

# TAPHS: Thermal-Aware Unified Physical-Level and High-Level Synthesis

Zhenyu (Peter) Gu<sup>1</sup>, Yonghong Yang<sup>2</sup>, Jia Wang<sup>1</sup>

Robert Dick<sup>1</sup>, Li Shang<sup>2</sup>

Northwestern Univ.<sup>1</sup>, Queen's Univ.<sup>2</sup>

01/27/2006

# Outline

- Introduction & past work
- Motivating example
- System infrastructure overview
- Thermal-aware techniques
  - Architecture-level technique
  - Physical-level technique
- Experimental results
- Conclusion

# Outline

- Introduction & past work
- Motivating example
- System infrastructure overview
- Thermal-aware techniques
  - Architecture-level technique
  - Physical-level technique
- Experimental results
- Conclusion

# Introduction & past work

- Thermal issues
  - Cooling costs
  - Reliability
  - Package costs
  - Performance
- Thermal-aware design
  - Requires a unified high-level and physical-level design optimization
  - Incremental synthesis is promising

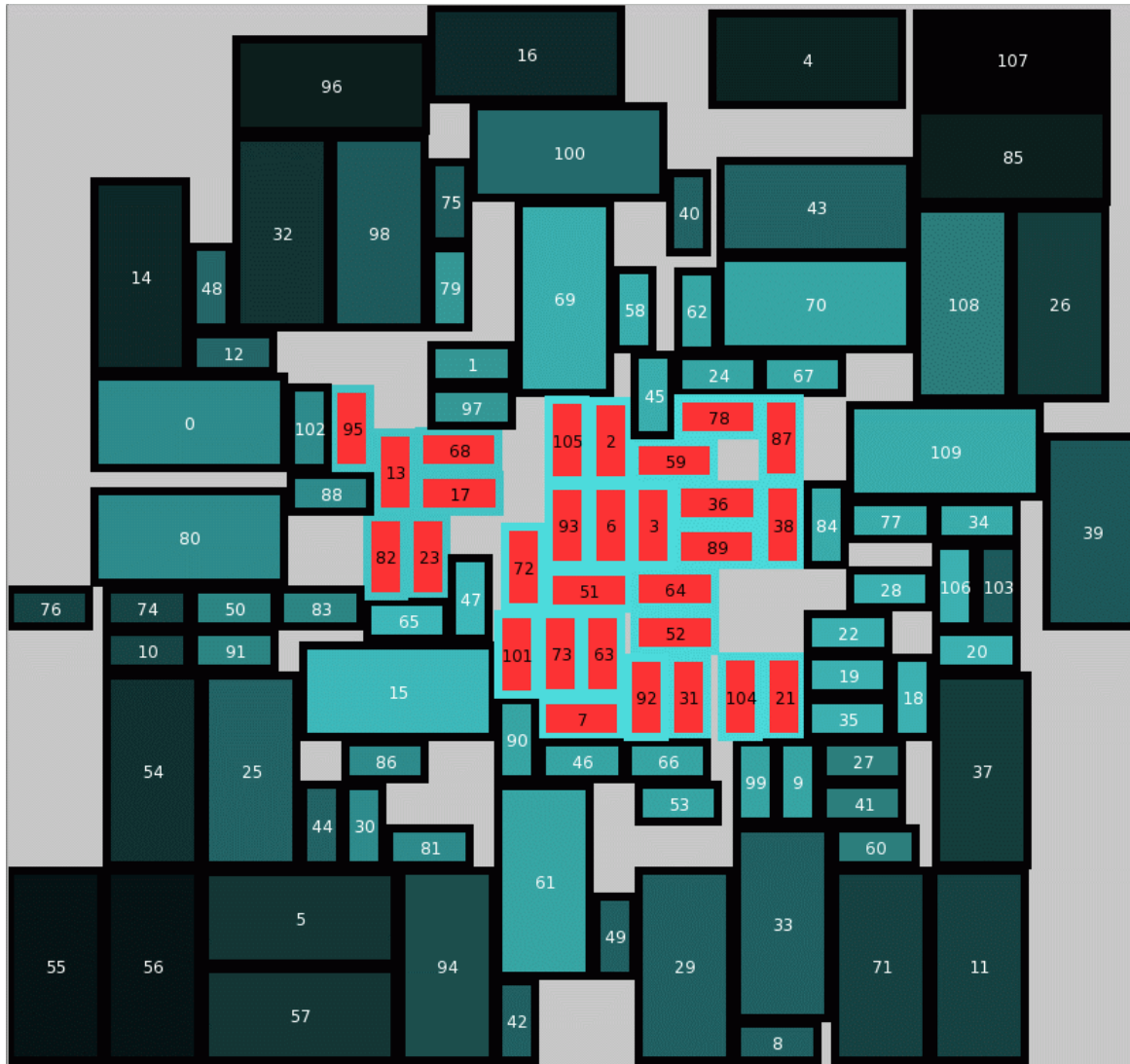
# Introduction & past work

- High-level and physical-level co-synthesis
  - J. P. Weng and A. C. Parker 93
  - D. Thomas 00
  - L. Zhong and N.K.Jha, 02
  - A. Stammermann, et al., 03
  - Z. P.Gu, et al., '05
- Thermal-aware analysis and design
  - K. Skadron, et al., 03
  - L. Shang, et al., 04
  - B. Goplen, et al., 03
  - J. Cong, et al., 04
  - R. Mukherjee, et al., 05

# Outline

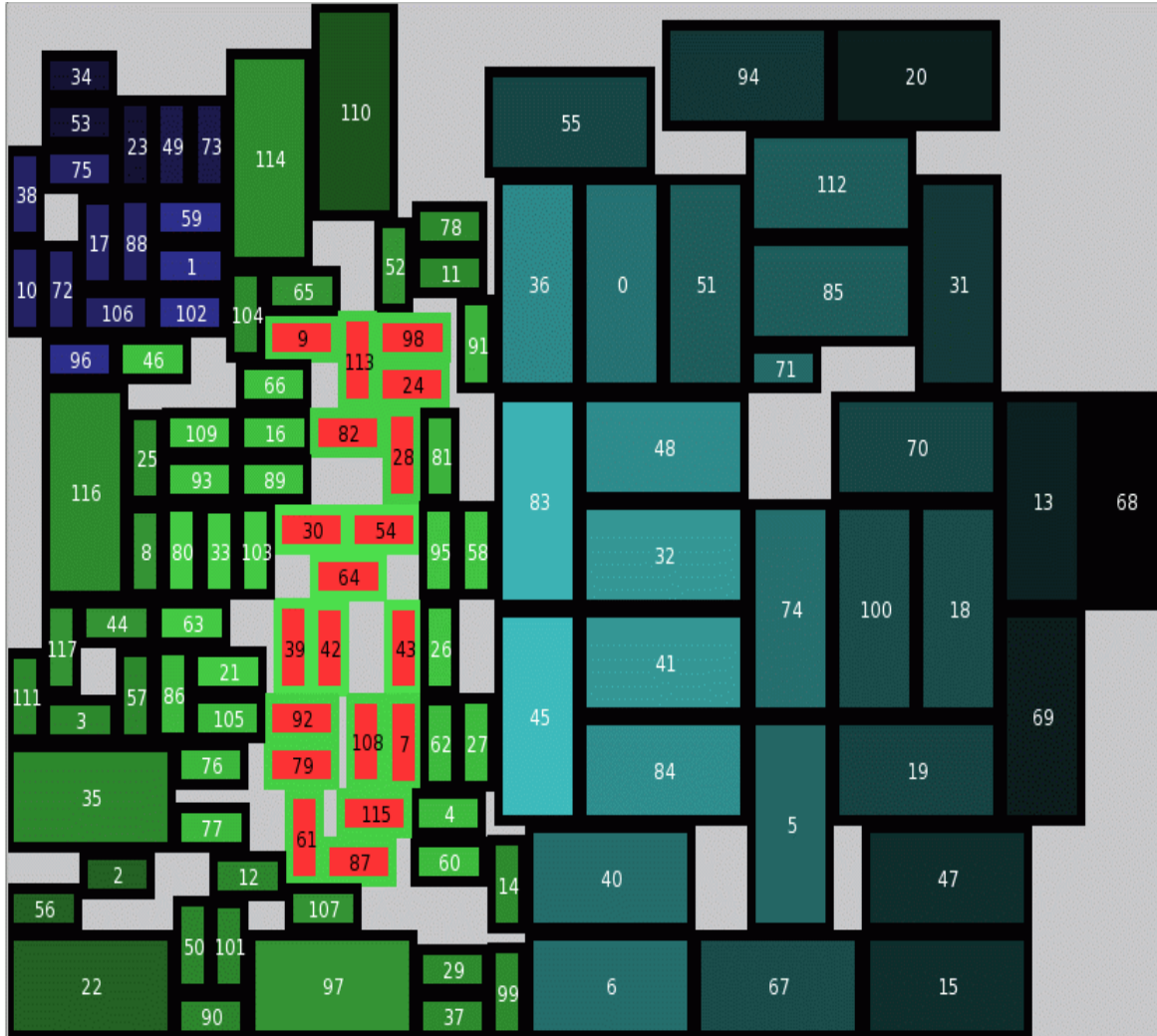
- Introduction & past work
- **Motivating example**
- System infrastructure overview
- Thermal-aware techniques
  - Architecture-level technique
  - Physical-level technique
- Experimental results
- Conclusion

# Motivating examples



**29 functional units have temperature higher than 85°C**

# Motivating examples



- **3 voltage islands**

- **19 functional units have temperature higher than 85°C**



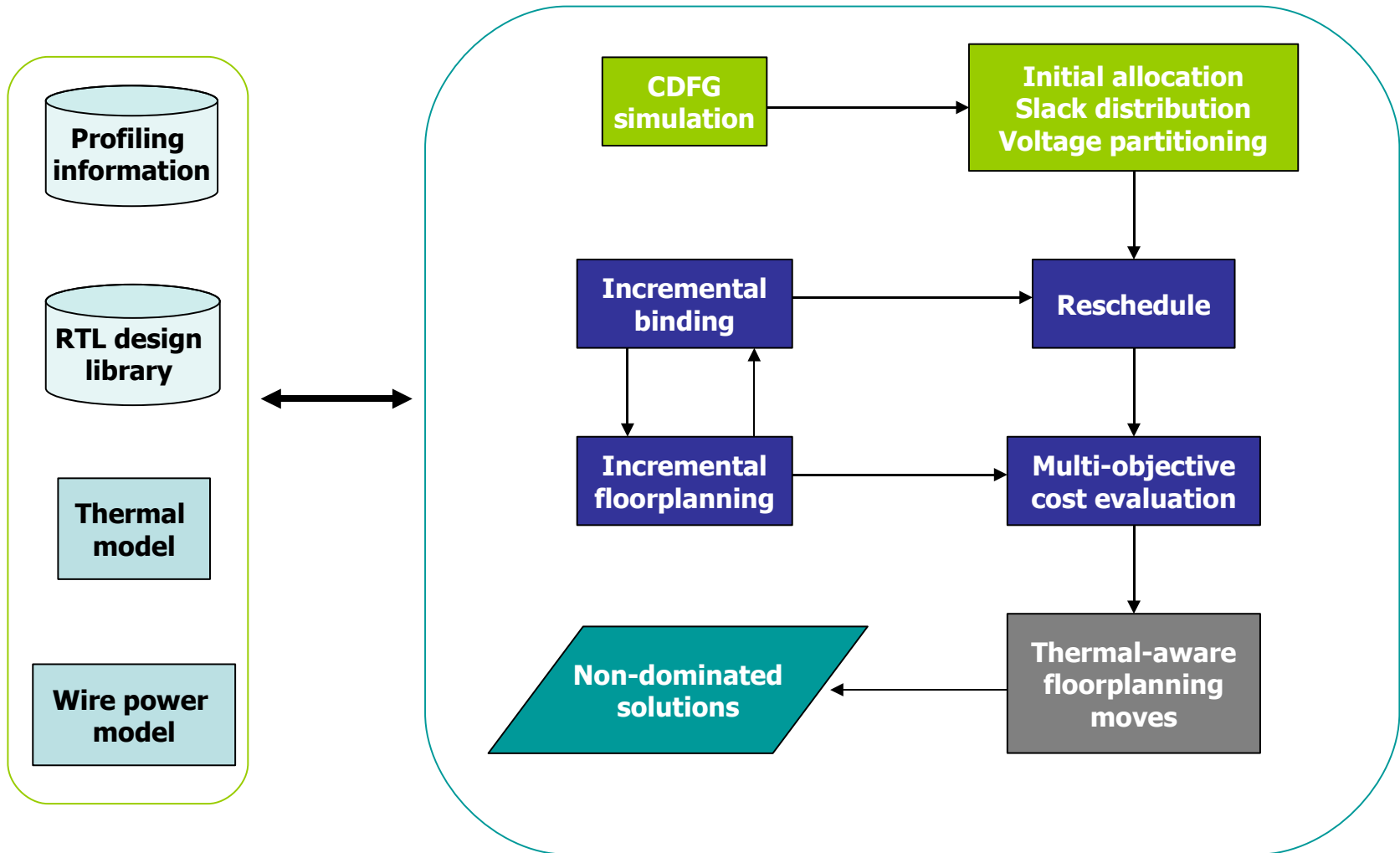
# Motivating examples

- Unified high-level and physical-level thermal optimization is necessary
- Voltage partitioning reduces power consumption
- Thermal-aware floorplanning further reduces peak temperature
- Incremental synthesis is essential to reduce synthesis time

# Outline

- Introduction & past work
- Motivating example
- System infrastructure overview
- Thermal-aware techniques
  - Architecture-level technique
  - Physical-level technique
- Experimental results
- Conclusion

# System infrastructure overview



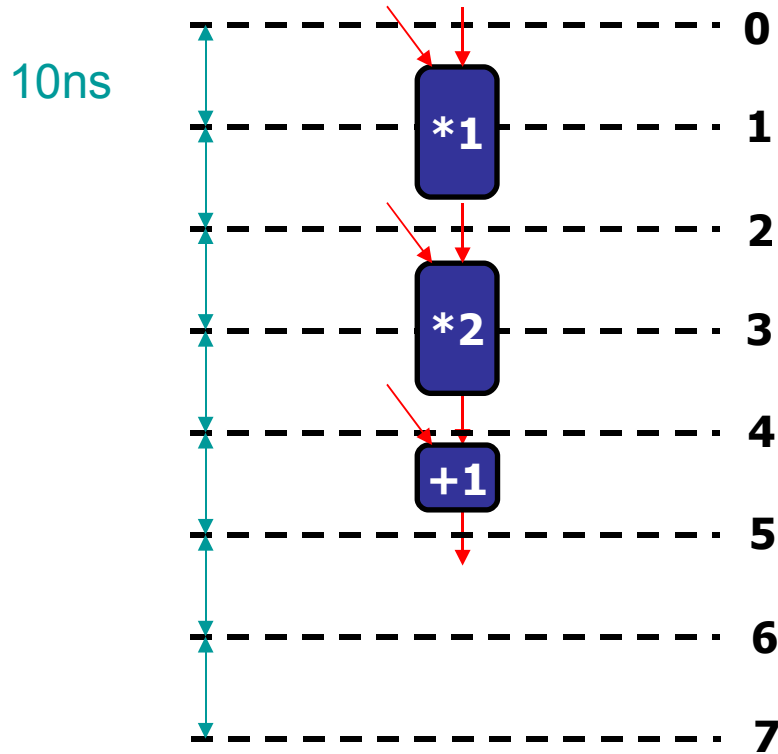
# Outline

- Introduction & past work
- Motivating example
- System infrastructure overview
- Thermal-aware techniques
  - Architecture-level technique
  - Physical-level technique
- Experimental results
- Conclusion

# Thermal-aware techniques

- Architecture-level
  - Voltage island
    - Reducing supply voltage increases circuit propagation delay
    - Overhead for voltage islands
  - Solution
    - Slack distribution
    - Voltage partitioning
- Physical-level
  - Voltage island generation
  - Thermal-aware swap operation

# Slack distribution



\*1:  $d=16\text{ns}$ ,  $c=4$

\*2:  $d=16\text{ns}$ ,  $c=2$

+1:  $d=8\text{ns}$ ,  $c=2$

$$e_i = C_i v_i^2 = C_i \left( \frac{d_i}{K_i} \right)^{\frac{2}{1-\alpha}}$$

for  $\alpha = 2$

$$E = \sum_{i \in p} C_i \left( \frac{d_i}{K_i} \right)^{\frac{2}{1-\alpha}}$$

for  $\alpha = 2$

Clock period = 10ns

Shared slack = 2

# Slack distribution

Slack: difference between latest and earliest start time.

$$\min_{\substack{\forall i \in p \\ v_i}} \sum_{i \in p} C_i \left( \frac{d_i}{K_i} \right)^{\frac{2}{1-\alpha}}$$

- D: bound on path execution time
- p: set of all operations on the path
- $d_i$ : delay of an operation's functional unit
- $v_i$ : voltage of an operation's functional unit
- $K_i$ : an execution time constant of an operation's functional unit
- $\alpha$ : alpha power law constant
- $C_i$ : switched capacitance of functional unit i

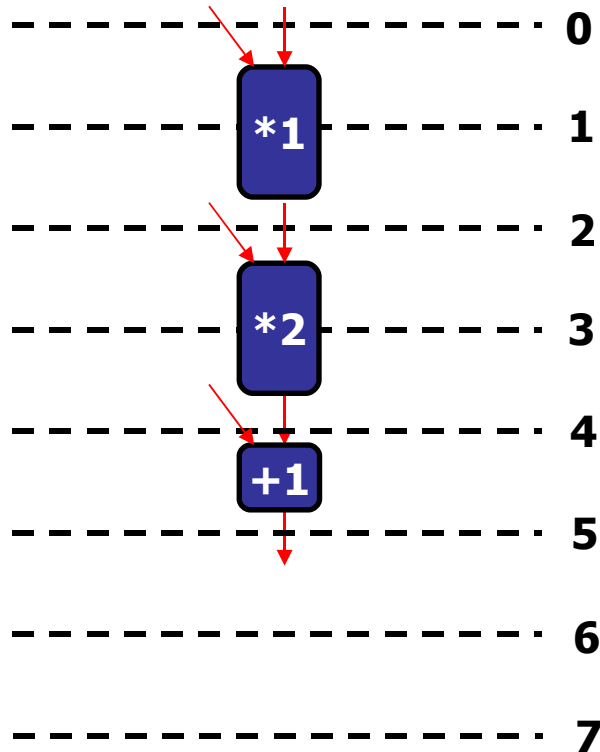
# Slack distribution

$$\forall_{i \in p} d_i = \frac{D}{N} \left( \frac{\frac{C_1}{K_1^{\frac{2}{1-\alpha}}}}{\frac{C_i}{K_i^{\frac{2}{1-\alpha}}}} \right)^{\frac{1-\alpha}{1+\alpha}} \quad \text{or} \quad \frac{D}{N} \sqrt[3]{\frac{C_i K_i^2}{C_1 K_1^2}} \quad \text{for } \alpha = 2$$

- $C_i$ : switched capacitance of functional unit  $i$
- $K_i$ : an execution time constant
- $D$ : bound on path execution time
- $p$ : set of all operations on the path
- $d_i$ : delay of an operation's functional unit
- $\alpha$ : alpha power law constant
- $N$ : sum of the optimal delay ratio



# Slack distribution



\*1:  $d=16\text{ns}$ ,  $c=4$

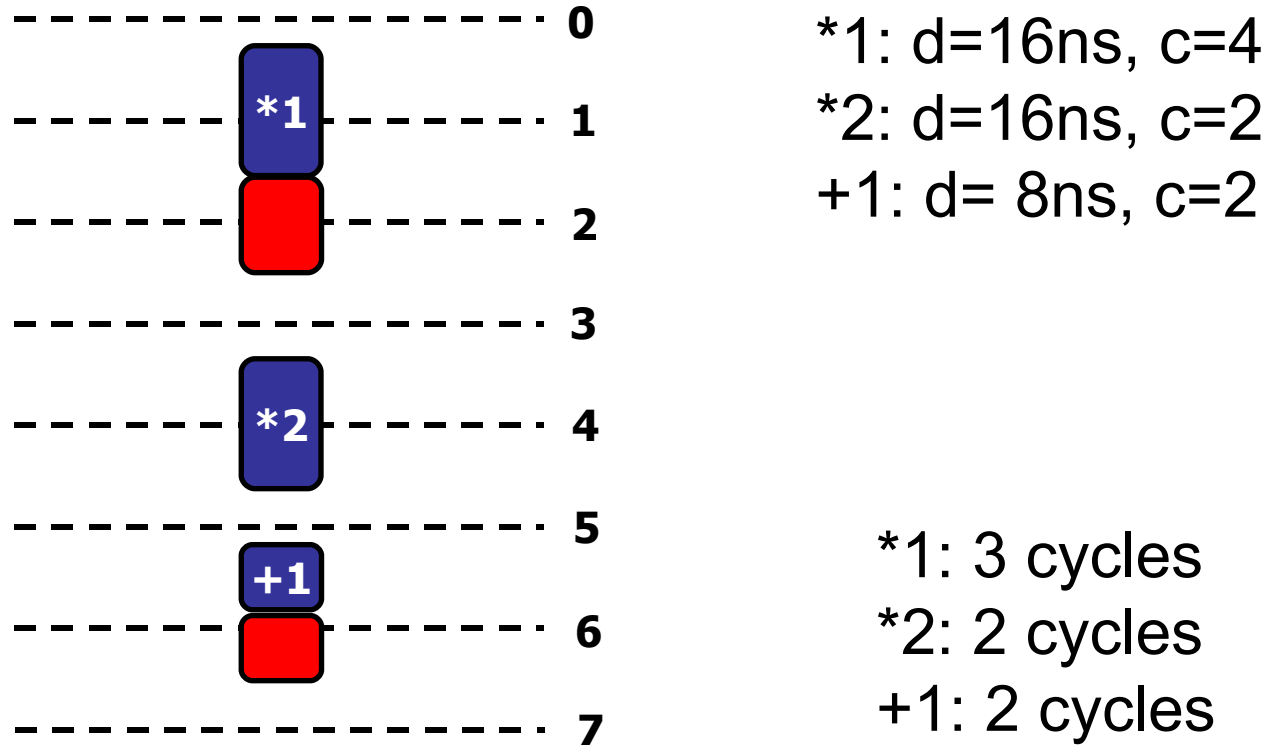
\*2:  $d=16\text{ns}$ ,  $c=2$

+1:  $d=8\text{ns}$ ,  $c=2$

Clock period = 10ns

Shared slack = 2

# Slack distribution



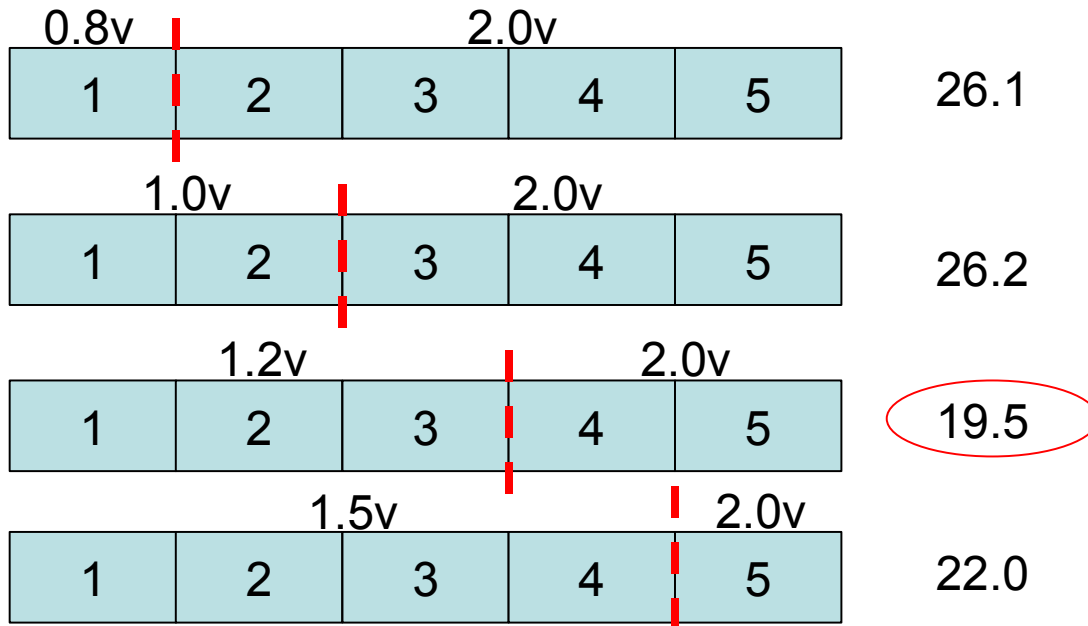
Clock period = 10ns  
Shared slack = 2

# Multiple voltage techniques

## Motivating example

FUs	FU1	FU2	FU3	FU4	FU5
Slack vs. delay	1.1	0.7	0.5	0.3	0
$V_{\min}$ (V)	0.8	1	1.2	1.5	2
C(pf)	2.0	0.2	3.0	1.0	2.0

# Multiple voltage techniques



Decreasing order of the slack

# Outline

- Introduction & past work
- Motivating example
- System infrastructure overview
- Thermal-aware techniques
  - Architecture-level technique
  - Physical-level technique
- Experimental results
- Conclusion

# Thermal-aware floorplanning

- Floorplan Representation
  - Adjacent Constraint Graph (ACG)  
(Zhou & Wang 'ICCD04)
    - A constraint graph with exactly one geometric relationship between every pair of modules
    - Operations have straightforward, local meaning in physical space
- Algorithm
  - Use simulated annealing for initial floorplan
  - Greedy iterative improvement for re-optimization
- Cost function
  - Area
  - Pair-wise edge within same voltage island
  - Wire power consumption

# Thermal-aware floorplanning

## Thermal-aware swap operation

- Motivation: minimizing average power consumption and minimizing peak temperature conflict with each other
- Solution: thermal-aware swap operation that exchanges hot (high-power density) functional units with cool (low-power density) functional units within the same voltage island

# Outline

- Introduction & past work
- Motivating example
- System infrastructure overview
- Thermal-aware techniques
  - Architecture-level technique
  - Physical-level technique
- Experimental results
- Conclusion



# Experimental results

## Thermal model

- Chip-package thermal model (Shang MICRO'03)
- Compared with COMSOL Multiphysics (FEMLAB): less than 2.5% estimation error on the Kelvin scale
- Two thermally conductive paths:
  - From the silicon die through the cooling package to the ambient environment
  - From the silicon die through the package to the printed circuit board
- Silicon thickness 200um
- Ambient temperature 45°C

# Experimental results

## **Benchmarks: 13 benchmarks**

- Technology: TSMC 0.18um
- Jacobi, the largest widely used benchmark:  
24 MULs, 8 DIVs, 8 ADDs, 16 SUBs
- 2 large random benchmarks generated by TGFF (Dick'98)
  - Random100: 20 ADDs, 15 SUBs, 19 MULs
  - Random200: 39 ADDs, 44 SUBs, 36 MULs

# Experimental results

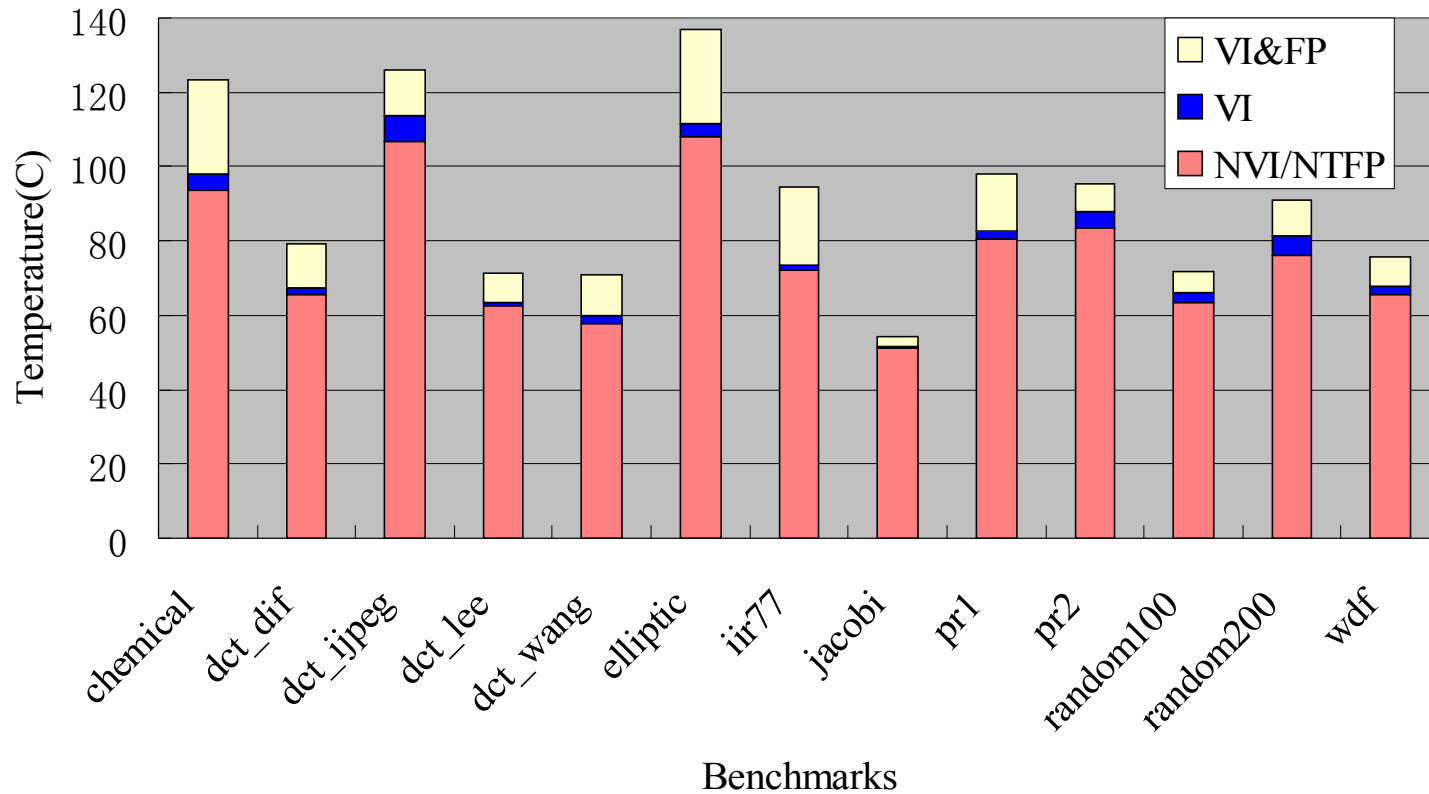
Example	No voltage islands			Voltage islands			Thermal FP	
	Peak T (°C)	Area (%)	Power (W)	Peak T (°C)	Area (%)	Power (W)	Peak T (°C)	Power (W)
dct_dif	79.0	87.9	0.85	67.3	92.5	0.60	65.6	0.55
	79.7	78.6	0.83	67.6	81.5	0.58	66.1	0.54
	80.3	83.7	0.85	69.8	83.4	0.61	67.4	0.57
	80.1	81.4	0.84	69.3	74.9	0.57	67.6	0.53
	81.7	80.7	0.86	69.9	80.0	0.60	68.4	0.56
	82.9	76.0	0.87	71.3	78.8	0.63	68.5	0.57
	84.5	68.8	0.87	71.4	75.8	0.62	68.7	0.57

# Experimental results

Example	No voltage islands			Voltage islands			Thermal FP	
	Peak T (°C)	Area (%)	Power (W)	Peak T (°C)	Area (%)	Power (W)	Peak T (°C)	Power (W)
dct_wang	70.7	101.3	0.70	59.8	109.8	0.42	57.6	0.39
	68.2	97.5	0.68	59.1	116.0	0.43	57.9	0.40
	68.5	108.1	0.68	60.1	108.0	0.42	58.1	0.39
	70.4	89.1	0.70	59.8	102.8	0.44	58.3	0.41
	71.3	100.5	0.69	61.1	100.1	0.45	59.3	0.42
	70.3	101.0	0.70	61.2	113.0	0.45	59.6	0.42
	72.0	85.1	0.72	63.1	109.8	0.48	60.7	0.43
	72.4	77.4	0.70	65.2	91.8	0.47	61.4	0.42
	72.0	88.9	0.72	66.3	90.8	0.48	63.6	0.43
	70.8	86.6	0.70	66.7	78.2	0.47	64.5	0.43

# Experimental results

Peak temperature comparison



# Outline

- Introduction & past work
- Motivating example
- System infrastructure overview
- Thermal-aware techniques
  - Architecture-level technique
  - Physical-level technique
- Experimental results
- Conclusion

# Conclusions

- Presented a thermal-aware high-level synthesis system, which supports tight integration with thermal model and physical design
- Experimental results indicate that TAPHS is able to trade off peak temperature, IC area, and power consumption
- Thermal-aware design needs tight interaction between high-level and physical-level synthesis
- Incremental algorithm can save synthesis time by reusing and building upon high-quality previous physical design solutions that required a huge amount of time and effort to produce

Q & A?



# Thermal-aware floorplanning

$$\underbrace{n\sqrt{A}}_{\text{Area}} + 2n \sum_{v \in V} L_v + \sum_{e \in E} C_e D_e$$

Voltage island

Wire power consumption

- A: area
- n: number of functional units
- v: a pair of functional units sharing the same voltage
- $L_v$ : separation between a pair of functional units sharing the same voltage
- e: an interconnect line
- $C_e$ : unit-length switched capacitance for the data transfer along e
- $D_e$ : length of interconnect e