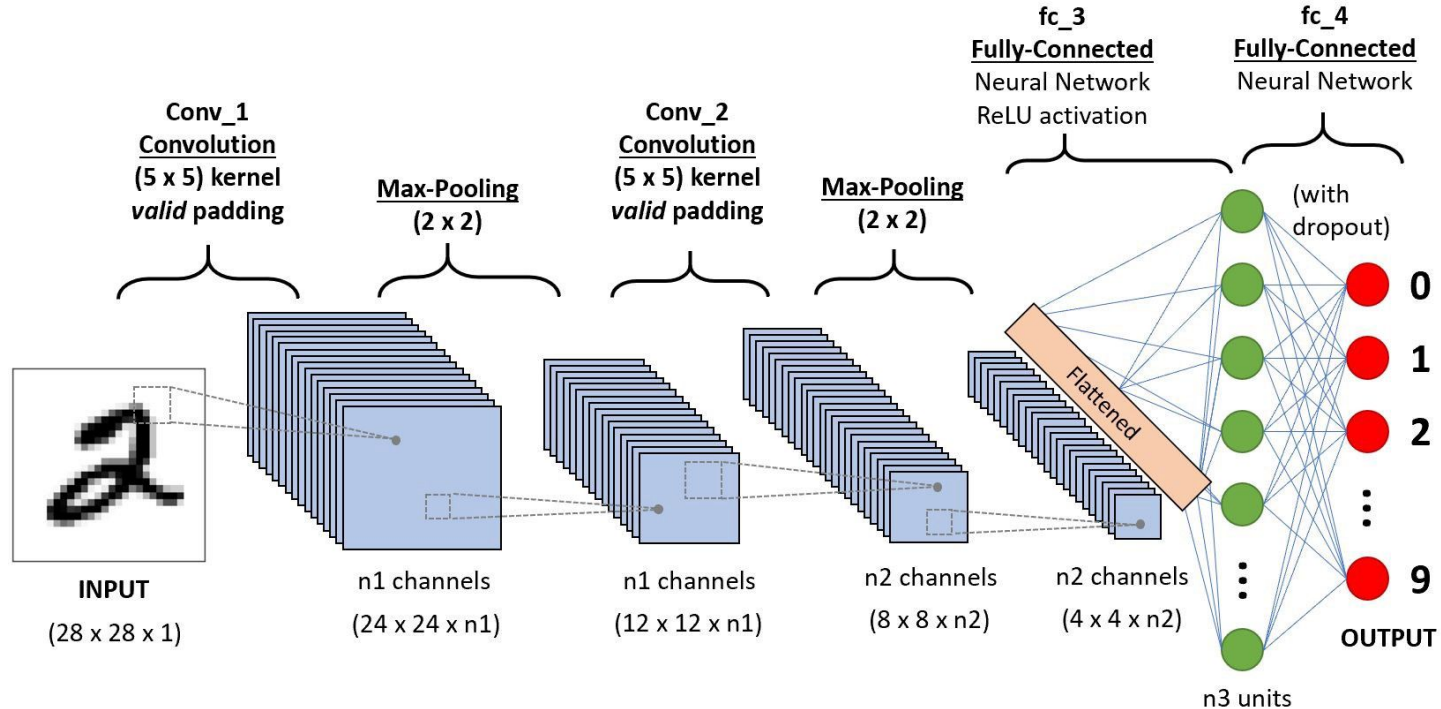# Restructuring Batch Normalization to Accelerate CNN training
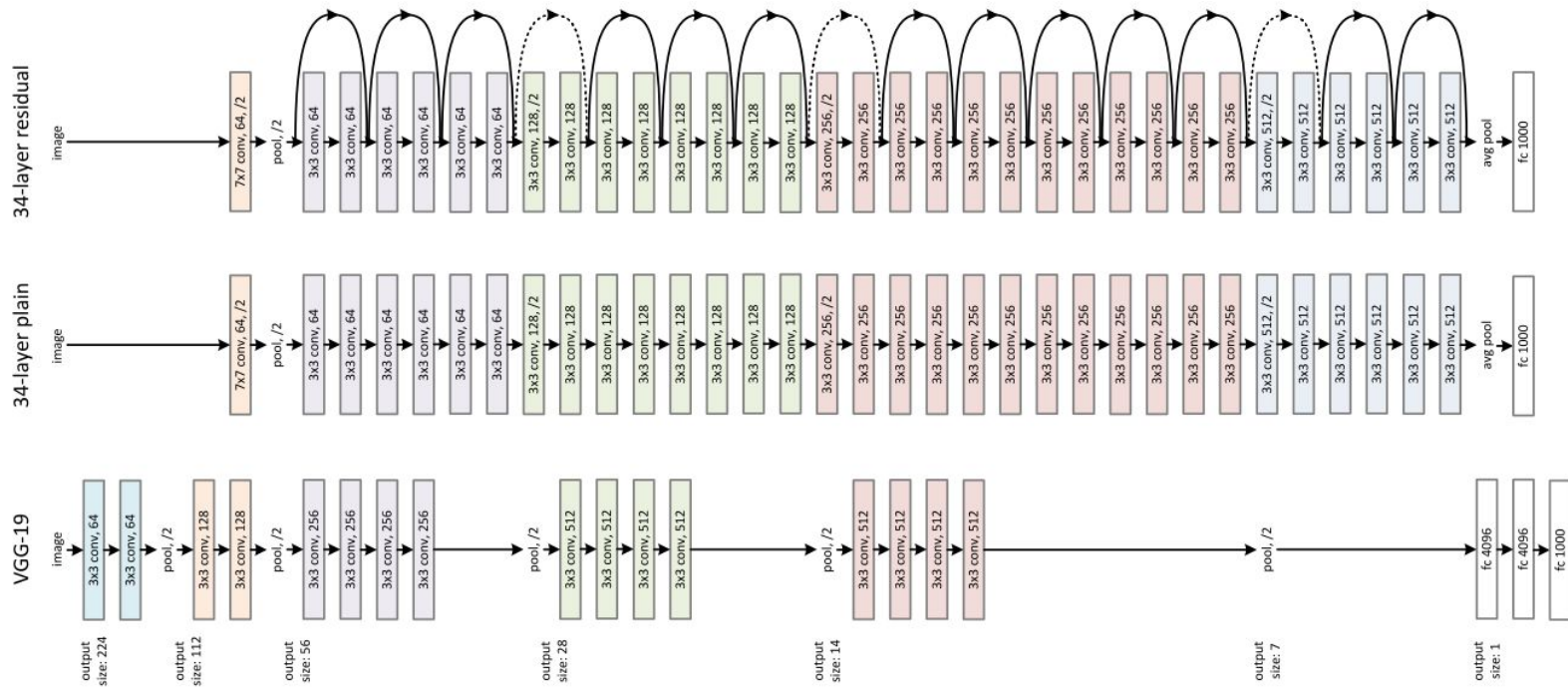
By Jung, Wonkyung, Daejin Jung, Byeongho Kim, Sunjung Lee, Wonjong Rhee, and Jung Ho Ahn in *Proceedings of Machine Learning and Systems* 1 (2019): 14-26.
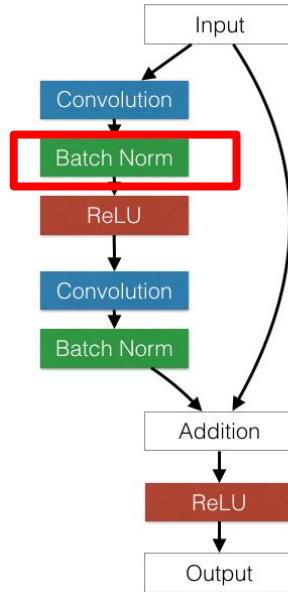
Presentation by: Saurav Sengupta
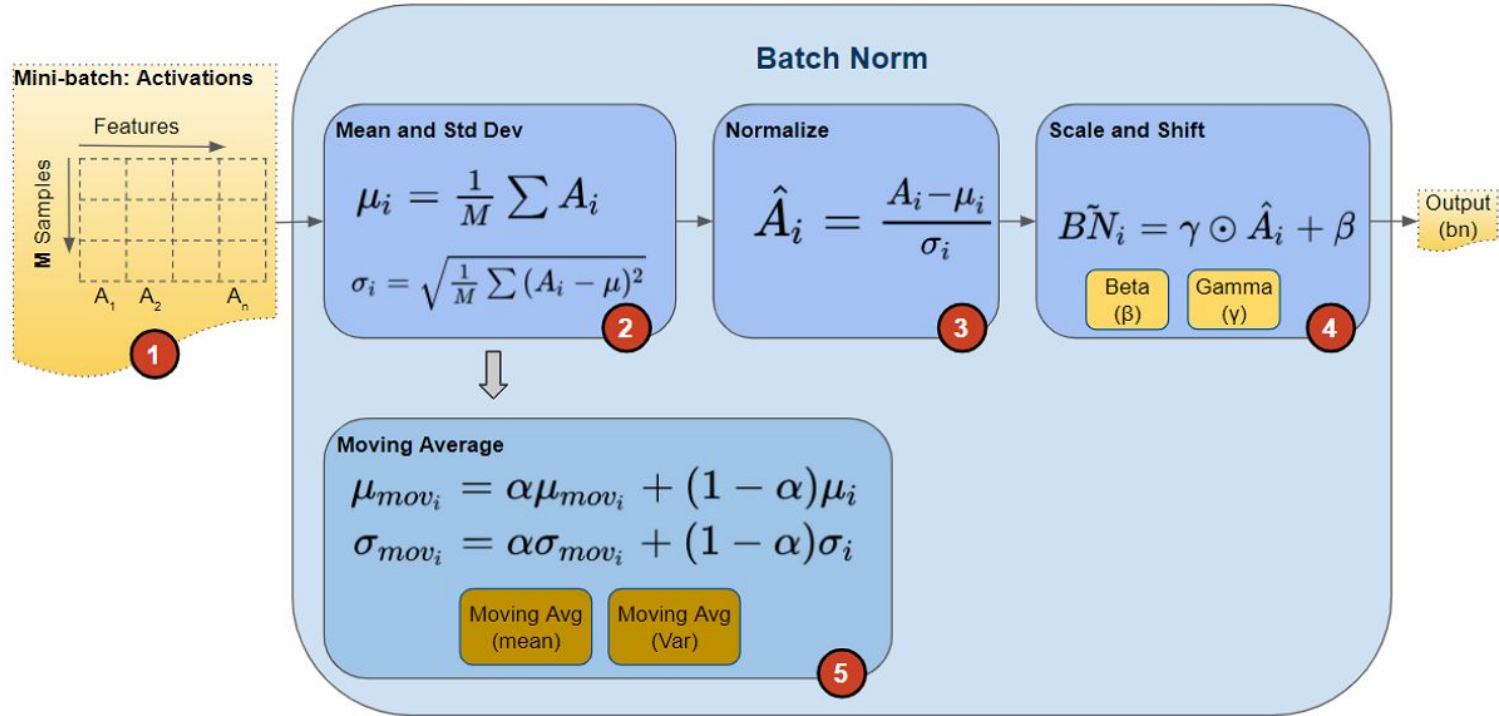
# Core Principles: CNN

# Residual Neural Networks (ResNets)

# ResNets with Batch Norm

# Batch Normalization

# Why does batch norm work?

Nobody knows for sure.

But there are conjectures as to why it might work:

1. BN makes the optimization landscape significantly smoother
2. Batch normalization also helps avoid large gradient updates that can result in diverging loss and activations growing uncontrollably.
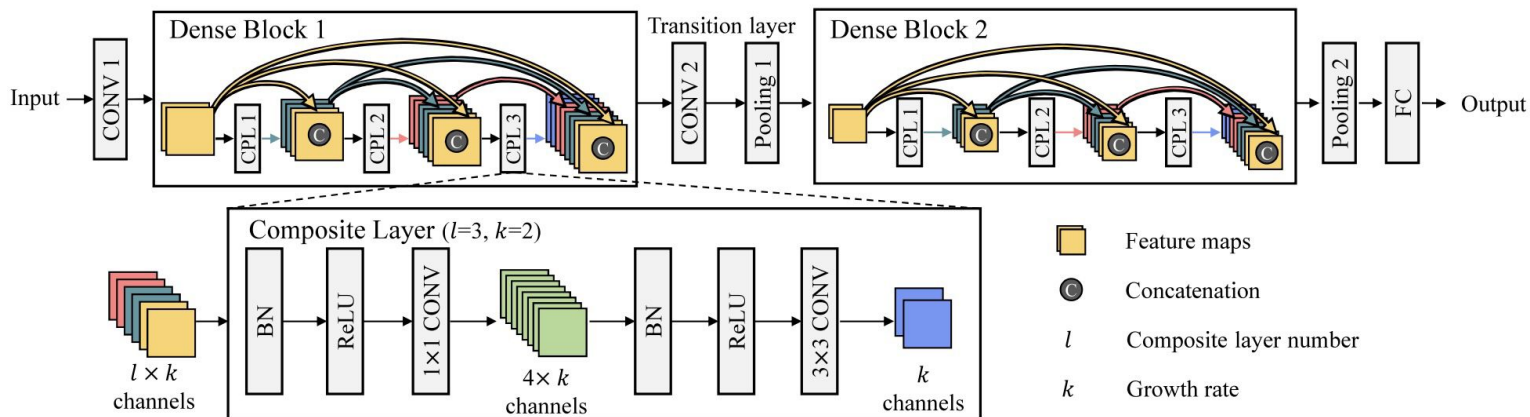
# DenseNet



*Figure 2.* An exemplar DensetNet structure with two Dense Blocks, connected through a transition layer which changes the number and size of the input feature maps to a Dense Block. Each Dense Block has multiple Composite Layers (CPLs), each of which is connected to every other CPL within a Dense Block in a feed-forward fashion. A CPL consists of six layers (BN, ReLU, $1\times1$ CONV, BN, ReLU, and $3\times3$ CONV). The $1\times1$ CONV layer in a CPL, called bottleneck layer, limits the number of input feature maps to $4\times k$ while the second CONV outputs $k$ channels that are concatenated to the input feature maps. Growth rate ($k$) is the variable for how many feature maps are concatenated per CPL; feature maps stack up as they go through CPLs.

# How is it different from ResNet?

- ResNet uses Element Wise Sum (EWS).

- DenseNets use concatenated feature maps.

- Less number of parameters to be optimized.

- Much more dense connectivity.

# Different Optimization Techniques for accelerators

Two choices:

1. Optimize Convolution Layers
2. Optimize the non-Convolution layers

Optimize Convolution layers by data reuse techniques for weights and input and output feature maps (ifmap and ofmaps), or pruning weights or adapting new memory technologies.
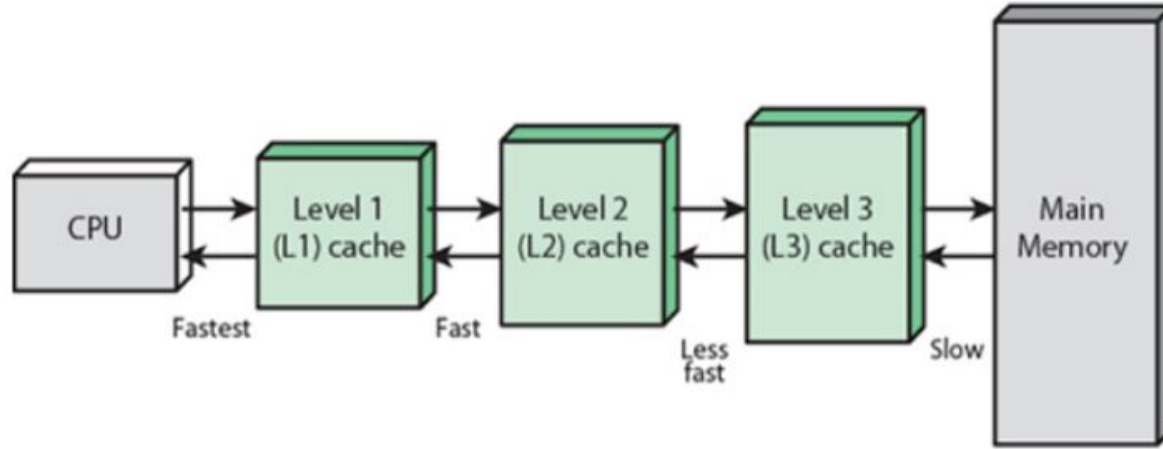
Most research has not focused on the non-CONV layers that most of the time outnumber the CONV layers.

DenseNet-121 (DenseNet with 120 CONV layers plus one FC layer) spends more than half of the execution time on the non-CONV layers.

# Key Problems

- CONV layers with smaller filter sizes have relatively high memory access rates compared to the total amount of computation. DenseNet has high number of these small filter sized CONV layers.
- non-CONV layers of DenseNet-121 are mostly bottlenecked by the peak main memory bandwidth of the system we use.
- CONV layers underutilize the available bandwidth.
- Memory bandwidths are not high enough to completely utilize all the computational resources.
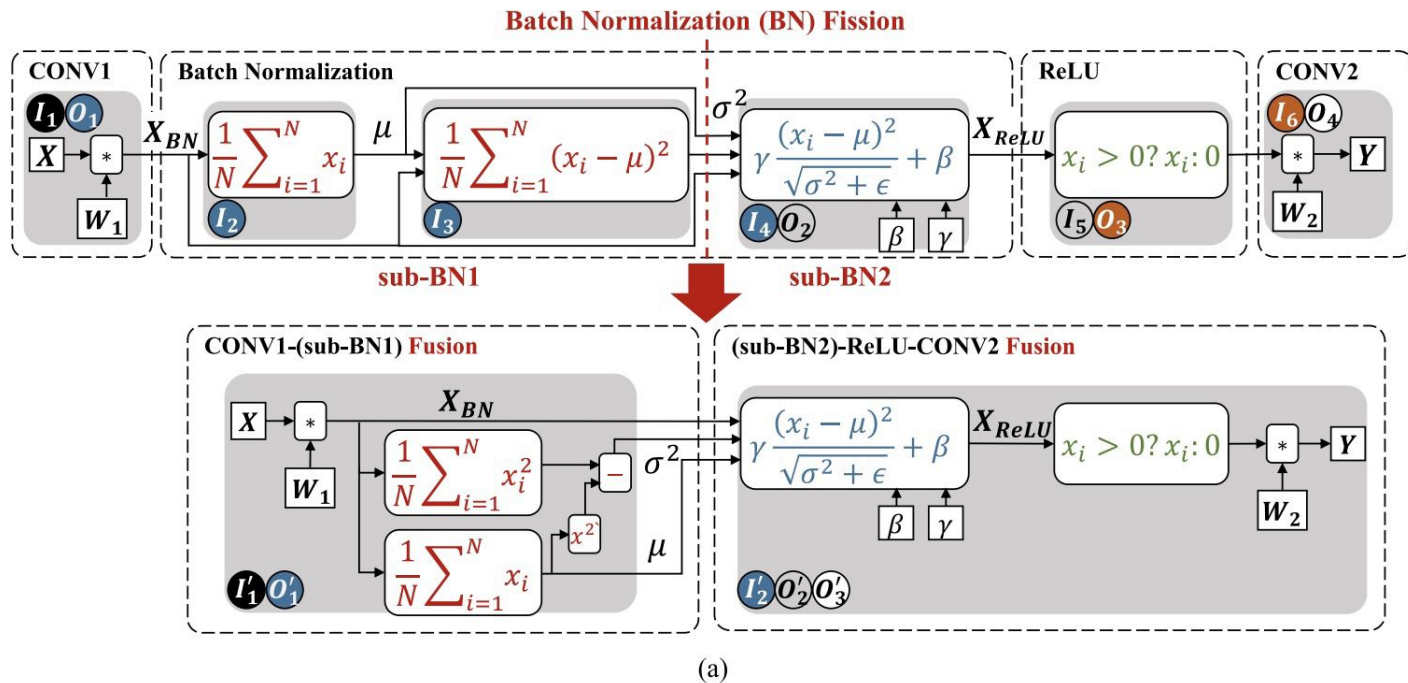
# Interlude: How memory works



Memory bandwidth: Rate at which data can be read from memory by the processor.
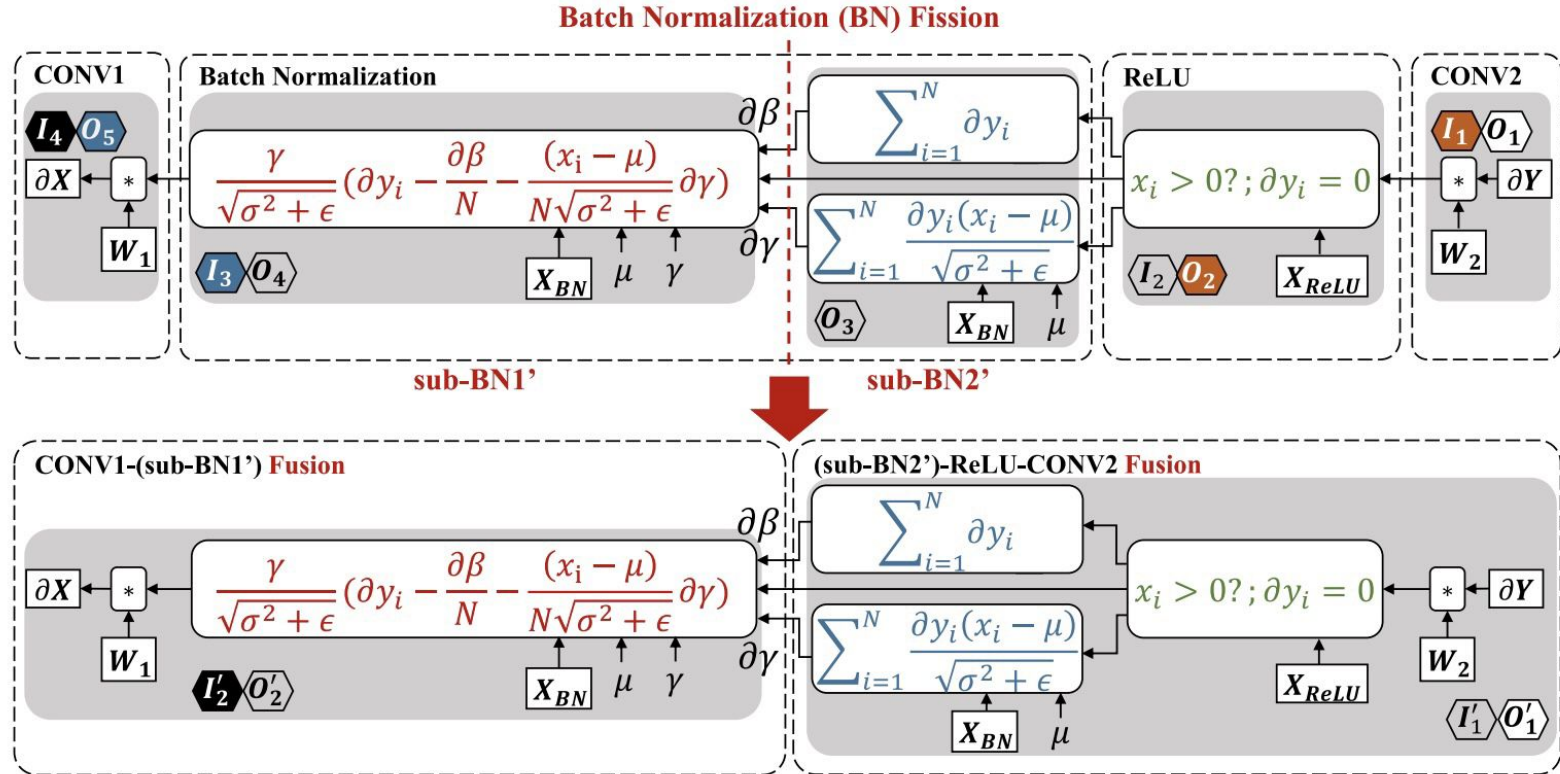
# Key Problems: Contd.

- Data reuse techniques do not work with Batch Norm layers.

- BN layers have strict data dependency, cross-layer data reuse is prohibited.

# Solution: Fission and Fusion (Forward)



(a)

$$V(X) = E((X - E(X))^2) = E(X^2) - E(X)^2$$

# Solution: Fission and Fusion (Backward)

# Results

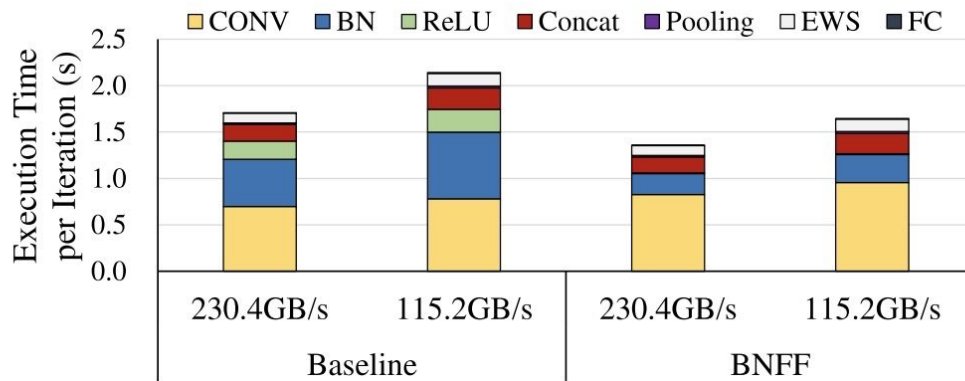- BNFF removes five memory sweeps per BN layer.



*Figure 8.* Execution time comparison of baseline and BNFF for 230.4GB/s and 115.2GB/s memory bandwidth using the Xeon processor with DenseNet-121.

# Results (Contd.)

"Experiments on a latest chip multiprocessor showed that the proposed BN restructuring can improve the training performance of DenseNet-121 by 47.9% for forward pass and by 15.4% for backward pass, leading to overall **25.7%** improvement. Applying the BN restructuring to GPU with an open-source linear algebra library also showed **17.4%** of performance improvement. The large improvement suggests that non-CONV layers are important candidates for acceleration and that future research should pay keen attention to how CNN models evolve"

# Limitations

- Only works on the various combinations of Conv+BatchNorm and BatchNorm+ReLU+Conv so works on DenseNets or ResNet but might not work on other implementations that do not have those sequences.
- Does not reduce the timing of backpropagation significantly because the number of reads and writes remain mostly the same.