

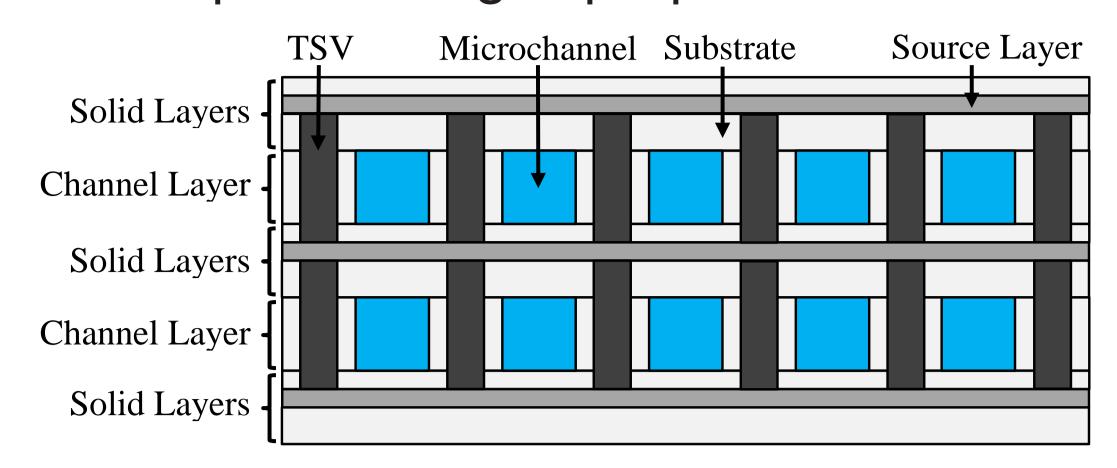


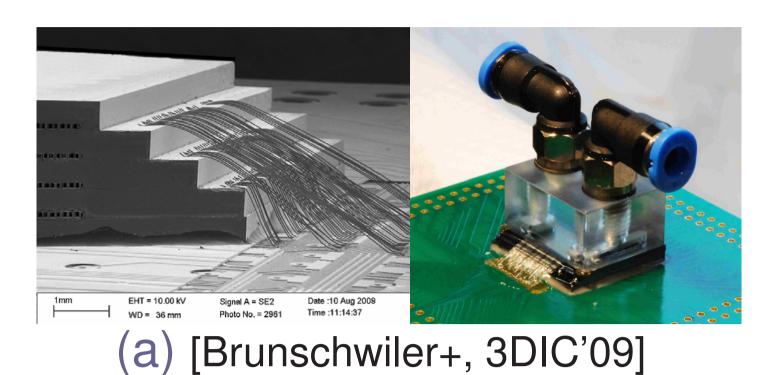
Gengjie Chen, Jian Kuang, Zhiliang Zeng, Hang Zhang, Evangeline F. Y. Young, Bei Yu Department of Computer Science and Engineering, The Chinese University of Hong Kong

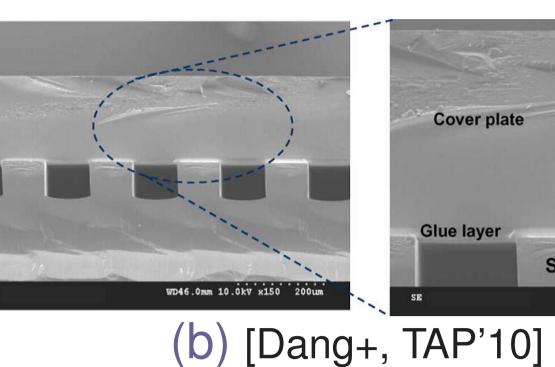
#### Introduction

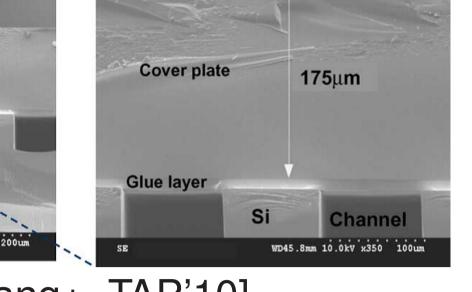
#### Why 3D IC Liquid Cooling?

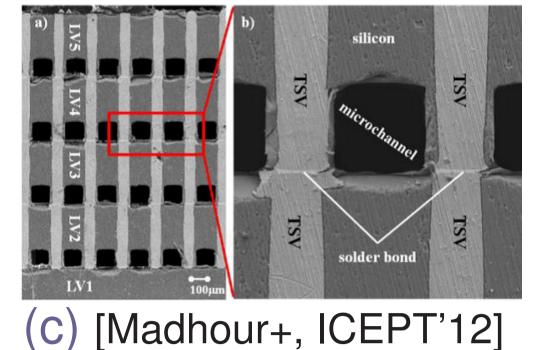
- Power is the number one problem in chip design
- ▶ 3D IC is promising for increasing computer performance
- But 3D IC worsens power problem by
- higher heat dissipation density
- larger thermal resistance from junction to ambient
- Microchannel-based liquid cooling is proposed as a solution











Challenges for 3D IC Liquid Cooling

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

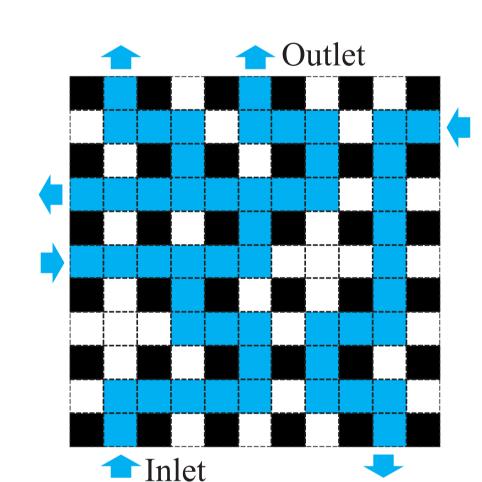
# 

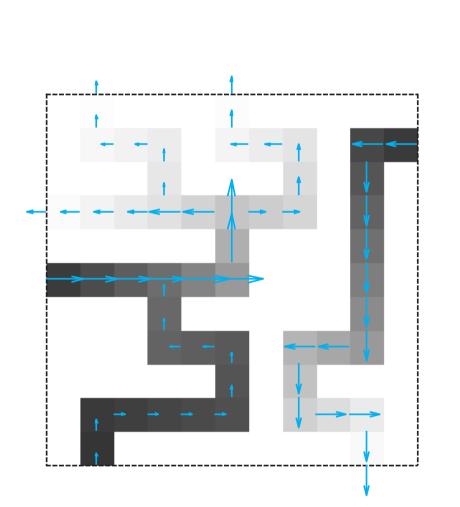
#### Thermal Modeling

- 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

- 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**
- $extstyle Q_{i,j} = g_{fluid,i,j} \cdot (P_i P_j) (g_{fluid,i,j}: fluid conductance)$  $\sum_{i \in N_i} Q_{i,j} = 0$  ( $N_i$ : neighboring cells, inlet/outlet)

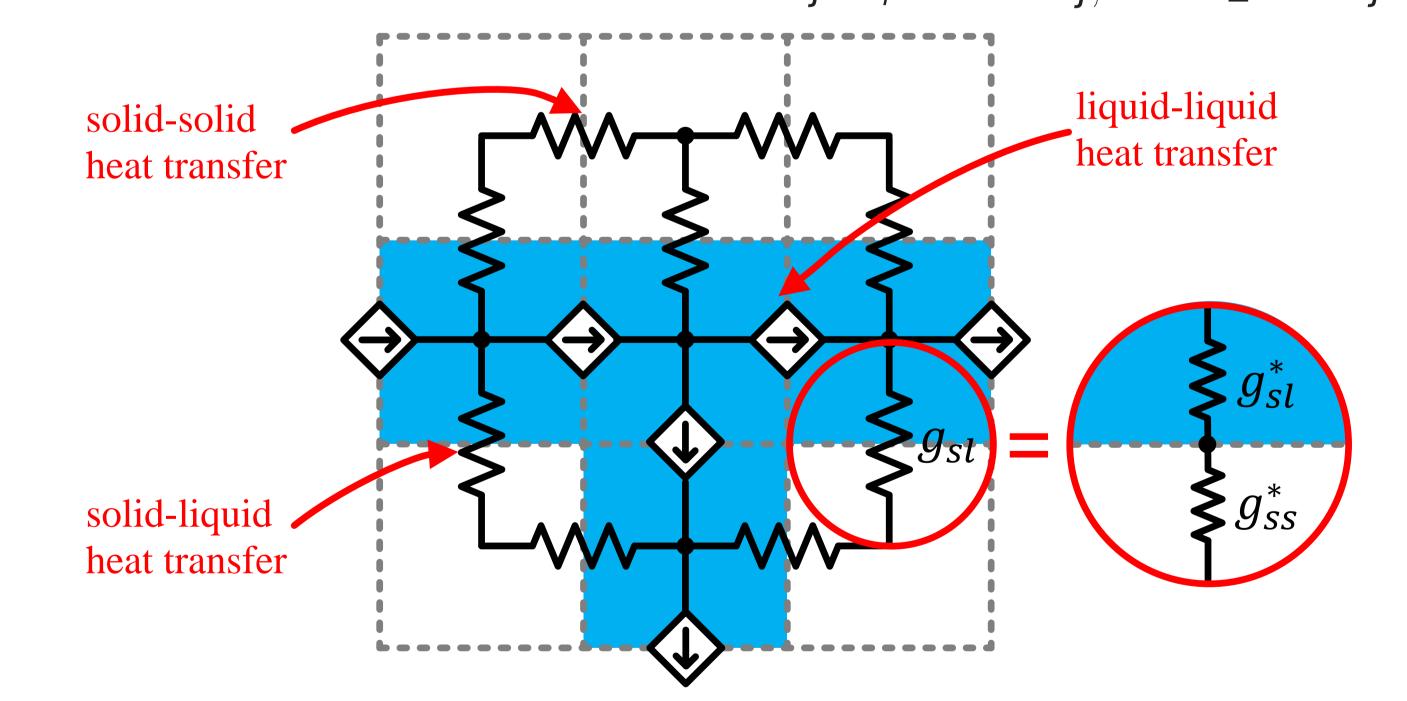




## 4RM Model

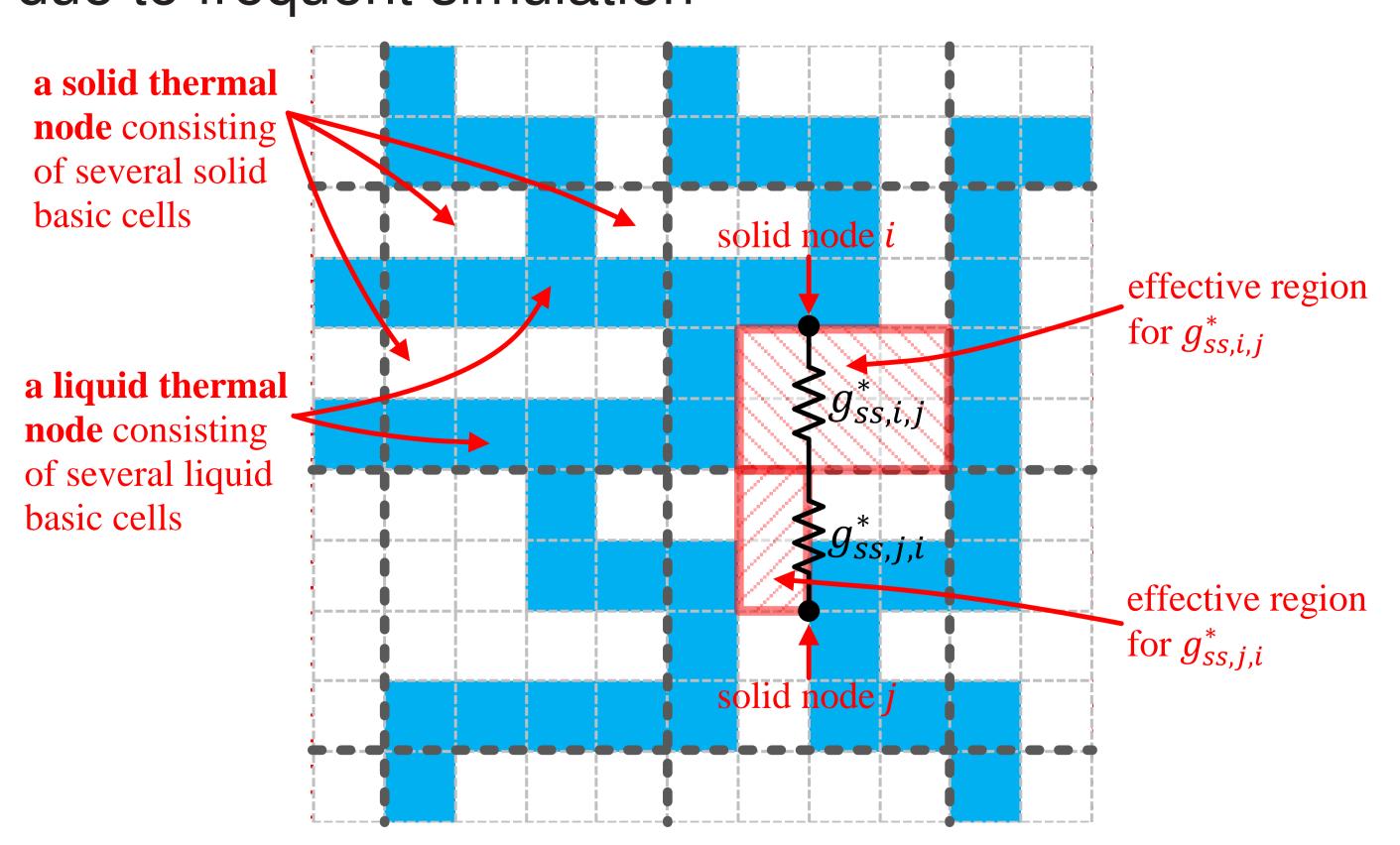
► Thermal cell = basic cell

- Solve temperature from a linear system considering three kinds of heat
- Solid-solid thermal conductance  $g_{ss} = \frac{q_{i,j}}{T_i T_i} = \frac{k_{solid} \cdot A_{i,j}}{l_{i,j}}$
- Solid-liquid thermal conductance  $g_{sl}=rac{q_{i,j}}{T_i-T_i}=g_{sl}^*\parallel g_{ss}^*=rac{g_{sl}^*\cdot g_{ss}^*}{g_{sl}^*+g_{ss}^*}$  with  $g_{sl}^* = h_{conv}A_{i,j}$
- Liquid-liquid heat transfer  $q_{ll} = C_V \cdot \sum_{j \in N_i} (Q_{j,i} \cdot T_{j,i}^*) = \frac{C_V}{2} \cdot \sum_{j \in N_i} (Q_{j,i} \cdot T_j)$



## Faster 2RM Model

- ightharpoonup No conforming channel geometry  $\Longrightarrow$  larger and fewer thermal cells  $\Longrightarrow$ speed-up
- 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

#### **Problem Formulations**

#### Decision variables

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

#### Metrics

- Pumping power  $W_{pump} = \frac{P_{sys} \cdot Q_{sys}}{r}$
- $ightharpoonup Q_{SVS}$ : system flow rate;  $\eta$ : efficiency term
- Thermal gradient  $\Delta T = \max_i (\Delta T_i)$
- $\triangle T_i$ : 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**

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

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

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

General considerations

 $\Delta T$  is most difficult to handle among all metrics ( $W_{pump}$ ,  $\Delta T$  and  $T_{max}$ )  $ightharpoonup W_{pump}$  vs.  $T_{max}$  is a simple trade-off under a specific N

- Liquid cooling alleviates  $T_{max}$  and worsens  $\Delta T$
- Three inducing factors for  $\Delta T$
- . 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

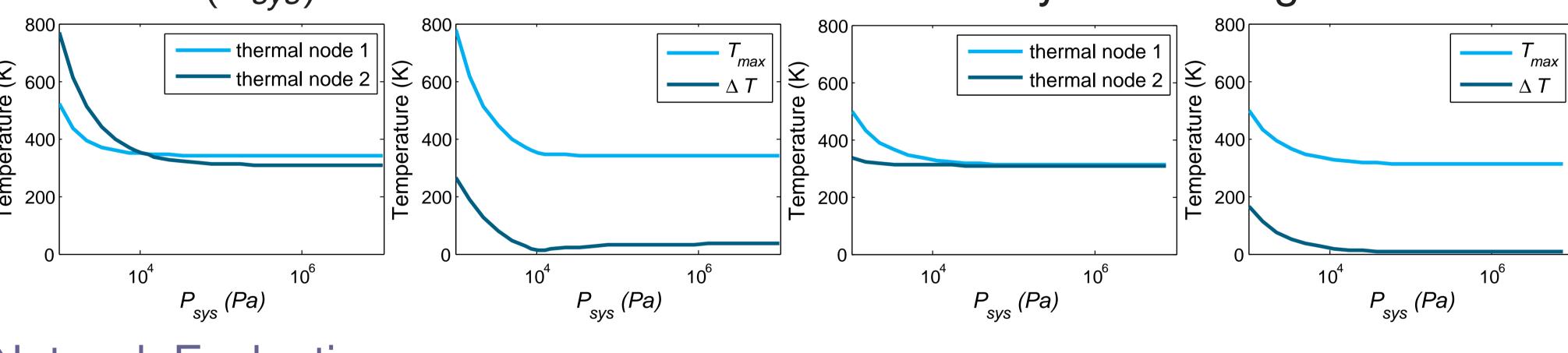
#### **Overall Flow of Pumping Power Minimization**

# **Input:** $N_{init}$ , $\Delta T^*$ , $T^*_{max}$ , stack description and floorplan files.

- Output:  $N, P_{sys}$ . 1:  $N \leftarrow N_{init}$ ; 2: while #iteration is within the limit do
  - Obtain neighboring network solution N';  $W'_{numn} \leftarrow \text{EVALUATENETWORK}(N', \Delta T^*, T^*_{max});$
- $N \leftarrow N'$  or not according to SA mechanism; if  $W'_{pump}$  converges then return N and  $P_{sys}$ ;
- : end while

Temperature vs. Pressure

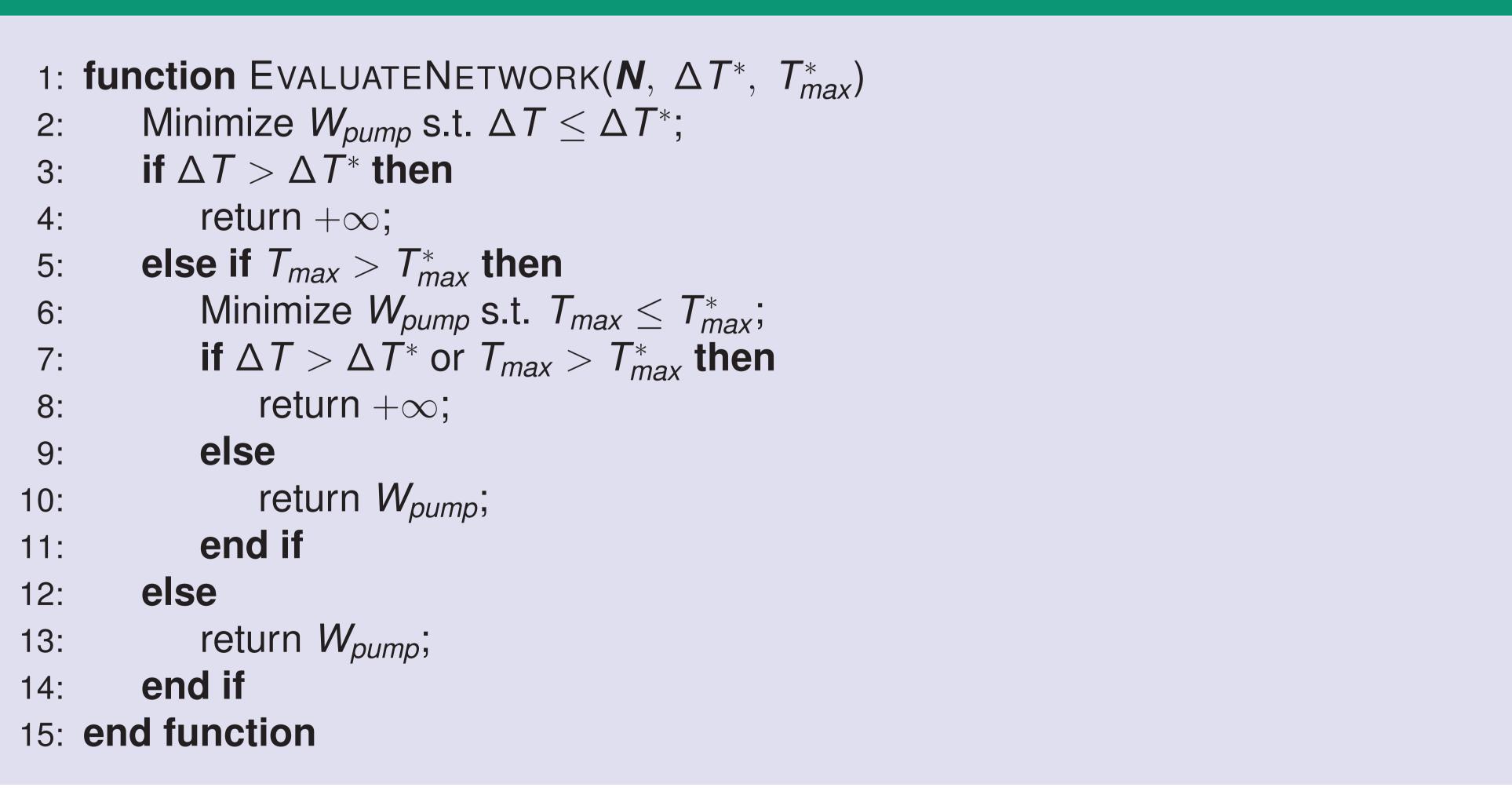
- As  $P_{sys}$  increases,  $T_{max}$  decreases and finally becomes approximately constant
- $ightharpoonup \Delta T = f(P_{svs})$  is either uni-modal or monotonically decreasing



## **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

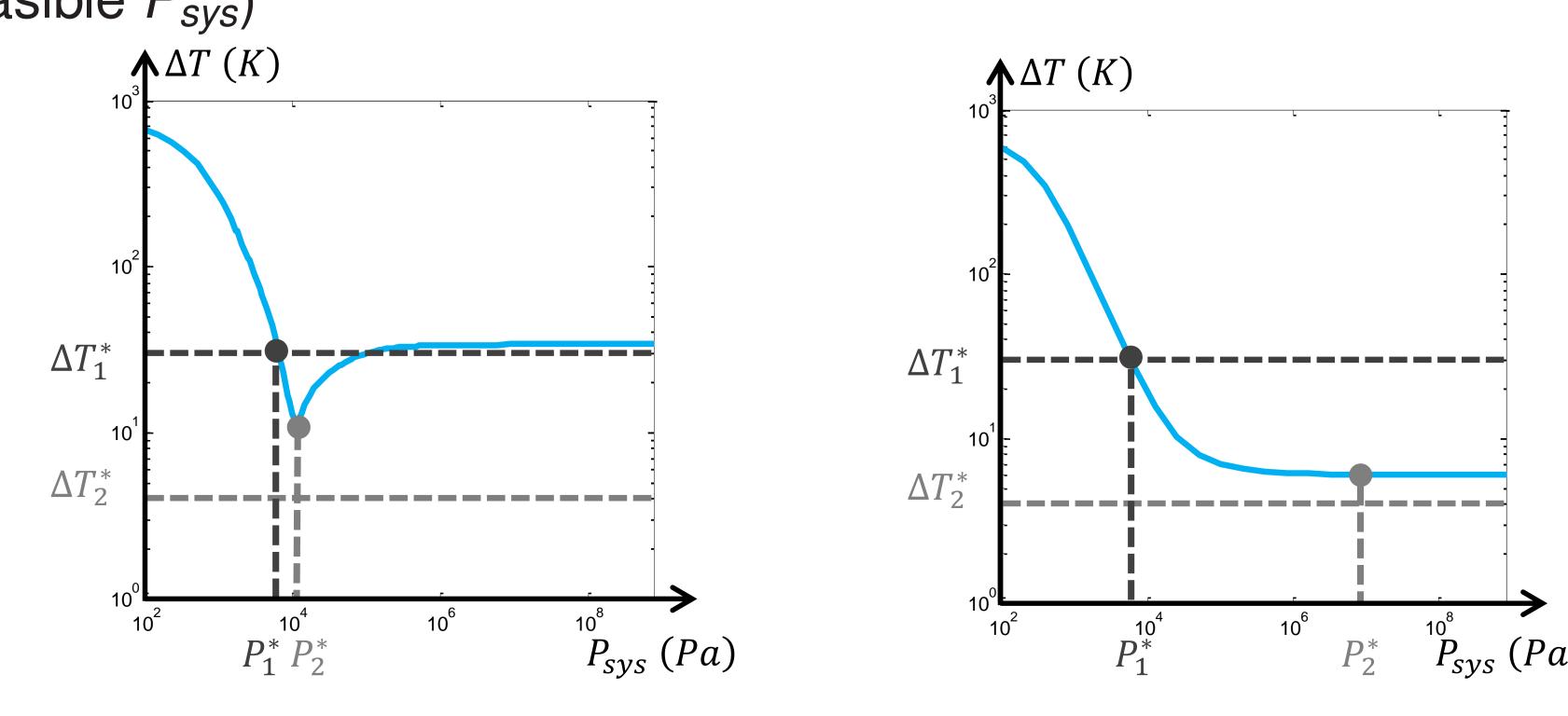
## Network Evaluation of Pumping Power Minimization



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

min 
$$P_{sys},$$
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  $P_{sys}$
- ightharpoonup Otherwise, return the  $P_{svs}$  for minimum f (show the nonexistence of feasible  $P_{svs}$ )

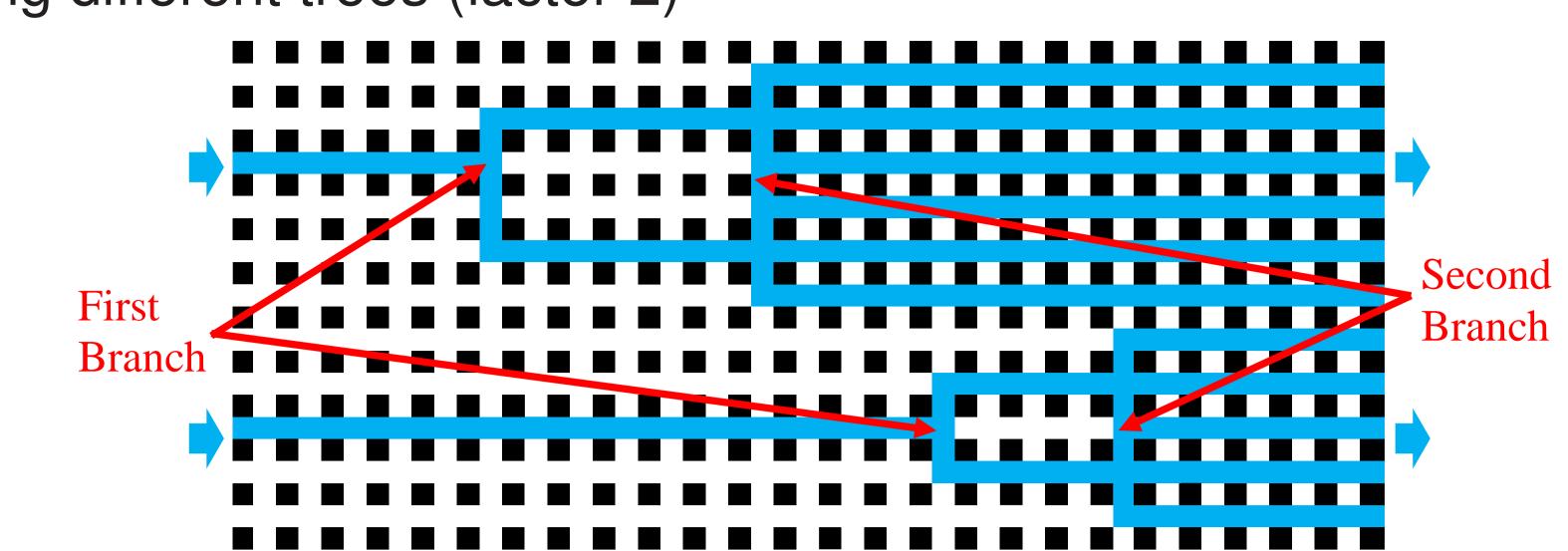


#### **Pumping Power Minimization**

#### Tree-like Cooling Network

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

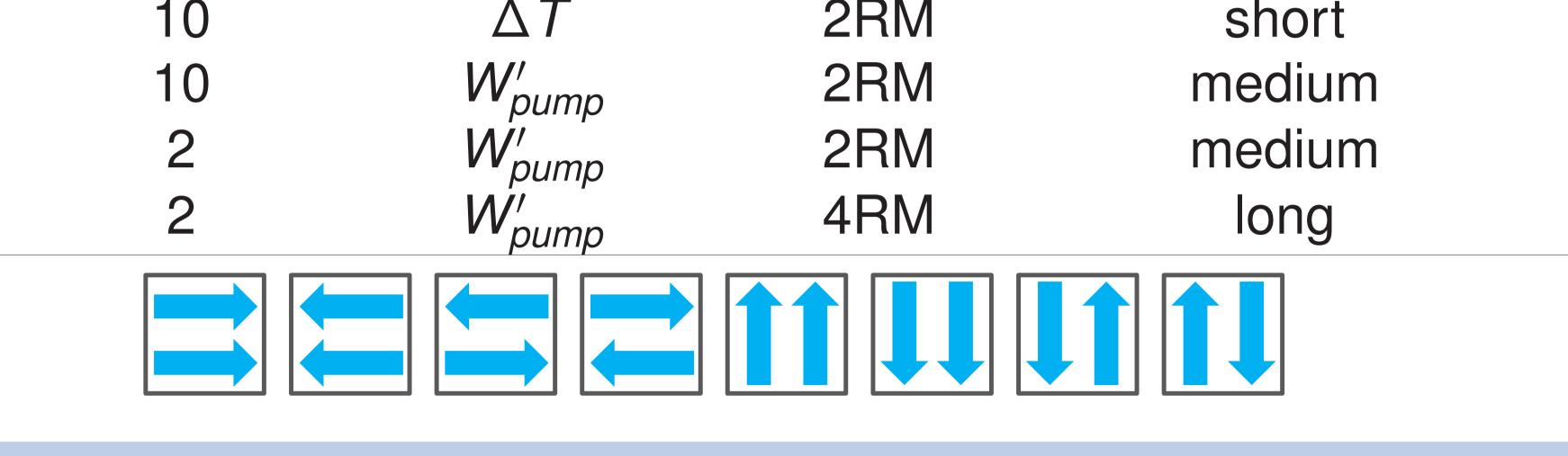
- Between upstream and downstream (factor 1)
- Among different trees (factor 2)



#### Network Topology Optimization

- In stage 1,  $\Delta T$  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 Simulator Runtime for an Iteration short



#### **Thermal Gradient Minimization**

Similar to solving pumping power minimization with some optimization Network Evaluation

Its simplified form becomes:

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

- 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  $\Delta T$  under a fixed  $P_{sys}$  is safe  $\implies$  **speed-up** 

- ightharpoonup 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

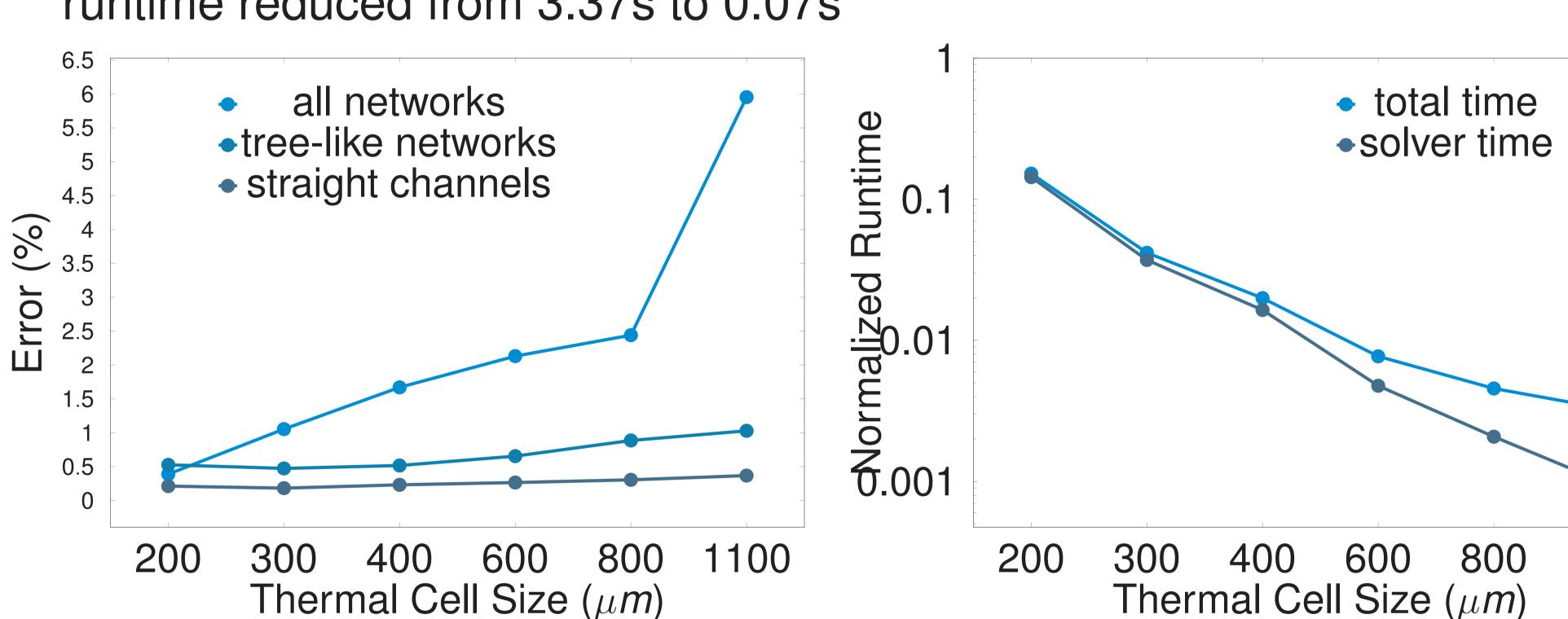
Stage # Step Size Objective Function Simulator Runtime for an Iteration short medium

## **Experimental Results**

- Faster 2RM Model 5 benchmarks, 40 network samples, 6 thermal cell sizes and 13 pressures
- ▶ Tree-like networks,  $400\mu m$  thermal cells: 0.52% errors (compared to 4RM), runtime reduced from 3.37s to 0.07s

4RM

medium



- Pumping Power Minimization
- ▶ 40 min for cases 1-3 and 240 min for case 4 79.61% better than baseline (unidirectional straight channels)

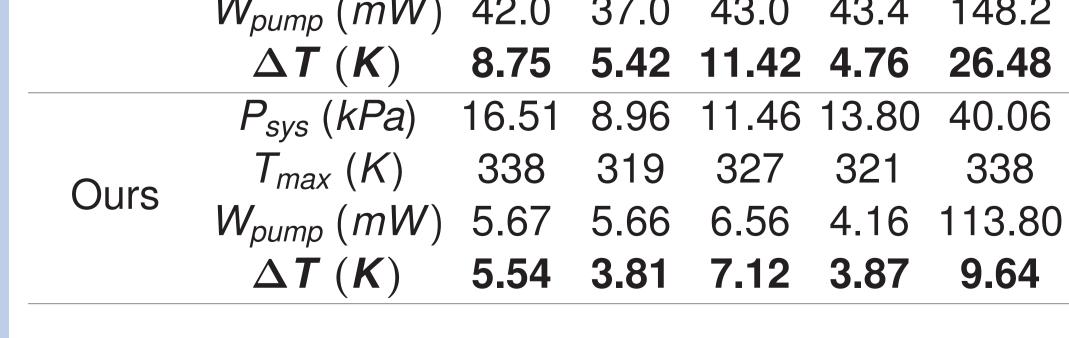
▶ 16.35% better than 1st place in ICCAD 2015 Contest 12.98 6.23 7.85 9.71 N/A 15.0 10.0 15.0 10.0 N/A  $W_{pump}$  (mW) 10.41 6.91 8.34 11.65 N/A 357 336 328 336 338  $\Delta T(K)$  15.0 10.0 15.0 10.0 10.0 Contest) W<sub>pump</sub> (mW) 1.72 1.51 3.36 2.96 113.96 bottom source layer  $P_{svs}$  (kPa) 8.72 5.13 5.81 8.27 40.10

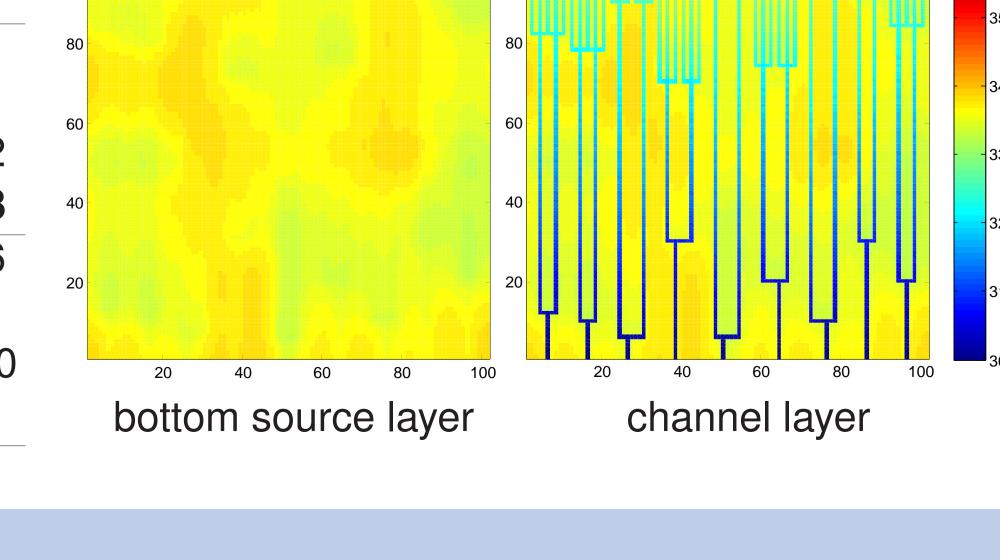
#### $W_{pump}$ (mW) 1.66 1.37 1.90 2.68 113.96 Thermal Gradient Minimization

► Constraint  $W_{pump}^*$  on  $W_{pump}$  is set to 0.1% of die power 37.27% better than baseline

15.00 10.0 15.0 10.00 10.00

P<sub>svs</sub> (kPa) 26.08 14.43 17.82 26.51 45.81  $W_{pump}$  (mW) 42.0 37.0 43.0 43.4 148.2





channel layer

E-Mail: gjchen@cse.cuhk.edu.hk