

Future Architecture Demands for More Aggressive Packaging



intel®

Josh Fryman, PhD
Office of the CTO, Intel Fellow
Feb 23, 2024



HETEROGENEOUS
INTEGRATION ROADMAP



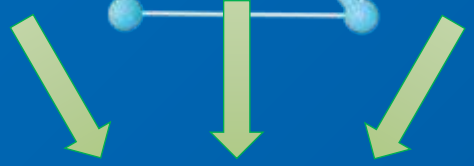
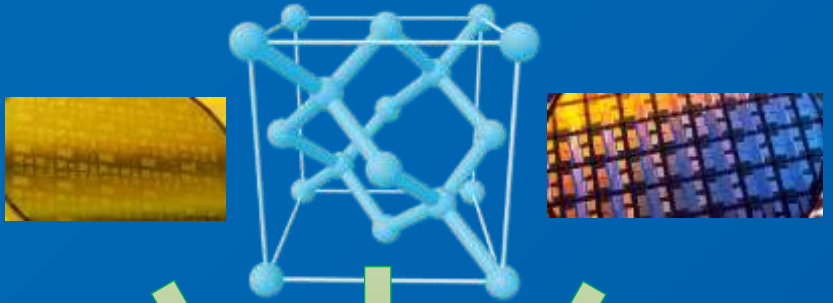
Moore's "Page 3" refocused upon by DARPA's ERI programs:

"The total cost of making a particular system function must be minimized. To do so, we could amortize the engineering over several identical items, or evolve flexible techniques for the engineering of large functions so that no disproportionate expense need be borne by a particular array.

[..]

It may prove to be more economical to build large systems out of smaller functions, which are separately **packaged and** interconnected. The availability of large functions, combined with functional design and construction, should allow the manufacturer of large systems to design and construct a considerable variety of equipment both rapidly and economically."

– *Gordon Moore, Electronics, No. 38, Vol. 8, April 19, 1965*



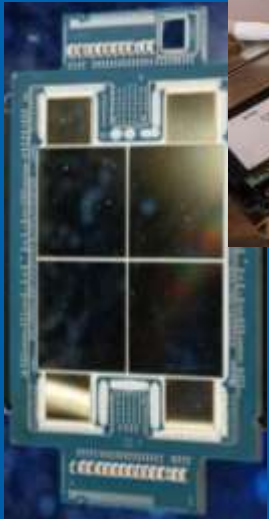
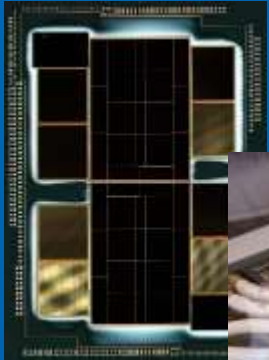
"... there is no spoon."

(There is no box... package... limit...)



HETEROGENEOUS
INTEGRATION ROADMAP

It's about systems, not just ingredients





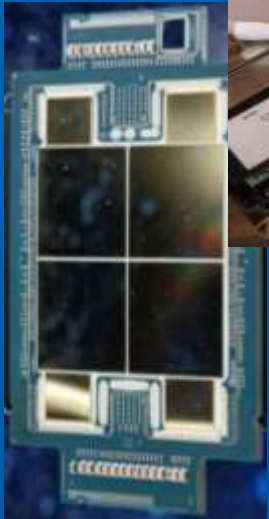
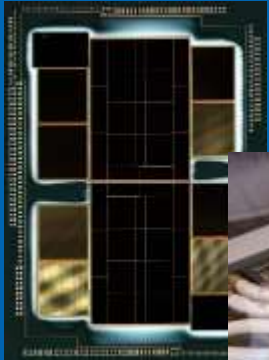
It's about systems, not just ingredients



Total Cost of Ownership →
System Technology Co-Optimization

AI, HPC, Mobile, Medical, AR/VR, Contact Lens, . . .

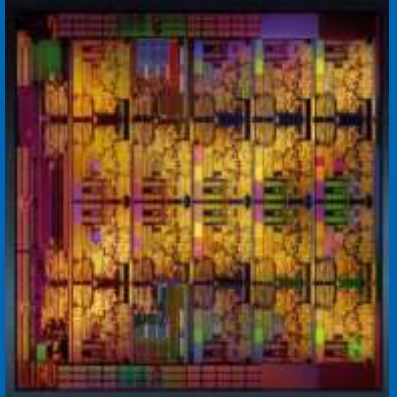
By 2030 scaling up to 50 TB/s DRAM BW,
20 TB/s IO BW, and 100 AI PFLOPS as a single module!



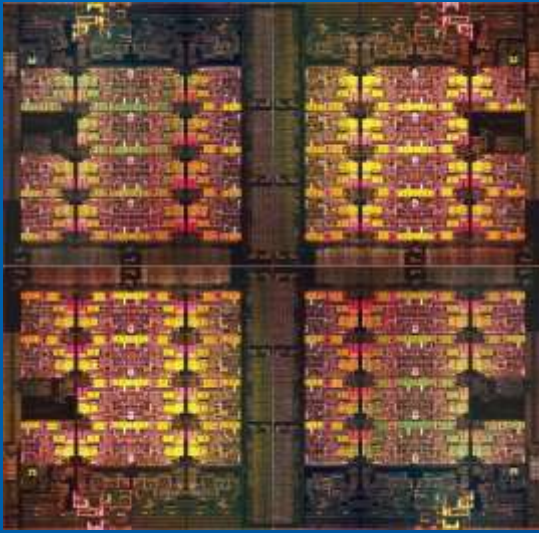


HETEROGENEOUS
INTEGRATION ROADMAP

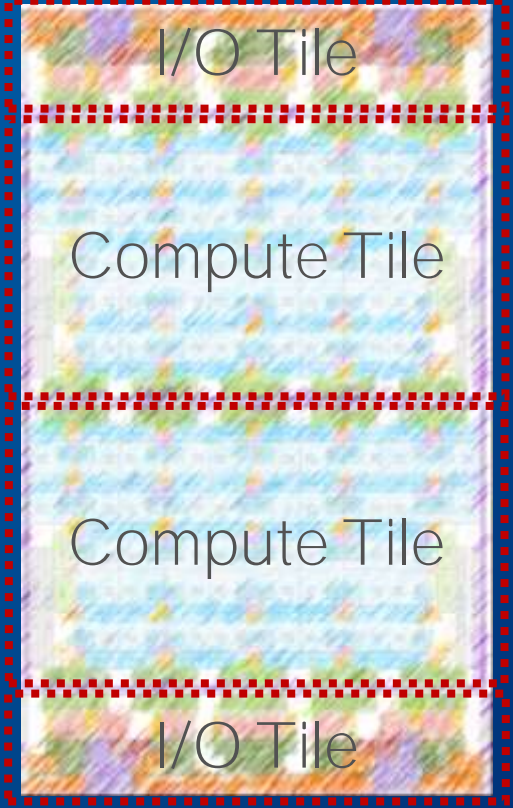
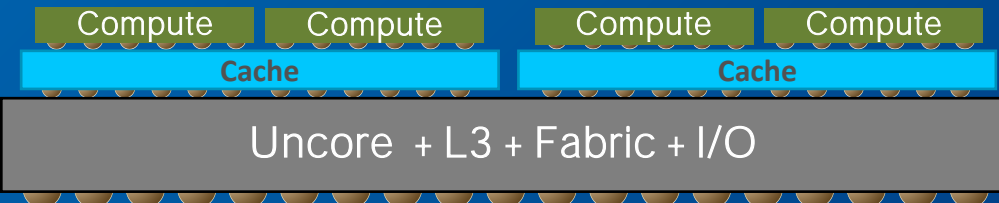
Are we really going all-in?



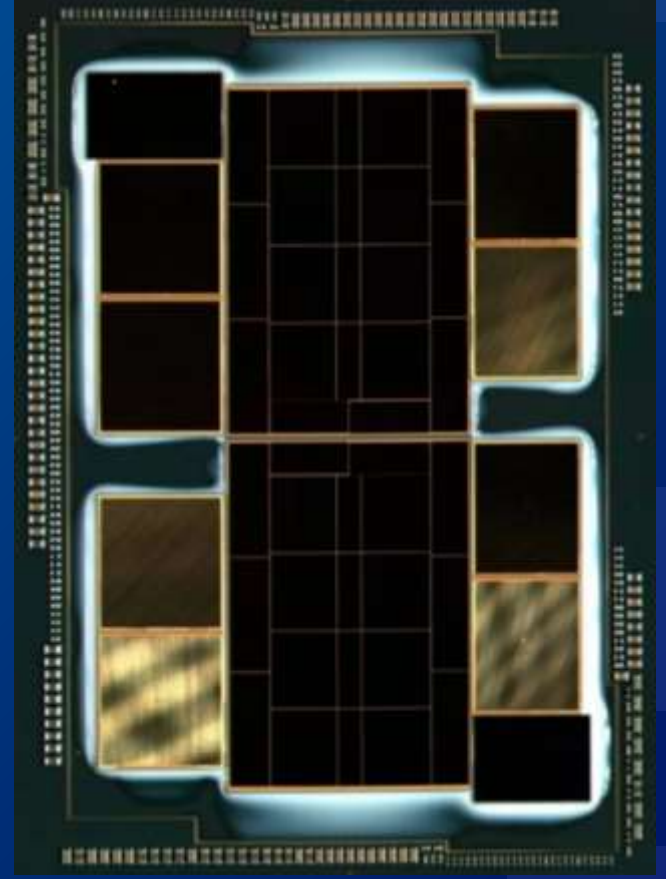
Monolithic



Polyolithic
multi-chip

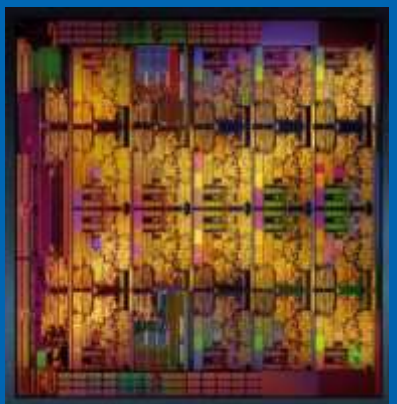


Polyolithic tiled



Heterogenous
polyolithic tiled

Are we really going all-in?



Monolithic

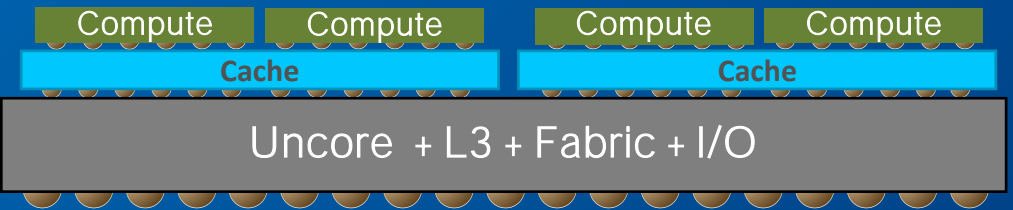


Multi-chip



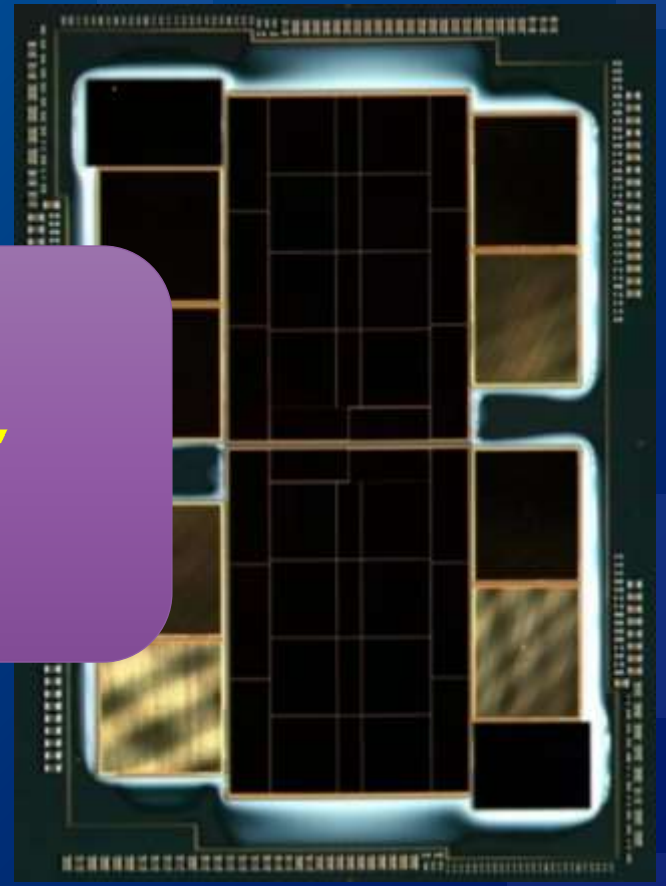
I/O Tile

Whether it's 2.5D or 3D in name ...
it's all the same old "planar thinking"
Is that bad?



I/O Tile

Polyolithic tiled



Heterogenous
polyolithic tiled

What about big-opportunity thinking?



"We need to think about packaging 60,000 mm² for processing all those LLMs ..."

– UCLA Prof. Subramanian Iyer, Director of NAPMP

Assume the equivalent of ~80 SotA GPUs at ~800 mm² each

sq mm	X,Y dims	pJ/b wires	pJ/b stops	pJ/b d2d	W intern	W extern	Notes
64,000	253	1,265	316	5	46	-	3 GHz, 12 TB/s mesh, 10% activity, no external switches, only counting fabrics
800	28	141	35	1	5	36	

If that 80-GPU part prices at \$2M, is that okay?

# units	W fabrics total	OpEx \$M annually	OpEx 5 yr life \$M
1	46	0.16	0.82
80	3,287	11.83	59.17



What about big-opp for optimization?



"V...
pro

2 for

Grossly simplistic model ...
but \$58M is a critical
opportunity for optimization!

Or is it?

(Coherency, Substrates,
Memory, ...)

at ~800 mm² each

sqm
64,000
800

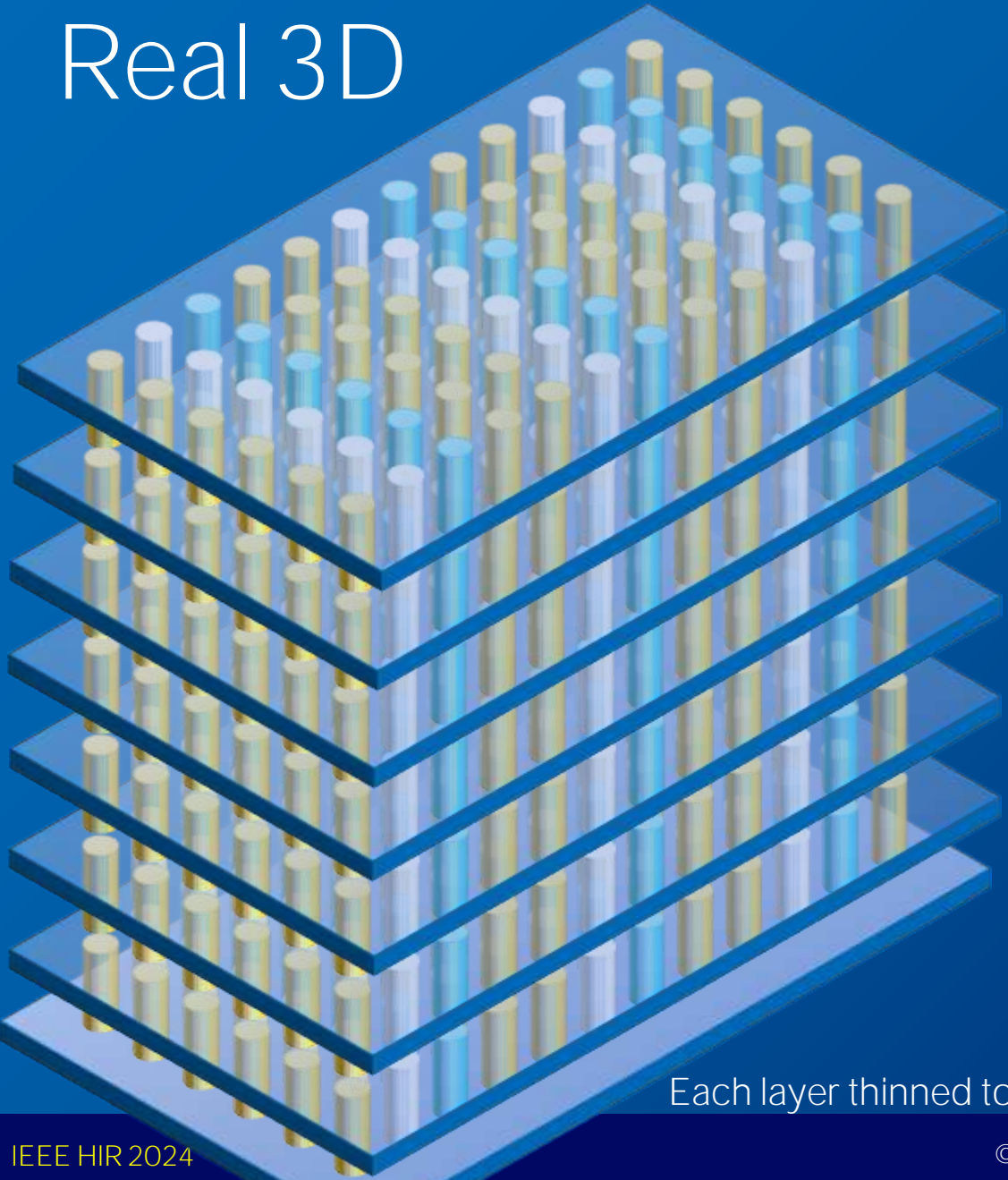
GHz, 12 TB/s, 90% activity, no external switches, only counting fabrics

... \$58M, is that okay?

# units	fabrics	OpEx \$M annually	OpEx 5 yr life \$M
46		0.16	0.82
8	3,287	11.83	59.17



Real 3D



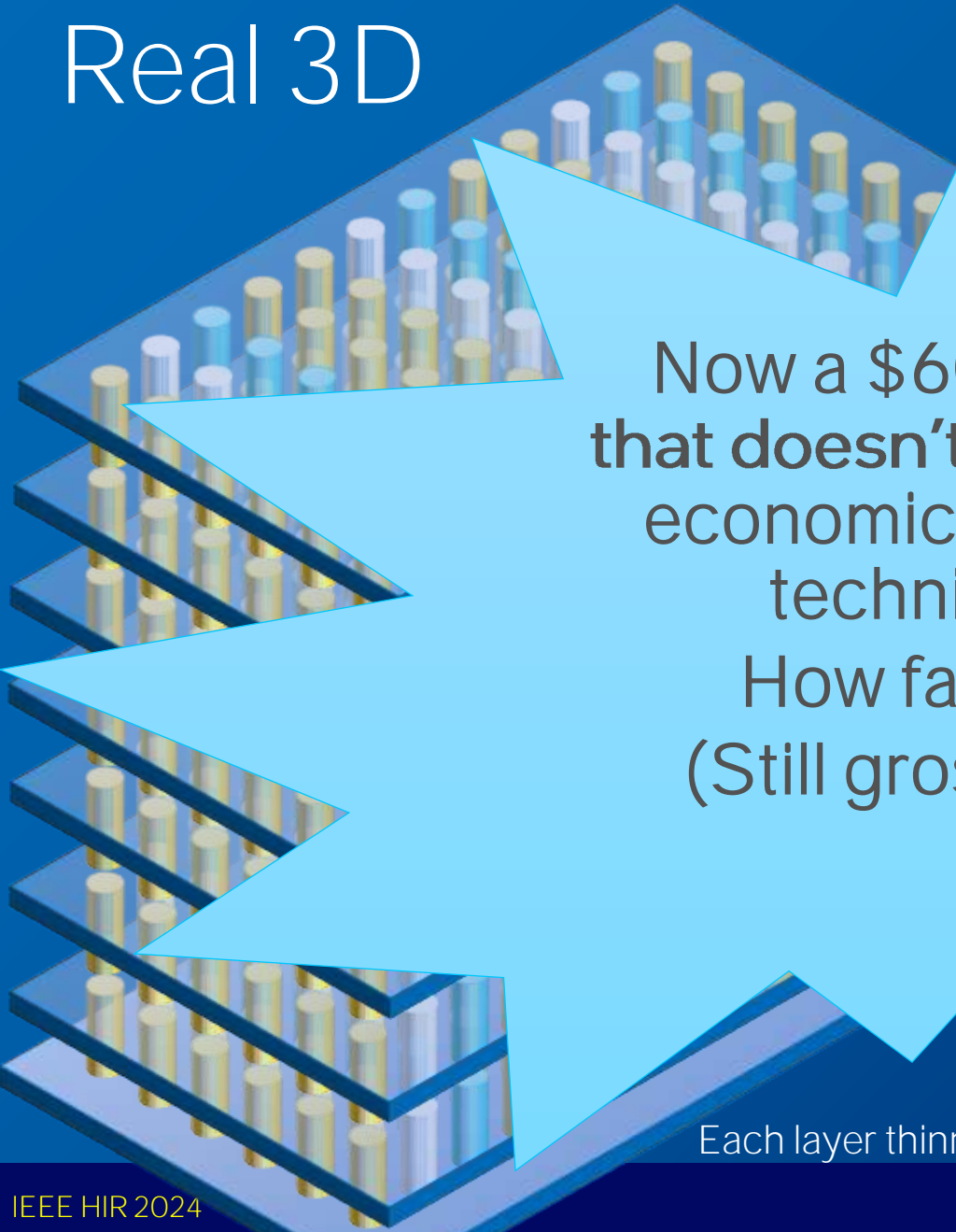
layers	sq mm	X,Y dims	pJ/b walk	pJ/b stops	pJ/b d2d	W intern	W extern
1	64,000	253	1,265	316	5	46	-
1	800	28	141	35	0	5	36
10	64,000	80	400	100	1	14	-
100	64,000	25	126	32	2	5	-
800	64,000	9	45	11	3	2	-

# units	W fabrics	\$M annually	5 yr life \$M
1	46	0.16	0.82
80	3,287	11.83	59.17
1	14	0.05	0.26
1	5	0.02	0.08
1	2	0.01	0.03

Each layer thinned to ~10um



Real 3D



Now a \$60M opportunity that doesn't risk fundamental economic hurdles ... only technical hurdles.
 How far can we go?
 (Still grossly simplistic)

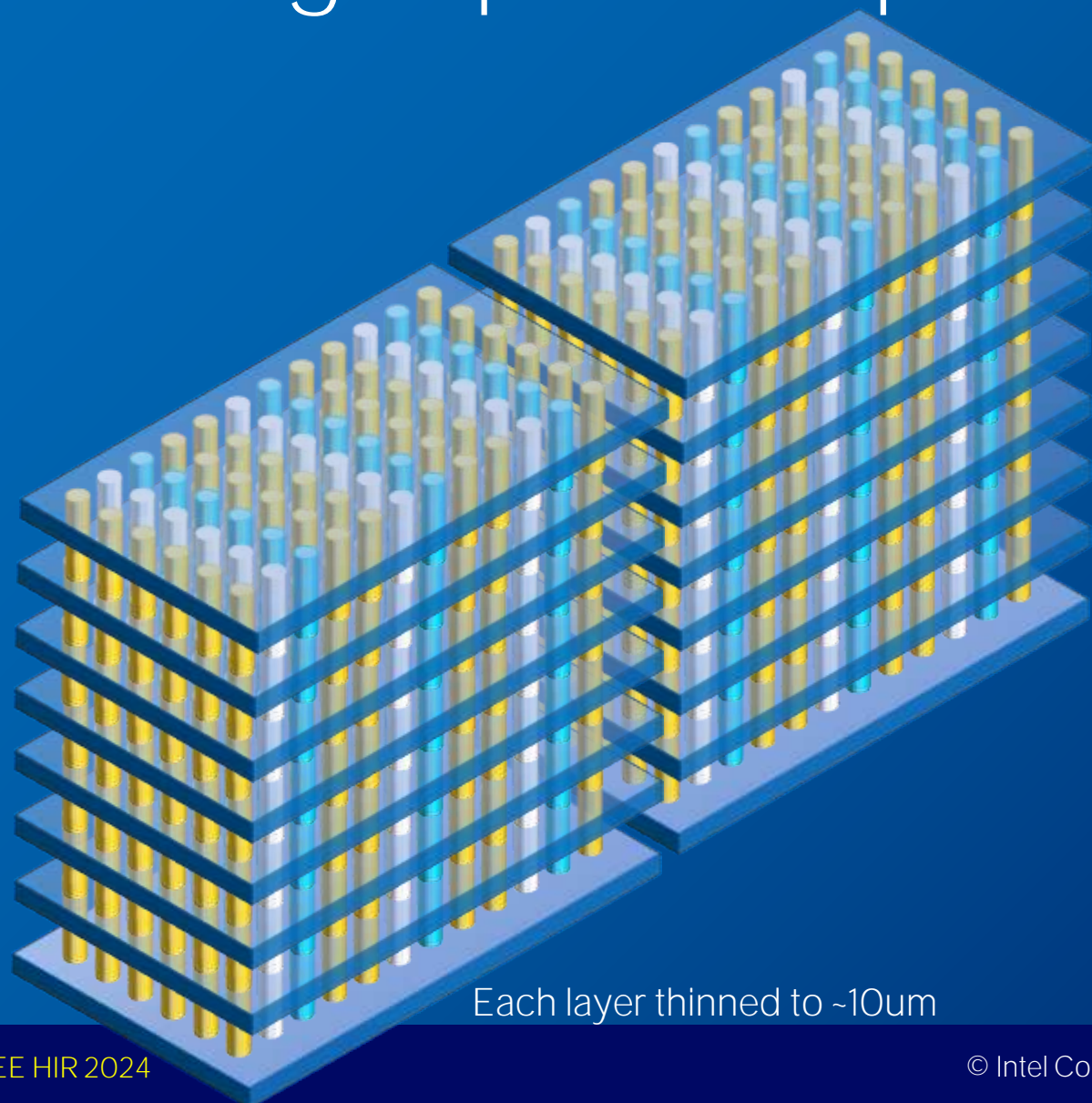
ms	pJ/b	stops	pJ/b d2d	W intern	W extern
216	5			46	-
				5	36
30	1			14	-
				5	-
				3	2

ics	\$/m	ly	5 yr life \$M
46		0.16	0.82
80		11.83	59.17
1		0.05	0.26
1		0.02	0.08
1		0.01	0.03

Each layer thinned to ~10um



Package sprawl vs. package scrapers

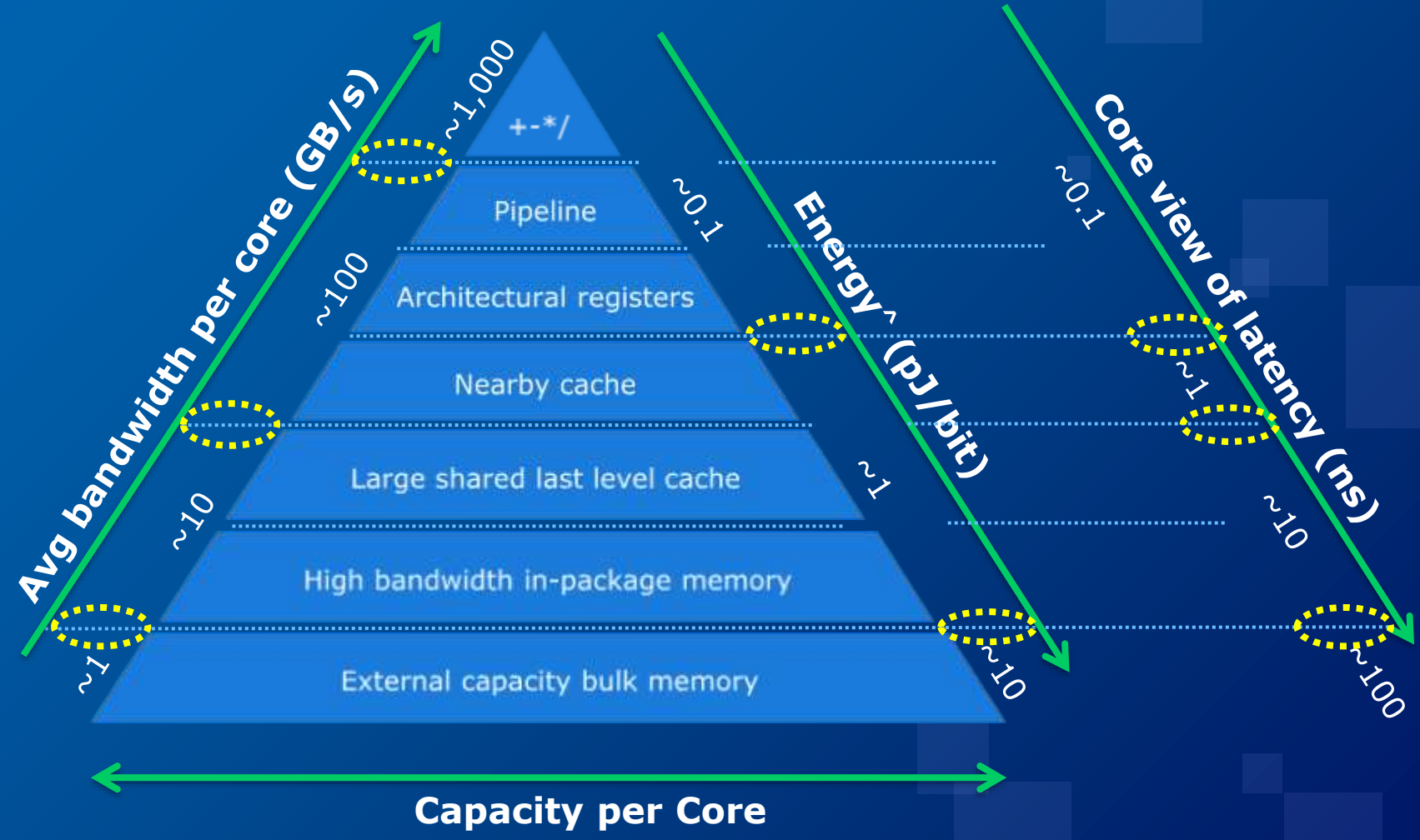


Each layer thinned to ~10um

- Assume ~8mm x ~8mm die size
 - 800 layers makes it ~8mm tall
- HBI-TSV density drives bandwidth by pitch
 - 9um \rightarrow 12k/mm² \rightarrow 800K total
 - 3um \rightarrow 111k/mm² \rightarrow 7.2M total
- Assume ~67% for power/gnd, ~33% for IO
 - ~0.27-2.4M IO signals per 1 GHz
 - ~0.260 – 2.4 Pbps total bidir BW per GHz
- 16 TB/s/dir – 150 TB/s/dir IO capable per GHz
 - Lots of floorplanning, thermal concerns
 - Future delivery of ≤ 1 um HBI-TSV resolves wiring challenges
 - Bottom die PHY area limits IO sustainable
- Some open hurdles
 - Time on tools
 - EDA, DFX, DFT, FA, etc.
 - Lack of lateral connectivity between high-rises

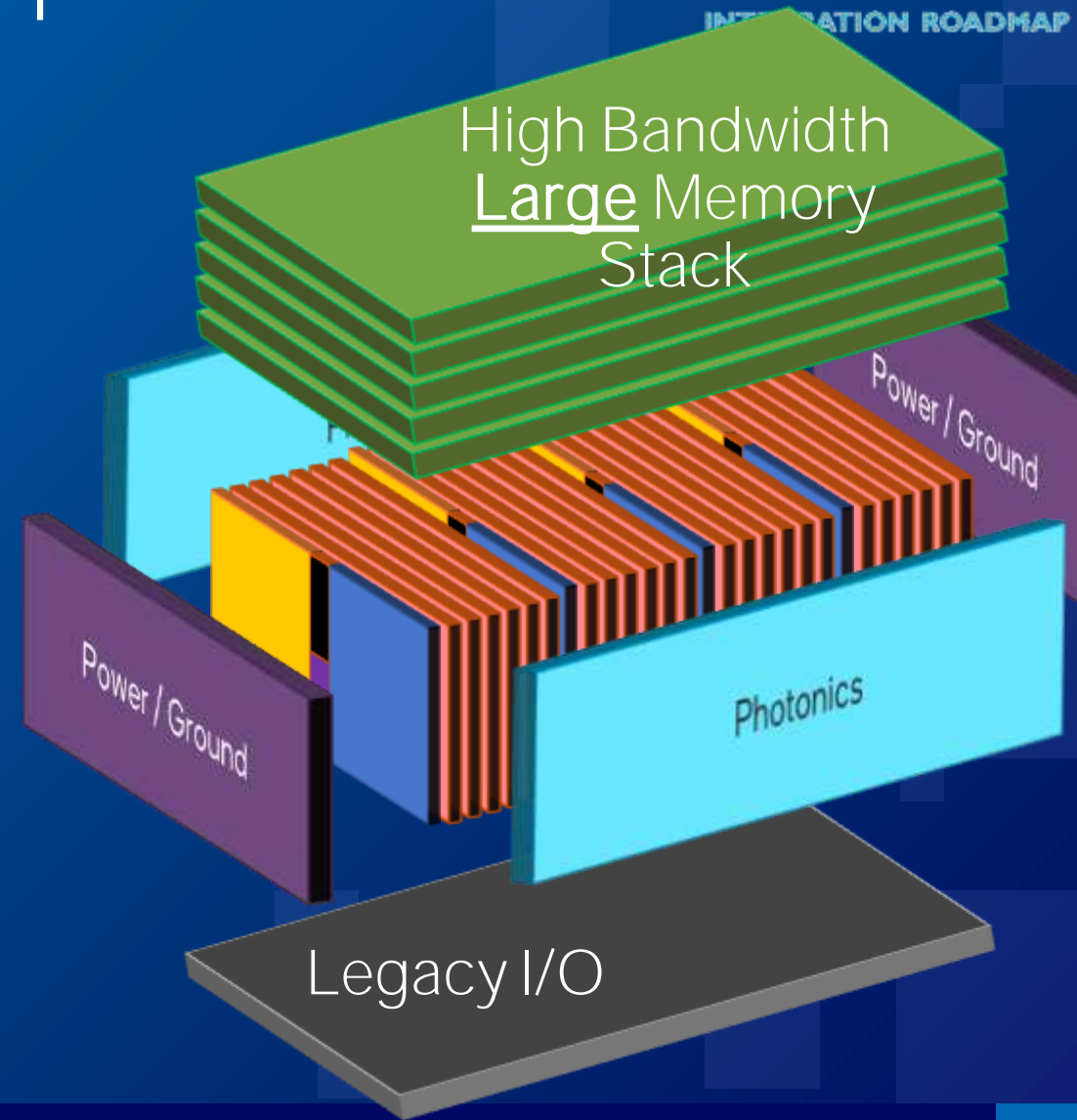
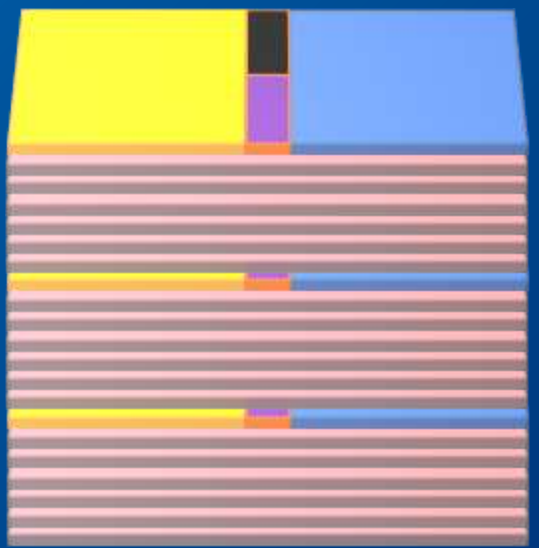
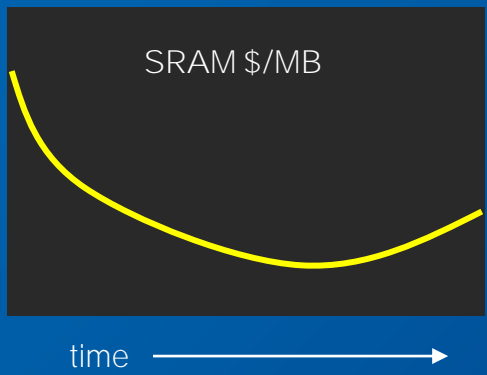
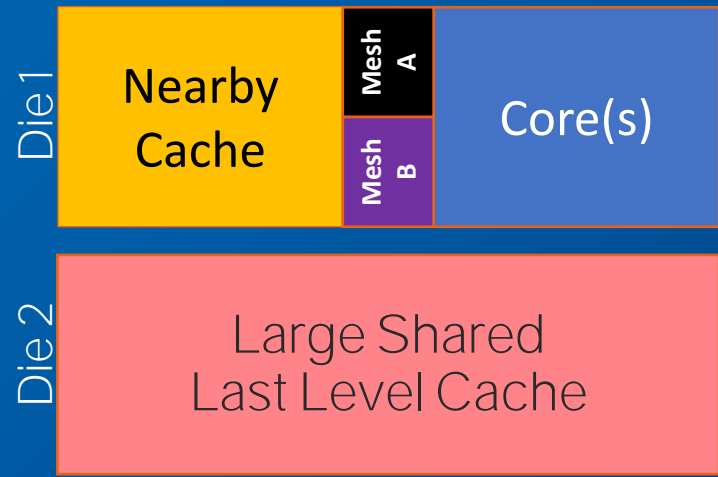
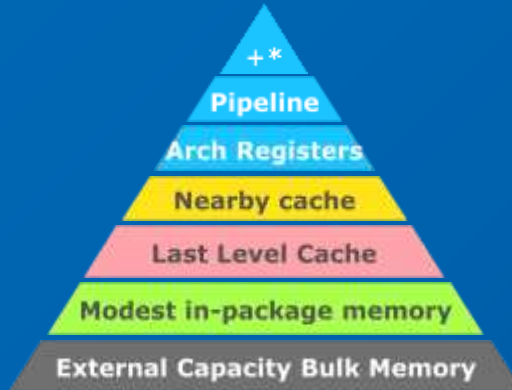
But are we exploiting the right things?

Tapering divisions in bandwidth, energy, and latency are opportunities to optimize

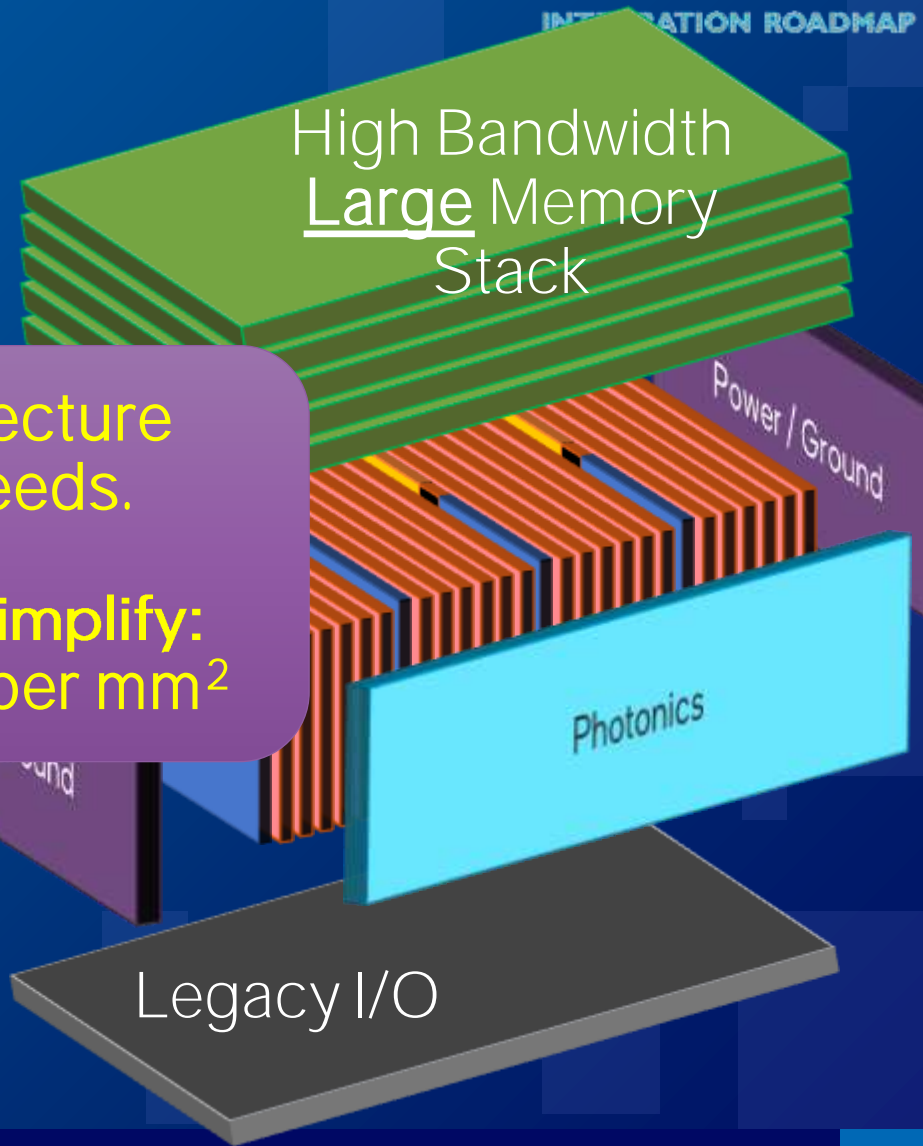
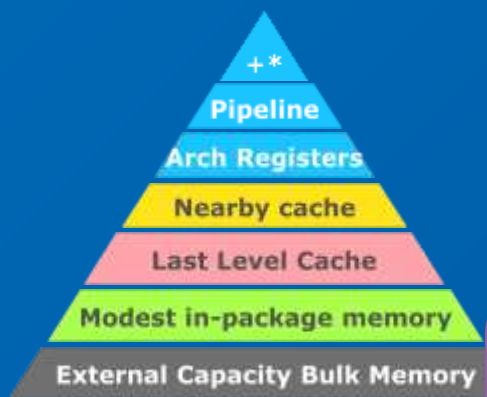


^Not counting overheads for access

3D as a spatial concept, not planar

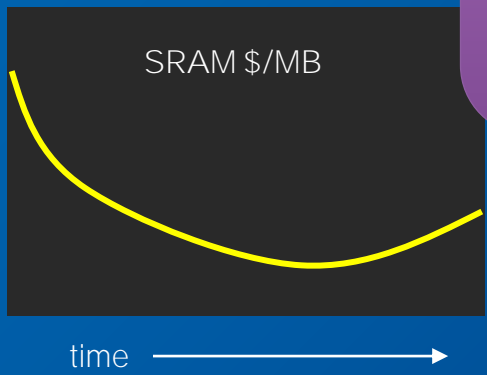


3D as a spatial concept, not planar



Match trade-offs to natural architecture inflection points – speeds and feeds.

Use all 6 faces of "cubic 3D" and simplify:
 $O(1k)$ hybrid bonding connections per mm^2



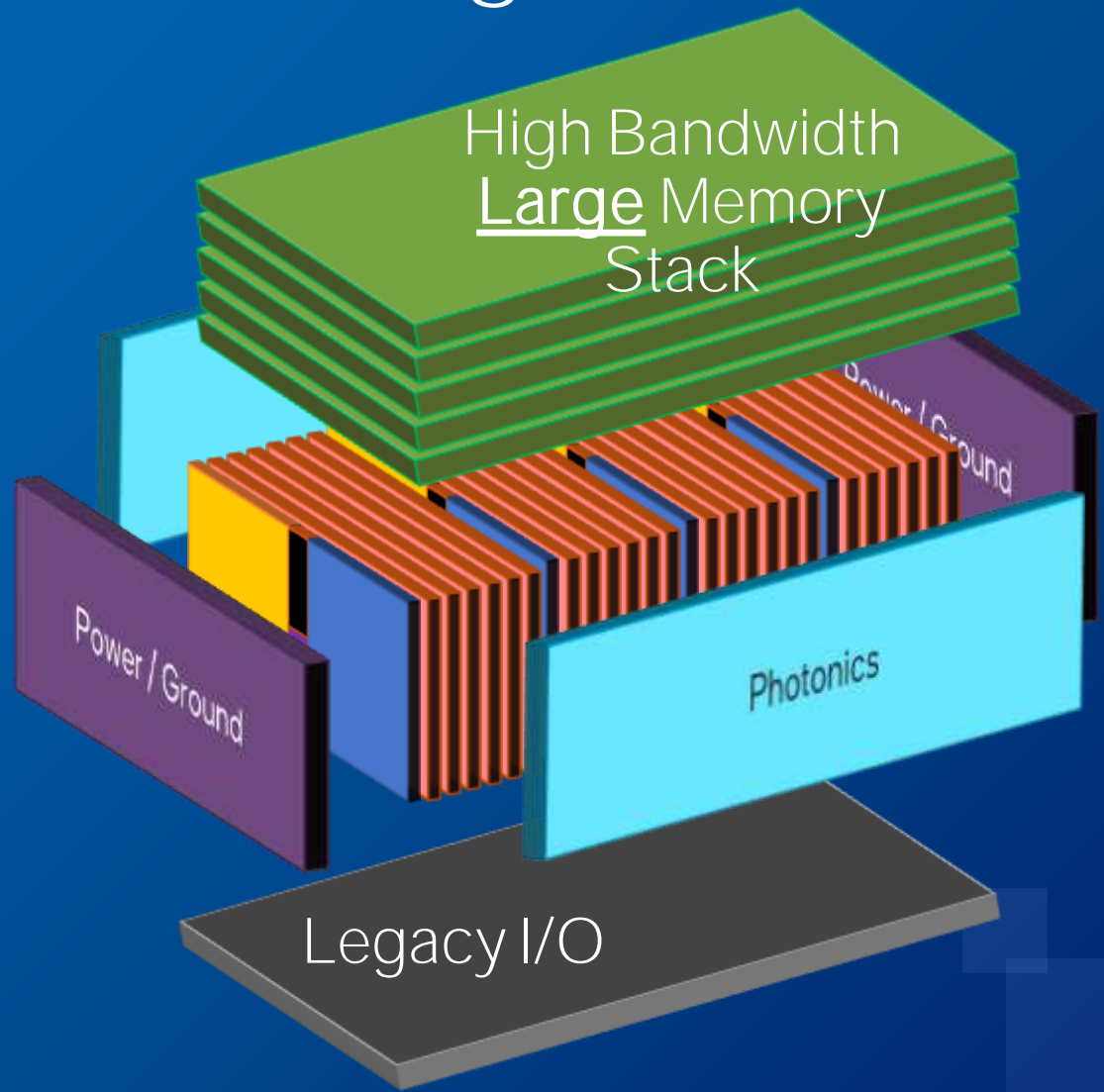
Challenges to achieving dense 3D

Manufacturing:

- edge polish
- right-angle attach
- via density
- bonding speed
- rework support
- redundancy
- tooling Z-height
- wafer thinning

Power and Thermals:

- thermal density
- power delivery
- cooling layers
- heterogeneous material



EDA:

- rotated die
- taper point isolation
- non-planar libraries
- formal verification
- DF<all> 1,000 layers
- scan chain time
- tests and coverage

Design:

- graceful degradation
- extreme interop
- built-in redundancy
- pluggable modules
- abstractions for all



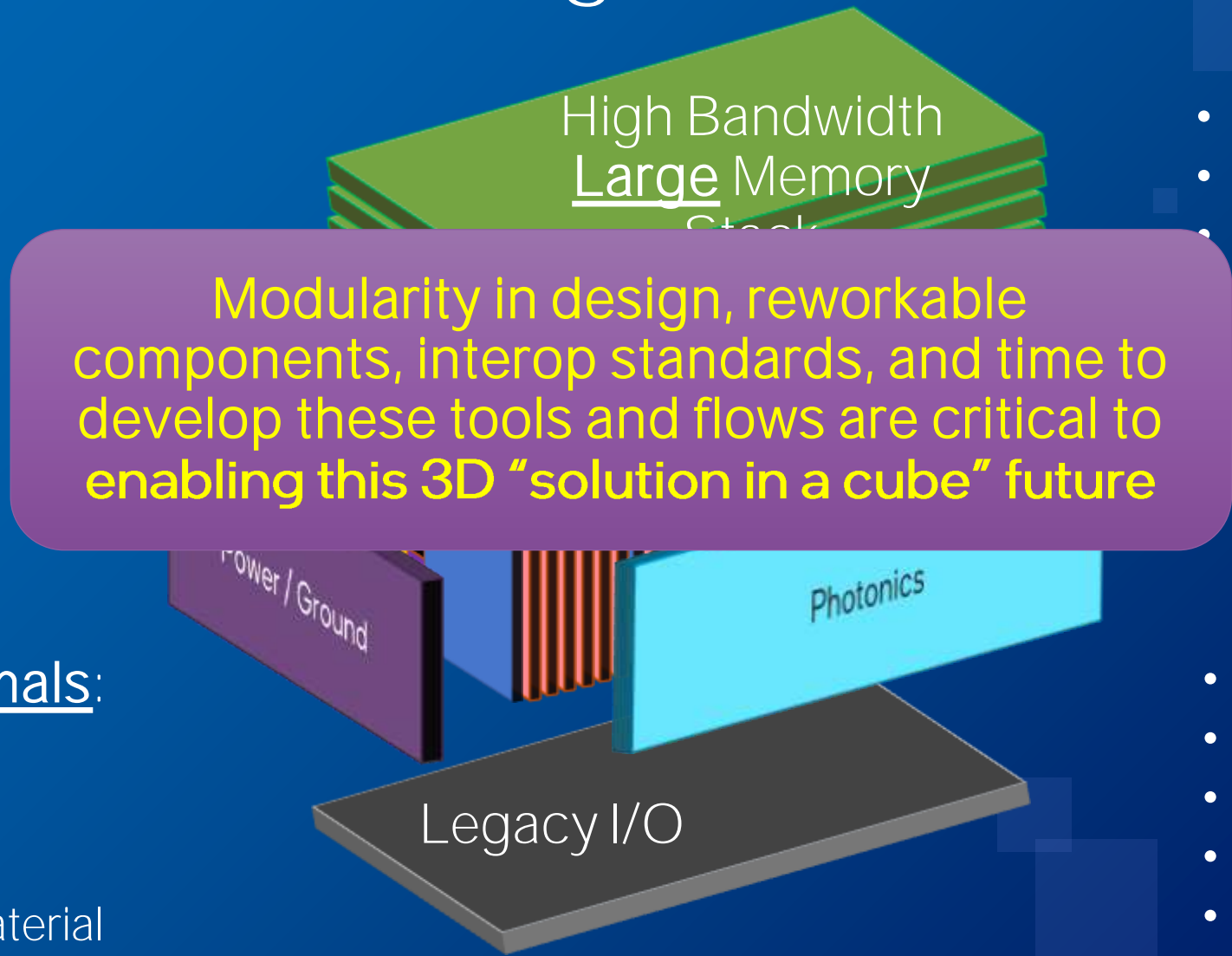
Challenges to achieving dense 3D

Manufacturing:

- edge polish
- right-angle attach
- via density
- bonding speed
- rework support
- redundancy
- tooling Z-height
- wafer thinning

Power and Thermals:

- thermal density
- power delivery
- cooling layers
- heterogeneous material



EDA:

- rotated die
- taper point isolation
- non-planar libraries
- formal verification
- DF<all> 1,000 layers
- scan chain time
- tests and coverage

Design:

- graceful degradation
- extreme interop
- built-in redundancy
- pluggable modules
- abstractions for all

About roadmaps ... hope is not a strategy



- AI will save us! Pixie dust on all!
 - AI will reduce design time
 - It will not innovate – yet
- Some packaging directions to 2035
 - Every 2 yr scaling factor targets
 - 10 yrs to “catch up” and then reassess
 - BW/mm (2x), /mm² (4x), /mm³ (8x)
 - IO escape BW/mm (4x), /mm² (8x)
 - X, Y dimension (1.5x, 1.5x)
 - Z dimension (4x)
 - pJ/b/d²d xing (0.7x)
 - Cooling W/mm² (4x), /mm³ (8x)



OPEN
Compute Project®

Building An Open Chiplet Economy
Monday, September 11, 2023 | Posted by [Sasi Vinnakota](#), Engineer, Lawrence Berkeley Laboratory

Introduction

Chiplet-based products have recently emerged from vendors and cloud providers, including AMD, to extend Moore's law. Their remarkable motivations: (1) to avoid very large die sizes for diversification through mix-and-match strategies; heterogeneous integration - integrating chiplet and/or performance optimization. One expects companies developing semiconductor product domain-specific accelerators is actually domain

High-performance, power-efficient three-dimensional system-in-package designs with universal chiplet interconnect express

[Debendra Das Sharma](#) ^{ORCID}, [Gerald Pascast](#), [Sathya Thangaraj](#) & [Kamal Avqiri](#)
Nature Electronics (2024) | [Cite this article](#)

421 Accesses | 1 Altmetric | [Metrics](#)

Abstract



intel®