# Imperial College London

# Exploiting Run-time Reconfiguration In Stencil Computation

Xinyu Niu, *Qiwei Jin, Wayne Luk, Qiang Liu and Oliver Pell*

Dept. of Computing, School of Engineering, Imperial College London, UK
School of Electronic Information Engineering, Tianjin University, China
Maxeler Technologies, UK

# Overview

- Background

- Partitioning algorithm

- Analytical model

- Reconfiguration scheduling

- Results

- Future work

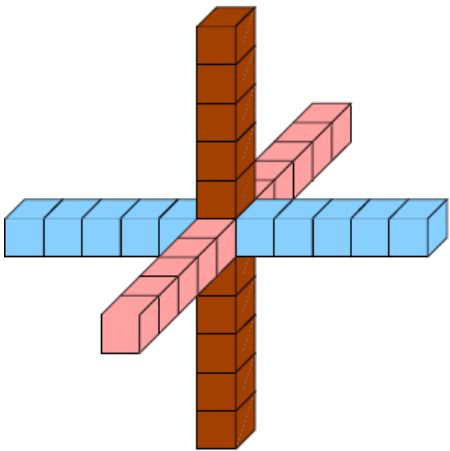- Conclusion

# Background: run-time reconfiguration

- Purna and Bhatia 1999

  - temporal partitioning and data flow graph scheduling

- Singhal and Bozorgzadeh 2006

  - floorplanning for reconfiguration

- Bruneel, Abouelella and Stroobandt 2009

  - mapping applications to self-reconfiguring platform

- Iskander 2010

  - accelerate design validation

- Koch and Torresen 2011

  - run-time reconfigurable sorter for large-scale problems

- Purna and Bhatia 1999

    - temporal partitioning and data flow graph scheduling

- Singhal and Bozorgzadeh 2006

    - floorplanning for reconfiguration

- Bruneel, Abouelella and Stroobandt 2009

    - mapping applications to self-reconfiguring platform

- Iskander 2010

    - accelerate design validation

- Koch and Torresen 2011

    - run-time reconfigurable sorter for large-scale problems

- Our focus: automate reconfigurable stencil computation

# Contributions

1. identify reconfiguration opportunities
   - −partitioning algorithm: generate configurations based on variations in time dimension
2. analyse runtime benefits for stencil computation
   - −analytical model: optimise generated partitions, to fully utilise available resources
3. evaluate run-time solutions
   - −Scheduling algorithm: evaluate run-time benefits and overhead, to provide optimal run-time solution
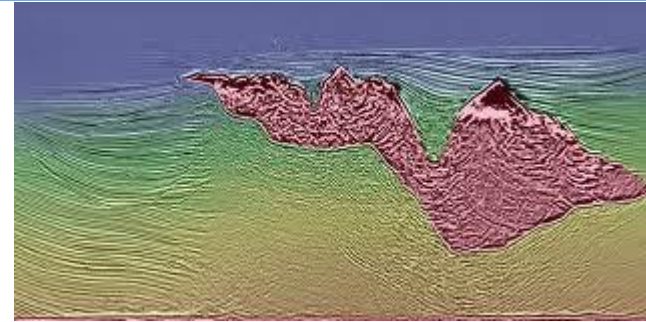
$$\frac{\partial^2 p}{\partial t^2} = v^2 \nabla^2 p + s(t)$$

Example: Oil & Gas Application

```
for t = 1 to tmax:
    for i = 0 to X*Y*Z-1:
        l = convolve(curr[i], stencil)
        next := 2 * curr[i] - prev[i] + vv[i] * l
        next[i] := next + source[i]
        apply_boundary_condition(curr, next)
        swap_buffers(prev, curr, next)
```
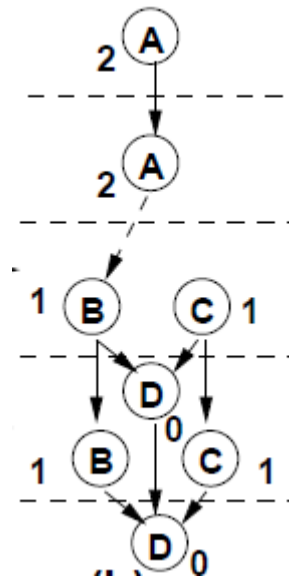
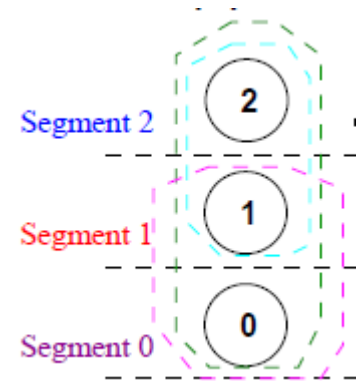Forward propagation

Forward propagation

Backward propagation

correlation

for

A

for

B    C

D

- segments
  - functions in same data dependency level combined as a segment
  - segment variations express algorithm variations in time dimension

- segments
  - functions in same data dependency level con[...] segment
  - segment variations express algorithm variati[...] dimension



- configurations
  - one or more segments are combined as a valid configuration to be executed
  - optimised configurations utilise run-time benefits

- segments
  - functions in same data dependency level combined as a segment
  - segment variations express algorith dimension

- configurations
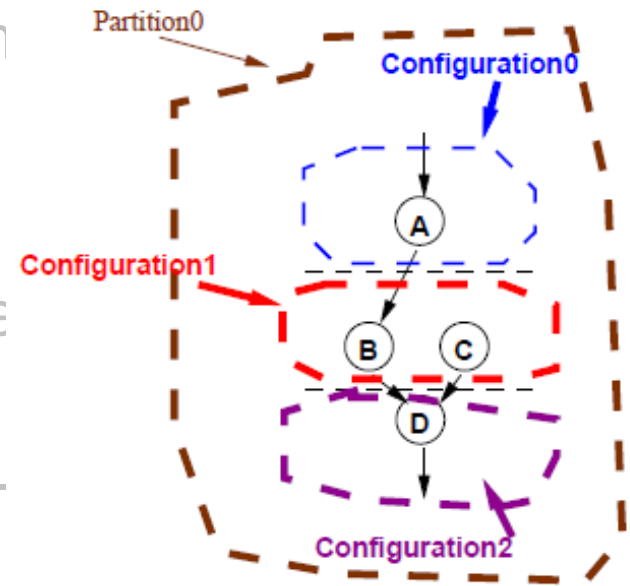  - one or more segments are combine configuration to be executed
  - optimised configurations utilise run-



- partitions
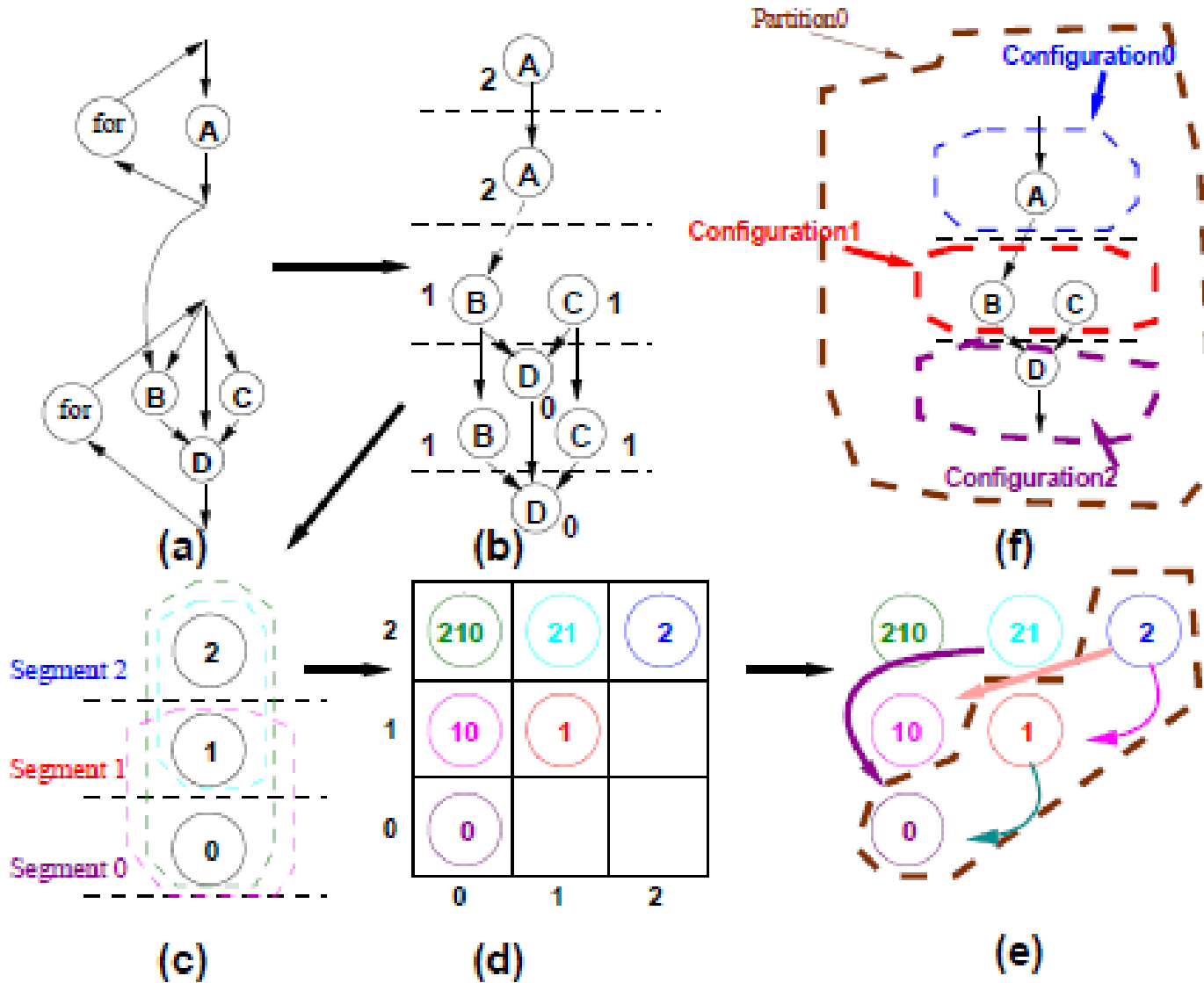  - one or more configurations are coordinated to accomplish application tasks
  - optimal partition balances run-time benefits and overhead

(a)　(b)　(c)　(d)　(e)　(f)

- fully pipelined data-paths
- bit-width optimisation, arithmetic operation transform

- fully pipelined data-paths

- bit-width optimisation, arithmetic operation transform

- fully pipelined data-paths
- bit-width optimisation and arithmetic operation transformation

- accumulating resource consumption

$$Ds = \frac{\sum N_{\mathrm{op,i}} \cdot B_{\mathrm{D}} \cdot T_{\mathrm{D,i}}}{A_{\mathrm{D}}}$$ ⟵ DSP consumption

$$Ls = \frac{\sum (N_{\mathrm{op,i}} \cdot B_{\mathrm{L}} \cdot T_{\mathrm{L,i}}) + I_{\mathrm{L}}}{A_{\mathrm{L}}}$$ ⟵ LUT consumption

$$Fs = \frac{\sum (N_{\mathrm{op,i}} \cdot B_{\mathrm{F}} \cdot T_{\mathrm{F,i}}) + I_{\mathrm{F}}}{A_{\mathrm{F}}}$$ ⟵ FF  consumption

- customised memory architecture



Data from Memory

(a)

(b)

- customised memory architecture
- data access blocking

- customised memory architecture
- analysing resource / bandwidth requirements

memory resource consumption

$$Bs = \frac{\sum_{1}^{P_{\text{knl}} \cdot P_{\text{t}} \cdot P_{\text{dp}}} m_{\text{dp}} \cdot (S \cdot (2 + N_{\text{c}}) + 1)}{A_{\text{B}} / (W_{\text{dp}} \cdot B_{\text{w}})}$$

$$m_{\text{dp}} = \frac{nx + (P_{\text{t}} - 1) \cdot 2S}{P_{\text{dp}}} \cdot (ny + (P_{\text{t}} - 1) \cdot 2S)$$

- customised memory architecture
- analysing resource / bandwidth requirements

$$Bs = \frac{\sum_1^{P_{\mathrm{knl}} \cdot P_{\mathrm{t}} \cdot P_{\mathrm{dp}}} m_{\mathrm{dp}} \cdot (S \cdot (2 + N_{\mathrm{c}}) + 1)}{A_{\mathrm{B}} / (W_{\mathrm{dp}} \cdot B_{\mathrm{w}})}$$

$$m_{\mathrm{dp}} = \frac{nx + (P_{\mathrm{t}} - 1) \cdot 2S}{P_{\mathrm{dp}}} \cdot (ny + (P_{\mathrm{t}} - 1) \cdot 2S)$$

off-chip bandwidth requirement

$$BW_{\mathrm{m}} \geq (W_{\mathrm{dp}} \cdot B_{\mathrm{w}} \cdot P_{\mathrm{dp}}) \cdot P_{\mathrm{knl}} \cdot f_{\mathrm{knl}}$$

$$W_{\mathrm{m}} = N \cdot (W_{\mathrm{dp}} \cdot B_{\mathrm{w}} \cdot P_{\mathrm{dp}}) \quad N \in \{1, 2, 3...\}$$

- parallel data-paths: multiple memory access
- serial: multiple time steps

- Objective

Minimise: $C_t = R_{c/c} \cdot \dfrac{D \cdot O_b \cdot O_t}{f_{knl} \cdot P_{knl} \cdot P_{dp} \cdot P_t}$

$$O_b = \frac{\frac{x - 2 \cdot S}{nx - 2 \cdot S} \cdot \frac{y - 2 \cdot S}{ny - 2 \cdot S} \cdot nx \cdot ny}{x \cdot y}$$

$$nx = \frac{x - 2 \cdot S}{\alpha} + 2 \cdot S \quad ny = \frac{y - 2 \cdot S}{\beta} + 2 \cdot S$$

$$O_t = \begin{cases} 1 & \alpha \cdot \beta = 1 \\ \frac{nx + (P_t - 1) \cdot 2 \cdot S}{nx} \cdot \frac{ny + (P_t - 1) \cdot 2 \cdot S}{ny} & \alpha \cdot \beta \neq 1 \end{cases}$$

- Objective

- Constraints

Subject to:

$$1 \geq Ds = \frac{\sum N_{\text{op,i}} \cdot B_{\text{D}} \cdot T_{\text{D,i}}}{A_{\text{D}}}$$

$$1 \geq Ls = \frac{\sum (N_{\text{op,i}} \cdot B_{\text{L}} \cdot T_{\text{L,i}}) + I_{\text{L}}}{A_{\text{L}}}$$

Minimise:
$$C_{\text{t}} = R_{\text{c/c}} \cdot \frac{D \cdot O_{\text{b}} \cdot O_{\text{t}}}{f_{\text{knl}} \cdot P_{\text{knl}} \cdot P_{\text{dp}} \cdot P_{\text{t}}}$$

$$1 \geq Fs = \frac{\sum (N_{\text{op,i}} \cdot B_{\text{F}} \cdot T_{\text{F,i}}) + I_{\text{F}}}{A_{\text{F}}}$$

$$O_b = \frac{\frac{x - 2 \cdot S}{nx - 2 \cdot S} \cdot \frac{y - 2 \cdot S}{ny - 2 \cdot S} \cdot nx \cdot ny}{x \cdot y}$$

$$1 \geq Bs = \frac{\sum_{1}^{P_{\text{knl}} \cdot P_{\text{t}} \cdot P_{\text{dp}}} m_{\text{dp}} \cdot (S \cdot (2 + N_{\text{c}}) + 1)}{A_{\text{B}} / (W_{\text{dp}} \cdot B_{\text{w}})}$$

$$nx = \frac{x - 2 \cdot S}{\alpha} + 2 \cdot S \quad ny = \frac{y - 2 \cdot S}{\beta} + 2 \cdot S \quad BW_{\text{m}} \geq (W_{\text{dp}} \cdot B_{\text{w}} \cdot P_{\text{dp}}) \cdot P_{\text{knl}} \cdot f_{\text{knl}}$$

$$O_t = \begin{cases} 1 & \alpha \cdot \beta = 1 \\ \frac{nx + (P_{\text{t}} - 1) \cdot 2 \cdot S}{nx} \cdot \frac{ny + (P_{\text{t}} - 1) \cdot 2 \cdot S}{ny} & \alpha \cdot \beta \neq 1 \end{cases}$$

- Objective
- Constraints
- Overhead
  - reconfiguration time
  - data transfer time

$$C_{\mathrm{re}} = \frac{\gamma \cdot Max(B_{\mathrm{s}}, D_{\mathrm{s}}, L_{\mathrm{s}}, F_{\mathrm{s}})}{\theta}$$

$$C_{\mathrm{m}} = \frac{2 \cdot D \cdot W_{\mathrm{dp}} \cdot (1 + N_{\mathrm{c}}) \cdot B_{\mathrm{w}}}{\phi}$$

**Algorithm 3** Partition Scheduling Algorithm.

**Labels:** $v_i$: nodes , $p_i$: partitions , Cur: current configuration
**Functions** $\text{Conf}(v_i, p_i)$: find the configuration in partition $p_i$ that $v_i \in c_i$

1: **for** $p_i \in$ Partitions **do**
2:     **for** $v_i \in$ Source Nodes **do**
3:         $\text{Cur} \leftarrow \text{Conf}(v_i, p_i)$
4:         **while** $v_i.\text{NextNode} \neq \emptyset$ **do**
5:           **if** $v_i \notin \text{Cur}$ **then**
6:             $\text{Cur} \leftarrow \text{Conf}(v_i, p_i)$
7:             $T_{exe} \mathrel{+}= \text{Cur}.C_{re} + \text{Cur}.C_m$
8:           **end if**
9:           $T_{exe} \mathrel{+}= \text{Cur}.C_t$
10:           $v_i \leftarrow v_i.\text{NextNode}$
11:         **end while**
12:         $p_i.\text{T} \leftarrow \text{Max}(T_{exe})$
13:     **end for**
14:     Partition $\leftarrow \text{Min}(p_i.\text{T})$
15: **end for**

- fully utilise available resources

- fully utilise available resources

- model accuracy

- fully utilise available resources

- model accuracy

- approximating peak performance

| data size[3] | CPU[1] | | | GPU[2] | | | | FPGA (static) | | | FPGA (optimal) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | s | m | l | | | | [14] | s | m | l | s | m | l |
| execution time (t) | 181.7 | 1458.2 | 5574.8 | [10] | [11] | [12] | | 3.6 | 18.0 | 147.9 | 2.9 | 13.4 | 110.59 |
| overhead time (t)[4] | 0.03 | 0.01 | 0 | | | | 0[7] | 0.14 | 0.58 | 3.88 | 0.22 | 0.82 | 5.8 |
| throughput (GFlop/s)[5] | 1.8 | 0.9 | 1.8 | 36 | 51.2 | 64.5 | 35.8[7] | 70.6 | 68.0 | 66.8 | 102.8 | 91.6 | 91.6 |
| speed-up | 1x | 1x | 1x | | n/a | | n/a | 39.2x | 76.4x | 37.11x | 57.1x | 102.9x | 50.9x |
| power (W)[5] | 182 | 185 | 183 | 461 | n/a | n/a | n/a | 128 | 129 | 124 | 131 | 128 | 127 |
| energy (10^3 J)[5] | 33 | 269 | 1020 | | n/a | | n/a | 0.5 | 2.3 | 18.3 | 0.4 | 1.7 | 14.6 |
| efficiency ((MFlop/s)/W) | 9.8 | 4.9 | 9.8 | 76.5 | n/a | n/a | n/a | 551.4 | 527.1 | 538.9 | 785.0 | 715.6 | 721.3 |
| efficiency gains | 1x | 1x | 1x | | n/a | | n/a | 56.7x | 108.4x | 55.4x | 80.8x | 145.1x | 71.4x |

- improved system performance
  - 1.59 times faster than the best published GPU and FPGA results

| data size[3] | CPU[1] | | | GPU[2] | | | FPGA (static) | | | | FPGA (optimal) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | s | m | l | | | | [14] | s | m | l | s | m | l |
| execution time (t) | 181.7 | 1458.2 | 5574.8 | [10] | [11] | [12] | | 3.6 | 18.0 | 147.9 | 2.9 | 13.4 | 110.59 |
| overhead time (t)[4] | 0.03 | 0.01 | 0 | | | | 0[7] | 0.14 | 0.58 | 3.88 | 0.22 | 0.82 | 5.8 |
| throughput (GFlop/s)[5] | 1.8 | 0.9 | 1.8 | 36 | 51.2 | 64.5 | 35.8[7] | 70.6 | 68.0 | 66.8 | **102.8** | 91.6 | 91.6 |
| speed-up | **1x** | **1x** | **1x** | | n/a | | n/a | 39.2x | 76.4x | 37.11x | 57.1x | **102.9x** | 50.9x |
| power (W)[5] | 182 | 185 | 183 | 461 | n/a | n/a | n/a | 128 | 129 | 124 | 131 | 128 | 127 |
| energy ($10^3$ J)[5] | 33 | 269 | 1020 | | n/a | | n/a | 0.5 | 2.3 | 18.3 | 0.4 | 1.7 | 14.6 |
| efficiency ((MFlop/s)/W) | 9.8 | 4.9 | 9.8 | 76.5 | n/a | n/a | n/a | 551.4 | 527.1 | 538.9 | **785.0** | 715.6 | 721.3 |
| efficiency gains | **1x** | **1x** | **1x** | | n/a | | n/a | 56.7x | 108.4x | 55.4x | 80.8x | **145.1x** | 71.4x |

- improved system performance
  - 1.59 times faster than the best published GPU and FPGA results
  - 1.45 times faster than optimized static implementation

| data size[3] | CPU[1] | | | GPU[2] | | | | FPGA (static) | | | FPGA (optimal) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | s | m | l | | | | [14] | s | m | l | s | m | l |
| execution time (t) | 181.7 | 1458.2 | 5574.8 | [10] | [11] | [12] | | 3.6 | 18.0 | 147.9 | 2.9 | 13.4 | 110.59 |
| overhead time (t)[4] | 0.03 | 0.01 | 0 | | | | 0[7] | 0.14 | 0.58 | 3.88 | 0.22 | 0.82 | 5.8 |
| throughput (GFlop/s)[5] | 1.8 | 0.9 | 1.8 | 36 | 51.2 | 64.5 | 35.8[7] | 70.6 | 68.0 | 66.8 | **102.8** | 91.6 | 91.6 |
| speed-up | **1x** | **1x** | **1x** | | n/a | | n/a | 39.2x | 76.4x | 37.11x | 57.1x | **102.9x** | 50.9x |
| power (W)[5] | 182 | 185 | 183 | 461 | n/a | n/a | n/a | 128 | 129 | 124 | 131 | 128 | 127 |
| energy (10[3] J)[5] | 33 | 269 | 1020 | | n/a | | n/a | 0.5 | 2.3 | 18.3 | 0.4 | 1.7 | 14.6 |
| efficiency ((MFlop/s)/W) | 9.8 | 4.9 | 9.8 | 76.5 | n/a | n/a | n/a | 551.4 | 527.1 | 538.9 | **785.0** | 715.6 | 721.3 |
| efficiency gains | **1x** | **1x** | **1x** | | n/a | | n/a | 56.7x | 108.4x | 55.4x | 80.8x | **145.1x** | 71.4x |

- improved system performance
  - 1.59 times faster than the best published GPU and FPGA results
  - 1.45 times faster than optimized static implementation
  - 1.42 times more energy efficient than the static designs

# Future work

- generalise design methods
- partial reconfiguration
- run-time scheduling
- further applications

# Conclusion

- design method to exploit run-time reconfiguration
  - partitioner
  - analytical model
  - scheduler
- improved system performance
  - 1.59 times faster than the best published GPU and FPGA results
  - 1.45 times faster than optimized static implementation
  - 1.42 times more energy efficient than the static designs