

Evolutionary Optimization of Convolutional Neural Networks

Vojtěch Čoupek

Supervisor: prof. Ing. Lukáš Sekanina, Ph.D.

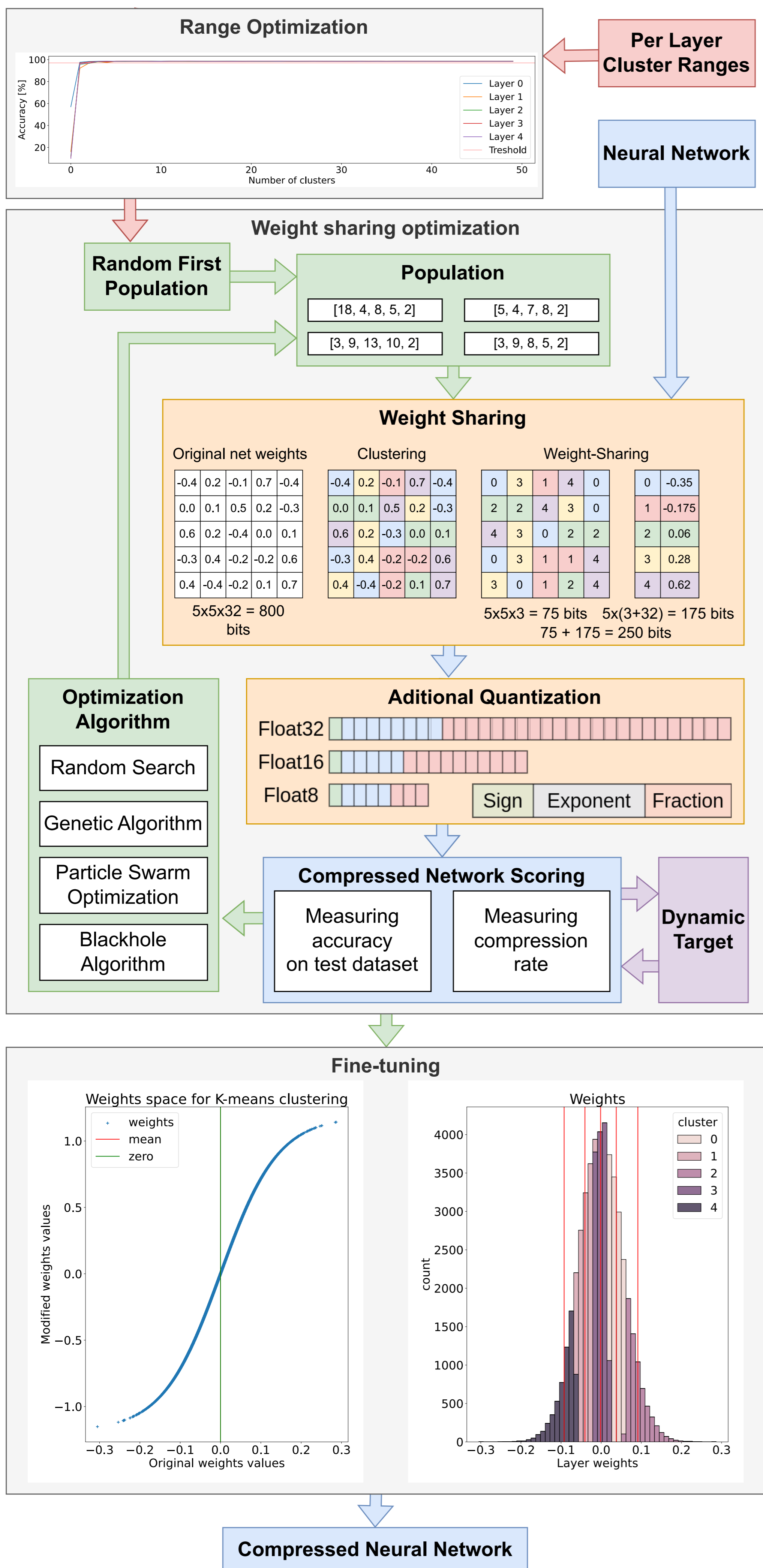


Figure 1: Weight-Sharing optimization computation flow

Weight-Sharing optimization

- Compressing neural network while retaining accuracy
 - Smaller weight memory → Shorter read time → Lower energy consumption
- 1 Range optimization – iterative sharing of layers by a given number of clusters. If the accuracy drops below the threshold, it is removed from the search space.
 - 2 Generate first population. Each member wraps a list of cluster numbers – position in the list refers to the layer in the model.
 - 3 Weight-Sharing – for each member it performs K-means clustering layer by layer. Weights are replaced by the cluster means. Additional Quantization to the means can be applied.
 - 4 Compressed network scoring – For each member, compression rate (CR) and accuracy (ACC) is computed. If CR or ACC of any member surpasses the ACC_{target} or CR_{target} values, they are updated, then the fitness is computed by equation 1.
 - 5 A new population is generated by the chosen optimization algorithm and the process is repeated from step 2 for the given number of iterations.
 - 6 The best found solution is fine-tuned by modulating K-means clustering space adding the second dimension using equation 3 – Searching for the best spread and focus values for each layer is the last step.

$$Fit = \frac{1}{\sqrt{\left(1 - \frac{ACC_{current}}{ACC_{target}}\right)^2 + \left(1 - \frac{CR_{current}}{CR_{target}}\right)^2}} \quad (1)$$

$$CR = \frac{n_w \cdot bits_{weight}}{n_w \cdot bits_{key} + n_k \cdot (bits_{red_weight} + bits_{key})} \quad (2)$$

$$y = spread \cdot (\max(w) - \min(w)) \cdot \tanh(\text{focus} \cdot (x - \text{mean}(w))) \quad (3)$$

Results – Le-Net-5

Activation	Quantization	Compress.	Acc. [%]	Acc. ft [%]	Acc. loss [%]
ReLU	float32	12.7529	97.56	97.80	0.66
	float16	14.1083	97.54	97.66	0.80
	float8	12.7291	97.48	97.78	0.68
Tanh	float32	19.9433	97.66	98.04	0.60
	float16	20.5451	97.68	98.02	0.62
	float8	19.7939	97.86	98.04	0.60

Table 1: Results of weight compression on Le-Net-5 with MNIST dataset

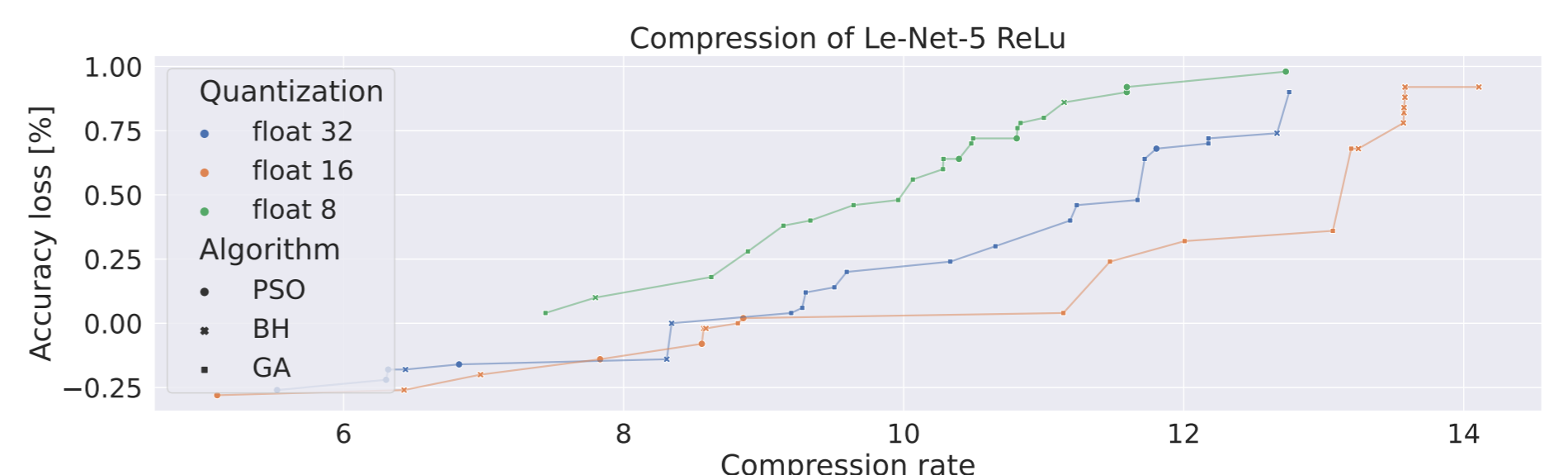


Figure 2: Different additional Quantizations on Weight-Sharing – Le-Net-5 ReLu

Results – Mobilenet_v2

Quantization	Compression	Top-1 Acc. [%]	Top-1 Acc. loss [%]
float32	4,8605	80,75	2,26
float16	4,7617	80,53	2,47
float8	–	–	–

Table 2: Results of weight compression on Mobilenet_v2 with Imagenette dataset