



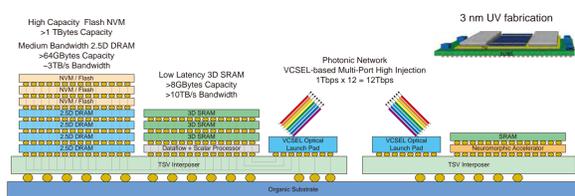
Infrastructure for Big Data and AI System Software and Application Research

Exploring Next-Gen Memory Hierarchy

To achieve higher performance and larger capacity on recent and future architectures, we need to explore next-gen memory hierarchy, including heterogeneous devices. Also placement of many-cores and memory devices can be reconsidered; 3D stacking of cores and memory chips may go mainstream in HPC/Big-data area.

Under such assumptions, understanding memory performance will be more challenging for data and address traffic routed among many-cores, which may often conform 2D mesh.

To analyze its effects, we have conducted preliminary measurements of memory latency from every core of Xeon Phi. We see 9% difference in the current architecture; the effect will be expanded in future architecture where chip cores are also stacked.



An Example of Future Architecture with 3D Stacking Technology and Heterogeneous Memory

33	32	25	24			5	4	1	0
170.9	171.0	168.1	168.1			161.2	161.2	163.1	163.2
35	34	27	26	19	18	13	12	7	6
171.3	171.4	165.2	165.3	165.5	165.5	162.9	162.9	161.5	161.5
		29	28	21	20	15		9	8
		168.9	168.9	166.8	166.8	163.6		162.0	162.0
		31	30	23	22	17	16	11	10
		169.0	169.0	167.8	167.8	164.7	164.7	162.9	162.9
63	62	59	58	53	52	47	46	41	40
174.1	174.1	171.8	171.8	168.8	168.8	165.8	165.8	165.1	165.1
65	64	61	60	55	54	49	48	43	42
175.1	175.1	172.0	172.0	170.5	170.5	167.3	167.3	166.2	166.2
67	66			57	56	51	50	45	44
176.7	176.8			171.6	171.6	168.4	168.4	167.5	167.5

max: 176.7 (ns)
min: 161.2 (ns)
delta: 15.5 (ns)

Memory access latency from every core of Xeon Phi 7285. As the target of access, one of MCDRAM modules, placed close to core 4, 5 is used.

Acknowledgments. This research is supported by the New Energy and Industrial Technology Development Organization (NEDO), Japan.

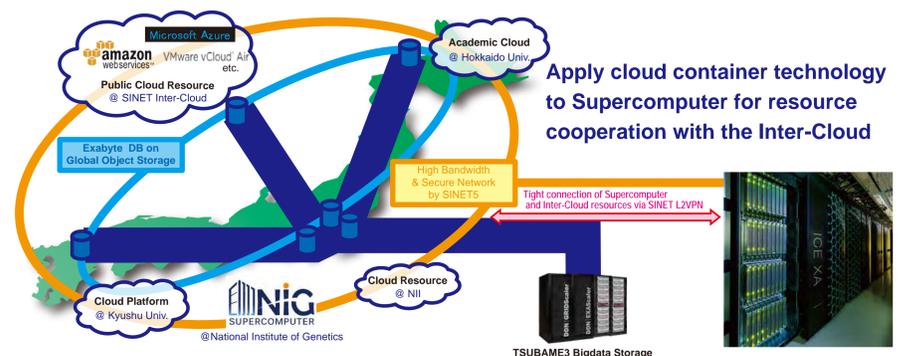
Inter-Cloud Infrastructure for Big Data Analysis

Building a Testbed Infrastructure on Overlay Cloud

- Using SINET5 network infrastructure (100Gbps Network)
- Cooperation with Cloud and Supercomputer
- Providing the Science Data Repository
 - Testbed: Petabytes class object storage
 - Real System: Exabytes class object storage

Cloud and Supercomputer Federation

- Development of collaboration technology between cloud and supercomputer using TSUBAME 3.0 as a test case

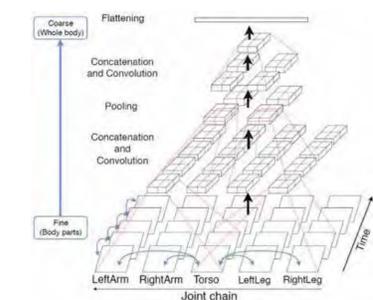


Acknowledgments. These researches are supported by CREST, JST CREST Grant Numbers JPMJCR1303 and JPMJCR1501. (Research Area: Advanced Core Technologies for Big Data Integration).

Efficient Deep Learning for Video Processing

Fine-to-Coarse Convolutional Neural Network for 3D Human Action Recognition

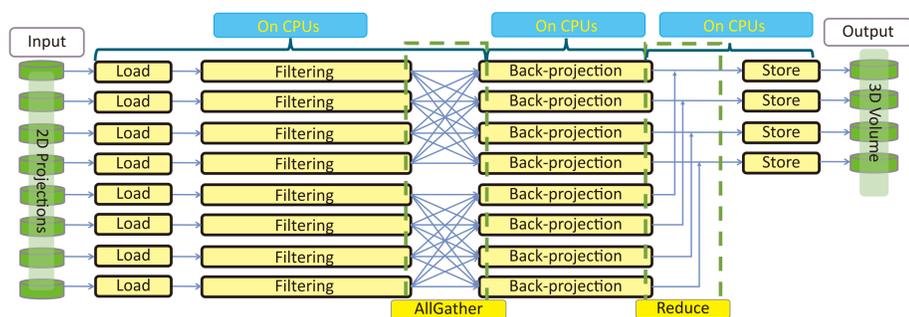
Current convolutional action recognition models struggle to fully capture the temporal dimension due to the required network depth. To avoid this issue, we investigate a coarse-to-fine network architecture to learn



long-distance temporal dependencies. The information of each body part is aggregated while respecting the structure of the body. We achieve state-of-the-art performance and improve it for two person interactions.

High-resolution Image Reconstruction on Supercomputers

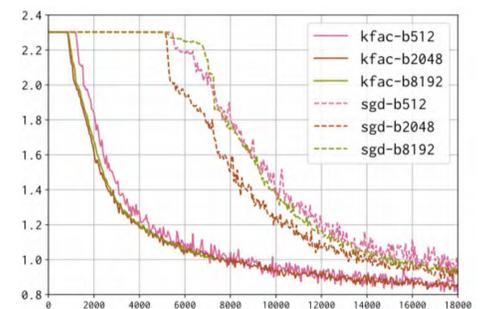
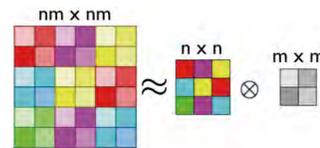
CT is a widely used technology that requires compute-intensive algorithms for image reconstruction. We exploit the heterogeneity of GPU-accelerated systems by overlapping the filtering and back-projection stages on CPUs and GPUs, respectively. We also propose a distributed framework for high-resolution image reconstruction on state-of-the-art GPU-accelerated supercomputers. We demonstrate the scalability and instantaneous CT capability of the distributed framework by using up to 2,048 V100 GPUs to solve 4K and 8K problems within 30 seconds and 2 minutes, respectively.



Presentation in SC19 Technical Program
iFDK: A Scalable Framework for Instant High-Resolution Image Reconstruction
Thursday 3:30pm-4pm @ Room 301-302-303

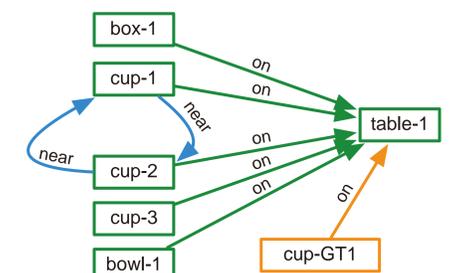
Second Order Optimization for Deep

Distributed deep learning on thousands of GPUs forces the batch stochastic descent methods to operate in a regime where the increasing batch size starts to have a detrimental effect on the convergence and generalization. We investigate the possibility of using second-order optimization methods with proper regularization as an alternative to conventional stochastic gradient descent methods.



Scene Graph Prediction

Understanding higher-level concepts such as recognizing a scene in a video stream requires algorithms that can leverage information beyond the pixel level. Graphs have been used to represent relationships between objects in a given image. A scene graph is a compact representation of an image that can represent both objects present and relationships between them. We investigate the application of graph neural networks based methods for designing efficient scene graph prediction models.



Acknowledgments. This research is supported by CREST, JST CREST Grant Number JPMJCR19F5.