

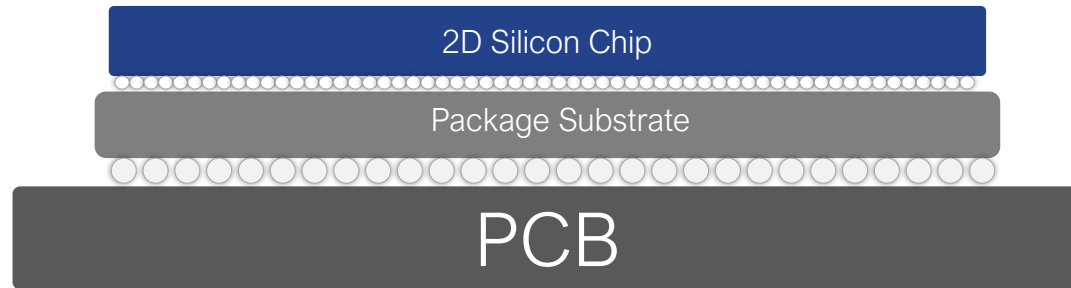


Automation and Optimization of Heterogeneous Systems

ISPD 2025

Henry Sheng
March 19, 2025

First, A Look at 2DIC



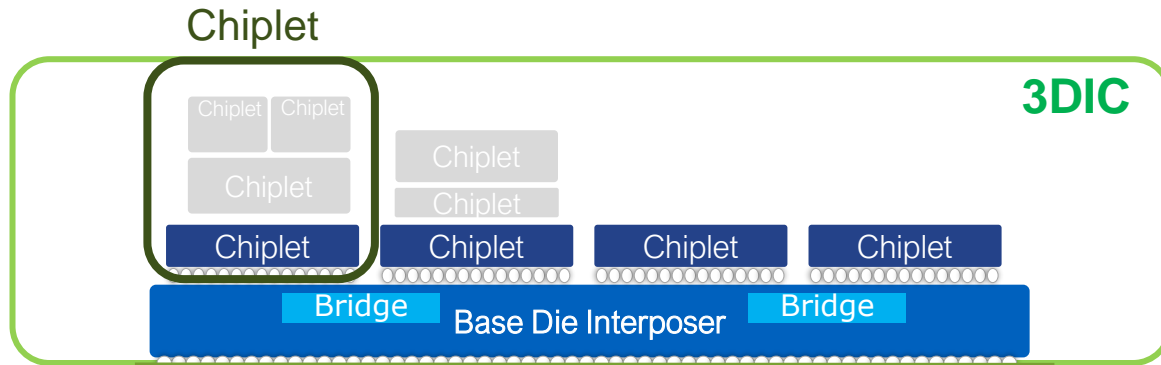
2D Flip Chip

- IC design traditionally fairly detached from lower package design
- Interface to package handled by flip chip and RDL routing solution at the chip level. Not usually visible at the block level.
- Lack of models, standards from Silicon to Package and PCB

Classic Package:

- Assembly focused, manual methods
- Classic capabilities SIP, discretes, wire bonds, etc.
- Can be single or multiple die
- Necessary piece for path to full systems solution

A Way to View 3DIC



3DIC/2.5DIC: SysMoore Scaling

- Technology densifying on annual basis
- Based on structure, not material
- High-density, part of overall logic system
- Material can co-mingle Si and Organic.
- Prompts need for unified design, implementation and analysis solution, along with automation

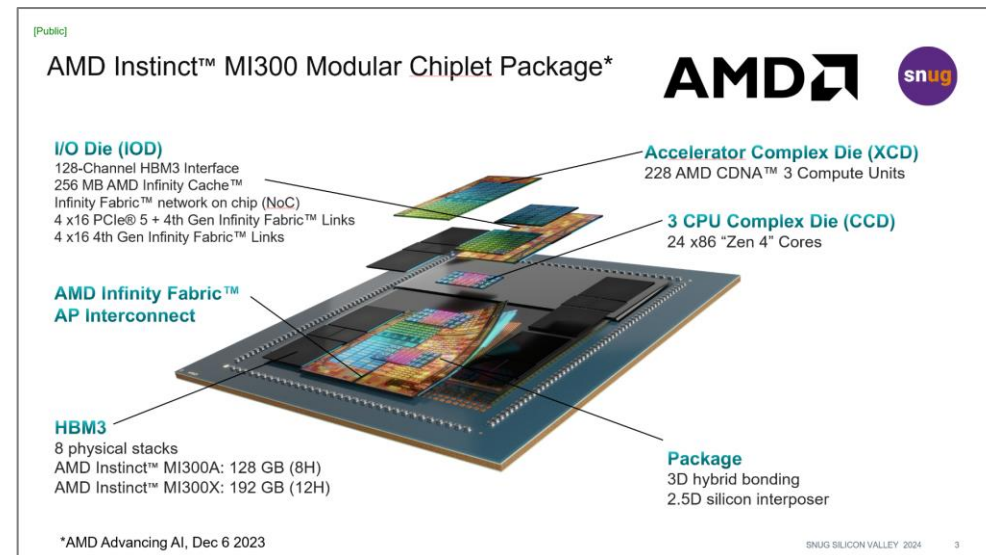
Classic Package: Form Factor Limited

- Assembly focused, manual methods
- Classic capabilities SIP, discretes, wire bonds, etc.
- Can be single or multiple die
- Densities fairly static, bound by form factor

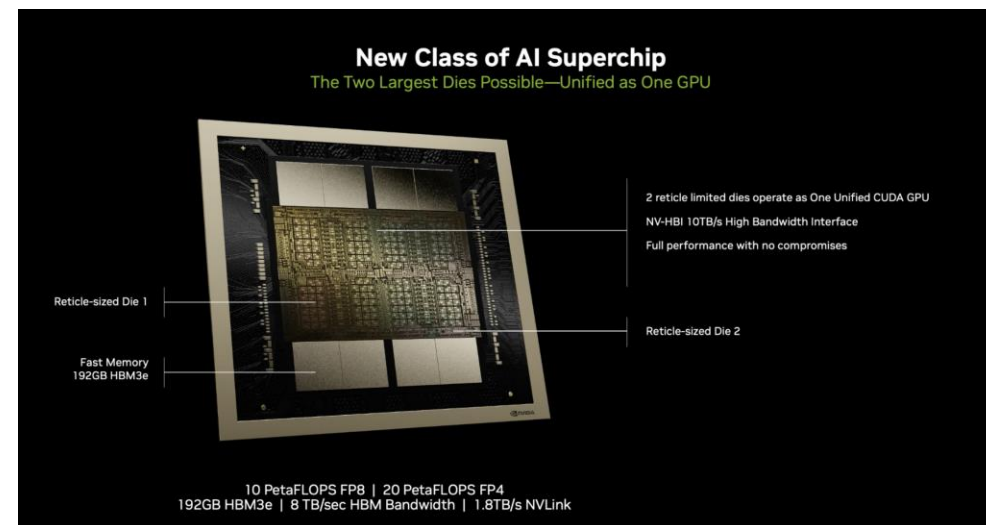


Multi-Die is Here, Now

Chip Name	Transistor Count (B)	Year	References
Cerebras CS-3 (WSE-3)	4000	2024	[Cerebras Official](https://www.cerebras.net), [TechCrunch](https://techcrunch.com)
Cerebras WSE-2	2600	2021	[Cerebras WSE-2 Details](https://www.cerebras.net)
NVIDIA Blackwell (GB200)	208	2024	[NVIDIA Blog](https://www.nvidia.com), [TechRadar](https://www.techradar.com)
AMD Instinct MI300	146	2023	[AMD MI300 Details](https://www.amd.com), [TechSpot](https://www.techspot.com)
Apple M2 Ultra	134	2023	[Apple Newsroom](https://www.apple.com), [TechCrunch](https://techcrunch.com)
Apple M1 Ultra	114	2022	[Apple M1 Ultra](https://www.apple.com), [TechRadar](https://www.techradar.com)
Intel Ponte Vecchio	100	2022	[Intel Ponte Vecchio](https://www.intel.com), [AnandTech](https://www.anandtech.com)
NVIDIA H100	80	2022	[NVIDIA H100 Details](https://www.nvidia.com)
Tesla Dojo	50	2024	[Trendforce](http://trendforce.com)
Broadcom Tomahawk 5	31	2022	[Broadcom Tomahawk 5](https://www.broadcom.com), [Network World](https://www.networkworld.com)
Intel Core i9-12900K (Alder Lake)	20	2022	[Intel Core i9-12900K](https://www.intel.com), [Tom's Hardware](https://www.tomshardware.com)
Intel Lakefield	10	2020	[Intel Lakefield](https://www.intel.com), [AnandTech](https://www.anandtech.com)

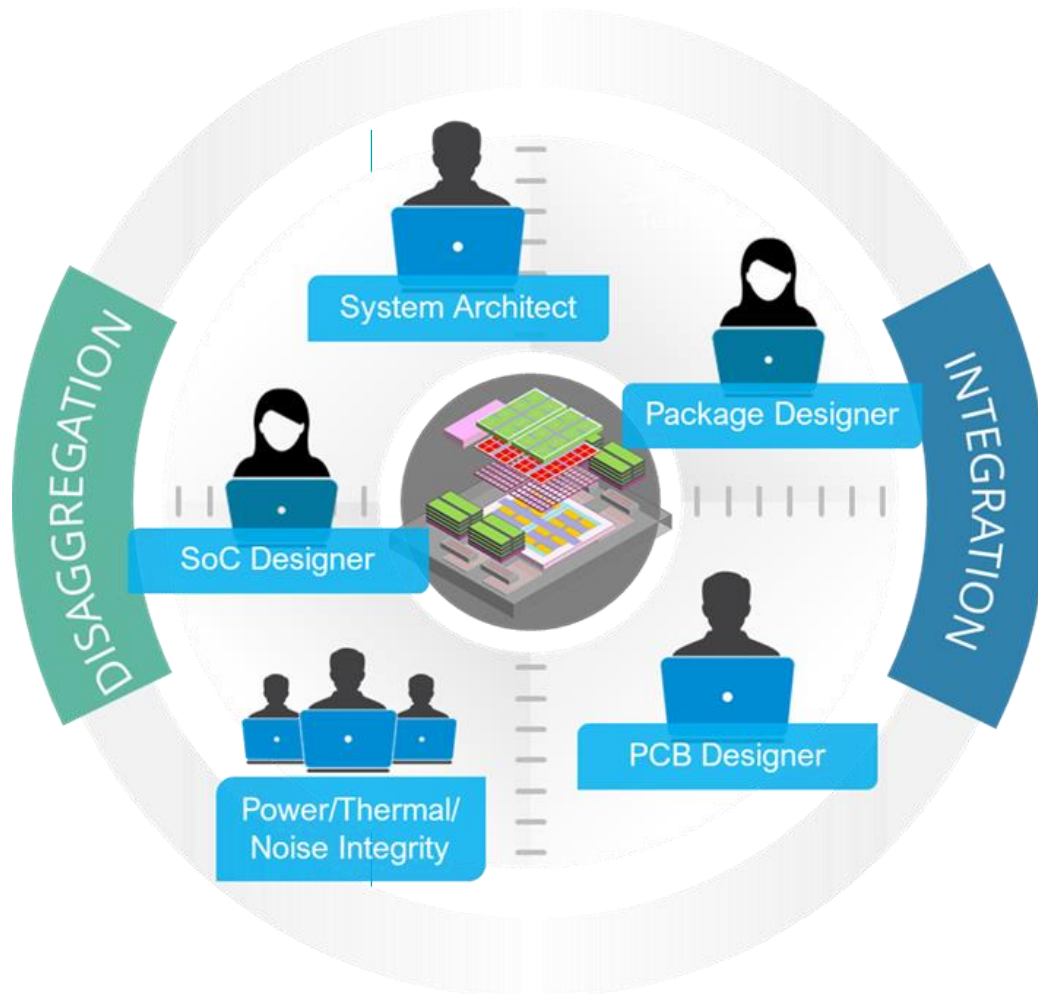


Source: AMD



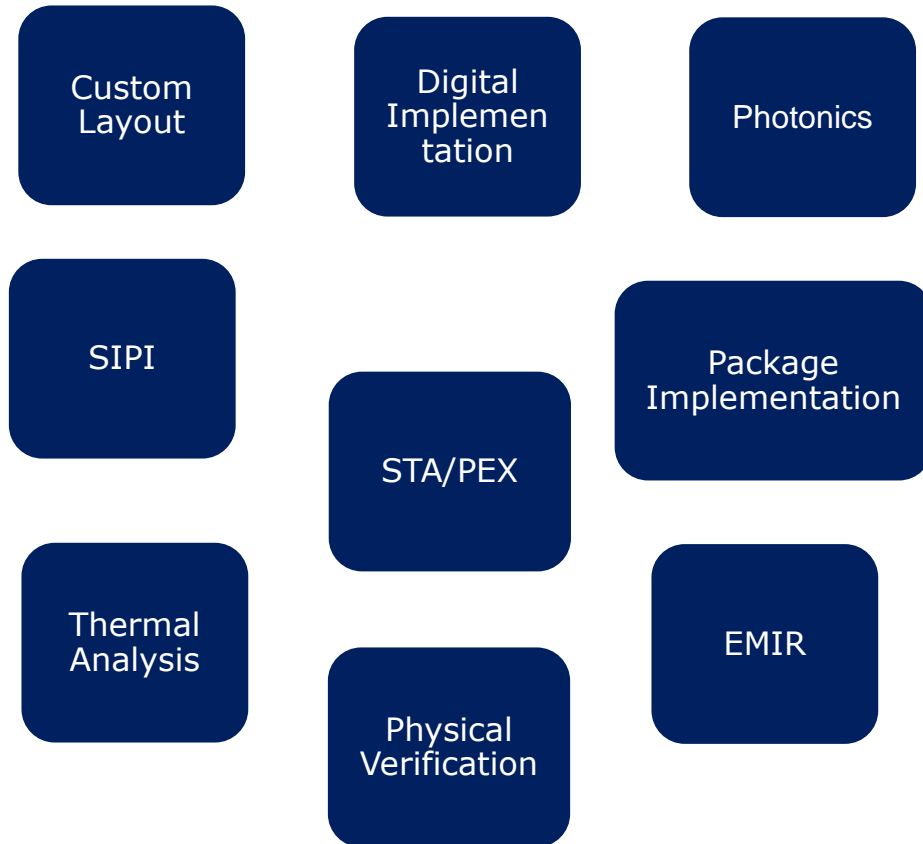
Source: NVIDIA

3D-IC Design: Multi-Team, Multi-Die, Multi-Node ...



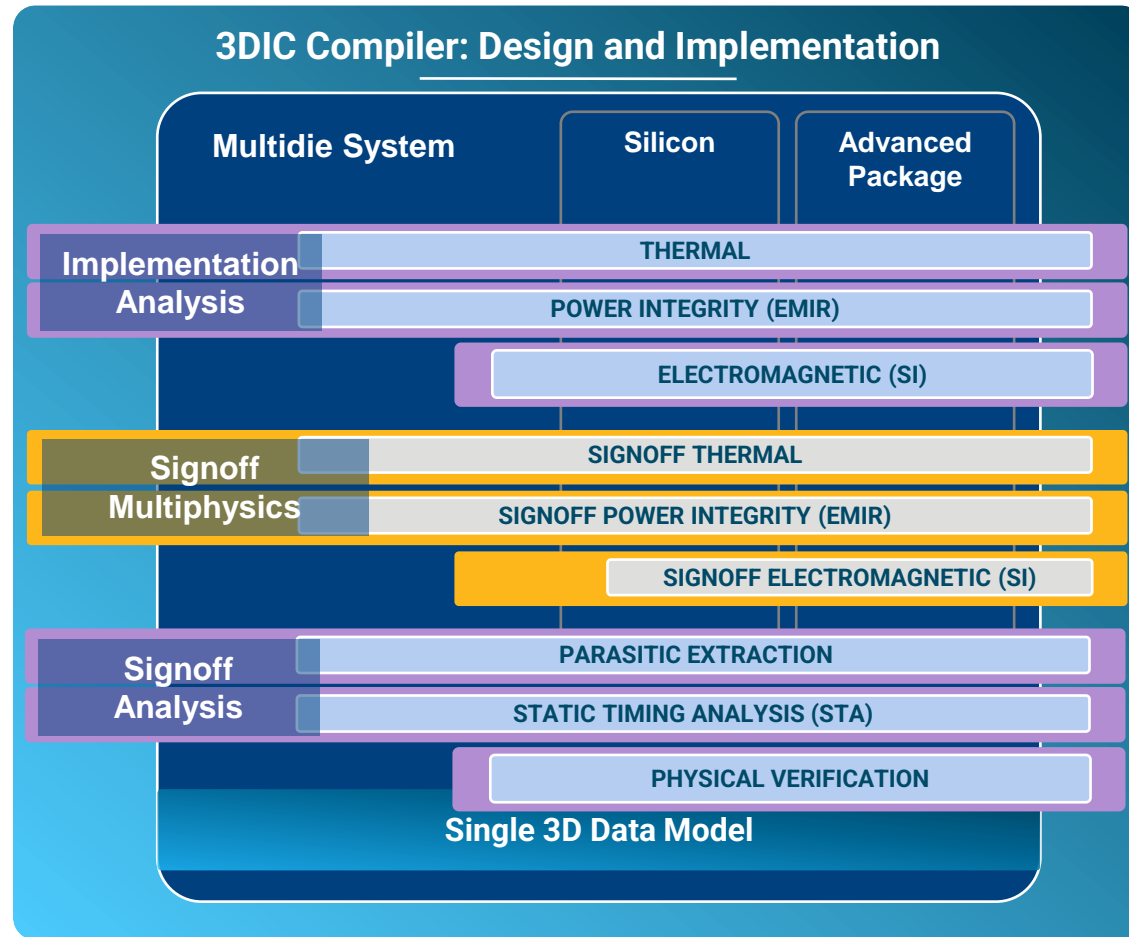
Multi-Die Not Only Integrates
Chiplets, But Also Integrates
Historically Independent Workflows

2010's Tooling for Chiplet EDA

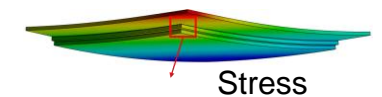
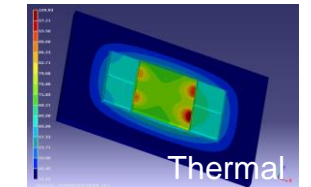
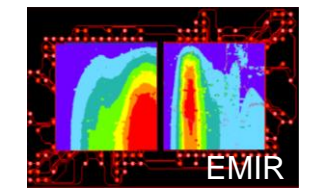
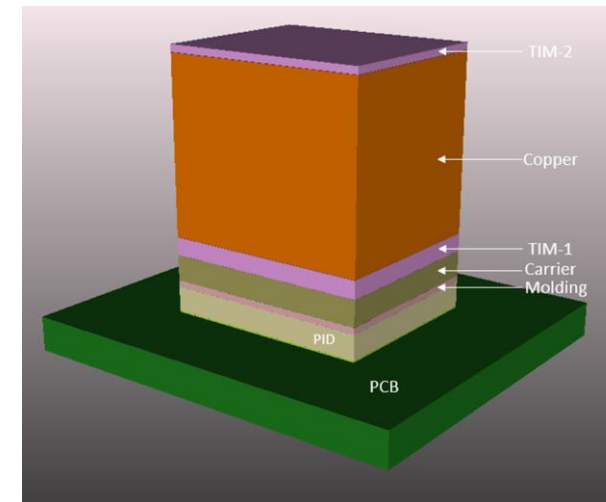


- Collection of disparate, unconnected single-die tools
- Different disciplines, different tools, different nomenclature
- Co-optimization is necessary to converge on designs, but prohibitively difficult

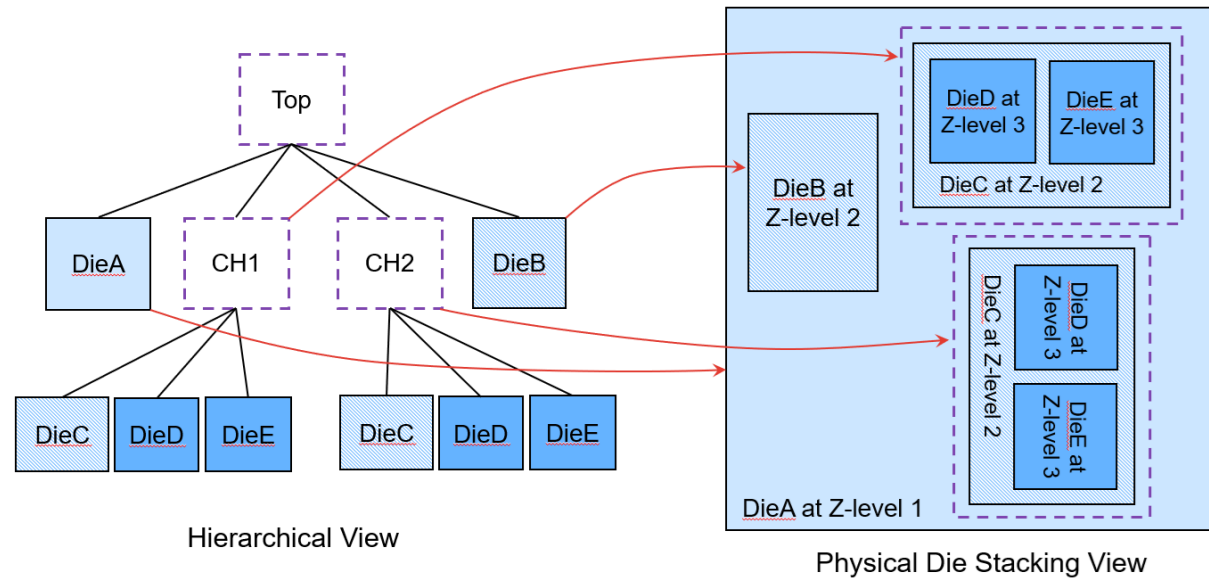
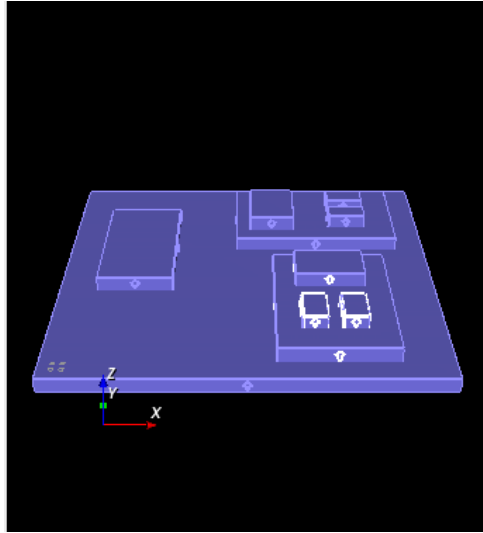
Modern Multi-Die EDA Platform



- Full Representation of Multi-Die Stack
- Co-Optimization → Optimization

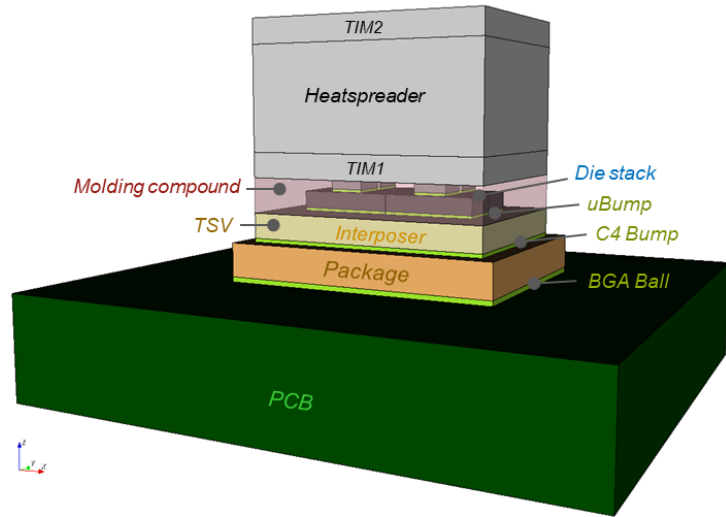


Heterogeneous Design Implementation

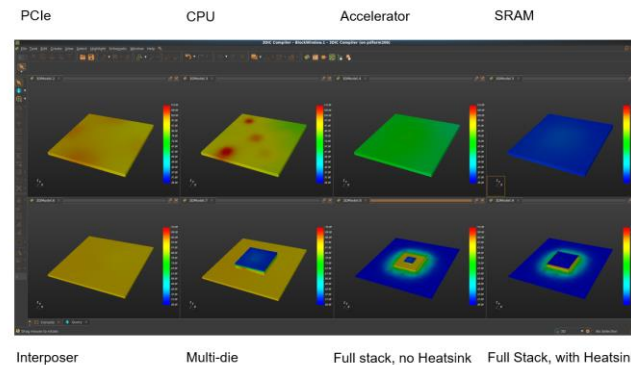
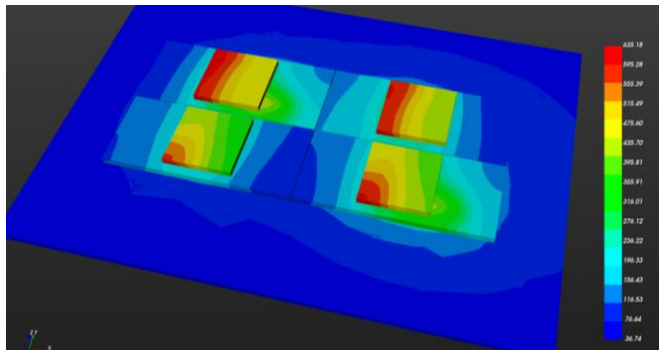


- Single environment for heterogeneous integration and design
 - Unified representation: Optical and Thermal shrinks, Orientations, Alignments
- Enable hierarchy handling – 3DIC's of 3DIC's
- Leaf components can be chiplet, interposer, package, optical, PCB, substrate, etc. each with well-defined, homogeneous process technology

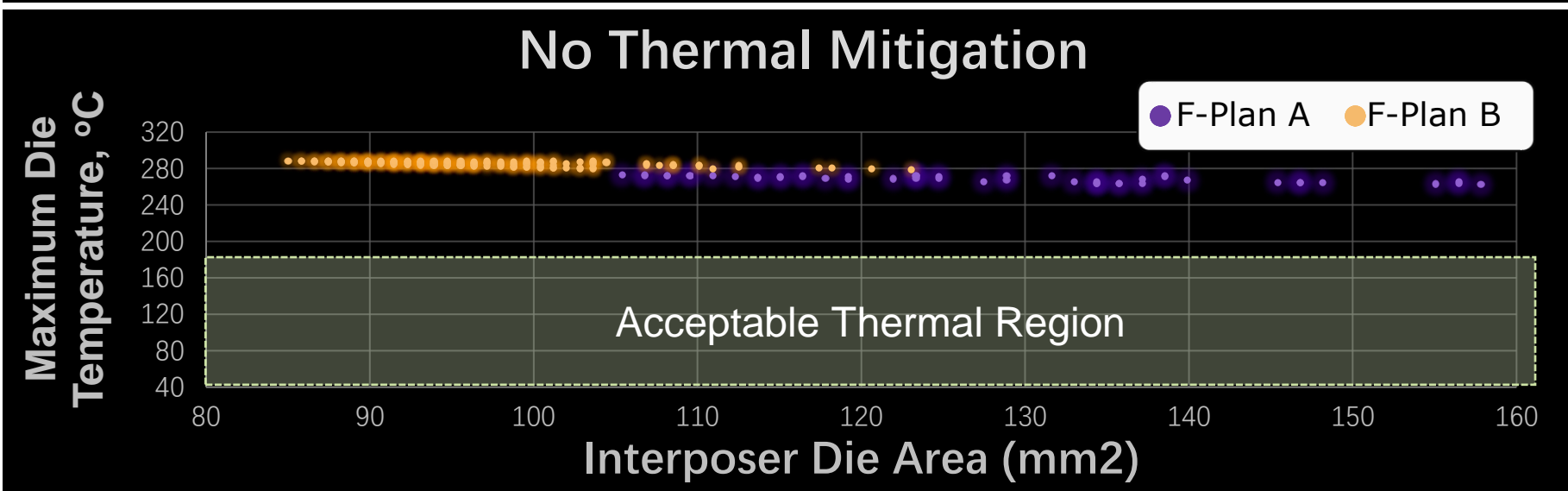
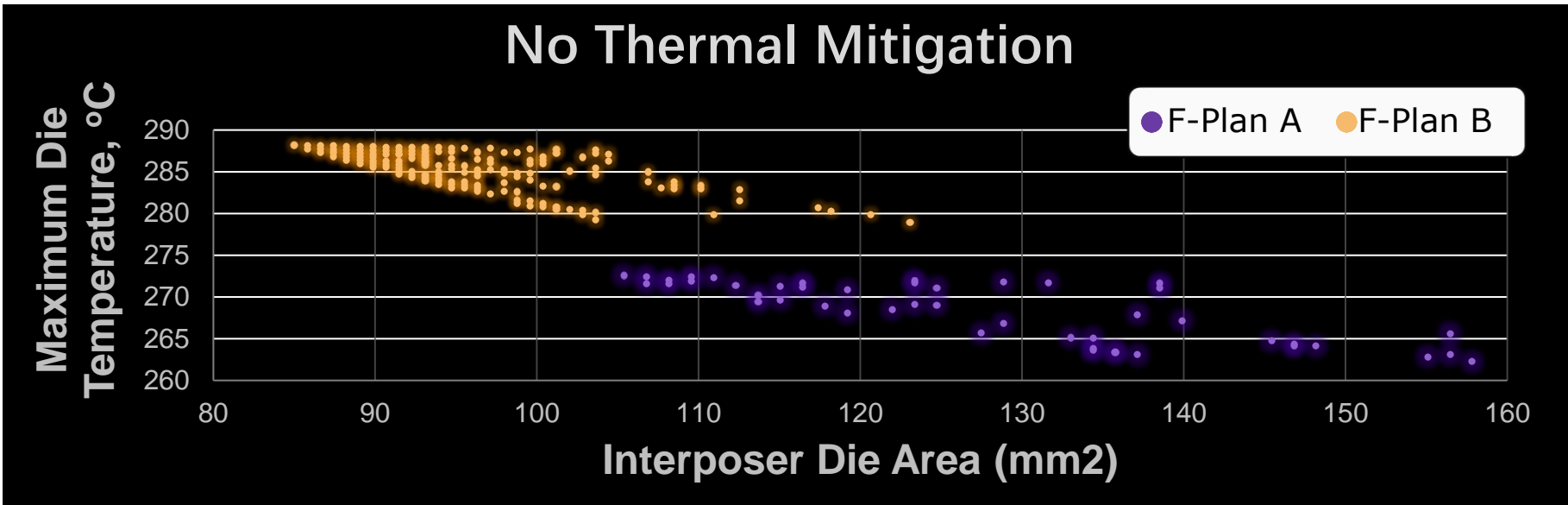
Multi-Die System Design Optimization



- Many possibilities on system optimization
- Degrees of Freedom: Design change, Die Orientation, Floorplan, Materials, Thermal components, Technology definition, IP design...
- Some Objectives: Thermal, Mechanical, EMIR, Physical Footprint

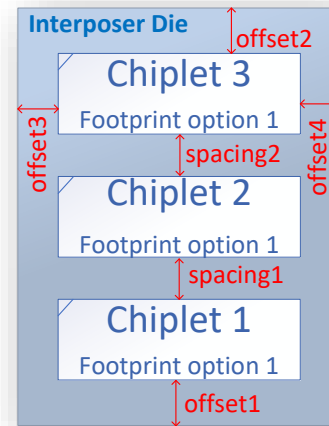


Architecture & Thermal Impact

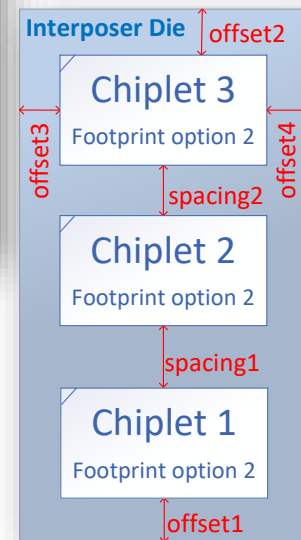


Floorplan Options

● F-plan A

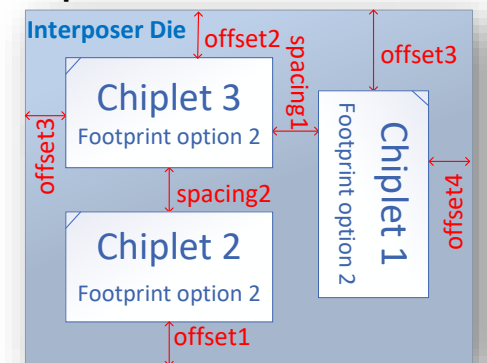


● F-Plan B

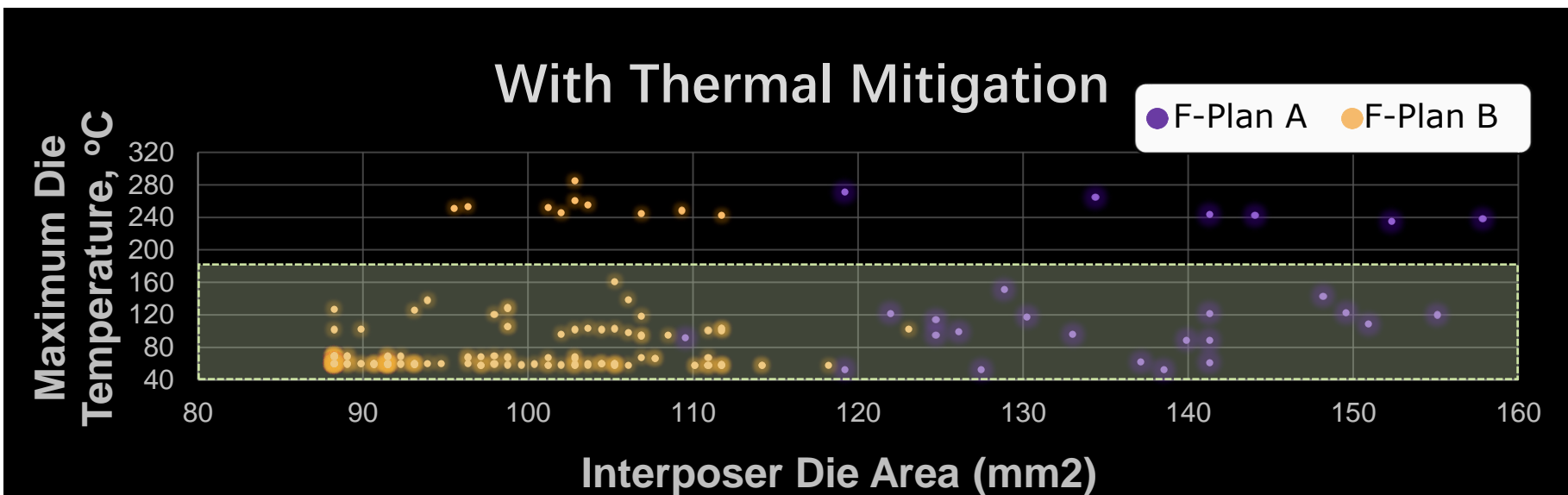
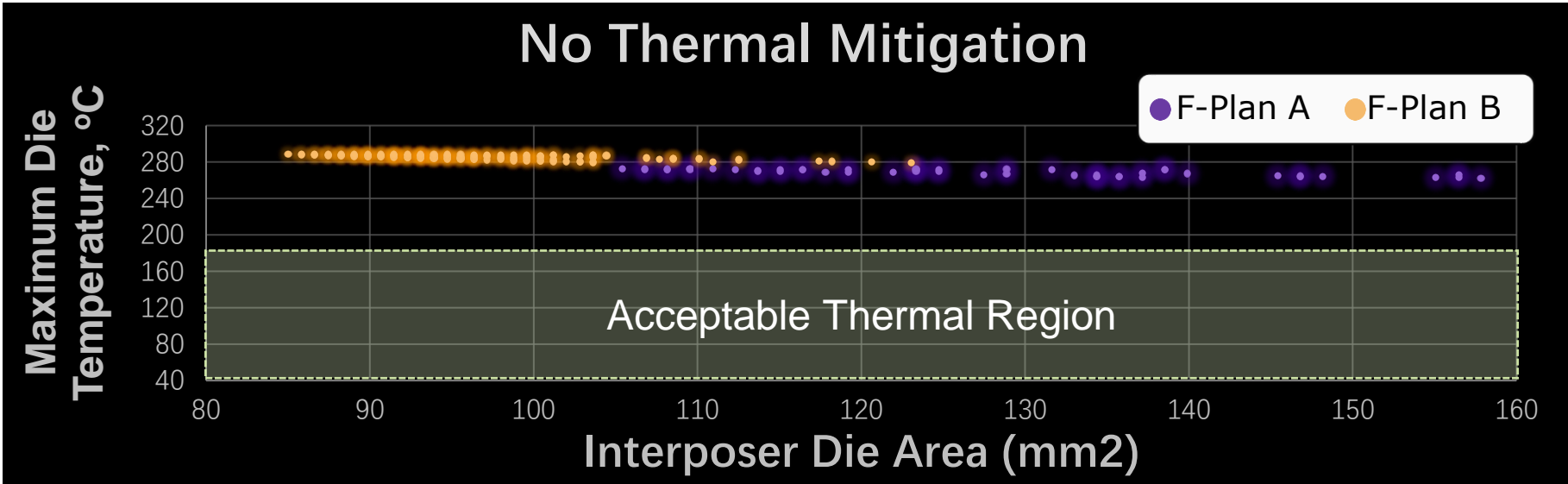


●
●
●

F-plan X



Architecture & Thermal Impact



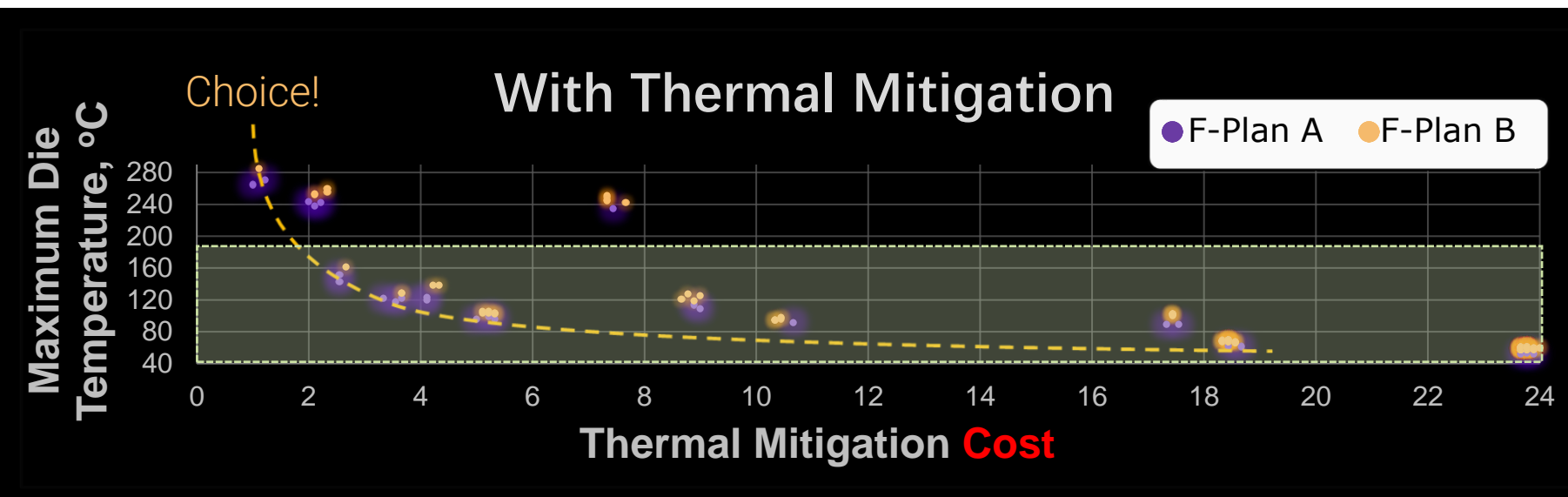
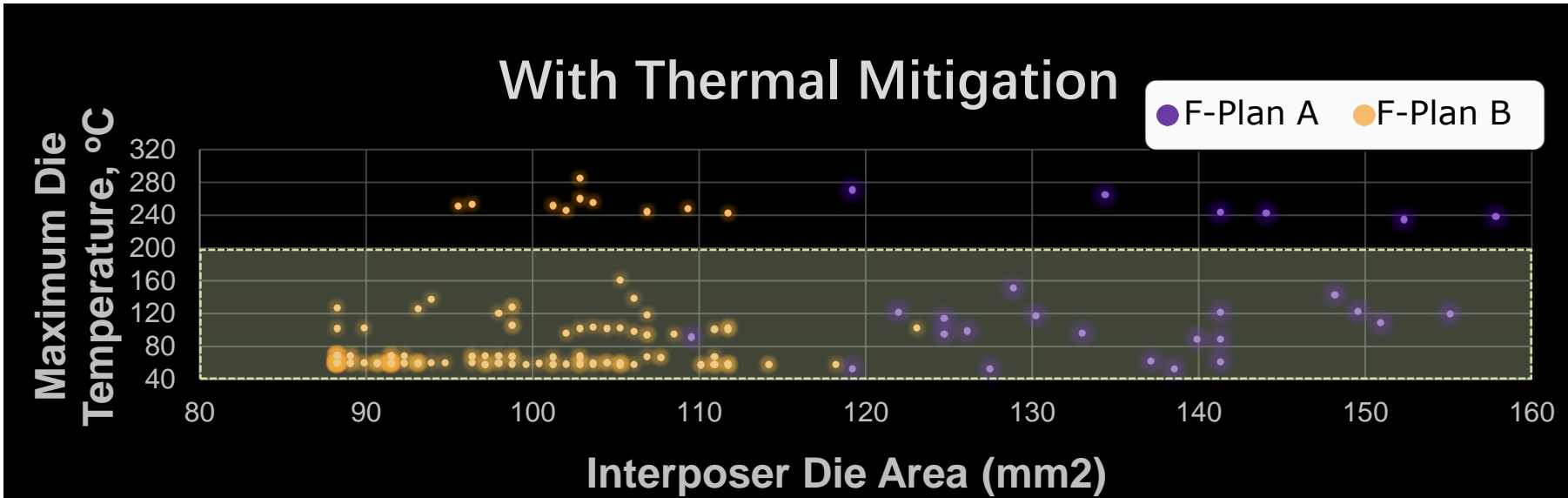
Thermal

- Thermal Insulation Material Index
- Heat Spreader
- Molding Compound
- Type of Cooling

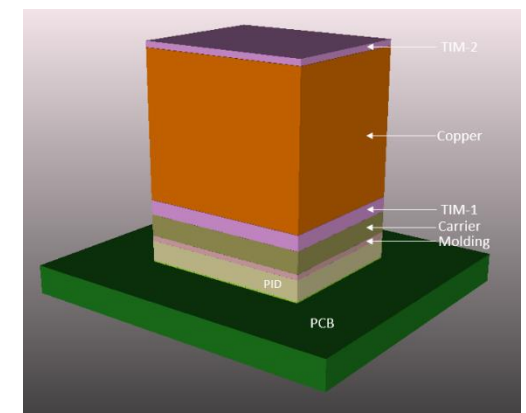
Architecture & Thermal Impact

Thermal

- Thermal Insulation Material Index
- Heat Spreader
- Molding Compound
- Type of Cooling

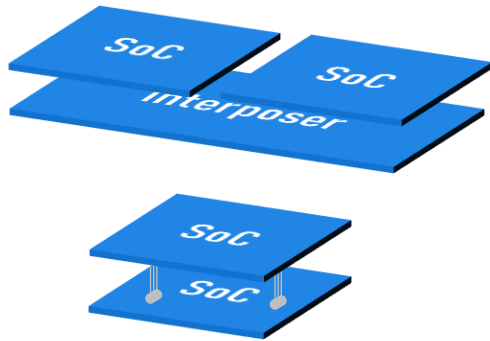


TechOnomics

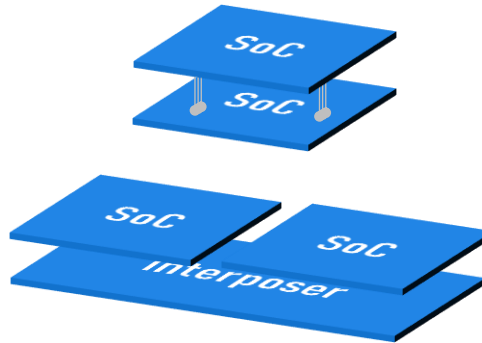


Implementation For 3DIC's

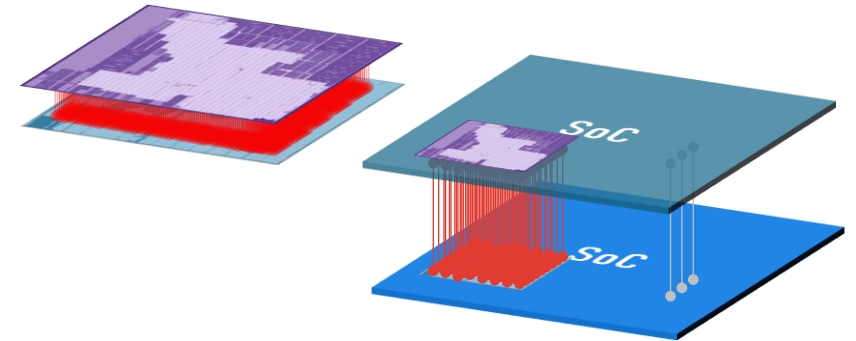
Coarse Connections
Source Synchronous



Coarse Connections
Clock Synchronous



Dense Connections
Clock Synchronous



Stack of coarsely bounded die
(protocol communication)

- Planned in 2.5D or 3D
- Implemented as multiple, 2D die
- Source synchronous D2D communication (UCIe, HBM)

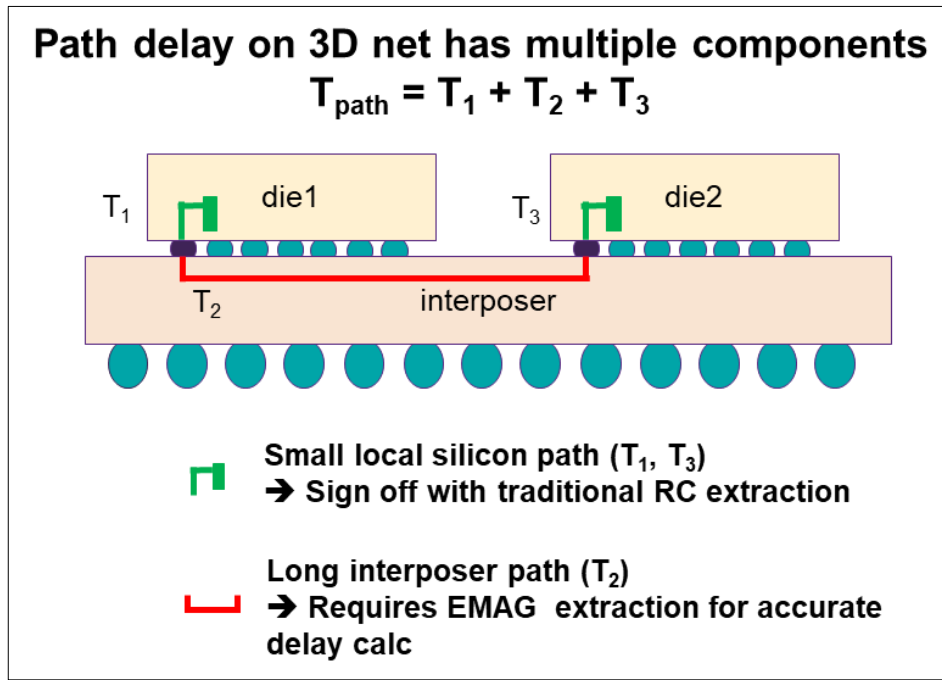
Stack of coarsely bounded die
(synch clk)

- Planned in 3D, Implemented as multiple 2D die
- Synchronous D2D communication
 - D2D optimization, clock...

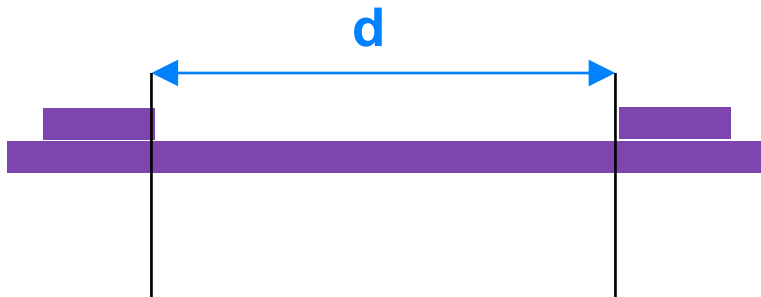
“Single chip” on multi-die

- Planned and implemented in 3D
- Similar flow as 2D, mapped to multiple die
- Synchronous D2D communication

Protocol → Clock Synchronous: When?

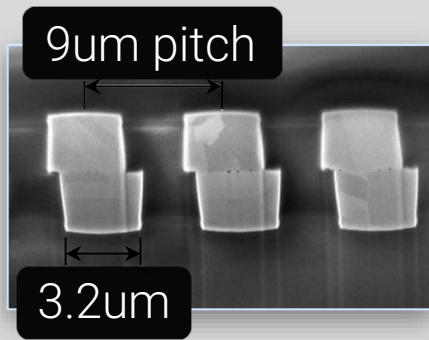


Source: AMD

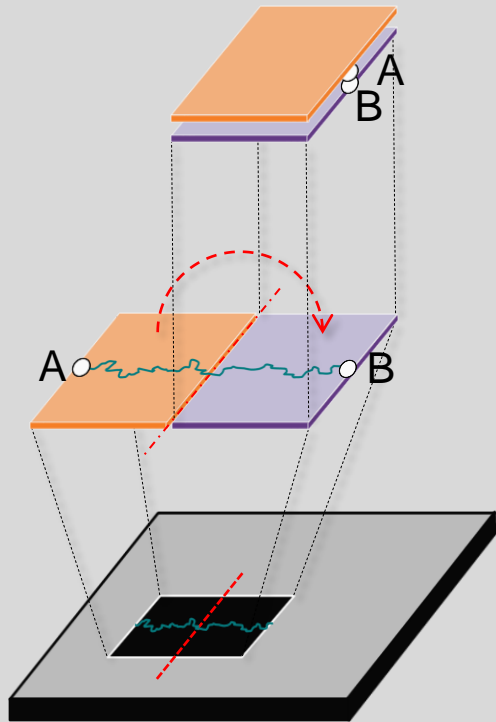


- Protocols have overheads
 - Logic
 - Latencies
 - Power
- Necessity for PCB due to distances
- At what D2D distance is it feasible to employ synchronous clock
 - Standard RC extractions for digital implementation not suitable for long nets
 - Static timing net delay calculation using full-wave Electromagnetic solver
- What happens at the limit where $d \rightarrow 0$?

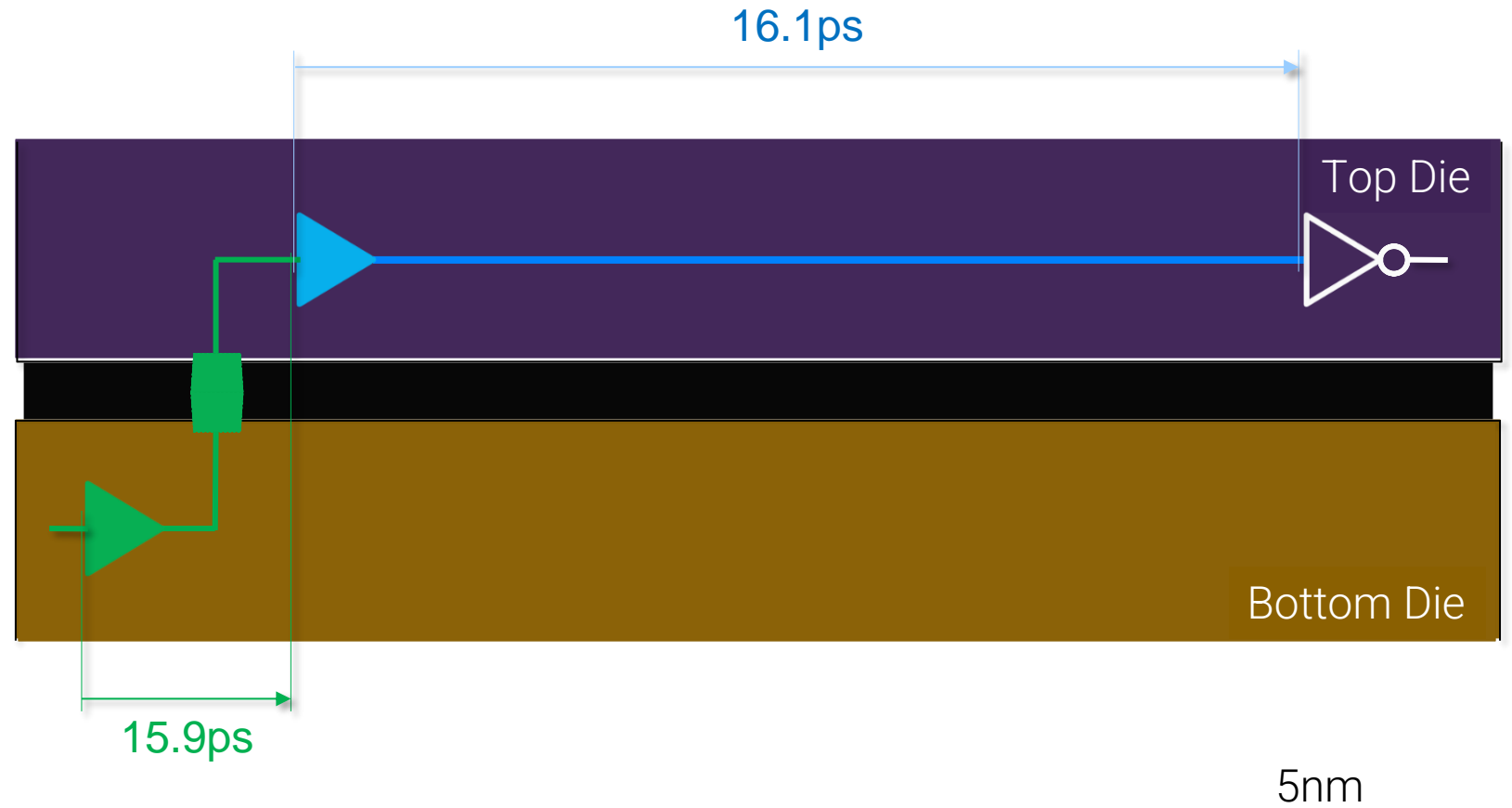
Hybrid Bonding



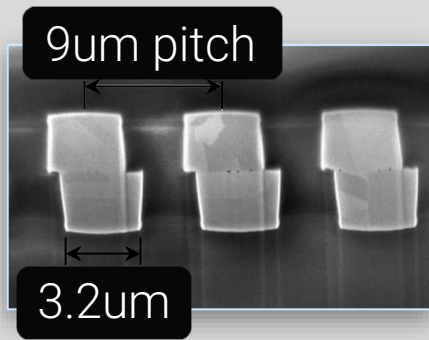
Die/Core Folding?



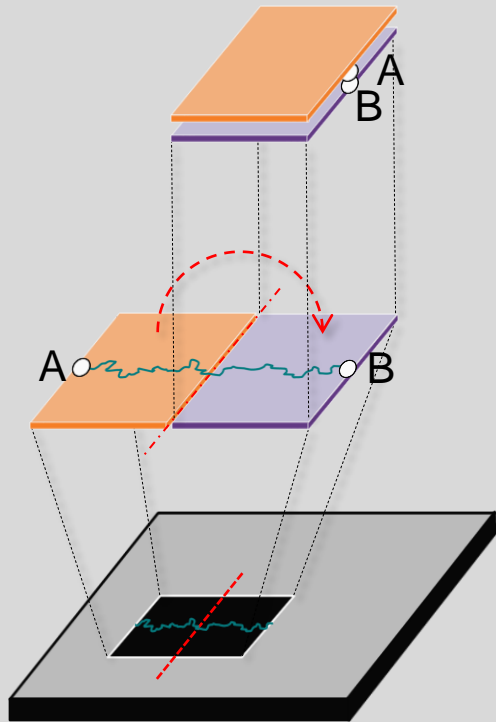
Shortest Delay: 'Across Chip?' or 'Between Chips?'



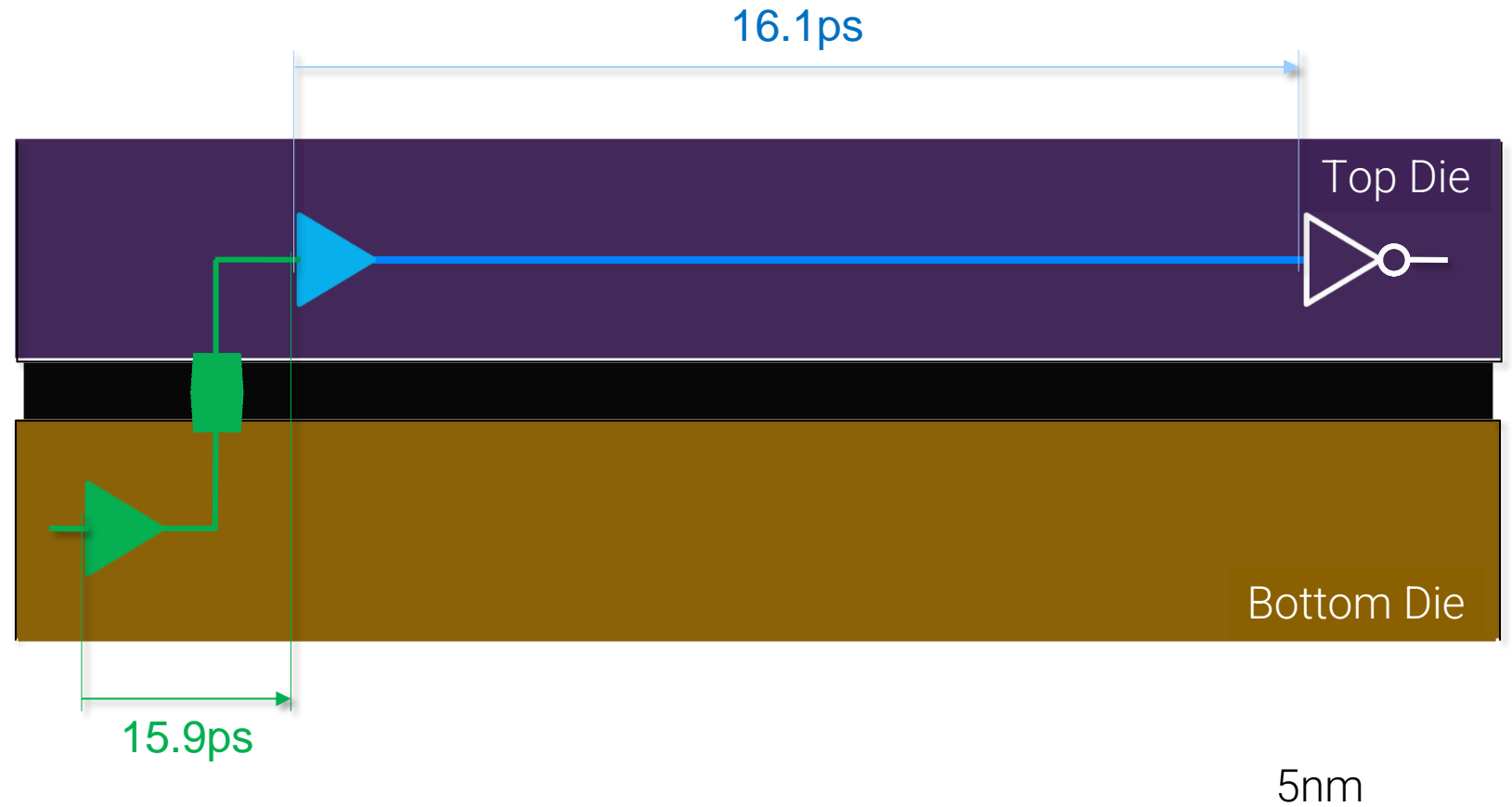
Hybrid Bonding



Die/Core Folding?



Shortest Delay: 'Across Chip?' or 'Between Chips?'



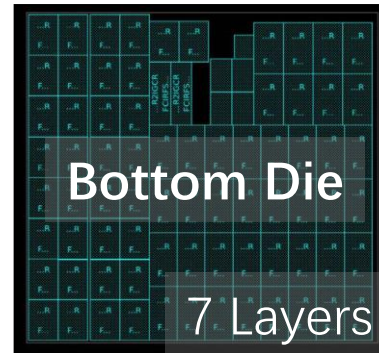
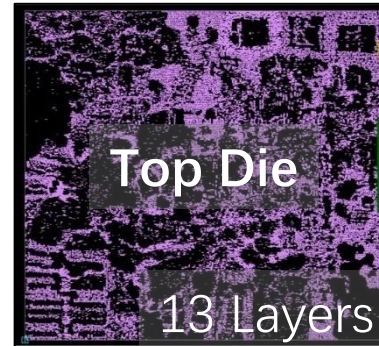
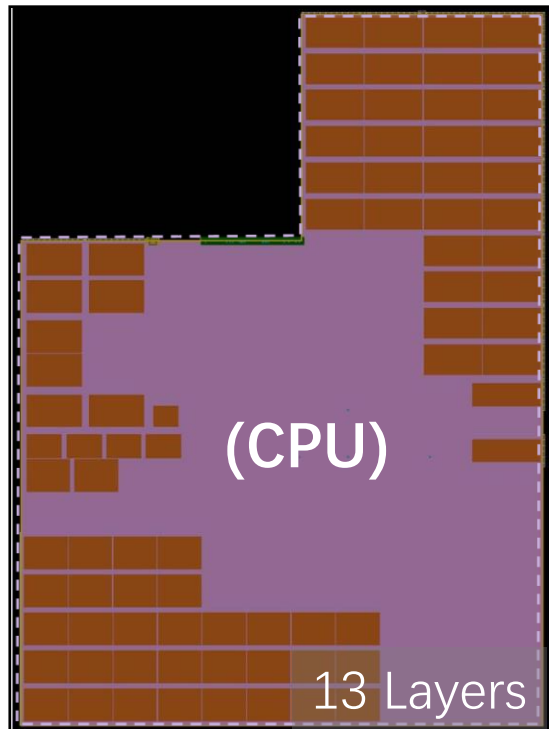
Hybrid Bond Pitch:

2023	2024	2025	2026
9um	6um	6um	4.5um

2D “Single Chip”



3D “Chiplets”



CPU Disaggregation

- ❑ Core Area: - 52%
- ❑ Frequency: + 2.8%
- ❑ TNS: - 19%
- ❑ Wire Length: - 2.5%

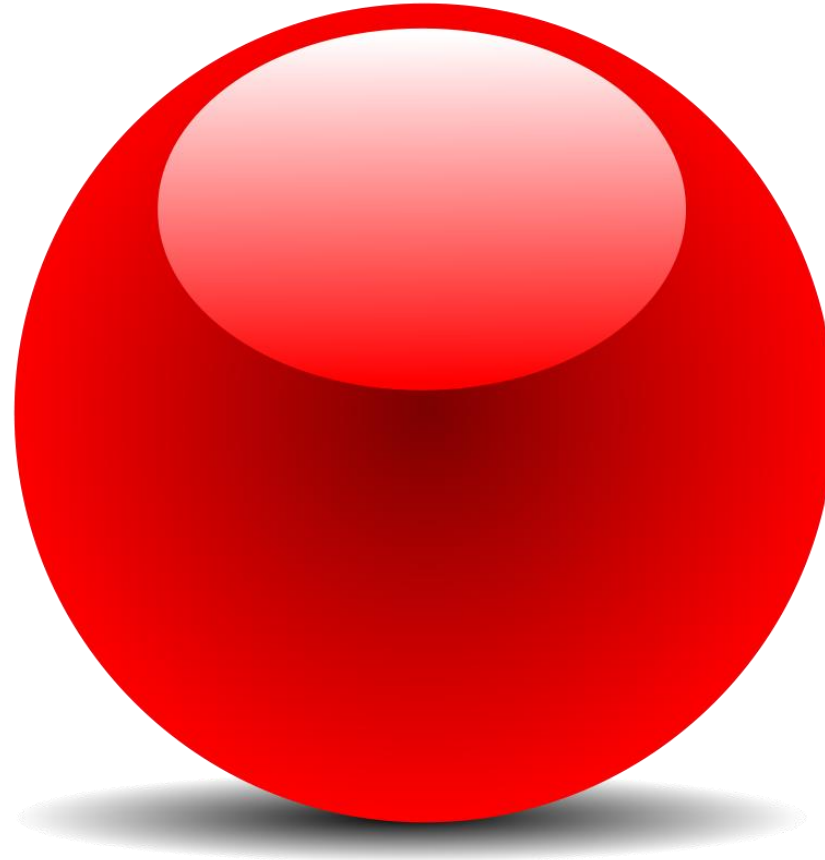
The Modern Bond



Hybrid Bond

Pitch 1 μ m, Density 1,000,000/mm²

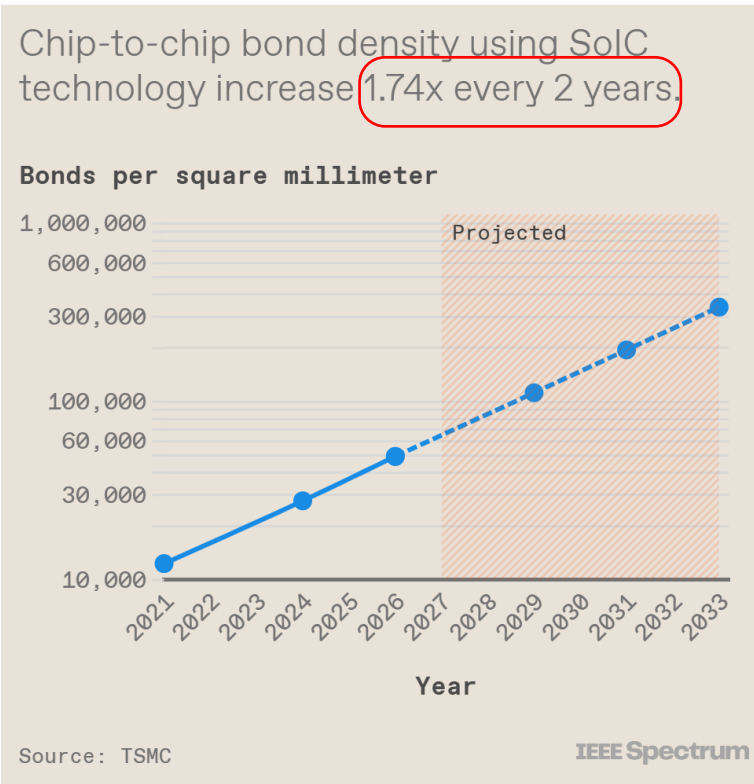
Kagawa et al, "Impacts of Misalignment on 1 μ m Pitch Cu-Cu Hybrid Bonding" IITC 2020



Classic Package Bump

Pitch 100 μ m, Density 100/mm²

Bond/Bump Pitch Scaling



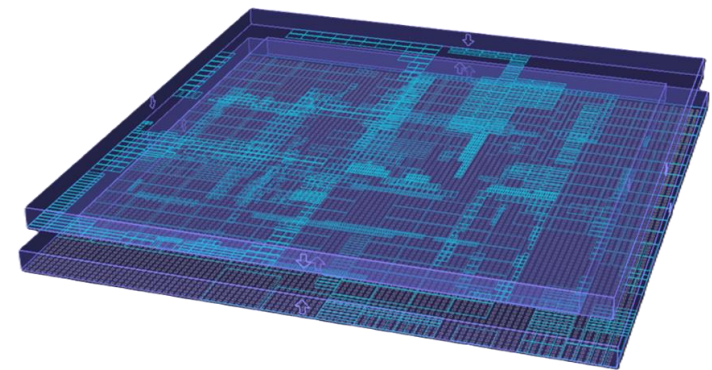
Mark Liu, Philip Wong, How We'll Reach a 1 Trillion Transistor GPU, IEEE Spectrum

- Inter-chiplet connections primarily done manually
- Bond densities growing at exponential rate – $O(1.32^n)$
- Reticle limits continue to increase: $3.3X \rightarrow 5.5X \rightarrow 8X \rightarrow 50-60X$
- Number of human fingers $O(1)$

	Bump Pitch	Relative Density
BGA	300u	1X
C4	140u	4.6X
ubump	40u	56X
Hybrid Bond	9u	1111X
Hybrid Bond+	1u	90,000X
Hybrid Bond++	0.1u	9,000,000X

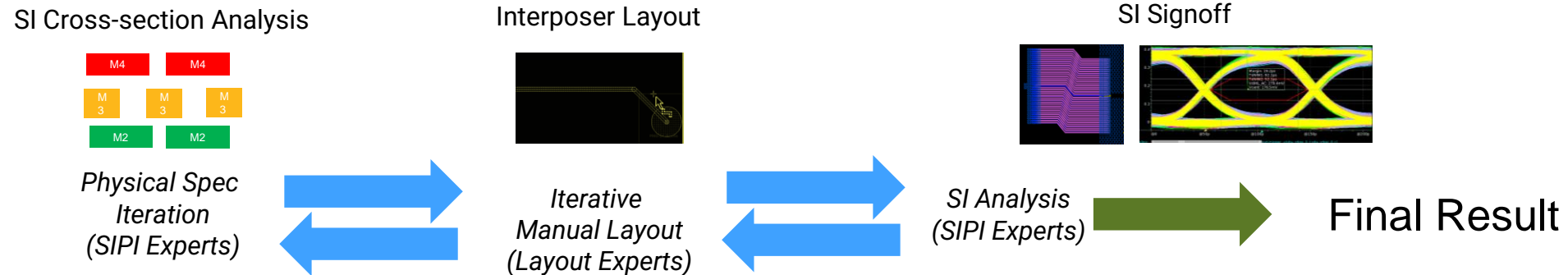
	HBM3/3e	HBM4
Interface Width	1024 bits	2048 bits

Source: Synopsys



New methods for 3D Bump Planning

Die-to-Die Interconnect Design

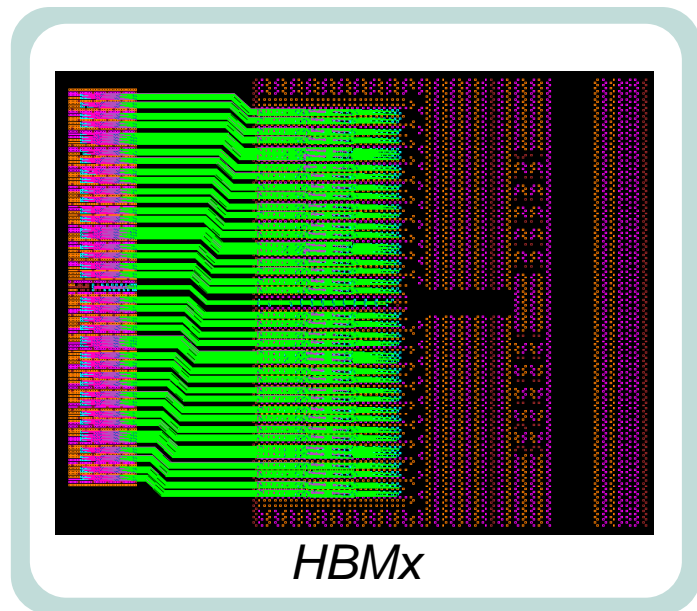


- “Spec-based” Methodology, ~6 months for a full channel implementation
- Routing typically manual, 3-4 weeks for a single implementation
- Real projects usually have multiple major design changes
- Cross-section Analysis: SI-based physical spec from simulating net patterns
- Interposer Layout: Based on SI spec and rules of thumb from PCB design.
- Tedious process, increasingly difficult to converge to spec, and sub-optimal

Routing Automation

- Types: High-speed (HBMx/UCIe/D2D), low-speed, perf/length matching
- High-speed chiplet routing is very different than general EDA routing
 - Almost 100% occupancy
 - No jogs – minimum wirelength and skew across bits
 - Chiplet offsets mandate ‘turns’ in D2D connections (45-degree, 90-degree)
 - Heterogeneity (source, interposer and destination are all different tech)
 - 100% shielding
- Implications on algorithmic solution
 - Limited incrementality
 - Heterogeneous constraints corresponding to heterogeneous composition
- Radical change in methodology from manual design entry
 - Manual design has high dependence on rules-of-thumb

HBMx Channel Routing



	Channel Width
HBM3/3e	1024 bits
HBM4	2048 bits

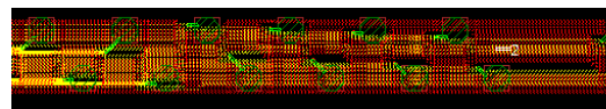
Increasing bitwidths with HBM4

- Synopsys 3DIC Compiler D2D Routing Pattern

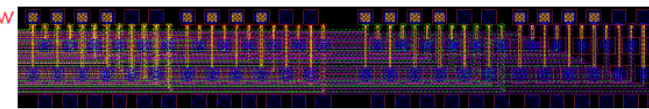
Source: GUC

- Work at 7.2Gbps through SI simulation
- Shorten channel signal routing runtime by 50%
- Easy for channel routing balance

Origin



New

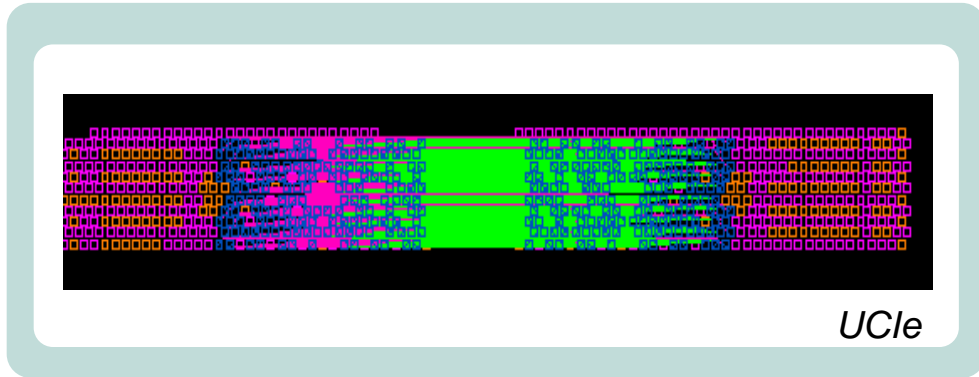


HBM3 Interposer	6% width increase StatEye opening 8% average timing margin improvement
-----------------	---

HBM3 Design	Custom Layout	Automated Routing	Compare Results
Eye Width@±60mV	59.8 ps	63.6 ps	6% better , for worst channel
Crosstalk (NEXT)	-30.9 dB	-29.6 dB	4% worse , still at good level, not hindering performance
Insertion Loss (IL)	-4.74 dB	-4.66 dB	2% better , less attenuation
Intra Channel Skew	6.5 ps	5.5 ps	15% better , less skew, worst case

Source: Synopsys

UCIe Channel Routing



Includes routing of signal (data, control, differential routing) and electromagnetic PG network (Shield, PG distribution...)

Virtually full occupancy

Absolute minimum wirelength/via

SI Signoff Results

UCIe IP Interposer Testchip: 17% better signal performance

Parameters	Spec	Custom [dB]	3DICC [dB]
IL	> -3db	-2.11	-1.96
PSFEXT	< -23db	-23.7	-23.91

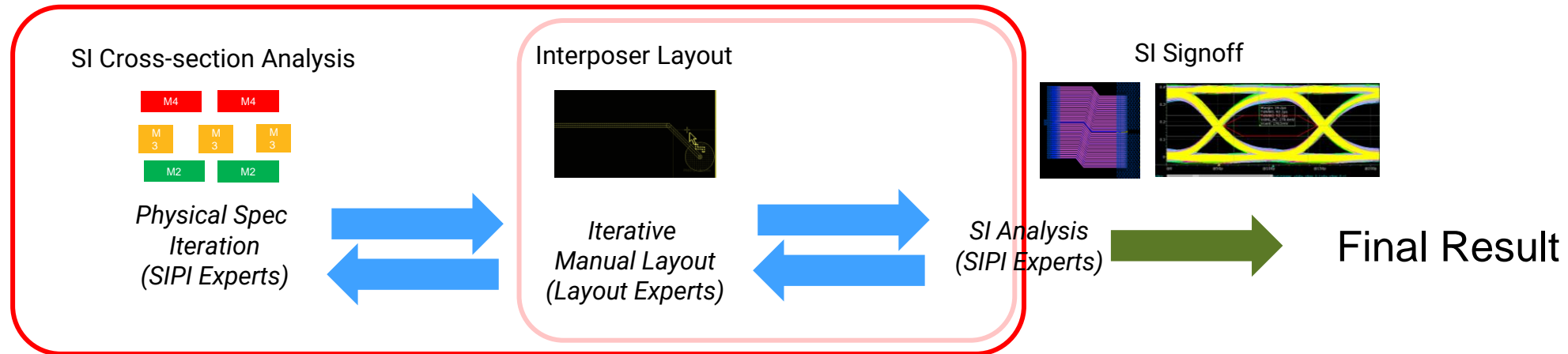
Source: Synopsys

UCIe SLM Test Vehicle Tapeout:

Parameters	Spec [ps]	Typical [ps]	Slow [ps]
ISI eye opening	50 (0.8UI)	59.1 (0.946UI)	58.0 (0.928UI)
ISI jitter	8.75	3.4	4.5
Xtalk eye open	50 (0.8UI)	52.3 (0.837UI)	51.7 (0.821UI)
Xtalk jitter	6.88	6.8	6.3
Total jitter[ps]	15.62	10.2	10.8
Total jitter[UI]	0.25UI	0.163UI	0.173UI

S. Kabir, TSMC OIP 2024

Going from Automation to Optimization



- Move from “Spec-based Layout” to “Analysis-Driven Interconnect Optimization”
 - Implementation and Analysis capabilities in same environment
- Enormous state space for optimality, typically highly simplified in practice
- Well-suited for Machine Learning

Results of Emag Optimization

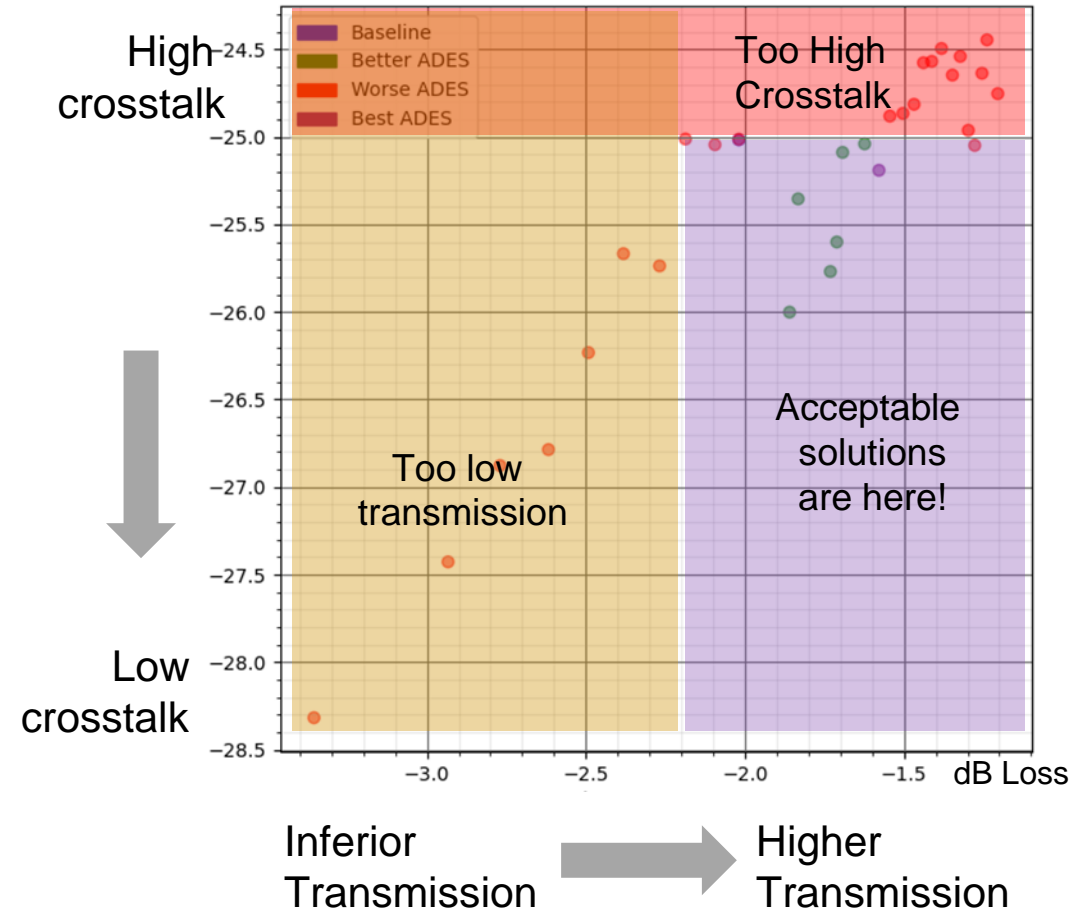
AI Machine Learning Based Optimization

Transmission Loss

	Baseline	Optimized	Improvement
Design 1 (HBM3)	-2.58 dB	-2.02 dB	0.56dB (-17%)
Design 2 (HBM3)	-2.11 dB	-1.88 dB	0.23dB (-8.7%)
Design 3 (UCIe)	-0.82 dB	-0.81 dB	0.01dB (-1.1%)
Design 4 (HBM3)	-1.27 dB	-1.21 dB	0.06dB (-4.1%)
Design 5 (UCIe)	-0.88 dB	-0.77 dB	0.11dB (-11.4%)

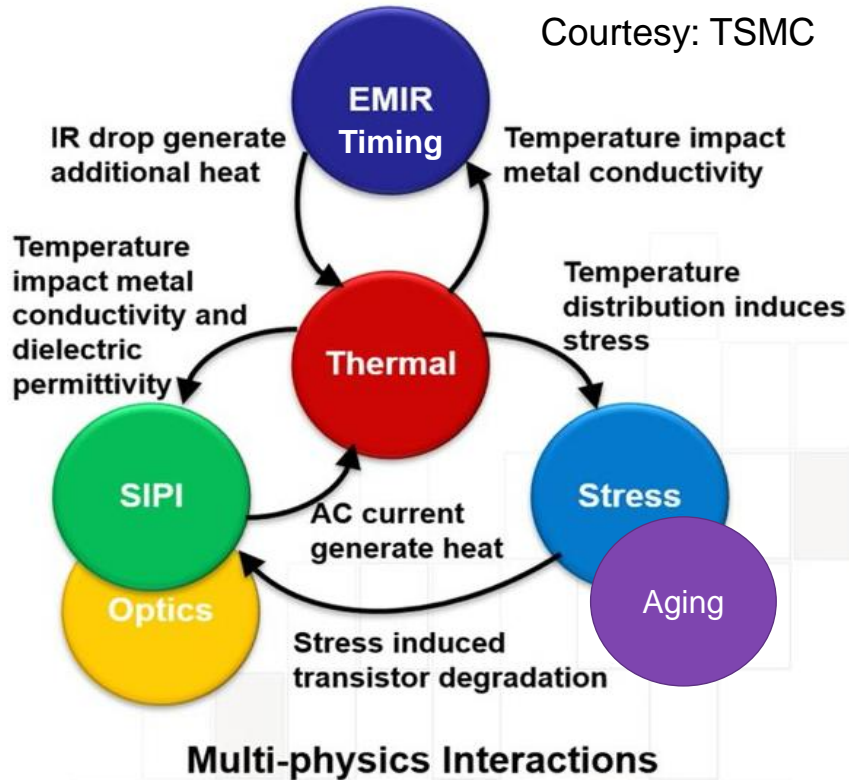
Source: Synopsys

Pareto Design Exploration

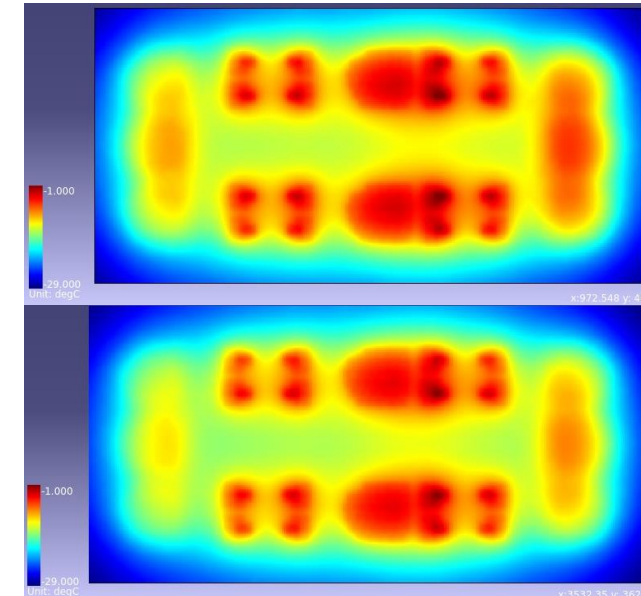
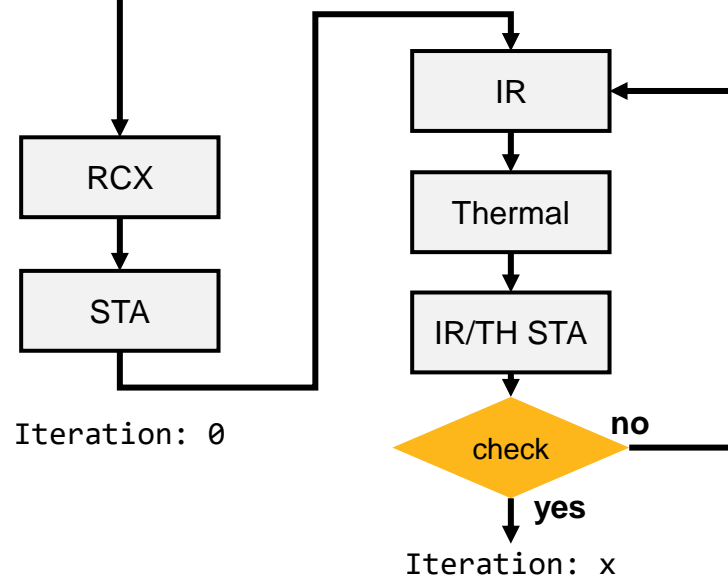


Multiphysics Analysis

Courtesy: TSMC



analyze_3d_multiphysics



Source: Ansys/Synopsys TSMC OIP NA 2024

- Accurate system analysis requires converging across multiple interacting effects
- Results impact both 2D and 3D design

Looking Forward...

- Design vs. Implementation: Different problem statements, user profiles
- Automation: Required for 3DIC's to go to scale
- Optimality: Goals of SI, PI, Thermal, EMIR, Area/Function. AI plays a role.
- Formalization: Cannot automate if you cannot codify the models, constraints and objectives a priori
 - vs. “I know it when I see it”
- Scaling: Abstractions, models, techniques. Bumps/bonds become via's or feedthroughs.
- Multiphysics Signoff – True multi-physics coming into the picture now – importance for 2D, 3D

SYNOPSYS®