

# MARKET LANDSCAPE FOR AI HARDWARE

KARL FREUND

KARL@MOORINSIGHTSSTRATEGY.COM

SR. ANALYST, AI AND HPC

MOOR INSIGHTS & STRATEGY

FOLLOW MY BLOGS COVERING MACHINE LEARNING HARDWARE  
ON FORBES: [HTTP://WWW.FORBES.COM/SITES/MOORINSIGHTS](http://www.forbes.com/sites/moorinsights) OR  
[HTTPS://MUCKRACK.COM/KARL-FREUND/ARTICLES](https://muckrack.com/karl-freund/articles)  
[WWW.MOORINSIGHTSSTRATEGY.COM](http://www.moorinsightsstrategy.com)

## GOALS OF THIS SESSION

1. Lay the foundation for the next two days
2. Review AI Strategy from the major vendors
3. Introduce you to a few new companies

# THE TWO WORLDS OF AI PROCESSING

## 1. Training a Network

- Runs take days or even weeks
- Trillions of billions of ops
- Massive data sets



## 2. Inference Processing

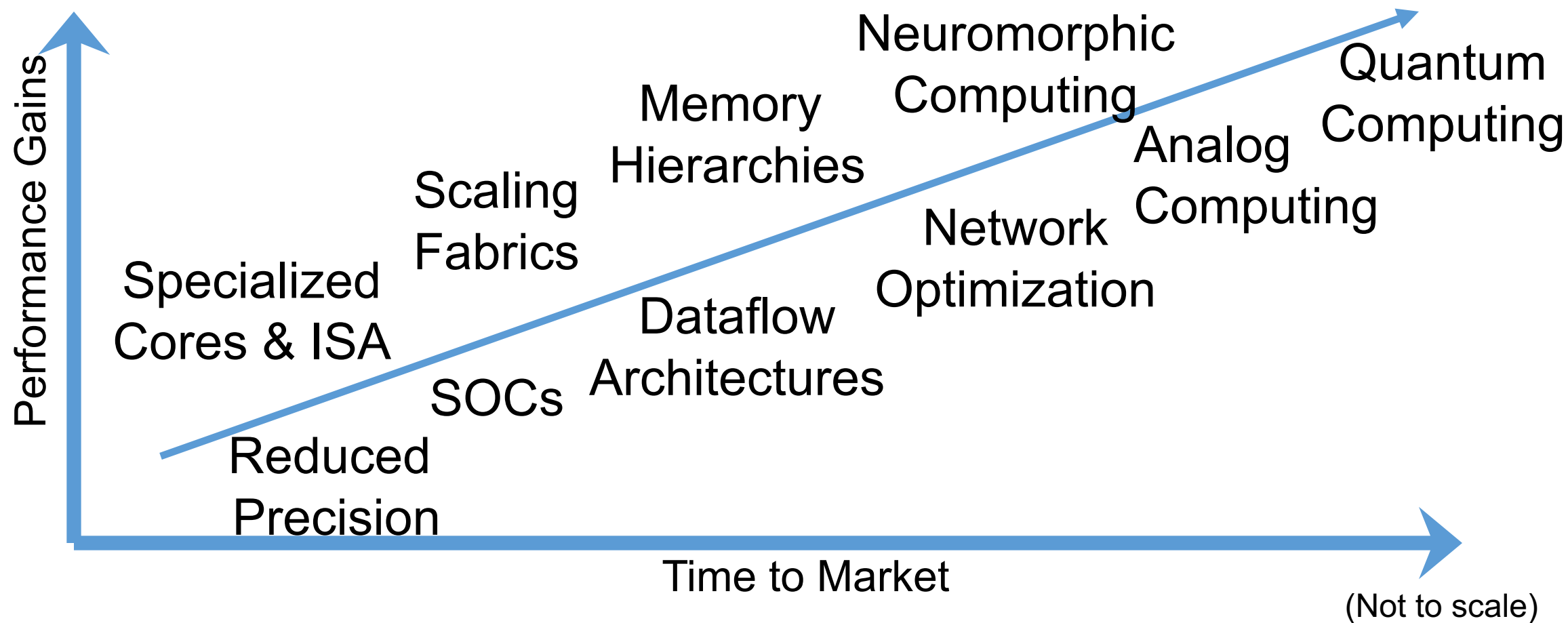
- Runs take milliseconds
- Huge range of requirements
- Will become the larger market



# HARDWARE CHOICES FOR AI (DEEP LEARNING)

- **CPU's:** Good for inference and model prototyping
  - Use the CPU's you already have
- **FPGA's:** Very good for Inference and other tasks (See Microsoft)
  - Difficult to master.
- **GPU's:** Very flexible accelerators.
  - Great for training and demanding inference workloads (Conversational AI)
  - Programmable for new models, new networks
  - Massive ecosystem already in place
- **ASIC's:** Great potential in AI, however ...
  - Requires significant R&D and time; need high volumes to break-even
  - Entails long development cycles in a fast-moving ecosystem
  - Each ASIC will need an ecosystem of software, models, researchers
  - Not a fungible asset in the cloud

# IMPROVING AI PROCESSING



# THE COMPETITIVE LANDSCAPE

## ***“A CAMBRIAN EXPLOSION IN CHIPS”***

[HTTPS://WWW.FORBES.COM/SITES/MOORINSIGHTS/2019/01/23/2019-A-CAMBRIAN-EXPLOSION-IN-DEEP-LEARNING-PART-1/#3CE8172F4DC1](https://www.forbes.com/sites/moorinsights/2019/01/23/2019-a-cambrian-explosion-in-deep-learning-part-1/#3CE8172F4DC1)

# NVIDIA: LEADER IN TRAINING, ALSO TARGETING INFERENCE

- GPU's benefit from programmability
  - Mature SW stack and development platform
- Volta delivers industry-leading 125 TOPS today
- T4 GPU is gaining inference traction (Google & AWS)
- Xavier SOC : flexible platform for edge AI
  - “DLA” is open-source ASIC for CNNs
- All clouds, OEM's, Universities, SW Startups



## INTEL HAS A LOT OF AI FIREPOWER

- New Xeons w/ 11x improvement in AI performance
- Nervana NNPs coming soon (training & inference)
- Acquired Mobileye for AVs and Movidius for edge vision
- Intel's future Exascale Engine (**X<sup>e</sup>**) could be game changer
  - Very few details have been disclosed.

<https://www.forbes.com/sites/moorinsights/2019/03/19/intel-and-cray-reaffirm-first-usa-exascale-supercomputer/#6bcd6fc29df>





# QUALCOMM

*(YES! QUALCOMM!)*



- Strong legacy of performance, power efficiency, volume and AI technology
  - Snapdragon now includes 4<sup>th</sup> generation AI Engine
- Qualcomm has announced their intentions to extend their presence in AI into the Data Center (Cloud AI100)
  - “Distributed Intelligence” Vision
- Watch this space

<https://www.forbes.com/sites/moorinsights/2019/04/10/qualcomm-dives-into-the-deep-end-of-the-data-center/#710d9163367f>

# FPGA'S: FLEXIBLE, PROGRAMMABLE INFERENCE

- HW programmable accelerators from Intel and Xilinx
- Microsoft champions FPGAs in datacenter
- Xilinx and AWS offer pre-built solutions
- AWS offers Xilinx as a Service
- FlexLogix is developing an Inference-specific FPGA (2020)



<https://www.forbes.com/sites/moorinsights/2019/08/28/xilinx-reveals-more-versal-details/#21b1147e330d>

# AI PROCESSING IN PUBLIC CLOUDS

- AWS has the large public cloud AI portfolio
  - Offers extensive suite of AI development tools & networks, many to support Alexa
  - Is building their own Inference chip (“Inferentia”)
  - Offers an app marketplace and F1 instances for Xilinx FPGAs
- Microsoft has invested heavily in GPUs and FPGA’s
  - Extensive API library with pre-trained neural networks
  - FPGA’s currently for internal use across MSFT infrastructure.
- Google has the TPU and a huge AI team
  - Tens of thousands engineers and scientists
- Alibaba and Baidu intend to build out indigenous Industry
  - Baidu has announced their own chips, Alibaba has founded a chip company



# COMPANY TARGET MARKETS (PARTIAL LIST)

## Training

Qualcomm  
Apple  
NVIDIA

NVIDIA  
Google  
Intel  
Baidu  
Huawei

Cerebras  
Groq  
Graphcore

## Inference

NVIDIA  
Intel  
Xilinx  
Qualcomm  
Apple  
Google  
Tesla  
NovuMind

Huawei  
Baidu  
FlexLogix  
Rain  
Tenstorrent  
GraphCore

Cambricon  
Mythic  
SambaNova  
Groq  
Gyr Falcon  
Cornami

Horizon Robotics  
Thinci  
Hailo  
Brainchip  
Syntient  
Eta Compute

NVIDIA  
Google  
Intel  
Xilinx  
Qualcomm  
~All Cloud Providers

Habana Labs  
Groq  
FlexLogix  
Cerebras  
Graphcore

Groq  
Wave  
AMD

## Edge

## Data Center

# CHIPS IN PRODUCTION\*

\* Does not Imply Adoption

Training

NVIDIA  
Google

Inference

NVIDIA  
Intel  
Xilinx  
Qualcomm  
Apple  
Google  
Tesla

Huawei  
Baidu

Cambricon

Horizon  
Robotics  
Thinci (?)

NVIDIA  
Google  
Intel  
Xilinx

Habana Labs

Edge

Data Center

## FINALLY, A FEW THOUGHTS...

1. **TOPS = Tremendously Overused Performance Stat**
  - Network specific benchmarks are better
  - Mlperf may become the solution, all vendors should publish!
2. Data Center AI demands programmability
3. Never, ever underestimate NVIDIA (or INTEL)
  - Nobody saw TensorCores coming.

STAY TUNED FOR WHAT'S NEXT!      @karlfreund

# THANK YOU!

KARL FREUND, SR. ANALYST, AI AND HPC

MOOR INSIGHTS & STRATEGY

[KARL@MOORINSIGHTSSTRATEGY.COM](mailto:KARL@MOORINSIGHTSSTRATEGY.COM)

@KARLFREUND

[WWW.MOORINSIGHTSSTRATEGY.COM](http://WWW.MOORINSIGHTSSTRATEGY.COM)