Two-stage Thermal-Aware Scheduling of Task Graphs on 3D Multi-cores Exploiting Application and Architecture Characteristics

 $\begin{array}{ccc} {\bf Zuomin} \ {\bf Zhu}^1 & {\rm Vivek} \ {\rm Chaturvedi}^2 & {\rm Amit} \ {\rm Kumar} \ {\rm Singh}^3 & {\rm Wei} \\ & {\rm Zhang}^1 & {\rm Yingnan} \ {\rm Cui}^2 \end{array}$ 

<sup>1</sup>Hong Kong University of Science and Technology

<sup>2</sup>Nanyang Technological University

<sup>3</sup>University of Southampton

January 12, 2017



#### Introduction

- Motivation
- Previous work

### Our Approach

- Two-stage scheduling (TSS)
- Communication-aware group stage
- Thermal-aware scheduling stage

#### **Experimental Results** 3

- Experimental Setup
- Peak Temperature Reduction
- Performance Evaluation
- Extension of TAS to multi-layer 3D CMP



- Motivation
- Previous work

#### Our Approach

- Two-stage scheduling (TSS)
- Communication-aware group stage
- Thermal-aware scheduling stage

### 3 Experimental Results

- Experimental Setup
- Peak Temperature Reduction
- Performance Evaluation
- Extension of TAS to multi-layer 3D CMP

# Challenges faced by 3D CMP



**Benefits:** 

- Multiple active layers stacked.
- Inergy efficiency and scalability.
- Short interconnect delay.

#### **Challenges:**

- High power density.
- 2 Thermal emergency.
- Reliability threat, device aging.

#### How to address thermal challenges?

Figure 1: Two-layer 3D CMP

Architecture

Based on thermal-aware or power-aware scheduling approaches:

#### Targeting independent tasks

- Stack power balance[1]
- Incremental update of thermal simulation[2]

### Targeting task graphs

- Swap tasks between hot and cold stacks[3]
- Ø Move low power tasks to top layer [4]
- Schedule tasks in the order of priority, thermal simulations for mapping[5]

- Neglect influence of neighboring cores and leakage power on temperature[1]
- High computational complexity[2]
- Incur huge swapping overhead[4, 3]
- High thermal simulation overhead at run-time [5]

Thus we propose a novel decoupling scheduling algorithm for peak temperature minimization while optimizing the makespan of the application.



- Motivation
- Previous work

#### 2 Our Approach

- Two-stage scheduling (TSS)
- Communication-aware group stage
- Thermal-aware scheduling stage

#### 3 Experimental Results

- Experimental Setup
- Peak Temperature Reduction
- Performance Evaluation
- Extension of TAS to multi-layer 3D CMP

# Two-stage scheduling algorithm

Decouple the optimization of makespan and temperature into two stages and optimize them separately.

#### Communication-aware group (CAG) stage: Design-time stage

- Bind individual tasks to super tasks assuming different number of available cores.
- Brute-force or genetic algorithm
- Inimize the makespan considering the worst-case communication.

### Thermal-aware scheduling (TAS) stage: Run-time stage

- Super tasks mapped to real available cores.
- ② Different strategies to handle different layers.
- Some the second second
- For peak temperature minimization.



Figure 2: An example DG model: Autoindust2

- Search the best mappings utilizing different number of cores.
- Exhaustive exploration when task number is small.
- GA is applied when the task number is bigger.
- Consider the maximum hop (communication delay) between cores.

Super task: a group of tasks (nodes) which are mapped to a particular core in a given scheduling interval.

#Cores	C1	C2	C3	C4	C5	C6	Period(cycles)
6	е	d	С	b	f	а	19018400
5	*	е	d	С	bf	а	19004900
4	*	*	de	С	bf	а	18991400
3	*	*	*	cde	bf	а	18977900
2	*	*	*	*	bcdef	а	18960000
1	*	*	*	*	*	bacdef	21360000



Figure 3: Best mappings generated at design-time for Autoindust2

Figure 4: Two super tasks with the minimum makespan

# Run-time peak temperature optimization



Figure 5: Flow chart of TSS

- Select the best mapping depending on the number of available cores.
- Classify super tasks based on their power.
- Different core selection strategy in different layers.
- A combination of two heuristics:
  - Thermal rank model in bottom layer (the layer closest to heat sink).
  - Combined power model in top layer (the layer far away from heat sink).

# Thermal rank model in bottom layer

- Dissipating heat to ambient faster, consider thermal efficiency.
- Thermal rank determines the likelihood of a core to receive new tasks for execution.
- Smaller thermal rank, higher thermal efficiency.

• 
$$TR = (T + Lw \cdot T_l + Vw \cdot T_v) \cdot \rho \cdot P_f$$

T: temperature of the candidate core;  $T_l, T_v$ : temperature of the lateral and vertical

neighboring cores;

 $\rho$ : the proximity to the heat sink;

 $P_f$ : position factor, the influence of absolute position of a core;

 $L_w, V_w$ : lateral and vertical thermal conductance.



Figure 6: 3D CMP Architecture

# Combined power model in top layer

$$P_{combined} = (R_{0\_amb} \cdot P_0 + R_{1\_amb} \cdot P_1) \cdot P_f$$

- Hotspot and temperature emergency in this layer.
- Limiting the total power in a core stack.
- Consider the stacked power in the top layer.



Figure 7: An example of target 3D CMP

ASPDAC 2017

## Extension to Multi-layer 3D CMP

• Consider thermal efficiency and the stacked power for different layers;

• Hot super task to hyper-bottom layer, cool super task to hyper-top layer. For a three-layer 3D CMP and four-layer 3D CMP.



Figure 8: A three-layer 3D CMP

Figure 9: A four-layer 3D CMP

# Run-time Mapping Algorithm



Zuomin Zhu (HKUST)

ASPDAC 2017



- Motivation
- Previous work

#### Our Approach

- Two-stage scheduling (TSS)
- Communication-aware group stage
- Thermal-aware scheduling stage

#### 3 Experimental Results

- Experimental Setup
- Peak Temperature Reduction
- Performance Evaluation
- Extension of TAS to multi-layer 3D CMP

Steps for comparison:

- CF and TR: targeting independent tasks. Temperate comparison to validate the run-time TAS algorithm.
- PTLS: targeting task graphs. Both temperature and performance comparison to validate our TSS algorithm.
- Vary scheduling interval from 1ms to 2ms to explore the scenarios under different workload pressure.

Methods	Abbreviation	References
Coolest First	CF	[4]
Thermal Profiling	TR	[1]
Peak Temperature List Scheduling	PTLS	[5]
Two-stage scheduling	TSS	proposed

Table 1: Methodologies considered for comparison

# Validate the TAS algorithm

- TR in [1] and CF in [4] are considered for only temperature comparison, validating the TAS stage.
- Sor fair comparison, we assume that the design time makespan optimization results are available to the approaches targeting independent tasks.
- 3 Seven benchmarks at fixed scheduling interval of 1ms are mapped to  $4 \times 4 \times 2$  CMP.
- Reduce peak temperature up to 5.0°C compared to CF, and up to 4.4°C compared to TR.



Figure 11: Absolute peak temperature of individual benchmarks

# Varying the scheduling interval

- Varying the scheduling interval to explore the scenarios under different workload pressure.
- Average reduction of 3.6°C, maximum reduction of 6.1°C.



# Comparing to PTLS on Two Layers

- $\blacksquare$  7 benchmarks on a two-layer 3D CMP arranged in the 4 imes 4 imes 2 pattern.
- ② CAG stage optimizes the makespan, and TAS stage reduces peak temperature.
- An average of 4.9% performance improvement, and 6.3°C peak temperature reduction.



# Comparing to PTLS on Four Layers

- A four layer 3D CMP arranged in the  $2 \times 4 \times 4$  pattern, the same total number of cores as two-layer 3D CMP.
- An average of 6.78% performance improvement, and 10.0°C peak temperature reduction.
- **O** Greater improvement is achieved with the increasing number of layers.



# TAS VS. Thermal rank only and Combined power only

- Compare our TAS approach which incorporates the consideration of both thermal efficiency and power stack against the ones considering only one factor.
- Ocheck the peak temperature at fixed scheduling interval of 1ms on two-layer, three-layer, and four-layer 3D CMPs.
- Solution  $\mathbf{S}$  Every layer is comprised of  $4 \times 4$  identical cores.
- TAS outperforms them by 7.3°C and 5.7°C reduction, respectively.



Figure 15: Relative peak temperature on 3D CMPs with different layers



- Motivation
- Previous work

#### Our Approach

- Two-stage scheduling (TSS)
- Communication-aware group stage
- Thermal-aware scheduling stage

### 3 Experimental Results

- Experimental Setup
- Peak Temperature Reduction
- Performance Evaluation
- Extension of TAS to multi-layer 3D CMP

#### Summary

- A two-stage thermal-aware task scheduling algorithm.
- Two steps: the communication-aware group stage at design-time, and the thermal-aware scheduling step at run-time.
- Average 4.9% performance improvement and 6.3°C peak temperature reduction.

#### Future Works

Apply DVFS into the thermal-aware scheduling algorithm to further reduce the peak temperature.

S. Liu et al., "Thermal-aware job allocation and scheduling for three dimensional chip multiprocessor," in *Proc. of Int. Symp. on Quality Electronic Design*, pp. 390–398, 2010.

C. H. Yu et al., "Thermal-aware on-line scheduler for 3-D many-core processor throughput optimization," IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems, vol. 33, no. 5, pp. 763–773, 2014.

V. Chaturvedi *et al.*, "Thermal-aware task scheduling for peak temperature minimization under periodic constraint for 3D-MPSoCs," in *Proc. of Int. Symp. on Rapid System Prototyping*, pp. 107–113, 2014.

Y. Cui *et al.*, "Thermal-aware task scheduling for 3D network-on-chip: A Bottom-to-Top scheme," in *Proc. Int. Symp. on Integrated Circuits*, pp. 224–227, 2014.

J. Li *et al.*, "Thermal-aware task scheduling in 3D chip multiprocessor with real-time constrained workloads," *ACM Trans. on Embedded Computing Systems*, vol. 12, no. 2, p. 24, 2013.

(日) (同) (三) (三)