

Memory CSD 2020

July 30th, 2020

ELT discussion

Thank you to the MANY who helped make this CSD Happen High Engagement / Inclusive / Corporate-Wide

Leads: Carolyn Duran / Frank Hady

CSO: Mark Pontarelli

Corporate Effort

DPG*	CCG
IAGS	TMG
NSG	CSO
SMG	CEG
ICAP	XPG

Finance
Supply Chain
Intel Labs
IP Engineering
Bain

*CESG, DPEA, IOTG, MIO, NPG, SBDO, XMG

Team of Experts and Advisors

Adiletta, Matt
Ahmad, Marisa
Appello, Adrienne
Ard, Jennifer
Bains, Kuljit
Balkan, Haluk
Chennupati, Srinivas
Chow, Gary
Coulson, Rick
Eylon, Eyran
Fazio, Al
Feghali, Wajdi

Forell, Kerry
Foucher, Yoann
Fryman, Josh
Galbi, Duane
Gardina, Jeff
Goldschmidt, Mark
Gomes, Wilfred
Hamzaoglu, Fatih
Hatzikos, George
Ihlefeld, Allen
Ilkbahar, Alper
Isbara, Melik

Karl, Eric
Kau, Derchang
Kottapalli, Sailesh
Lal, Manoj
Leta, Eric
Loop, Becky
Luhmann, Fiona
Mather, Nate
Megit, Jason
Moga, Adrian
Mueller, Steve
Nagisetty, Ramune

Nikonov, Dmitri E
Noonan, Tim
Onufryk, Peter
Osborne, Randy
Peschka, Brian
Pickett, Jay
Rajput, Pushpdeep
Rangarajan, Madhu
Rao, Radhika
Royer, Bob
Roytman, Eduard
Schulz, Matthew

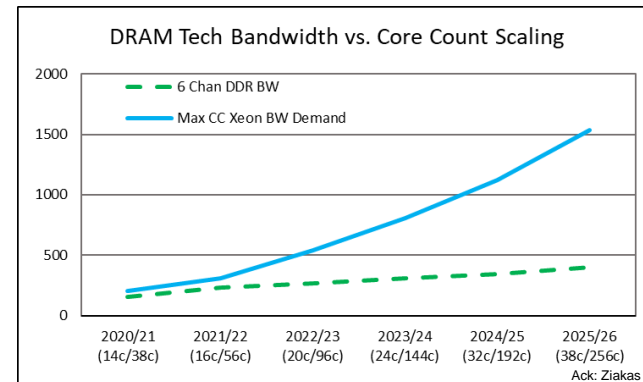
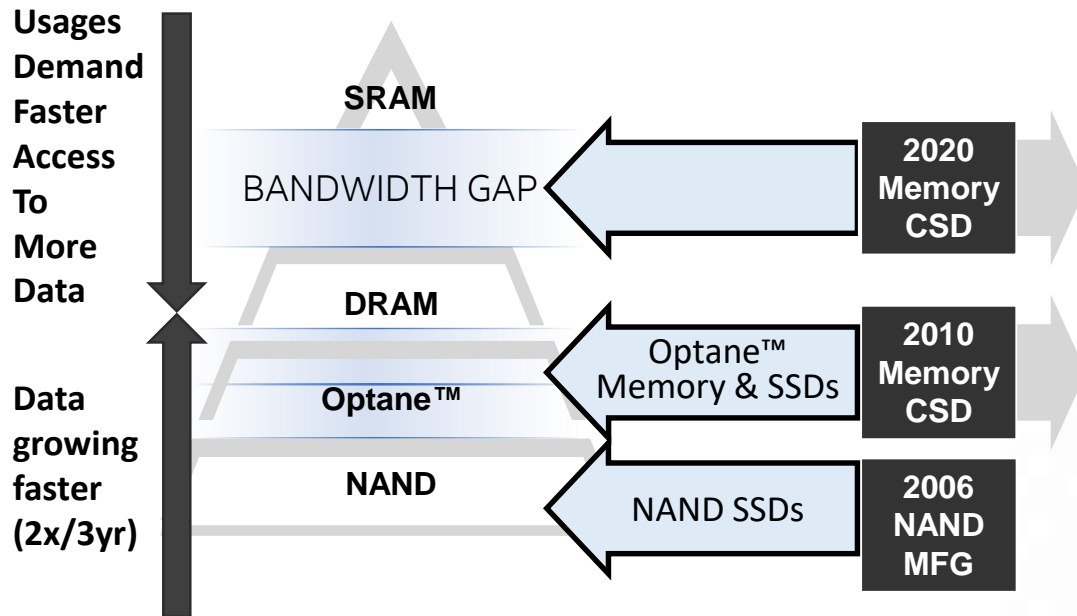
Sell, Ben
Singhal, Ronak
Skarpness, Mark
Soto, Percy
Sury, Samantika
Tomishima, Shigeki
Tripathi, Brijesh
Wechsler, Ofri
Ziakas, Dimitrios

Sponsors: Rob Crooke, Raja Koduri, Murthy Renduchintala, Navin Shenoy

Expected Outcome

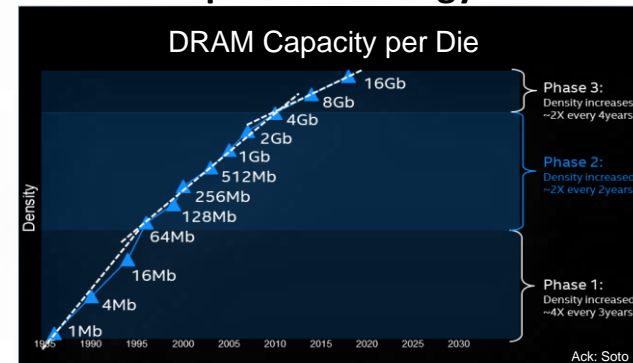
- **Align: We *must* do something to address the fundamental physics limitations of high bandwidth data movement in order to scale compute.**
 - We are NOT suggesting Intel get into the DRAM business
- **Discuss/Decide: Strengthen XPU product leadership by establishing a differentiated cache strategy**
 - Invest and deliver Adamantine technology roadmap for a differentiating 1s of Gigabytes cache
 - Commit to employing that solution for AI/Graphics/Xeon products
- **Act: As a result of this CSD we will also**
 - Work with the DRAM industry and key customers to ensure Intel is competitive for 10s of Gigabyte high bandwidth working sets
 - Initiate an *ELT empowered* Memory Initiative to ensure cross BU alignment in memory requirements and systems innovations

Memory CSD Scope: SRAM to DRAM

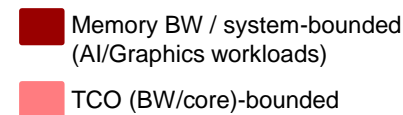


DRAM is Falling Behind in Performance:
This CSD

DRAM is Falling Behind in Capacity:
Optane™ Strategy

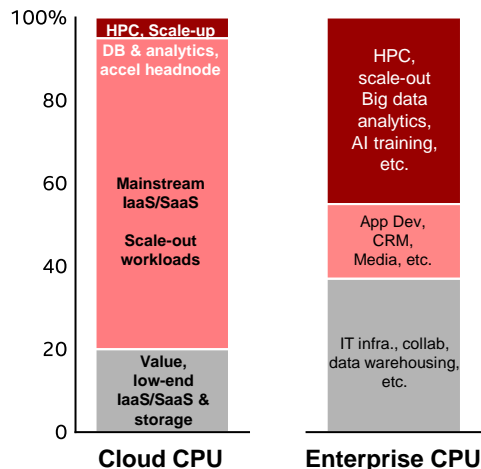


Today's DRAM challenges across segments: problem statements



Data center: Core count & AI are scaling faster than DRAM bandwidth/capacity/power

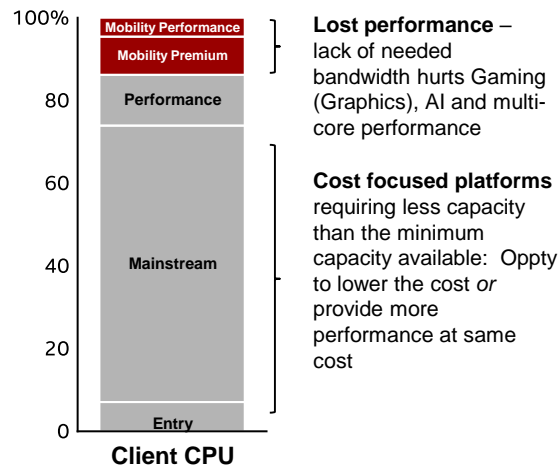
% cores and workloads (today)



- Creating unsustainable increases in our customer's Capex (cost) and OpEx (power)
- Not delivering performance customers can charge for (relative to GPU, accelerators)
- Decreasing the relevance of Intel products because consumers cannot access full benefit

Client: Application needs are increasingly misaligned with DRAM technology

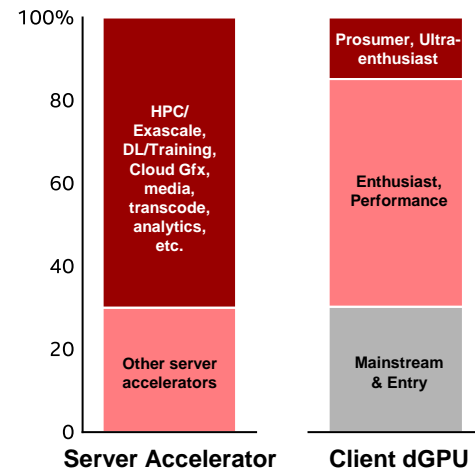
Client CPU volume (today)



- Cost focused platforms paying for more capacity than needed
- Performance focused clients suffering low graphics performance due to DRAM bandwidth shortfalls

Accelerators: Memory performance bottlenecking system performance

Accelerator revenue TAM (today)



- DC Accelerator: Customers want max BW, 10's of GB capacity per chip, and min power and are prepared to pay for memory architecture innovations.
- Client Accelerator: high Graphics performance depends on high memory bandwidth to a small working set (<1 to 10GB). Solutions are cost sensitive.

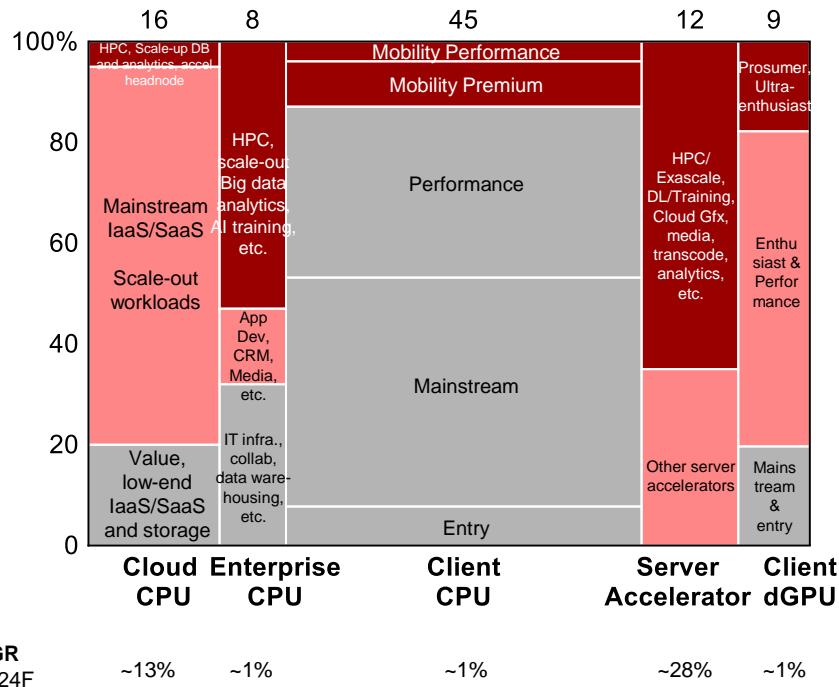
Future: Large and growing share of our XPU TAM is memory bandwidth bound

2024 Server, Client CPU and accelerator TAM

Total = ~\$90B

■ Memory BW Bounded = ~\$20B

■ TCO Bounded = ~\$25B

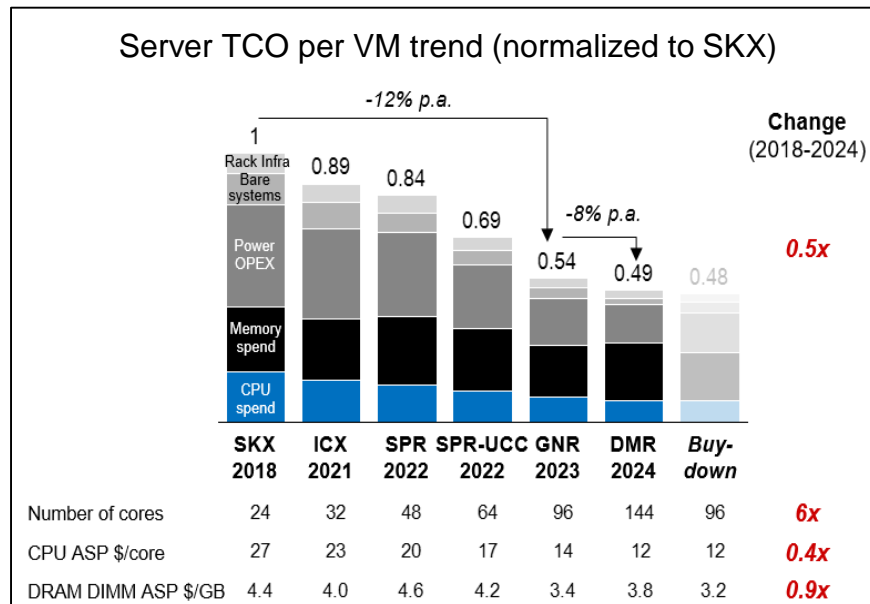


Two workload challenges

- AI/graphics workloads across CPU and accelerators demanding **high-bandwidth** memory:
 - \$20B TAM
 - Fastest growing segments
 - Highest gross margin (>70%)
 - **Major growth bets for Intel (e.g. dGPU), challenging strong incumbent**
 - Underpinning market cap
- Large portion of datacenter scale-out workloads bounded by **bandwidth per core**, **leading to a TCO wall as our core count grows to >100 in 2024**

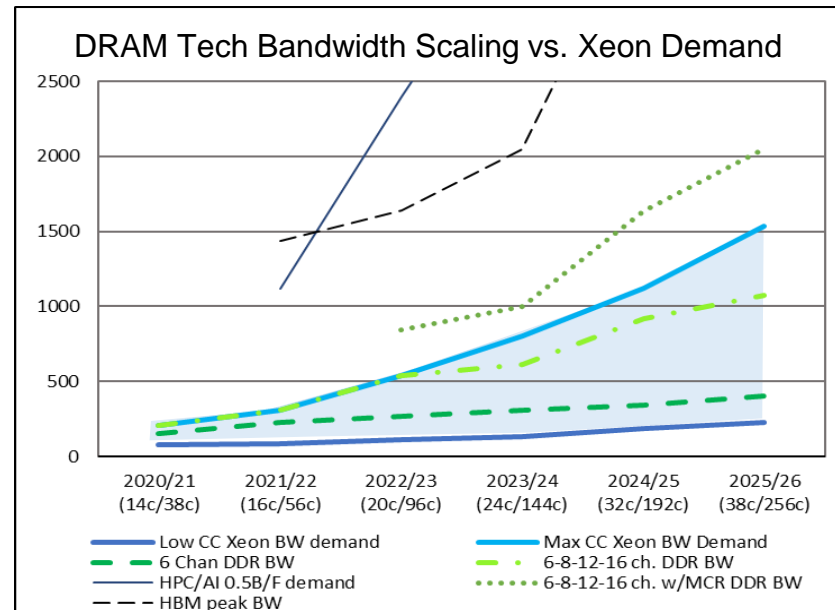
Our DC customers are impacted by the DRAM bandwidth gap

CSP TCO reduction limited by DRAM as CPU core count grows; client min. DRAM capacity drives BOM cost



Source: Madhu Rangarajan

CSP driven increased Xeon core counts and new Xeon AI accelerators requiring costly system changes to avoid DRAM bandwidth bottleneck

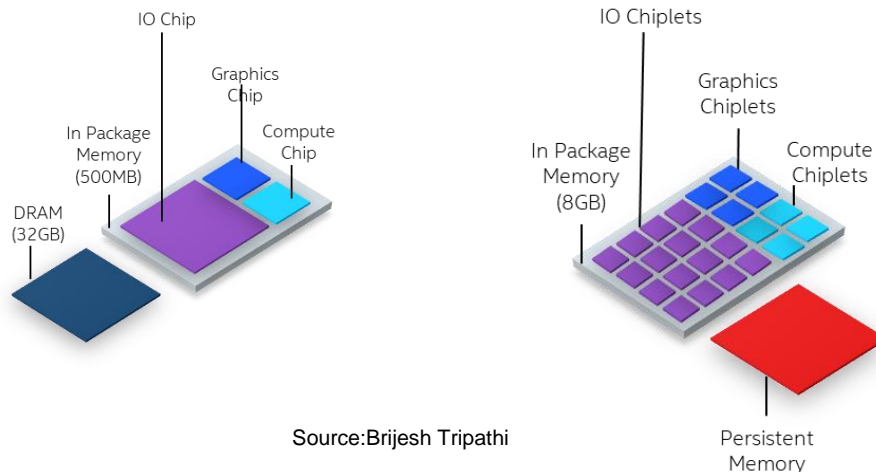


Source: Dimi Ziakas

A subset of server applications can benefit significantly from an ADM-based cache*

*Some benefit from ADM's latency advantage, while others benefit from ADM's bandwidth advantage. With ADM at the right price, an optional ADM cache can provide advantage to server customers.

Client usages require lower power with high bandwidth



Source: Brijesh Tripathi

MTL (2022/2023)

- ❑ **Improved Gaming Experience w/ 30%-50% better graphics perf.**
(Integrated Graphics bandwidth)

Future (LNL+, 2023+)

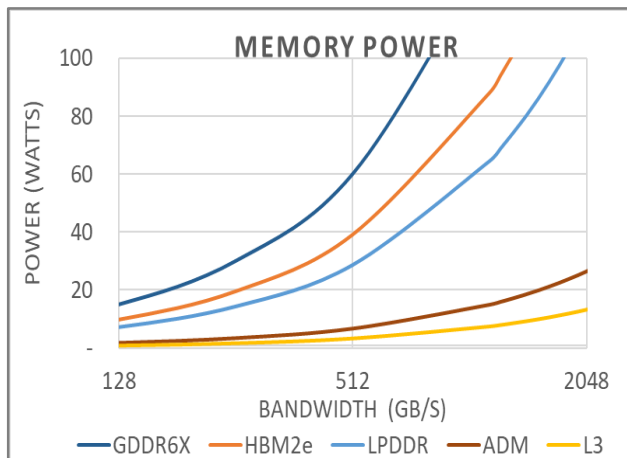
- ❑ **Better Gaming Experience**
(Integrated Graphics bandwidth)
- ❑ **Better Content creation**
(Bandwidth for Multicore workloads)
- Better Form Factor**
(DRAM replacement)
- Better battery life**
(DRAM replacement for lower pwr bw)

Integration of a on SoC differentiated Cache (ADM*) across our Client products would deliver leadership gains in Gaming, Content Creation and Battery life.

Technical Fundamentals Drive High Bandwidth Memory Selection

Memory Power:

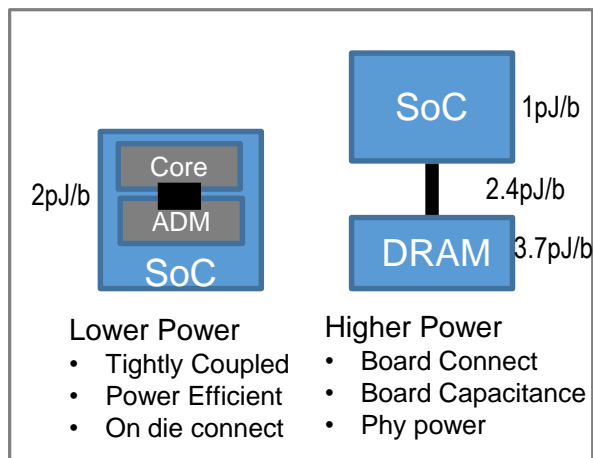
Changing Memory type required eventually as bandwidth grows but power budget remains



Source: Randy Osborne

Interconnect Power:

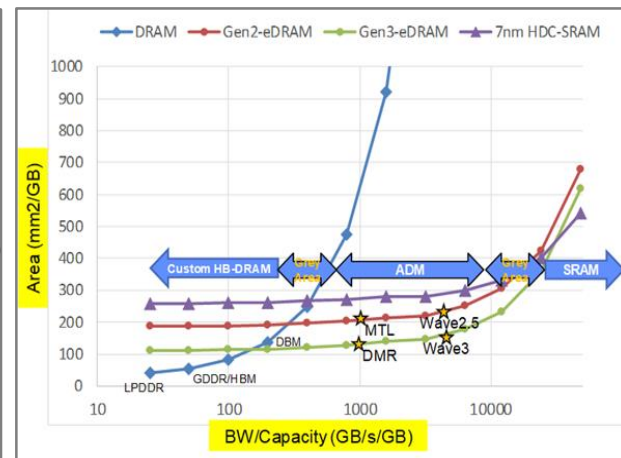
Lower capacitance memory interconnect required as bandwidth grows but power budget remains



Source: Fatih Hamzaoglu, Bob Royer

Bandwidth/Capacity:

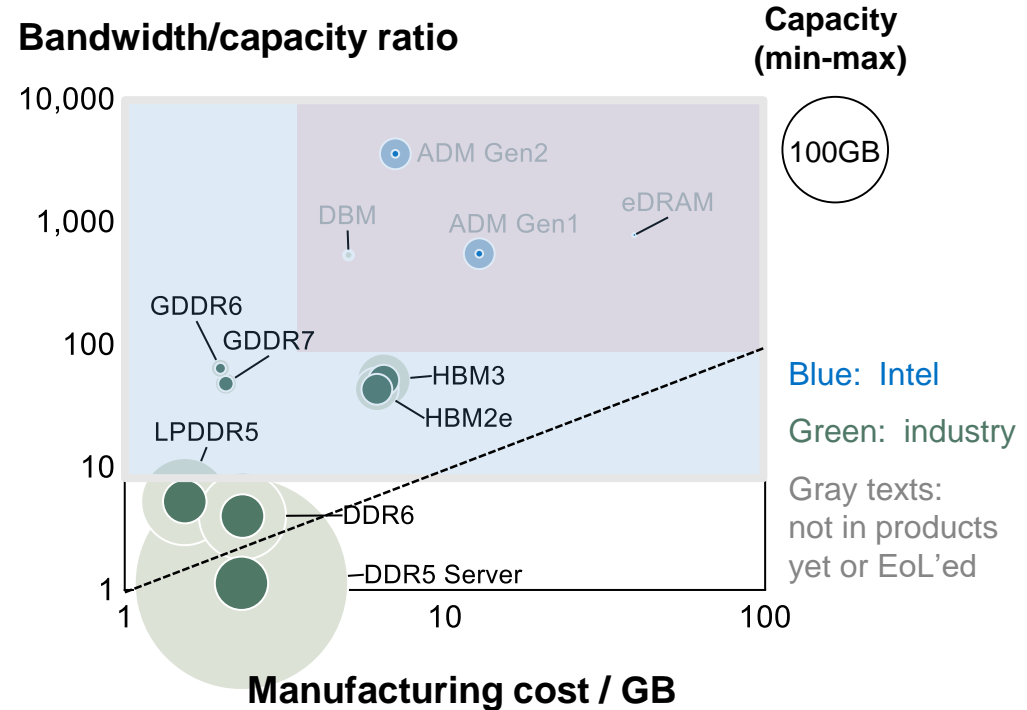
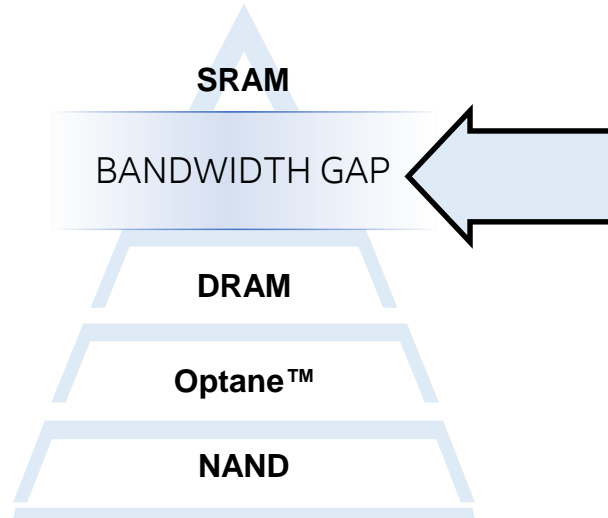
The area efficient memory depends on the combination of capacity / bandwidth needed



Source: Fatih Hamzaoglu

**With low pJ/b, tight integration, and high bandwidth/capacity
ADM promises a unique opportunity for Intel**

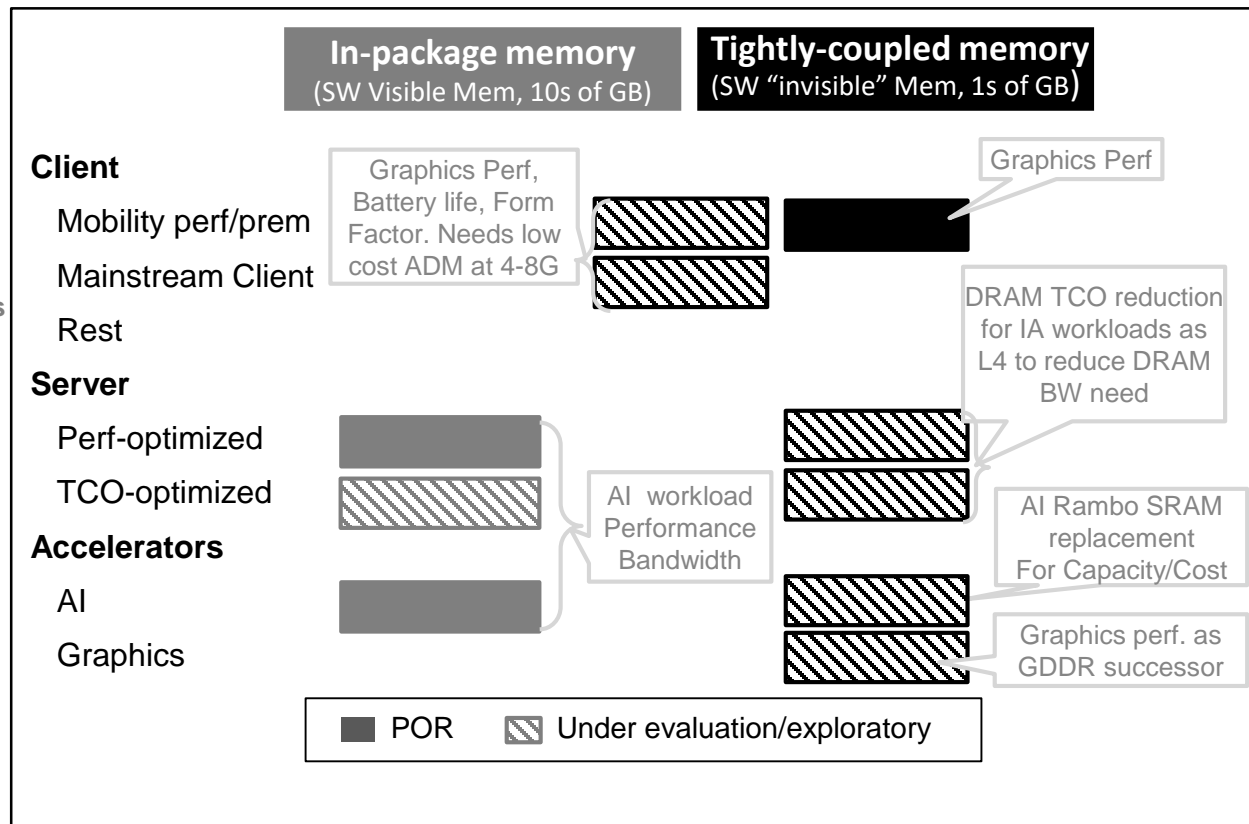
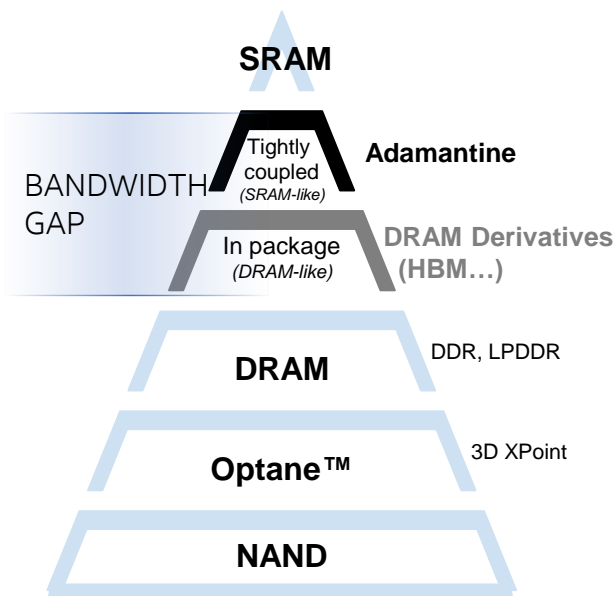
Traditional DRAM solutions cannot address bandwidth/capacity/cost simultaneously



ADM is best for several TB/s bandwidth @ 1s of GBs in capacity

High bandwidth DRAM derivatives (HBM) are best for 10s of GB capacity with 100s of GB/s bandwidth

Usage Needs Drive Two Technology Choices



What are others doing to address memory challenges?

In-package memory (DRAM-like, 10s of GB)

Competitors and customers are investing in industry-standard IPM, with pathfinding efforts in customization

Examples

- ❑ **Logic and memory vendors** driving **HBM** standard to next-gen, in accelerators and CPU
- ❑ **Apple** A13 **LPDDR** integrated with logic using **TSMC InFo packaging**
- ❑ **Google, TSMC, and Samsung or Hynix** investigating **HBLL** (High Bandwidth Low Latency) memory (1s of GB)
- ❑ **Nvidia and Micron** working on **unique memory** (limited info known)
- ❑ **TSMC** and IP ecosystem investing in **High-Bandwidth Interconnect PHY** to increase compute and memory intimacy in chiplet designs
- ❑ ...

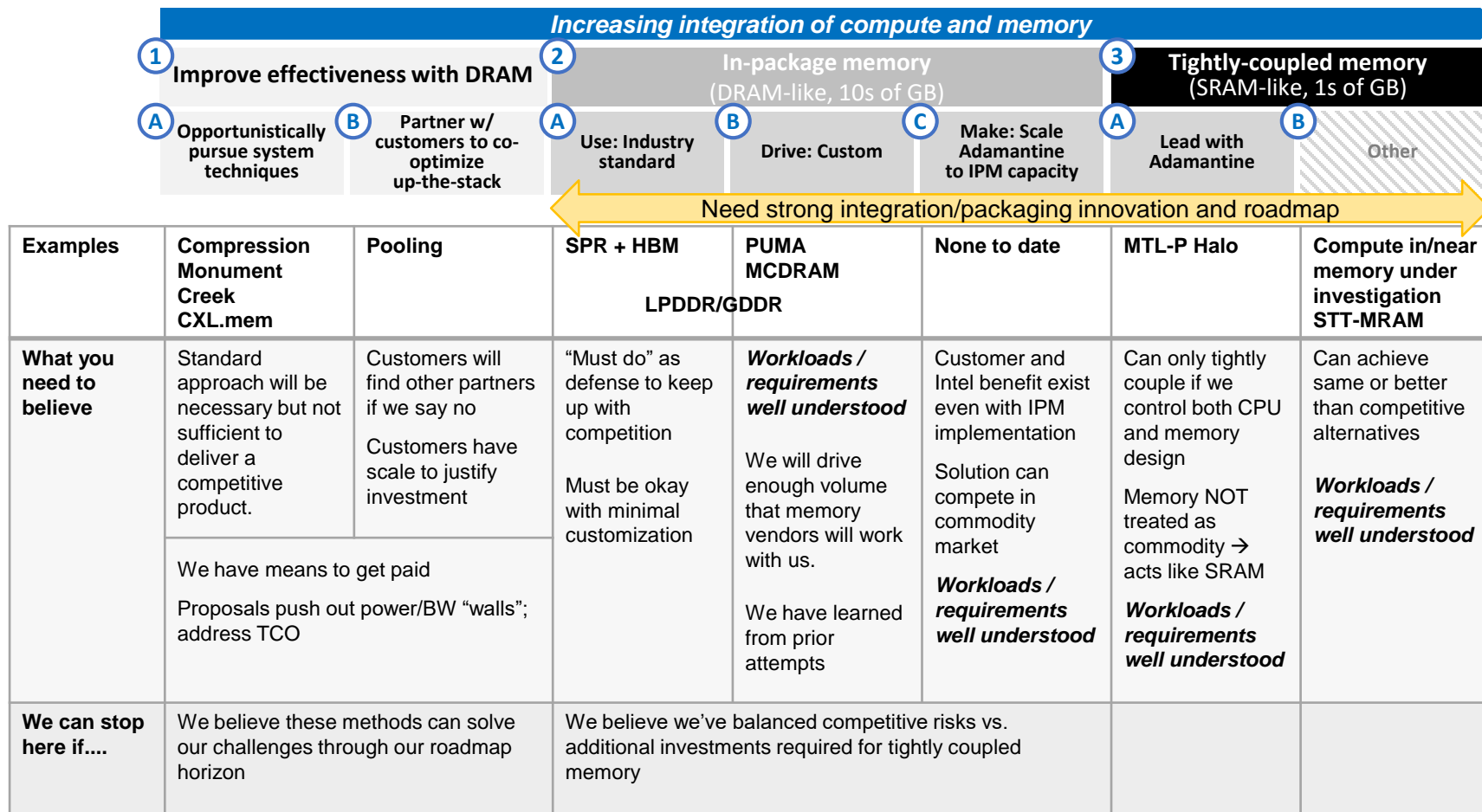
Tightly-coupled memory (SRAM-like, 1s of GB)

A variety of efforts to further increase memory-compute intimacy. (No direct substitute to Adamantine)

Examples

- ❑ **Logic vendors** working on larger/better **SRAM** cache
- ❑ **AMD** working on **Foveros-like 3D** stacking, likely with **TSMC**
- ❑ **Apple TSMC** co-innovating on **next-gen 3D stacking** patents
- ❑ **TSMC, Samsung, Hynix, GloFo** etc. investing in **STT-MRAM** with embedded NVM, potential to expand to tightly-coupled memory
- ❑ **Many** investing in **Compute in Memory** pathfinding, limited commercialization success (Mythic AI early mover)
- ❑ ...

No silver bullet: Framing the Intel Options

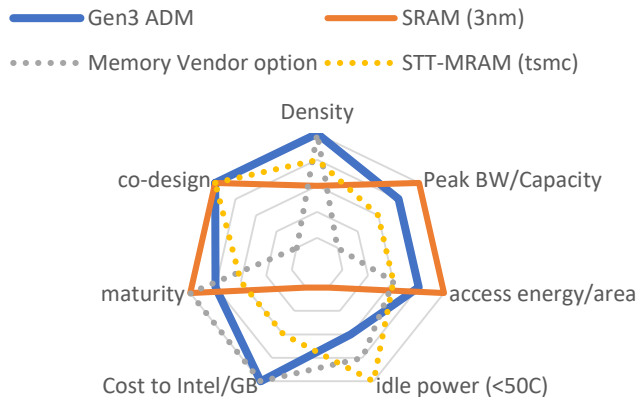


Discussion Flow

- **Tightly Coupled** → Adamantine
- In Package Memory → high bandwidth DRAM variants
- Systems → Many

Solving for high bandwidth needs in the 1s of GB memory capacity: Tightly-coupled ADM is a differentiator

Technology Trends for 1s of GB High BW memory (2024 estimate)



Driving Factors

- Graphics and AI workloads are driving the need for extremely high bandwidths at 1s of GB capacity
- Tightly coupling compute and memory is the only way to meet bandwidth without hitting power/TCO walls.
- Tightly coupled solutions require co-design of compute and memory

Adamantine Advantage

- Only known solution in desired timeframe (~2yr adv?)
- Best means to address BW, power, capacity challenges
- Intel controls all input to maximize benefit to customer
- “invisible” to higher level SW (like SRAM)

Not included:

IPM: Traditional solutions over-index on capacity and are power hungry

In-memory, Near-memory compute: Research should continue → to date, point solutions are likely

Staging ADM to achieve long term XPU differentiation

All available info indicates we have ~2yr head start relative to competition – use it!

ADM cost and Platform volume need to commit to scale together (key learning from eDRAM)

- Technology Scaling and Cost: \$/GB – TMG commitment
- Product Volume: GB shipped – BU commitment

Phase	What We Mean	What it Takes	Success
<u>Disrupt</u> and Commit by 2022/2023	Invest up front in tech roadmap Disrupt in graphics with MTL-P Invest in targeted Xeon use case → where we hit the physics wall first	Committed investment to roadmap (100s of \$M, 22nm) Targeted, small volume (~10 MU for MTL-P; 1-2 MU for Xeon) Margin compromise	Design wins for MTL-P Early learning on ADM in DC applications – we figure out to do it right at the system level
<u>Innovate</u> and Establish by 2023/2024	Select targeted expansion at corporate level	Continued TD roadmap Committed BU / architecture resources Deep working set /use case knowledge Margin compromise	Maturing technology to drive scale; developing common requirements Design wins in key areas
<u>Scale</u> by 2024/2025	Utilize ADM across the business as core business opportunity	Continued TD roadmap must hit cost metrics Committed BU / architecture resources Deep working set /use case knowledge Must compete against alternatives	Scale enables standard business model and margins

Opportunities to scale are many, need resolution on each with aligned strategy

Proposed stage	BU	Area	Motivation	Maturity	Comments
Disrupt	CCG	Perf/Prem Mobility	As cache: Disrupt in graphics	POR for MTL-P	Compete and win against entry discrete graphics
	DPG	Performance optimized	PoC in BHS or EGS?	Concept	Test value proposition with side attach
Innovate	Accelerators	AI/Graphics	As cache: RAMBO SRAM replacement for capacity/cost	Under investigation	PVC-next, Elasti-Pro, Elasti-Sound
	CCG	Premium mobility	As memory: Enable form factor & compute differentiator over dGx	Concept	Target same perf advantage as an ADM cache but lower system power
	DPG	TCO optimized	As cache: DRAM TCO reduction for IA workloads	Under investigation for DMR	ADM L4 cache reduces DRAM BW needs, allows higher core count
Scale	DPG	Perf optimized	Performance and TCO reduction	Under investigation for DMR	Initial assessment shows 6% perf boost, equivalent TCO/VM at 50% margins; many assumptions
	CCG	Mainstream mobility	As memory: Enable form factor & compute differentiator over dGx	Concept	As above but mainstream cannot tolerate cost premium
	CCG	Prem/Perf mobility; Mainstream	As cache: Cost neutral with longer battery life and higher graphics performance	Under investigation	Need 4-8GB ADM, addresses power/thermal issues
	Accelerators	AI/Graphics	As memory: Performance	Concept	GDDR successor
	DPG	NPG	Virus scanning; small capacity	Concept	Need to address ECC, security

Discussion / Decision

- Despite a very critical look ADM → team became convinced ADM will best strengthen XPU product leadership by establishing a differentiated cache strategy
- This IS the POR plan to disrupt Graphics with MTL-P

Strategic Decision

Go: Phased Implementation of ADM

No Go on ADM

or

- ~~Require ADM “attach” volume necessary to hit expected core business ROI and margins out of the gate~~
- ~~Decommit client; suffer competitive risk~~

Intel Confidential
20-07-30 Memory CSD_ELT_v10.PPTX

Staging ADM for long term success; learning from eDRAM

- Commitment: Upfront to end goal
- Expediency: All available info indicates we have ~2yr head start relative to competition – use it!
- Execution: Logic, ADM, subsystems and package all need to come together in concert
- What does it look like?

Phase	What We Mean	What It Takes	Success
Disrupt and Commit by 2022/2023	Invest up front in tech roadmap Disrupt in graphics with MTL-P Invest in targeted Xeon use case → where we hit the physics wall first	Committed investment to roadmap (100s of \$M, 22nm) Targeted, small volume (~10 MU for MTL-P; 1-2 MU for Xeon) Margin compromise	Design wins for MTL-P Early learning on ADM in DC applications – we figure out to do it right at the system level
Innovate and Establish by 2023/2024	Select targeted expansion at corporate level	Continued TD roadmap Committed BU / architecture resources Deep working set / use case knowledge Margin compromise	Maturing technology to drive scale; developing common requirements Design wins in key areas
Scale by 2024/2025	Utilize ADM across the business as core business opportunity	Continued TD roadmap must hit cost metrics Committed BU / architecture resources Deep working set / use case knowledge Must compete against alternatives	Scale enables standard business model and margins

18

Intel Options to address 10s of GB needs

	In-package memory (10s of GB)		
	Off-the-shelf (i.e. HBM)	Custom	Scale Adamantine to IPM capacity
Pros	Minimal investment	Clear customer need Workload intelligence System optimization Can provide differentiation Opp for new business models	Reduces reliance on DRAM Drives value to Intel System optimization Can provide differentiation
Cons	Might not solve the problem We need to adapt to adopt We start from behind Absorb memory "price"	Custom solution may limit scale Absorb memory "price" Multi-party alignment req'd	Might not solve the problem ADM benefits degraded as IPM Single/sole source to customer
Investment required to get to a memory to use	1's of \$M	10's \$M	\$100's of \$M
Intel key contributions	Industry Influencing	Compute/memory Optimization Industry Influencing Packaging Innovation	ADM process Compute/Memory Optimization Packaging innovation

IPM key messages

• Guiding beliefs (a.k.a. lessons learned)

- We have to focus on solving the customer problem first → need deep understanding of workloads, their use of memory and their future trajectory.
- Any IPM implementation requires optimization of the software stack.
- New innovations must be judged against realistic internal and memory ecosystem business models and affordability assumptions.
- 10s of GB capacities need to start from a DRAM process; we have had decades of success guiding the industry and can continue to push capabilities
- AI workloads drive high bandwidth memory needs

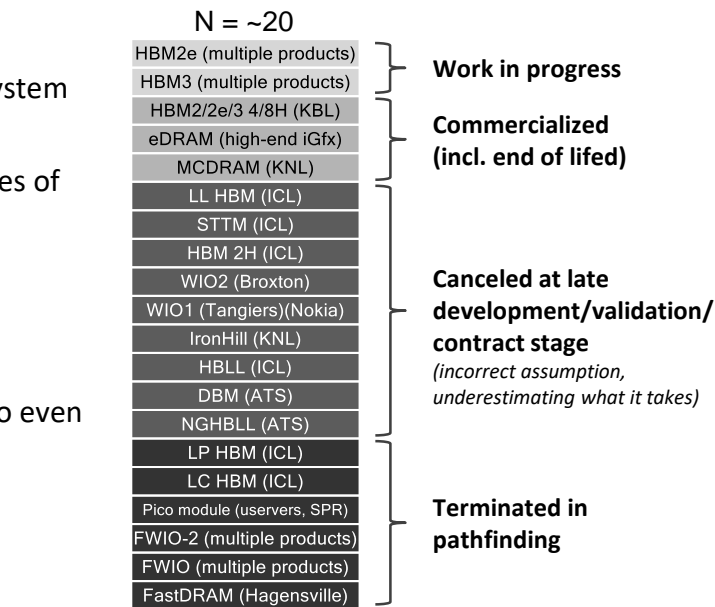
• Plan

- Continue to drive next gen memory with industry, deepening collaboration to even more strongly influence new-to-Intel memories (i.e. HBM, GDDR)
- Determine opportunity to create a strategic partnership with a definitional customer to develop a differentiated solution.

• Alternative

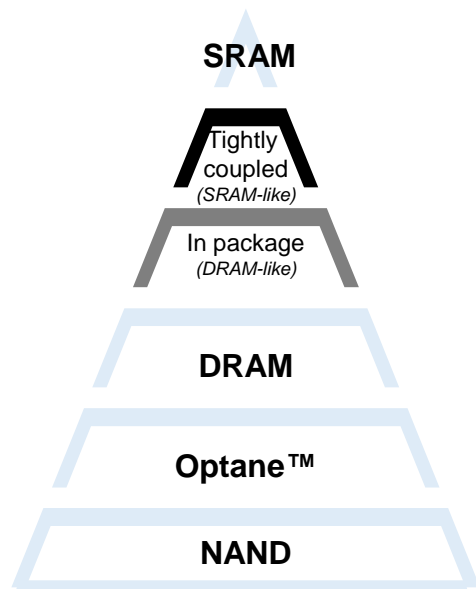
- Adapt to adopt external IPM offerings as they become available

In-package memory technologies attempted & products landing zone



Systems innovations remain crucial to mitigate DRAM's limitations

Memory Tech Bounds the Experience



This CSD:

A new memory/level is necessary for multigenerational bandwidth improvements

Systems Innovations Optimize Memory Use

Innovation	Intercept		BW	Cap.	Pwr
Memory Side Cache –	GNR or DMR	DC	X		
Monument Creek	GNR	DC	X		
Cache compression/DeDup		DC	X	X	
Increase IP Memory Efficiency	2022, MTL	Client	X		X
Memory Compression	Now-2023	All	X	X	X
Additional Channels	Ongoing	All	X	X	
Page level Memory compression		DC		X	
Optane™ Memory Mode/App Direct	now	DC		X	
CXL capacity expansion (2LM)	GNR for CXL	All		X	
Flat2LM	GNR (2023)	DC		X	
Memory pooling	GNR (2023)	DC		X	

Per BU:

Systems Innovations for one-time Memory bandwidth advantages

Recommendation: BU's keep portfolio of programs to mitigate customer pain points, use Memory Initiative, solution should meet standard business model expectations (margin, etc.)

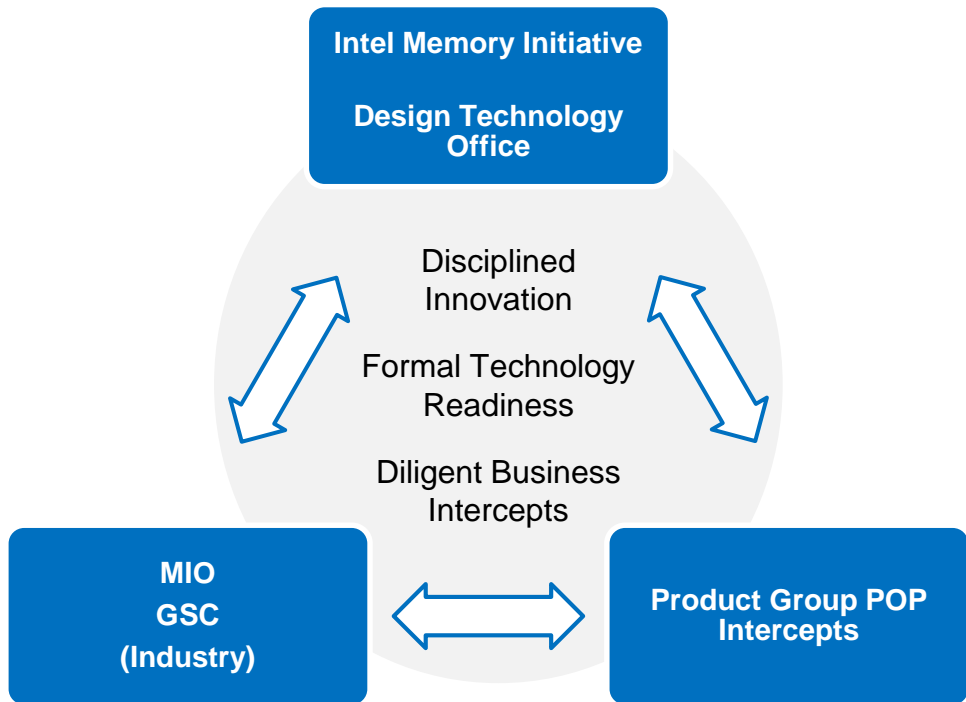
Plan: Proceed with creation of an Intel Memory Initiative

Expected Impact:

- Robust technology portfolio
- Anticipate/shape industry transitions
- Reduce late stage roadmap changes
- Reduce roadmap TTM
- Improve solution stacks/SW @ launch

Next Steps:

- Identify Initiative lead & ELT sponsor
- Appoint ADM roadmap integration lead
- Initiative lead and PM to align Initiative, DTO, MIO, and technical communities
- Work with BUs to create tech roadmap <> product roadmap interfaces



Expected Outcome / Next Steps

- **Align: We *must* do something to address the fundamental physics limitations of high bandwidth data movement in order to scale compute.**

- We are NOT suggesting Intel get into the DRAM business

→ *Agreed*

- **Discuss/Decide: Strengthen XPU product leadership by establishing a differentiated cache strategy**

- Invest and deliver Adamantine technology roadmap for a differentiating 1s of Gigabytes cache

- Commit to employing that solution for AI/Graphics/Xeon products

→ *Agreement to move forward with product-level assessment of ADM solution(s)*

- **Act: As a result of this CSD we will also**

- Work with the DRAM industry and key customers to ensure Intel is competitive for 10s of Gigabyte high bandwidth working sets

- Initiate a an *ELT empowered* Memory Initiative to ensure cross BU alignment in memory requirements and systems innovations

→ *Agreed*