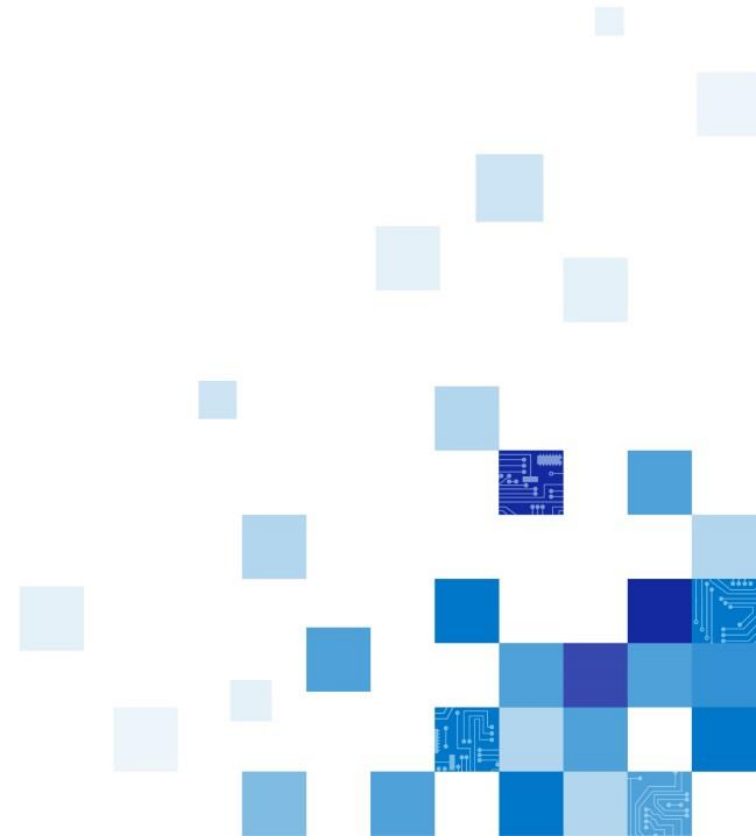


LPW(LP WIDE IO) Introduction

Feb. 2024 | Samsung Memory

Confidential



Caution for Document License

THIS DOCUMENT AND ALL INFORMATION PROVIDED HEREIN (COLLECTIVELY, “INFORMATION”) IS PROVIDED ON AN “AS IS” BASIS AND REMAINS THE SOLE AND EXCLUSIVE PROPERTY OF SAMSUNG ELECTRONICS CO., LTD. CUSTOMER MUST KEEP ALL INFORMATION IN STRICT CONFIDENCE AND TRUST, AND MUST NOT, DIRECTLY OR INDIRECTLY, IN ANY WAY, DISCLOSE, MAKE ACCESSIBLE, POST ON A WEBSITE, REVEAL, REPORT, PUBLISH, DISSEMINATE OR TRANSFER ANY INFORMATION TO ANY THIRD PARTY. CUSTOMER MUST NOT REPRODUCE OR COPY INFORMATION, WITHOUT SPECIFIC WRITTEN CONSENT FROM SAMSUNG. CUSTOMER MUST NOT USE, OR ALLOW USE OF, ANY INFORMATION IN ANY MANNER WHATSOEVER, EXCEPT FOR CUSTOMER’S INTERNAL EVALUATION PURPOSE. CUSTOMER MUST RESTRICT ACCESS TO INFORMATION TO THOSE OF ITS EMPLOYEES WHO HAVE A BONA FIDE NEED TO KNOW FOR SUCH PURPOSE AND ARE BOUND BY OBLIGATIONS AT LEAST AS RESTRICTIVE AS THIS CLAUSE. BY RECEIVING THIS DOCUMENT, IT IS UNDERSTOOD THAT CUSTOMER AGREES TO THE FOREGOING AND TO INDEMNIFY SAMSUNG FOR ANY FAILURE TO STRICTLY COMPLY THEREWITH. IF YOU DO NOT AGREE TO ANY PORTION OF THIS CLAUSE, PLEASE RETURN ALL INFORMATION AND ALL COPIES (IF ANY) WITHIN 24 HOURS OF RECEIPT THEREOF.

Agendas

1) Market Introduction

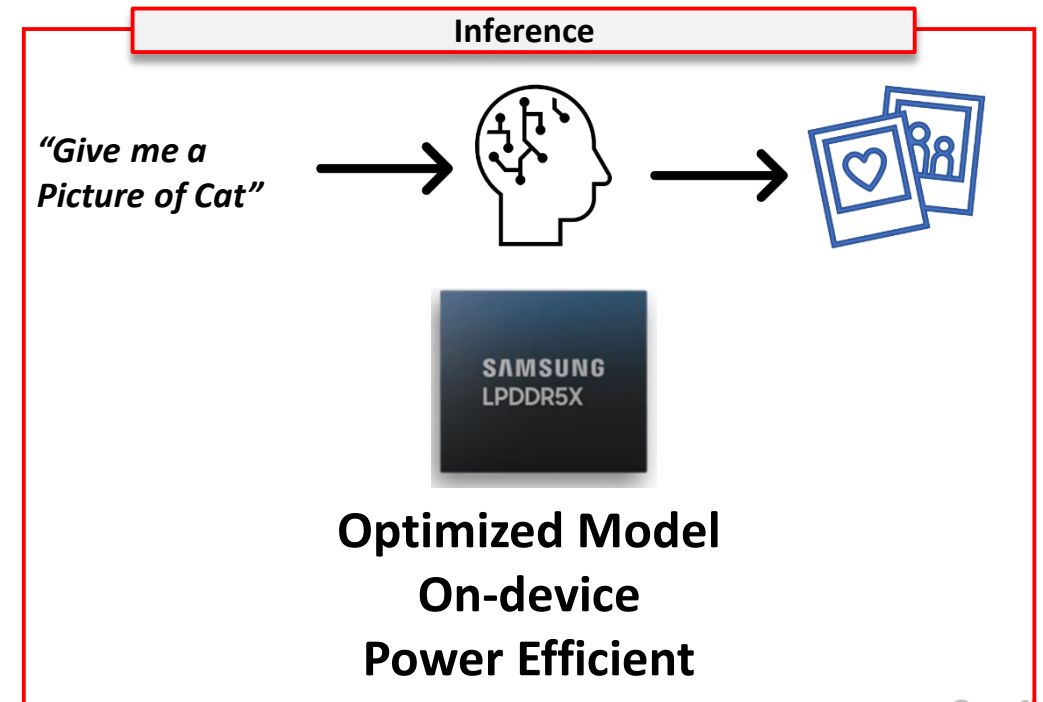
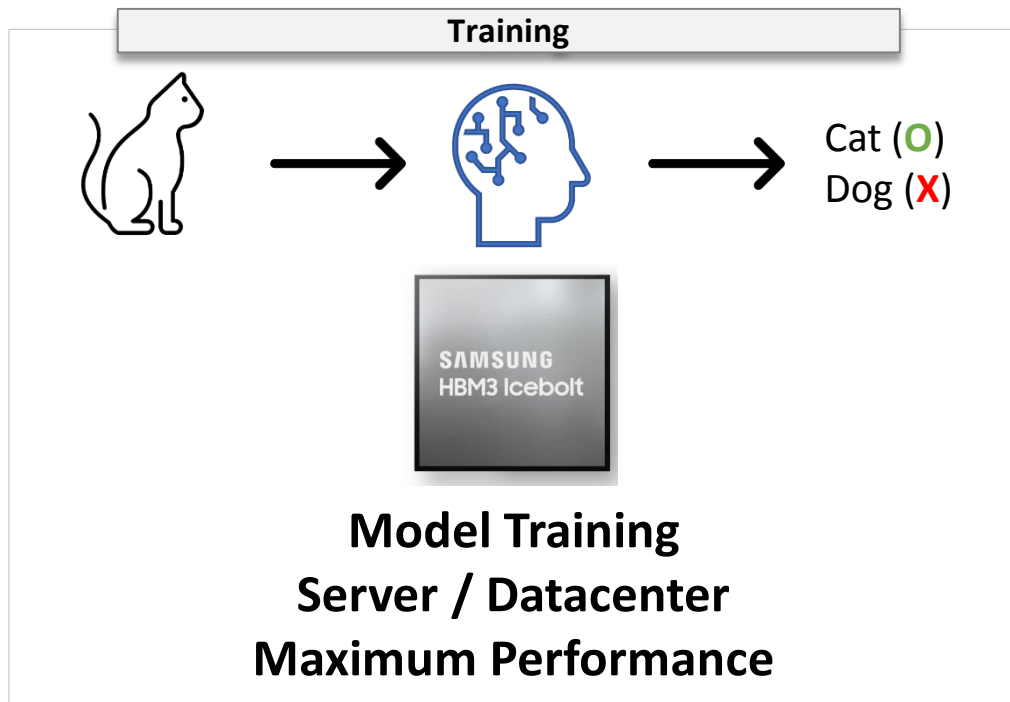
- AI Market Trend
- On-Device Generative AI Analysis and Memory Roadmap

2) AI Memory Solution for Mobile

- LPW(LPDDR Wide IO) Concept Introduction
- LPW(LPDDR Wide IO) Architecture
- LPW(LPDDR Wide IO) Packaging

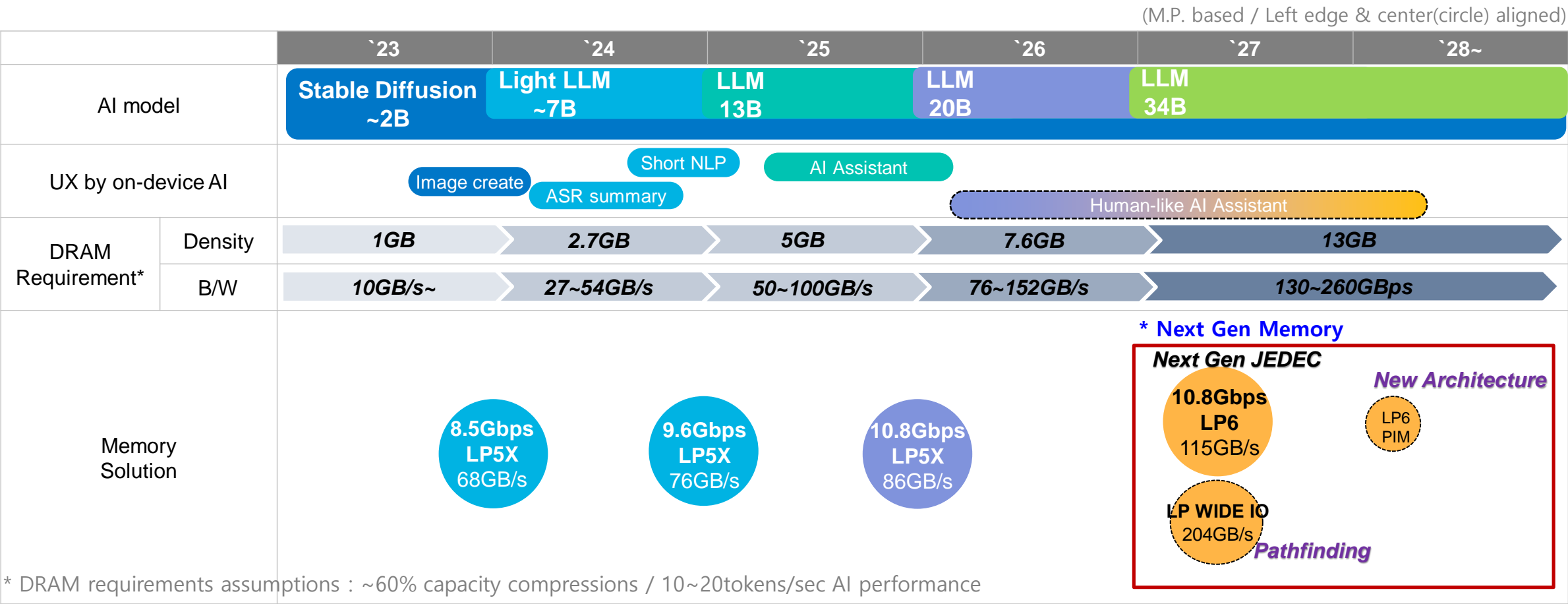
Market Trends for Artificial Intelligence

1. [~2023] New Models, Learning new capability from existing data → Training
 - 1) Objective: Build Deep Learning Framework (e.g. GPT, PaLM, LLaMA)
2. [2023~] Introducing real-life applications, Applying capability to new data → Inference
 - 1) Objective: Real life Apps or Services (e.g. ChatGPT, Bard)
 - Generative Models & On-device support



Mobile Market Driven by On-device Generative AI

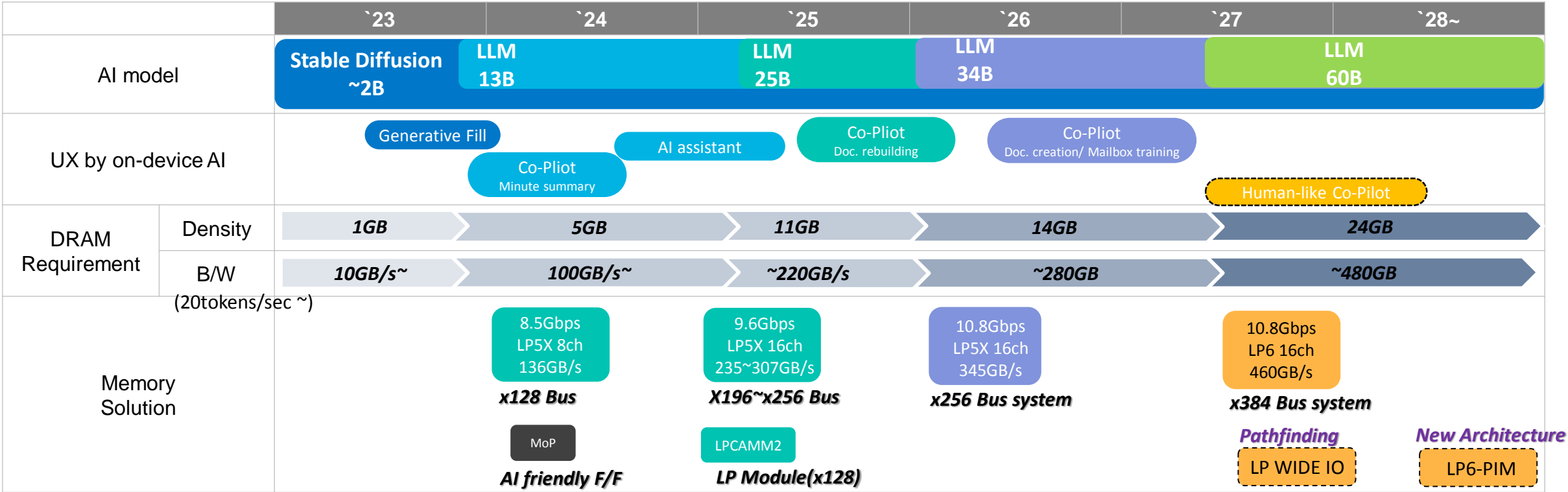
1. Mobile: Simpler queries & less accuracy → ASR for Voice assistant, light-weight photo editing for social media
- 1) Large Language Model for Mobile : Compromised model for target UX
2. LP Mem. Solution : LP5X (D1b 9.6Gbps → 10.8Gbps) → Next Gen and Pathfinding (LP6 & LP Wide IO, LP-PIM)



PC Market Driven by On-device Generative AI

1. AI-enabled PC is expected to be deployed for commercial area first and then gradually spread to consumer area
- 1) Commercial : Productivity & Creativity (e.g. Paperwork automation, Auto code creation, AI generated meeting minutes)

2) Consumer : Convenience & Entertainment (e.g. Personal assistant, Real-time avatars, AI enhanced gaming)
2. Multi-Channel system(x128↑) for sufficient BW and High-Density solution for utilizing larger model



• LLM Memory requirements assumptions
- 60% Memory footprint reduction : Quantization(Int4) + HW Compression/Accelerating logic

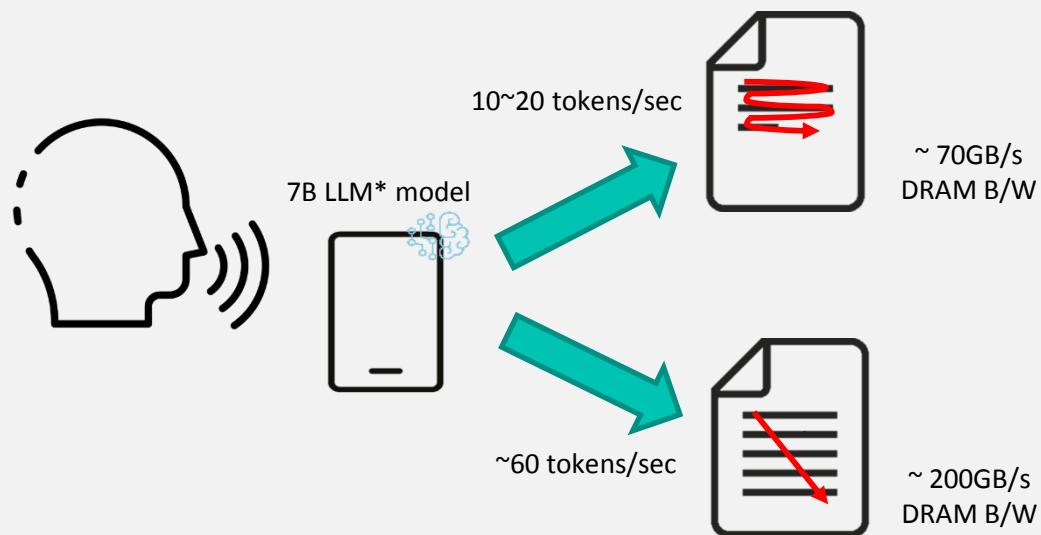
Post LP5X Memory Solutions for Generative AI

1. On-Device Gen AI (Memory bounded UX) → Higher BW & Lower power LP Memory sub-system is crucial

- 1) LPDDR6 : ~120GB/s BW, Lower power feature(VDD2D 0.875v / Efficiency mode) JEDEC standardization
→ Flagship SOCs will support LPDDR6 targeting '26.2H mass production
- 2) LP Wide IO: 204.8GB/s Total B/W, Leveraging Wide IO architecture to attain lower power (lower pin speed, 3.2Gbps)
→ Need to explore packaging solutions (e.g Wide IO level vertical wire-bonding & Fine Pitch bonding option)

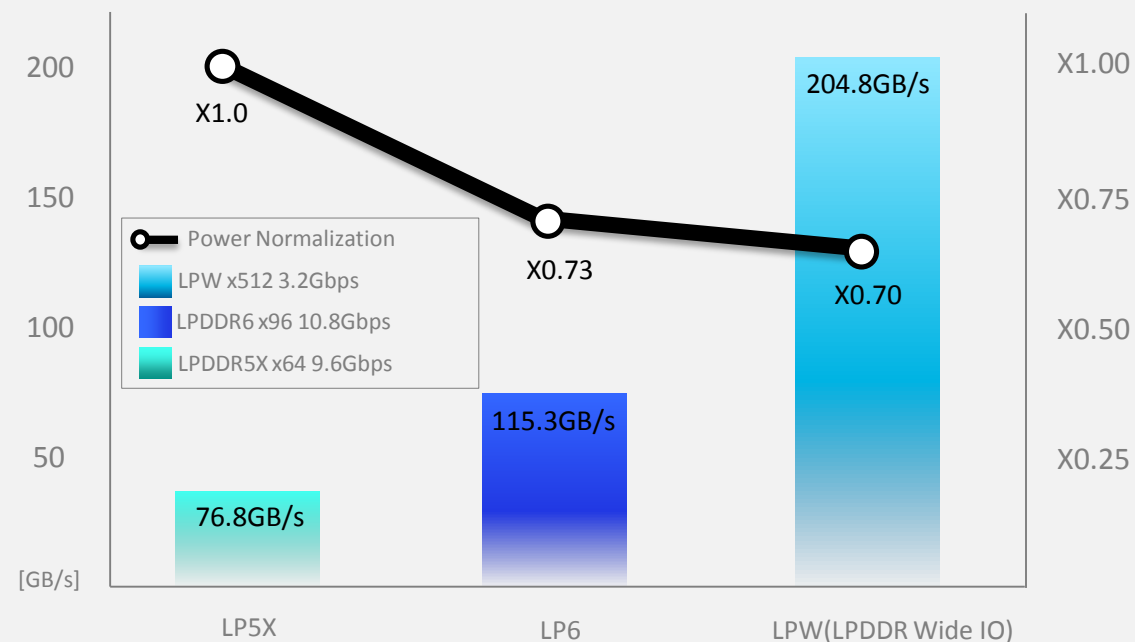
Higher Memory B/W Requirements

- ✓ Paragraph-by-paragraph fast reading and response
- Current GPT engine token generation speed : 4~6 words/sec



(* Large Language Model)

Power Efficiency& Total Bandwidth per PKG



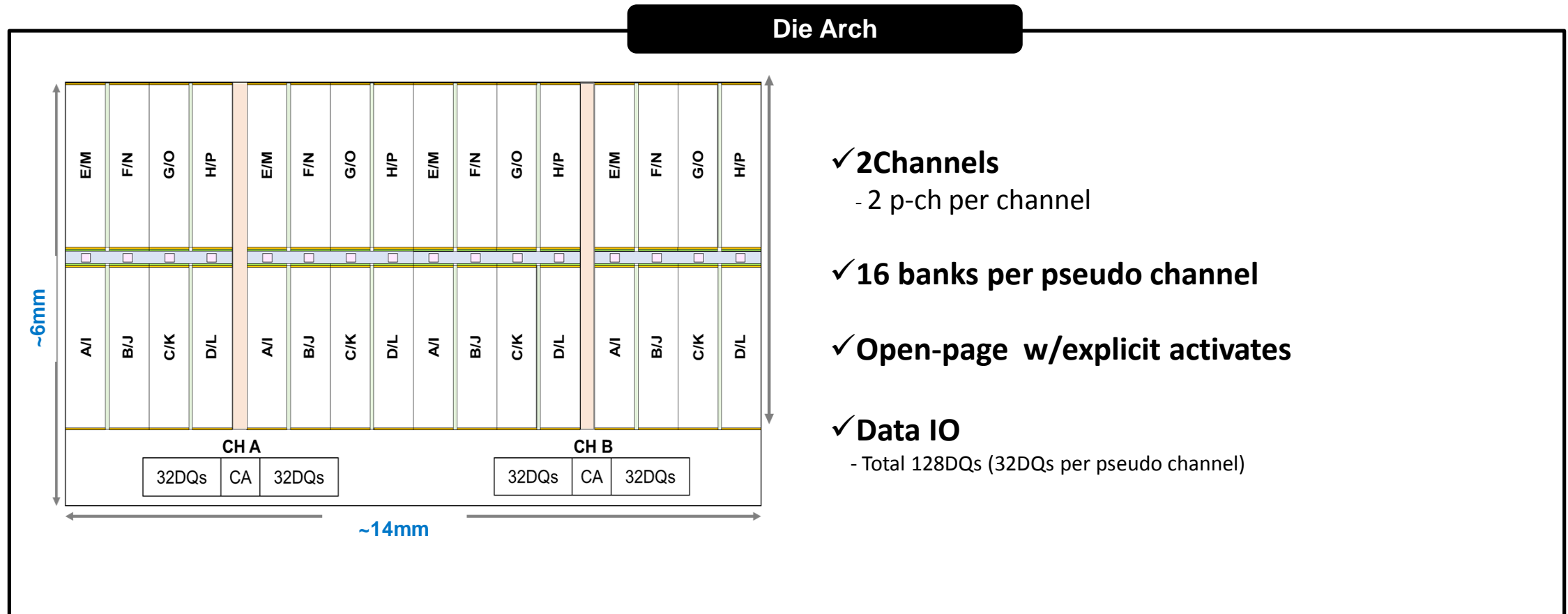
(* Compared with same process node)

Confidential

LPW(LP Wide IO) introduction

1. Samsung is now reviewing a large-capacity LP wide-IO as main memory solution

1) Capacity : 24Gb, Peak BW : 51.2GB/sec, Power efficiency : ~1.9pJ/b and DQ Count : 128 w/Max. 3.2Gbps



LPW(LP Wide-IO) : Die architecture

1. LPW die architecture is similar to LLW

But larger capacity considering main memory solution with stacking

1) Other key differences : Channel count IO, Data rate, Core timing and use of explicit activates

	LLW (Low Latency Wide IO)	LPW(LP Wide IO)	LPDDR6
Capacity and BW	1Gb, 128GB/s	24Gb, 51.2GB/s	16Gb, 28.8GB/s
Organization	4 channel 2 pseudo channels (pCH) per channel 8 banks per pCH (No bank groups) x64 DQ per pCH w/Burst length 8, 16, 32 → Min 64B Closed-page operation w/Implicit activates	2 channel 2pseudo channels (pCH) per channel 16 banks per pCH (no Bank groups) x32 DQ per pCH w/burst length 16 and 32 → Min 64B (* BL8 can be considered for Min 32B) Open-page operation w/explicit activates	2 channel 2 Sub channels (Sub-CH) per die 16 banks per Sub-CH (4Banks/4BG) x12 DQ per Sub-CH w/burst length 24 and 48 → Min 32B Open-page operation w/explicit activates
Signaling and clocking	DDR Signaling for DQ and CA 96 Mbps, 1, 2Gbps data rates 1 diff. clock per channel 1 diff. DQS per x32 DQ	DDR Signaling for DQ and CA 0.8, 1.6, 3.2Gbps data rates 1 diff. clock per channel 1 diff. DQS per x32 DQ	DDR Signaling for DQ and CA 10.8Gbps data rates (1 st target) 1 diff. clock per Sub-channel 1 diff. DQS per Sub-channel
Critical Core Timing	tRCmin=28/32ns, Latency=30ns	tRCmin=60ns, Latency(tRCDr + RL + BL16) < 42 ns	tRCmin=60ns
Refresh	Only all-bank refresh per pseudo channel	Both per-bank and all-bank refresh	Per dual-bank and all-bank refresh
Calibrations	Only at cold boot + background periodic ZQC	Only at cold boot + background periodic ZQC	background periodic ZQC
Testability	Direct Access Port, Boundary Scan	Direct Access Port, Boundary Scan	N/A
Reparability	Lane Repair, Post-Package Repair	Lane Repair, Post-Package Repair	Post-Package Repair,

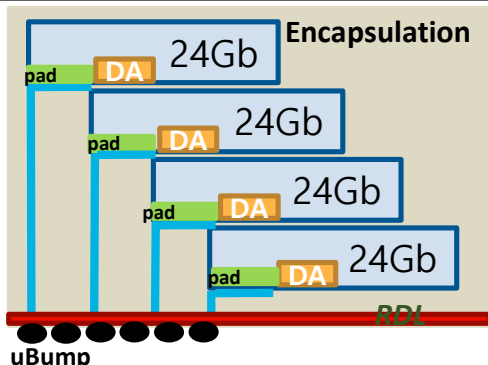
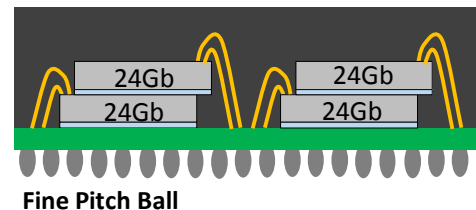
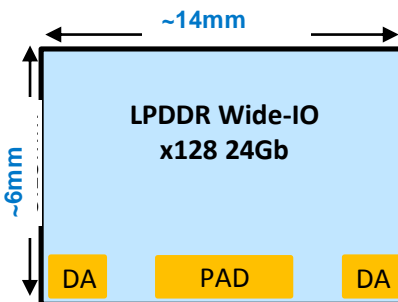
Confidential

Case study for LPW PKG options : VWB +RDL, Fine-pitch PKG, Wafer Biz

**Vertical Wire Bonding*

1. Currently, LPW is pathfinding status & now reviewing technical feasibility on packaging options

1) Case : 1. VWB + Wafer level RDL w/new equipment, 2. Fine-Pitch PKG w/ existing infra, 3. Wafer Biz

Case study			1. VWB + RDL	2. Fine-pitch PKG	3. Wafer Biz
LP Wide-IO	Concept				
	Chip die		LPW architecture	←	
Tech. feasibility	PKG	Stacking	New solution -VWB@4H stacking	D-DDP structure - x256/512PKG @wire bonding	N/A
		Pad Pitch	60+@μm	80μm -Staggered	.
		Size	Reviewing	~16.0mm x ~14.0mm @0.27ball pitch assumption	.
Risk point			VWB feasibility Testability @DA Pad is consideration	PKG Size increase • Reviewing 4H stack option to minimize Y-axis (Tentative size : ~16mm x ~9mm)	.

Confidential

A journey shared takes us beyond