

# CoDeQ: End-to-End Joint Model Compression with Dead-Zone Quantizer for High-Sparsity and Low-Precision Networks

Jonathan Wenshøj  
Department of Computer Science  
University of Copenhagen  
jowe@di.ku.dk

Bob Pepin  
Department of Computer Science  
University of Copenhagen  
bob@pepin.io

Tong Chen  
Department of Computer Science  
University of Copenhagen  
toch@di.ku.dk

Raghavendra Selvan  
Department of Computer Science  
University of Copenhagen  
raghav@di.ku.dk

## Abstract

While joint pruning–quantization is theoretically superior to sequential application, current joint methods rely on auxiliary procedures outside the training loop for finding compression parameters. This reliance adds engineering complexity and hyperparameter tuning, while also lacking a direct data-driven gradient signal, which might result in sub-optimal compression. In this paper, we introduce **CoDeQ**, a simple, fully differentiable method for joint pruning–quantization. Our approach builds on a key observation: the dead-zone of a scalar quantizer is equivalent to magnitude pruning, and can be used to induce sparsity directly within the quantization operator. Concretely, we parameterize the dead-zone width and learn it via backpropagation, alongside the quantization parameters. This design provides explicit control of sparsity, regularized by a single global hyperparameter, while decoupling sparsity selection from bit-width selection. The result is a method for Compression with Dead-zone Quantizer (CoDeQ) that supports both fixed-precision and mixed-precision quantization (controlled by an optional second hyperparameter). It simultaneously determines the sparsity pattern and quantization parameters in a single end-to-end optimization. Consequently, CoDeQ does not require any auxiliary procedures, making the method architecture-agnostic and straightforward to implement. On ImageNet with ResNet-18, CoDeQ reduces bit operations to  $\sim 5\%$  while maintaining close to full precision accuracy in both fixed and mixed-precision regimes<sup>1</sup>

<sup>1</sup>Source code available at <https://github.com/saintslab/CoDeQ>

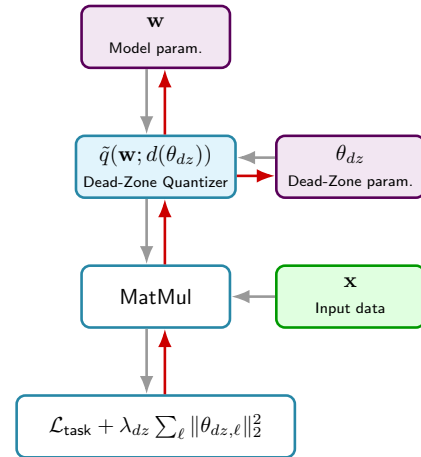


Figure 1. Schematic overview of the proposed CoDeQ method for a linear layer with a matrix multiplication (MatMul). Grey arrows show the forward pass; red arrows show gradient flow. CoDeQ parameterizes the dead-zone width  $d$ , which determines the magnitude range mapped to zero, enabling model weights, pruning– and quantization– parameters to be learned jointly via backpropagation. By regularizing the dead-zone width to be large, CoDeQ promotes sparse solutions controlled by a single global hyperparameter  $\lambda_{dz}$ .

## 1. Introduction

Modern neural networks continue to grow in scale and computational cost, creating significant barriers for deployment in latency-, memory-, and energy-constrained environments. Model compression has therefore become essential, with pruning and quantization being the most widely adopted techniques for reducing computational load and

Property	Ours	GETA	FITC	SQL	QST	CLIPQ	DJPQ	BB
Fully Learnable	✓	✗	✗	✗	✗	✗	✓	✓
Hyperparameters	2	5	3	2	3	2	5	3
Fixed-bit	✓	✓	✗	✓	✗	✗	✗	✗
Architecture Agnostic	✓	✓	✓	✗	✗	✗	✗	✗

Table 1. Comparison of the relative ease of adopting joint pruning–quantization methods. *Fully Learnable*: The pruning and quantization parameters are obtained directly through gradient-based updates in training or fine-tuning, without auxiliary search or allocation procedures. *Hyperparameters*: The number of hyperparameters that require tuning to use the method. *Fixed-bit*: The method allows for and is evaluated in a setting with a single, user-specified bit-width across layers, without MP allocations. *Architecture Agnostic*: The method is readily applicable and empirically validated on heterogeneous architectures.

model size. In practice, these techniques are typically applied sequentially (e.g., pruning a model and then quantizing it), each with its own fine-tuning phase. This pipeline increases training cost and complexity, while potentially yielding sub-optimal results because pruning and quantization decisions are made independently. These limitations have motivated extensive research on *joint* pruning–quantization methods, where sparse and quantized weights are found simultaneously [13–18, 20, 21].

Despite substantial effort, current joint methods remain underused and difficult to deploy in practice. Typical machine learning pipelines consist of an inner loop that updates model parameters with minimal to no user intervention needed and an outer loop that configures the inner loop, for example, by selecting hyperparameters. Extending the outer loop, however, tends to introduce substantial additional demands on the user. In existing literature, the joint pruning–quantization problem is framed as determining a layer-wise allocation of bit-widths and sparsity patterns. This is often handled by introducing an auxiliary procedure in the outer loop that runs either before or during training to compute the layer-wise allocations. However, this general approach introduces practical challenges for adopting joint methods and may also result in sub-optimal compression outcomes. Tab. 1 summarizes the relative ease of adopting relevant methods.

### 1.1. Challenges

**Dependence on Auxiliary Procedures.** Many joint pruning–quantization methods rely on non-differentiable heuristics, or discrete search procedures to determine layer-wise sparsity and bit-width allocations [13–15, 17, 20, 21]. These auxiliary procedures often run in multiple stages or alternate with gradient updates, increasing training cost and pipeline complexity. Additionally, the compression parameters lack a direct gradient signal from the inner loop, meaning the resulting allocations may not be optimal given the particular model–dataset pair. Furthermore, procedure designs tend to depend on architectural assumptions, limiting generality and complicating integration into existing workflows [13, 16, 18, 20].

**Hyperparameter Burden.** Joint pruning–quantization methods naturally introduce hyperparameters related to compression, such as desired global sparsity levels, bit-width or bit budgets. However, additional procedure-specific hyperparameters are also introduced, which users must configure [13–15, 17, 20, 21]. These values are rarely known *a priori*, and exhaustive sweeps over both compression and procedure hyperparameters are typically infeasible. As a result, achieving good accuracy–compression trade-offs demands considerable user effort.

**Mixed-precision Quantization.** Many reported results focus on mixed-precision (MP) quantization, often employing non-power-of-two bit-widths [13–18, 20, 21]. However, commodity hardware predominantly supports uniform 4-bit and 8-bit quantization, complicating deployment. Furthermore, many methods couple MP search and sparsity selection via a single auxiliary procedure, which means that bit-width cannot be fixed while finding only layer-wise sparsity parameters [13, 15–18, 21]. This coupling is problematic since bit-width often is hardware-dictated, and *sparsity is the principal degree of freedom* for improving efficiency.

### 1.2. Contributions

To address these challenges, we introduce **CoDeQ**, a simple and fully differentiable approach for joint pruning–quantization, as illustrated in Fig. 1. CoDeQ builds on the observation that the dead-zone (zero bin) of a scalar quantizer naturally implements magnitude pruning. By parameterizing this dead-zone width and learning it via back-propagation, CoDeQ can simultaneously find layer-wise sparsity patterns and quantized weights, at desired bit-width (or with an optional learnable MP bit-width allocations), eliminating the need for auxiliary procedures (as shown in Fig. 2). Because CoDeQ comprises a differentiable quantizer, the method is inherently architecture-agnostic and can be applied to any architecture supported by standard quantization-aware training (QAT) workflows. Furthermore, pruning sparsity is governed by a single global hyperparameter, with an optional second hyperparameter enabling MP. Finally, as CoDeQ is built on a standard uni-

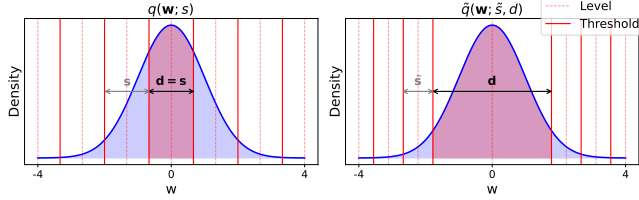


Figure 2. Visualization of ( $b$ -bit mid-tread) uniform symmetric quantizer Eq. (1) (left) and with adjustable dead zone Eq. (10) (right). Increasing the width of the dead-zone increases the sparsity by including more weights into the zero bin, while also making the step-size smaller in the non-zero region.

form symmetric quantizer, it remains straightforward to implement in existing QAT and deployment frameworks. Our key contributions are:

- **Unified pruning via quantization:** In Sec. 4, we show theoretically that the dead-zone of a scalar quantizer is equivalent to magnitude pruning. We show a formulation for a quantizer with adjustable dead-zone width, allowing for control over sparsity. Additionally, with an appropriate scale factor, no quantization levels are wasted in the pruned region, preserving the full non-zero grid for unpruned weights.
- **A simple, differentiable joint pruning–quantization method:** In Sec. 5 we propose **CoDeQ**, which learns layer-wise dead-zone widths and quantized weights during QAT, with either fixed-bit or optionally MP. Sparsity is governed by a single global hyperparameter (with an additional one for MP), eliminating the need for auxiliary procedures.
- **Strong trade-offs with minimal engineering overhead:** In Sec. 6 we empirically demonstrate that CoDeQ produces highly sparse, low-precision models with competitive accuracy in both fixed-bit and MP settings, while avoiding the extensive implementation and tuning overhead required by prior joint approaches.

## 2. Related work

**Joint Pruning–Quantization Pipelines.** Early work such as Deep Compression [7] applies pruning and quantization in sequential stages: a dense model is first pruned using hand-tuned thresholds, then quantized by selecting the bit-width per layer, followed by additional fine-tuning. In these pipelines, the compression policy is largely manual: users decide pruning ratios and bit-widths, and the training loop simply enforces those choices. Subsequent joint pruning–quantization methods retains a similar reliance on user-specified compression targets, but also introduce a non-trivial amount of outer-loop implementations and configurations. CLIP-Q [15] uses Bayesian optimization in an outer loop to select a pruning rate and bit-width for each layer, and then runs a separate joint fine-tuning phase to re-

alize these allocations. APQ [17] couples a supernet with an evolutionary search over architecture, channel counts, and per-block bit-widths under latency or energy constraints. Automatic Neural Network Compression (SQL) [20] integrates pruning and mixed-precision quantization and employs ADMM-style updates with periodic projection and discretization steps to meet a global target size. GETA [14] constructs a dependency graphs and uses a custom projected optimizer to jointly decide which channel or attention-head groups to prune and layer-wise bit-width allocations, given a global sparsity targets, bit-width ranges, and scheduler periods. FITCompress (FITC) [21] achieves combined pruning and mixed-precision via iterative search and optimization over a design space.

**Quantization as Pruning.** A different line of joint methods explicitly embed pruning into the quantizer. Quantized Sparse Training (QST) [13] defines a quantizer and pruning operator that first thresholds weights by magnitude and then applies a quantizer with step size tied to the pruning threshold. QST takes as input a global target pruning rate, a target bit-width, and a coefficient  $\alpha$  that scales the step size. Because the pruning threshold is reused as the step size, high sparsity naturally results in lower effective precision.

Few works avoid auxiliary procedures entirely, by learning all necessary compression parameters via the quantizer. Differentiable Joint Pruning and Quantization (DJ PQ) [18] jointly optimizes compression parameters, via a single bit operation (BOP)- aware objective. It uses channel gates with a regularizer to induce structured sparsity, together with a learnable non-linear quantizer, parameterized by ranges and step sizes from which an effective bit-width is derived. Sparsity and bit-width are coupled, and compression behavior is governed by two global regularization weights and several gate- and quantizer-specific hyperparameters. In reported layer-wise pruning ratios, sparsity is concentrated in early layers while other layers stay dense, possibly indicating sub-optimal pruning patterns. Bayesian Bits (BB) [16] decomposes quantization into gated power-of-two quantizers. The effective bit-width and whether a group is pruned are determined by which gates are active. Gate activations are optimized via a variational objective with sparsity-inducing priors and a global regularization parameter. At test time, gates probabilities are thresholded, followed by additional fine-tuning. Our work continue this line of using the quantizer for determining compression parameters.

## 3. Background

### 3.1. Quantization

A commonly used form of neural network weight quantization is the *mid-tread uniform symmetric quantizer*  $q(\mathbf{w}; s)$

defined as

$$\begin{aligned}\bar{\mathbf{w}} &= \text{clip}\left(\left\lfloor \frac{\mathbf{w}}{s} \right\rfloor, -Q_b, Q_b\right), \\ \hat{\mathbf{w}} &= s \cdot \bar{\mathbf{w}} = q(\mathbf{w}; s),\end{aligned}\quad (1)$$

where  $s > 0$  is the scale factor,  $Q_b = 2^{b-1} - 1$  and  $-Q_b, Q_b$  being the minimum and maximum quantization index, respectively,  $\lfloor \cdot \rfloor$  represents the element-wise round-to-nearest operator, and  $\text{clip}(\mathbf{w}, \alpha, \beta)$  limits each element of  $\mathbf{w}$  to the range  $[\alpha, \beta]$ . A common choice for  $s$  is the absmax scale factor

$$s = \frac{\max(|\mathbf{w}|)}{Q_b}. \quad (2)$$

For the quantizer in Eq. (1), the scale factor  $s$  coincides with the step size; the distance between two adjacent quantization levels.

### 3.2. Pruning

Neural network pruning can be expressed as the element-wise multiplication (shown as  $\odot$ ) of a binary mask  $\mathbf{m} \in \{0, 1\}^{|\mathbf{w}|}$  to the weight vector  $\mathbf{w}$ , producing pruned weights

$$\mathbf{m} \odot \mathbf{w}, \quad (3)$$

where each entry  $m_i \in \mathbf{m}$  indicates whether a weight is retained ( $m_i = 1$ ) or discarded ( $m_i = 0$ ). The pattern of zeros in the binary mask determines the sparsity structure: an *unstructured* mask places zeros arbitrarily throughout the weight vector, whereas a *structured* mask zeros out entire rows, columns, channels, or blocks. A common rule for unstructured pruning is *magnitude pruning*, which removes all weights whose absolute value falls below a threshold  $\tau > 0$ . In this case, the mask is

$$m_i(\tau) = 1_{\{|w_i| \geq \tau\}}, \quad (4)$$

so that  $\mathbf{m}(\tau) \odot \mathbf{w}$  with  $1_{\{\cdot\}}$  as the indicator function.

### 3.3. Quantization/Pruning Aware Training

QAT refers to methods that directly optimize model parameters under their quantized approximation [9, 10]. Concretely, the forward pass is performed using the quantized weights  $\hat{\mathbf{w}}$ , thereby simulating the behavior of low-precision inference:

$$\mathcal{L}_{\text{QAT}} = \mathcal{L}(q(\mathbf{w})), \quad (5)$$

where  $\mathcal{L}(\cdot)$  denotes the training objective and  $q(\cdot)$  an arbitrary scalar quantizer. Analogously, we define *Pruning-Aware Training (PAT)* as directly optimizing the loss under a pruning mask. In this case, the forward pass uses the pruned parameters  $\mathbf{w} \odot \mathbf{m}$ , where  $\mathbf{m}$  is a binary pruning mask:

$$\mathcal{L}_{\text{PAT}} = \mathcal{L}(\mathbf{w} \odot \mathbf{m}). \quad (6)$$

The earliest instance of PAT can be traced to [12], where a network is iteratively retrained under a fixed pruning mask that is recomputed at each iteration. Later works [5, 6] extended PAT to an end-to-end formulation, allowing the mask  $\mathbf{m}$  itself to be learnable.

### 3.4. Straight-Through Estimator

Consider the chain rule for a function  $F$  which is either non-differentiable or with a derivative of zero almost everywhere:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \frac{\partial \mathcal{L}}{\partial F} \cdot \frac{\partial F}{\partial \mathbf{w}}. \quad (7)$$

Here, the term  $\frac{\partial F}{\partial \mathbf{w}}$  prevents gradient-based optimization through  $F$ . The Straight-Through Estimator (STE) [1] replaces  $\frac{\partial F}{\partial \mathbf{w}}$  with a simple proxy  $g(x) = 1$  in the backward pass, allowing gradients to flow through  $F$ . For a uniform quantizer (1), STE is often applied to the rounding operation  $\lfloor \cdot \rfloor$ , replacing its almost-everywhere zero derivative with 1, while differentiating the remaining operations normally. Here, the clip operation contributes a gradient mask, so gradients are only propagated for inputs within the representable range. Consequently, following Eq. (7), the gradient for Eq. (1) is approximated as

$$\frac{\partial \mathcal{L}_{\text{QAT}}}{\partial \mathbf{w}} \approx \frac{\partial \mathcal{L}_{\text{QAT}}}{\partial \hat{\mathbf{w}}} \cdot \mathbf{1}_{\{|\mathbf{w}/s| \leq Q_b\}}. \quad (8)$$

Under this approximation, the quantizer behaves as an identity mapping during the backward pass within the representable range  $[-Q_b, Q_b]$ , while remaining exact in the forward pass.

## 4. Joint Pruning and Quantization with Dead-Zone Quantizer

In this section, we show how magnitude pruning is directly related to the dead-zone of a quantizer. We then provide a concrete formulation of a dead-zone quantizer with adjustable dead-zone width and uniform step size in the remaining quantization grid. Lastly we show how to set the step-size.

### 4.1. Dead-Zone Quantizer

Let  $q$  denote a quantizer mapping  $\mathbf{w}$  to its quantized counterpart  $\hat{\mathbf{w}}$ , i.e.,  $\hat{\mathbf{w}} = q(\mathbf{w})$ . The quantizer  $q$  may induce sparsity if 0 is a possible reconstruction level of  $q$ , so that entries of  $\hat{\mathbf{w}}$  can take the value  $\hat{w}_i = 0$ . Suppose the scalar interval mapped to zero is symmetric, i.e.,  $\hat{w}_i = 0$  for all  $w_i \in [-d/2, d/2]$  for some  $d > 0$ . We refer to  $d$  as the *dead-zone width* of the quantizer  $q$ . A quantizer with dead-zone width  $d$  induces the same sparsity pattern as magnitude pruning at threshold  $\tau = d/2$ . In particular,

$$\tilde{q}(\mathbf{w}; s, d) = \tilde{q}(\mathbf{w}; s, d) \odot \mathbf{m}(d/2), \quad (9)$$



since  $\tilde{q}$  maps all weights with  $|w_i| \leq d/2$  to zero, while the remaining weights are quantized with uniform step-size  $s$ . The uniform quantizer defined in Eq. (1) is a special case with dead-zone width  $d = s$  which gives  $q(\mathbf{w}; s) = \tilde{q}(\mathbf{w}; s, s) = q(\mathbf{w}; s) \odot \mathbf{m}(s/2)$ , which is the approach used in QST [13]. Here we introduce an inherent trade-off: for aggressive pruning we need a large step-size which then increases the quantization error, whereas increasing the precision reduces quantization error but wastes quantization levels. Thus, to obtain a dead-zone quantizer that remains usable across bit-widths and step sizes, we seek a formulation that decouples the dead-zone width  $d$  from the step size  $s$ , enabling  $d$  to be tuned freely. The  $b$ -bit *mid-tread uniform symmetric quantizer with adjustable dead-zone* (which we henceforth refer to as the *dead-zone quantizer*), denoted by  $\hat{\mathbf{w}} = \tilde{q}(\mathbf{w}; s, d)$ , is defined as:

$$\begin{aligned}\bar{\mathbf{w}} &= \text{clip} \left( \left\lfloor \frac{\text{sign}(\mathbf{w}) \cdot \text{ReLU}(|\mathbf{w}| - \delta)}{s} \right\rfloor, -Q_b, Q_b \right) \\ \hat{\mathbf{w}} &= \text{sign}(\bar{\mathbf{w}}) \cdot \delta + s \cdot \bar{\mathbf{w}},\end{aligned}\quad (10)$$

where  $\delta = \frac{d}{2} - \frac{s}{2}$  is the difference between the half-widths of the desired dead-zone width  $d$  and the dead-zone width  $s$  in the underlying uniform quantization grid.

This formulation enables independent control over the dead-zone width  $d$  while preserving uniformly spaced quantization levels elsewhere. Intuitively, subtracting  $\delta$  inside the ReLU shifts the effective quantization thresholds around zero. For all inputs with  $|w_i| \leq d/2$  we have

$$0 \leq \text{ReLU}(|w_i| - \delta) \leq \frac{d}{2} - \delta = \frac{s}{2}, \quad (11)$$

with ReLU ensuring negative values are set to zero. Consequently, all weights  $w_i \in [-d/2, d/2]$  are mapped to zero, so the dead-zone width is  $d$ , while the remaining quantization levels still lie on a uniform grid with step size  $s$ . If we set  $d = s$ , then  $\delta = 0$ , and we have:

$$\text{sign}(\mathbf{w}) \cdot \text{ReLU}(|\mathbf{w}| - \delta) = \mathbf{w}, \quad \text{sign}(\bar{\mathbf{w}}) \delta = 0,$$

with the  $\tilde{q}$  reducing to the uniform quantizer Eq. (1).

## 4.2. Pruning-Aware Absmax Scale Factor

With a symmetric quantizer, we have  $2Q_b + 1$  quantization levels. The non-zero levels then number  $2Q_b$ , which means we need to calculate the step size in the non-zero region for a total of  $2Q_b - 1$  steps. Additionally to calculate the step size, we note that while the total dynamic range of the weights is  $2 \max(|\mathbf{w}|)$ , only  $2 \max(|\mathbf{w}|) - d$  is available for the non-zero levels once the dead-zone of width  $d$  is carved out. Thus, the pruning-aware scale factor  $\tilde{s}$  is given by:

$$\tilde{s} = \frac{2 \max(|\mathbf{w}|) - d}{2Q_b - 1} = \frac{\max(|\mathbf{w}|) - d/2}{Q_b - 1/2}. \quad (12)$$

If we set  $d = \tilde{s}$ , we obtain

$$d = \frac{\max(|\mathbf{w}|) - d/2}{Q_b - 1/2} \implies d = \frac{\max(|\mathbf{w}|)}{Q_b} = s, \quad (13)$$

which recovers the absmax scale factor Eq. (2) for all steps. In summary, the adaptive scale factor in Eq. (12) ensures that all quantization levels fall within the weight range and also decreases the quantization error of non-pruned weights, as we increase  $d$ . We illustrate this difference in Fig. 2.

## 5. Compression with Dead-zone Quantizer (CoDeQ)

In this section we show how to achieve joint pruning-quantization by training with the dead-zone quantizer. We do this by parameterizing the width of the dead-zone and then regularize that width. Lastly we extend the method to also include learnable bit.

### 5.1. Gradient Optimization through Dead-Zone Quantizer

To enable gradient-based learning through the dead-zone quantizer Eq. (10), we follow the common practice of applying the STE to the rounding operation  $\lfloor \cdot \rfloor$ . Our quantizer also contains sign operations in  $\bar{\mathbf{w}}$  and  $\hat{\mathbf{w}}$ , but we do not apply STE to these. Allowing gradients to pass through the sign operators creates an additional gradient path and distort the gradient magnitude of  $\tilde{q}$  wrt.  $w_i$  away from 1. We additionally use ReLU to shift the input before quantization. Since ReLU has zero gradient in its negative region, no learning signal would reach weights that have been pruned by the ReLU. To preserve a signal, we apply STE to the ReLU. For similar reasons, we also STE the clipping operation so that weights in the saturated region can still move during training. This gives us the effective gradient

$$\frac{\partial \mathcal{L}(\tilde{q}(\mathbf{w}))}{\partial \mathbf{w}} \approx \frac{\partial \mathcal{L}(\tilde{q}(\mathbf{w}))}{\partial \hat{\mathbf{w}}} \cdot \mathbf{1}, \quad (14)$$

with the indicator function always being true.

### 5.2. Learnable Dead-Zone

To achieve differentiable sparsification with the dead-zone quantizer Eq. (10) we make the dead-zone width learnable, by parameterizing it as a differentiable function of a trainable parameter  $\theta_{dz}$ :

$$d(\theta_{dz}) = 2 \max(|\mathbf{w}|) \cdot (1 - \tanh(|\theta_{dz}|)). \quad (15)$$

This formulation ensures that  $d$  remains positive and smoothly bounded within the dynamic range of the weights  $[0, 2 \max(|\mathbf{w}|)]$ . The  $\tanh(\cdot)$  function provides a convenient way to map the unbounded parameter  $\theta_{dz} \in \mathbb{R}$  to a normalized scaling factor, allowing gradient-based optimization while preventing  $d$  from exceeding the dynamic range of the weights.

Intuitively, when  $|\theta_{dz}|$  is small,  $\tanh(|\theta_{dz}|) \approx 0$ , leading to a large dead-zone and thus stronger pruning. As  $|\theta_{dz}|$  increases,  $\tanh(|\theta_{dz}|) \rightarrow 1$ , shrinking the dead-zone and reducing the pruning effect. This parameterization therefore enables adaptive, data-driven control of the pruning strength during training.

### 5.3. Optimization objective for CoDeQ

With the quantizer in Eq. (10) now fully differentiable and with a learnable dead-zone, we can optimize with QAT directly. To incentivize high sparsity solutions, we add weight decay to  $\theta_{dz}$ . The overall regularized loss for updating the model parameters and the dead-zone width can then be written as

$$\mathcal{L}_{\text{CoDeQ}} = \mathcal{L}(\tilde{q}(\mathbf{w}; \tilde{s}, d(\theta_{dz}))) + \lambda_{dz} \|\theta_{dz}\|_2^2, \quad (16)$$

where  $\lambda_{dz}$  controls the regularization strength for the dead-zone parameter. Optimizing a dead-zone quantizer  $\tilde{q}(\mathbf{w}; \tilde{s}, d(\theta_{dz})) = \tilde{q}(\mathbf{w}; \tilde{s}, d(\theta_{dz})) \odot \mathbf{m}(d(\theta_{dz})/2)$  implicitly corresponds to optimizing both the quantization-aware loss in Eq. (5) and the pruning-aware loss in Eq. (6). A high-level schematic of the joint learning of the model parameters and the dead-zone parameters for CoDeQ are shown in Figure 1, and pseudo-code for the algorithm in Algorithm 1.

### 5.4. CoDeQ with Learnable Bit-Width

The formulation of our joint pruning and quantization method, CoDeQ, assumes fixed bit-width quantization. MP can be incorporated into CoDeQ using learnable bit-widths in the scale factor  $\tilde{s}$  in Eq. (12). To enable this, we parameterize the bit-width value  $b$  using a continuous latent variable  $\theta_{bit} \in \mathbb{R}$ :

$$\tilde{s}(\theta_{bit}) = \frac{\max(|\mathbf{w}|) - d/2}{2^{b(\theta_{bit})-1} - 1 - 1/2}, \quad (17)$$

recalling we have  $Q_b = 2^{b-1} - 1$  and with the learnable bit-width  $b(\theta_{bit})$  given by:

$$b(\theta_{bit}) = \lfloor \tanh(|\theta_{bit}|)(b_{\max} - b_{\min}) + b_{\min} \rfloor. \quad (18)$$

This formulation constrains the effective bit-width to lie within the valid range  $[b_{\min}, b_{\max}]$  while allowing smooth, differentiable updates of  $\theta_{bit}$  during training. The  $\tanh(\cdot)$  function provides a bounded and continuous mapping from the unbounded real parameter  $\theta_{bit} \in \mathbb{R}$  to a normalized interval, ensuring stable gradient flow and preventing invalid bit-width values. Intuitively, this allows learning the extent of the bit-width budget  $b_{\max} - b_{\min}$  to use. When  $|\theta_{bit}|$  is small,  $\tanh(|\theta_{bit}|) \approx 0$ , meaning the bit width becomes  $b_{\min}$ .

During the forward pass, we discretize the continuous bit-width to an integer value using rounding  $\lfloor \cdot \rfloor$ , and apply STE on the rounding operator during the backward pass.

---

#### Algorithm 1: CoDeQ with dead-zone quantizer

---

**Input:** Model with parameters  $\{W_\ell\}$ , Data  $\mathcal{D}$ , Optimizer OPT, reg  $\lambda_{dz}$ , Quantizer  $\tilde{q}(\cdot; \cdot)$ , constants  $b, C$

**Output:**  $\{W_\ell\}, \{\tilde{s}_\ell\}, \{\theta_{dz,\ell}\}$

**Parameters:** *Weights*  $\{W_\ell\}$ , *Dead-Zones*  $\{\theta_{dz,\ell}\}$

**Initialize:**

$Q \leftarrow 2^{b-1} - 1$

$\{\theta_{dz,\ell}\} \leftarrow C$

**Forward pass:**

**foreach**  $\ell \in \text{layers}$  **do**

$R_\ell \leftarrow \max(|W_\ell|)$

$d_\ell \leftarrow 2R_\ell \cdot (1 - \tanh(|\theta_{dz,\ell}|))$

$\tilde{s}_\ell \leftarrow \frac{R_\ell - d_\ell/2}{Q - 1/2}$

$\hat{W}_\ell \leftarrow \tilde{q}(W_\ell; \tilde{s}_\ell, d_\ell)$

$\mathcal{L} \leftarrow \mathcal{L}_{\text{task}}(\{\hat{W}_\ell\}) + \lambda_{dz} \sum_\ell \|\theta_{dz,\ell}\|_2^2$

**Backward pass:**

**foreach**  $\ell \in \text{layers}$  **do**

$\lfloor \text{update } W_\ell, \theta_{dz,\ell} \text{ with OPT} \rfloor$

---

The overall loss for CoDeQ for jointly learning the dead-zone width and bit-width can be written as

$$\mathcal{L}(\tilde{q}(\mathbf{w}; \tilde{s}(\theta_{bit}), d(\theta_{dz}))) + \lambda_{bit} \|\theta_{bit}\|_2^2 + \lambda_{dz} \|\theta_{dz}\|_2^2, \quad (19)$$

where  $\lambda_{bit}$  controls regularization strength for the bit-width parameter.

## 6. Experiments

In this section we empirically evaluate CoDeQ on a range of architectures and datasets, using the loss formulations in Eqs. (16) and (19). We focus on two questions: (i) how CoDeQ compares to recent joint pruning-quantization methods in terms of accuracy and BOPs, and (ii) whether CoDeQ can achieve competitive compression in hardware-friendly fixed-bit regimes.

### 6.1. Experimental Set-Up

**Models and Data:** We perform our experiments on ResNet-18, ResNet-20 and ResNet-50 [8]. We follow the training regime from the original work unless otherwise noted. We fine-tune ResNet-18 and ResNet-50 pretrained models from TorchVision [2] on ImageNet [3] for 120 epochs with cosine annealing on the learning rates. For ResNet-20 we train from scratch on CIFAR-10 [11] for 300 epochs with cosine annealing. For ResNet-18 and ResNet-20 we use a batch-size of 512 and for ResNet-50 a batch-size 256. The Tiny Vision Transformer [19] is trained from

scratch on CIFAR10 for 200 epochs with an AdamW optimizer,  $3 \times 10^{-4}$  learning rate with cosine annealing and a batch size of 128.

**Quantization Configuration:** For mixed-precision (MP) experiments we make both bit-width and dead-zone learnable. For 4-bit experiments, we keep bit-width fixed and only learn the dead-zone. In all cases we use layer-wise quantization granularity, meaning each layer has one scale factor. Using the absmax scale factor we follow the conventions from [9] of not tracking the gradient for the operation  $\text{absmax} = \max(|w|)$ . Additionally we use the 99<sup>th</sup> quantile of the absmax for stability with outliers. To avoid division by zero in the case where the whole layer is pruned, i.e.  $d = 2 \max(|w|)$  we add an  $\epsilon = 10^{-8}$  to the scale factor in Eq. (12).

**Initialization.** In all experiments we initialize  $\theta_{dz}$  and  $\theta_{bit}$  to 3, which gives  $\tanh(3) \approx 0.995$  meaning we initialize with a full bit budget and a dead-zone width close to 0. For all experiments with MP we set minimum bit  $b_{\min} = 2$  and maximum bit  $b_{\max} = 8$  in Eq. (18).

**Learning Rates.** For fine-tuning we use a learning rate of  $10^{-4}$  on both  $\theta_{dz}$  and  $\theta_{bit}$ , while we use a higher learning rate of  $10^{-3}$  when training from scratch.

**Lambdas.** For ResNet-20 we set  $\lambda_{bit} = 0.01$  and  $\lambda_{dz} = 0.01$ , ResNet-18  $\lambda_{bit} = 0.1$  and  $\lambda_{dz} = 0.02$  and ResNet-50  $\lambda_{bit} = 0.05$  and  $\lambda_{dz} = 0.02$ . For TinyVit we set  $\lambda_{dz} = 0.5$  and  $\lambda_{bit} = 0.75$ . When bit is not learned we use the  $\lambda_{dz}$  values also used in the MP experiments.

**Methods for Comparison** We compare CoDeQ against several joint pruning-quantization methods: QST [13], SQL [20], CLIP-Q [15], and GETA [14]. When available, we report numbers from the original papers.

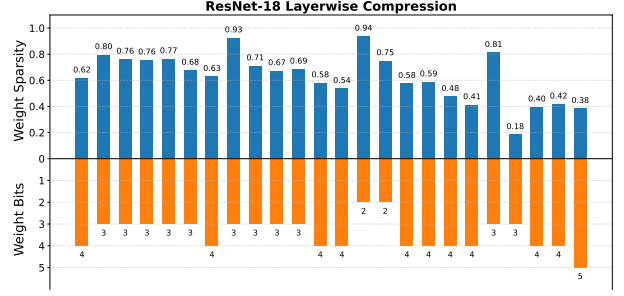
**Compression Metric Bit-Operations (BOPs):** For evaluating the computational efficiency of the compressed models, we adopt the BOP formulation from [18], but adapt their MAC accounting to reflect unstructured sparsity rather than structured (channel-wise) pruning. Specifically, we scale the dense MAC cost by the weight density  $d_W$ .

For a convolutional layer  $\ell$  with  $c_{\ell-1}$  input channels,  $c_\ell$  output channels, kernel size  $k_h \times k_w$ , and spatial output size  $m_{h,\ell} \times m_{w,\ell}$ , the dense MAC count is

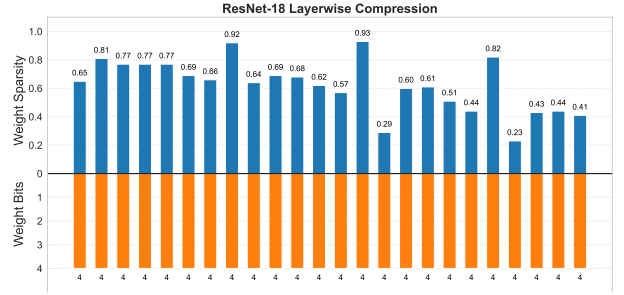
$$\text{MACs}_\ell^{\text{dense}} = c_{\ell-1} \cdot c_\ell \cdot k_h \cdot k_w \cdot m_{h,\ell} \cdot m_{w,\ell}. \quad (20)$$

To incorporate unstructured weight sparsity, we let  $d_W \in [0, 1]$  denote the fraction of non-zero weights in layer  $\ell$ . We ignore activation sparsity and any structured reduction in input channels. The expected MAC count then becomes

$$\text{MACs}_\ell^{\text{unstruct}} = d_W \cdot \text{MACs}_\ell^{\text{dense}}. \quad (21)$$



(a) Layer-wise compression for a ResNet-18 model with MP. Negative  $y$ -axis shows bit-width (descending) and positive  $y$ -axis shows sparsity (ascending).



(b) Layer-wise compression for a ResNet-18 model with fixed 4-bit weights. Negative  $y$ -axis shows bit-width (descending) and positive  $y$ -axis shows sparsity (ascending).

Figure 3. Layer-wise compression patterns learned by CoDeQ for ResNet-18 on ImageNet.

Given weight and activation precisions  $w_{\text{bits}}$  and  $a_{\text{bits}}$ , respectively, the resulting BOPs for layer  $\ell$  are

$$\text{BOPs}_\ell = \text{MACs}_\ell^{\text{unstruct}} \cdot w_{\text{bits}} \cdot a_{\text{bits}}. \quad (22)$$

We set  $a_{\text{bits}} = 32$  in all experiments, as we do not quantize activations. We note that this BOP metric reflects arithmetic cost under idealized support for unstructured weight sparsity.

## 6.2. Results

We evaluate CoDeQ across CIFAR-10 and ImageNet benchmarks, focusing on both Top-1 accuracy and BOPs. Across all experiments, CoDeQ delivers competitive or superior accuracy-compression trade-offs with a small number of compression hyperparameters and *without auxiliary procedures*.

### 6.2.1. CIFAR-10 Experiments

**ResNet-20.** We first assess CoDeQ on CIFAR-10 using ResNet-20, a widely used benchmark for pruning and quantization. As shown in Tab. 2, CoDeQ achieves the highest accuracy among all compared approaches while also delivering the lowest BOPs. Notably, the fixed-bit variant (despite being constrained to uniform 4-bit precision), outper-

forms the next-best method, QST [13], in both accuracy and BOPs. MP bit yields slightly lower BOPs for CoDeQ, but the fixed-bit model performs nearly identically, providing evidence that CoDeQ generalizes effectively to hardware-friendly fixed-bit settings.

**Vision Transformers.** To assess generality beyond convolutional architectures, we evaluate CoDeQ on the TinyViT vision transformer [19] on CIFAR-10, summarized in Tab. 5. CoDeQ achieves strong compression with minimal accuracy degradation, preserving performance close to the uncompressed model while substantially reducing BOPs. These results show that CoDeQ can be effectively and readily applied to transformer-based models.

### 6.2.2. ImageNet Experiments

**ResNet-18.** We next validate CoDeQ on ImageNet using ResNet-18, a standard compression benchmark. As reported in Tab. 3, CoDeQ matches the accuracy of QST within run-to-run variation while achieving a further reduction in BOPs. Fixed-bit and MP variants again exhibit nearly identical accuracy–compression behavior, confirming that fixed-bit training remains stable at larger scale. The learned layer-wise compression patterns, visualized in Fig. 3, show that CoDeQ naturally assigns higher precision and lower sparsity to the first and last layers, consistent with established heuristics, while allocating greater sparsity and lower precision to deeper layers.

**ResNet-50.** On ResNet-50, shown in Tab. 4, CoDeQ delivers substantial efficiency gains while maintaining competitive accuracy. Compared to QST, CoDeQ incurs a small top-1 accuracy drop of approximately 0.81% but reduces BOPs from 4.5% to 2.62% relative to the full-precision baseline. This represents a significant improvement in computational efficiency and highlights the practical value of CoDeQ for large-scale deployments where modest accuracy changes must be balanced against large reductions in compute.

### 6.3. Ablation

In CoDeQ, sparsity is primarily controlled by the dead-zone parameter  $\theta_{dz}$ . Regularizing  $\theta_{dz}$  encourages a wider dead-zone, which in turn prunes more weights. However, standard weight decay on the model weights can also indirectly promote sparsity, as it pulls weights toward zero and thus into the dead-zone.

To study this interaction, we compare two configurations on Figs. 4a and 4b; (i) CoDeQ while varying the dead-zone regularization but without additional weight decay on model weights, and (ii) CoDeQ with fixed dead-zone regularization but with weight decay on the weights. As seen

Table 2. Performance of ResNet-20 on CIFAR-10. Results are reported as the mean and standard deviation over 3 runs.

Method	Weights	Pruning	Acc (%)	Rel. BOPs (%)
Baseline	32bit	None	91.70	100
SQL [20]	MP	Unstruc.	90.90	6.1
QST-P [13]	MP	Unstruc.	91.80	5.0
GETA [14]	MP	Struct.	91.42	4.5
QST-Q [13]	MP	Unstruc.	91.60	3.3
CoDeQ (Ours)	4bit	Unstruc.	<b>91.88</b> $\pm$ 0.15	<b>2.95</b> $\pm$ 0.12
CoDeQ (Ours)	MP	Unstruc.	<b>91.86</b> $\pm$ 0.10	<b>2.61</b> $\pm$ 0.08

in Fig. 4b, applying weight decay alone can reach a sparsity level comparable to the dead-zone-regularized model by progressively shrinking weights into the zero bin, rather than expanding the bin itself as in Fig. 4a. In contrast, dead-zone regularization explicitly increases the width of the zero bin, allowing large-magnitude weights to remain relatively unaffected while aggressively pruning small ones.

This comparison highlights an important interplay between dead-zone regularization and weight decay: both can induce sparsity, but via different mechanisms.

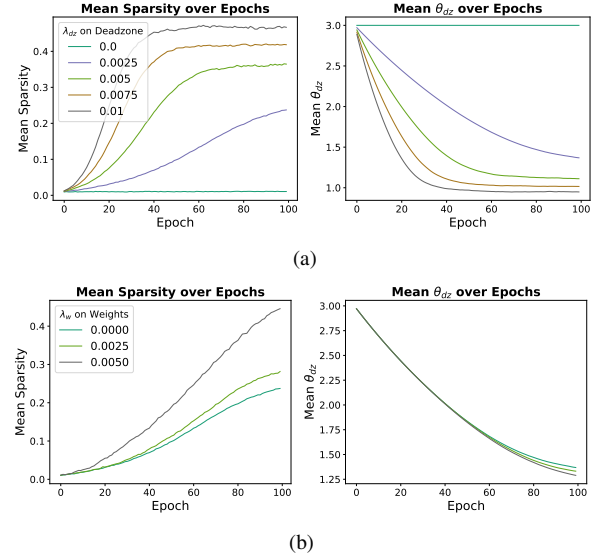


Figure 4. ResNet-20 on CIFAR-10 trained for 100 epochs with varying levels of  $\lambda_{dz}$  in plot (a) we note how increasing the regularization strength, reduces  $\theta_{dz}$  and increases sparsity. In plot (b), we instead add weight decay on the weights while keeping  $\lambda_{dz}$  fixed.  $\lambda_w$  is the strength of weight decay on the model parameters. Note how the dead-zone remains close, for different levels of sparsity.

## 7. Discussion

We have shown empirically that CoDeQ is a competitive joint pruning–quantization method that matches or surpasses recent baselines across architectures and datasets, while remaining simple to configure and deploy.



Table 3. Performance of ResNet-18 on ImageNet. Results are reported as the mean and standard deviation over 3 runs.

Method	Weights	Pruning	Acc (%)	Rel. BOPs (%)
Baseline	32bit	None	70.30	100
SQL [20]	MP	Unstruc.	68.60	6.20
QST-B [13]	MP	Unstruc.	<b>69.90</b>	5.00
CoDeQ (Ours)	4bit	Unstruc.	<b>69.83</b> $\pm$ 0.14	<b>4.75</b> $\pm$ 0.19
CoDeQ (Ours)	MP	Unstruc.	<b>69.81</b> $\pm$ 0.09	<b>4.42</b> $\pm$ 0.05

Table 4. Performance of ResNet-50 on ImageNet. Results are reported based on one run.

Method	Weights	Pruning	Acc (%)	Rel. BOPs (%)
Baseline	32bit	None	76.30	100
CLIP-Q [15]	MP	Unstruc.	73.70	6.30
GETA [14]	MP	Struct.	74.40	5.38
QST-B [13]	MP	Unstruc.	<b>76.10</b>	4.50
CoDeQ (Ours)	MP	Unstruc.	75.29	<b>2.62</b>

Table 5. Mean and std. of 3 runs with a TinyVIT on CIFAR10. Results are reported as the mean and standard deviation over 3 runs.

Method	Weights	Pruning	Acc (%)	Rel. BOPs (%)
Baseline	32bit	None	87.99 $\pm$ 0.09	100
CoDeQ (Ours)	4bit	Unstruc.	87.86 $\pm$ 0.19	2.15 $\pm$ 0.04
CoDeQ (Ours)	MP	Unstruc.	87.85 $\pm$ 0.21	1.66 $\pm$ 0.01

**Low BOPs at Comparable Accuracy.** Across all benchmarks, CoDeQ attains the lowest BOPs among methods with similar accuracy to existing methods. Learning dead-zone widths (and, optionally, MP bit-widths) via a single differentiable operator proves to be a strong alternative to approaches that rely on auxiliary procedures to determine compression parameters.

**Sparsity With Minimal Hyperparameters.** Existing joint pruning–quantization methods often require substantial user effort to adopt. In CoDeQ, sparsity emerges from learned dead-zone widths and is governed by a single global regularization parameter. This substantially reduces hyperparameter tuning and shifts sparsity allocation to a data-driven inner-loop mechanism that applies uniformly across architectures.

**Decoupling Bit-Width and Sparsity selection.** Our experiments show that CoDeQ maintains strong accuracy–efficiency trade-offs under fixed 4-bit quantization, and that fixed-bit and mixed-precision variants perform similarly. Thus, CoDeQ can provide competitive compression without mixed-precision search, offering a direct path to hardware deployment where bit-width is fixed and sparsity is the principal degree of freedom.

**Limitations and Future Work.** This work has three main limitations. First, we use layer-wise quantization and do not explore finer granularities (e.g., channel- or block-wise), which could further improve accuracy. Second, we rely on a classical absmax scale factor rather than learning scales (e.g., via learned step-size quantization [4]). Third, CoDeQ currently induces unstructured sparsity, which is generally harder to exploit efficiently on existing hardware than structured pruning.

## 8. Conclusion

We have shown, both theoretically and empirically, that the dead-zone of a scalar quantizer naturally implements magnitude pruning. Building on this observation, we introduced CoDeQ, a simple, fully differentiable method that unifies pruning and quantization within a single parameterized quantizer. Our analysis show how the dead-zone induce sparsity, and our experiments demonstrate that CoDeQ yields sparse, low-precision models with competitive accuracy across convolutional and transformer architectures, in both fixed-bit and mixed-precision regimes. By removing the need for auxiliary procedures and explicit pruning schedules, CoDeQ offers a practical path toward user-friendly, fixed-bit compatible model compression, where sparsity and precision are learned directly from data.

## Acknowledgments

Authors thank members of Machine Learning Section and [SAINTS Lab](#) for useful discussions throughout. Authors acknowledge the funding received from the European Union’s Horizon Europe Research and Innovation Action programme under grant agreements No. 101070284, No. 101070408 and No. 101189771. RS also acknowledges funding received under Independent Research Fund Denmark (DFF) under grant agreement number 4307-00143B.

## References

- [1] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013. 4
- [2] TorchVision Contributors. Torchvision: Pytorch computer vision library. <https://github.com/pytorch/vision>, 2016. 6
- [3] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009. 6
- [4] Steven K. Esser, Jeffrey L. McKinstry, Deepika Bablani, Rathinakumar Appuswamy, and Dharmendra S. Modha. Learned step size quantization. In *International Conference on Learning Representations*, 2020. 9

- [5] Yiwen Guo, Anbang Yao, and Yurong Chen. Dynamic network surgery for efficient dnns. *Advances in neural information processing systems*, 29, 2016. 4
- [6] Song Han, Jeff Pool, John Tran, and William Dally. Learning both weights and connections for efficient neural networks. In *Advances in Neural Information Processing Systems*, pages 1135–1143, 2015. 4
- [7] Song Han, Huizi Mao, and William J. Dally. Deep compression: Compressing deep neural network with pruning, trained quantization and huffman coding. In *ICLR*, 2016. 3
- [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778. IEEE, 2016. 6
- [9] Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for efficient integer-arithmetic-only inference. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2704–2713, 2018. 4, 7
- [10] Raghuraman Krishnamoorthi. Quantizing deep convolutional networks for efficient inference: A whitepaper, 2018. 4
- [11] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009. 6
- [12] Yann LeCun, John Denker, and Sara Solla. Optimal brain damage. In *Advances in Neural Information Processing Systems*. Morgan-Kaufmann, 1989. 4
- [13] Jun-Hyung Park, Kang-Min Kim, and Sangkeun Lee. Quantized sparse training: A unified trainable framework for joint pruning and quantization in dnns. *ACM Trans. Embed. Comput. Syst.*, 21(5), 2022. 2, 3, 5, 7, 8, 9
- [14] Xiaoyi Qu, David Aponte, Colby Banbury, Daniel P Robinson, Tianyu Ding, Kazuhito Koishida, Ilya Zharkov, and Tianyi Chen. Automatic joint structured pruning and quantization for efficient neural network training and compression. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 15234–15244, 2025. 3, 7, 8, 9
- [15] Frederick Tung and Greg Mori. Clip-q: Deep network compression learning by in-parallel pruning-quantization. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7873–7882, 2018. 2, 3, 7, 9
- [16] Mart Van Baalen, Christos Louizos, Markus Nagel, Rana Ali Amjad, Ying Wang, Tijmen Blankevoort, and Max Welling. Bayesian bits: Unifying quantization and pruning. *Advances in neural information processing systems*, 33:5741–5752, 2020. 2, 3
- [17] Tianzhe Wang, Kuan Wang, Han Cai, Ji Lin, Zhijian Liu, and Song Han. APQ: joint search for network architecture, pruning and quantization policy. *CoRR*, abs/2006.08509, 2020. 2, 3
- [18] Ying Wang, Yadong Lu, and Tijmen Blankevoort. Differentiable joint pruning and quantization for hardware efficiency. In *European Conference on Computer Vision*, pages 259–277. Springer, 2020. 2, 3, 7
- [19] Kan Wu, Jinnian Zhang, Houwen Peng, Mengchen Liu, Bin Xiao, Jianlong Fu, and Lu Yuan. Tinyvit: Fast pretraining distillation for small vision transformers. In *European conference on computer vision*, pages 68–85. Springer, 2022. 6, 8
- [20] Haichuan Yang, Shupeng Gui, Yuhao Zhu, and Ji Liu. Automatic neural network compression by sparsity-quantization joint learning: A constrained optimization-based approach. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2178–2188, 2020. 2, 3, 7, 8, 9
- [21] Ben Zandonati, Glenn Bucagu, Adrian Alan Pol, Maurizio Pierini, Olya Sirkin, and Tal Kopetz. Towards optimal compression: Joint pruning and quantization, 2023. 2, 3