

Minimizing Thermal Gradient and Pumping Power in 3D IC Liquid Cooling Network Design Gengjie Chen, Jian Kuang, Zhiliang Zeng, Hang Zhang, Evangeline F. Y. Young, Bei Yu Department of Computer Science and Engineering, The Chinese University of Hong Kong

Introduction

Why 3D IC Liquid Cooling?

- Power is the number one problem in chip design
- **3D IC** is promising for increasing computer performance
- But 3D IC worsens power problem by
- higher heat dissipation density
- Iarger thermal resistance from junction to ambient
- Microchannel-based liquid cooling is proposed as a solution



WD46.0mm 10.0kV x150 200





Problem Formulations

Decision variables

- Cooling network topology N
- System pressure drop P_{svs}

Metrics

- Pumping power $W_{pump} = \frac{P_{sys} \cdot Q_{sys}}{m}$
- Q_{SVS} : system flow rate; η : efficiency term
- ► Thermal gradient $\Delta T = \max_i (\Delta T_i)$
- $\land \Delta T_i$: range of node temperatures in *i*-th source layer
- **Peak temperature** T_{max}

Design Rules

- TSV positions are at alternating basic cells in both dimensions
- Inlets and outlets can only occur at edges of channel layer
- At most one "continuous" inlet and outlet on each side
- **Problem 1: Pumping Power Minimization**

min W_{pump} , s.t. $P_{sys} \in \mathbb{R}^+$, $N \in \mathcal{N}$, $T_{max} \leq T^*_{max}$, $\Delta T \leq \Delta T^*$.

 $(\mathcal{N} \text{ is the set of all legal cooling networks})$ **Problem 2: Thermal Gradient Minimization**

Pumping Power Minimization

Tree-like Cooling Network

Hierarchical tree-like structure is simple and can balance cooling:

- Between upstream and downstream (factor 1)
- Among different trees (factor 2)

First

Network Topology Optimization

(1)

(2)

- In stage 1, ΔT under a **fixed** P_{sys} is used as cost function to accelerate
- In earlier stages, more rounds are performed to fully explore solution space
- Eight types of global flow directions are attempted

| Stage # | Step Size | Objective Function | n Simulator | Runtime for an Iteration |
|---------|-----------|---------------------------|-------------|--------------------------|
| 1 | 10 | ΔT | 2RM | short |

Challenges for 3D IC Liquid Cooling

- Hot downstream and cool upstream \implies
- large thermal gradient \implies
- reliability and timing issues
- \blacktriangleright limited channel diameter \Longrightarrow
- high pumping requirement \implies
- overhead to whole system
- Limitations of previous work
- No considering thermal gradient
- Assuming unidirectional straight channels
- Assuming unrealistic constant-temperature heat source

Thermal Modeling

- Most existing models assume unidirectional straight channels
- 4-register model (4RM) in 3D-ICE [Sridhar+, TOC'14]
- Accurate
- Has been extended for flexible topology
- ► Slow
- We construct a fast 2-register model (2RM) for cooling network Basics
- Divide channel layer into basic cells with a 2D grid Either solid (white/black, black reserved for TSV) or liquid (blue) Solve local pressure P_i and flow rate $Q_{i,j}$ from a **linear system** $P_{i,j} = g_{fluid,i,j} \cdot (P_i - P_j) (g_{fluid,i,j}: fluid conductance)$ \blacktriangleright $\sum_{i \in N_i} Q_{i,j} = 0$ (N_i : neighboring cells, inlet/outlet)

min ΔT s.t. $P_{sys} \in \mathbb{R}^+$, $N \in \mathcal{N}$, $T_{max} \leq T^*_{max}$, $W_{pump} \leq W^*_{pump}$.

General considerations

- ΔT is most difficult to handle among all metrics (W_{pump} , ΔT and T_{max}) \blacktriangleright W_{pump} vs. T_{max} is a simple trade-off under a specific **N**
- Liquid cooling alleviates T_{max} and worsens ΔT
- Three inducing factors for ΔT
- Temperature rise of coolant
- 2. Non-uniform power source distribution
- 3. Non-uniform channel distribution
- Factor 3 can be used to compensate for factors 1 & 2

Pumping Power Minimization

The problem is divided into two levels:

- ▶ Inner: P_{sys} is varied to minimize W_{pump} for a specific N, which evaluates N
- Outer: simulated annealing (SA) searches for a good N

Overall Flow of Pumping Power Minimization

Input: N_{init} , ΔT^* , T^*_{max} , stack description and floorplan files. **Output:** N, P_{SVS} . 1: $\mathbf{N} \leftarrow \mathbf{N}_{init}$; 2: while #iteration is within the limit do

- Obtain neighboring network solution N';
- $W'_{pump} \leftarrow \text{EVALUATENETWORK} (N', \Delta T^*, T^*_{max});$
- $N \leftarrow N'$ or not according to SA mechanism;
- if W'_{pump} converges then return N and P_{sys} ; : end while



Thermal Gradient Minimization

Similar to solving pumping power minimization with some optimization Network Evaluation

Its simplified form becomes:

min $f(P_{sys})$, s.t. $P_{sys} \in \mathbb{R}^+$, $P_{sys} \leq P_{sys}^*$.

(4)

- Solving (4) is simpler:
- If P_{svs}^* locates on falling side of f, it is optimal already
- Otherwise, adopt golden section search

Network Topology Optimization

Minimizing W_{pump} under a fixed P_{sys} is unrelated to temperature and meaningless, but minimizing ΔT under a fixed P_{sys} is safe \implies speed-up

- Some iterations are evaluated by one simulation under a fixed P_{svs}
- The original stage 1 is no longer needed
- Another stage with 4RM is affordable to replace the original stage 3

| Stage # | Step Size | Objective Functic | on Simulator Ru | untime for an Iteration |
|---------|-----------|--------------------------|-----------------|-------------------------|
| 1 | 10 | $\Delta T'$ | 2RM | short |
| 2 | 10 | $\Delta T'$ | 4RM | medium |
| 3 | 2 | $\Delta T'$ | 4RM | medium |





4RM Model

- ► Thermal cell = basic cell
- Solve temperature from a linear system considering three kinds of heat transfer

Solid-solid thermal conductance $g_{ss} = \frac{q_{i,j}}{T_i - T_i} = \frac{k_{solid} \cdot A_{i,j}}{I_{i,j}}$ Solid-liquid thermal conductance $g_{sl} = \frac{q_{i,j}}{T_i - T_i} = g_{sl}^* \parallel g_{ss}^* = \frac{g_{sl}^* \cdot g_{ss}^*}{g_{sl}^* + g_{ss}^*}$ with $g_{sl}^* = h_{conv} A_{i,j}$

Liquid-liquid heat transfer $q_{II} = C_V \cdot \sum_{j \in N_i} (Q_{j,i} \cdot T_{j,i}^*) = \frac{C_V}{2} \cdot \sum_{j \in N_i} (Q_{j,i} \cdot T_j)$



Faster 2RM Model

No conforming channel geometry \implies larger and fewer thermal cells \implies speed-up

Temperature vs. Pressure

► As P_{sys} increases, T_{max} decreases and finally becomes approximately constant

 $\land \Delta T = f(P_{svs})$ is either uni-modal or monotonically decreasing



Network Evaluation

- ► Replace W_{pump} by P_{sys} , as W_{pump} vs. P_{sys} is monotonic for a specific **N**
- ▶ Ignore T_{max} first, as it is easier to handle
 - Step 1: solve the problem without constraint T_{max}^*
 - Step 2: check T_{max} and find optimal solution by binary search

Network Evaluation of Pumping Power Minimization



Experimental Results

Faster 2RM Model

5 benchmarks, 40 network samples, 6 thermal cell sizes and 13 pressures

Free-like networks, $400 \mu m$ thermal cells: 0.52% errors (compared to 4RM), runtime reduced from 3.37s to 0.07s



Pumping Power Minimization

- ▶ 40 min for cases 1-3 and 240 min for case 4
- 79.61% better than baseline (unidirectional straight channels)
- 16.35% better than 1st place in ICCAD 2015 Contest





Important due to frequent simulation



Thermal nodes

- In solid layers, $m \times m$ basic cells = a thermal node
- In channel layers, $m \times m$ basic cells = a solid thermal node + a liquid one Heat transfer
- Solid-solid: only consider complete conducting paths
- Solid-liquid: project horizontal heat transfer to vertical direction
- Liquid-liquid: sum heat transfer over multiple channel connections

In step 1, by further substituting $\Delta T = f(P_{svs})$, Problem 1 becomes single-variable:

s.t. $P_{sys} \in \mathbb{R}^+$, $f(P_{sys}) \leq \Delta T^*$.

- Solve (3) by searching (with three probing points):
- ► If a feasible P_{sys} exists, return optimal P_{sys}
- Otherwise, return the P_{svs} for minimum f (show the nonexistence of feasible P_{sys},



Gengije Chen – CSE Department – The Chinese University of Hong Kong – Hong Kong

