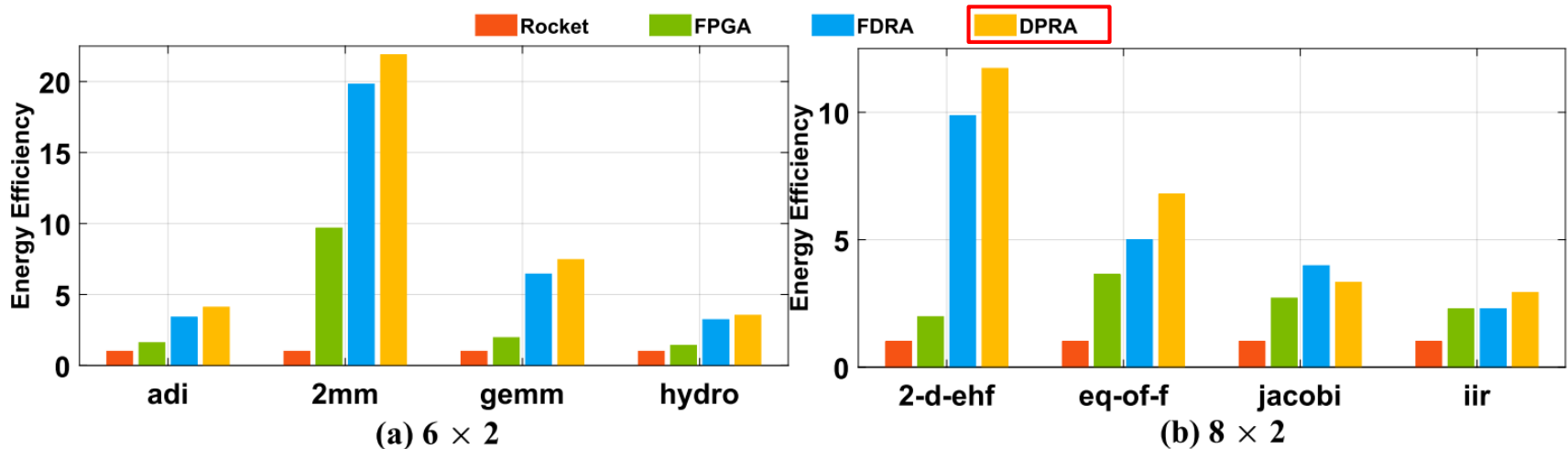


The area of each SoC

System Level Evaluation

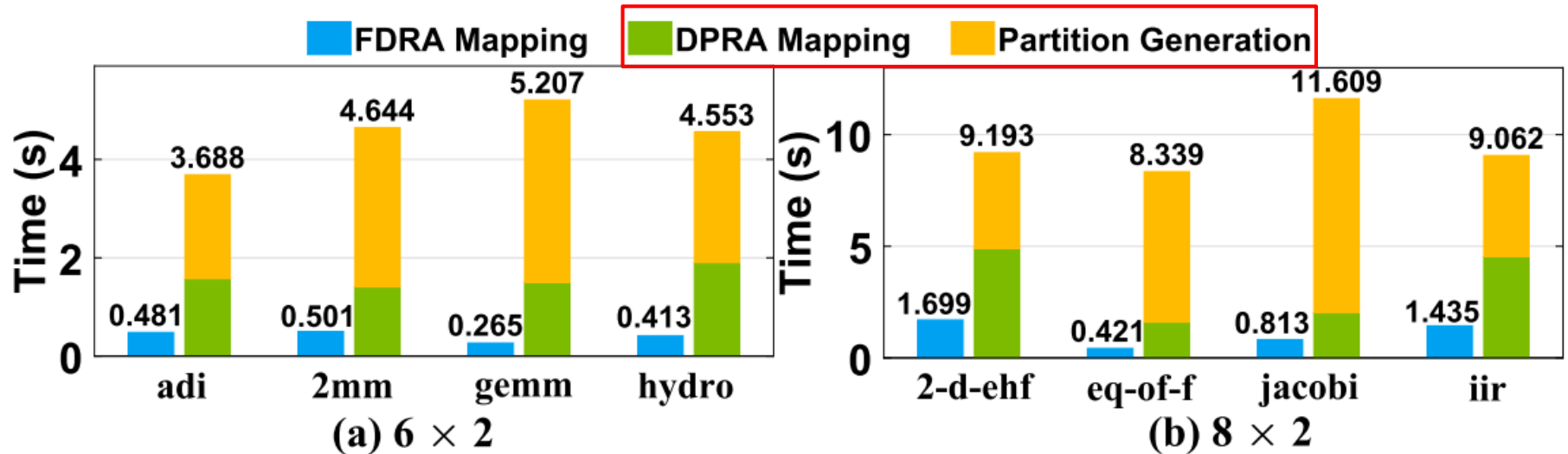


Comparison of performance



Comparison of energy efficiency

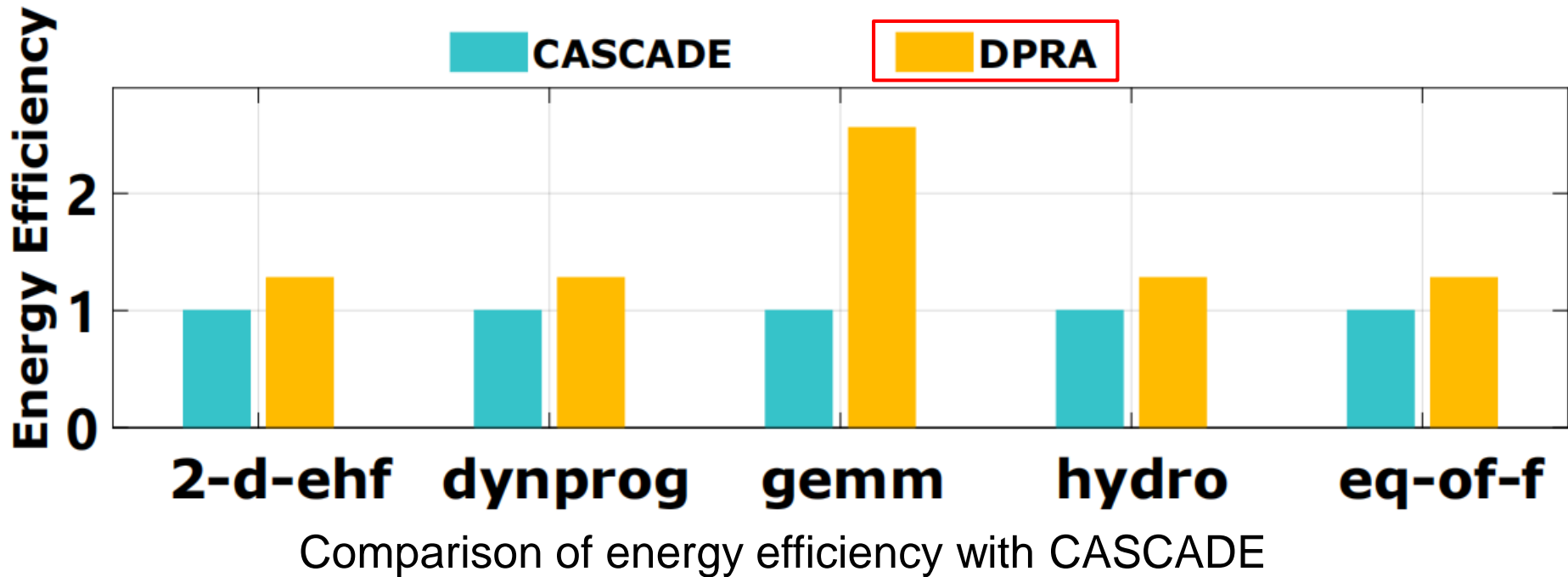
System Level Evaluation



Comparison of compilation time

CGRA Level Evaluation

- $Energy\ efficiency = \frac{Throughput}{Power\ consumption}$
- $Throughput = \frac{Frequency}{Mapped_II}$



Contents

- Introduction
- Framework of our work
- Analyzing & transforming passes
- SoC runtime configuration of CGRA
- Experiment
- **Conclusion**

Conclusion

- A CSP-based memory conflict detection for memory partition
- A graph coloring-based approach for memory access scheduling during the partition process to improve performance
- An end-to-end CGRA framework, including the data parallel architecture and compiler with efficient data parallelism support



Thanks

Yuan Dai, Xuchen Gao, Chen Shen, Bingbing Peng,
Wenbo Yin, Wai-Shing Luk*, Lingli Wang*

State Key Laboratory of Integrated Chips and Systems
Fudan University, Shanghai, China
Email: xcgao22@m.fudan.edu.cn, *llwang@fudan.edu.cn