# Optimal Gradient Checkpoint Search for Arbitrary Computation Graphs

# Jianwei Feng, Dong Huang

Robotics Institute Carnegie Mellon University Pittsburgh, PA 15213 jfeng1@andrew.cmu.edu, donghuang@cmu.edu

#### **Abstract**

Deep Neural Networks(DNNs) require huge GPU memory when training on modern image/video databases. Unfortunately, the GPU memory in off-the-shelf devices is always finite, which limits the image resolutions and batch sizes that could be used for better DNN performance. Existing approaches to alleviate memory issue include better GPUs, distributed computation and gradient checkpointing. Among them, gradient checkpointing is a favorable approach as it focuses on trading computation for memory and does not require any upgrades on hardware. In gradient checkpointing, during forward, only a subset of intermediate tensors are stored, which are called Gradient Checkpoints (GCPs). Then during backward, extra local forwards are conducted to compute the missing tensors. The total training memory cost becomes the sum of (1) the memory cost of the gradient checkpoints and (2) the maximum memory cost of local forwards. To achieve maximal memory cut-offs, one needs optimal algorithms to select GCPs.

Existing gradient checkpointing approaches rely on either manual input of GCPs or heuristics-based GCP search on linear computation graphs (LCGs), and cannot apply to arbitrary computation graphs(ACGs). In this paper, we present theories and optimal algorithms on GCP selection that, for the first time, apply to ACGs and achieve maximal memory cut-offs. Extensive experiments show that our approach constantly outperforms existing approaches on LCGs, and can cut off up-to 80% of training memory<sup>1</sup> with a moderate time overhead (around 40%) on LCG and ACG DNNs, such as Alexnet, VGG, Resnet, Densenet and Inception Net.

## 1 Introduction

Deep Neural Networks(DNNs) require huge GPU memory when training on modern image/video databases. For popular backbone DNNs used in feature extraction of images, such as AlexNet (Krizhevsky, Sutskever, and Hinton 2012), VGG (Simonyan and Zisserman 2014) and ResNet (He et al. 2016), the memory cost increases quadratically with the input image resolution and network depth. For example, given an median size input

Copyright © 2020, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

<sup>1</sup>Cutting off 80% of training memory means one can double the input image size or quadruple the batch size on the same GPUs.

tensor of  $[BatchSize \times Channel \times Width \times Height] =$ [32, 3, 224, 224], ResNet101 requires around 5 GB memory in training. In more challenging tasks, DNNs that detect small objects and large number of object categories require input image resolution of more than  $600 \times 600$  (Ren et al. 2015; Singh et al. 2017; Redmon and Farhadi 2018) and can easily consume more than 10 GB with just a small batch size. The memory issue is even worse for video-based DNNs, such as CDC (Shou et al. 2017), C3D (Ji et al. 2013) and 3D-ResNet (Hara, Kataoka, and Satoh 2017). To recognize complex activities in video, the input video clips would be as long as 64 frames and could easily go beyond 10 GB using a moderate network. Memory issue also occurs in training DNN compositions, such as Generative adversarial networks (GANs). Multiple generator and discriminator networks are simultaneously stored in GPU memory.

Existing efforts to address memory issues presented three main approaches: (1) Better single GPUs. Recent GPUs provide larger memory at the expense of exponentially growing price and power consumption. For instance, from TitanXp, Quadro P6000 to Tesla V100, for 1-2.7 times increase in memory, the prices increase 2.8-8.5 times. (2) Parallelization among multiple GPUs (Dean et al. 2012; Shi et al. 2009; Langford, Smola, and Zinkevich 2009; Mcdonald et al. 2009; McDonald, Hall, and Mann 2010; Zinkevich et al. 2010; Agarwal et al. 2014; Agarwal and Duchi 2011), which requires expensive clusters, introduces substantial I/O cost, and does not reduce the total memory cost. (3) Gradient checkpointing (Chen et al. 2016; Gruslys et al. 2016), which focuses on trading computation for memory and reduces the total memory cost without any upgrade in hardware. Note that recent affordable GPUs (e.g., Nvidia GTX 1080-Ti, RTX 2080 Ti), although limited in memory (around 11GB), provide exceptional improvement in GPU cores and FLOPS. Trading computation costs for memory is a very attractive solution that make it possible to train very heavy DNNs with finite GPU memory.

The regular DNN training process consists of two alternated stages: forward and backward. Fig. 1 (a) illustrates an example of feed-forward neural networks. In the forward stage, the network takes an input tensor, and computes the tensors at each layer until producing the output. In the back-

ward stage, the difference between the output and ground truth is passed back along the network to compute the gradients at each layer. The regular training approach saves tensors at all layers during forward, because they are needed to compute gradients during backward. The total memory cost is the sum of cost over all these intermediate tensors.

Gradient checkpointing is a high-level training approach that trade extra computation time for substantial saving of GPU memory. Fig. 1 (b) illustrates its main idea. During gradient checkpoint training, only a subset of intermediate tensors (which are called gradient checkpoints (GCPs)) are stored in the first forward, and the missing tensors needed during backward are computed via extra local re-forwards. The total memory cost is the sum of the cost at the subset of intermediate tensors and the maximum memory cost among local re-forwards. Training with gradient checkpointing can lead to substantial memory reduction, with the time overhead of local re-forwards. To achieve maximal memory cut-offs, one needs optimal algorithms to select GCPs. Note that given the computation graph of a network, the GCP algorithm only needs to run once before the gradient checkpointing training and can be viewed as a preprocessing step.

In this paper, We propose sophisticate theories and efficient algorithms that automatically find the **optimal** GCPs in DNN with **arbitrary** computation graph, Using these GCPs gradient checkpoint training leads to the smallest memory cost. Comparing to existing GCP searching approaches, the optimality of our approach does not pose any assumption on computation graph and is the first optimal algorithm that applies to arbitrary computation graphs.

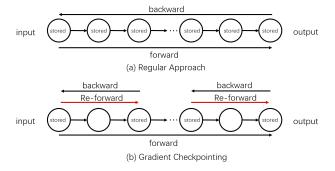


Figure 1: Regular Training Approach vs. Gradient Checkpointing. (a) The regular approach saves all tensors during forward, and uses these tensors to compute gradients during backward. (b) Gradient Checkpointing saves a subset of tensors during the first forward, and conducts extra local reforwards to compute tensors and gradients during backward.

#### 2 Related Work

To alleviate the memory pressure from a single GPU processor, many researchers utilized the well-established techniques for distributed computation (Dean et al. 2012; Shi et al. 2009; Langford, Smola, and Zinkevich 2009; Mcdonald et al. 2009; McDonald, Hall, and Mann 2010; Zinkevich et al. 2010; Agarwal et al. 2014; Agarwal and

Duchi 2011). These techniques distribute memory pressure to possibly infinite GPUs or server clusters, but do not reduce the total memory cost of DNNs.

Some researchers reduced the memory on finite hardware by optimizing computation graph of DNN and performing liveness analysis. The computation graph of DNNs describes the dependencies of tensors among layers. Liveness analysis recycles garbage to manage memory. These ideas were originated from compiler optimization (Aho, Sethi, and Ullman 1986) and has been widely adopted by deep learning frameworks: Theano (Bastien et al. 2012; Bergstra et al. 2010), MXNet (Chen et al. 2015), Tensorflow (Abadi et al. 2016) and CNTK (Yu et al. 2014). Some other techniques efficiently swap data between CPU and GPU (Wang et al. 2018; Rhu et al. 2016). These techniques usually cost extra I/O time and still do not actually reduce the total memory cost.

Other approaches focus on trading computation for memory with the idea of gradient checkpointing. Popular deep learning frameworks such as Pytorch (Paszke et al. 2017) and Tensorflow (Abadi et al. 2016) provide functions for users to manually define GCPs in computation graph and perform gradient checkpoint training. These functions are user-dependent and their performance highly relies on the selected GCPs.

There are also algorithms to solve for GCPs automatically. Chen et al. (Chen et al. 2016) develop algorithms to solve for GCPs based on heuristics in a simple case of linear computation graph (will also be discussed in section 4). Chen's approach is only applicable and not even optimal for linear computation graph (LCG). For arbitrary computation graph (ACG), Chen's approach does not apply.

Gruslys et al. (Gruslys et al. 2016) targets at gradient checkpoint problem for recurrent neural network (RNN). In recurrent neural network, the hidden state of each time step has the same size and thus consume equal amount of memory. Gruslys utilizes this characteristic and develops dynamic programming algorithm to solve for optimal GCPs for RNN given a memory budget. Gruslys's approach is based on a strong assumption on linear computation graph that the memory cost of all intermediate tensors in the computation graph is identical. Thus his approach is not even applicable when this assumption does not hold.

The main contribution of this paper is proposing algorithms to solve **optimal** GCPs for **arbitrary computation graph**. The difference between our approach and other approaches is summarized in Table.1.

#### 3 Overview

The gradient checkpoint approach consists of two steps: preprocessing and training. In the preprocessing step, we run our GCP algorithms to solve optimal GCPs given the computation graph of a network. Then in the training step, we only store tensors at the optimal GCPs during the first forward. During backward, the tensors and gradients at missing vertices are recovered by local re-forward operations. Our algorithms focus on solving optimal GCPs in the preprocessing step and is thus an **one-time effort** conducted before training. This configuration is the same as

Table 1: ✓ ✓ is both applicable and optimal, ✓ X is applicable but not optimal, XX is not applicable nor optimal.

Approach	applicable & optimal in identical cost LCG	applicable & optimal in arbitrary cost LCG	applicable & optimal in ACG	automatic	with budget
manual input	<b>√</b> X	<b>√</b> X	<b>√</b> X	Х	Х
Chen's approach	<b>//</b>	✓X	XX	✓	×
Gruslys's approach	<b>//</b>	XX	XX	✓	<b>✓</b>
ours	11	<b>//</b>	<b>11</b>	✓	X

other gradient checkpoint algorithms (Chen et al. 2016; Gruslys et al. 2016).

In section 4, we start with the Linear Computation Graph (LCG) and formulate the optimization problem of solving optimal GCPs. We first discuss a special case of LCGs, where we can easily compute an optimal solution in analytic form and understand the effectiveness of gradient checkpointing. Then we present our algorithms to solve for optimal GCPs in arbitrary LCGs.

In section 5, we present our approach on Arbitrary Computation Graphs (ACGs). Section 5 is organized by a bottom-up manner. We first introduce all the basic components, including definitions and sub-algorithms, and then the final solver based on these components.

In section 6, we present extensive experiments on networks with both linear and non-linear computation graphs. Due to space limit, we cannot put all illustrative examples in the paper. Extra illustrative examples are included in the "Extra Examples" section of the supplementary material.

In section 7, we present our conclusion for this paper.

# 4 Linear Computation Graph (LCG)

Denote a computation graph of a DNN as G=(E,V).  $E=\{e_i\}$  and  $V=\{v_i\}$  are the edges and vertices in the computation graph, respectively. The vertices represent the intermediate tensors and the edges represent DNN operations. Denote function  $l(\cdot)$  as a measure of memory cost.  $V^R$  is the subset of vertices selected as GCPs during the first forward.  $l(v_i^R)$  is defined as the memory cost of the ith gradient checkpoint in  $V^R$ . For two adjacent gradient checkpoint  $v_i^R$  and  $v_{i+1}^R$  in set  $V^R$ , suppose the ith gradient checkpoint  $v_i^R$  corresponds to vertex  $v_j$  in the original computation graph, and  $v_{i+1}^R$  corresponds to  $v_k$ , the memory cost during re-forwards from  $v_i^R$  to  $v_{i+1}^R$  is defined as  $l(v_i^R, v_{i+1}^R) = \sum_{t=j+1}^{k-1} l(v_t)$ , which is the sum of cost over all the vertices between  $v_j$  and  $v_k$  in the computation graph. Using these notations, solving the optimal GCPs is formulated as an optimization problem:

$$\min_{V^R} (\sum_i l(v_i^R) + \max_i l(v_i^R, v_{i+1}^R)), \tag{1}$$

where the  $\sum_i l(v_i^R)$  is the sum of the memory cost over all the GCPs, and  $\max_i l(v_i^R, v_{i+1}^R))$  is the maximal cost among the local re-forwards. Eq. 1 describes the peak memory during gradient checkpoint training. Solution to Eq. 1 produces the optimal GCPs in  $V^R$ .

For easy illustration, we start by solving Eqn. 1 on Linear Computation Graphs (LCG) (Fig. 2 (a)). For LCGs, Eqn. 1 can be solved in two cases.

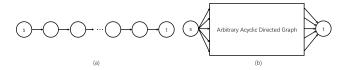


Figure 2: (a) Linear Computation Graph (LCG). "s" denotes the source vertex, "t" denotes the target vertex. (b) Arbitrary Computation Graph (ACG). The structure between "s" and "t" vertices may contain arbitrary branches and connections.

Case(1) LCG with Identical Vertex Cost: Suppose a LCG has N vertices, each of which has the same cost as  $l(v_i) = \frac{1}{N}$  and the total cost of these N vertices is 1. Obviously, the optimal solution is reached when vertices in  $V^R$  are distributed evenly in the LCG. Suppose the number of vertices in  $V^R$  is k. The total cost is then  $\frac{k}{N} + \frac{1}{k}$ . The optimal solution of Eqn. 1 is achieved when  $k = \sqrt{N}$ , and the optimal total cost is  $\frac{2}{\sqrt{N}}$ .

From Case(1), we can get a sense of the effectiveness of gradient checkpointing. The original memory cost is 1, and can be reduced to  $\frac{2}{\sqrt{N}}$  at the time overhead of extra local forwards. When the network is deep, i.e. N is large, huge amount of memory cost can be cut off. For example, when N=100, we can reduce the memory cost to  $\frac{1}{5}$  of the original cost. Chen's approach (Chen et al. 2016) is developed exactly from this observation and thus is only optimal in this case.

Case (2) LCG with Non-identical Vertex Cost: When the assumption of identical cost does not hold, the solution to Eqn. 1 does not have an analytic form. Denote the maximal Re-forward cost  $\max_i l(v_i^R, v_{i+1}^R))$  as a constant C, and the solution to Eqn. 1 is reduced to solving for  $\min_{V^R} \sum_i l(v_i)$ , such that all the re-forward memory costs satisfy the constraint  $l(v_i^R, v_{i+1}^R)) \leq C$ .

Given a constant C as this constraint, we can solve the reduced problem by constructing a new graph, called Accessibility Graph  $G^A = (E^A, V)$ . The edges of  $G^A$ , called Accessibility Edge  $e^A_{ij}$ , exists between vertex  $v_i$  and  $v_j$  if and only if  $l(v_i, v_j) \leq C$ . Now the constraints are all encoded in the accessibility graph, we can solve the unconstrained problem  $\min_{V^R} \sum_i l(v^R_i)$ , which is equivalent to finding the shortest path from the source vertex and the target vertex in the Accessibility Graph. Notice that in the optimal solution

of Eqn. 1,  $\max_i l(v_i^R, v_{i+1}^R)) = C = l(v_j, v_k)$ . C would be the cost  $l(v_j, v_k)$  between a vertex pair. Therefore, to determine C of an optimal solution, we can simply traverse all possible C by using the loss of every vertex pair, and find optimal solution under each C. The best of it would then be the optimal solution of Eqn. 1. Algorithm 1 summarizes the steps for searching an optimal solution for LCGs. For a computation graph with N vertices, the time complexity of Algorithm 1 is  $O(N^4)$ .

Algorithm 1 Linear Computation Graph (LCG) Solver

**Input:** a linear computation graph G **Output:** optimal GCPs  $V^R$ 

- 1: **for** each vertex pair  $(v_i, v_j)$  in G **do**
- 2: Set the maximal term as  $l(v_i, v_j)$
- 3: Construct Accessibility Graph
- 4: Find the shortest path in the Accessibility Graph as a candidate solution  $V^R$
- 5: Compute the total cost of candidate solution  $V^R$
- 6: Save the solution  $V^R$  if the total cost is smaller.

# 5 Arbitrary Computation Graph(ACG)

As the generalization of LCGs, we present theory and algorithms for DNNs with Arbitrary Computation Graphs (ACG), in particular the acyclic directed graphs (Fig. 2 (b)).

For ACGs, we follow the same idea in LCGs: traverse all possible max term C and solve a constrained problem for each C. The following subsections are organized in a bottom-up manner: we first introduce all the basic components and then the final algorithm. Due to space limit, proofs and further analysis are in the supplementary material.

#### 5.1 Definition and Theorem



Figure 3: Closed Set Examples: (a) Closed set in a graph. there cannot exist a closed set between  $v_2$  and  $v_4$  because  $v_3$  depends on  $v_1$ . There can exist a closed set between  $v_1$  and  $v_3$  because  $v_2$  doesn't depend on any other vertex. (b) Splittable Closed Set (Type 1).  $v_2$  is the splitting vertex of  $s_{13}$ . (c) Branched Closed Set (Type 2). (d) Non-branched Closed Set (Type 3).

**Definition 1** Closed Set: A set s containing vertices and edges is a closed set if and only if it satisfies the following three properties: 1. All the vertices of s have a common ancestor  $v_i$  and a common descendent  $v_j$ ; 2. Denote the vertex subset of s as V, edge subset as E, and the set of edges between two arbitrary vertices of  $V \cup \{v_i, v_j\}$  is E', the edge from  $v_i$  to  $v_j$  (if exists) as  $e_{ij}$ . E must either be E'

or  $E' - \{e_{ij}\}$ ; 3. An arbitrary  $v_1 \in V$  doesn't have edge with another arbitrary  $v_2 \notin V \cup \{v_i, v_j\}$ . For multiple valid closed sets between  $v_i$  and  $v_j$ , we denote the largest one as  $s_{ij}$ 

**Definition 2** 
$$[s_{ij}] = s_{ij} \cup \{v_i, v_j\}.$$
  $[s_{ij}) = s_{ij} \cup \{v_i\}.$   $(s_{ij}] = s_{ij} \cup \{v_j\}$ 

Closed Set can be viewed as an independent sub-graph that has no cross edge with the other parts in the graph. For convenience, in the definition we exclude its source and target vertex. In the definition, Property 1 confines the set only has one source vertex and one target vertex. Property 2 confines the edge subsets of s to be one of two cases: E' or  $E' - \{e_{ij}\}$ . Both cases are valid although they have different edges. Property 3 guarantees the independence of such a set s, meaning that the vertices within s have no connections with other vertices outside  $s \cup \{v_i, v_i\}$ . This property is for the independent backward and re-forward between GCPs. As there might be multiple valid closed sets between  $v_i$  and  $v_i$ , which corresponds to the **Branched Closed Set** in Definition 5, we denote the largest closed set between  $v_i$  and  $v_i$  as  $s_{ij}$  and denote smaller closed set with an extra superscript, such as  $s_{ij}^1$ .

**Definition 3** Splitting Vertex: A vertex  $v_t \in s_{ij}$  is a splitting vertex of  $s_{ij}$  if and only if  $s_{it}$  exists,  $s_{tj}$  exists and  $s_{ij} = s_{it} \cup s_{tj} \cup \{v_t\}$  and  $s_{it} \cap s_{tj} = \emptyset$ 

**Definition 4** *Splittable Closed Set (Type 1):* A closed set with at least one splitting vertex.

**Definition 5** *Branched Closed Set (Type 2):* A closed set is branched if it has 0 splitting vertex and can be divided into branches:  $s_{ij} = s_{ij}^1 \cup s_{ij}^2$  and  $s_{ij}^1 \cap s_{ij}^2 = \emptyset$ 

**Definition 6** Non-branched Closed Set (Type 3): A closed set  $s_{ij}$  is non-branched if it has 0 splitting vertex and no branch:  $\exists s_{ij}^1 \subseteq s_{ij}$ 

The definition of **Splitting Vertex** is to describe whether a closed set can be divided into two linearly arranged closed set. A closed set is splittable if it has at least one splitting vertex and is defined as **Closed Set Type 1**. Among closed sets with no splitting vertex, we categorize the closed sets with branches as **Closed Set Type 2**, and the closed set without branches as **Closed Set Type 3**. The examples of different types of closed sets are shown in Fig. 3.

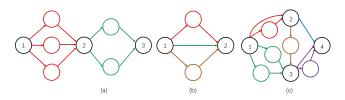


Figure 4: Examples on Division of Closed Sets. Members of different Closed Sets are colored. (a) Division of closed set Type 1. The division is  $\{[s_{12}], [s_{23}]\}$ . (b) Division of closed set Type 2. The division is  $\{[s_{12}], [s_{12}^3], [s_{12}^3]\}$ . (c) Division of closed set Type 3. The division is  $\{[s_{12}], [s_{13}], [s_{23}], [s_{24}], [s_{34}]\}$ 

All closed sets can be further decomposed into a set of smaller closed sets, which is called **the Division of Closed Set**. Closed set type 1 can be divided into linearly arranged segments connected by the splitting vertices. Closed set type 2 can be divided into its branches. Closed set type 3 requires closer investigation. We don't want trivial division, for example, division that is formed by every edge in the closed set. We define **Maximal Split** to describe the division such that each member of the division is as large as possible. An example of maximal split is shown in Fig. 4 (c). In the definition of maximal split, the term maximal is implied by saying that any subset of this split cannot be combined into a single closed set. If it can, then the maximal split will be formed by this larger closed set and all the rest of the previous split. For closed set type 3, we use its maximal split as its division.

**Definition 7** *Maximal Split:*  $\{[s_{pq}]\}$  *is a maximal split of non-branched*  $s_{ij}$  *if*  $[s_{ij}] = \cup \{[s_{pq}]\}$  *and*  $\forall s_{ab}, s_{cd} \in \{[s_{pq}]\}, s_{ab} \cap s_{cd} = \emptyset$  *and*  $\exists \{[s'_{pq}]\} \subseteq \{[s_{pq}]\}$  *such that*  $\cup \{[s'_{pq}]\} = [s_{kt}] \subseteq [s_{ij}]$ 

**Definition 8** *Division of Closed Set:* For a Closed set type 1, its division is the linear segments separated by all its splitting vertices; for Type 2, its division is all its branches, any of which cannot be divided into more branches; for Type 3, its division is its maximal split.

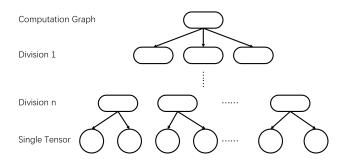


Figure 5: Division tree of a computation graph. The root node is the whole computation graph (largest closed set). All the leaf nodes are single tensors (smallest closed set). Every other node except root and leaves is a member of the division of its parent.

**Definition 9** *Division Tree:* Division tree is a representation of a computation graph, where the root node is the whole computation graph, the leaf nodes are all the single tensors in the computation graph, and for a non-leaf node, its children is the members of its division.

With the division of 3 types of closed sets, the computation graph can be re-organized into a division tree (Figure 5) where a non-leaf node would be a closed set and its children would be its corresponding division. The root node is the whole computation graph, the largest closed set, and the leaf nodes would be single tensors in the computation graph. With the division tree, we can use divide-and-conquer method to search for optimal solution of Eqn.1.

**Theorem 1** The division tree of a computation graph is unique and complete.

The uniqueness of the division tree indicates that the optimal solution of the division tree would also be the optimal solution of the whole computation graph. The completeness indicates that the division tree has included all the possible members of solutions and represents the whole search space for finding the optimal solution. This theorem also indicates that the optimal solution in the division tree is equivalent to the optimal solution of the original problem. Theorem 1 is proved in the supplementary material.

## 5.2 Algorithm

We search optimal GCP solutions for ACGs by solving several sub-problems using Algorithm 2-4 respectively. Based on these components, we present our final solver as Algorithm 5.

**Algorithm 2** judges whether a vertex is a splitting vertex of a closed set. This algorithm mainly follows the Definition 3 and uses vertex set to check the property of a splitting vertex. With this algorithm, we can judge whether a closed set is type 1 and get its division if it is, as the division of type 1 is simply closed sets seperated by the splitting vertices. Suppose there are N vertices in closed set  $s_{ij}$ , the time complexity of **Algorithm 2** is  $O(N^2)$ .

**Algorithm 2** Judge whether a vertex  $v_t$  is a splitting vertex of closed set  $s_{ij}$ 

**Input:** closed set  $s_{ij}$ , vertex  $v_t$ 

Output: True/False

- 1: Let  $\{v_{in}\}$  be the vertices of all the vertices within  $[s_{ij}]$  that have paths to  $v_t$ . Let  $\{v_{out}\}$  be the vertices of all the vertices within  $[s_{ij}]$  that have paths from  $v_t$ .
- 2: **if**  $\{v_{in}\} \cup \{V_{out}\} \cup \{v_t\} = \{v|v \in [s_{ij}]\}$  and  $\{v_{in}\} \cap \{V_{out}\} = \emptyset$  and  $\not\exists v_1 \in \{v_{in}\}, v_2 \in \{v_{out}\}, v_1, v_2$  have connections **then**
- 3: Return True
- 4: else
- 5: Return False

**Algorithm 3** examines whether a closed set is branched. It uses a growing algorithm to check whether an independent subpart of this closed set can form a closed set. If a non-trivial closed set  $s_{ij}$  has an edge from  $v_i$  to  $v_j$ , then it is branched because this edge itself can be treated as a closed set. Combined with Algorithm 2, we can know the type of a closed set and get its division if it is type 2. Suppose there are N vertices in  $s_{ij}$ , the time complexity of **Algorithm 3** is  $O(N^2)$ .

**Algorithm 4** addresses the problem of finding the maximal split, the division of a closed set type 3  $s_{ij}$ . First get all the possible closed sets within  $s_{ij}$  and use a property of maximal split to judge whether this closed set is a member of the maximal split. The property is: there cannot exist another closed set  $s_{ab} \subsetneq s_{ij}$  but contains any member of this maximal split. This property is proved in Lemma 6 of the supplementary material. Suppose there are N vertices in  $s_{ij}$ , the time complexity of **Algorithm 4** is  $O(N^4)$ .

## Algorithm 3 Judge whether $s_{ij}$ is branched

if  $s = \{v \in s_{ij}\}$  then

Return False

Return True

**Input:** closed set  $s_{ij}$ 

Output: True/False 1: **if**  $s_{ij}$  has at least 1 vertex **then** if  $s_{ij}$  includes an edge from  $v_i$  to  $v_j$  then 3: Return True 4: 5: Initialize a vertex set  $s = \{v_k\}$ .  $v_k \in s_{ij}$  is a randomly chosen vertex. while True do 6: 7: For any  $v_t \in s_{ij}, v_t \not\in s$  that has connection to any  $v_k \in s$ , add  $v_t$  to s. 8: if No more vertex can be added to s then 9: Break

Algorithm 4 Find the maximal split of a non-branched  $s_{ij}$  with 0 splitting vertex

**Input:** closed set  $s_{ij}$ 

else

Return False

10:

11: 12:

13:

15:

14: **else** 

**Output:** maximal split of  $s_{ij}$  (a set of closed sets)

1: for each vertex pair  $(v_k, v_t)$  except  $(v_i, v_j)$  in  $[s_{ij}]$  do

- 2: For all the vertices  $\{v\}$  that have paths from  $v_k$  and have paths to  $v_t$ .
- 3: **if**  $\not\supseteq v_2 \not\in \{v\}$  and  $v_2 \neq v_k, v_t, v_2$  has connection to a  $v_1 \in \{v\}$  **then**
- 4: Form a closed set  $s_{kt}$  with all these vertices.
- 5: **for** each formed closed set  $s_{kt}$  **do**
- 6: If there doesn't exist a  $s_{ab}$  such that  $s_{kt} \subsetneq s_{ab} \subsetneq s_{ij}$ , put  $s_{kt}$  into the maximal split.

**Algorithm 2-4** are the sub-components of our final solver. With them, we can categorize and get the division of a closed set, and reform the computation graph as division tree to set up recursions easily.

**Algorithm 5** describes how we do recursion in division tree. Given a max term C as the constraint, we propose a greedy idea: for a closed set, never expand to its division unless the its cost exceed the constraint. In other word, if the constraint doesn't allow a leap over this closed set, we expand it to its division and look for more GCPs in the next level. Otherwise, there's no need to expand it since the closed set already satisfies the constraint and doesn't need more GCPs inside. Once a closed set is expanded, its source and target vertex will be added to GCPs. For closed set type 1, if some children of it are expanded, the rest reforms a few linear segments and can be further optimized by the LCG solver under the constraint C (inside the loop of Algorithm 1). If some children of the closed set type 2 or 3 are expanded, there is no optimization for the unexpanded closed

sets.

**Algorithm 6** is the final solver of Eqn. 1 for ACGs. First, the division tree of the computation graph is built with Algorithms 2-4. Similar to the LCG solver, a list of max term is formed to contain the costs of all the possible closed sets for traverse. Then for each max term C, the recursion function in Algorithm 5 is called with the whole computation graph (the largest closed set) as the input. Suppose there are N vertices in computation graph, the overall time complexity of **Algorithm 6** is  $O(N^4)$ . Note that given an ACG, ACG Solver is a pre-procession step and only needs to run once before the gradient checkpointing training.

```
Algorithm 5 Recursion in division tree: V^R \leftarrow recur(s, V^R, C)
```

**Input:** a closed set s, current GCPs  $V^R$ , max term C **Output:** new GCPs  $V^R$  with GCPs in s added

- 1: for each child closed set  $s_{ij}$  in the division of s do
- 2: **if** cost of  $s_i$  greater than C **then**
- 3: Add source and target vertex  $v_i$  and  $v_j$  to  $V^R$
- 4:  $V^R = recur(s_{ij}, V^R, C)$
- 5: **if** closed set s is type 1 **then**
- for each segment s<sub>k</sub> separated by the expanded children closed set do
- 7: Solve  $s_k$  with LCG Solver under constraint C and add to  $V^R$ :  $V^R = V^R + LCGSolver(s_k, C)$

Algorithm 6 Arbitrary Computation Graph (ACG) Solver

**Input:** an arbitrary computation graph G **Output:** optimal GCPs  $V^R$ 

- 1: Get all possible closed set and their costs. Use their costs to form the max term list  $\{c\}$ .
- 2: Reorganize the computation graph into a division tree: from the root node (the computation graph), build its children from its division, until all the leaf nodes are single tensors.
- 3: **for** each possible max term C in max term list  $\{c\}$  **do**
- 4: Set  $V^{R}$  empty
- 5:  $V^R = recur(G, V^R, C)$
- 6: Summarize the total loss, save the current solution  $V^R$  if it's better.

## 6 Experiment

We evaluated our approach on (1) networks with linear computation graphs, such as Alexnet (Krizhevsky, Sutskever, and Hinton 2012) and Vgg (Simonyan and Zisserman 2014). (2) networks with non-linear computation graphs, such as Resnet (He et al. 2016), Densenet (Huang et al. 2017) and Inception net (Szegedy et al. 2016). In Table 2, We compared our approach with Chen's approach (Chen et al. 2016) and a random baseline and the regular training approach. Note that (Chen et al. 2016) only works on linear computation graphs and is not applicable to non-linear computation

Table 2: Training memory cut-offs and time overheads of gradient checkpointing training with respect to regular training. The GCPs used in gradient checkpointing training are provided by random baseline, Chen's (Chen et al. 2016) and our GCP algorithm, respectively. Note that random baseline reports the **best** number over 10 random trials. (Chen et al. 2016) does not apply to non-linear networks (or ACGs). During gradient checkpointing training, random, Chen's and our approach has the same time overhead, therefore they share the "Checkpointing Time" and "Checkpointing Time Overhead" columns in the table.

	Regular	Random	Chen's	Ours	Ours	Regular	Checkpointing	Checkpointing
Linear network	Memory	Memory	Memory	Memory	Memory	Time	Time	Time
	(MB)	(MB)↓	(MB)↓	(MB)↓	Cut-offs↑	(Sec)	(Sec)↓	Overhead↓
Alexnet batch 1024	3550	2944	3108	2620	26%	1.295	1.816	40%
Vgg11 batch 64	2976	2314	2292	1802	39%	0.606	0.819	35%
Vgg13 batch 64	4152	2720	2586	2586	38%	1.020	1.333	31%
Vgg16 batch 64	4470	3210	2894	2586	42%	1.307	1.696	30%
Vgg19 batch 64	4788	3098	2894	2502	48%	1.593	2.060	29%
Non-linear network	Regular	Random	Chen's	Ours	Ours	Regular	Checkpointing	Checkpointing
	Memory	Memory	Memory	Memory	Memory	Time	Time	Time
	(MB)	(MB)↓	(MB)↓	(MB) ↓	Cut-offs↑	(Sec)	(Sec)↓	Overhead↓
Resnet18 batch 256	5402	3636	N/A	2898	46%	1.144	1.599	40%
Resnet34 batch 128	3900	2108	N/A	1544	60%	1.041	1.419	36%
Resnet50 batch 64	5206	2714	N/A	1798	65%	0.740	1.027	40%
Resnet101 batch 32	3812	1500	N/A	970	75%	0.624	0.853	37%
Resnet152 batch 16	2810	1024	N/A	564	80%	0.450	0.628	39%
Densenet121 batch 32	3984	2132	N/A	776	81%	0.558	0.789	42%
Densenet161 batch 16	3658	1534	N/A	616	83%	0.511	0.708	39%
Densenet169 batch 32	4826	2128	N/A	848	82%	0.714	1.022	43%
Densenet201 batch 16	3164	1440	N/A	582	82%	0.449	0.651	45%
Inceptionv3 batch 32	2976	1244	N/A	910	69%	0.563	0.763	35%

graphs. Our approach directly works on arbitrary computation graphs. For random baseline, we randomly select 1-5 GCPs among all vertices in the computation graph. We repeat this random selection for 10 times and report the **best** solution (i.e. the solution with minimal memory consumption) among 10 trials. For non-linear networks, random selection can yield invalid solution (unable to do independent forward and backward between GCPs). In this case, we repeat random selection process until we have 10 valid solutions and report the **best** among them.

All experiments were conducted in Pytorch 1.0. GPU memory costs (MB) are measured in Float32. To remove irrelevant cost, such as model weights and Pytorch CUDA interface, training memory costs were computed as the memory difference under two input sizes. For example, for Alexnet, we first measure the training memory under input size [BatchSize, Channel, Width, Height] =[16, 3, 224, 224] as  $r_1$  and that under input [32, 3, 224, 224]as  $r_2$ . The Alexnet memory cost under input [16, 3, 224, 224]is reported as  $r_2 - r_1$ . To make the best use of public codes, the input to Inception net is [BatchSize, 3, 300, 300], and the input to all other networks is [BatchSize, 3, 224, 224]. We also measured the training time per iteration (Sec) averaging over 20 iterations. As Random baseline, Chen's approach and our approach all conduct one extra forwarding, these three approaches have the same time overhead and share the "Checkpoint Time" and "Checkpoint Time Overhead" columns in the Table. 2.

Table. 2 shows that our approach cuts down great amount of memory from the regular approach at reasonable time overheads. For instance, for linear network Vgg19, 48% memory was cut down at the expense of 29% time overhead.

Due to our optimal solution on computation graphs, gradient checkpointing outperforms Chen's approach and also constantly outperforms the best solution of 10 random trials. For deeper and non-linear networks, Chen's approach does not apply, while our approach can still give substantial memory cut and constantly outperform the best solution of 10 random trials. On the deepest Resnet152, 80% memory cut was achieved with only 39% time overhead. For Densenet series, more than 80% memory cuts were achieved with around 40% time overhead.

## 7 Conclusion

Gradient checkpointing is a fundamental approach that makes it possible to train very heavy DNNs on finite GPU hardware. However, existing efforts on this approach are stagnant at heuristic GCP searching and LCGs. To our knowledge, our theoretical and algorithmic results are the first top-down work that achieve an optimal memory GCP solution for DNNs with arbitrary computation graphs. Our advance of gradient checkpointing is general and can be further integrated with any low-level techniques such as distributed computing, GPU/CPU swapping, computation graph optimization and liveness analysis.

#### References

Abadi, M.; Agarwal, A.; Barham, P.; Brevdo, E.; Chen, Z.; Citro, C.; Corrado, G. S.; Davis, A.; Dean, J.; Devin, M.; et al. 2016. Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *arXiv* preprint *arXiv*:1603.04467.

Agarwal, A., and Duchi, J. C. 2011. Distributed delayed

- stochastic optimization. In *Advances in Neural Information Processing Systems*, 873–881.
- Agarwal, A.; Chapelle, O.; Dudík, M.; and Langford, J. 2014. A reliable effective terascale linear learning system. *The Journal of Machine Learning Research* 15(1):1111–1133.
- Aho, A. V.; Sethi, R.; and Ullman, J. D. 1986. Compilers, principles, techniques. *Addison Wesley* 7(8):9.
- Bastien, F.; Lamblin, P.; Pascanu, R.; Bergstra, J.; Goodfellow, I.; Bergeron, A.; Bouchard, N.; Warde-Farley, D.; and Bengio, Y. 2012. Theano: new features and speed improvements. *arXiv preprint arXiv:1211.5590*.
- Bergstra, J.; Breuleux, O.; Bastien, F.; Lamblin, P.; Pascanu, R.; Desjardins, G.; Turian, J.; Warde-Farley, D.; and Bengio, Y. 2010. Theano: A cpu and gpu math compiler in python. In *Proc. 9th Python in Science Conf*, volume 1.
- Chen, T.; Li, M.; Li, Y.; Lin, M.; Wang, N.; Wang, M.; Xiao, T.; Xu, B.; Zhang, C.; and Zhang, Z. 2015. Mxnet: A flexible and efficient machine learning library for heterogeneous distributed systems. *arXiv preprint arXiv:1512.01274*.
- Chen, T.; Xu, B.; Zhang, C.; and Guestrin, C. 2016. Training deep nets with sublinear memory cost. *arXiv* preprint *arXiv*:1604.06174.
- Dean, J.; Corrado, G.; Monga, R.; Chen, K.; Devin, M.; Mao, M.; Senior, A.; Tucker, P.; Yang, K.; Le, Q. V.; et al. 2012. Large scale distributed deep networks. In *Advances in neural information processing systems*, 1223–1231.
- Gruslys, A.; Munos, R.; Danihelka, I.; Lanctot, M.; and Graves, A. 2016. Memory-efficient backpropagation through time. In *Advances in Neural Information Processing Systems*, 4125–4133.
- Hara, K.; Kataoka, H.; and Satoh, Y. 2017. Learning spatio-temporal features with 3d residual networks for action recognition. In *Proceedings of the ICCV Workshop on Action, Gesture, and Emotion Recognition*, volume 2, 4.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 770–778.
- Huang, G.; Liu, Z.; Weinberger, K. Q.; and van der Maaten, L. 2017. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, volume 1, 3.
- Ji, S.; Xu, W.; Yang, M.; and Yu, K. 2013. 3d convolutional neural networks for human action recognition. *IEEE transactions on pattern analysis and machine intelligence* 35(1):221–231.
- Krizhevsky, A.; Sutskever, I.; and Hinton, G. E. 2012. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, 1097–1105.
- Langford, J.; Smola, A. J.; and Zinkevich, M. 2009. Slow learners are fast. *Advances in Neural Information Processing Systems* 22:2331–2339.
- Mcdonald, R.; Mohri, M.; Silberman, N.; Walker, D.; and Mann, G. S. 2009. Efficient large-scale distributed training

- of conditional maximum entropy models. In Advances in Neural Information Processing Systems, 1231–1239.
- McDonald, R.; Hall, K.; and Mann, G. 2010. Distributed training strategies for the structured perceptron. In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, 456–464. Association for Computational Linguistics.
- Paszke, A.; Gross, S.; Chintala, S.; Chanan, G.; Yang, E.; DeVito, Z.; Lin, Z.; Desmaison, A.; Antiga, L.; and Lerer, A. 2017. Automatic differentiation in pytorch.
- Redmon, J., and Farhadi, A. 2018. Yolov3: An incremental improvement. *CoRR* abs/1804.02767.
- Ren, S.; He, K.; Girshick, R.; and Sun, J. 2015. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, 91–99.
- Rhu, M.; Gimelshein, N.; Clemons, J.; Zulfiqar, A.; and Keckler, S. W. 2016. vdnn: Virtualized deep neural networks for scalable, memory-efficient neural network design. In *Microarchitecture (MICRO)*, 2016 49th Annual IEEE/ACM International Symposium on, 1–13. IEEE.
- Shi, Q.; Petterson, J.; Dror, G.; Langford, J.; Smola, A.; Strehl, A.; and Vishwanathan, V. 2009. Hash kernels. In *Artificial intelligence and statistics*, 496–503.
- Shou, Z.; Chan, J.; Zareian, A.; Miyazawa, K.; and Chang, S.-F. 2017. Cdc: convolutional-de-convolutional networks for precise temporal action localization in untrimmed videos. In 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 1417–1426. IEEE.
- Simonyan, K., and Zisserman, A. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv* preprint arXiv:1409.1556.
- Singh, B.; Li, H.; Sharma, A.; and Davis, L. S. 2017. R-FCN-3000 at 30fps: Decoupling detection and classification. *CoRR* abs/1712.01802.
- Szegedy, C.; Vanhoucke, V.; Ioffe, S.; Shlens, J.; and Wojna, Z. 2016. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2818–2826.
- Wang, L.; Ye, J.; Zhao, Y.; Wu, W.; Li, A.; Song, S. L.; Xu, Z.; and Kraska, T. 2018. Superneurons: Dynamic gpu memory management for training deep neural networks. *arXiv* preprint arXiv:1801.04380.
- Yu, D.; Eversole, A.; Seltzer, M.; Yao, K.; Huang, Z.; Guenter, B.; Kuchaiev, O.; Zhang, Y.; Seide, F.; Wang, H.; et al. 2014. An introduction to computational networks and the computational network toolkit. *Microsoft Technical Report MSR-TR-2014–112*.
- Zinkevich, M.; Weimer, M.; Li, L.; and Smola, A. J. 2010. Parallelized stochastic gradient descent. In *Advances in neural information processing systems*, 2595–2603.