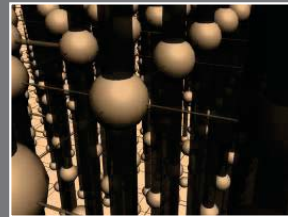
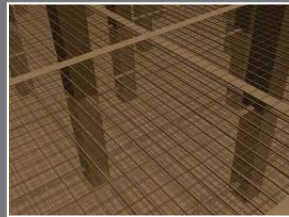
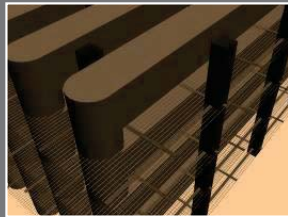


# Snap-3D: A Constrained Placement-Driven Physical Design Methodology for Face-to-Face-Bonded 3D ICs



Pruek Vanna-iampikul, Chengjia Shao, Yi-Chen Lu,  
Sai Pentapati, and Sung Kyu Lim  
Georgia Institute of Technology



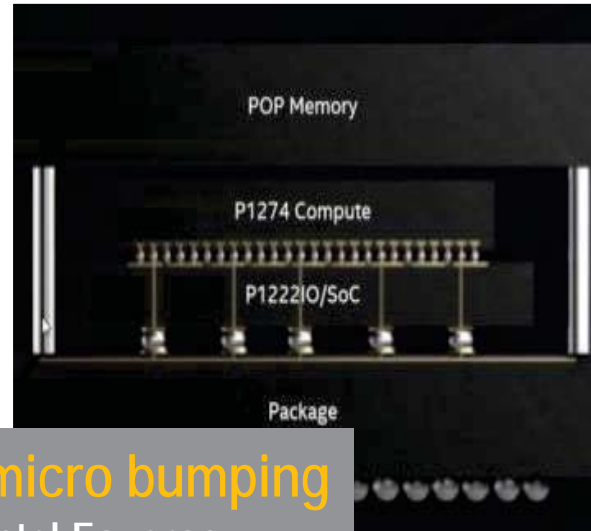
- Pseudo-3D vs. True-3D physical design flows
- Snap-3D flow
  - Overview
  - Details
  - Strengths
- Experimental Results
- Conclusions



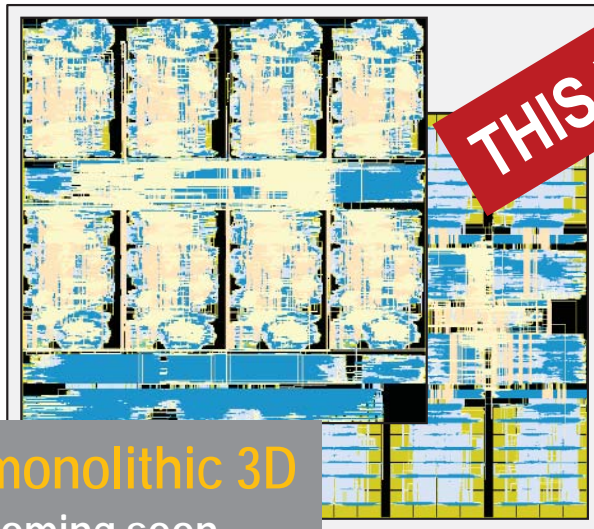
# Heterogenous Integration Technologies



2.5D interposer  
TSMC CoWoS

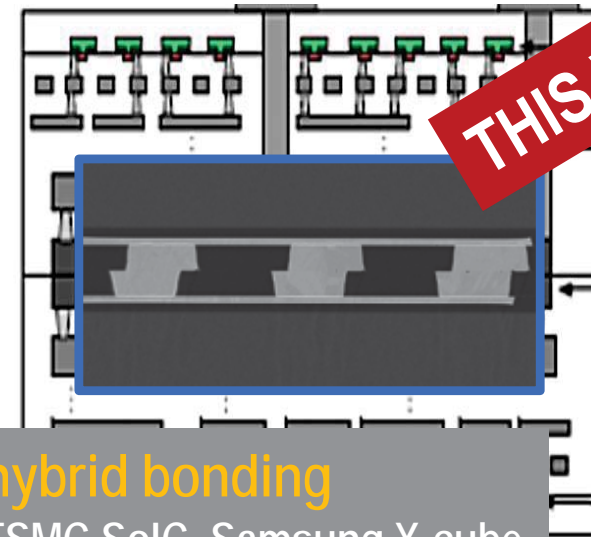


micro bumping  
Intel Foveros



monolithic 3D  
coming soon

**THIS WORK**

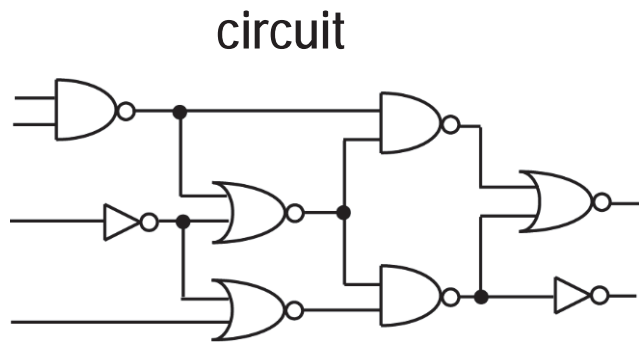


hybrid bonding  
TSMC SoIC, Samsung X-cube

**THIS WORK**



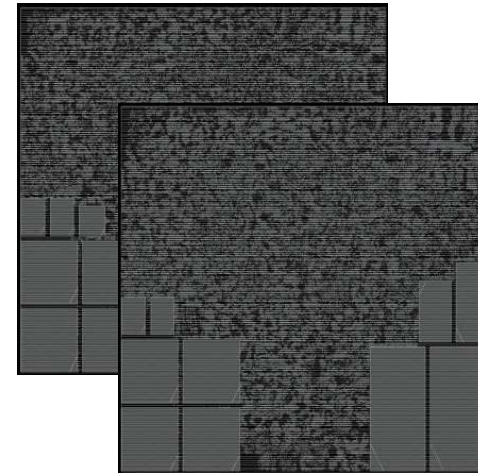
# Pseudo-3D vs. True-3D EDA Tools



1



true-3D  
place/route



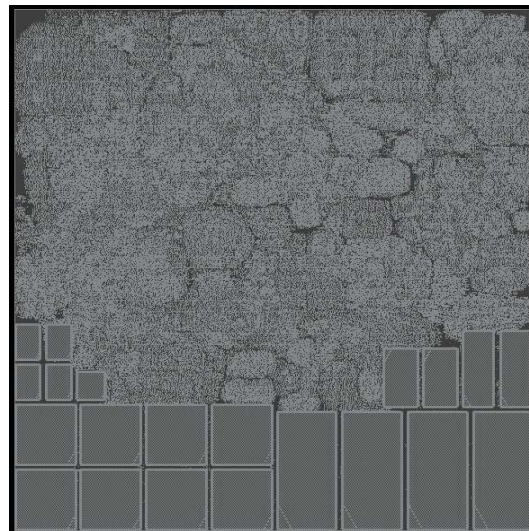
commercial tool **NOT READY**

2



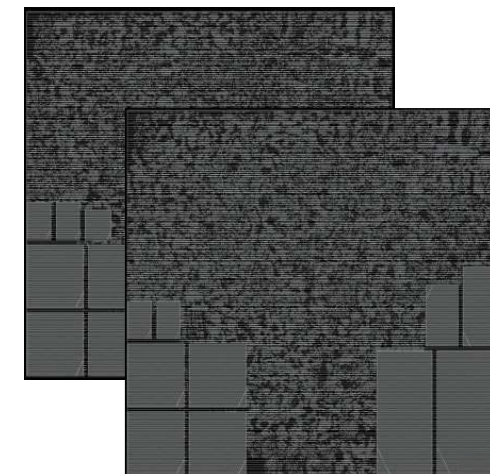
pseudo-3D  
place/route

**THIS WORK**



intermediate 2D (commercial tool **READY**)

transformation

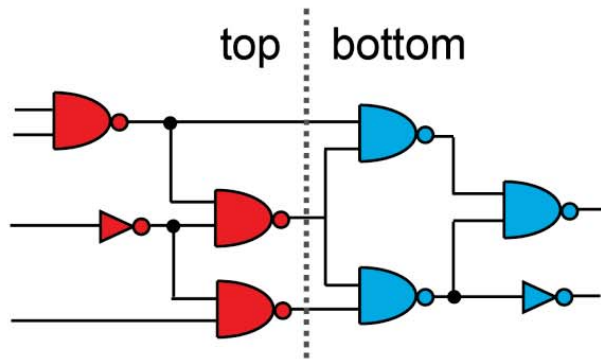


final 3D

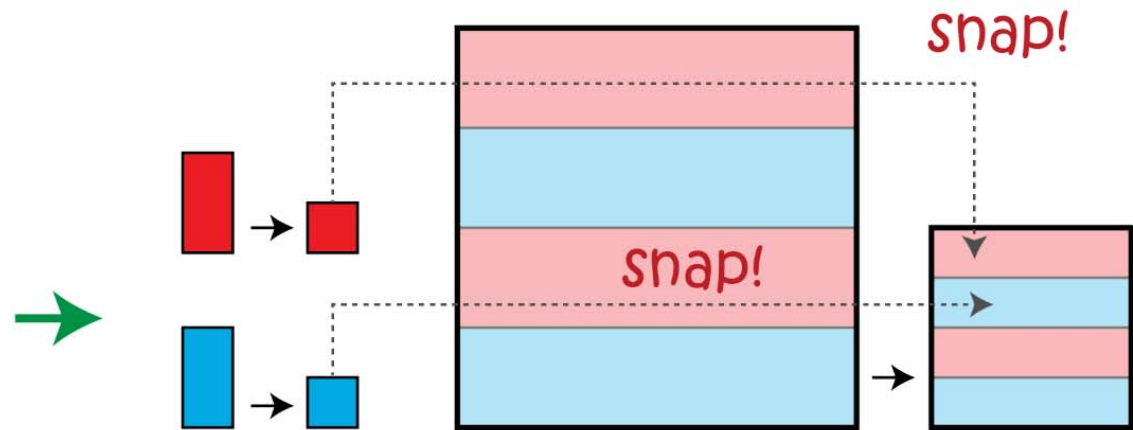


# Snap-3D: Overview

5/16

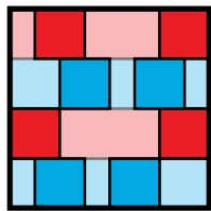


1. partition circuit

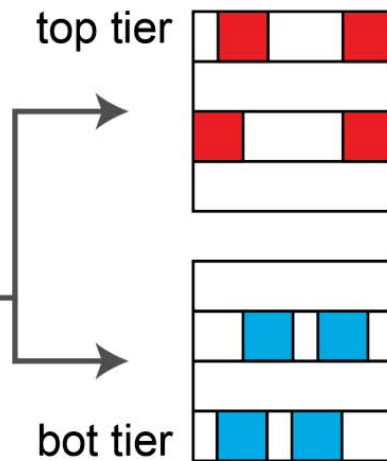


2. halve cell height + row height/width

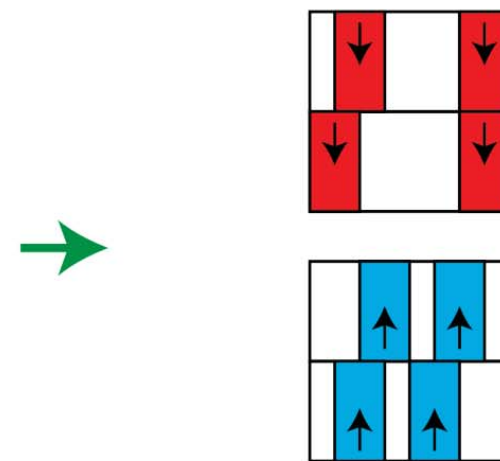
We "snap" cells to rows!



3. 2D placement with half height



4. tier partitioning



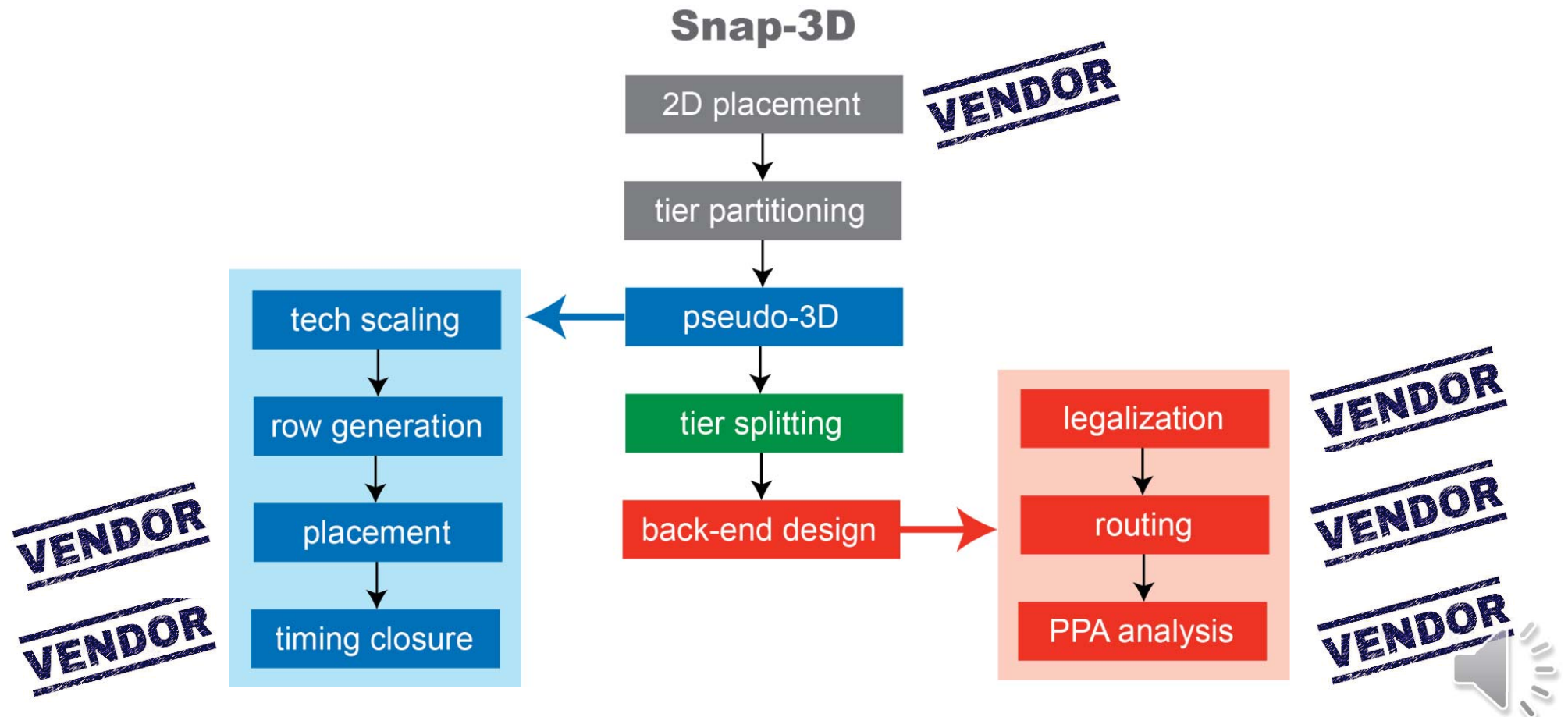
5. restore height



# Snap-3D: Design Flow

6/16

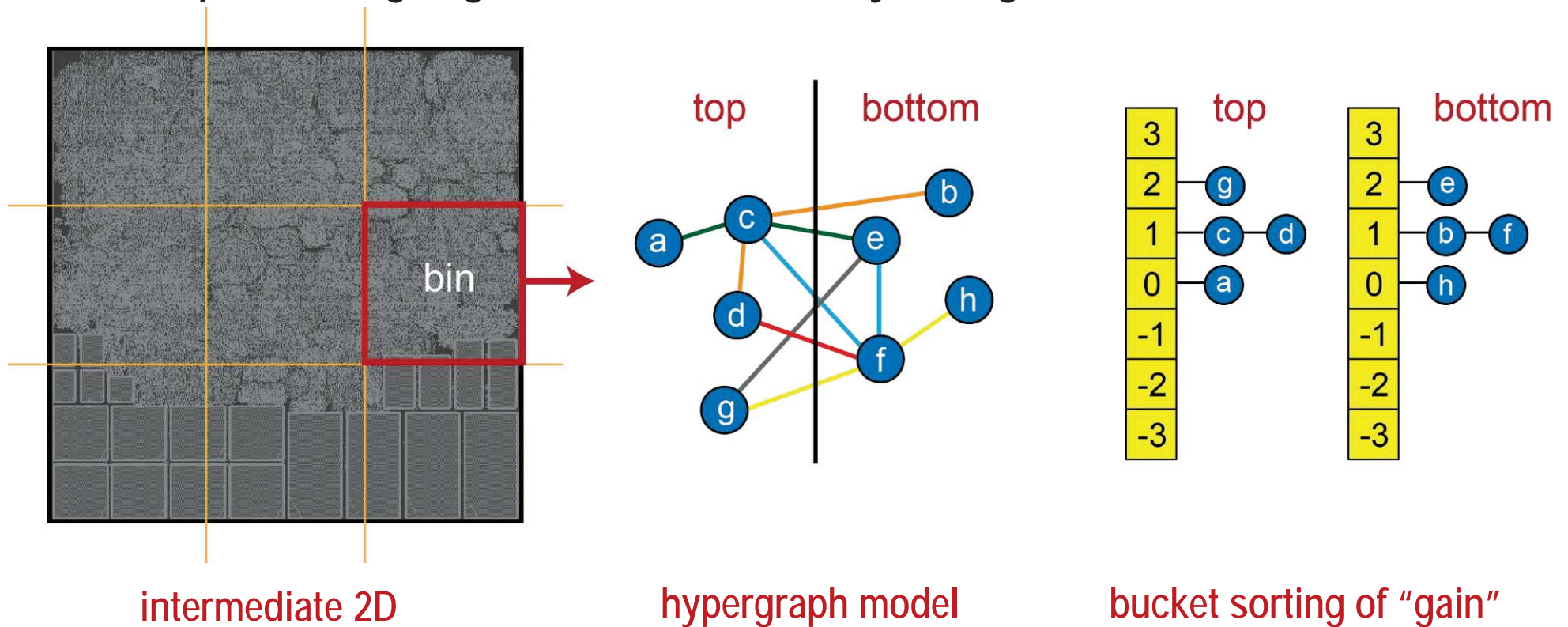
- Goal
  - Use EDA vendor tools as much as possible
  - Then add key missing engines and seamlessly integrate



# Our Automatic Tier Partitioner

7/16

- Bin-based hypergraph partitioning
  - Divide 2D into bins, and partition each bin
  - Bi-partitioning engine is Fiduccia-Matheyses algorithm [1982]



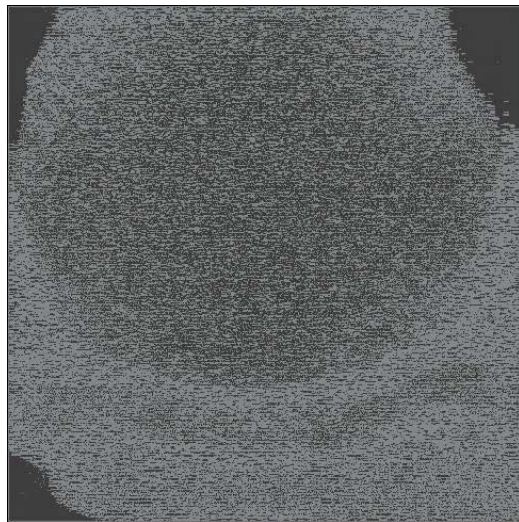
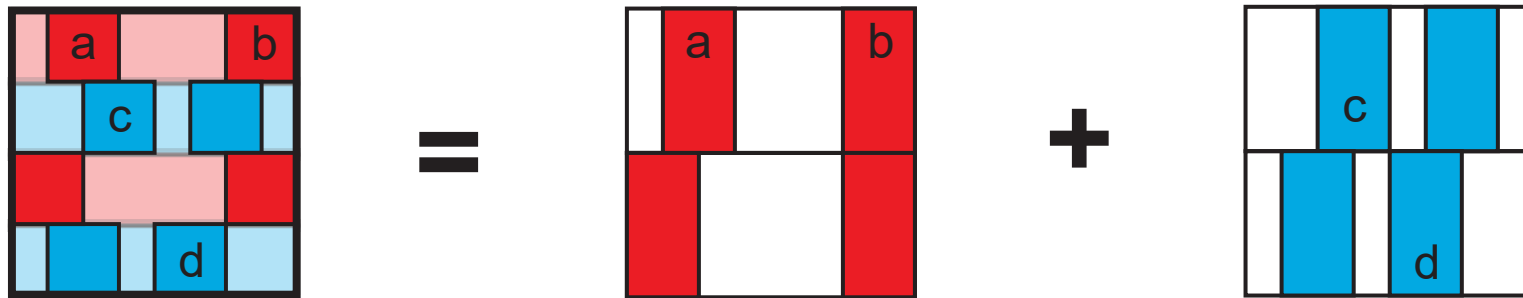
Why binning? Bin size determines F2F usage!



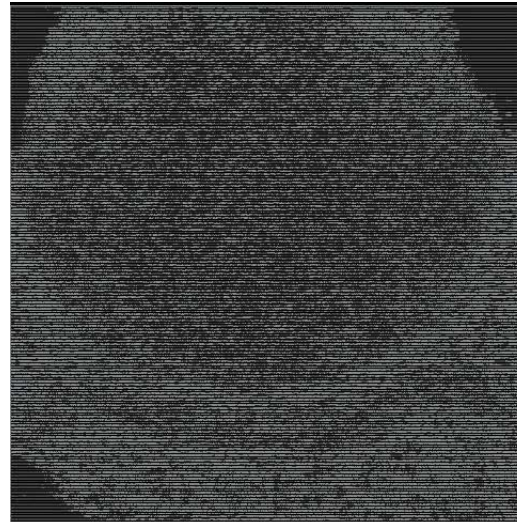
# Snap-3D: Key Benefit (1/2)

8/16

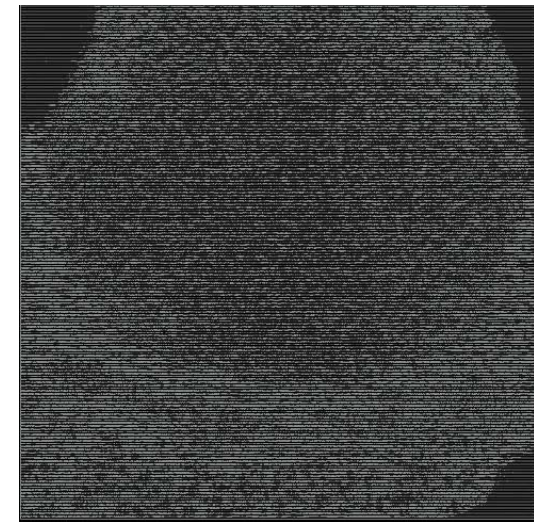
- Commercial placement quality
  - 2D placement **preserved** in 3D placement!



Snap-3D placement



top tier placement



bottom tier placement





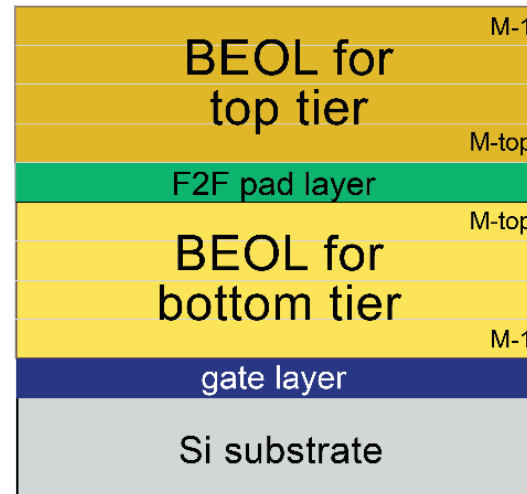


# Snap-3D: Key Benefit (2/2)

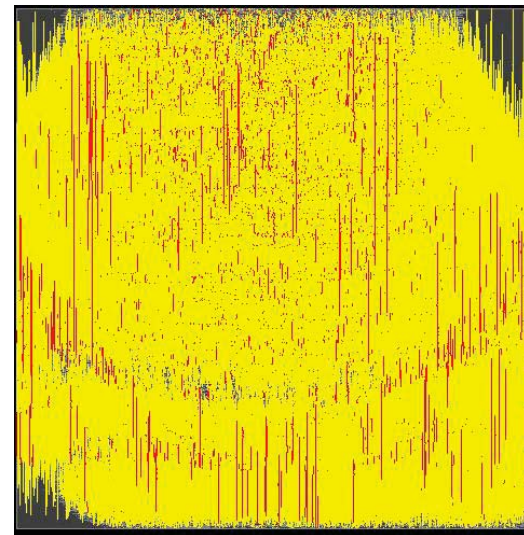
9/16

- Commercial routing quality
  - We route both tiers **simultaneously** with double metal stack
  - This allows **metal layer sharing!**

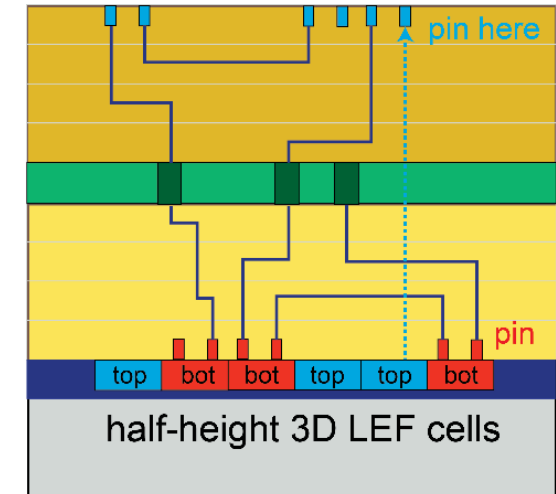
-  connecting cells in the bottom tier
-  **connecting cells in the top tier!!!**



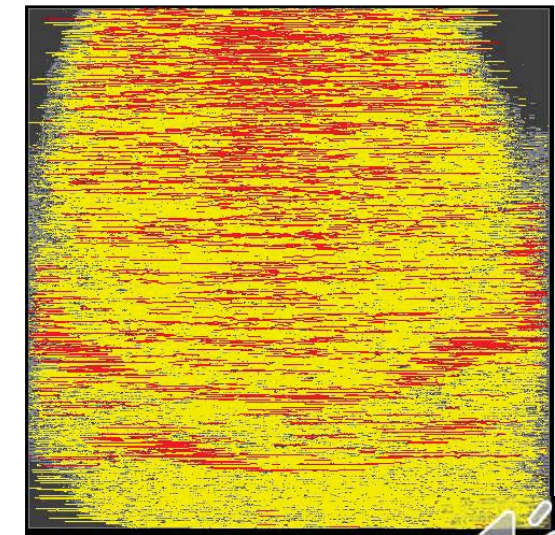
cross-sectional view



M5 Bottom



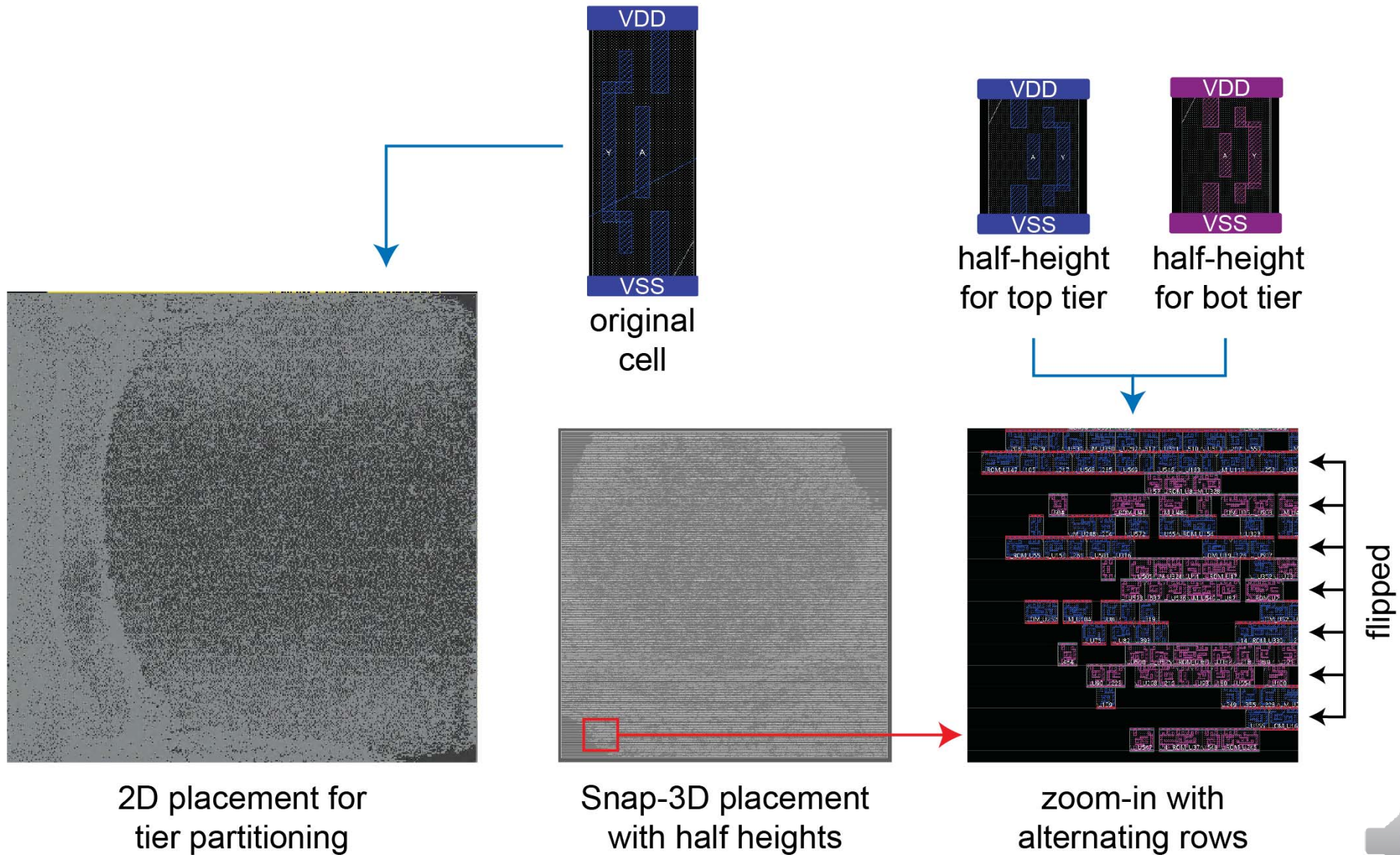
cross-sectional view



M6 Bottom



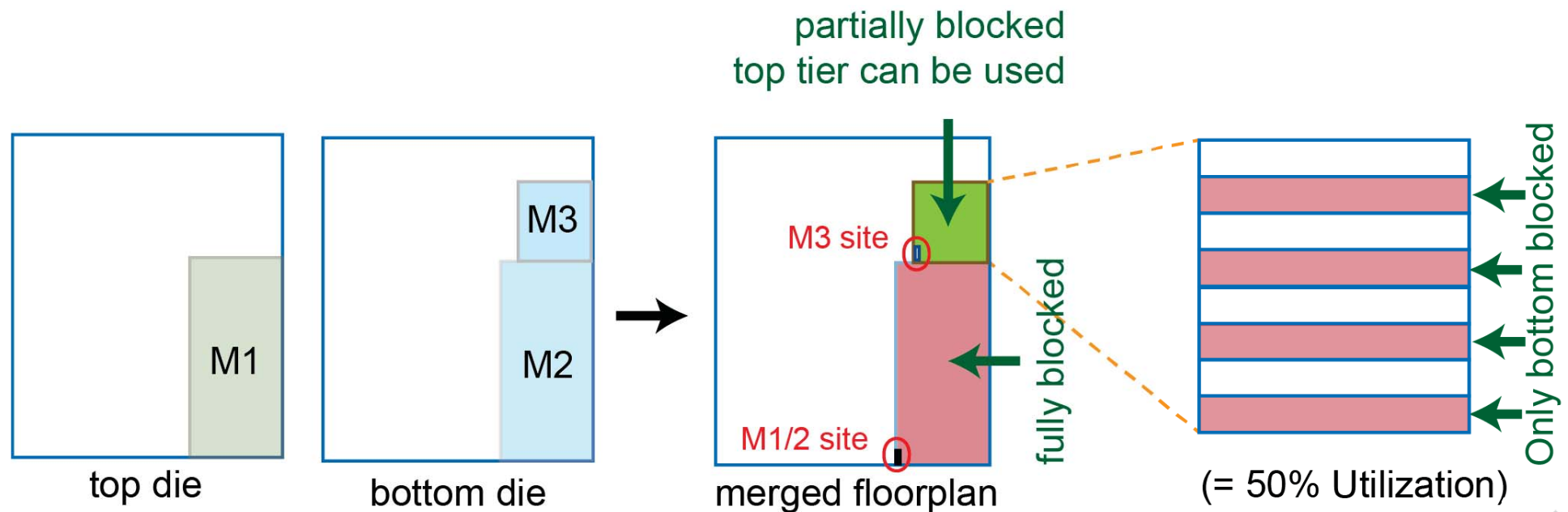
# Snap-3D: Placement Sample



# Handling Memory Macros

11/16

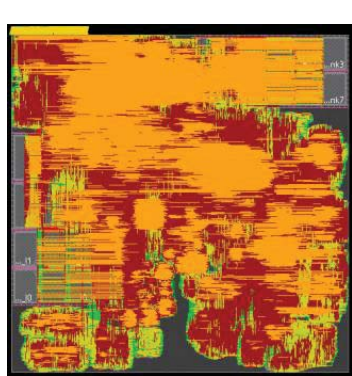
- Memory macros are used in processor designs
  - Mostly placed manually: become placement blockages in Snap-3D
  - If both tiers are blocked: gate placement not allowed
  - If one tier is blocked: corresponding rows are not used



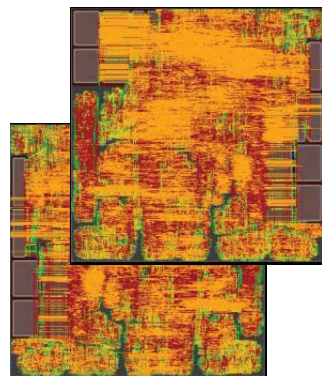
# Full-Chip GDS Layouts

12/16

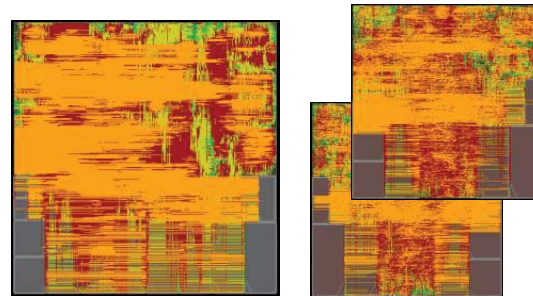
- Snap-3D using TSMC 28nm
  - Not just placement: does routing, timing closure, and PPA simulations
  - High-quality layouts: OUTPERFORMS COMMERCIAL 2D PPA



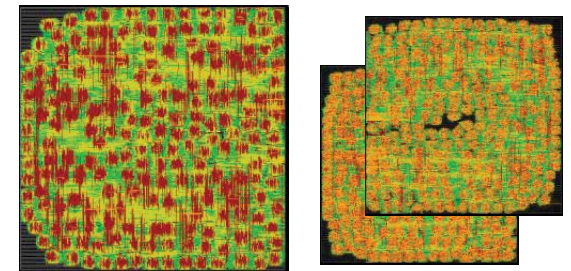
Cortex A53 2D vs. 3D



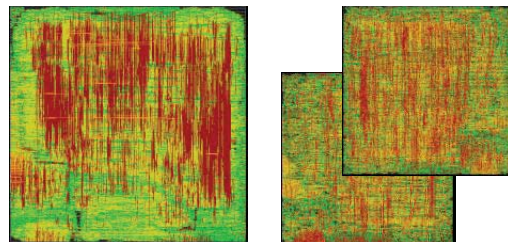
Cortex A7 2D vs. 3D



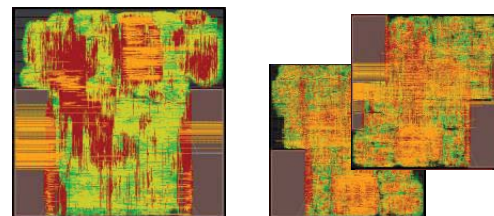
AES\_128 2D vs. 3D



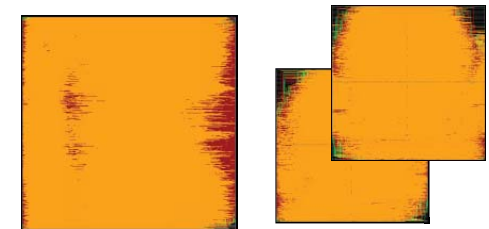
TATE 2D vs. 3D



RocketCore 2D vs. 3D



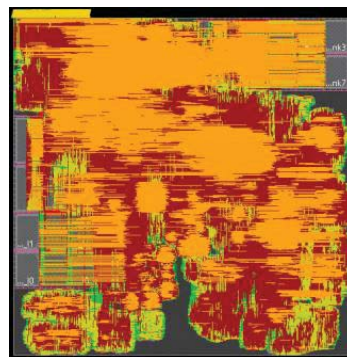
LDPC 2D vs. 3D



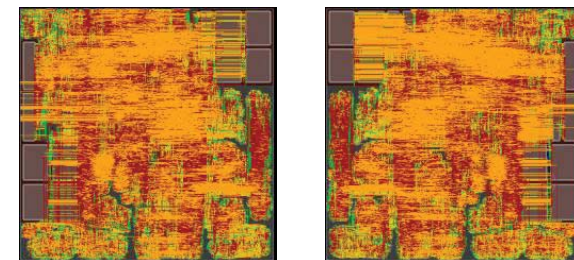
# A53 Full-Chip PPA

13/16

	2D Innovus	Shrunk-2D [2]	Compact-2D [3]	Snap-3D
target freq (GHz)	same			
footprint (mm <sup>2</sup> )	1.0	0.5	0.5	0.5
# F2F pads	-	1.0	1.01	1.15
wirelength (m)	1.0	0.69	0.70	0.73
power (mW)	1.0	0.67	0.66	0.67
WNS (ns)	1.0	0.57	1.12	0.33
<b>power × delay</b>	<b>2.10</b>	<b>1.12</b>	<b>1.46</b>	<b>0.97</b>



Innovus 2D full-chip GDS, A53



Snap-3D full-chip GDS, A53



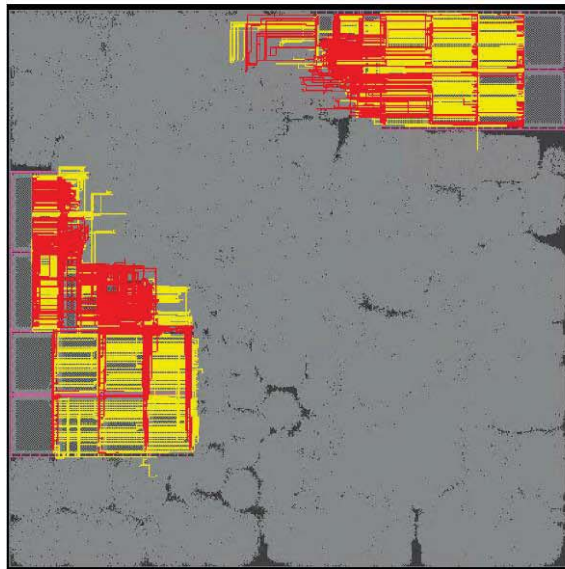
# A53 Memory Latency/Energy

14/16

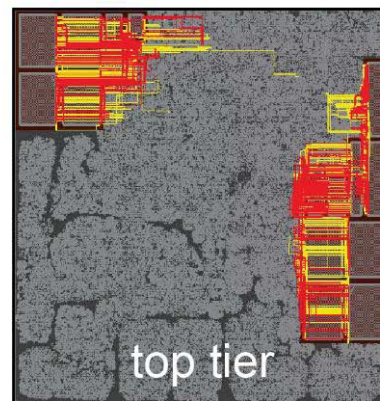
- Shorter WL in 3D
  - Helps reduce memory access latency and power!

yellow: input to memory  
red: output from memory

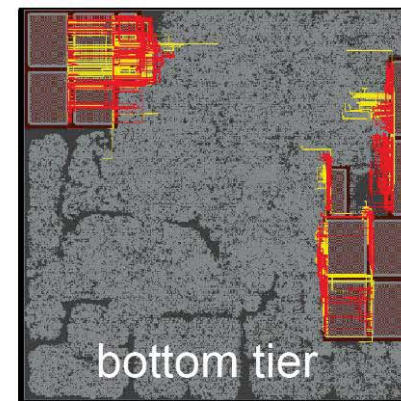
metric	2D	3D	3D gain
Energy/cycle (pJ)	3.73	2.57	30.8
Input latency (max, ps)	209	202	3.4
Input latency (ave, ps)	70	44	37.1
Output latency (max, ps)	272	125	54.0
Output latency (ave, ps)	57	28	50.9



Cortex A53 2D



top tier



bottom tier

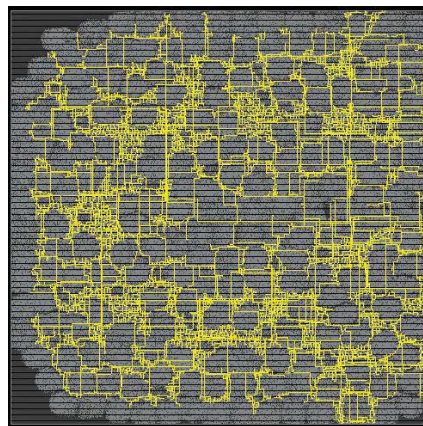
Cortex A53 3D



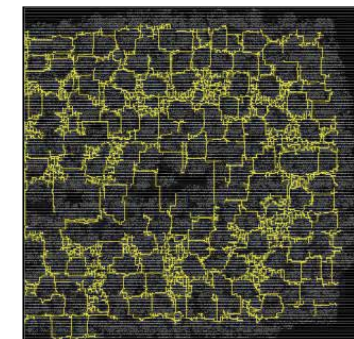
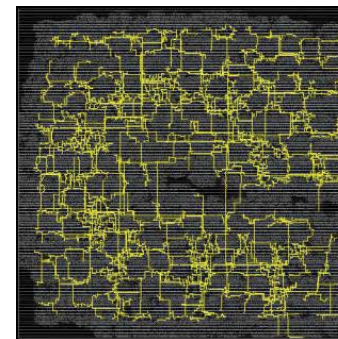
# Clock Comparison : AES @ 28nm

15/16

Clock Metrics	2D Innovus	Shrunk-2D [2]	Compact-2D [3]	Snap-3D
Clock Latency (ps)	211.8	181.5	177.6	<b>166.1</b>
Clock Skew (ps)	9.9	11.7	11.3	<b>8.5</b>
Clock WL. (mm)	43.42	42.15	41.33	<b>38.99</b>
# Clk. F2F pads	0	674	671	731
# Clock Buffer	875	910	849	<b>862</b>



clock tree for AES, 2D Innovus



clock tree for AES, Snap-3D



- **Snap-3D key ideas**
  - Use half heights (for cells and rows)
  - Do tier partitioning first and snap cells to rows (= constrained placement)
  - Use double metal stack for routing
- **Snap-3D key benefits**
  - 2D placement = 3D placement
  - Metal layer borrowing is supported
  - Outperforms Innovus 2D, Shrunk-2D [2] and Compact-2D [3]

