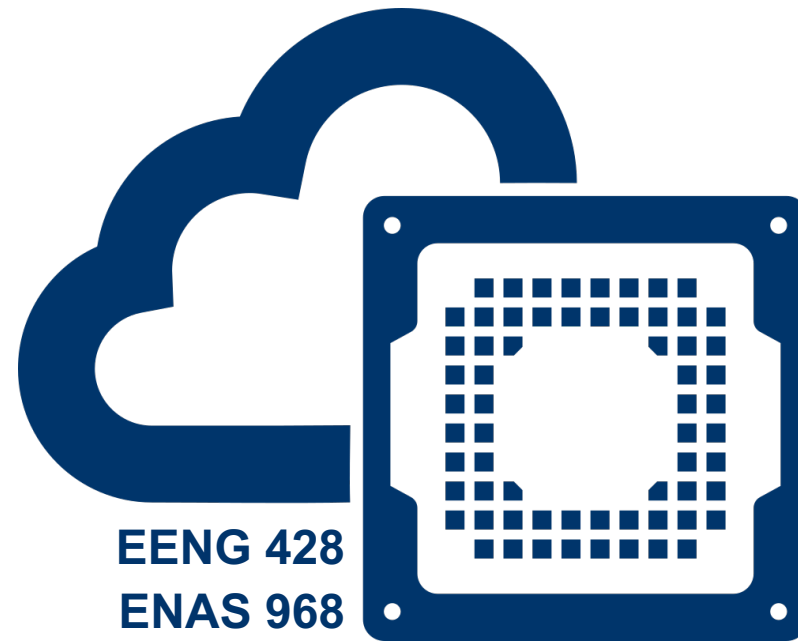


# Cloud FPGA



[bit.ly/cloudfpga](https://bit.ly/cloudfpga)



## Lecture: Intel FPGAs in Cloud FPGAs

Prof. Jakub Szefer

Dept. of Electrical Engineering, Yale University

EENG 428 / ENAS 968

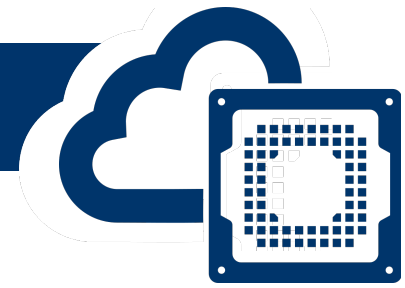
Cloud FPGA



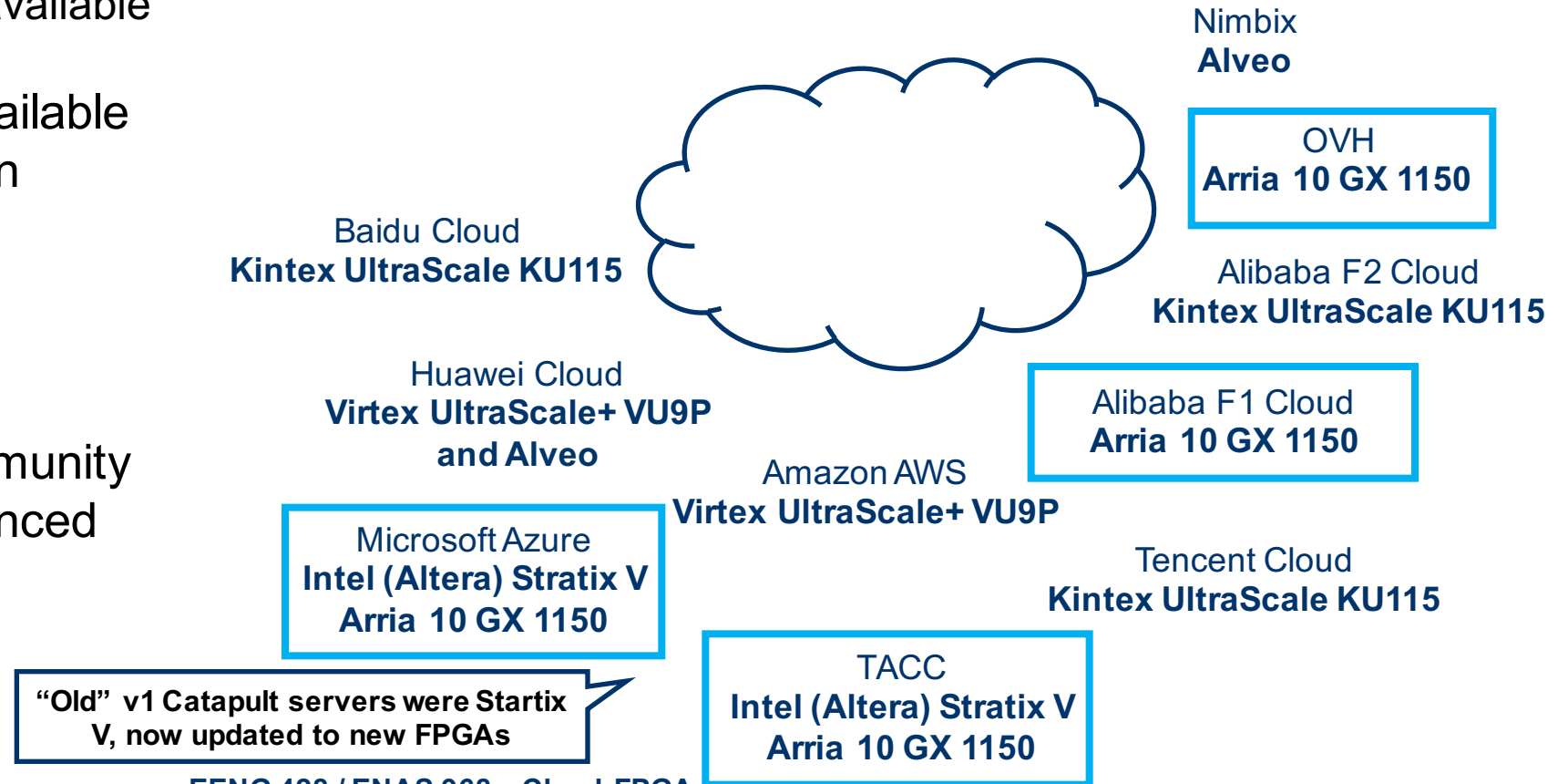
Share:  
[bit.ly/cloudfpga](http://bit.ly/cloudfpga)

EENG 428 / ENAS 968 – Cloud FPGA  
© Jakub Szefer, Fall 2019

# FPGA Chips Used in Cloud FPGAs

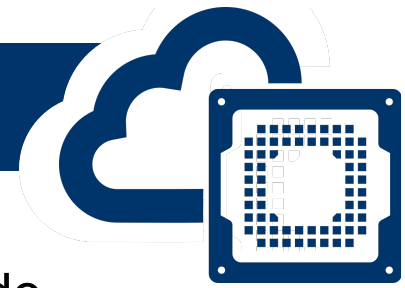


- Different Cloud FPGA providers make various FPGAs available for remote access
  - Major FPGA vendors are: **Xilinx** and **Intel** (Altera)
  - Typically, Cloud FPGAs only provide one type of FPGA board available
- Intel (Altera) FPGAs are available in public clouds as well as in Microsoft's data centers where they accelerate search and AI operations
- Available to academic community through TACC (Texas Advanced Computing Center)
  - Part of Chameleon Cloud [www.chameleoncloud.org](http://www.chameleoncloud.org)



Share:  
[bit.ly/cloudfpga](http://bit.ly/cloudfpga)

# Intel (Altera) FPGA Families

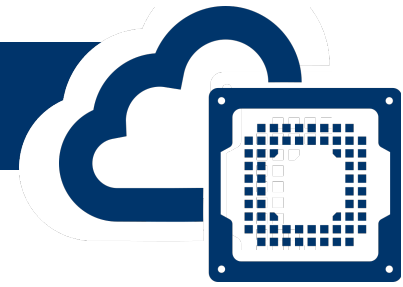


Intel acquired Altera Corporation in 2015, the Altera (now Intel) FPGA offerings include

- Cyclone – low-resource FPGAs
  - **Cyclone IV** (2009)
  - **Cyclone V** (28nm, circa 2011)
  - **Cyclone 10** (2017)
- Arria – mid-range FPGAs
  - **Aria II** (40nm, circa 2009)
  - **Aria V** (28nm, circa 2011)
  - **Aria 10** (20nm, circa 2013)
- Stratix – high-performance FPGAs
  - **Stratix IV** (40nm, circa 2008)
  - **Stratix V** (28nm, circa 2010)
  - **Stratix 10** (14nm, circa 2013)
- Embedded Arm cores:
  - Many of the FPGAs have variants with a hard IP Arm core or cores



# Metrics for Comparing FPGAs

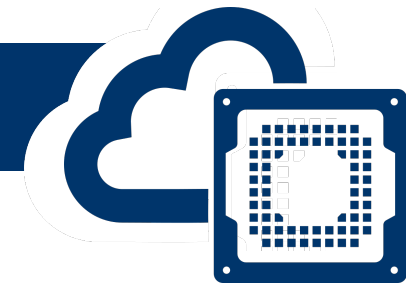


Main metrics to consider when comparing FPGAs are:

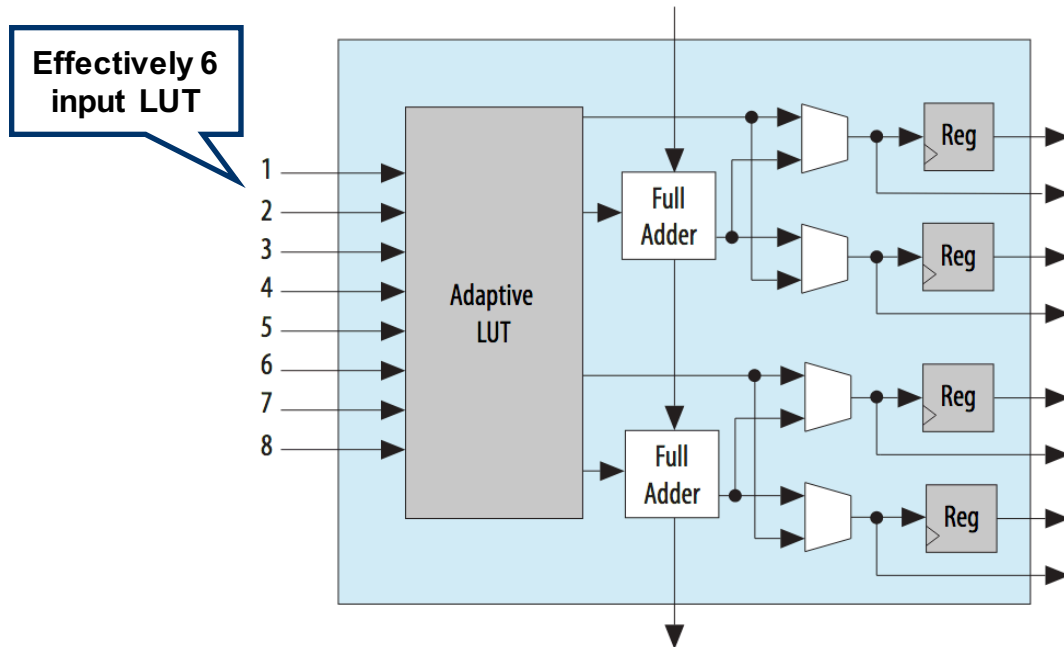
- **Number of Look-Up Tables (LUTs)** → Logic Elements (**LEs**)
- **Size of LUTs**
- **Number of Flip-Flops** → **Registers**
- **Size of Block RAM storage** → “**M20K**” blocks
- **Logic Cells** → **System Logic Elements**, rough estimate of how much logic can fit on FPGA, abstracting away details
- **Speed grade** → **Lower #** is faster, e.g., -1 is fastest and -10 is slowest, roughly corresponds to clock frequency in ns
- **I/O resources**
- **Digital Signal Processing (DSP) blocks, Floating Point blocks, etc.**
- **Other embedded Hard IP blocks or CPUs**

**Adaptive Logic Modules (ALMs)** contain logic elements (**LEs**) and **registers**

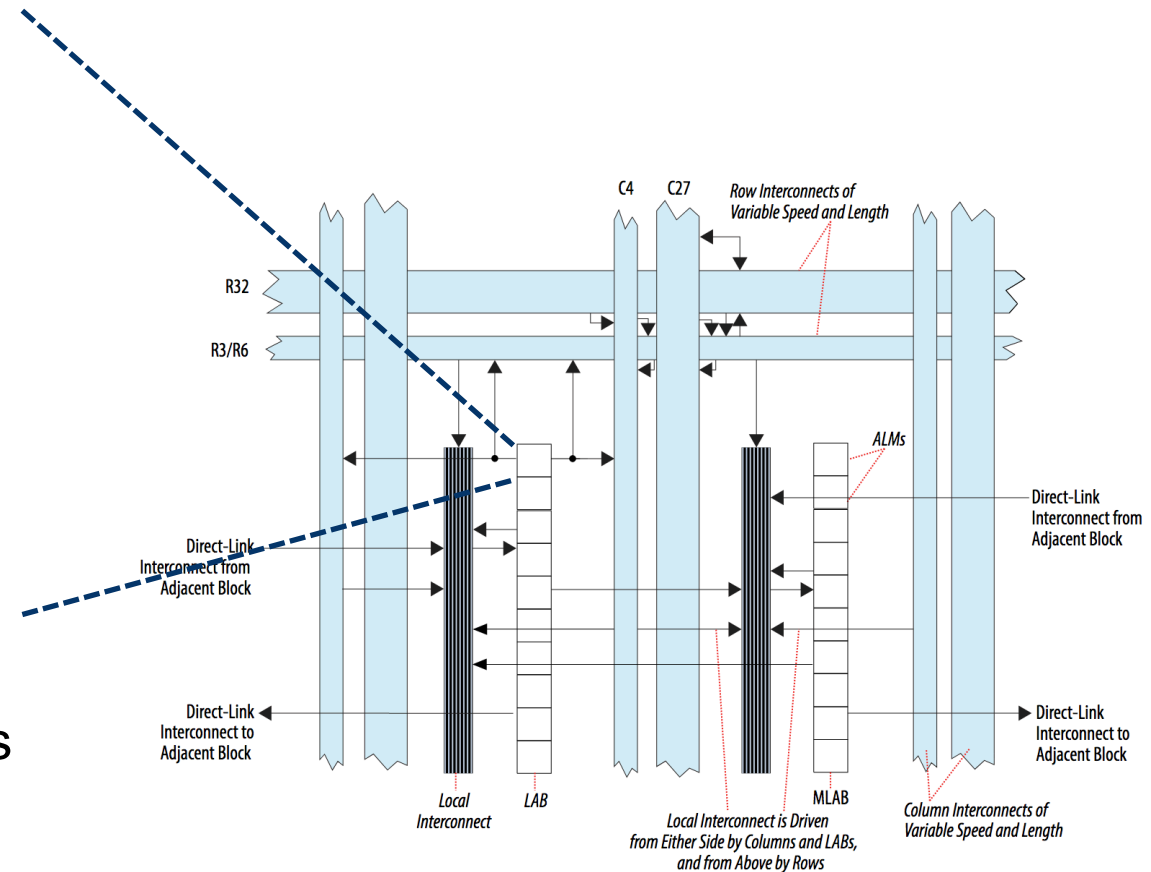
# LABs, ALMs, and LEs



- Example from Arria 10: an ALM contains LUT, adders, and registers



- ALMs are arranged into Logic Array Blocks (LABs), similar to Slices in Xilinx

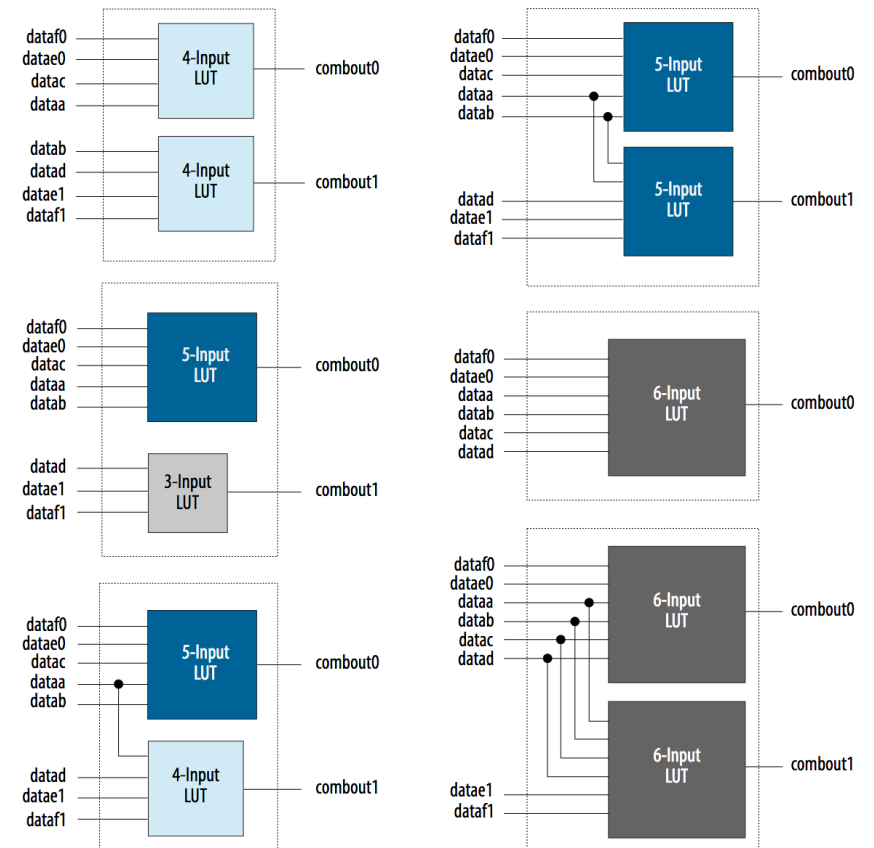
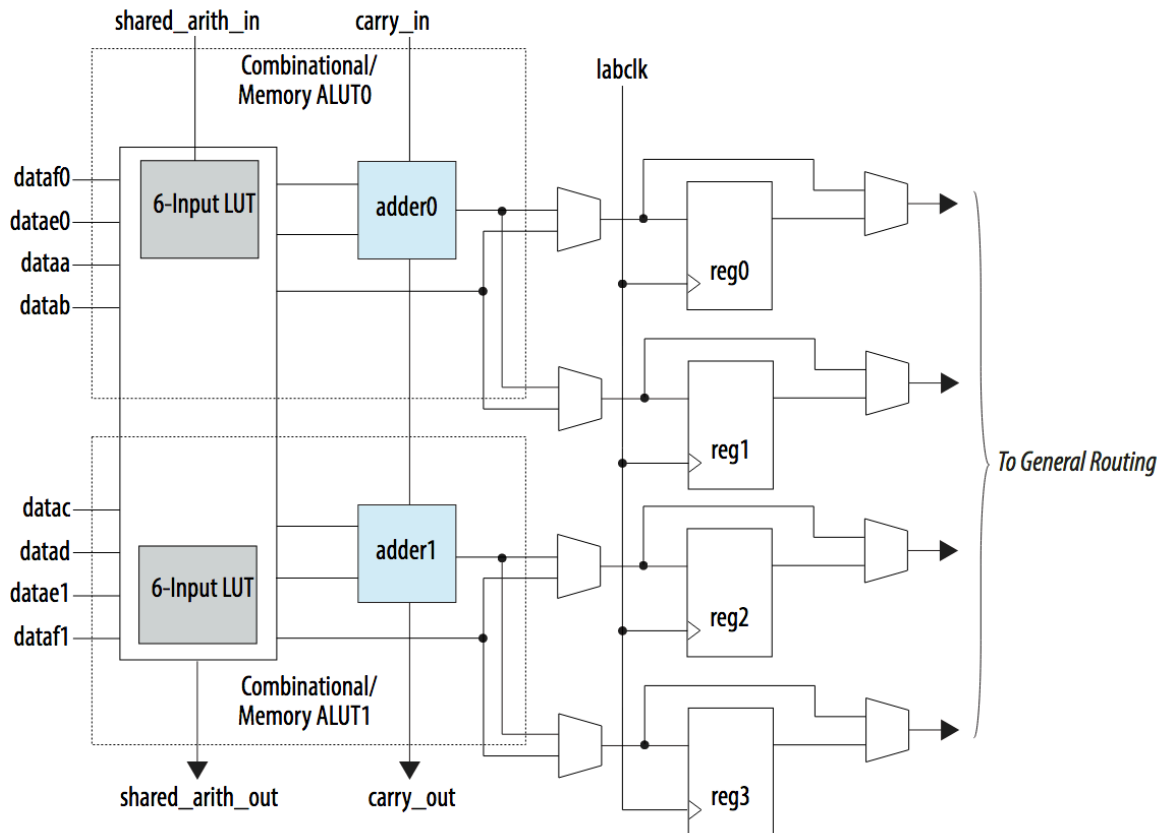


Share:  
[bit.ly/cloudfpga](http://bit.ly/cloudfpga)

# ALM Details



- Example from Arria 10: up to two functions can be implemented in one Intel Arria 10 ALM, or a single function of up to six inputs → hence effectively 6 input LUT

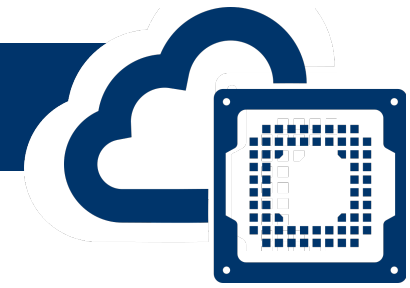


Share:  
[bit.ly/cloudfpga](http://bit.ly/cloudfpga)





# Cyclone 10 FPGAs



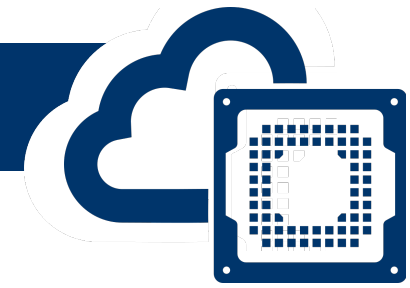
- Cyclone 10 FPGAs are lowest-end FPGA form Intel
  - GX more powerful
  - LP are smallest, about 1x to 10x smaller than GX

PRODUCT LINE		10CX085	10CX105	10CX150	10CX220
Resources	Logic elements (LEs) <sup>1</sup>	85,000	104,000	150,000	220,000
	Adaptive logic modules (ALMs)	31,000	38,000	54,770	80,330
	ALM registers	124,000	152,000	219,080	321,320
	M20K memory blocks	291	382	475	587
	M20K memory size (Kb)	5,820	7,640	9,500	11,740
	MLAB memory size (Kb)	653	799	1,152	1,690
	Variable-precision digital signal processing (DSP) blocks	84	125	156	192
	18 x 19 multipliers	168	250	312	384
	Peak fixed-point performance (GMACS) <sup>2</sup>	151	225	281	346
	Peak floating-point performance (GFLOPS) <sup>3</sup>	59	88	109	134
I/O and Architectural Features	Global clock networks	32	32	32	32
	Regional clocks	8	8	8	8
	Maximum user I/O pins	192	284	284	284
	Maximum LVDS pairs 1.4 Gbps (RX or TX)	72	118	118	118
	Maximum transceiver count (12.5 Gbps)	6	12	12	12
	Maximum 3V I/O pins	48	48	48	48
	PCI Express* (PCIe*) hard IP blocks (Gen2 x4) <sup>4</sup>	1	1	1	1
	Memory devices supported	DDR3, DDR3L, LPDDR3			



Share:  
[bit.ly/cloudfpga](https://bit.ly/cloudfpga)

# Arria 10 FPGAs



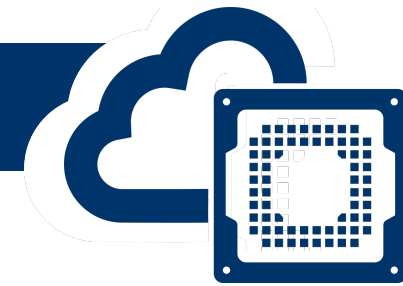
- Arria 10 are mid-range FPGAs
  - GT, GX, SX variants
  - About 1x to 10x bigger than biggest Cyclone 10

PRODUCT LINE		GX 160 SX 160	GX 220 SX 220	GX 270 SX 270	GX 320 SX 320	GX 480 SX 480	GX 570 SX 570	GX 660 SX 660	GX 900	GX 1150	GT 900	GT 1150	
Resources	LEs (K)	160	220	270	320	480	570	660	900	1,150	900	1,150	
	System logic elements (K)	210	288	354	419	629	747	865	1,180	1,506	1,180	1,506	
	Adaptive logic modules (ALMs)	61,510	83,730	101,620	118,730	181,790	217,080	250,540	339,620	427,200	339,620	427,200	
	Registers	246,040	334,920	406,480	474,920	727,160	868,320	1,002,160	1,358,480	1,708,800	1,358,480	1,708,800	
	M20K memory blocks	440	588	750	891	1,438	1,800	2,133	2,423	2,713	2,423	2,713	
	M20K memory (Mb)	9	11	15	17	28	35	42	47	53	47	53	
	MLAB memory (Mb)	1.0	1.8	2.4	2.8	4.3	5.0	5.7	9.2	12.7	9.2	12.7	
	Hardened single-precision floating-point multipliers/adders	156/156	191/191	830/830	985/985	1,368/1,368	1,523/1,523	1,688/1,688	1,518/1,518	1,518/1,518	1,518/1,518	1,518/1,518	
	18 x 19 multipliers	312	382	1,660	1,970	2,736	3,046	3,376	3,036	3,036	3,036	3,036	
	Peak fixed-point performance (GMACS) <sup>1</sup>	343	420	1,826	2,167	3,010	3,351	3,714	3,340	3,340	3,340	3,340	
Peak floating-point performance (GFLOPS)	140	172	747	887	1,231	1,371	1,519	1,366	1,366	1,366	1,366		
Clocks, Maximum I/O Pins, and Architectural Features	Global clock networks	32	32	32	32	32	32	32	32	32	32	32	
	Regional clocks	8	8	8	8	8	8	16	16	16	16	16	
	Hard processor system (available in SX devices only)	Dual-core Arm* Cortex*-A9 MPCore* processor. See the following page for details.							-				
	Maximum LVDS channels (1.6 G)	120	120	168	168	222	324	270	384	384	312	312	
	Maximum user I/O pins	288	288	384	384	492	696	696	768	768	624	624	
	Transceiver count (17.4 Gbps)	12	12	24	24	36	48	48	96	96	72	72	
	Transceiver count (25.78 Gbps)	-	-	-	-	-	-	-	-	-	6	6	
	PCIe* hardened IP blocks (Gen3 x8) <sup>2</sup>	1	1	2	2	2	2	2	4	4	4	4	
	Maximum 3 V I/O pins	48	48	48	48	48	96	96	-	-	-	-	



Share:  
[bit.ly/cloudfpga](http://bit.ly/cloudfpga)

# Stratix 10 FPGAs



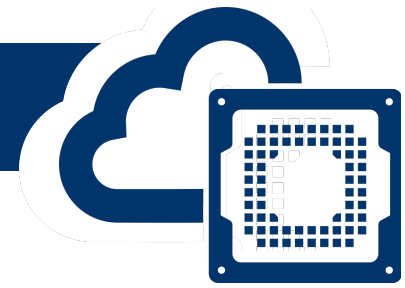
- Stratix 10 are high-end FPGAs
  - GX (baseline), SX (with quad-core Arm processors), TX (fastest I/O),
  - MX (3D stacked HBM), DX (coherence bus with host CPU)

PRODUCT LINE		GX 400 SX 400	GX 650 SX 650	GX 850 SX 850	GX 1100 SX 1100	GX 1650 SX 1650	GX 2100 SX 2100	GX 2500 SX 2500	GX 2800 SX 2800	GX 1660	GX 2110	GX 10M
Resources	Logic elements (LEs) <sup>1</sup>	378,000	612,000	841,000	1,325,000	1,624,000	2,005,000	2,422,000	2,753,000	1,679,000	2,073,000	10,200,000
	Adaptive logic modules (ALMs)	128,160	207,360	284,960	449,280	550,540	679,680	821,150	933,120	569,200	702,720	3,466,080
	ALM registers	512,640	829,440	1,139,840	1,797,120	2,202,160	2,718,720	3,284,600	3,732,480	2,276,800	2,810,880	13,864,320
	Hyper-Registers from Intel® Hyperflex™ FPGA architecture	Millions of Hyper-Registers distributed throughout the monolithic FPGA fabric										
	Programmable clock trees synthesizable	Hundreds of synthesizable clock trees										
	M20K memory blocks	1,537	2,489	3,477	5,461	5,851	6,501	9,963	11,721	6,162	6,847	12,950
	M20K memory size (Mb)	30	49	68	107	114	127	195	229	120	134	253
	MLAB memory size (Mb)	2	3	4	7	8	11	13	15	9	11	55
	Variable-precision digital signal processing (DSP) blocks	648	1,152	2,016	2,592	3,145	3,744	5,011	5,760	3,326	3,960	3,456
	18 x 19 multipliers	1,296	2,304	4,032	5,184	6,290	7,488	10,022	11,520	6,652	7,920	6,912
Peak fixed-point performance (TMACS) <sup>2</sup>	2.6	4.6	8.1	10.4	12.6	15.0	20.0	23.0	13.3	15.8	13.8	
Peak floating-point performance (TFLOPS) <sup>3</sup>	1.0	1.8	3.2	4.1	5.0	6.0	8.0	9.2	5.3	6.3	5.5	
Secure device manager	AES-256/SHA-256 bitstream encryption/authentication, physically unclonable function (PUF), ECDSA 256/384 boot code authentication, side channel attack protection											
I/O and Architectural Features	Hard processor system <sup>4</sup>	Quad-core 64-bit ARM® Cortex®-A53 up to 1.5 GHz with 32KB I/D cache, NEON coprocessor, 1 MB L2 Cache, direct memory access (DMA), system memory management unit, cache coherency unit, hard memory controllers, USB 2.0 x2, 1G EMAC x3, UART x2, SPI x4, I2C x5, general purpose timers x7, watchdog timer x4										
	Maximum user I/O pins	SX 400	SX 650	SX 850	SX 1100	SX 1650	SX 2100	SX 2500	SX 2800	-	-	-
	Maximum LVDS pairs 1.6 Gbps (RX or TX)	392	392	688	688	704	704	1160	1160	688	688	2,304
	Maximum LVDS pairs 1.6 Gbps (RX or TX)	192	192	336	336	336	336	576	576	336	336	1152 <sup>5</sup>
	Total full duplex transceiver count	24	24	48	48	96	96	96	96	48	48	48
	GXT full duplex transceiver count (up to 28.3 Gbps)	16	16	32	32	64	64	64	64	32	32	-
	GX full duplex transceiver count (up to 17.4 Gbps)	8	8	16	16	32	32	32	32	16	16	48
	PCI Express® (PCIe®) hard intellectual property (IP) blocks (Gen3 x16)	1	1	2	2	4	4	4	4	2	2	4 <sup>6</sup>
	Memory devices supported	DDR4, DDR3, DDR2, DDR, QDR II, QDR II+, RLDram II, RLDram 3, HMC, MoSys										

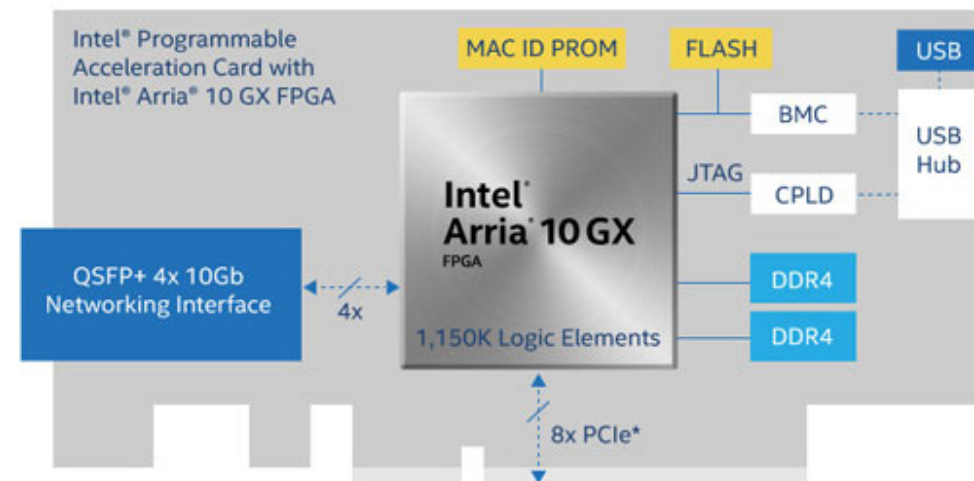


Share:  
[bit.ly/cloudfpga](https://bit.ly/cloudfpga)

# Programmable Accelerator Cards (PACs)

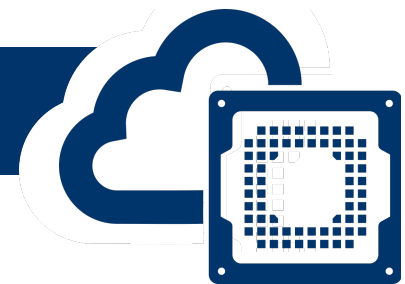


- Intel D5005 Programmable Acceleration Card (2019)
  - Uses Intel's 14nm Stratix 10 SX FPGA architecture
  - 2,800,000 logic elements
  - Thermal Design Power (TDP) of 215W
  - PCIe 3.0 and DDR4
  - QSFP networking
- Different card versions with GT (fastest, biggest), GX, or SX (slowest, smallest)
- Embedded Arm processor  
Hard IP core



Share:  
[bit.ly/cloudfpga](http://bit.ly/cloudfpga)

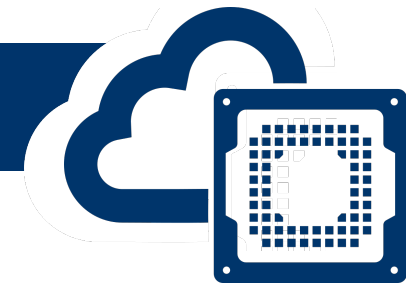
# Stratix V and Aria 10 FPGAs in Cloud FPGAs



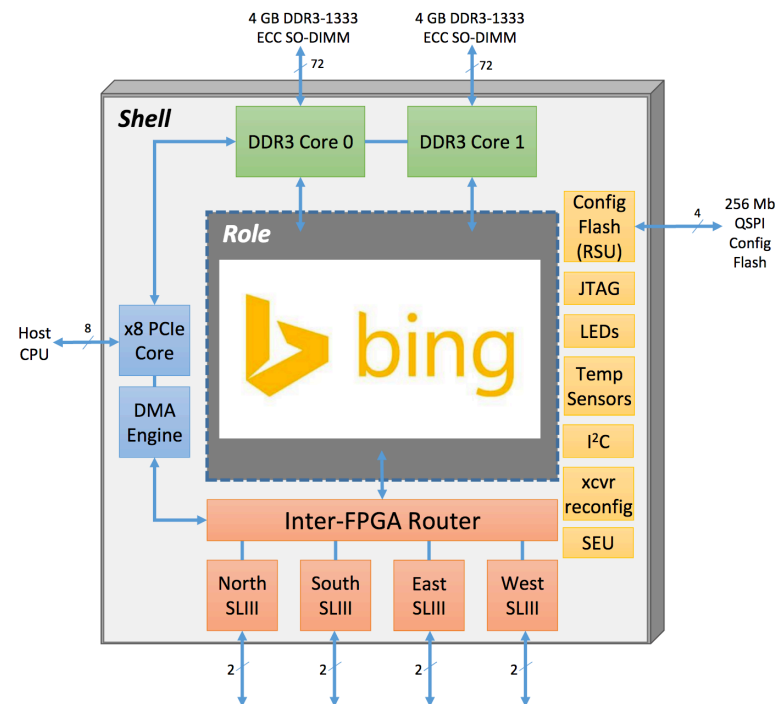
Share:  
[bit.ly/cloudfpga](https://bit.ly/cloudfpga)

EENG 428 / ENAS 968 – Cloud FPGA  
© Jakub Szefer, Fall 2019

# Project Catapult from Microsoft



- One of the first Cloud FPGA deployments
- Catapult FPGA Accelerator Card (Microsoft + Intel FPGAs)
  - Altera Stratix V GS D5
  - 172,000 ALMs
  - PCIe 3.0 and DDR3
- Used to accelerate Bing search, circa 2011
  - Defined 'shell' and 'role' ideas for splitting logic into two parts, where users control 'role'

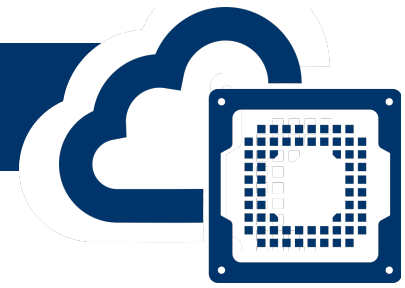


Not a public Cloud FPGA in same sense as Amazon F1, but one of first offerings of FPGAs in data centers



Share:  
[bit.ly/cloudfpga](http://bit.ly/cloudfpga)

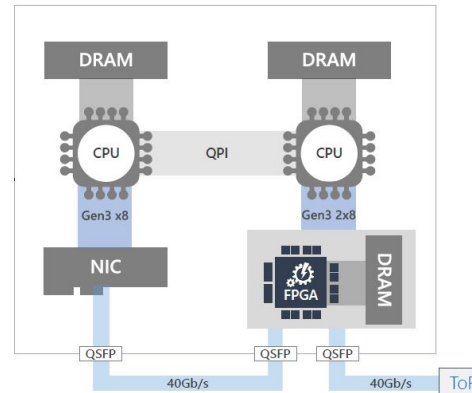
# Projects Catapult and Brainwave



- Project Catapult transformed into project Brainwave for accelerating AI on FPGAs in data centers
- Catapult v1 – original FPGA accelerator card, Stratix V
- Catapult v2 – FPGA card with direct access to network, Stratix V
  - FPGA can process network packets
  - FPGA can bypass CPU and get data more quickly from network



Catapult v2 Mezzanine card

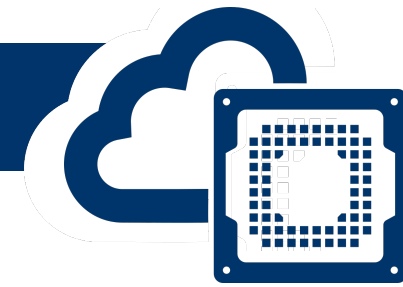


- Brainwave project with Arria 10 boards (and Stratix 10)

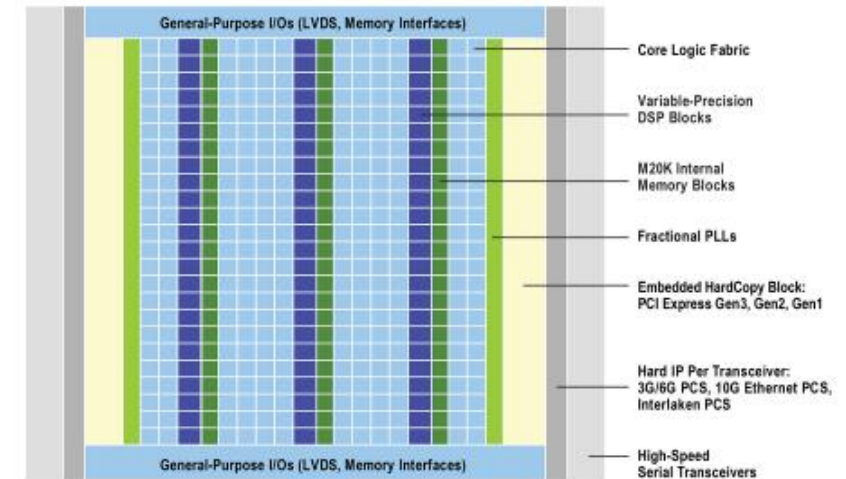




# Stratix V FPGAs



- Catapult v2 uses Stratix V D5 FPGAs
  - Column-style FPGA like most modern chips
  - Mid-range Stratix V FPGA
  - Relatively large M20K size
  - High number of DSP and multiplier blocks



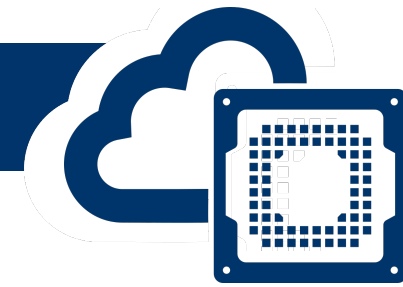
PRODUCT LINE	STRATIX V GS FPGAs <sup>1</sup>					STRATIX V GX FPGAs <sup>1</sup>										STRATIX V E FPGAs <sup>1</sup>	
	5SGSD3	5SGSD4	5SGSD5	5SGSD6	5SGSD8	5SGXA3	5SGXA4	5SGXA5	5SGXA7	5SGXA9	5SGXAB	5SGXB5	5SGXB6	5SGXB9	5SGXBB	5SEE9	5SEEB
LEs (K)	236	360	457	583	695	340	420	490	622	840	952	490	597	840	952	840	952
ALMs	89,000	135,840	172,600	220,000	262,400	128,300	158,500	185,000	234,720	317,000	359,200	185,000	225,400	317,000	359,200	317,000	359,200
Registers	356,000	543,360	690,400	880,000	1,049,600	513,200	634,000	740,000	938,880	1,268,000	1,436,800	740,000	901,600	1,268,000	1,436,800	1,268,000	1,436,800
M20K memory blocks	688	957	2,014	2,320	2,567	957	1,900	2,304	2,560	2,640	2,640	2,100	2,660	2,640	2,640	2,640	2,640
M20K memory (Mb)	13	19	39	45	50	19	37	45	50	52	52	41	52	52	52	52	52
MLAB memory (Mb)	2.72	4.15	5.27	6.71	8.01	3.92	4.84	5.65	7.16	9.67	10.96	5.65	6.88	9.67	10.96	9.67	10.96
Variable-precision DSP blocks	600	1,044	1,590	1,775	1,963	256	256	256	256	352	352	399	399	352	352	352	352
18 x 18 multipliers	1,200	2,088	3,180	3,550	3,926	512	512	512	512	704	704	798	798	704	704	704	704



Share:  
[bit.ly/cloudfpga](http://bit.ly/cloudfpga)



# Arria 10 GX FPGAs



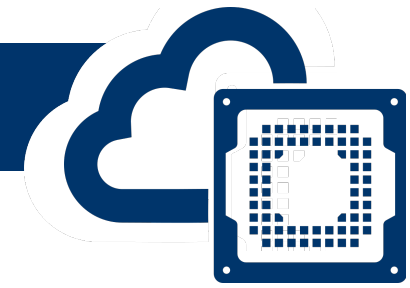
- OVH, Alibaba, and Microsoft cloud deployments use Arria 10 GX
  - 2x ~ 3x resources of Stratix V
  - Bigger multipliers, FP hard IP cores,

PRODUCT LINE		GX 160 SX 160	GX 220 SX 220	GX 270 SX 270	GX 320 SX 320	GX 480 SX 480	GX 570 SX 570	GX 660 SX 660	GX 900	GX 1150	GT 900	GT 1150	
Resources	LEs (K)	160	220	270	320	480	570	660	900	1,150	900	1,150	
	System logic elements (K)	210	288	354	419	629	747	865	1,180	1,506	1,180	1,506	
	Adaptive logic modules (ALMs)	61,510	83,730	101,620	118,730	181,790	217,080	250,540	339,620	427,200	339,620	427,200	
	Registers	246,040	334,920	406,480	474,920	727,160	868,320	1,002,160	1,358,480	1,708,800	1,358,480	1,708,800	
	M20K memory blocks	440	588	750	891	1,438	1,800	2,133	2,423	2,713	2,423	2,713	
	M20K memory (Mb)	9	11	15	17	28	35	42	47	53	47	53	
	MLAB memory (Mb)	1.0	1.8	2.4	2.8	4.3	5.0	5.7	9.2	12.7	9.2	12.7	
	Hardened single-precision floating-point multipliers/adders	156/156	191/191	830/830	985/985	1,368/1,368	1,523/1,523	1,688/1,688	1,518/1,518	1,518/1,518	1,518/1,518	1,518/1,518	
	18 x 19 multipliers	312	382	1,660	1,970	2,736	3,046	3,376	3,036	3,036	3,036	3,036	
	Peak fixed-point performance (GMACS) <sup>1</sup>	343	420	1,826	2,167	3,010	3,351	3,714	3,340	3,340	3,340	3,340	
Peak floating-point performance (GFLOPS)	140	172	747	887	1,231	1,371	1,519	1,366	1,366	1,366	1,366		
Clocks, Maximum I/O Pins, and Architectural Features	Global clock networks	32	32	32	32	32	32	32	32	32	32	32	
	Regional clocks	8	8	8	8	8	8	16	16	16	16	16	
	Hard processor system (available in SX devices only)	Dual-core Arm* Cortex*-A9 MPCore* processor. See the following page for details.											
	Maximum LVDS channels (1.6 G)	120	120	168	168	222	324	270	384	384	312	312	
	Maximum user I/O pins	288	288	384	384	492	696	696	768	768	624	624	
	Transceiver count (17.4 Gbps)	12	12	24	24	36	48	48	96	96	72	72	
	Transceiver count (25.78 Gbps)	–	–	–	–	–	–	–	–	–	6	6	
	PCIe* hardened IP blocks (Gen3 x8) <sup>2</sup>	1	1	2	2	2	2	2	4	4	4	4	
	Maximum 3 V I/O pins	48	48	48	48	48	96	96	–	–	–	–	

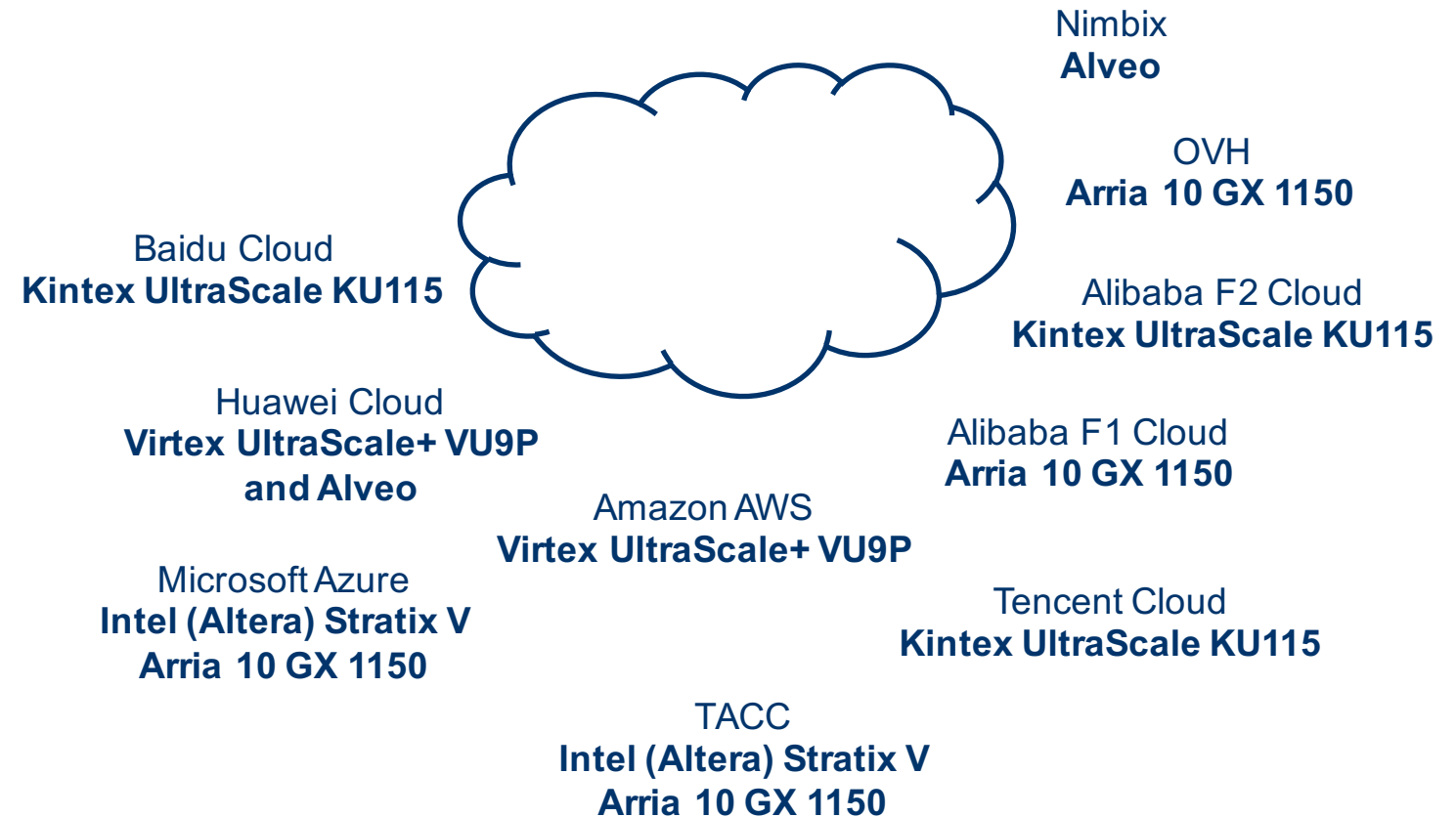


Share:  
[bit.ly/cloudfpga](http://bit.ly/cloudfpga)

# Summary

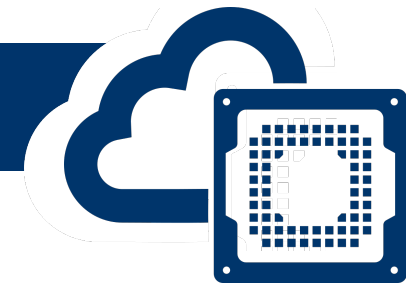


- Different Cloud FPGA vendors give access to different types of FPGAs
- Intel (Altera) FPGAs are available in public clouds as well as in Microsoft's data centers where they accelerate search and AI operations



Share:  
[bit.ly/cloudfpga](http://bit.ly/cloudfpga)

# References



1. “Intel Arria 10 Product Table”. Available at: <https://www.intel.com/content/dam/www/programmable/us/en/pdfs/literature/pt/arria-10-product-table.pdf>.
2. “Intel Arria 10 Device Overview”. Available at: [https://www.intel.com/content/dam/www/programmable/us/en/pdfs/literature/hb/arria-10/a10\\_overview.pdf](https://www.intel.com/content/dam/www/programmable/us/en/pdfs/literature/hb/arria-10/a10_overview.pdf).
3. “Intel Arria 10 Core Fabric and General Purpose I/Os Handbook”. Available at: [https://www.intel.com/content/dam/www/programmable/us/en/pdfs/literature/hb/arria-10/a10\\_handbook.pdf](https://www.intel.com/content/dam/www/programmable/us/en/pdfs/literature/hb/arria-10/a10_handbook.pdf).
4. “Intel Programmable Acceleration Card (PAC) with Arria 10 GX FPGA”. Available at: <https://www.digikey.com/en/product-highlight/i/intel/programmable-acceleration-card>.
5. “Intel Cyclone 10 GX FPGAs Product Table”. Available at: <https://www.intel.com/content/dam/www/programmable/us/en/pdfs/literature/pt/cyclone-10-gx-product-table.pdf>.
6. “Large-Scale Reconfigurable Computing in a Microsoft Datacenter”. Available at: <https://www.microsoft.com/en-us/research/uploads/prod/2014/06/HC26.12.520-Recon-Fabric-Pulnam-Microsoft-Catapult.pdf>.
7. “Microsoft Shows off Project BrainWave Persistent Inferencing from FPGA Cache”. Available at: <https://www.servethehome.com/microsoft-shows-off-project-brainwave-persistent-inferencing-from-fpga-cache/>.
8. “Intel Stratix 10 GX/SX Product Table”. Available at: <https://www.intel.com/content/dam/www/programmable/us/en/pdfs/literature/pt/stratix-10-product-table.pdf>.
9. “Intel Stratix V”. Available at: <https://www.intel.com/content/www/us/en/products/programmable/fpga/stratix-V.html>.

