## A. Gradient of Quantization Mapping

In this section, we present the formula for the gradient of the quantization mapping with respect to $d$, $t$, and $q_m$.

$$\nabla_d x^Q = \text{sgn}(x) \cdot \begin{cases} \left( \lfloor \frac{|x|^t}{d} \rceil - \frac{|x|^t}{d} \right), & |x| \le q_m, \\ \left( \lfloor \frac{(q_m)^t}{d} \rceil - \frac{(q_m)^t}{d} \right), & |x| > q_m, \end{cases} \tag{15}$$

$$\nabla_t x^Q = \text{sgn}(x) \cdot \begin{cases} |x|^t \log(|x|), & |x| \le q_m, \\ (q_m)^t \log(q_m), & |x| > q_m, \end{cases} \tag{16}$$

$$\nabla_{q_m} x^Q = \begin{cases} 0, & |x| \le q_m, \\ \text{sgn}(x) t (q_m)^{t-1}, & |x| > q_m. \end{cases} \tag{17}$$

**Remark.** The computation involving $x$ in this section represents element-wise operations.

## B. Proof for Proposition 5.1

In this section, we present the proof for Proposition 5.1. For convenience, we restate the proposition as follows.

**Proposition B.1.** *Let $\hat{\nabla}_x f$ be the full gradient of function $f(x, d, q_m, t)$ with respect to $x$. With forget rate $\gamma$ selection rule Eq. (13) and quantization step size $d$ selection rule Eq. (14), the search direction $s(x)$ is a descent direction for the function $f$ with respect to $x$ at $x$.*

*Proof.* Denote the full gradient of function $f(x, d, q_m, t)$ with respect to $x$ as $\nabla_x f$. The search direction $s(x)$ is rewritten as

$$s(x) = \begin{cases} -\alpha [\nabla_x f]_g, & g \in \mathcal{G}_I, \\ -\alpha [\nabla_x f]_g - \gamma [x^Q]_g, & g \in \mathcal{G}_R. \end{cases} \tag{18}$$

Since $-\alpha [\nabla_x f]_g^T [\nabla_x f]_g < 0$ for $g \in \mathcal{G}_I$, it suffices to show that for $g \in \mathcal{G}_R$,

$$[\nabla_x f]_g^T \left[ -\alpha [\nabla_x f]_g - \gamma [x^Q]_g \right] < 0.$$

It follows from (9) that for $g \in \mathcal{G}_R$,

$$-\alpha [\nabla_x f]_g - \gamma [x^Q]_g = \underbrace{-\alpha [\nabla_x f]_g - \gamma [\text{sgn}(x) \cdot \text{clip}_{q_m}^t (|x|)]_g}_{[s_{\text{clip}}(x)]_g} - \gamma \cdot d [\text{sgn}(x) \cdot R(x)]_g.$$

Denote the angle between $-[\nabla_x f]_g$ and $-[\text{sgn}(x) \cdot \text{clip}_{q_m}^t (|x|)]_g$ as $\theta_\gamma$. It follows that the vector $[s_{\text{clip}}(x)]_g$ can be decomposed into two orthogonal vectors, i.e.,

$$[s_{\text{clip}}(x)]_g = [\hat{s}_{\text{clip}}(x)]_g + [\tilde{s}_{\text{clip}}(x)]_g,$$

where $[\hat{s}_{\text{clip}}(x)]_g$ is orthogonal to vector $[\nabla_x f]_g$ and $[\tilde{s}_{\text{clip}}(x)]_g$ is parallel to vector $[\nabla_x f]_g$. Since $[\hat{s}_{\text{clip}}(x)]_g^T [\nabla_x f]_g = 0$, we have that

$$\| [\hat{s}_{\text{clip}}(x)]_g \| = \gamma \sin \theta_\gamma \| [\text{sgn}(x) \cdot \text{clip}_{q_m}^t (|x|)]_g \|.$$

Using the orthogonality between vector $[\hat{s}_{\text{clip}}(x)]_g$ and vector $[\tilde{s}_{\text{clip}}(x)]_g$, we have that

$$\begin{aligned}
\| [\tilde{s}_{\text{clip}}(x)]_g \|^2 &= \| [s_{\text{clip}}(x)]_g \|^2 - \| [\hat{s}_{\text{clip}}(x)]_g \|^2 \\
&= \| -\alpha [\nabla_x f]_g - \gamma [\text{sgn}(x) \cdot \text{clip}_{q_m}^t (|x|)]_g \|^2 - \gamma^2 \sin^2 \theta_\gamma \| [\text{sgn}(x) \cdot \text{clip}_{q_m}^t (|x|)]_g \|^2 \\
&= \alpha^2 \| [\nabla_x f]_g \|^2 + 2\alpha\gamma [\nabla f(x)]_g^T [\text{sgn}(x) \cdot \text{clip}_{q_m}^t (|x|)]_g + \gamma^2 \cos^2 \theta_\gamma \| [\text{sgn}(x) \cdot \text{clip}_{q_m}^t (|x|)]_g \|^2 \\
&= \left[ \alpha \| [\nabla_x f]_g \| + \gamma \cos \theta_\gamma \| [\text{sgn}(x) \cdot \text{clip}_{q_m}^t (|x|)]_g \| \right]^2.
\end{aligned}$$

Given the norm and direction of the vector $[\tilde{s}_{\text{clip}}(x)]_g$, we have $[\tilde{s}_{\text{clip}}(x)]_g$ expressed as, for $g \in \mathcal{G}_R$,

$$[\tilde{s}_{\text{clip}}(x)]_g = -\frac{\alpha \| [\nabla_x f]_g \| + \gamma \cos(\theta_\gamma) \| [\text{sgn}(x) \cdot \text{clip}_{q_m}^t (|x|)]_g \|}{\| [\nabla_x f]_g \|} [\nabla_x f]_g. \tag{19}$$

Combining the forget rate selection rule (13) and the expression (19) allows us to have that for $g \in \mathcal{G}_R$,

$$
\begin{aligned}
[\nabla_x f]_g^T [s_{\text{clip}}(x)]_g &= [\nabla_x f]_g^T \left[ [\hat{s}_{\text{clip}}(x)]_g + [\tilde{s}_{\text{clip}}(x)]_g \right] \\
&= [\nabla_x f]_g^T [\tilde{s}_{\text{clip}}(x)]_g \\
&= -\alpha \|[\nabla_x f]_g\|^2 - \gamma \cos(\theta_\gamma) \|[\nabla_x f]_g\| \| [\text{sgn}(x) \cdot [\text{clip}_{q_m}^t(|x|)]_g] \| \\
&< -\eta \alpha \|[\nabla_x f]_g\|^2 .
\end{aligned}
\tag{20}
$$

Further, our quantization step size $d$ selection rule (14) guarantees that

$$
-\eta \alpha \|[\nabla_x f]_g\|^2 - \gamma d [\nabla_x f]_g^T [\text{sgn}(x) \cdot R(x)]_g < 0.
\tag{21}
$$

Combining Eq. (20) and Eq. (21) allows us to have that

$$
\begin{aligned}
[\nabla_x f]_g^T \left[ -[\nabla_x f]_g - \gamma [x^Q]_g \right] &= [\nabla_x f]_g^T \left[ [s_{\text{clip}}(x)]_g - \gamma \cdot d[\text{sgn}(x) \cdot R(x)]_g \right] \\
&< -\eta \alpha \|[\nabla_x f]_g\|^2 - \gamma d[\nabla_x f]_g^T [\text{sgn}(x) \cdot R(x)]_g < 0,
\end{aligned}
$$

which completes the proof. □

## C. Joint Stage Implementation Details

The update rule in Eq. (14) alone is insufficient to ensure that the bit width constraint in Eq. (7b) is consistently satisfied. To address this issue, we introduce an algorithm, outlined in Algorithm 4, to adaptively adjust the forget rate $\gamma$ and quantization step size $d$ such that the computed bit width stays within the target bit width range $[b_l, b_u]$. Meanwhile, with the adaptive algorithm in place, the search direction $s(x)$ continues to be a descent direction when stochastic gradient $\hat{\nabla}_x f$ is assumed to be full gradient, as demonstrated in Proposition C.1. In addition, there are three hyperparameters that appear in Eq. (13) and Eq. (14) and they are selected as $\eta = 0.9$, $\xi = 0.999$, and $\epsilon = $1e-8.

**Proposition C.1.** *Let $\hat{\nabla}_x f$ be the full gradient of function $f(x, d, q_m, t)$ with respect to $x$. With the Algorithm 4 in place (applied immediately after Line 17 in Algorithm 1), the search direction $s(x)$ is still a descent direction with respect to function $f$ at the point $x$.*

*Proof.* Denote the full gradient of function $f(x, d, q_m, t)$ with respect to $x$ as $\nabla_x f$. Let's first consider the following three simple cases. When $clip > \epsilon$, the forget rate is equal to 0 and therefore, $s(x) = -\nabla f(x)$, which is a descent direction with respect to function $f(x, d, q_m, t)$ with respect to $x$ at $x$. When $\cos(\theta_\gamma) \geq 0$ or $\cos(\theta_d) \geq 0$, $s(x)$ is guaranteed to be descent direction with respect to function $f(x, d, t, q_m)$ with respect to $x$ when $\gamma$ and $d$ are positive values.

Now, it remains to consider the following two cases: $\cos(\theta_\gamma) < 0, clip > \epsilon$ and $\cos(\theta_d) < 0$. As indicated in Eq. (13) and Eq. (14), we have that $s(x)$ is a descent direction with respect to function $f(x, d, q_m, t)$ with respsect to $x$ if when $\cos(\theta_\gamma) < 0$ and $clip > \epsilon$,

$$
\gamma \in \left( 0, -\frac{\alpha \|[\nabla_x f]_g\|}{\cos(\theta_\gamma) \| [\text{sgn}(x) \cdot \text{clip}_{q_m}^t(|x|)]_g \|} \right),
\tag{22}
$$

and when $\cos(\theta_d) < 0$,

$$
d \in \left( 0, -\frac{\xi \eta \alpha \|[\nabla_x f]_g\|}{\gamma \cos(\theta_d) \| [\text{sgn}(x) \cdot R(x)]_g \|} \right).
\tag{23}
$$

When $\cos(\theta_\gamma) < 0$, we guarantee that with the Algorithm 4 in place, the forget rate $\gamma$ always lie in the range specified in Eq. (22) since $\gamma$ either decreases by a factor of $\beta$ (see Line 4) or remains the same (see Line 6). When $\cos(\theta_d) < 0$, we consider two cases based on if the forget rate decreases. If forget rate decreases (see Line 4), then the range given in Eq. (14) is changed to

$$
\left( 0, -\frac{\xi \eta \alpha \|[\nabla_x f]_g\|}{\beta \gamma \cos(\theta_d) \| [\text{sgn}(x) \cdot R(x)]_g \|} \right)
\tag{24}
$$

It follows that increasing $d$ by a factor of $\frac{1}{\beta}$ guarantees that $d$ lies within the range Eq. (24). If forget rate remains the same, then $d$ always lie in the range Eq. (23) since $d$ decreases by a factor of $\beta$. □

---

**Algorithm 4** Adaptive update rule for $\gamma$ and $d$.

---

1: **Inputs.** Variables: $\gamma$, $d$, bit width range: $[b_l, b_u]$, $\beta \in (0, 1)$, fixed quantization variables: $q_m, t$.

2: **while** $\log_2 \left( \frac{(q_m)^t}{d} + 1 \right) + 1 \notin [b_l, b_u]$ **do**

3:     **if** $\log_2 \left( \frac{(q_m)^t}{d} + 1 \right) + 1 > b_u$ **then**

4:         $\gamma \leftarrow \beta\gamma, d \leftarrow d/\beta$.

5:     **else if** $\log_2 \left( \frac{(q_m)^t}{d} + 1 \right) + 1 < b_l$ **then**

6:         $\gamma \leftarrow \gamma, d \leftarrow \beta d$.

7:     **end if**

8: **end while**

9: **Outputs.** $\gamma$, $d$.

---

## D. Numerical Experiment Setup

First, we provide details on how we initialize quantization parameters. For each layer that contain quantization parameters, the exponential $t = 1$ and the maximum of quantization range $q_m$ is set to the layerwise maximum of the weight tensor. For experiments on ResNet20, VGG7, and ResNet50, the quantization step size $d$ is chosen such that the resulting bit width is 32 bits while for Bert, the quantization step size $d$ is selected to achieve a bit width of 8 bits.

Next, we provide details on how we select the optimizer and the learning rate for different experiments. For ResNet20, we use the SGD optimizer and the initial learning rate is set to 1e-1 with StepLR learning rate scheduler. For experiments of VGG7, we use the optimizer ADAM and the learning rate is set to 1e-3 with StepLR learning rate scheduler. For ResNet50, we use the optimizer SGD and the learning rate is set to 1e-1 with StepLR learning rate scheduler. For Bert, we use AdamW with learning rate as constant 3e-5. For all four experiments, the learning rate for quantization parameters is set as constant 1e-4. For details on how we set hyperparameters related to projection stage and pruning stage, one can find them in Tab. 7.

Table 7. Experiment setup for all four experiments. In the following table, the unit for projection steps and pruning steps is the number of epochs. As for Bert, the experiment setups are same under all sparsity ratios (10%, 30%, 50%, 70%).

| Model | Sparsity level | Total epochs | Projection periods $B$ | Projection steps $K_b$ | Pruning periods $P$ | Pruning steps $K_p$ | Bit width reduction $b_r$ | Bit width range $[b_l, b_u]$ |
|---|---|---|---|---|---|---|---|---|
| ResNet20 | 0.35 | 350 | 7 | 35 | 5 | 30 | 2 | [4,16] |
| VGG7 | 0.7 | 200 | 5 | 20 | 10 | 30 | 2 | [4,16] |
| ResNet50 | 0.4,0.5 | 120 | 5 | 5 | 10 | 10 | 2 | [4,16] |
| Bert | 0.1,0.3,0.5,0.7 | 10 | 4 | 1 | 6 | 6 | 2 | [4,16] |

## E. Ablation Study

Our proposed QASSO consists of four distinct stages: warm-up stage, projection stage, joint stage, and cool-down stage. To evaluate the contribution of each stage, we conduct an ablation study on two benchmarks, ResNet56 trained from scratch on CIFAR10 and Phi2 fine-tuned from a pre-trained model on the Common-Sense task. The results demonstrate that each stage positively contributes to the model's performance, as measured by test accuracy. As shown in Fig. 4a, removing any of the four stages, especially the joint stage and cool-down stage, results in a noticeable decline in test accuracy. The significance of the joint stage and cool-down stage stems from the fact that a significant knowledge transfer is conducted to retain the information lost when applying pruning and quantization.

Moreover, each stage's contribution varies over downstream applications. For instance, the joint stage plays a more critical role when fine-tuning a pre-trained model compared to training from scratch. This can be attributed to the fact that pre-trained models inherently possess a wealth of useful knowledge, and the joint stage helps preserve performance by effectively transferring this knowledge under quantization constraints.

In addition, we perform an ablation study (See Fig. 4b) using ResNet56 on CIFAR10 benchmark to study the limit of each compression technique within **GETA** framework. As highlighted in [27], structured pruning methods typically achieve sparsity greater than 80%. However, under joint setup, accuracy begins to degrade significantly beyond 60% sparsity. This

suggests quantization error constrains aggressive pruning, lowering the achievable sparsity threshold from 80% to 60% for ResNet56-CIFAR10. For quantization, satisfactory accuracy is typically retained with bit width $\geq$ 2bits when sparsity $\leq$ 60%. When sparsity exceeds 60%, model becomes less tolerant to lower bit width, requiring at least 4-bit to retain performance.

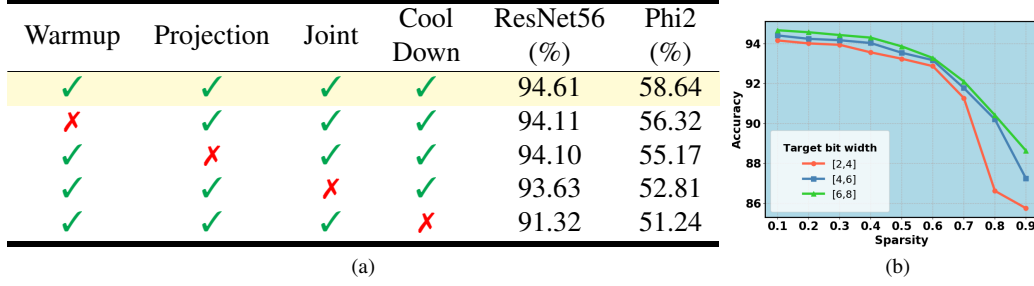| Warmup | Projection | Joint | Cool Down | ResNet56 (%) | Phi2 (%) |
|:---:|:---:|:---:|:---:|:---:|:---:|
| ✓ | ✓ | ✓ | ✓ | 94.61 | 58.64 |
| ✗ | ✓ | ✓ | ✓ | 94.11 | 56.32 |
| ✓ | ✗ | ✓ | ✓ | 94.10 | 55.17 |
| ✓ | ✓ | ✗ | ✓ | 93.63 | 52.81 |
| ✓ | ✓ | ✓ | ✗ | 91.32 | 51.24 |

(a)



(b)

Figure 4. The Fig. 4a presents an ablation study evaluating the necessity of the four distinct stages of the QASSO optimizer using ResNet56 on the CIFAR10 benchmark and Phi2-2.7B on a common-sense task. The last two columns indicate the model's test accuracy. The Fig. 4b illustrates the limits and boundaries of various compression techniques applied to ResNet56 on the CIFAR10 dataset.

# F. Quantization-Aware Dependency Graph

For more intuitive illustration, we present quantization-aware dependency graphs of Bert1 (mini-Bert with one transformer block) and VGG7. Both the original and post-analysis versions of these graphs are shown. To enhance readability of the graph's finer details, we recommend zooming in to a scale of 1500% or higher using Adobe PDF Reader.
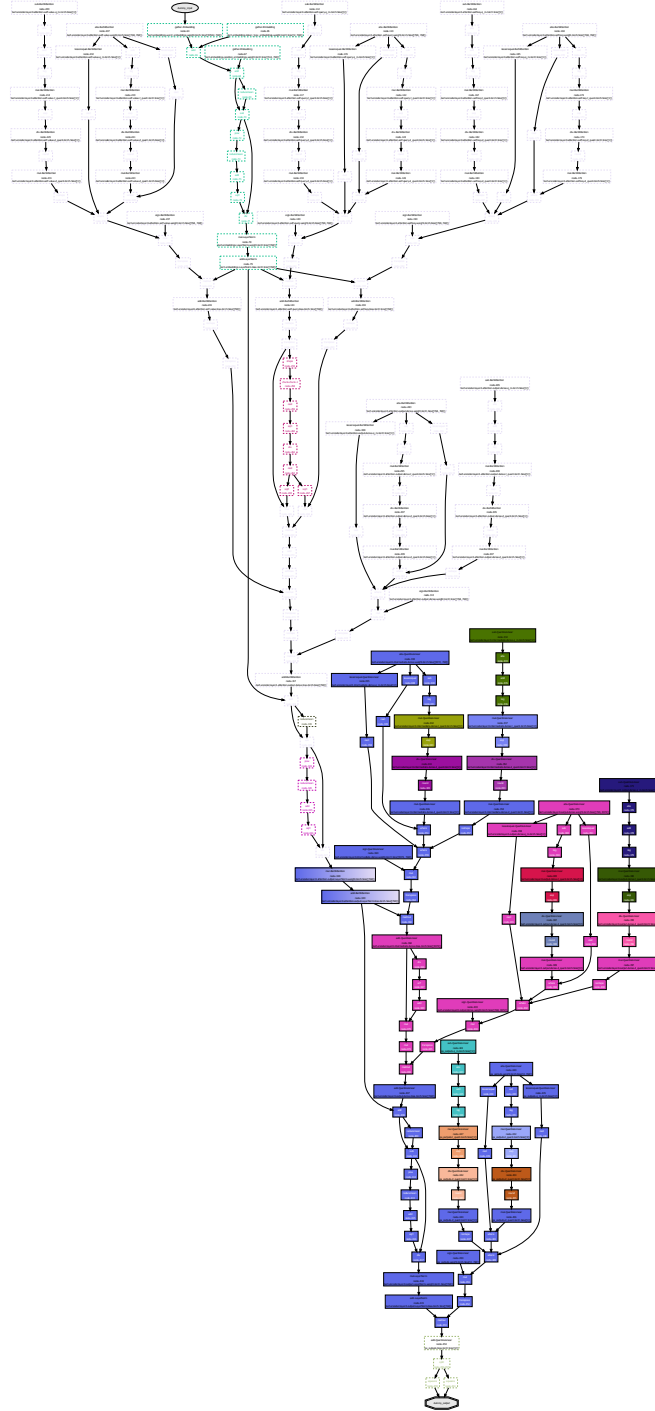
**Bert.**



Figure 5. Bert1 before performing quantization-aware dependency graph analysis.
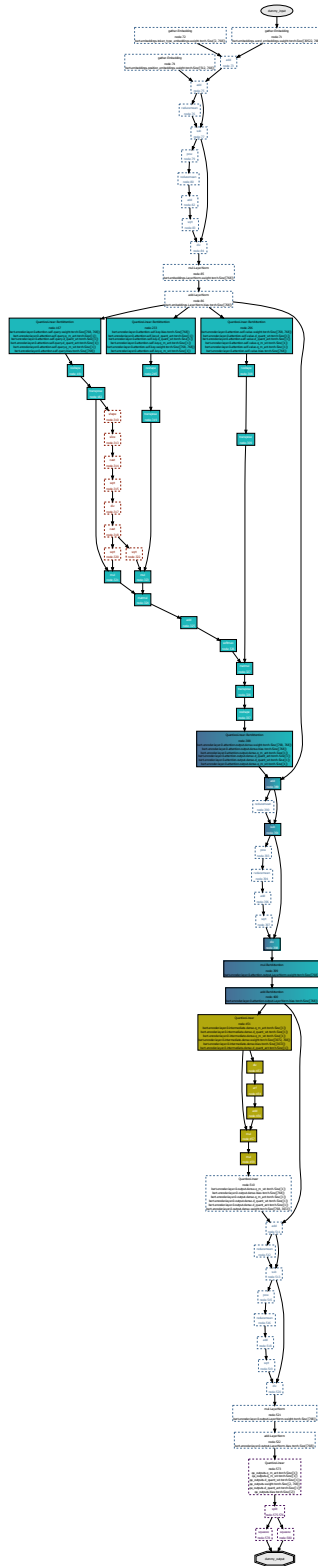
BERT.



Figure 6. Bert1 after performing quantization-aware dependency graph analysis.
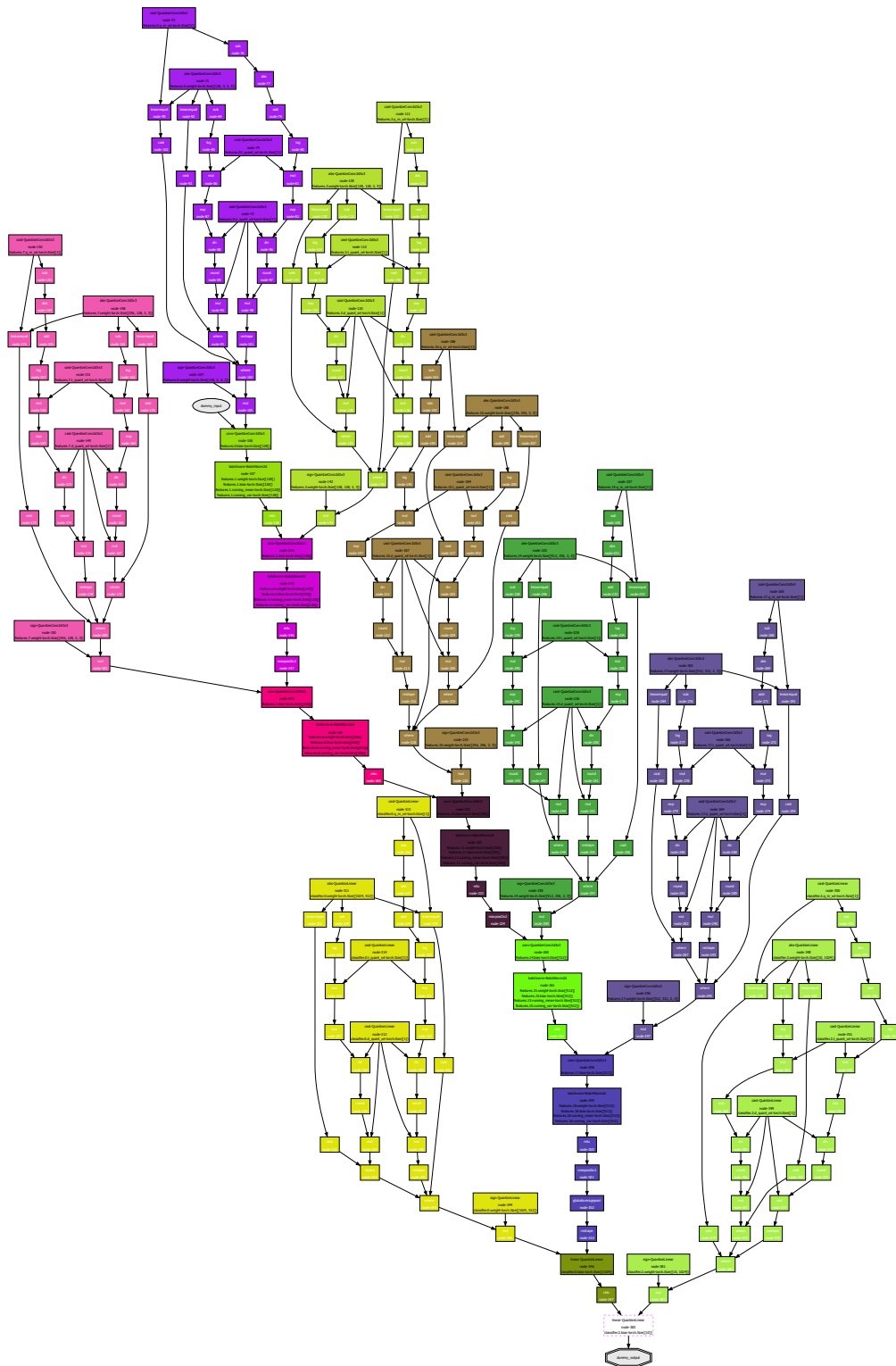
VGG7.



Figure 7. VGG7 after performing quantization-aware dependency graph analysis.

VGG7.



Figure 8. VGG7 after performing quantization-aware dependency graph analysis.