

sconce v0.99

Feature Support

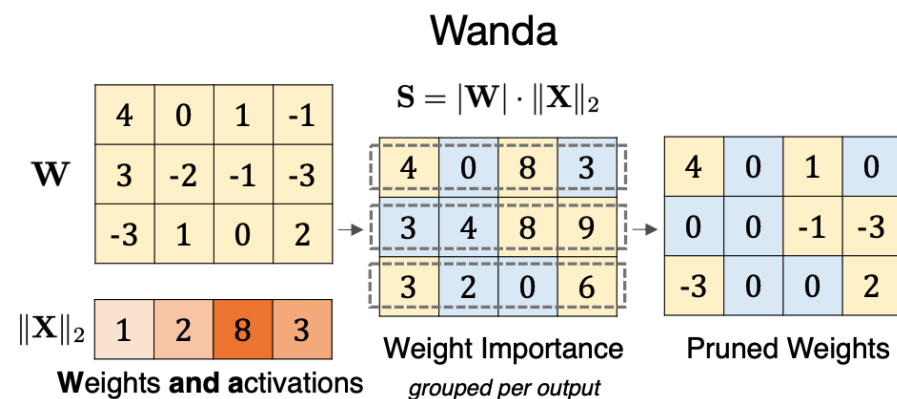
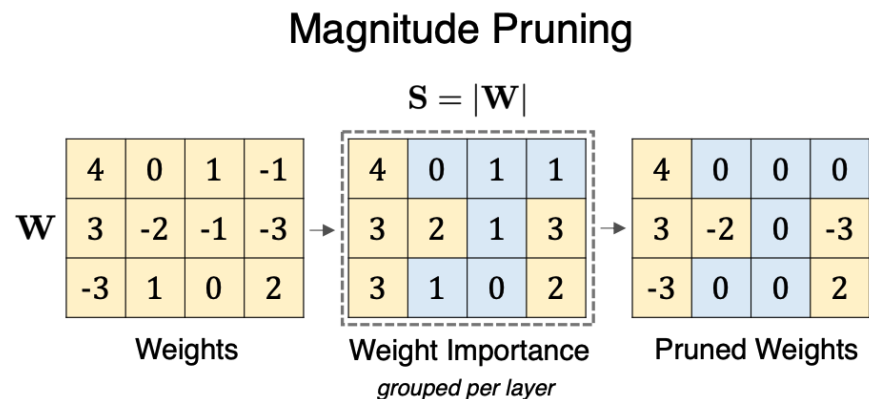
- Auto Sensitivity Scan for Pruning -> Finds Best Sparsity Ratio for Pruning [Least Performance Degradation and Max Performance]
- Supports CWP, GMP Pruning. Room for WANDA, GPTQ, etc..
- QAT
- Auto-Layer Fusion

sconce v1.1

Altruism is all you need !!! Don't just be self attentive

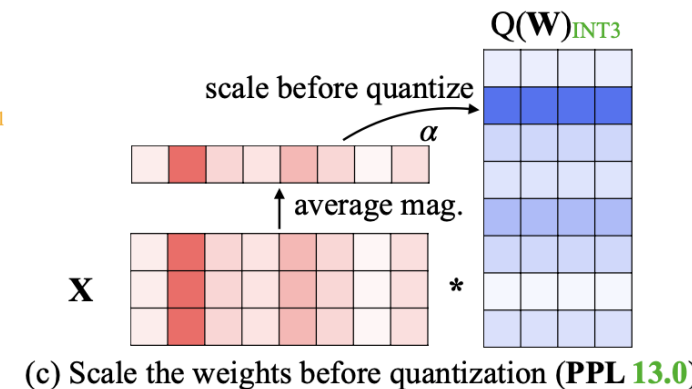
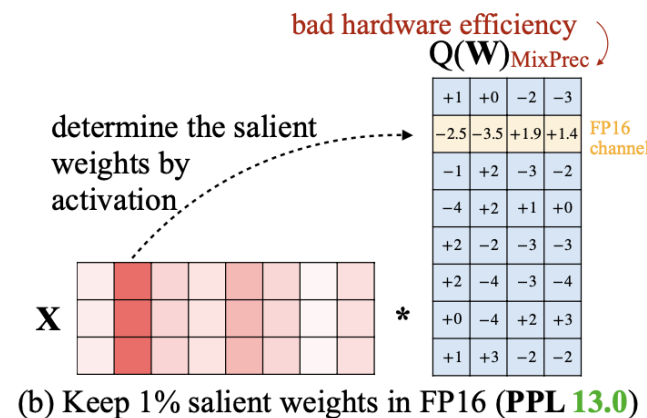
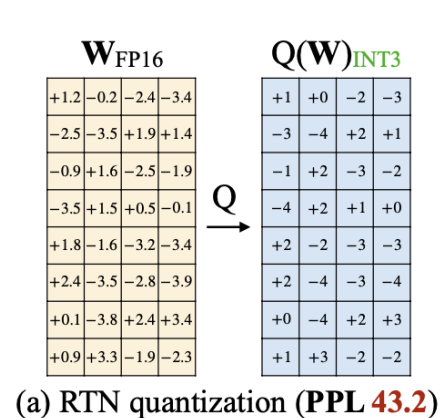
A SIMPLE AND EFFECTIVE PRUNING APPROACH FOR LARGE LANGUAGE MODELS

Mingjie Sun^{1*} Zhuang Liu^{2*} Anna Bair¹ J. Zico Kolter^{1,3}
¹Carnegie Mellon University ²Meta AI Research ³Bosch Center for AI



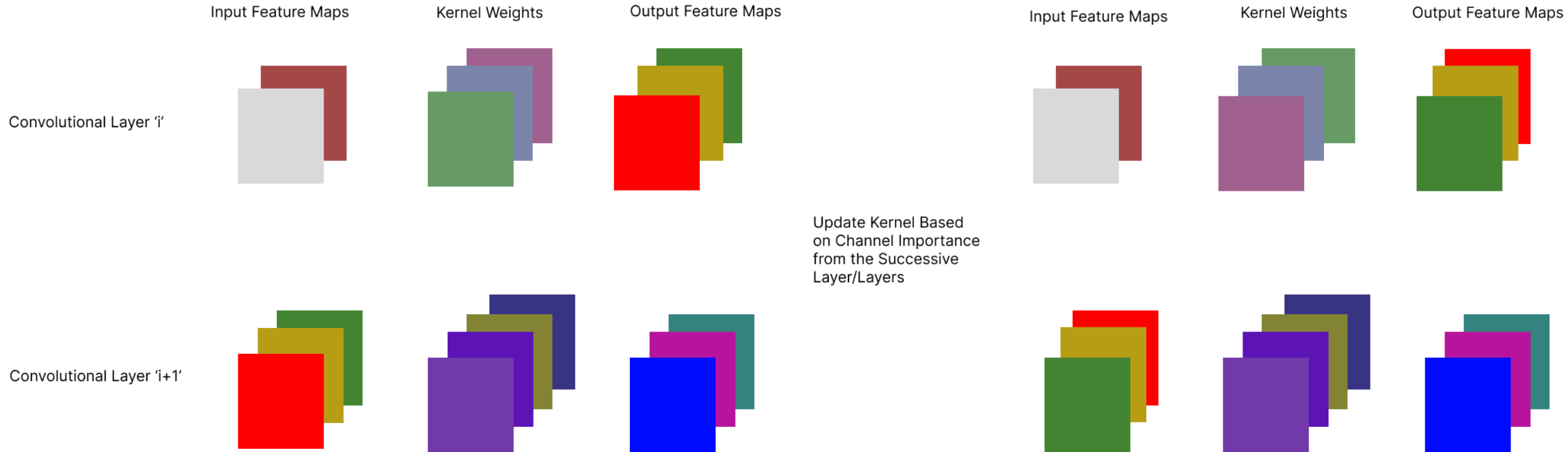
AWQ: Activation-aware Weight Quantization for LLM Compression and Acceleration

Ji Lin^{1*} Jiaming Tang^{1,2*} Haotian Tang¹ Shang Yang¹ Xingyu Dang³ Chuang Gan¹ Song Han¹
¹MIT ²SJTU ³Tsinghua University
<https://github.com/mit-han-lab/llm-awq>



- WANDA: <https://github.com/locuslab/wanda>
- AWQ: <https://github.com/mit-han-lab/llm-awq>

Make Kernels Aware of the Future Kernel Spaces

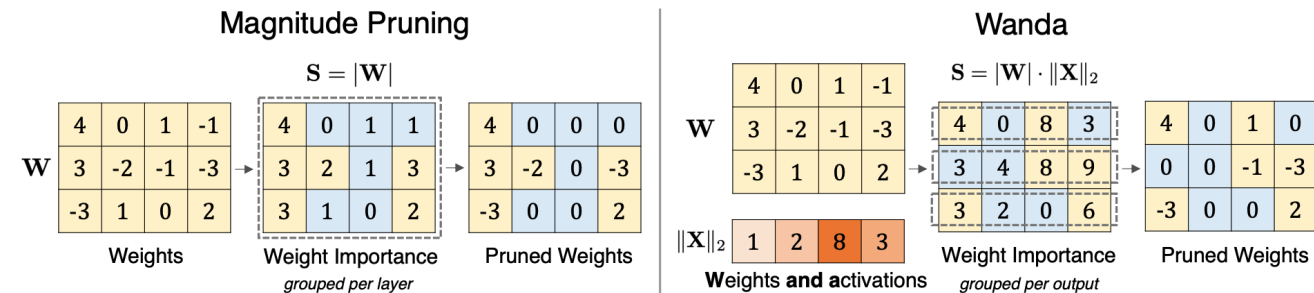


Code: <https://github.com/satabios/sconce/blob/main/tutorials/Pruning.ipynb>

Citations:

- EIE: <https://arxiv.org/abs/1602.01528>
- Multi-scale channel importance sorting and spatial attention mechanism for retinal vessels segmentation
- EACP: An effective automatic channel pruning for neural networks

Channel-Based Activation Aware Pruning



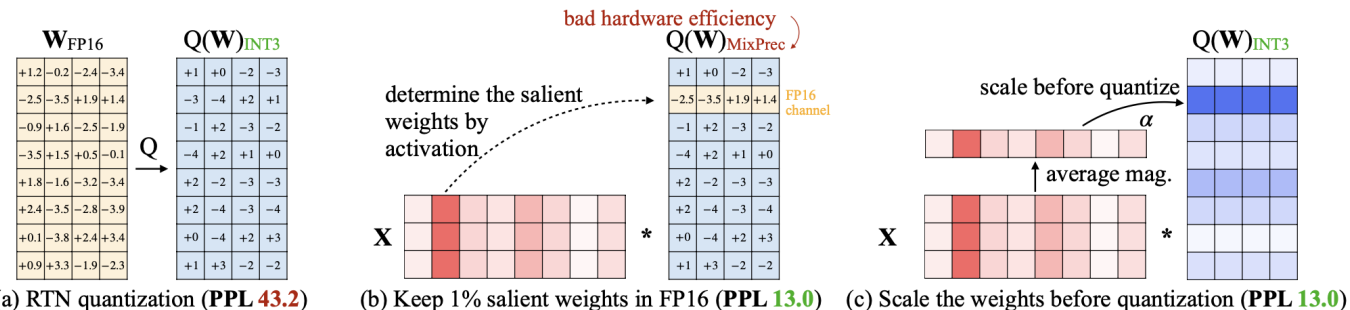
But Channel Wise!!

| Method | Weight Update | Calibration Data | Pruning Metric S_{ij} | Complexity |
|-----------|---------------|------------------|---|--------------------------|
| Magnitude | ✗ | ✗ | $ \mathbf{W}_{ij} $ | $O(1)$ |
| SparseGPT | ✓ | ✓ | $[\mathbf{W} ^2 / \text{diag}[(\mathbf{X}\mathbf{X}^T + \lambda \mathbf{I})^{-1}]]_{ij}$ | $O(d_{\text{hidden}}^3)$ |
| Wanda | ✗ | ✓ | $ \mathbf{W}_{ij} \cdot \ \mathbf{X}_j\ _2$ | $O(d_{\text{hidden}}^2)$ |

| Method | Weight Update | Sparsity | LLaMA | | | | LLaMA-2 | | |
|-----------|---------------|----------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | | | 7B | 13B | 30B | 65B | 7B | 13B | 70B |
| Dense | - | 0% | 59.99 | 62.59 | 65.38 | 66.97 | 59.71 | 63.03 | 67.08 |
| Magnitude | ✗ | 50% | 46.94 | 47.61 | 53.83 | 62.74 | 51.14 | 52.85 | 60.93 |
| SparseGPT | ✓ | 50% | 54.94 | 58.61 | 63.09 | 66.30 | 56.24 | 60.72 | 67.28 |
| Wanda | ✗ | 50% | 54.21 | 59.33 | 63.60 | 66.67 | 56.24 | 60.83 | 67.03 |
| Magnitude | ✗ | 4:8 | 46.03 | 50.53 | 53.53 | 62.17 | 50.64 | 52.81 | 60.28 |
| SparseGPT | ✓ | 4:8 | 52.80 | 55.99 | 60.79 | 64.87 | 53.80 | 59.15 | 65.84 |
| Wanda | ✗ | 4:8 | 52.76 | 56.09 | 61.00 | 64.97 | 52.49 | 58.75 | 66.06 |
| Magnitude | ✗ | 2:4 | 44.73 | 48.00 | 53.16 | 61.28 | 45.58 | 49.89 | 59.95 |
| SparseGPT | ✓ | 2:4 | 50.60 | 53.22 | 58.91 | 62.57 | 50.94 | 54.86 | 63.89 |
| Wanda | ✗ | 2:4 | 48.53 | 52.30 | 59.21 | 62.84 | 48.75 | 55.03 | 64.14 |

- Register hooks to fetch O/P Feature Maps
- Run through a Calibration Dataset
- Prune Channels(Kernel Weights) Based on Activations

Activation Aware - QAT



$$\mathbf{s}^* = \arg \min_{\mathbf{s}} \mathcal{L}(\mathbf{s}), \quad \mathcal{L}(\mathbf{s}) = \|\mathbf{Q}(\mathbf{W} \cdot \mathbf{s})(\mathbf{s}^{-1} \cdot \mathbf{X}) - \mathbf{W}\mathbf{X}\| \quad (3)$$

$$\mathbf{Q} \left(\begin{array}{cccc} \mathbf{W} & & & \\ +1.2 & -0.2 & -2.4 & -3.4 \\ -2.5 & -3.5 & +1.9 & +1.4 \\ -0.9 & +1.6 & -2.5 & -1.9 \\ -3.5 & +1.5 & +0.5 & -0.1 \\ +1.8 & -1.6 & -3.2 & -3.4 \\ +2.4 & -3.5 & -2.8 & -3.9 \\ +0.1 & -3.8 & +2.4 & +3.4 \\ +0.9 & +3.3 & -1.9 & -2.3 \end{array} \begin{array}{l} \times 1 \\ \times 2 \\ \times 1 \\ \times 1 \\ \times 1 \\ \times 1 \\ \times 1 \\ \times 1 \end{array} \right) \xrightarrow{\text{fuse to previous op}} \mathbf{W}\mathbf{X} \rightarrow \mathbf{Q}(\mathbf{W} \cdot \mathbf{s})(\mathbf{s}^{-1} \cdot \mathbf{X})$$

But Channel Wise!!

- Apply Scaling on Feature Maps based on Activation Saliency
- QAT

| PPL↓ | | Llama-2 | | | LLaMA | | | |
|--------------|--------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | | 7B | 13B | 70B | 7B | 13B | 30B | 65B |
| FP16 | - | 5.47 | 4.88 | 3.32 | 5.68 | 5.09 | 4.10 | 3.53 |
| INT3 g128 | RTN | 6.66 | 5.52 | 3.98 | 7.01 | 5.88 | 4.88 | 4.24 |
| | GPTQ | 6.43 | 5.48 | 3.88 | 8.81 | 5.66 | 4.88 | 4.17 |
| | GPTQ-R | 6.42 | 5.41 | 3.86 | 6.53 | 5.64 | 4.74 | 4.21 |
| | AWQ | 6.24 | 5.32 | 3.74 | 6.35 | 5.52 | 4.61 | 3.95 |
| INT4 g128 | RTN | 5.73 | 4.98 | 3.46 | 5.96 | 5.25 | 4.23 | 3.67 |
| | GPTQ | 5.69 | 4.98 | 3.42 | 6.22 | 5.23 | 4.24 | 3.66 |
| | GPTQ-R | 5.63 | 4.99 | 3.43 | 5.83 | 5.20 | 4.22 | 3.66 |
| | AWQ | 5.60 | 4.97 | 3.41 | 5.78 | 5.19 | 4.21 | 3.62 |

Complete Flow

- Sort Channels Based on Successive Channels
- Activation Aware Pruning (WANDA- like)
- Activation Aware Quantization (AWQ- like)

Possible Additions

- **Layer-Wise Neural Network Compression via Layer Fusion:**
<https://proceedings.mlr.press/v157/o-neill21a/o-neill21a.pdf>
- **Layer-Selective Rank Reduction:**
<https://github.com/pratyushasharma/laser>