



Inshrinkerator: Compressing Deep Learning Training Checkpoints via Dynamic Quantization

Amey Agrawal
Georgia Institute of Technology
ameyagrwal@gatech.edu

Venkata Prabhakara
Sarath Nookala*
Meta Inc.
sarathnookala@meta.com

Sameer Reddy*
Cisco Inc.
samredd@cisco.com

Vidushi Vashishth*
Google Inc.
vvashishth@google.com

Satwik Bhattamishra
University of Oxford
satwik.bmishra@cs.ox.ac.uk

Kexin Rong
Georgia Institute of Technology
krong@gatech.edu

Alexey Tumanov
Georgia Institute of Technology
atumanov@gatech.edu

ABSTRACT

The likelihood of encountering in-training failures rises substantially with larger Deep Learning (DL) training workloads, leading to lost work and resource wastage. Such failures are typically offset by checkpointing, which comes at the cost of storage and network bandwidth overhead. State-of-the-art approaches involve lossy model compression mechanisms, which induce a tradeoff between the resulting model quality and compression ratio. We make a key enabling observation that the sensitivity of model weights to compression varies during training, and different weights benefit from different quantization levels, ranging from retaining full precision to pruning. We propose (1) a non-uniform quantization scheme that leverages this variation, (2) an efficient search mechanism that dynamically finds the best quantization configurations, and (3) a quantization-aware delta compression mechanism that rearranges weights to minimize checkpoint differences and thereby improving compression.

We instantiate these contributions in Inshrinkerator, an in-training checkpoint compression system for DL workloads. Our experiments show that Inshrinkerator consistently achieves a better tradeoff between accuracy and compression ratio compared to prior works, enabling a compression ratio up to 39x and withstanding up to 10 restores with negligible accuracy impact in fault-tolerant training. Inshrinkerator achieves at least an order of magnitude reduction in checkpoint size for failure recovery and transfer learning without any loss of accuracy.

CCS CONCEPTS

• **Computing methodologies** → **Neural networks; Distributed computing methodologies.**

KEYWORDS

Deep Learning Training Systems, Checkpoint Compression, Quantization, Delta Compression

ACM Reference Format:

Amey Agrawal, Sameer Reddy, Satwik Bhattamishra, Venkata Prabhakara Sarath Nookala, Vidushi Vashishth, Kexin Rong, and Alexey Tumanov. 2024. Inshrinkerator: Compressing Deep Learning Training Checkpoints via Dynamic Quantization. In *ACM Symposium on Cloud Computing (SoCC '24)*, November 20–22, 2024, Redmond, WA, USA. ACM, New York, NY, USA, 20 pages. <https://doi.org/10.1145/3698038.3698553>

1 INTRODUCTION

Large-scale Deep Learning (DL) training workloads, increasingly require weeks or months of compute time on GPU clusters [64]. As the duration and scale of these training efforts grows, so does the frequency of failures. For example, the DL training workloads in a Microsoft cluster encounter a

*Work done while at Georgia Institute of Technology.



This work is licensed under a Creative Commons Attribution International 4.0 License.

SoCC '24, November 20–22, 2024, Redmond, WA, USA

© 2024 Copyright held by the owner/author(s).

ACM ISBN 979-8-4007-1286-9/24/11.

<https://doi.org/10.1145/3698038.3698553>

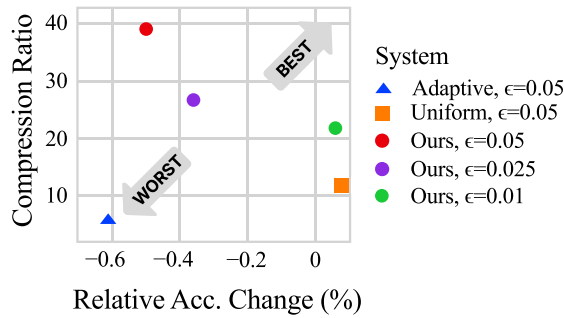


Figure 1: Inshrinkerator provides better accuracy-storage tradeoffs compared to baselines for ResNet152 training.

failure every 45 minutes on average (excluding early failures) due to various system and user errors [40]. This makes the role of checkpointing—creating periodic snapshots of a DL model during its training—increasingly crucial. When failure occurs, training can be recovered from the most recent checkpoint to reduce lost work.

During the model development phase, developers often resume from a stable mid-training checkpoint instead of starting from scratch for each bug fix or modification [1]. Checkpoints are also used in transfer learning, where a saved model state can be adapted for a different task. Due to these varied uses, there is a growing research interest in checkpointing systems [11, 23, 37, 41].

Frequent checkpoints ensure minimal loss of progress in the event of failures but also increase the storage cost. For instance, checkpoints for language models like Pythia [9], OPT [64] require tens of gigabytes of storage for each checkpoint. During the course of training, hundreds of such checkpoints are generated. As models grow in complexity and size, and as organizations seek to maintain a multitude of checkpoints, compressing these checkpoints becomes necessary.

As a result, several model compression systems have been proposed in the recent past [11, 23, 37, 41], however, existing checkpoint compression systems face shortcomings in terms of suboptimal trade-off between compression ratio and accuracy degradation as shown in Figure 1. We find that this issue originates due to the following factors. First, model parameters exhibit different levels of sensitivity to compression. Systems such as Check-N-Run [23] and QD-Compressor [41] adopt uniform quantization strategies that provide all parameters with the same level of resolution – oblivious to the level of their impact on the model quality. Second, it’s observed that models, as they accrue knowledge over the course of training, exhibit higher quantization error. However, existing systems adhere to fixed quantization configurations, e.g., a preset number of quantization bins, throughout training.

To overcome these limitations, we introduce Inshrinkerator—an efficient, transparent in-training model checkpoint compression system for DL workloads (Figure 2). Inshrinkerator is premised on the observation that not all model parameters contribute equally to model quality, and this contribution varies as the model trains. This observation informs three main aspects of Inshrinkerator.

First, Inshrinkerator uses a comprehensive quantization configuration space that partitions model parameters into three groups based on their significance: parameters preserved with full precision (unchanged), pruned parameters (set to zero), and quantized parameters. For the quantized parameters, Inshrinkerator proposes a novel *non-uniform* quantization approach using a sketch-based approximate K-Means clustering, achieving up to 65× speed up compared to a naïve implementation. This enables more effective compression by reducing the number of required quantization bins by 2-3× compared to uniform quantization.

Second, Inshrinkerator includes an efficient dynamic quantization configuration search component. This component automatically adapts the quantization configuration as parameters’ sensitivity towards compression changes over training time. This has the combined benefit of improving compression performance and alleviating the cognitive burden on ML practitioners to manually choose these configurations.

Finally, Inshrinkerator implements a novel quantization-aware delta encoding scheme, based on the observation that most parameters remain in the same quantization bin across consecutive checkpoints. Our delta encoding scheme rearranges the model parameters such that parameters for each quantization bucket are stored separately, improving delta encoding efficiency with compression ratios 3-4× higher than existing delta compression schemes.

Inshrinkerator builds on the following key contributions:

- A novel *non-uniform* quantization algorithm using approximate K-Means clustering.
- A mechanism for efficient, automatic quantization configuration search space navigation during training.
- A delta compression algorithm with effective run length encoding enabled by quantization-aware model parameter rearrangement.

We evaluate Inshrinkerator on 7 different model families, including tasks in vision and language modeling. Inshrinkerator reduces checkpoint size by 26-39× for fault-tolerant training and sustains up to ten failures (for multi-day training jobs) with negligible impact on the final accuracy of the trained model, achieving a 1.3-3.3× improvement over state-of-the-art methods. Inshrinkerator can also reduce the storage overhead of snapshots of pre-trained models used for transfer learning by 10× with no impact on the performance of the fine-tuned model. Overall, Inshrinkerator can achieve at least

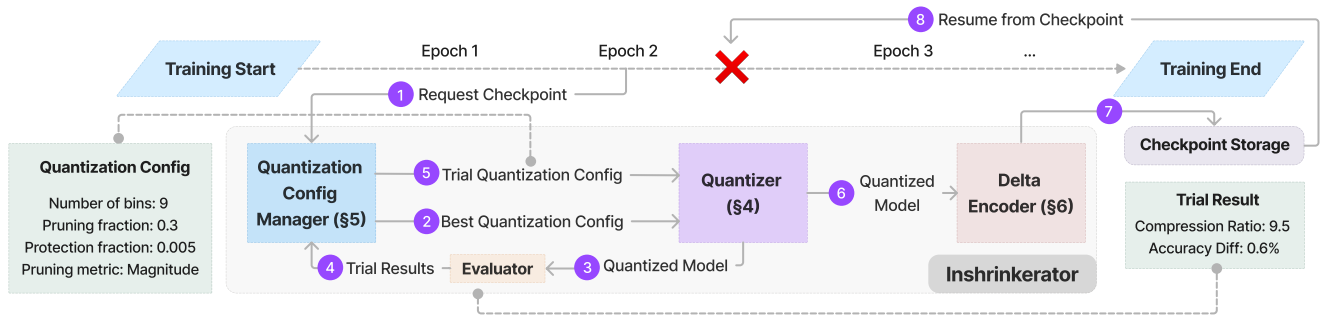


Figure 2: High-level workflow of a checkpoint compression and saving request, and resumption of training on failure from a saved compressed checkpoint in Inshrinkerator.

an order of magnitude checkpoint overhead reduction on both use cases with minimal accuracy loss.

2 BACKGROUND AND USE CASES

We start by introducing the notion of checkpointing in deep learning systems and highlight a variety of scenarios where model checkpoint compression is beneficial. Model checkpoints serve as snapshots of a deep learning model at a certain stage in its training process. They encompass the model’s architecture, learned parameters, and optimizer state. One can think of checkpoints as analogous to code commits in version control systems, where each checkpoint represents a version of the model at a specific stage of its development. In the rest of this section, we describe several use cases for model checkpoint compression.

Failure Recovery In Training Workloads. Large-scale deep-learning training workloads are highly susceptible to various hardware and software failures. Any failure during this period prompts recovery from the most recent checkpoint. Joen et al. [40] show that on average DL training workloads in a Microsoft cluster encounter a failure every 45 minutes (excluding early failures) due to various system and user errors. Zhang et al. [64] also observe that the training of the OPT-2 model required 100+ restarts due to failures across the span of two months. As we move toward developing not only larger but also more complex models and distributed training systems, the likelihood of experiencing system failures is expected to rise. Simultaneously, constraints on bandwidth and storage capacity, especially in shared multi-tenant environments, limit the frequency of checkpoints. This creates tension between frequent checkpointing to minimize wasted work and managing the storage and network checkpointing overheads. Addressing this challenge calls for a compression mechanism that enables more frequent checkpoints with minimal overhead.

Iterative Model Development. In the development of deep learning models, checkpoints are useful for multiple use cases including version control and continual learning. DL model development is an iterative process with bug fixes, hyperparameter tuning, and architectural adjustments. Given the substantial GPU hours invested in training, it is wasteful to discard progress just to make these adjustments. As a result, practitioners make these changes mid-way through training and resume the training from a stable checkpoint [1]. This prompts the need for an effective git-like version control system for models, to manage their development cycle. Beyond the training phase, production environments often require continuous updates to models to accommodate new data and maintain performance. In such scenarios, model checkpoints are preserved for extended periods for reasons such as debugging, reliability, and provenance. Delta-encoded checkpoints, where only the changes between successive checkpoints are stored, can help minimize storage requirements and streamline recovery processes in these use cases.

Model Hubs & Transfer Learning. Model hubs, such as Hugging Face [3], provide pre-trained models for a variety of tasks. These pre-trained models can be used for transfer learning, a process in which a pre-trained model is fine-tuned for a new task or domain. With the growth of the number of models on such hubs, the bandwidth required to transfer model checkpoints becomes a significant concern. For example, the BERT Base [19] model alone was downloaded more than 41 million times from the HuggingFace model hub [2] in a 30-day window as of August 2023. Model checkpoint compression can reduce the bandwidth required for transferring models between users and model hubs.

3 OVERVIEW

This section outlines the design goals and key components of Inshrinkerator.

	CNR [23]	QD [41]	LC [11]	DDNN [37]	GOBO [63]	Inshrinkerator
In-training	✓	✓	✓	✓	✗	✓
Low Overhead	✓	✓	✓	✓	✗	✓
Model Agnostic	✗	✓	✓	✓	✓	✓
Scalable	✓	✓	✓	✓	✗	✓
Algorithmically Transparent	✓	✗	✓	✓	✓	✓
Dynamic Configuration Management	✗	✗	✗	✓	✗	✓
Non-uniform Quantization	✗	✗	✓	✗	✓	✓
Quantization-aware Delta Encoding	✗	✓	✓	✓	✗	✓

Table 1: Comparison of checkpoint compression systems.

Design Goals. Inshrinkerator aims to preserve model accuracy under multiple restores from compressed checkpoints while minimizing storage overhead. Additional design goals include:

- **Low Overhead:** Compression should impose minimal overhead ($< 0.5\%$) on training runtime.
- **Model Agnostic:** The quantization algorithm should handle diverse model architectures without requiring model-specific adjustments.
- **Scalable:** Ability to process models with hundreds of millions to billions of parameters efficiently.
- **Algorithmically Transparent:** No interference with the original training algorithm, unlike quantization-aware training methods.

Table 1 compares the main features of Inshrinkerator with existing checkpoint compression systems. A detailed discussion is provided in § 11.

System Overview. Inshrinkerator processes checkpoint requests in three stages, as illustrated in Figure 2. First, the **QUANTIZATION CONFIG MANAGER** performs a distributed search to identify an optimal quantization configuration. Next, the **QUANTIZER** compresses the model using the identified configuration. Finally, the **DELTA ENCODER** processes the compressed checkpoint in CPU memory.

QUANTIZATION CONFIG MANAGER (§ 5): The quantization configuration (Figure 2, Table 2) summarizes all quantization related parameters. Over the course of training, Inshrinkerator dynamically updates the quantization configuration throughout training to maximize the compression

while satisfying to user-defined quality threshold (ϵ) for maximum allowed quantization error per checkpoint. An efficient search algorithm identifies a configuration maximizing compression ratio while adhering to ϵ .

QUANTIZER (§ 4): Given a specific quantization configuration, Inshrinkerator divides all the model parameters into three categories accordingly: prune (set to zero if below the pruning threshold), protect (preserve in full precision if above the protection threshold), or quantize (apply non-uniform quantization). Inshrinkerator uses a novel efficient non-uniform quantization scheme using sensitivity-aware approximate K-Means clustering, where quantization granularity is determined by the number of bins used in K-Means (e.g., 8-bins=3-bits). This approach offers superior quantized models with reduced runtime overhead compared to traditional methods and allows for any number of quantization bins.

DELTA ENCODER (§ 6): Inshrinkerator performs lossless compression using quantization-aware delta compression, exploiting similarities between consecutive checkpoints. **DELTA ENCODER** employs a novel parameter rearrangement technique enabling efficient run-length encoding, reducing storage overhead by up to two orders of magnitude. After delta-compression and run-length encoding, model parameters are encoded in a combined byte stream format.

4 QUANTIZER

This section outlines the design and implementation of the **QUANTIZER**. An overview of **QUANTIZER**'s workflow is shown in Figure 3. First, the **QUANTIZER** ranks model parameters according to their magnitude and sensitivity (§ 4.1). Based on the ranking, the parameters are then partitioned into three groups (§ 4.2). The least important parameters are pruned, the most important ones are protected with high precision, while the remaining parameters are quantized using our novel non-uniform quantization approach (§ 4.4).

4.1 Metrics for Parameter Ranking

First, we need to rank parameters by importance in order to assign them for protection, pruning, or quantization.

The two popular metrics used to determine the importance of parameters are magnitude and sensitivity. Magnitude-based ranking approaches simply use the parameter magnitude as an importance score [33]. The magnitude score $I_m(w)$ for parameter w is defined as $I_m(w) = |w|$. Sensitivity-based ranking approaches attempt to estimate the impact on the end-to-end loss value when a parameter is altered [21, 44]. In this paper, we use the first-order Taylor series approximation formulation of sensitivity for the sake of computational efficiency. The sensitivity score is defined as $I_s(w) = |\nabla \mathcal{L}(w)w|$,

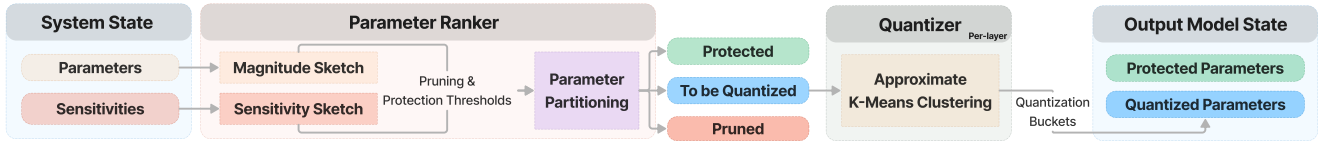
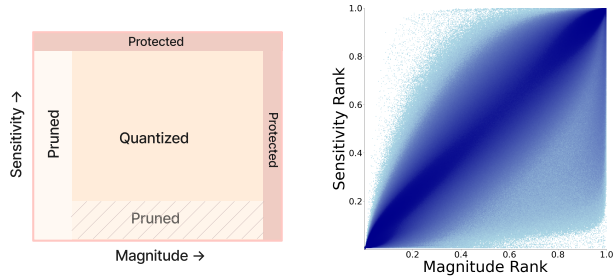


Figure 3: A high-level illustration of the quantization process in Inshrinkerator.



(a) Parameter partitioning schema in Inshrinkerator. (b) Parameter distribution for ResNet152.

Figure 4: Model parameters are partitioned into three groups for pruning, protection, and quantization by jointly evaluating magnitude and sensitivity scores.

where \mathcal{L} represents the loss function and $\nabla \mathcal{L}(w)$ is the gradient of weight w with respect to the loss function. It has been shown that magnitude-based parameter ranking approaches perform poorly in the early stage of training [56]. On the other hand, as we approach convergence, gradients start to diminish, making the sensitivity metric less reliable.

To address the above limitations, we propose to use a combination of sensitivity and magnitude to get a reliable ranking throughout the training process. Figure 4b shows the distribution of parameters across the two ranking dimensions for a ResNet152 model trained on Imagenet data. While there is a strong correlation between magnitude and sensitivity, we observe a significant number of outliers, which have a high magnitude score but a low sensitivity score or vice versa. By ranking parameters across both the metrics, Inshrinkerator can identify and preserve these parameters.

For efficient sensitivity computation, we reuse gradients computed during the training process for sensitivity estimates and asynchronously copy the gradients to pinned CPU memory to avoid GPU memory overhead. We discuss implementation details in Section 4.5.

4.2 Parameter Partitioning

Given the magnitude and sensitivity score for each model parameter, the QUANTIZER partitions the parameters into one of the three groups (Figure 4a).

Protection. Inshrinkerator preserves the most important model parameters in the high precision `bf16` format.

Recent quantization studies [15, 18] show a significant reduction in quantization error when a few important weights are protected. Our experiments suggest that protecting a small fraction (0.05-0.1%) of both the highest magnitude and sensitivity parameters provides the best performance. The precise protection fraction is determined dynamically by the QUANTIZATION CONFIG MANAGER. Based on the fraction (F_{prot}) of parameters to be protected, we identify the minimum magnitude and sensitivity thresholds (quantiles) for a parameter to qualify for protection.

Pruning. We prune the least significant model parameters, based on the pruning fraction F_{prun} (typically between 10-40%) and the pruning metric, magnitude or sensitivity. An important design consideration is the granularity level at which the pruning is performed. Layer-level pruning can degrade model quality due to different redundancy levels across the network layers. Conversely, global pruning may disproportionately affect different layer types (e.g., Convolution, Attention, Linear) due to varied weight distributions. To tackle this, we perform *per-layer type* pruning, i.e. using the same pruning fraction across all layers of a given type. We observed that in GPT-2 Medium, we can eliminate up to 30% parameters with a minimal effect on the model’s quality.

4.3 Efficient Parameter Partitioning via Quantile Sketches

Computing quantiles required for determining pruning and protection thresholds can be costly, especially for large models due to the log-linear time complexity of the operation. Moreover, in cases where models are trained via model parallelism, and parameters are distributed across different worker groups, computing global quantiles can be quite challenging.

To facilitate efficient partitioning of model parameters, we use a sketch-based quantile estimation algorithm that is a simplified and parallelized version of DDSketch [48]. This enables us to compute approximate quantile estimates in linear time with a configurable α -relative-error bound, that is, the sketch produces an estimated quantile value corresponding to the quantile value x_q such that $|\hat{x}_q - x_q| \leq \alpha x_q$. For computing protection and pruning thresholds, that require tail quantiles, DDSketch’s relative error bounds outperform alternative quantile sketches that guarantee rank error. Furthermore, the quantile sketches are mergeable, which

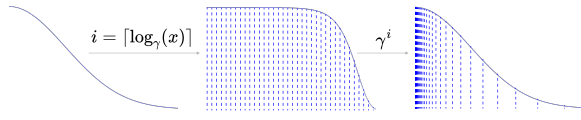


Figure 5: DDSketch’s log space transformation creates a histogram with value-dependent bucket width.

simplifies the computation of global quantile estimates from individual model shards.

The sketching algorithm transforms the input values into a logarithmic space and then divides them into an equal-width histogram (Figure 5). Due to the use of a logarithmic scale, we achieve varying levels of detail for different ranges of the original values. Specifically, smaller original values will land in narrower bins, while larger values will fall into broader bins. This allows the sketch to provide the α -relative error property of the sketch. To produce quantile estimates, the sketch sums up the buckets until it finds the bucket containing the desired quantile value.

Algorithm 1 describes the parallelizable version of DDSketch [48] optimized for GPU execution. We forgo the bucket merging step in the original DDSketch for ease of parallelization. This has limited impact in practice, as we observe that the number of buckets does not exceed a few thousand even for models with hundreds of millions of parameters with a strict relative error bound of 1%. We implement a highly parallelized version of this algorithm using Pytorch that can run on GPUs. We observe a 3-4 \times speed-up compared to off-the-shelf GPU-enabled implementation provided by CuPy [51] while using significantly less memory.

Algorithm 1 GetQuantileSketchHistogram

Require: X {Input list}

Require: α {Relative Error Bound}

- 1: $\gamma \leftarrow (1 + \alpha)/(1 - \alpha)$
 - 2: Define X_q as an empty array
 - 3: **for all** x in X **do**
 - 4: Append $\lceil \log_\gamma(x) \rceil$ to X_q
 - 5: **end for**
 - 6: (bucket_vals, bucket_counts) \leftarrow unique(X_q)
 - 7: **return** bucket_vals, bucket_counts
-

4.4 Non-uniform Quantization via Approximate K-Means Clustering

The remaining model parameters that are not pruned or protected go through quantization. Quantization reduces the precision of model parameters to lower-bit representations

(e.g., from 32-bit floating-point to 8-bit integers), leveraging redundancy and noise tolerance in large networks. Uniform quantization strategies create equally spaced quantization levels across the parameter space, while non-uniform quantization approaches use varying intervals between levels to adapt to the parameter distribution.

Non-uniform quantization techniques, such as via K-Means clustering [23, 32, 63], generally provide better compression and model quality. However, they have seen limited adoption in checkpoint compression systems due to their computational complexity. For instance, GOBO [63], that utilizes a variant of K-Means clustering, takes roughly 50 minutes to quantize a small model like ResNet-18 (11.7M parameters).

Algorithm 2 Approximate K-Means Clustering

Require: X {Input data points}, k {Number of clusters}

Require: α {Relative error}, σ {Linear combination factor}

- 1: $X_q, C_q \leftarrow$ GetQuantileSketchHistogram(X, α) {Obtain log sketch histogram Algorithm 1 for coarse grain clustering. X_q are the histogram bucket centers and C_q are the corresponding frequencies.}
 - 2: $\hat{X}_q, \hat{C}_q \leftarrow$ Normalize(X_q), Normalize(C_q)
 - 3: $W = \sigma * \hat{C}_q + (1 - \sigma) * |\hat{X}_q|$ {Compute sample weight for histogram buckets according to frequency and magnitude}
 - 4: centers \leftarrow WeightedKMeans(X_q, W, k) {Perform weighted k-means++ on histogram keys and values}
 - 5: **return** centers
-

Approximate K-Means Clustering. We introduce a novel non-uniform quantization method leveraging approximate K-means clustering, which strikes a balance between quantization quality and computational efficiency. Our key insight is to perform K-means clustering on histogram buckets of the parameters, rather than on the raw parameters themselves. This significantly reduces the computational burden of K-means clustering, resulting in up to 65 \times speedup compared to the standard implementation in CuML (§ 9.3).

Algorithm 2 describes the two-step quantization algorithm. First, we group model parameters into coarse-grain buckets, utilizing the same log-space projection approach in DDSketch. We then perform weighted k-means clustering, where the histogram bins serve as the data and their heights as sample weights, to establish quantization boundaries. Two key optimizations enhance this clustering performance:

First, we introduce an additional parameter, σ , to our weighted K-Means to help balance the allocation of resolution between high-importance parameters with low frequency and high-frequency parameters with low importance. Usually, non-uniform clustering offers increased resolution to denser parameter spaces. However, for neural network

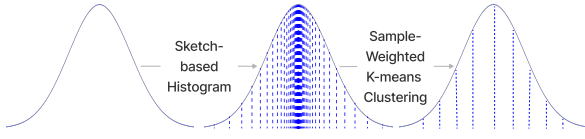


Figure 6: Two-step quantization: (1) group model parameters into coarse-grained buckets with sketch-based histogram computation. (2) cluster parameter buckets using sample weighted K-means clustering.

Algorithm 3 Weighted K-Means++ Initialization

Require: X {Input data points}
Require: W {Input sample weights}
Require: k {Number of clusters}

- 1: $C \leftarrow X[\text{random_choice}(W)]$ {Initialize the first centroid randomly according to weights}
- 2: **for** i in 2 to k **do**
- 3: $D \leftarrow \text{min_distance}(X, C)$ {Compute distance of each point to nearest centroid}
- 4: $D \leftarrow D \cdot W$ {Weight the distances}
- 5: $P \leftarrow D / \text{sum}(D)$ {Compute probabilities}
- 6: $C \leftarrow \text{CUX}[\text{random_choice}(P)]$ {Select new centroid}
- 7: **end for**
- 8: **return** C

parameters, this would allocate greater resolution to values near zero. To address this, we calculate the sample weight of a bucket, b^i , as a linear combination of the bucket’s normalized frequency and magnitude: $w^i = \sigma * b_{\text{freq}}^i + (1 - \sigma) * b_{\text{mag}}^i$ (see line 2 in Algorithm 2). Our observations suggest that values of σ in the range $0.1 - 0.4$ generally provide the best compression results. We use $\sigma = 0.2$ for all experiments.

Second, the success of the K-Means clustering algorithm relies heavily on proper initialization. Inappropriate initialization can lead to suboptimal solutions and slower convergence. K-means++ [6] addresses this by ensuring the initial centroids are uniformly distributed across the space, often leading to near-optimal solutions with a single round of clustering. We extend the K-means++ [6] algorithm for initialization with support for sample weights, Algorithm 3 provides a sketch of this process.

Error Analysis. In the first step of the quantization algorithm, we create a histogram using DDSketch’s log-space projection approach, which provides a configurable α -relative error bound. Using this approach, we formally show that the difference between optimal loss of K-Means clustering over raw inputs X and its histogram counterpart \tilde{X} is bounded

by α . We can further show that when the clustering is performed with sample-weighted k-means++, then the expected loss over \tilde{X} is bounded by the optimal loss over X and an additional term related to α^2 (Appendix A.3).

4.5 Efficient Sensitivity Computation

To save computational resources, we reuse gradients computed during the training process for sensitivity estimates instead of calculating them separately. To avoid any GPU memory overhead, we asynchronously copy the gradients to pinned CPU memory as soon as the backward pass is completed. The gradient copy operation is overlapped with the optimizer step, data loading, and forward pass of the next batch, making the overhead of this operation negligible as shown in Figure 7. The copied gradient values are accumulated on the CPU using an exponential moving average $\text{EMA}_{g_w}^t = \beta * g_w^t + (1 - \beta) * \text{EMA}_{g_w}^{t-1}$, where $\text{EMA}_{g_w}^t$ is the exponential moving average of the gradient g_w at timestep t . We set the value of β to be 0.9 in all our experiments to emphasize the recent gradients. The gradient copy mechanism is only invoked for a limited preconfigured number of batches before checkpointing is scheduled to be performed to further avoid any impact on steady-state throughput. We identify that typically recording gradients for fifty batches is sufficient to get good sensitivity estimates.

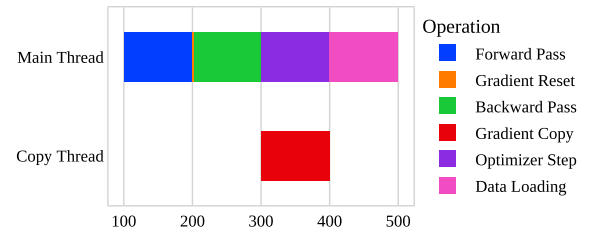


Figure 7: Schematic representation of asynchronous gradient offloading for efficient sensitivity computation in Inshrinkerator.

5 QUANTIZATION CONFIG MANAGER

In this section, we describe QUANTIZATION CONFIG MANAGER’s search space of quantization configurations and its efficient search algorithm. Unlike checkpoint compression systems that use a predefined quantization configuration throughout training, the manager dynamically searches for a configuration that maximizes compression ratio while satisfying the quality constraint ϵ at every checkpoint.

Search Space. Table 2 describes the quantization configurations that comprise the search space of QUANTIZATION CONFIG MANAGER. To keep the search space manageable, we choose to use the same number of quantization bins across

Algorithm 4 Guided Exhaustive Search

Require: Quality threshold T , Configuration cube C_{cube} , Compression Ratio CR_{max}

Ensure: Optimal configuration C_{opt}

```

1:  $C_{\text{opt}} \leftarrow \text{Null}$ 
2:  $CR_{\text{max}} \leftarrow 0$ 
3:  $\text{Diagonal} \leftarrow \text{GetDiagonal}(C_{\text{cube}})$ 
4:  $C_{\text{opt}}, CR_{\text{max}} \leftarrow \text{ParallelSearch}(\text{Diagonal}, T)$ 
5:  $\text{SubCubes} \leftarrow \text{GetFeasibleSubCubes}(C_{\text{cube}}, C_{\text{opt}})$ 
6: for each  $\text{SubCube}$  in  $\text{SubCubes}$  do
7:    $C_{\text{temp}}, CR_{\text{temp}} \leftarrow \text{GuidedExhaustiveSearch}(T, \text{SubCube})$ 
8:   if  $CR_{\text{temp}} > CR_{\text{max}}$  then
9:      $C_{\text{opt}} \leftarrow C_{\text{temp}}$ 
10:     $CR_{\text{max}} \leftarrow CR_{\text{temp}}$ 
11:   end if
12: end for
13: return  $C_{\text{opt}}$ 

```

all layers except embedding layers. We use separate quantization parameters for embedding layers, since we have observed that embedding layers are considerably more sensitive to quantization, in line with findings from other research on Transformer model quantization [57, 63].

Parameter	Values
Number of bins	4, 6, 8, 12, 16, 32
Pruning Fraction	0, 0.1, 0.2, 0.3, 0.4, 0.5
Pruning Metric	Magnitude, Sensitivity
Protection Fraction	0.0005, 0.005, 0.01

Table 2: Inshrinkerator’s quantization configuration space.

Overview of the Search Algorithm. Naïvely evaluating each configuration in the search space is costly – the search space contains over 200 configurations for non-embedding layers alone. To speed up the search, our key insight is that the optimal quantization configuration remains largely similar between adjacent checkpoints. At the first checkpoint, we perform a guided exhaustive search to identify the optimal configuration. Subsequently, we greedily evaluate the configurations that are close to the previous best configuration and only fall back to the exhaustive search if we can not identify a configuration that satisfies the quality constraint. We observe that we need to resort to exhaustive search at most 2-3 times during model training.

Guided Exhaustive Search. We use domain knowledge to reduce the cost of exhaustive search. If there was only one parameter, like the number of quantization bins, model quality would increase monotonically with the number of bins, while the compression ratio decreases. In this setting,

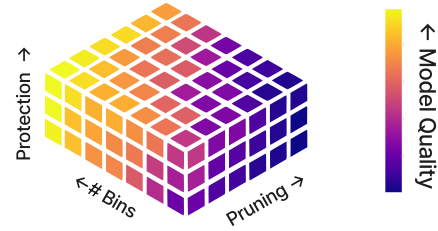


Figure 8: Configuration search space as a function of pruning fraction, protection fraction, and number of bins. The color gradient indicates the direction of higher quality.

we can find the optimal number of bins using a variation of binary search. We apply this approach to multiple configuration parameters, organizing them in a ‘configuration cube’ where each axis represents a configuration knob as shown in Figure 8. This layout ensures a monotonic increase in model quality across each axis, allowing us to use a divide-and-conquer method (Algorithm 4). We repeat this search procedure with the two pruning metrics, magnitude and sensitivity, and select the better of the two.

Delta-Neighbourhood Search. In the steady-state, we obtain a suitable configuration by greedily selecting the best configuration within the ϵ -delta neighborhood of the configuration chosen at the previous checkpoint. The ϵ -delta neighborhood is defined as all configurations that are at most ϵ steps away in the configuration cube from the given configuration. At the start of each neighborhood search, we evaluate a configuration identical to the last best configuration, except with a different pruning metric. If we observe a significant improvement in the quantized model’s quality, we change the pruning metric and proceed with the neighborhood search using the new metric. As training progresses, the model becomes denser and more sensitive to compression. Hence, subsequent neighborhood searches intentionally rule out configurations more aggressive than the prior setting.

Parallelized Implementation. Since each data parallel replica maintains an identical model state, we can parallelize the search process. Configuration search trials are evaluated in groups of m , where m is the degree of data parallelism. While scheduling the trials, we arrange them in a way that allows for an early exit. For instance, in the delta neighborhood search, we order the configurations in decreasing order of compression ratio. If in a given round of evaluation, we find a solution that satisfies the quality constraint, we can exit the search without evaluating the rest of the configurations. Our experiments show that for a data parallelism degree of eight, a high-quality configuration can be found within a single round of evaluation.

6 DELTA ENCODER

To enhance the storage efficiency of compressed models, we introduce a DELTA ENCODER that stores the differences between successive checkpoints. Delta encoding works well with quantized models, tracking changes between quantization buckets rather than exact values.

Calculating Delta. To compute the delta between the adjacent checkpoints, we transform quantized values into integers and calculate the difference between the quantized integer values of the successive checkpoints for each layer. We employ the technique proposed in QD-compressor [41] that treats the quantization range as a cyclic buffer. Concretely, for integer-mapped quantized checkpoints C_q^{i-1} and C_q^i , we compute the delta D_i as $D_i = (C_q^{i-1} - C_q^i) \bmod B$, where B is the number of quantization bins.

Due to the dynamic changes in the system configuration, the number of quantization bins may vary between checkpoints. We address this by adjusting the size of the cyclic buffer to the maximum number of bins in either checkpoint. Additionally, pruned and protected values in our quantization scheme are treated as unique quantization levels with the same delta calculation. The protected values from checkpoint C^i are stored separately to aid reconstruction.

Encoding Delta. After delta calculation, we use a combination of run-length encoding [29] and Huffman [39] encoding to reduce its storage footprint. Given that parameter changes are minimal in the later training stages, run-length encoding is particularly effective for compression. Standard run-length encoding can incur significant overhead due to the values with a run-length of one. We can reduce this overhead by only storing the run-lengths only when it is greater than one. However, this requires a mechanism to distinguish between the encoded values and their run-lengths. We achieve this by simply negating the sign of all the values, such that all the run-lengths are positive while all delta values are less than or equal to zero. Huffman encoding is then applied to the run-length encoded data for further compression.

Optimization: Parameter Rearrangement. The efficiency of run-length encoding depends on the arrangement of the parameters. We observe that the *migration rate* (i.e., the fraction of parameters that change the quantization bucket between adjacent checkpoints) of parameters can differ significantly between different quantization buckets. In the default arrangement where parameters of all quantization buckets are interleaved, the run lengths get interrupted by the buckets with high migration rates, resulting in poor compression.

To prevent this, we rearrange parameters to isolate the parameters from different quantization buckets. After calculating the delta, we split delta values corresponding to each quantization bucket in the source checkpoint into separate groups, as shown in the second step of Figure 9. Each group

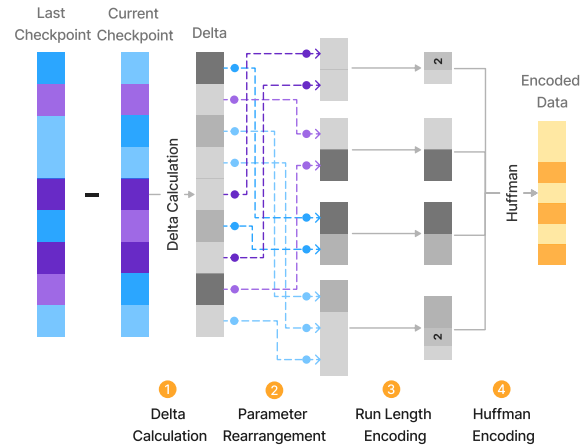


Figure 9: Rearrange and split computed delta based on the quantization bucket each parameter belonged to in the last checkpoint.

Algorithm 5 Rearranging parameters for efficient encoding

Require: Δ (Delta values calculated from adjacent checkpoints)

Require: C_q^{i-1} (The quantized integer values of the last checkpoint)

Ensure: Delta groups (A dictionary containing delta values for each unique quantization bin)

- 1: bins \leftarrow unique_bins(C_q^{i-1})
- 2: delta_groups \leftarrow create_empty_dict(bins)
- 3: **for** $j \in [0, \text{length}(C_q^{i-1})]$ **do**
- 4: delta_groups[$C_q^{i-1}[j]$].append($\Delta[j]$)
- 5: **end for**
- 6: **return** delta_groups

is then encoded separately. This ensures that quantization bins with high migration rates don't affect others. The process can easily be reversed during reconstruction and does not require any additional storage to obtain the parameter arrangement. Algorithms 5 and 6 provide a sketch of this algorithm. The approach is particularly useful when the number of quantization bins changes between checkpoints. As demonstrated in § 9.3, this optimization can significantly improve the compression ratio.

7 IMPLEMENTATION

We have developed Inshrinkerator and all the baseline systems discussed in § 8 using Pytorch, amounting to 7377 lines of Python code. To enhance the performance of delta encoding, we have implemented Huffman encoding and run-length encoding using 300 lines of C++ code. Inshrinkerator can be effortlessly integrated into user applications with less than 10

Algorithm 6 Reconstructing parameters after rearrangement

Require: `delta_groups` (Delta groups obtained from the rearranging step)

Require: C_q^{i-1} (The quantized integer values of the previous checkpoint)

Ensure: C_q^i (The quantized integer values of the current checkpoint)

```

1:  $C_q^i \leftarrow \text{create\_zero\_array}(\text{length}(C_q^{i-1}))$ 
2: for  $j \in [0, \text{length}(C_q^{i-1})]$  do
3:    $\text{bin} \leftarrow C_q^{i-1}[j]$ 
4:   if  $\text{delta\_groups}[\text{bin}] \neq \emptyset$  then
5:      $C_q^i[j] \leftarrow C_q^{i-1}[j] + \text{delta\_groups}[\text{bin}].\text{pop}(0)$ 
6:   end if
7: end for
8: return  $C_q^i$ 

```

```

from inshrinkerator import CompressorRegistry
# Before train loop
compressor = CompressorRegistry.get_compressor(
    system, model, eval_batches, search_config,
    search_metric, threshold)
...
while global_step < max_training_steps:
    ...
    # Inside train loop
    if should_checkpoint(global_step):
        compressor.compress(global_step)
    ...
    # Before flushing gradients
    compressor.before_gradient_flush()
    optimizer.zero_grad()
    ...
    # After backwards pass
    compressor.on_backward_pass_end()
    optimizer.step()
    ...

```

Listing 1: Inshrinkerator API usage with native Pytorch applications

lines of code modifications (Listing 1). Additionally, we provide a custom callback for the MosaicML composer library, enabling single-line integration of our system into applications developed with Composer (Listing 2). This strategy can be readily extended to frameworks such as Pytorch-lightning and Keras, both of which offer a high-level trainer interface.

8 EVALUATION SETUP

In this section, we describe the setup used in our evaluations. Our experimental environment comprises servers

```

from composer import Trainer
from inshrinkerator.integrations.composer import
    InshrinkeratorCallback
# Add InshrinkeratorCallback to list of callbacks
trainer = Trainer(
    ...
    callbacks=[InshrinkeratorCallback()],
)
trainer.fit()

```

Listing 2: MosaicML Composer application with Inshrinkerator callback integration.

Model	Task	Parameters
ResNet18 [34]	Vision	11.7 M
ResNet50 [34]	Vision	25.6 M
ResNet152 [34]	Vision	60.2 M
MobileNet V3L [36]	Vision	5.5 M
VGG19 [58]	Vision	143.7 M
ViT B32 [22]	Vision	88.2 M
ViT L32 [22]	Vision	306.5 M
BERT Base [19]	Language	110 M
BERT Large [19]	Language	345 M
GPT-2 Medium [52]	Language	335 M
Pythia 1B [9]	Language	1 B

Table 3: Summary of models used in the evaluation.

equipped with 8 NVIDIA A40s GPUs connected with peer-wise NVLINK, 128 AMD CPU cores, and 504 GB memory.

We compare checkpoint compression systems along the following three metrics:

Storage Efficiency. We report the end-to-end compression ratio, calculated by comparing the total size of compressed model checkpoints to their uncompressed counterparts.

Model Quality. We compare the quality of machine learning models trained from compressed checkpoints against a baseline model trained without checkpoint compression and report the *relative* quality degradation on accuracy or loss.

Runtime Overhead. We report time spent on compression as a fraction of the total training time.

8.1 Evaluation Scenarios, Models and Datasets

We evaluate Inshrinkerator across 7 different model families, including tasks in vision and language modeling. These models vary in complexity, with the number of parameters ranging from 5.5 million to 1 billion. Details on all models and datasets used are reported in Table 3 and Table 4.

Dataset	Task	Training Mode
Imagenet [17]	Image Classification	Pre-training
C4 [53]	Masked Language modeling	Pre-training
Openwebtext [28]	Next Token Prediction	Pre-training
STS-B [62]	Semantic Textual Similarity	Transfer Learning
MNLI [62]	Natural Language Inference	Transfer Learning
Alpaca [60]	Instruction Fine-Tuning	Transfer Learning

Table 4: Summary of datasets used for evaluation

We evaluate Inshrinkerator for two use cases:

Fault-tolerant Training. We simulate multiple failures during the training process and perform the recovery using compressed checkpoints. Our evaluation involves four models: ResNet-152, ViT-32L, BERT Base, and GPT2 Medium. The first two, representing the vision models, are trained on the Imagenet 1K datasets for classification tasks. The latter two, BERT and GPT models, undergo training for language modeling tasks using C4 and OpenWebText datasets, respectively. Training of GPT2 Medium and ViT-32L requires \approx 650 NVIDIA A40 GPU hours. The hyper-parameters and implementation details are provided in Appendix A.1.

Transfer Learning. We perform transfer learning tasks by training from a compressed checkpoint of pre-trained models. We use BERT Large and Pythia models. The BERT model is fine-tuned on GLUE tasks, STS-B, SST-2, and MNLI, while the Pythia model is fine-tuned using the Alpaca dataset.

8.2 Baseline Systems

We compare the performance of Inshrinkerator against two state-of-the-art checkpoint compression systems, Check-N-Run [23] and QD-compressor [41], as well as a post-training quantization system GOBO [63] that uses a similar quantization technique. Due to a lack of openly accessible functional implementations, we re-implemented their described quantization and delta compression algorithms within the Inshrinkerator framework. All re-implementations benefit from GPU-accelerated quantization algorithms. None of the systems supports dynamic configurations, so we extended Inshrinkerator’s search capability to all baselines, with each baseline having a different search space defined by its quantization algorithm. Unless otherwise mentioned, we use a default per-checkpoint quality degradation threshold of $\epsilon = 0.05$ in the configuration search. Appendix A.2 provides additional implementation details.

- **Uniform:** An entropy-based variable bit-width uniform quantization scheme and delta compression with Huffman encoding, used in QD-compressor.
- **Adaptive:** Adaptive uniform quantization [30] algorithm used in Check-N-Run. The delta compression algorithm was designed for recommendation models and does not directly extend to general DL workloads.
- **KMeans:** Naive K-Means non-uniform quantization with L1 distance used in GOBO. GOBO does not have a delta compression algorithm.

9 EVALUATION

In this section, we evaluate the empirical performance of Inshrinkerator. Experiments show that:

- Inshrinkerator achieves up to 26-39 \times compression ratio and less than 0.6% runtime overhead in fault-tolerant training, consistently offering the best accuracy-storage tradeoff compared to alternatives (§ 9.1).
- Inshrinkerator achieves better or similar performance on downstream transfer learning tasks using compressed checkpoints that are 10 \times smaller than the size of the pre-trained models, outperforming the closest state-of-the-art compression system by 2 \times (§ 9.2).
- Inshrinkerator’s non-uniform quantization algorithm and delta encoding scheme contribute meaningfully to the overall performance (§ 9.3).

9.1 End-to-end Evaluation: Fault-Tolerant Training

We conduct end-to-end training of four models under simulated failure conditions. Each experiment has 10 evenly distributed failures (every 3-8 hours of training) throughout the training duration. All models are trained using data parallelism across 8 NVIDIA A40 GPUs, with the largest models taking 3.5 calendar days (27 GPU days).

Table 5 compares different systems in terms of model quality, compression ratio, and runtime overhead on the failure recovery task. Overall, Inshrinkerator achieves between 26-39 \times reduction in storage requirements across various model families, marking a 3.5-6.5 \times improvement over Adaptive, and 1.3-3.3 \times improvement over Uniform. This is achieved while keeping the end-to-end relative degradation below 1%, and the overhead of our system remains less than 0.6% of the total training time. For ViT-L32, we observe that restoring from compressed checkpoints can act as a regularization mechanism, and result in better end-to-end performance.

While Uniform achieves a higher compression ratio compared to Adaptive due to its delta compression ability, we find the quality of models generated by the system to be inconsistent. In particular, we see a large drop in accuracy (to 4.43%) after 8 restores for BERT-Base training. KMeans is

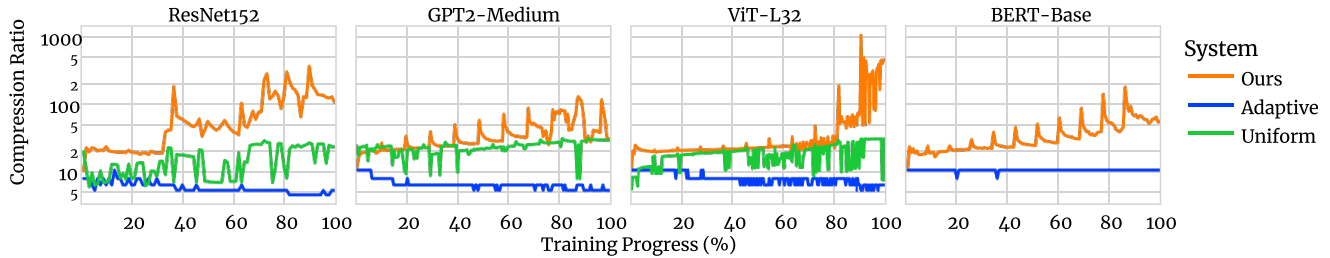


Figure 10: Compression ratio at every checkpoint during the training of various models.

	System	Deg. %	CR	Ovr. %
ResNet-152 (Base Acc: 77.41%)	Adaptive	0.36	5.71	0.61
	Uniform	-0.09	11.82	0.11
	Ours	0.50	39.09	0.35
ViT L32 (Base Acc: 71.36%)	Adaptive	-0.21	7.82	0.96
	Uniform	0.12	16.31	0.18
	Ours	-0.46	26.19	0.51
BERT Base (Base Acc: 68.49%)	Adaptive	3.97	6.78	1.03
	Uniform	×	×	×
	Ours	0.80	29.25	0.58
GPT-2 M (Base Loss: 2.72)	Adaptive	0.25	6.35	0.55
	Uniform	1.10	22.15	0.10
	Ours	0.61	29.81	0.32

Table 5: Performance of checkpoint compression systems for failure-recovery. Deg. % = Relative quality degradation, CR = Compression ratio, Overhead % = Runtime Overhead.

omitted in this experiment as it was originally designed for post-training quantization and is significantly slower than other in-training methods (details in § 9.3).

Figure 10 shows the change of the compression ratio throughout training, where the peaks in the compression ratio correspond to checkpoints directly following restores from the compressed state. Inshrinkerator’s delta encoding scheme is particularly effective during the later stages of training when model updates decelerate. Inshrinkerator’s peak compression performance is above 100× for all models.

Figure 1 illustrates the trade-off between storage overhead and model quality under different quality thresholds (ϵ) in ResNet-152 training. At $\epsilon = 0.01$, we achieve 21.82× compression, with no end-to-end quality degradation. At $\epsilon = 0.05$, there is a small end-to-end relative accuracy drop (0.5%) while the compression ratio goes up to 39.09×. For all thresholds, Inshrinkerator demonstrates an end-to-end degradation of less than 1% and achieves a better compression-quality tradeoff when compared to Adaptive and Uniform. Notably, both baselines struggle to reach high compression ratios.

Finally, Inshrinkerator can withstand up to 20 restores for ResNet152, which is equivalent to a restore every 4.5 epochs, while still maintaining a relative error of 1.1%; for reference, the error is only 0.16% for 5 restores.

Model	Task	Metric	Base	Adap	Uni	Ours
BERT-L	SST-2	CR	-	4.57	4.00	11.32
		Acc	93.2	92.9	93.0	93.3
		MNLI	Acc	85.9	86.0	86.0
Pythia 1B	Alpaca	PC	86.1	87.6	87.8	88.3
		CR	-	5.33	4.57	9.99
		CE	0.894	0.910	0.906	0.919

Table 6: Performance comparison on the transfer learning task. CR = Compression ratio, Acc = Accuracy, PC = Pearson correlation, CE = Cross-entropy Loss.

9.2 End-to-end Evaluation: Transfer Learning

For the transfer learning experiments, we perform the compression process on pre-trained model checkpoints, then deploy the compressed models for downstream task fine-tuning. The checkpoints are obtained with a quality threshold $\epsilon = 0.05$. We evaluate the BERT-Base and Pythia 1B models on a suite of downstream tasks. Table 6 shows results for all three fine-tuning tasks. We observe that for both the models, Inshrinkerator can achieve compression ratios up to 11× even without delta encoding, 2× higher compared to other systems. Notably, we find that the use of compressed checkpoints does not harm the performance of the downstream model, and we achieve performance comparable to baseline across all tasks and models.

9.3 Microbenchmarks & Ablation Studies

Quality of Non-uniform Quantization. To evaluate different quantization algorithms, we collect ten checkpoints spaced throughout the training of different models. For each

Model	Mean Number of Quantization Bins (\downarrow)											
	0.1% Degradation				1% Degradation				5% Degradation			
	CNR+	QD+	GB+	Ours	CNR+	QD+	GB+	Ours	CNR+	QD+	GB+	Ours
ResNet18	42.84	×	×	21.25	33.75	35.51	33.12	16.85	23.55	23.25	17.50	10.19
ResNet50	46.36	×	×	21.39	35.07	53.65	31.78	14.36	23.06	33.69	20.12	8.24
ResNet152	59.32	×	×	20.18	40.71	×	×	11.92	26.74	39.36	×	7.20
MobileNet V3 Large	65.24	×	×	27.92	54.44	×	×	32.48	36.68	47.96	30.79	17.91
VGG19	50.46	×	×	13.06	31.01	45.36	×	8.99	18.48	25.15	×	6.26
ViT B32	37.01	×	×	19.88	28.37	38.43	×	14.33	20.85	28.79	×	9.75
ViT L32	37.40	×	×	13.31	32.95	×	×	10.79	25.36	35.34	×	8.46
BERT Base	×	×	×	54.30	81.45	×	×	22.75	43.18	53.10	×	12.74
GPT-2 Medium	106.86	×	×	54.71	47.71	60.17	×	16.19	24.60	31.14	×	9.08

Table 7: One-shot compression performance under quality degradation constraints. We consider the mean number of quantization bins required to achieve the quality degradation constraint across ten evenly spaced checkpoints during training. System timeouts (30 mins for each evaluation of each configuration) are represented by ×, while scenarios, where no suitable configuration is found, are marked by ×.

Model	Runtime in Milliseconds (\downarrow)			
	Quantile		Clustering (K=32)	
	CuPy [51]	Ours	CuML [54]	Ours
ResNet152	71.9	15.2	10200	1160
ViT L32	364	80.9	50100	1140
BERT Base	155	33.9	20700	411
GPT-2 M	390	86.2	53400	819

Table 8: Runtime of various operations in Inshrinkerator compared with off-the-shelf implementations.

algorithm, we determine the minimum number of quantization bins needed to meet per-checkpoint thresholds of $\epsilon = 0.001, 0.01, 0.05$, then calculate the mean number of bins across all checkpoints. We disable parameter pruning in our system for this experiment to allow for easier comparison. The results are presented in Table 7.

We observe that Inshrinkerator consistently outperforms across all the models and quality thresholds. KMeans, which uses a naïve implementation of k-means clustering algorithm for quantization, takes over 30 minutes for models with more than 30 million parameters, while Inshrinkerator usually takes a few seconds. Uni form fails to find valid quantization configurations for strict quality thresholds. Among baselines, only Adaptive manages to constantly find quantization configurations that satisfy the quality thresholds. However, it requires a 2-4× higher number of quantization bins to achieve the same quality as Inshrinkerator due to their use of uniform quantization.

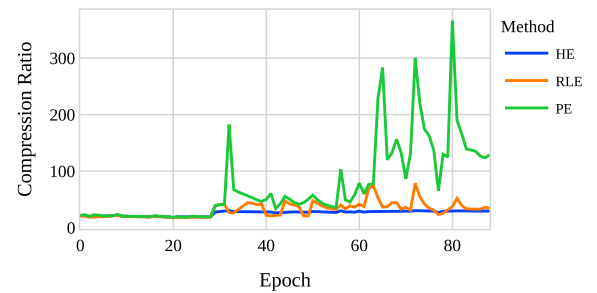


Figure 11: Parameter rearrangement significantly improves the performance of delta encoding schemes (ResNet-152). HE = Huffman Encoding, RLE = Run Length Encoding + Huffman Encoding, PE = Parameter Rearrangement-based Encoding.

Quantization Runtime. Table 8 reports the runtime of two key operations in our quantization algorithm: computing qualities for pruning and protection, and the K-Means clustering for quantization. We implement our approach in Pytorch with support for both CPU and GPU execution. We compare the performance of our implementation against off-the-shelf GPU-enabled implementations for quantile estimation in CuPy and K-Means clustering in CuML respectively. Our quantile estimation algorithm achieves 3-4× higher performance compared to CuPy, while we note a speed-up of up to 65× in quantile computation.

Parameter Rearrangement-based Delta Encoding. The key optimization in Inshrinkerator’s delta encoding scheme is parameter rearrangement, which enables it to decrease the entropy in each compressed chunk. To assess the contribution of this optimization, we compare three variants of our delta encoding scheme: the Huffman encoding variant used in Uni form (HE), run-length plus Huffman encoding (RLE),

and parameter rearrangement-based delta encoding (PE) on the ResNet-152 training job. The results are illustrated in Figure 11. With parameter rearrangement, we achieve much higher compression in the later phase of training. Notably, we obtain a maximum compression ratio of 366 with rearrangement compared to 79 with RLE and 30 with HE. This translates to overall lower storage overhead. The Huffman and run-length encoding-based schemes only achieve overall storage overhead reduction of 25× and 28× respectively, as opposed to 39× with parameter rearrangement.

10 DISCUSSION

Inshrinkerator demonstrates significant reduction in checkpoint storage overhead with minimal impact on model accuracy. However, several aspects of the system design warrant further investigation and development:

- **Delta Compression Optimization:** While delta compression optimizes storage, it can increase restore latency as the chain of delta checkpoints grows. Future work should explore adaptive policies for balancing storage efficiency and restore time, potentially incorporating periodic full checkpoints or intermediate compaction. These policies could dynamically adjust based on training progress, failure rates, and available resources.
- **Integration with High-Frequency Checkpointing:** Recent work [49] has proposed dynamic, high-frequency, iteration-level checkpointing to minimize work loss due to failures. As model sizes increase, the storage and network costs of such frequent checkpoints become significant. Future research could explore how Inshrinkerator can be integrated with these systems to alleviate storage and network bottlenecks while maintaining the benefits of frequent checkpointing. This integration may require adapting Inshrinkerator’s configuration search for short intervals between checkpoints and considering the impact on I/O bandwidth and overall system performance.
- **Adaptive Application-Specific Tuning:** The quality degradation threshold ϵ should be dynamically tuned based on application requirements and environmental factors. Future work could develop a framework for automatically determining optimal ϵ values using characteristics such as model architecture, dataset properties, and training environment. Online learning techniques could be employed to adapt ϵ during training, balancing compression efficiency with application-specific quality requirements.
- **Quantization Heuristics:** Our current approach of monotonically increasing quantization precision during training prioritizes model quality. Future work could explore more

sophisticated strategies, to optimize the quantization strategy throughout training. Multi-objective optimization approaches could be developed to balance compression ratio, model quality, and computational overhead dynamically.

These considerations open avenues for further research in adaptive checkpoint compression strategies, potentially leading to even more efficient and versatile checkpoint storage solutions for large-scale deep learning training.

11 RELATED WORK

Our work intersects with three main areas of research: model compression techniques, checkpoint compression systems, and checkpointing at scale. We discuss each of these below.

11.1 Model Compression Techniques

Model compression techniques aim to reduce the storage, memory, and computational requirements of deep learning models without compromising their performance.

Quantization. Quantization has been widely studied, particularly for efficient model inference. Post-Training Quantization (PTQ) [7, 10, 13, 26, 35] and Quantization-Aware Training (QAT) [14, 31, 38, 46, 55, 59] are two main approaches. While QAT incorporates quantization errors into the loss function during training, our proposed in-training checkpoint quantization operates on intermediate model states without affecting the training optimization function.

Variable bit-width quantization assigns different bit-widths to model layers based on their importance [12, 20, 21, 25, 42, 45]. This approach enables better compression without compromising performance, although determining the optimal bit-width for each layer can be time-consuming.

Pruning. Pruning is another popular model compression technique in which unimportant parameters are removed from the model. Magnitude-based pruning techniques remove parameters with the smallest absolute values [33], while sensitivity-based pruning methods consider the impact of parameter removal on the model’s loss [16, 24, 61, 65].

Recent work on the Lottery Ticket Hypothesis with Rewinding [27] has shown that simple magnitude-based heuristics can identify extremely sparse subnetworks early in training that yield performance similar to the full network when trained in isolation. Our work leverages this insight for in-training pruning, applying it to checkpoints during training.

11.2 Checkpoint Compression Techniques

Several studies have explored compression techniques specifically for model checkpointing systems to minimize storage overhead [11, 23, 37, 41]. These systems typically employ a combination of model quantization and delta encoding, leveraging the high similarity between adjacent checkpoints.

Quantization in Checkpoint Compression. Recent checkpoint compression systems have employed various quantization techniques. LC-Checkpoint [11] uses an exponent-based non-uniform quantization scheme, choosing the non-linear function heuristically. Delta-DNN [37] proposes a uniform quantization approach with an exhaustive configuration search strategy. QD-Compressor [41] employs an entropy-based variable bit-width uniform quantization scheme, while Check-N-Run [23] applies an adaptive uniform quantization method [30] for large deep recommendation models.

Delta Compression. Delta compression is a key technique in reducing storage overhead for model checkpoints. Each compressed checkpoint only stores the differences since the previous checkpoint, allowing for reconstruction of the final checkpoint by iteratively applying these deltas.

While QD-Compressor [41] proposes a delta encoding algorithm, it doesn't account for changes in quantization bit widths induced by their variable bit-width technique. Check-N-Run [23] introduces a delta compression scheme specifically for recommendation models, but its applicability to traditional deep learning models is limited.

11.3 Checkpointing at Scale

As machine learning models grow in size and complexity, efficient checkpointing strategies become crucial.

Distributed Checkpointing. Distributed checkpointing systems [23, 50] save partial model states across multiple nodes within a distributed system. This approach is especially useful when the model state cannot fit on a single device and reduces the overhead associated with storing and retrieving large checkpoints. Our proposed method is fully compatible with distributed checkpointing, enabling checkpoint compression even when the model state is distributed across different processes.

Frequent Checkpointing. Frequent checkpointing minimizes the risk of losing progress due to unexpected failures and allows for more granular control over the training process. CheckFreq [49] facilitates frequent checkpointing through a two-phase process that overlaps computation with checkpointing operations. It dynamically determines checkpointing intervals to balance trade-offs between overhead and wasted work.

12 CONCLUSION

Inshrinkerator combines non-uniform quantization, dynamic quantization configuration search, and quantization-aware delta compression to enable efficient in-training checkpointing. In our experiments, Inshrinkerator demonstrates the ability to induce a tradeoff space between the end-to-end model quality and checkpoint storage overhead, achieving

up to $10\times$ compression without loss in model quality for fault-tolerant training and transfer learning.

ACKNOWLEDGMENTS

This material is based upon work partially supported by the National Science Foundation under grant number CNS-2420977 and IIS-2335881. We would also like to express our sincere gratitude to the reviewers, the PC panel, and especially our shepherd Dr. Saurabh Bagchi for their insightful comments and thoughtful consideration, which significantly improved the quality of this paper.

Disclaimer: Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation.

REFERENCES

- [1] 2022. OPT 175B Training Logbook. https://github.com/facebookresearch/metaseq/blob/main/projects/OPT/chronicles/OPT175B_Logbook.pdf
- [2] 2023. bert-base-uncased · Hugging Face. <https://huggingface.co/bert-base-uncased>
- [3] 2023. Hugging Face. <https://huggingface.co>
- [4] 2023. MosaicML BERT. <https://www.mosaicml.com/blog/mosaicbert>
- [5] Margareta Ackerman and Shai Ben-David. 2009. Clusterability: A theoretical study. In *Artificial intelligence and statistics*. PMLR, 1–8.
- [6] David Arthur and Sergei Vassilvitskii. 2007. K-means++ the advantages of careful seeding. In *Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*. 1027–1035.
- [7] Ron Banner, Yury Nahshan, and Daniel Soudry. 2019. Post training 4-bit quantization of convolutional networks for rapid-deployment. *Advances in Neural Information Processing Systems* 32 (2019).
- [8] Shai Ben-David and Nika Haghtalab. 2014. Clustering in the presence of background noise. In *International Conference on Machine Learning*. PMLR, 280–288.
- [9] Stella Biderman, Hailey Schoelkopf, Quentin Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. 2023. Pythia: A suite for analyzing large language models across training and scaling. *arXiv preprint arXiv:2304.01373* (2023).
- [10] Yaohui Cai, Zhewei Yao, Zhen Dong, Amir Gholami, Michael W Mahoney, and Kurt Keutzer. 2020. Zeroq: A novel zero shot quantization framework. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 13169–13178.
- [11] Yu Chen, Zhenming Liu, Bin Ren, and Xin Jin. 2020. On efficient constructions of checkpoints. *arXiv preprint arXiv:2009.13003* (2020).
- [12] Yoojin Choi, Mostafa El-Khomy, and Jungwon Lee. 2016. Towards the limit of network quantization. *arXiv preprint arXiv:1612.01543* (2016).
- [13] Yoni Choukroun, Eli Kravchik, Fan Yang, and Pavel Kisilev. 2019. Low-bit quantization of neural networks for efficient inference. In *2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW)*. IEEE, 3009–3018.
- [14] Matthieu Courbariaux, Yoshua Bengio, and Jean-Pierre David. 2015. Binaryconnect: Training deep neural networks with binary weights during propagations. *Advances in neural information processing systems* 28 (2015).
- [15] Saurabh Dash and Saibal Mukhopadhyay. 2020. Hessian-driven unequal protection of dnn parameters for robust inference. In *Proceedings*

- of the 39th International Conference on Computer-Aided Design. 1–9.
- [16] Pau De Jorge, Amartya Sanyal, Harkirat S Behl, Philip HS Torr, Gregory Rogez, and Puneet K Dokania. 2020. Progressive skeletonization: Trimming more fat from a network at initialization. *arXiv preprint arXiv:2006.09081* (2020).
- [17] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 248–255.
- [18] Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. 2022. Llm. int8 (0): 8-bit matrix multiplication for transformers at scale. *arXiv preprint arXiv:2208.07339* (2022).
- [19] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [20] Zhen Dong, Zhewei Yao, Daiyaan Arfeen, Amir Gholami, Michael W Mahoney, and Kurt Keutzer. 2020. HAWQ-V2: Hessian Aware trace-Weighted Quantization of Neural Networks. In *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (Eds.), Vol. 33. Curran Associates, Inc., 18518–18529.
- [21] Zhen Dong, Zhewei Yao, Amir Gholami, Michael W Mahoney, and Kurt Keutzer. 2019. Hawq: Hessian aware quantization of neural networks with mixed-precision. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 293–302.
- [22] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020).
- [23] Assaf Eisenman, Kiran Kumar Matam, Steven Ingram, Dheevatsa Mudigere, Raghuraman Krishnamoorthi, Krishnakumar Nair, Misha Smelyanskiy, and Murali Annavaram. 2022. Check-N-Run: a checkpointing system for training deep learning recommendation models. In *19th USENIX Symposium on Networked Systems Design and Implementation (NSDI 22)*. 929–943.
- [24] A.P. Engelbrecht. 2001. A new pruning heuristic based on variance analysis of sensitivity information. *IEEE Transactions on Neural Networks* 12, 6 (2001), 1386–1399. <https://doi.org/10.1109/72.963775>
- [25] Fartash Faghri, Iman Tabrizian, Ilia Markov, Dan Alistarh, Daniel M Roy, and Ali Ramezani-Kebrya. 2020. Adaptive gradient quantization for data-parallel sgd. *Advances in neural information processing systems* 33 (2020), 3174–3185.
- [26] Jun Fang, Ali Shafiee, Hamzah Abdel-Aziz, David Thorsley, Georgios Georgiadis, and Joseph H Hassoun. 2020. Post-training piecewise linear quantization for deep neural networks. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*. Springer, 69–86.
- [27] Jonathan Frankle and Michael Carbin. 2018. The lottery ticket hypothesis: Finding sparse, trainable neural networks. *arXiv preprint arXiv:1803.03635* (2018).
- [28] Aaron Gokaslan and Vanya Cohen. 2019. OpenWebText Corpus. <http://Skylion007.github.io/OpenWebTextCorpus>.
- [29] Solomon Golomb. 1966. Run-length encodings (corresp.). *IEEE transactions on information theory* 12, 3 (1966), 399–401.
- [30] Hui Guan, Andrey Malevich, Jiyan Yang, Jongsoo Park, and Hector Yuen. 2019. Post-training 4-bit quantization on embedding tables. *arXiv preprint arXiv:1911.02079* (2019).
- [31] Philipp Gysel, Mohammad Motamedi, and Soheil Ghiasi. 2016. Hardware-oriented approximation of convolutional neural networks. *arXiv preprint arXiv:1604.03168* (2016).
- [32] Song Han, Huizi Mao, and William J Dally. 2015. Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding. *arXiv preprint arXiv:1510.00149* (2015).
- [33] Song Han, Jeff Pool, John Tran, and William Dally. 2015. Learning both weights and connections for efficient neural network. *Advances in neural information processing systems* 28 (2015).
- [34] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015. Deep residual learning for image recognition. *arXiv 2015. arXiv preprint arXiv:1512.03385* 14 (2015).
- [35] Xiangyu He and Jian Cheng. 2018. Learning Compression from Limited Unlabeled Data. In *Proceedings of the European Conference on Computer Vision (ECCV)*.
- [36] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, et al. 2019. Searching for mobilenetv3. In *Proceedings of the IEEE/CVF international conference on computer vision*. 1314–1324.
- [37] Zhenbo Hu, Xiangyu Zou, Wen Xia, Sian Jin, Dingwen Tao, Yang Liu, Weizhe Zhang, and Zheng Zhang. 2020. Delta-DNN: Efficiently compressing deep neural networks via exploiting floats similarity. In *Proceedings of the 49th International Conference on Parallel Processing*. 1–12.
- [38] Itay Hubara, Matthieu Courbariaux, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. 2016. Binarized Neural Networks. In *Advances in Neural Information Processing Systems*, D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett (Eds.), Vol. 29. Curran Associates, Inc. https://proceedings.neurips.cc/paper_files/paper/2016/file/d8330f857a17c53d217014ee776bfd50-Paper.pdf
- [39] David A Huffman. 1952. A method for the construction of minimum-redundancy codes. *Proceedings of the IRE* 40, 9 (1952), 1098–1101.
- [40] Myeongjae Jeon, Shivaram Venkataraman, Amar Phanishayee, Junjie Qian, Wencong Xiao, and Fan Yang. 2019. Analysis of Large-Scale Multi-Tenant GPU Clusters for DNN Training Workloads.. In *USENIX Annual Technical Conference*. 947–960.
- [41] Haoyu Jin, Donglei Wu, Shuyu Zhang, Xiangyu Zou, Sian Jin, Dingwen Tao, Qing Liao, and Wen Xia. 2023. Design of a Quantization-Based DNN Delta Compression Framework for Model Snapshots and Federated Learning. *IEEE Transactions on Parallel and Distributed Systems* 34, 3 (2023), 923–937. <https://doi.org/10.1109/TPDS.2022.3230840>
- [42] Prad Kadambi, Karthikeyan Natesan Ramamurthy, and Visar Berisha. 2020. Comparing fisher information regularization with distillation for dnn quantization. In *NeurIPS 2020 Workshop: Deep Learning through Information Geometry*.
- [43] Jnangpt Karpathy. [n. d.]. Karpathy/nanogpt: The simplest, fastest repository for training/finetuning medium-sized gpts. <https://github.com/karpathy/nanoGPT>
- [44] Yann LeCun, John Denker, and Sara Solla. 1989. Optimal brain damage. *Advances in neural information processing systems* 2 (1989).
- [45] Guiying Li, Chao Qian, Chunhui Jiang, Xiaofen Lu, and Ke Tang. 2018. Optimization based layer-wise magnitude-based pruning for DNN compression.. In *IJCAI*. 2383–2389.
- [46] Zhouhan Lin, Matthieu Courbariaux, Roland Memisevic, and Yoshua Bengio. 2015. Neural networks with few multiplications. *arXiv preprint arXiv:1510.03009* (2015).
- [47] TorchVision maintainers and contributors. 2016. *TorchVision: PyTorch’s Computer Vision library*.
- [48] Charles Masson, Jee E Rim, and Homin K Lee. 2019. DDSketch: A fast and fully-mergeable quantile sketch with relative-error guarantees. *arXiv preprint arXiv:1908.10693* (2019).
- [49] Jayashree Mohan, Amar Phanishayee, and Vijay Chidambaram. 2021. CheckFreq: Frequent, Fine-Grained DNN Checkpointing.. In *FAST*, Vol. 21. 203–216.
- [50] Bogdan Nicolae, Jiali Li, Justin M. Wozniak, George Bosilca, Matthieu Dorier, and Franck Cappello. 2020. DeepFreeze: Towards Scalable Asynchronous Checkpointing of Deep Learning Models. In *2020 20th*

- IEEE/ACM International Symposium on Cluster, Cloud and Internet Computing (CCGRID)*. 172–181. <https://doi.org/10.1109/CCGrid49817.2020.00-76>
- [51] Ryosuke Okuta, Yuya Unno, Daisuke Nishino, Shohei Hido, and Crissman Loomis. 2017. CuPy: A NumPy-Compatible Library for NVIDIA GPU Calculations. In *Proceedings of Workshop on Machine Learning Systems (LearningSys) in The Thirty-first Annual Conference on Neural Information Processing Systems (NIPS)*. http://learningsys.org/nips17/assets/papers/paper_16.pdf
- [52] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog* 1, 8 (2019), 9.
- [53] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2019. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *arXiv e-prints* (2019). arXiv:1910.10683
- [54] Sebastian Raschka, Joshua Patterson, and Corey Nolet. 2020. Machine Learning in Python: Main developments and technology trends in data science, machine learning, and artificial intelligence. *arXiv preprint arXiv:2002.04803* (2020).
- [55] Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. 2016. Xnor-net: Imagenet classification using binary convolutional neural networks. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV*. Springer, 525–542.
- [56] Maying Shen, Pavlo Molchanov, Hongxu Yin, and Jose M Alvarez. 2022. When to prune? a policy towards early structural pruning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12247–12256.
- [57] Sheng Shen, Zhen Dong, Jiayu Ye, Linjian Ma, Zhewei Yao, Amir Gholami, Michael W Mahoney, and Kurt Keutzer. 2020. Q-bert: Hessian based ultra low precision quantization of bert. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 8815–8821.
- [58] Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014).
- [59] Shyam A Tailor, Javier Fernandez-Marques, and Nicholas D Lane. 2020. Degree-quant: Quantization-aware training for graph neural networks. *arXiv preprint arXiv:2008.05000* (2020).
- [60] Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford Alpaca: An Instruction-following LLaMA model. https://github.com/tatsu-lab/stanford_alpaca.
- [61] Enzo Tartaglione, Andrea Bragagnolo, Francesco Odierna, Attilio Fian-drotti, and Marco Grangetto. 2022. SeReNe: Sensitivity-Based Regularization of Neurons for Structured Sparsity in Neural Networks. *IEEE Transactions on Neural Networks and Learning Systems* 33, 12 (2022), 7237–7250. <https://doi.org/10.1109/TNNLS.2021.3084527>
- [62] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. 2018. GLUE: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461* (2018).
- [63] Ali Hadi Zadeh, Isak Edo, Omar Mohamed Awad, and Andreas Moshovos. 2020. Gobo: Quantizing attention-based nlp models for low latency and energy efficient inference. In *2020 53rd Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*. IEEE, 811–824.
- [64] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068* (2022).
- [65] Cankun Zhong, Yang He, Yifei An, Wing W. Y. Ng, and Ting Wang. 2022. A Sensitivity-based Pruning Method for Convolutional Neural Networks. In *2022 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*. 1032–1038. <https://doi.org/10.1109/SMC53654.2022.9945569>

A EVALUATION SETUP

A.1 Model Implementations

For pre-training the ResNet-152 and ViT-32L models on Imagenet1k, we utilize the implementation and hyperparameters available in the Torchvision package [47], while BERT-Base and GPT-2 Medium training leverages implementations provided by MosaicML [4] and NanoGPT [?] repositories respectively. Notably, we maintain the original training duration and model hyperparameters, avoiding any potential skewing of evaluation results, with the exception of GPT-2 Medium where we truncate the training procedure to 30 billion tokens to restrict the training duration. In ViT-32L we find that the default parameters result in unstable training, so we reduce the learning from 0.003 to 0.001.

A.2 Baseline Implementation

In this section, we describe modifications made to Check-N-Run [23], QD-compressor [41], and GOBO [63] algorithms during our implementation of the baselines.

For QD-compressor, we address an oversight in the original description which restricted the applicability of the delta compression algorithm to cases where the number of bins assigned to a layer changes between successive checkpoints.

While GOBO was originally proposed for post-training quantization, we implement it within the Inshrinkerator framework to allow in-training compression. We also modify the GOBO clustering algorithm to include a soft termination condition which significantly reduces the runtime.

Since Check-N-Run was originally developed with the intent to support deep learning recommendation models, the delta compression algorithm proposed in the system cannot be used for regular deep learning workloads. Thus we only implement the quantization algorithm [30] used by Check-N-Run for comparison.

Recall that we also added Inshrinkerator’s dynamic configuration search capability for all baselines. Tables 9 to 12 list the range of values used for the configuration parameters for Inshrinkerator, QD, CNR, and GOBO respectively.

Parameter	Values
Number of bins	4, 6, 8, 12, 16, 32
Number of bins (Embedding Layers)	16, 32
Pruning Fraction	0, 0.1, 0.2, 0.3, 0.4, 0.5
Pruning Metric	Magnitude, Sensitivity
Protection Fraction	0.0005, 0.005, 0.01

Table 9: Search space for Inshrinkerator quantization configuration

Parameter	Values
Minimum number of bins	8, 16, 32, 64, 128, 256
Maximum number of bins	8, 16, 32, 64, 128, 256

Table 10: Search space for QD quantization configuration

Parameter	Values
Number of bins	8, 16, 32, 64, 128, 256
Number of step bins	10, 25, 50, 100
Range	0.1, 0.2, 0.3, 0.4, 0.5

Table 11: Search space for CNR quantization configuration

Parameter	Values
Number of bins	8, 16, 32, 64, 128, 256
Number of bins (Embedding Layers)	16, 32
Outlier threshold	-4
Max Iteration	1000

Table 12: Search space for GOBO quantization configuration

A.3 Additional Details for Approximate K-Means

A.3.1 Proof for Error-bound Guarantees. We restate the notations and propositions here for clarity. Let $X = \{x_1, \dots, x_n\}$ and $\tilde{X} = \{\tilde{x}_1, \dots, \tilde{x}_n\}$ denote a set of n points where $x_i, \tilde{x}_i \in [0, 1]$ and for all i , $|x_i - \tilde{x}_i| \leq \alpha x_i$ for some $\alpha \in (0, 1)$. Let $M = \{\mu_1, \dots, \mu_k\}$ be the cluster centers of X and the loss of M over inputs X is defined as $L(M) = \frac{1}{n} \sum_{x \in X} \min_{\mu \in M} |x - \mu|$. The cluster centers \tilde{M} and loss L_q over \tilde{X} are defined similarly. The optimal loss $L(M^*)$ is the loss of the optimal cluster centers $M^* = \arg \min_M L(M)$. The optimal loss $L_q(\tilde{M}^*)$ and cluster centers \tilde{M}^* over \tilde{X} are defined similarly. Let $\mu_M(i) = \arg \min_{\mu} |\mu - x_i|$ denote the cluster center in M that is closest to the i th input x_i ($\tilde{\mu}_{\tilde{M}}(i)$ is defined similarly for \tilde{x}_i).

Note that the loss $L(\tilde{M})$ over the original inputs X with any clustering \tilde{M} over \tilde{X} is at most α greater than the loss $L_q(\tilde{M})$,

$$L_q(\tilde{M}) = \frac{1}{n} \sum_{i=1}^n |\tilde{\mu}_{\tilde{M}}(i) - \tilde{x}_i| \leq \frac{1}{n} \sum_{i=1}^n (|\tilde{\mu}_{\tilde{M}}(i) - x_i| + \alpha x_i)$$

$$= L(\tilde{M}) + \frac{\alpha}{n} \sum_{i=1}^n x_i \leq L(\tilde{M}) + \alpha$$

where the first inequality follows from triangle inequality and the second follows from the fact that $\frac{1}{n} \sum_{i=1}^n x_i \leq 1$. Using similar arguments, one can verify that $L(\tilde{M}) \leq L_q(\tilde{M}) + \alpha$.

PROPOSITION A.1. *Given any two sets X, \tilde{X} of n points as defined above such that $|x_i - \tilde{x}_i| \leq \alpha x_i \forall i$. For any $k \geq 1$, the difference between the optimal loss of clusterings over X and \tilde{X} are bounded as,*

$$|L(M^*) - L_q(\tilde{M}^*)| \leq \alpha \quad (1)$$

PROOF. We first show that $L_q(\tilde{M}^*) \leq L(M^*) + \alpha$.

$$\begin{aligned} L_q(\tilde{M}^*) &= \frac{1}{n} \sum_{i=1}^n |\tilde{\mu}_{\tilde{M}^*}(i) - \tilde{x}_i| \leq \frac{1}{n} \sum_{i=1}^n |\mu_{M^*}(i) - \tilde{x}_i| \\ &\leq \frac{1}{n} \left(\sum_{i=1}^n |\mu_{M^*}(i) - x_i| + \alpha \sum_{i=1}^n x_i \right) \leq L(M^*) + \alpha \end{aligned} \quad (2)$$

In Equation (2), the first inequality follows from the fact that \tilde{M}^* is the optimal clustering over \tilde{X} and hence replacing the cluster centers $\tilde{\mu}_{\tilde{M}^*}(i)$ with $\mu_{M^*}(i)$ leads to a quantity which is greater than equal to $L_q(\tilde{M}^*)$. The next inequality follows from triangle inequality and uses the fact that $|x_i - \tilde{x}_i| \leq \alpha x_i$. Using similar arguments, one can verify that $L(M^*) \leq L_q(\tilde{M}^*) + \alpha$. Combining the two equations, we have that $|L(M^*) - L_q(\tilde{M}^*)| \leq \alpha$ which is the main statement of the proposition.

The proof for $L(M^*) \leq L_q(\tilde{M}^*) + \alpha$ follows similar arguments leading to the main statement of the result. \square

REMARK A.1. *The results show that the difference between optimal loss of k -means clustering over inputs \tilde{X} and its grouped counterpart X is bounded by α . Building upon this, we will show that when the clustering is performed with sample-weighted k -means++, the expected loss over \tilde{X} is bounded by the optimal loss over X and α^2 .*

To keep notations close to [6], we denote the loss over squared norm as $\phi(M) = \frac{1}{n} \sum_{x \in X} \min_{\mu \in M} |x - \mu|^2$. The loss $\phi_q(M)$ is defined similarly over \tilde{X} : $\phi_q(\tilde{M}) = \frac{1}{n} \sum_{\tilde{x} \in \tilde{X}} \min_{\tilde{\mu} \in \tilde{M}} |\tilde{x} - \tilde{\mu}|^2$.

PROPOSITION A.2. *Given any two sets X, \tilde{X} of n points as defined above such that $|x_i - \tilde{x}_i| \leq \alpha x_i \forall i$. For any $k \geq 1$, let \hat{M} be the clustering obtained by applying weighted k means++ (Algorithms 3 and 2) over the set \tilde{X} . Then,*

$$E[\phi_q(\hat{M})] \leq 16(\ln k + 2)(\phi(M^*) + \alpha^2) \quad (3)$$

PROOF. We first upper bound the optimal loss $\phi_q(\tilde{M}^*)$ with $\phi(M^*)$ and α . See that

$$\begin{aligned} \phi_q(\tilde{M}^*) &= \frac{1}{n} \sum_{i=1}^n |\tilde{\mu}_{\tilde{M}^*}(i) - \tilde{x}_i|^2 \leq \frac{1}{n} \sum_{i=1}^n |\mu_{M^*}(i) - \tilde{x}_i|^2 \\ &\leq 2 \left(\frac{1}{n} \sum_{i=1}^n |\mu_{M^*}(i) - x_i|^2 + \frac{\alpha^2}{n} \sum_{i=1}^n x_i^2 \right) \leq 2(\phi(M^*) + \alpha^2) \end{aligned} \quad (4)$$

Applying weighted-kmeans++ (Algorithms 3 and 2) over the set \tilde{X} produces \hat{M} . Using Theorem 1.1 in [6], we have that $E[\phi_q(\hat{M})] \leq 8(\ln k + 2)(\phi_q(\tilde{M}^*))$ and using Eq. 4 to replace $\phi_q(\tilde{M}^*)$, we obtain the main statement of the proposition. \square

REMARK A.2. *The result implies that Algorithms 2 and 3 is $O(\ln k)$ -competitive with $\phi(M^*) + \alpha^2$ while being orders of magnitude faster in comparison with the vanilla k -means++ which is $O(\ln k)$ -competitive with $\phi(M^*)$.*

We now analyze certain conditions where the clustering memberships of every point in the two cases (X and \tilde{X}) would be identical. Let $C = \{X_1, \dots, X_k\}$ be a clustering (or partition) of $X \subseteq [0, 1]$ with centers $M = \{\mu_1, \dots, \mu_k\}$. For all $x \in X_i$, and $i \neq j$, $|x - \mu_i| \leq |x - \mu_j|$. Let the minimum distance between any two points in two different clusters in C be the split of the clustering denoted as $split_C(X)$. Let the maximum distance between any two points in the same cluster be the width of the clustering denoted as $width_C(X)$.

DEFINITION A.1. [8] *A clustering $C = \{X_1, \dots, X_k\}$ of $X \subseteq [0, 1]$ is σ -separable for $\sigma \geq 1$ if $split_C(X) > \sigma \cdot width_C(X)$.*

In other words, any clustering C is σ -separable if the minimum distance between any two points in separate clusters is greater than the maximum distance between two points in the same cluster.

We now show that if the optimal clustering over the set X is σ -separable, then for small enough α depending on the value of σ , the optimal clustering over X and \tilde{X} will be the same.

First note that, if a σ -separable clustering C exists for a given set of inputs, then only one such partitioning of the inputs.

LEMMA A.1. [5] *If there exists a k -clustering C of X , for $k \geq 2$, such that C is σ -separable, then there is only one such partitioning of X .*

The proof of the above Lemma can be found in [5]. We now show that if $\alpha < \frac{split_C(X) - width_C(X)}{4}$, then the same partitioning over \tilde{X} is also σ -separable and hence the optimal clustering for both X and \tilde{X} are identical.

Let s_1 and s_2 be the two points in X such that the distance between them is minimum among all pairs of points in separate clusters ($|s_2 - s_1| = \text{split}_C(X)$). Let w_1 and w_2 be two points in X such that the distance between them is maximum among all pairs of points in the same cluster.

For the points in X to have a σ -separable clustering in \tilde{X} , the following condition should hold in the worst case,

$$s_2(1 - \alpha) - s_1(1 + \alpha) > w_2(1 + \alpha) - w_1(1 - \alpha)$$

$$\implies (s_2 - s_1) - \alpha(s_1 + s_2) > (w_2 - w_1) + \alpha(w_1 + w_2)$$

$$\implies \alpha < \frac{(s_2 - s_1) - (w_2 - w_1)}{4} = \frac{\text{split}_C(X) - \text{width}_C(X)}{4}$$

This implies that if the original clusters are σ -separable and α is small enough then the optimal cluster memberships over X and \tilde{X} will be identical.