# Virtuous Cycle of AI

Pradeep K Dubey

Intel Senior Fellow and Fellow of IEEE

Director Parallel Computing Lab, Intel Labs

3rd Annual Meeting of Heterogenous Integration Roadmap (IEEE EPS Chapter)
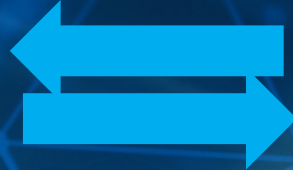Milpitas, California Feb 20, 2020

Machines:
Number Crunching
AND
Decision Making

Division of Labor Between Man and Machine Is Getting Disrupted:
*Faster than Anyone Predicted!*

# Virtuous Demand Cycle of AI Compute

# Virtuous Demand Cycle of Compute: A Functional View



Sense

Act

Reason

# AI Compute Needs

AlphaGo Zero needs:
2 EF-days to train →
Need a 100 ExaFlop machine
to train within an hour *



**AlexNet to AlphaGo Zero: A 300,000x Increase in Compute**

(intel)

# Illustrating population-scale AI in health sector

## Explosion of trans-OMICS data in the coming decade

➢ From the molecular level to the genomic level to the organ level
➢ Interactions with medications, nutrients, therapeutic devices, and the environment.



DNA
(**gen**omics)

**transcription**  translation  folding

RNA
(**transcript**omics)

amino acid chain

Protein
(**prote**omics)
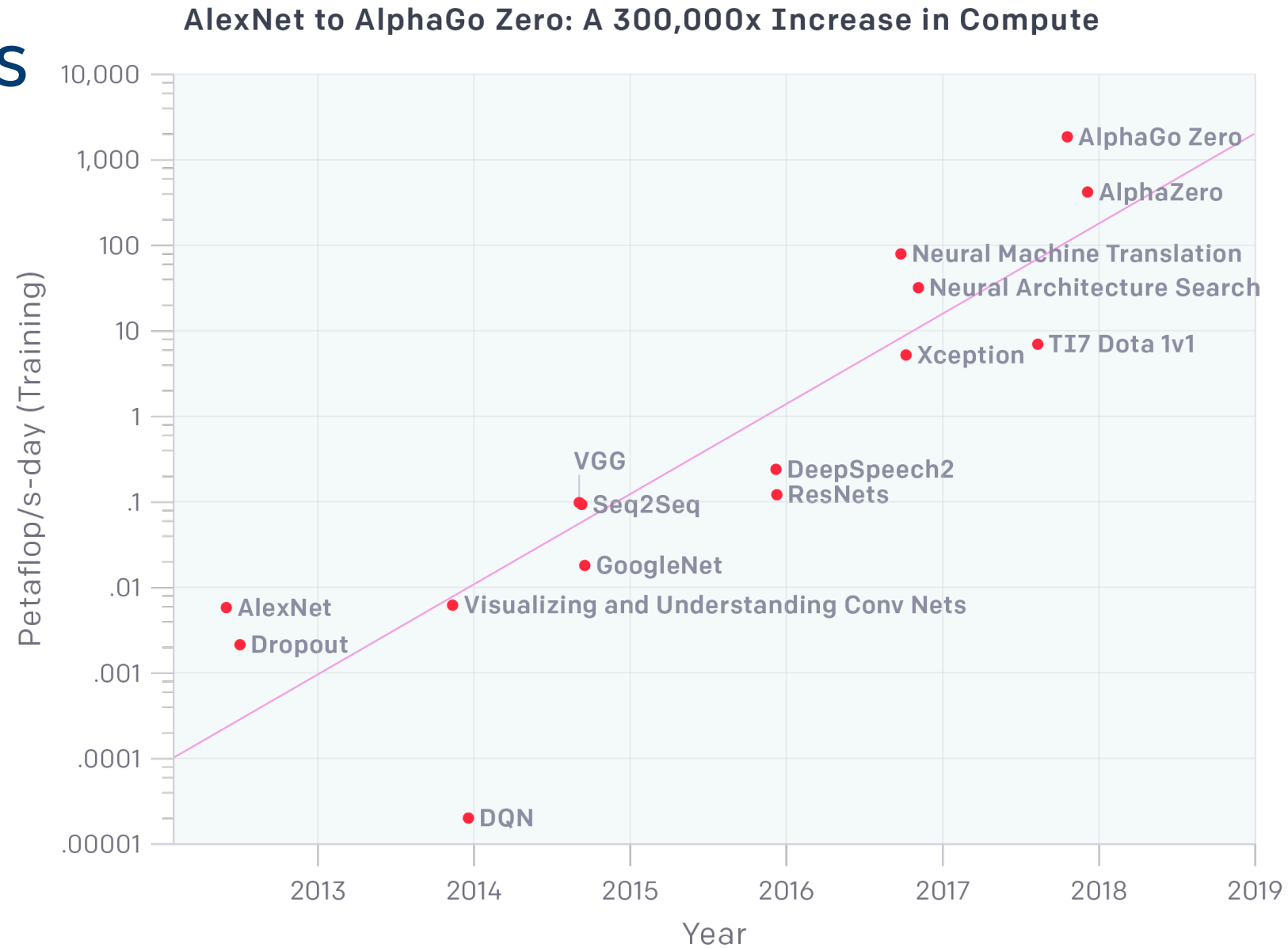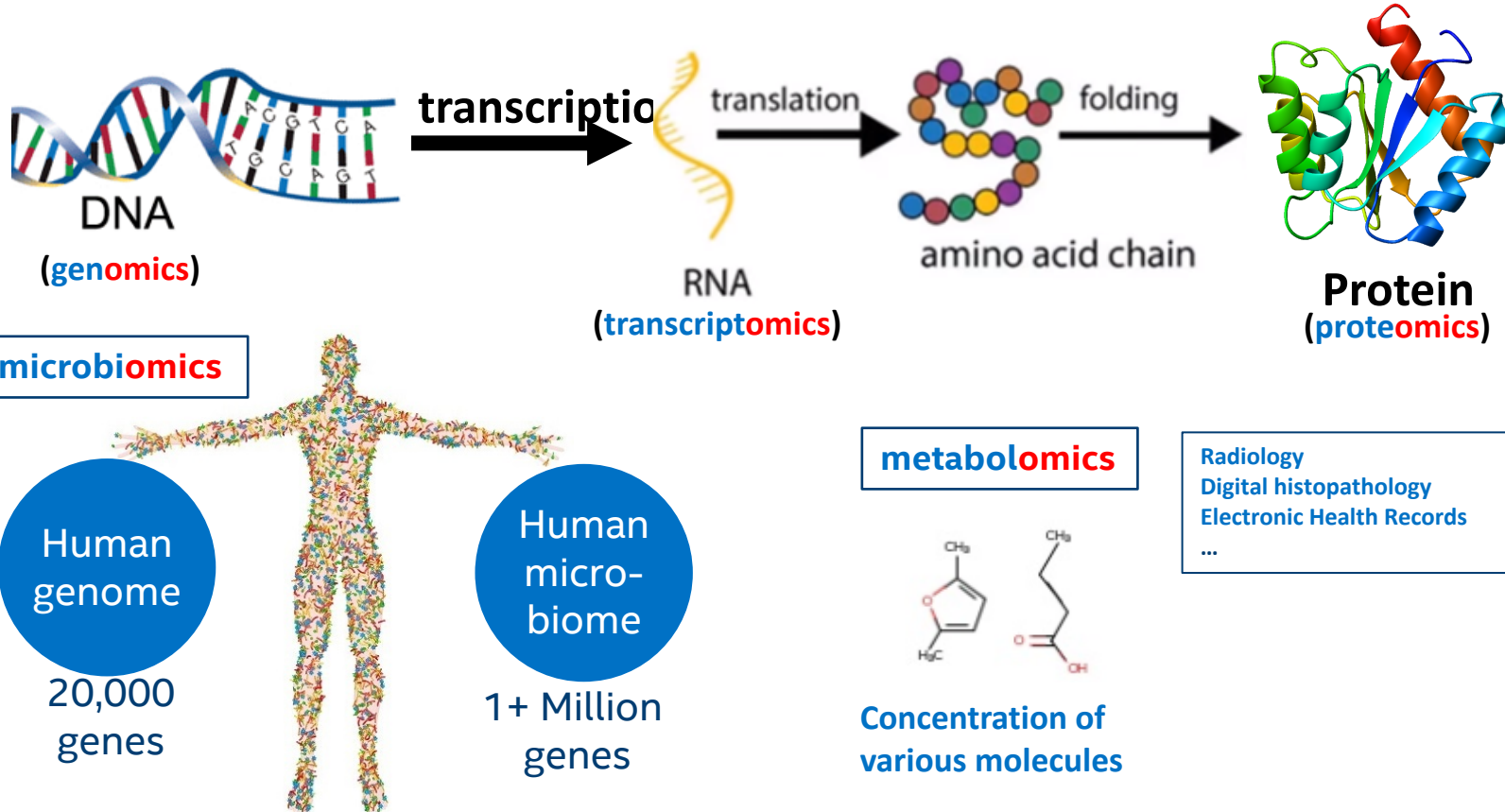
**microbiomics**

Human genome

20,000 genes

Human micro-biome

1+ Million genes

**metabolomics**

Radiology
Digital histopathology
Electronic Health Records
...

Concentration of various molecules

## Platform Characteristics

➢ **High throughput** omics pipeline over O(Billion) record data sets

➢ A multi-modal **AI pipeline** to correlate omics data across multiple datasets & generates actionable recommendations

➢ An open, modular **data-centric-computing platform** optimized for security, scalability and performance across many workloads

➢ 100s of exabytes of **storage capacity, processed in minutes** for results

➢ Best-in-class TCO with **easy to use/extend SW-Stack**

(intel)

# High-dimensional learning Infrastructure

# Graph Analytics
Improve social-network analysis, fraud-ring detection, anomaly detection



## CHALLENGES

- Sparse and irregular memory accesses
- Small data accesses with frequent synchronization
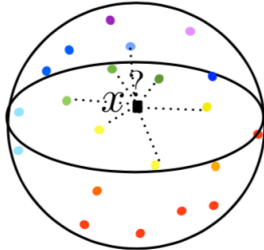- Scaling to very large datasets

**Real-Time Decision Making**
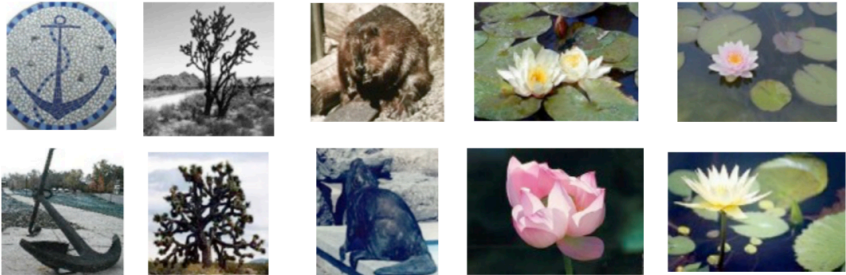
(intel)

# Simplified ML Theory World View

- If we only solve a **single** problem and design world-class hardware-software stack for it, what would that be:

  - *High Dimensional Learning*

  - *"Learn a super-compact,deep (hierarchical) approximation of dynamic graphs – computable in polynomial time, and evolving very slowly in time"*



**Curse of Dimensionality**

- $f(x)$ can be approximated from examples $\{x_i, f(x_i)\}_i$ by local interpolation if $f$ is regular and there are close examples:

- Need $\epsilon^{-d}$ points to cover $[0, 1]^d$ at a Euclidean distance $\epsilon$
  $\Rightarrow \|x - x_i\|$ is always large

Friday, October 3, 14

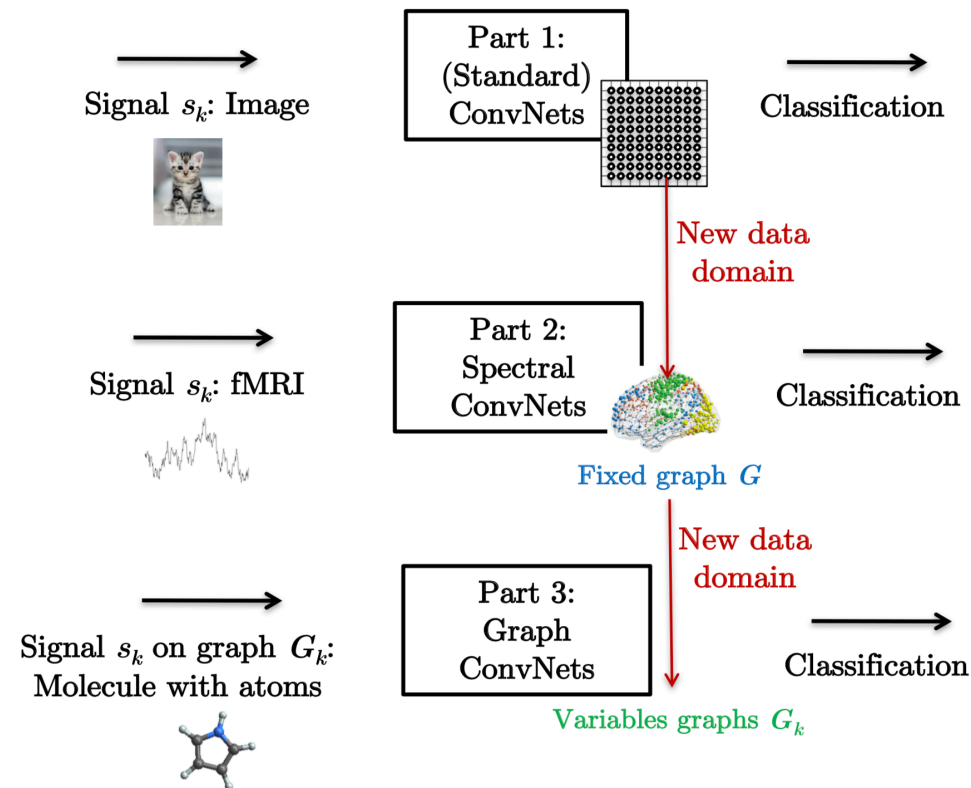**Credit for this slide above goes to Prof. Stephane Mallat ***

## Recent (2017) Algorithmic Breakthrough: Geometric Deep Learning **

# Will convolutions + downsampling prove effective in non-Euclidiean space as well?



Graph ConvNet architectures

# DARPA Graph Analytics Challenge



DARPA Taps Intel for Graph Analytics Chip Project

Michael Feldman | June 7, 2017 04:22 CEST

Intel shook hands with DARPA to craft HIVE Big Data platform a reality

**INTEL NAMED TO DARPA PROJECT FOCUSED ON MACHINE LEARNING AND ARTIFICIAL INTELLIGENCE**

Intel has been selected by DARPA, a U.S. Department of Defense agency, to collaborate on the development of a powerful new data-handling and computing platform that will leverage machine learning and other artificial intelligence (AI) techniques.

SR Siliconreview Team
2017-06-09

**PROGRAM TO DEVELOP NEW TECHNOLOGIES TO REALIZE 1,000X PERFORMANCE-PER-WATT GAINS IN THE ABILITY TO HANDLE GRAPH ANALYTICS**

intel

# Key Technologies and Scalability

## RE-IMAGINED ARCHITECTURE



**CPU Support for Small, Irregular Memory Accesses**

**Near-Memory Atomics**

## FULLY INTEGRATED



**Global Memory Model**

**Packaging for High I/O and Memory bandwidth**

## SEAMLESS SCALING



**Network as 1st –class Citizen**

**Flatten Latency Hierarchy**

**Point-to-point Messages**

(intel)

# The Future: Third Wave of AI



Deep Neural Networks Getting Augmented: NN + X + Memory
Such As: CNN + Bayes Net + Sparse Embeddings

# LOOKING AHEAD ...

## CONVENTIONAL
- Known procedures
- Generate answers

## DEEP LEARNING
- Known answers
- Generate procedures with training

## GRAPH ANALYTICS
- Graph edges and vertices represent relationships
- Big, sparse data structures

## NEUROMORPHIC
- Many procedures
- Adapt the answer with reinforcement

## QUANTUM
- Answers superimposed
- Select and measure the answer

# What makes delivering AI hard and fun!

## Better model building

LEARNING WITH LESS DATA AND SUPERVISION

DEEP NEURAL NETWORKS GETTING AUGMENTED: DATA-DRIVEN + ANALYTICAL + MEMORY

LEARNING MODELS THAT ARE EASIER TO REASON AND BETTER TO GENERALIZE

CONTINUOUS LEARNING FOR MISSION-CRITICAL AI

## More efficient and pervasive model deployment

THROUGHPUT, ACCURACY, AND MODEL SIZE TRADEOFFS: SPARSIFICATION AND PRUNING

SELF-LEARNING AND PERSONALIZATION AT THE EDGE

## Compute architecture needs of AI

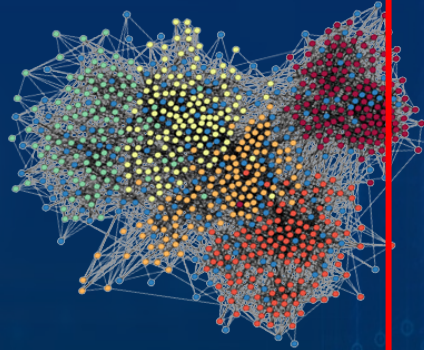REDUCING ARITHMETIC PRECISION WHILE PRESERVING ACCURACY

FEEDING THE COMPUTE ← MEMORY AND NETWORK; COMPUTE NEAR NETWORK

DOMAIN-SPECIFIC ARCHITECTURES → TRADITIONAL, SPATIAL, NEUROMORPHIC, QUANTUM

## Productivity and Scaling needs of AI

STRONG-SCALING AI TO HPC SCALE ON CLOUD INFRASTURCURE: HIGHER-ORDER METHODS

DELIVERING PERFORMANCE-PRODUCTIVITY: FaaS AND NEW LANGUAGES AND ABSTRACTIONS

We have a problem though …

# Where Are All the Ninjas?

2019-01-17 / DEVELOPMENT

## The Talent Shortage of Software Developers in 2019

1. According to Code.org, there are less than 50,000 Computer Science graduates in 2017. But, there are over 500,000 open computing positions in the United States. This could mean that in 2020, the available seats for this position will exceed qualified applicants by a million which could widen the gap even more.

**Only 10% of programmer pool have academic training to be *ninjas***
**New additions to programing pool are data scientists**
**Hardware is getting increasingly complex for extracting high performance**
**Hence, *ninja gap* is not likely to go down!**

Time to switch gear ….

From:
Designing hardware/software to meet AI needs

To:
AI helping us with our software-hardware needs

This implies a virtuous cycle for realized AI performance growth!

# Virtuous Supply Cycle of AI Compute

Design Software-Hardware for AI

AI for Software-Hardware Design

# MACHINE PROGRAMMING

## ANY TECHNIQUE THAT AUTOMATES SOME ASPECT OF SOFTWARE (AND HARDWARE) DEVELOPMENT.



## ENABLE THE WORLD TO CREATE SOFTWARE.

# The Machine Programming Inflection Point

- **ML/DL algorithms**
  - Transformer, neural recursion, DeepCoder, time series precision and recall, ...

- **Parallel hardware**
  - CPU, GPU, NNP, TPU, FPGAs, ...

- **Rich & dense programming data**
  - 33,000 github repos in 2008
  - 180,000,000 github repos in 2019

- *And programs are graphs after all*

**Inflection Point**

Business grows

Business declines

# Three Pillars of Machine Programming *



* Justin Gottschlich, Intel
Armando Solar-Lezama, MIT
Nesime Tatbul, Intel
Michael Carbin, MIT
Martin Rinard, MIT
Regina Barzilay, MIT
Saman Amarasinghe, MIT
Joshua B Tenenbaum, MIT
Tim Mattson, Intel

ACM SIGPLAN Workshop on Machine Learning and Programming Languages (MAPL), PLDI'18, arxiv.org/pdf/1803.07244.pdf

# Learning to Optimize Halide with Tree Search and Random Programs *

# Halide Learned System Design

# *Significantly Outperforms Human Ninjas*



A new automatic scheduling algorithm for Halide

Speed-up (higher is better)

Prior work (Mullapudi 2016) — Expert Humans — This paper

**Larger search space**
- includes more Halide scheduling features
- extensible

**Hybrid cost model**
- Mix of machine learning and hand-designed terms
- Can model complex architectures

# *Neo: A Learned end-to-end Query Optimizer ***

## Traditional Optimizers vs. Neural Optimizer

|  | Traditional Optimizers | Neural Optimizer (Neo) |
|---|---|---|
| **Creation** | Human developers | Learning from demonstration, Reinforcement learning |
| **Query Representation** | Operator tree | Tree with feature encodings |
| **Cost Model** | Hand-crafted model | Learned DNN model (value network) |
| **Plan Space Search** | Enumeration, Heuristics, Dynamic programming | DNN-guided search strategy (best-first) |
| **Cardinality Estimation** | Histograms, Hand-crafted models | Histograms, Learned embeddings (row vectors) |

Neo: first fully learned end-to-end solution to DB query optimization:
Matches/exceeds state of the art open-source & commercial optimizers within 24 hrs of training

# Neo's Learning Framework



Key ingredients:

- Learning from demonstration

- Value Model: DNN based on tree convolution

- Reinforcement learning

- Query and plan level feature encoding techniques (e.g., row-vectors that capture predicate semantics)

# Neo vs. Traditional Optimizers: Overall Performance



**Performance of each system's optimizer**

Neo's query plans are better

- Neo's query plans are:
  15-40% faster than Postgres
  20-25% faster than SQLite

- Neo also beats commercial
  optimizers except for TPC-H.

JOB — (purple) Join Order Benchmark over the IMDB database, complex predicates, lots of column correlations
TPC-H — (green) Standard OLAP benchmark from TPC at scale factor 10
Corporation — (blue) 2 TB dataset with 8000 unique queries from a large DB corporation

# Neo vs. Oracle Optimizer: Learning Curve

DSAIL Launched Jan 2019 @ MIT in collaboration with Google-Microsoft <link>

# Summary

We are at an unprecedented convergence of massive compute with massive data …
This confluence will have a lasting impact on both how we do computing and what computing can do for us!

# Questions?

# Notice and Disclaimers

Notice: This document contains information on products in the design phase of development. The information here is subject to change without notice. Do not finalize a design with this information. Contact your local Intel sales office or your distributor to obtain the latest specification before placing your product order.

INFORMATION IN THIS DOCUMENT IS PROVIDED IN CONNECTION WITH INTEL® PRODUCTS. EXCEPT AS PROVIDED IN INTEL'S TERMS AND CONDITIONS OF SALE FOR SUCH PRODUCTS, INTEL ASSUMES NO LIABILITY WHATSOEVER, AND INTEL DISCLAIMS ANY EXPRESS OR IMPLIED WARRANTY RELATING TO SALE AND/OR USE OF INTEL PRODUCTS, INCLUDING LIABILITY OR WARRANTIES RELATING TO FITNESS FOR A PARTICULAR PURPOSE, MERCHANTABILITY, OR INFRINGEMENT OF ANY PATENT, COPYRIGHT, OR OTHER INTELLECTUAL PROPERTY RIGHT. Intel products are not intended for use in medical, life saving, or life sustaining applications. Intel may make changes to specifications, product descriptions, and plans at any time, without notice.

All products, dates, and figures are preliminary for planning purposes and are subject to change without notice.

Designers must not rely on the absence or characteristics of any features or instructions marked "reserved" or "undefined." Intel reserves these for future definition and shall have no responsibility whatsoever for conflicts or incompatibilities arising from future changes to them.

Performance tests and ratings are measured using specific computer systems and/or components and reflect the approximate performance of Intel products as measured by those tests.  Any difference in system hardware or software design or configuration may affect actual performance.

The Intel products discussed herein may contain design defects or errors known as errata which may cause the product to deviate from published specifications. Current characterized errata are available on request.

Copies of documents which have an order number and are referenced in this document, or other Intel literature, may be obtained by calling 1-800-548-4725, or by visiting Intel's website at http://www.intel.com.

# Notice and Disclaimers Continued ...

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products.  For more information go to http://www.intel.com/performance

Intel's compilers may or may not optimize to the same degree for non-Intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include SSE2, SSE3, and SSE3 instruction sets and other optimizations. Intel does not guarantee the availability, functionality, or effectiveness of any optimization on microprocessors not manufactured by Intel. Microprocessor-dependent optimizations in this product are intended for use with Intel microprocessors. Certain optimizations not specific to Intel microarchitecture are reserved for Intel microprocessors. Please refer to the applicable product User and Reference Guides for more information regarding the specific instruction sets covered by this notice.  Notice revision #20110804