# FreeLunch: Compression-based GPU Memory Management for Convolutional Neural Networks

**Shaurya Patel**
**University of British Columbia**
**University of Massachusetts, Amherst**

**Tongping Liu**
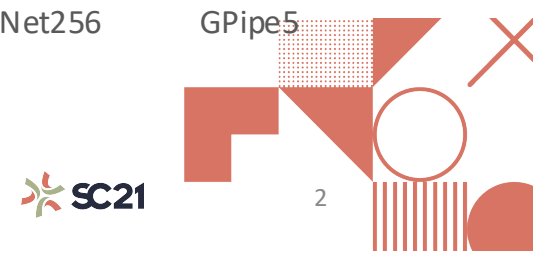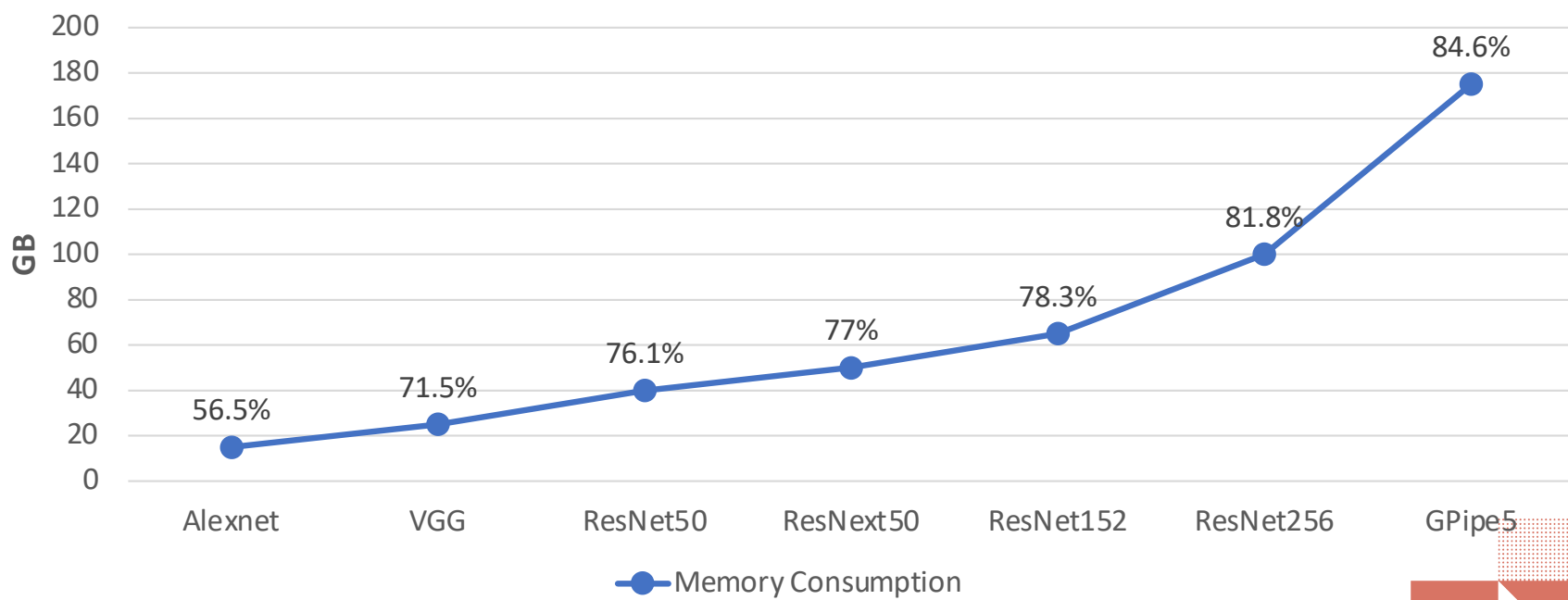**University of Massachusetts, Amherst**
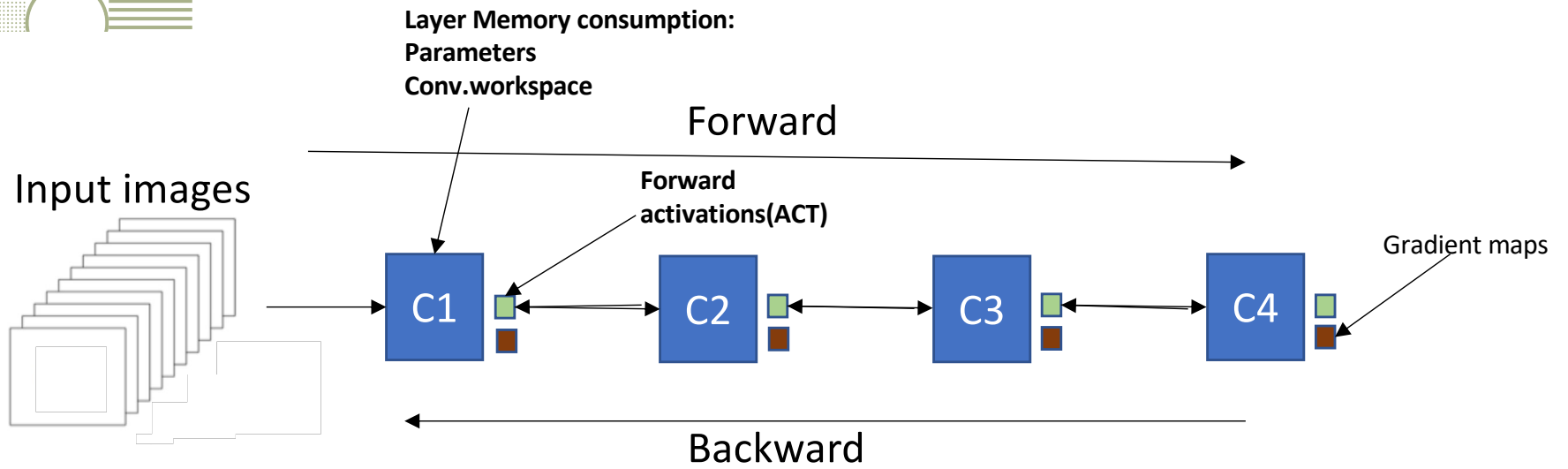
**Hui Guan**
**University of Massachusetts, Amherst**

# CNN memory consumption trend



Memory Consumption

CNN training

Layer Memory consumption:
Parameters
Conv.workspace

Forward

Input images

Forward activations(ACT)

Gradient maps

C1    C2    C3    C4

Backward
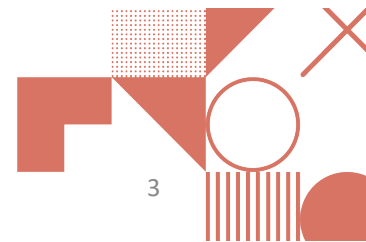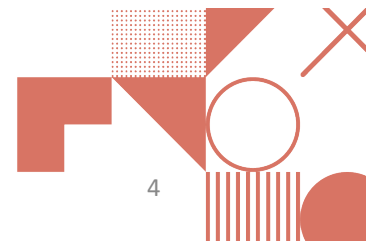
Forward Activations(much larger in size than model params) need to persist in memory until the gradient updates in backward phase!
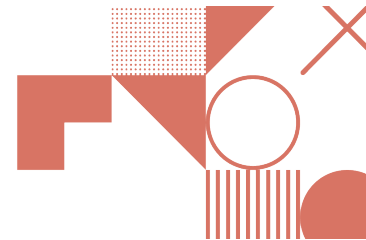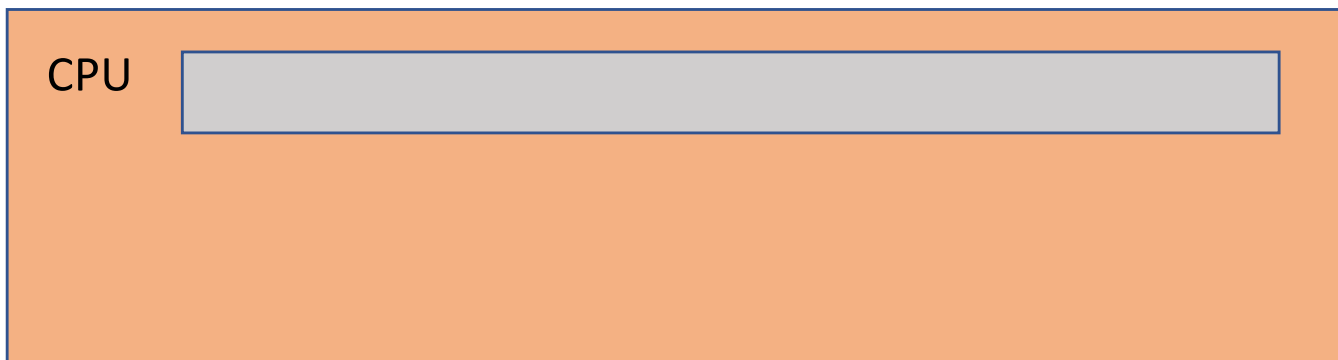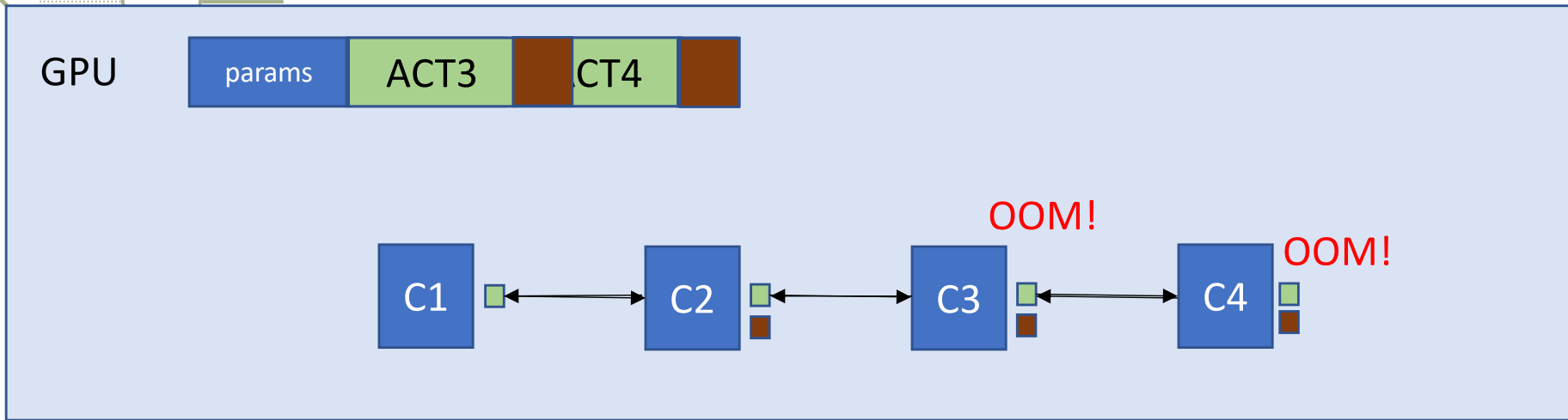
Policies for memory management

- Swapping
  - Capuchin [X. Peng et al., 2020]
  - SwapAdvisor [C-C. Huang et al., 2020]
  - Superneurons [L. Wang et al., 2018]
  - …

- Recomputation
  - Capuchin [X. Peng et al., 2020]
  - Superneurons [L. Wang et al., 2018]
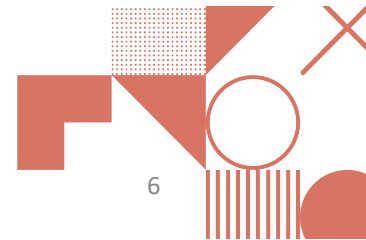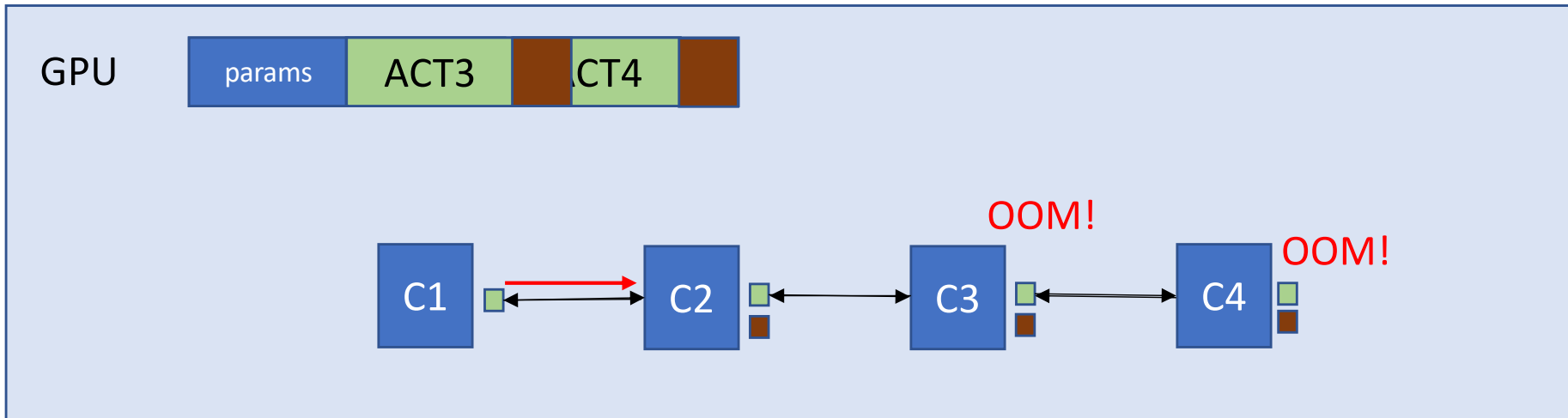  - …

Swapping

CPU-GPU bandwidth is a bottleneck!
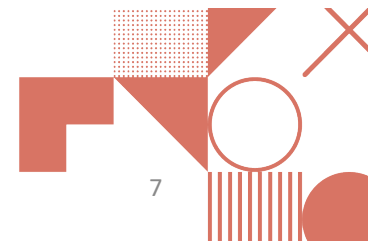
Recomputation is complex and has lineage dependencies!

- A compression based policy for CNN training.
  - basic Idea: **compress and keep the tensors on GPU memory**.
  - avoids the bandwidth issue introduced by swapping.
  - avoids the computation complexity of recomputation.

- Challenges:
  - How to reduce the compression overhead?
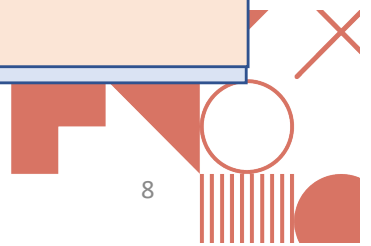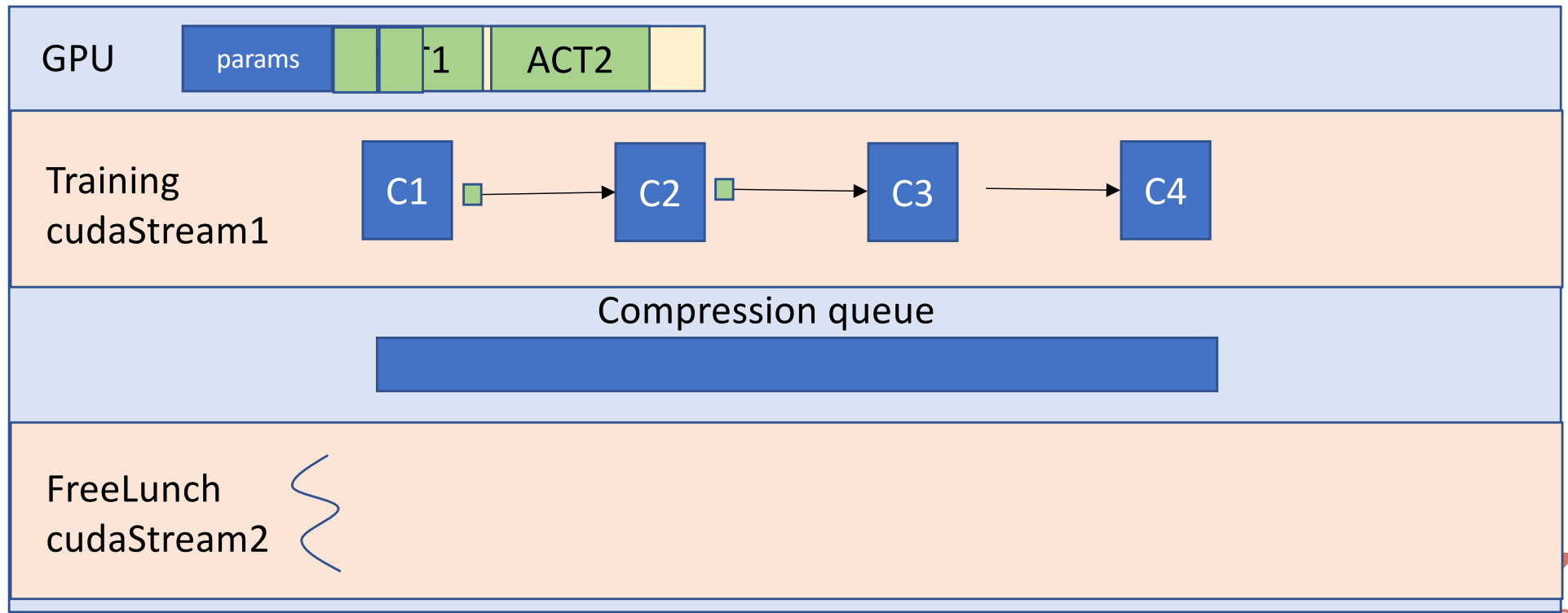
    Parallel workflow

    Optimizations:
    - Sliding Compression Workspace
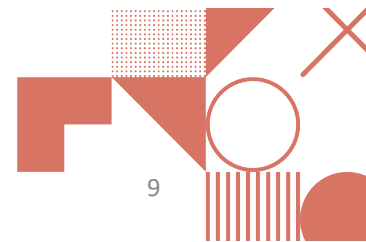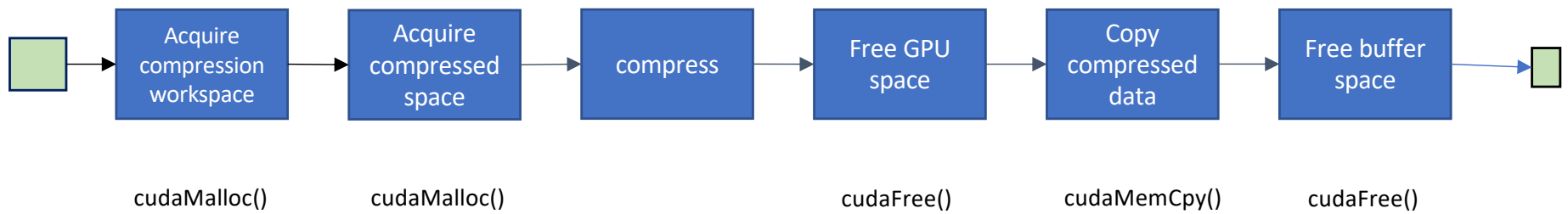    - Persistent Tensor Buffers

# Parallel workflow

**GPU**

| params | | | T1 | ACT2 | |

**Training cudaStream1**

C1 → C2 → C3 → C4

**Compression queue**

**FreeLunch cudaStream2**

**Memory operations synchronize all cuda streams!**



| | Acquire compression workspace | | Acquire compressed space | | compress | | Free GPU space | | Copy compressed data | | Free buffer space | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | cudaMalloc() | | cudaMalloc() | | | | cudaFree() | | cudaMemCpy() | | cudaFree() | |

Typical workflow implementation
Sliding compression workspace workflow
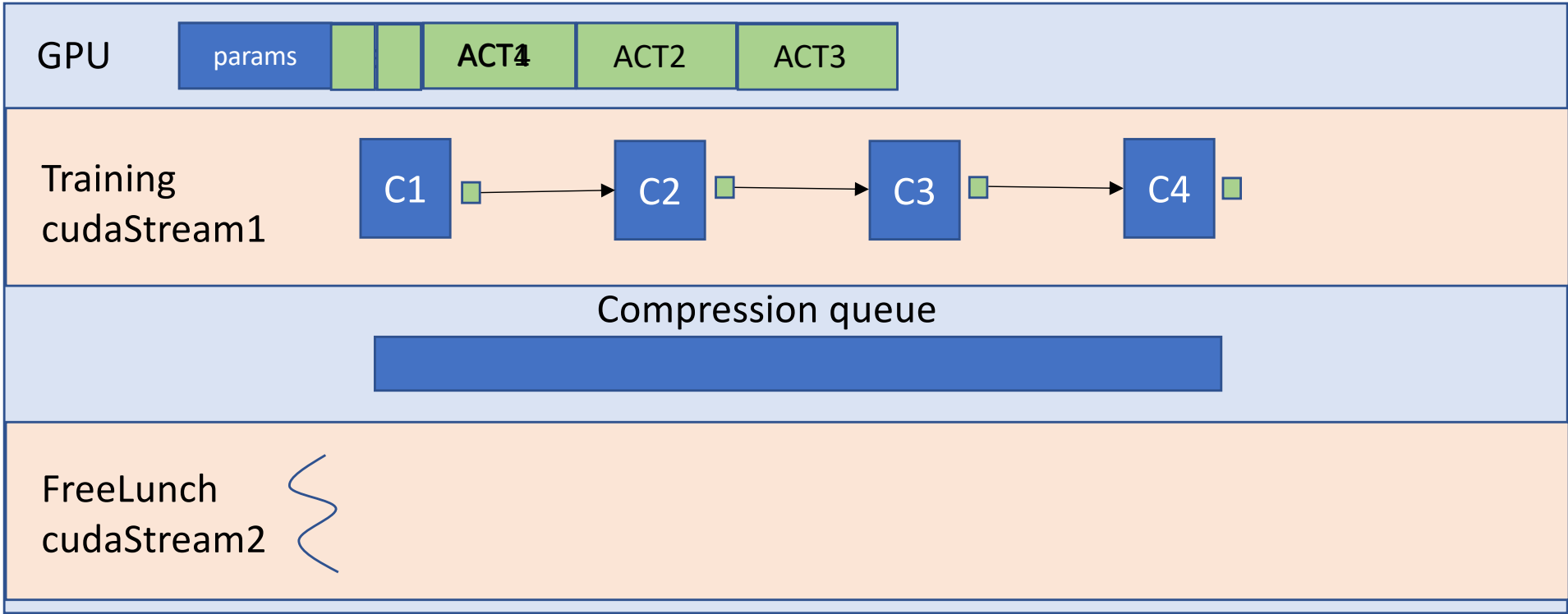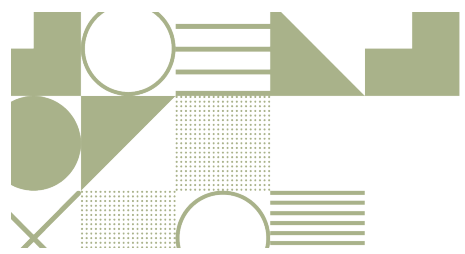


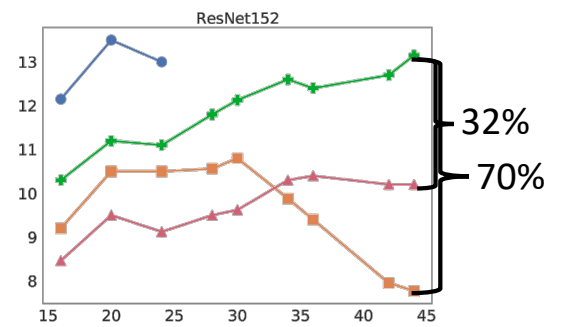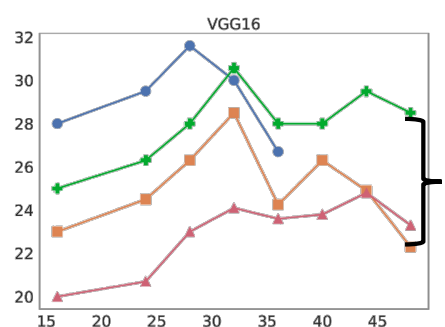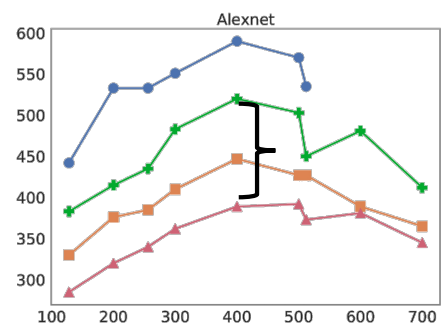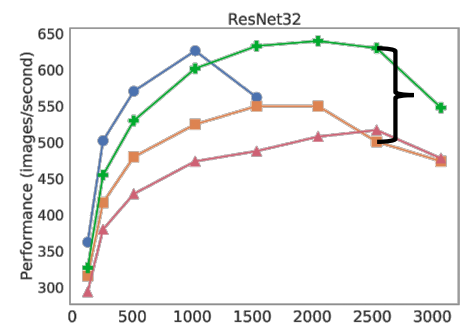This workflow introduces multiple blocking operations!

Persistent tensor buffers

- Can FreeLunch improve training throughput while reducing memory consumption of CNN training?

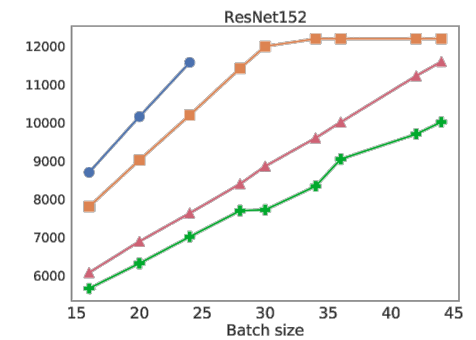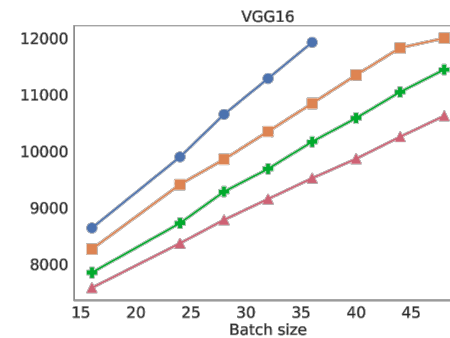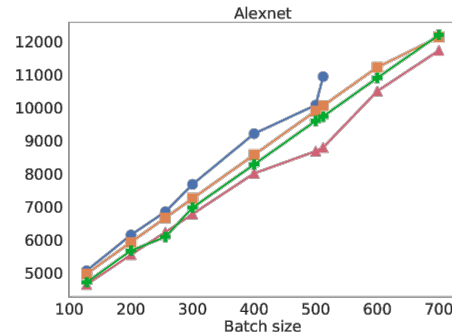- How effective are the optimizations in FreeLunch compared with other compression-based baselines?
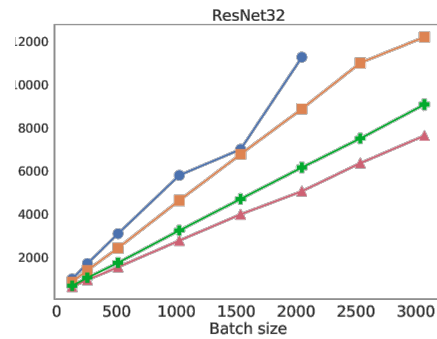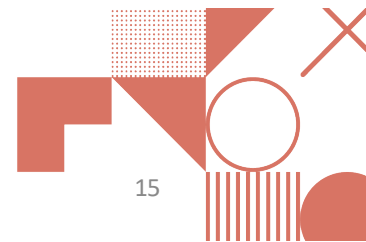
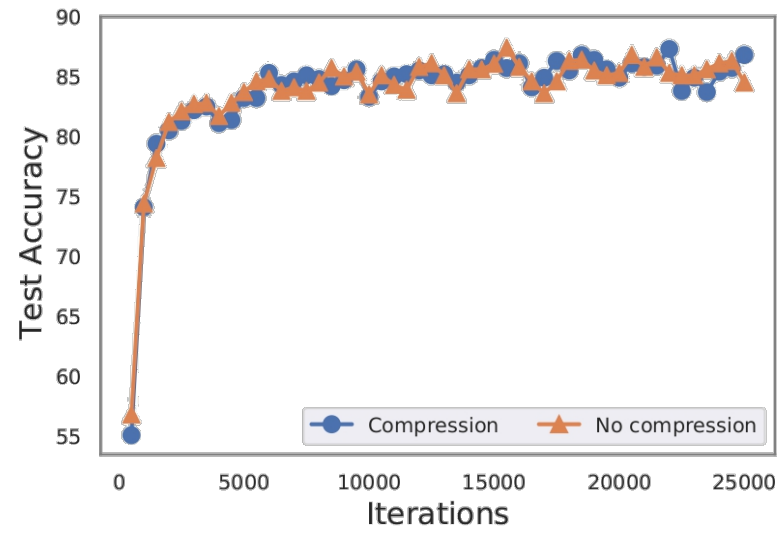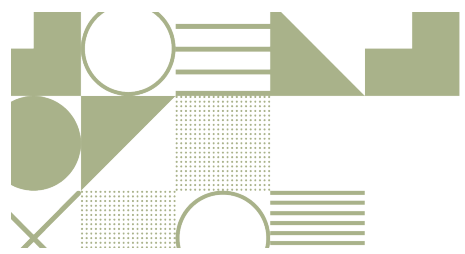# Throughput as compared to other policies
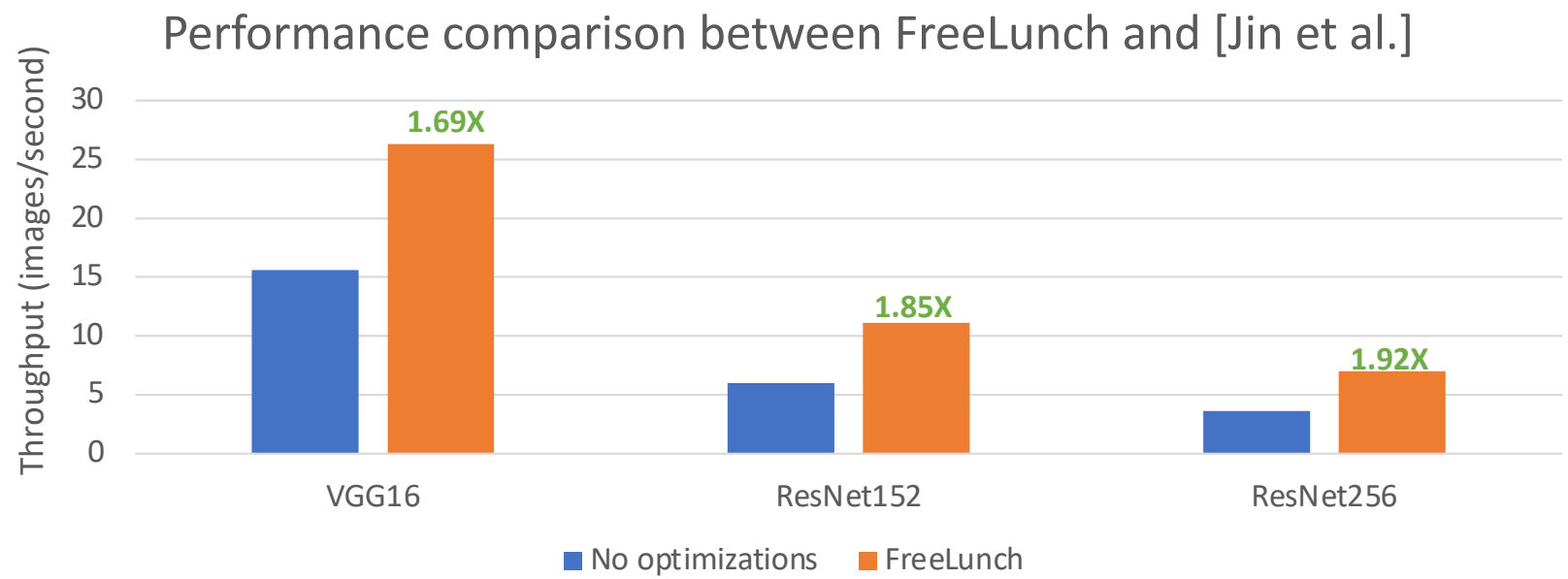
# Memory consumed as compared to other policies

# No observed impact on accuracy of model

Impact of optimizations

Performance comparison between FreeLunch and [Jin et al.]
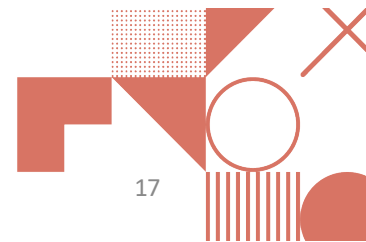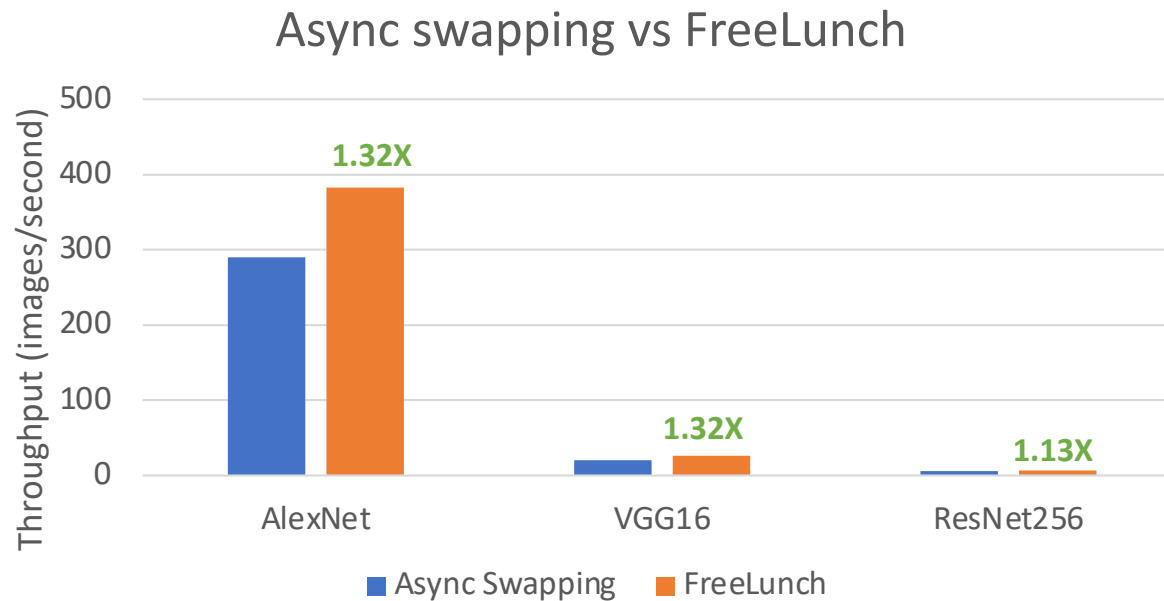
# Throughput comparison with async swapping

- Capuchin and SwapAdvisor use swapping in an async manner.
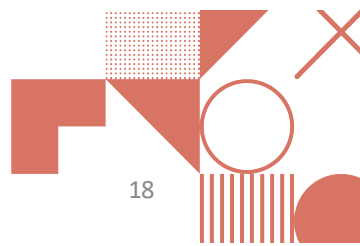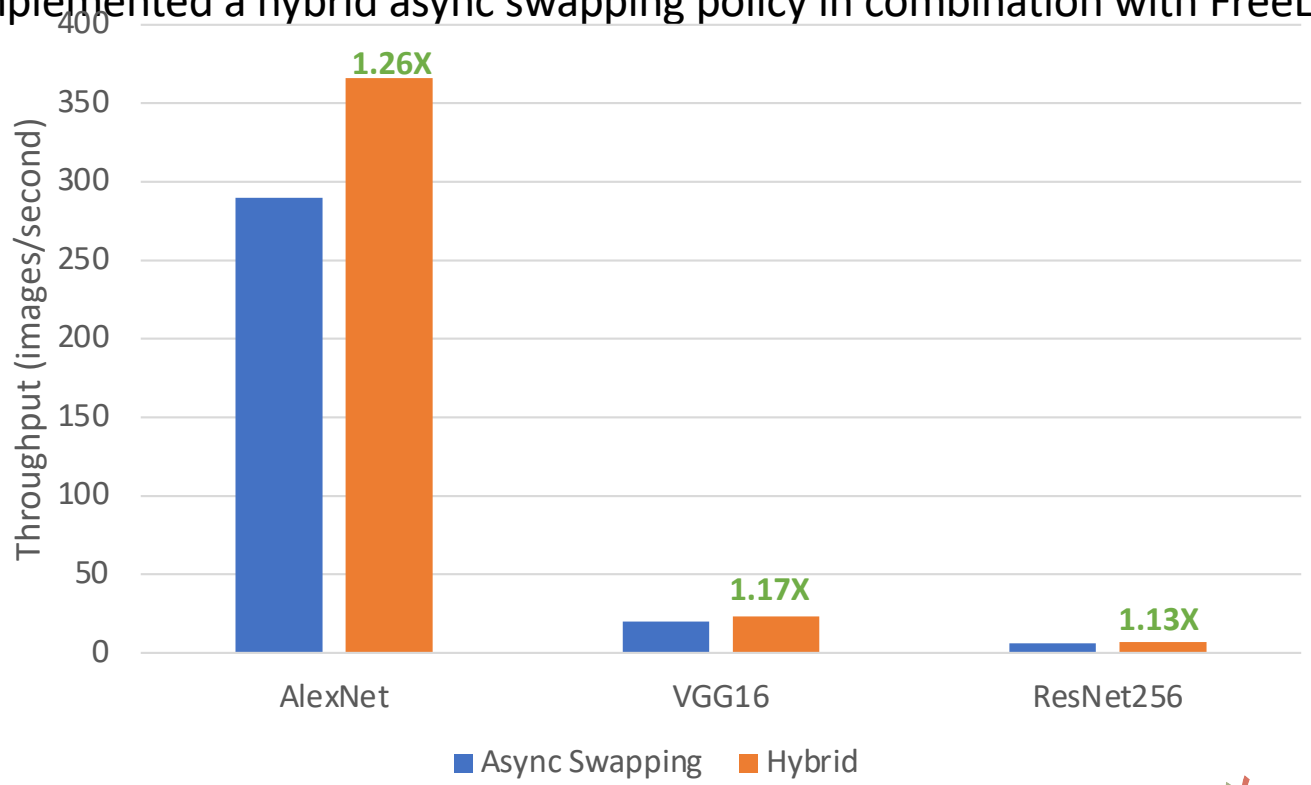- We implement async swapping and compare it to FreeLunch.

## Async swapping vs FreeLunch

Hybrid policy

• We implemented a hybrid async swapping policy in combination with FreeLunch

Async swapping vs Hybrid



18

Summary

- We introduce FreeLunch that effectively avoids the bandwidth and concurrent execution that swapping and recomputation face.

- We incorporate two optimizations as part of FreeLunch to make compression parallelizable and improve performance.

- We show that FreeLunch achieves up to 70% better throughput and up to 32% better memory consumption.