Memory Allocation on
[Base / RDMA / GR + NIC / GR + RNIC]

# INDEX

# BASE vs RDMA

# BASE vs RDMA



CPU 효율 차이

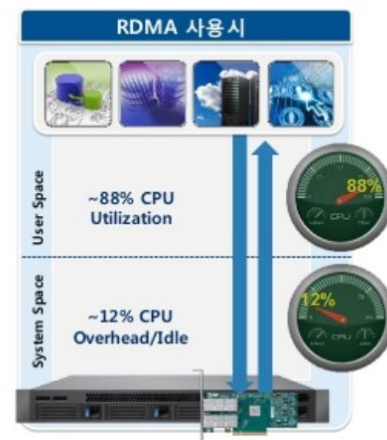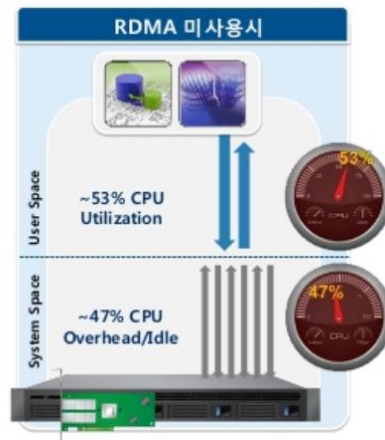# GR + NIC

RDMA requires RDMA-enable NIC
     RNICs : iWARP, ROCE NICs
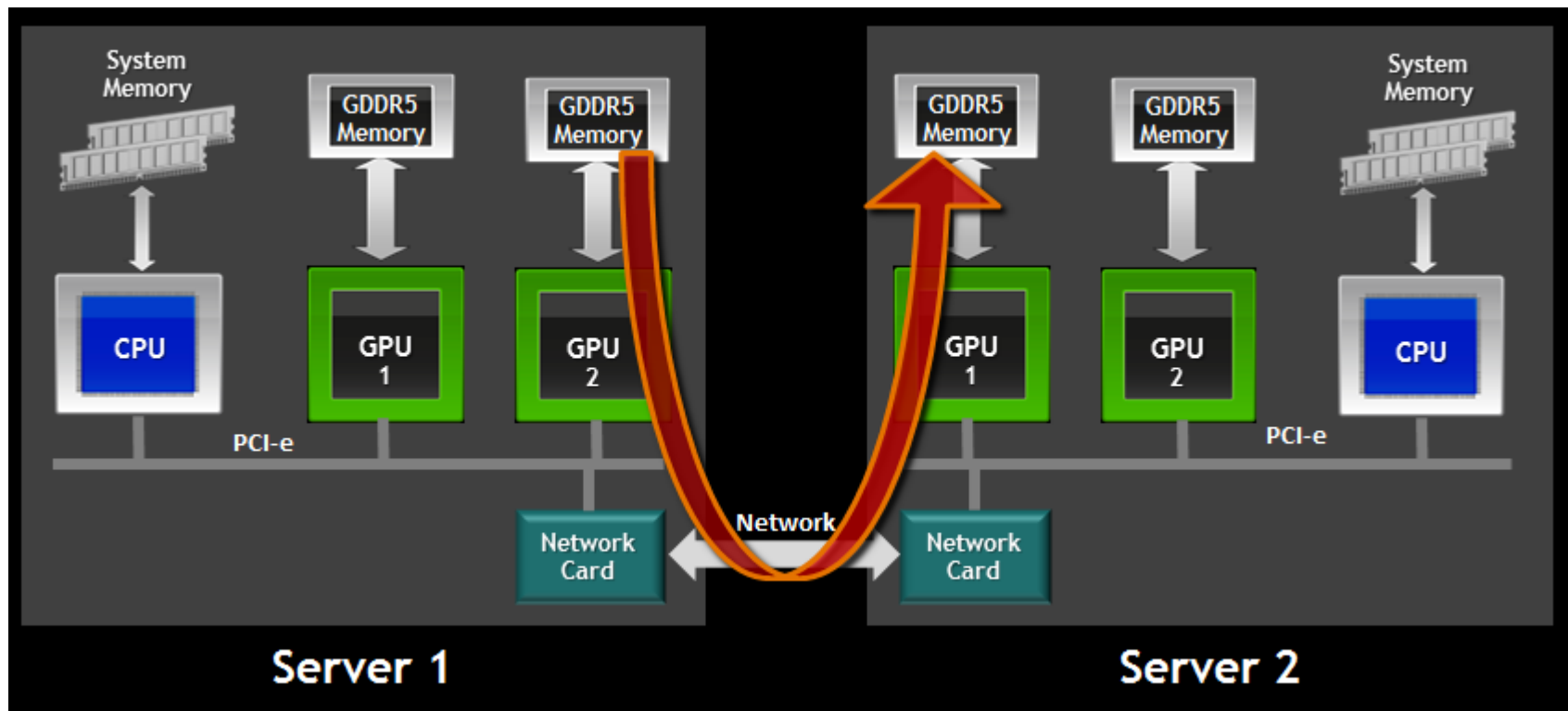     Infiniband : already infused in the IB networks

∴ 일반 NIC을 사용한 GPU RDMA 관련된 연구 자료를 찾기 어렵다
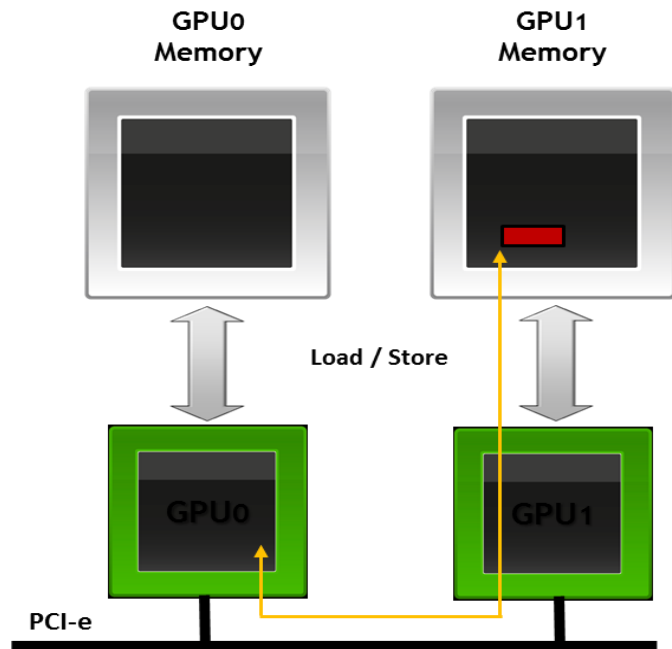
# Mellanox Requirements

**Table 2 - GPUDirect RDMA System Requirements**

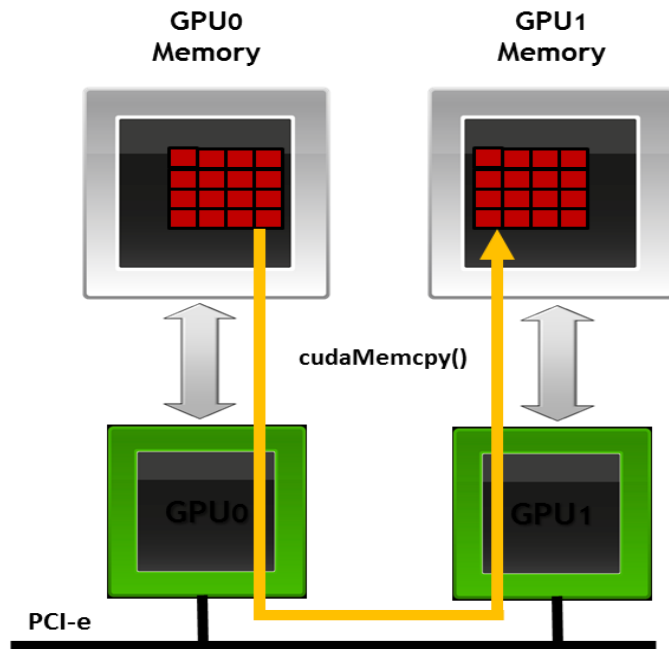| Platform | Type and Version |
|---|---|
| HCAs | • ConnectX®-3 (VPI/EN)<br>• ConnectX®-3 Pro<br>• Connect-IB®<br>• ConnectX®-4 (VPI/EN)<br>• ConnectX®-4 Lx<br>• ConnectX®-5 (VPI/EN)<br>• ConnectX®-6 (VPI/EN)<br>• NVIDIA® Tesla™ / Quadro K-Series or Tesla™ / Quadro™ P-Series GPU |
| Software/Plugins | • MLNX_OFED v2.1-x.x.x or later<br>www.mellanox.com -> Products -> Software - > InfiniBand/VPI Drivers -> Linux SW/ Drivers<br>• Plugin module to enable GPUDirect RDMA<br>www.mellanox.com -> Products -> Software - > InfiniBand/VPI Drivers -> GPUDirect RDMA (on the left navigation pane)<br>• NVIDIA Driver<br>http://www.nvidia.com/Download/index.aspx?lang=en-us<br>• NVIDIA CUDA Runtime and Toolkit<br>https://developer.nvidia.com/cuda-downloads<br>NVIDIA Documentation<br>http://docs.nvidia.com/cuda/index.html#getting-started-guides |

# GR + RNIC
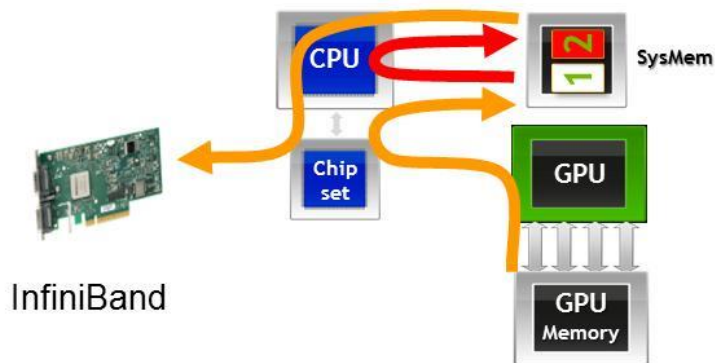
# GR + RNIC



**P2P Direct Access**

**P2P Direct Transfers**

# GR + RNIC



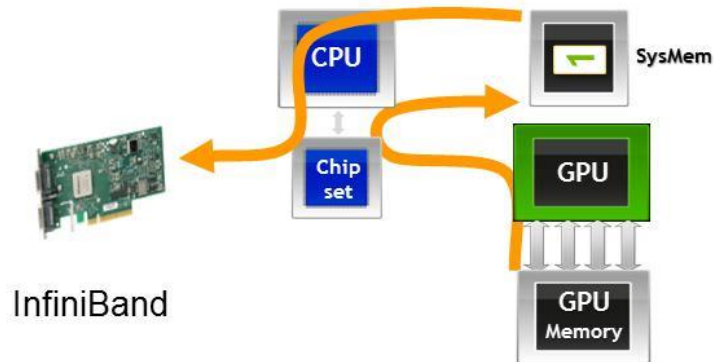## Without GPUDirect

Same data copied three times:
1. GPU writes to pinned sysmem1
2. CPU copies from sysmem1 to sysmem2
3. InfiniBand driver copies from sysmem2

## With GPUDirect

Data only copied twice
Sharing pinned system memory makes sysmem-to-sysmem copy unnecessary

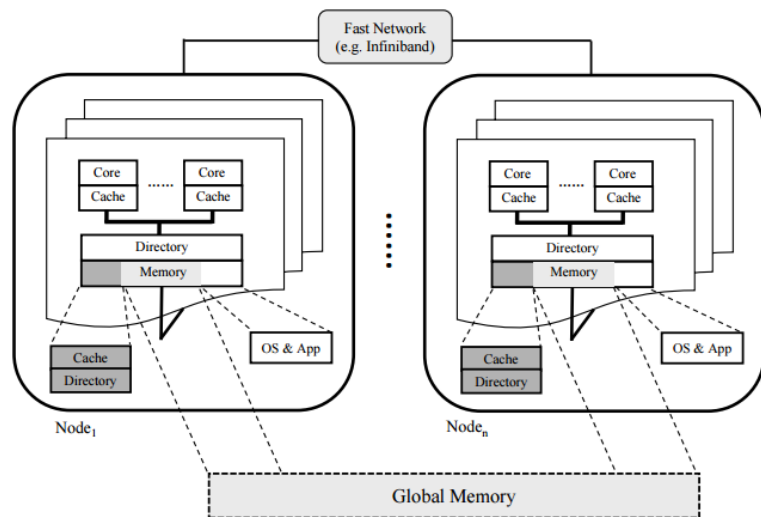# GAM : Efficient Distributed Memory Management with RDMA and Caching



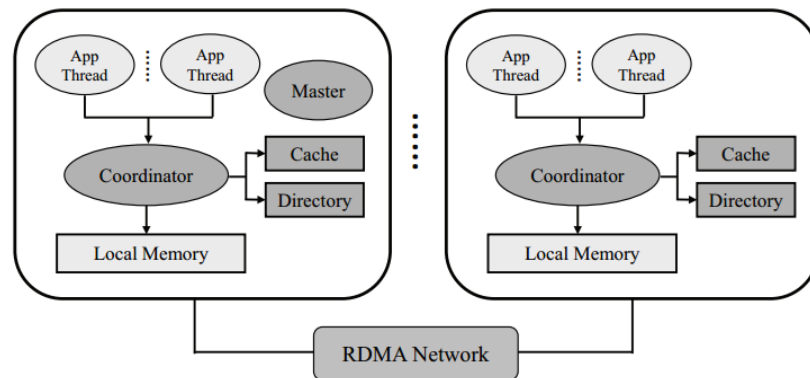**Figure 1:** Overview of GAM

**Figure 5:** Architecture of GAM

GAM manages the free memory distributed among multiple nodes to provide a unified memory model.

# Accelerating TensorFlow with RDMA for high-performance deep learning

## OSU AR-gRPC Architecture

- Adaptive RDMA gRPC
- Features
  - Hybrid Communication engine
    - Adaptive protocol selection between eager and rendezvous
  - Message pipelining and coalescing
    - Adaptive chunking and accumulation
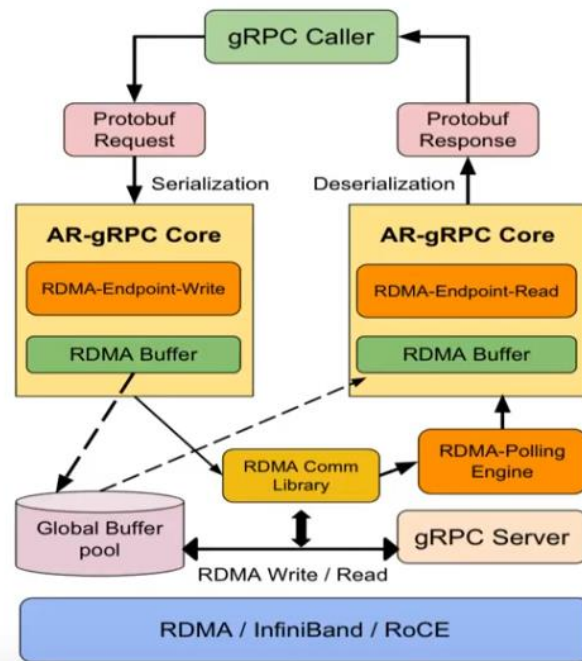    - Intelligent threshold detection
  - Zero copy transmission
    - Zero copy send/recv

# 그림 출처

Mellanox 블로그랑 white paper

논문들