

Figure 4: Genetic Loss Optimization (GLO) overview. A genetic algorithm constructs candidate loss functions as trees. The best loss functions from this set then has its coefficients optimized using CMA-ES. GLO loss functions are able to train models more quickly and more accurately.

245 **Initial population, recombination, and mutation**

246 The initial population is composed of randomly generated trees with a maximum depth of 2. Re-
 247 cursively starting from the root, nodes are randomly chosen from the allowable operator and leaf
 248 nodes using a weighting (where $\log(\dots)$, x , y are three-times as likely and $\sqrt{\dots}$ is two-times as likely
 249 as $+$, $*$, $-$, \div , 1 , -1), this can impart a bias and prevent, for example, the integer 1 from occurring
 250 too frequently. The genetic algorithm has a population size of 80, incorporates elitism with 6 elites
 251 per generation, and uses roulette-sampling.

252 Recombination is accomplished by randomly splicing two trees together. For a given pair of parent
 253 trees, a random element is chosen in each as a crossover point. The two subtrees, whose roots are the
 254 two crossover points, are then swapped with each other. Both resultant trees become part of the next
 255 generation. Recombination occurs with a probability of 80%.

256 To introduce variation into the population, the genetic algorithm has the following mutations, applied
 257 in a bottom-up fashion: integer scalar nodes are incremented or decremented with a 5% probability,
 258 nodes are replaced with a weighted-random node with the same number of children with a 5%
 259 probability, nodes (and their children) are deleted and replaced with a weighted-random leaf node
 260 with a $5\% * 50\% = 2.5\%$ probability, leaf nodes are deleted and replaced with a weighted-random
 261 element (and weighted-random leaf children if necessary) with a $5\% * 50\% = 2.5\%$ probability.

262 **MNIST experimental setup**

263 The first target task used for evaluation was the MNIST Handwritten Digits dataset [I5], a widely
 264 used dataset where the goal is to classify 28×28 pixel images as one of ten digits. The MNIST
 265 dataset has 55,000 training samples, 5,000 validation samples, and 10,000 testing samples.

266 A simple CNN architecture with the following layers is used: (1) 5×5 convolution with 32 filters,
 267 (2) 2×2 stride-2 max-pooling, (3) 5×5 convolution with 64 filters, (4) 2×2 stride-2 max-pooling,
 268 (5) 1024-unit fully-connected layer, (6) a dropout layer [I0] with 40% dropout probability, and (7) a
 269 softmax layer. ReLU [I9] activations are used. Training uses stochastic gradient descent (SGD) with
 270 a batch size of 100, a learning rate of 0.01, and, unless otherwise specified, for 20,000 steps.

271 **CIFAR-10 experimental setup**

272 To further validate GLO, the more challenging CIFAR-10 dataset [I3] (a popular dataset of small,
 273 color photographs in ten classes) was used as a medium to test the transferability of loss functions
 274 found on a different domain. CIFAR-10 consists of 50,000 training samples, and 10,000 testing
 275 samples.

276 A simple CNN architecture, taken from [6] (and itself inspired by AlexNet [I4]), with the following
 277 layers is used: (1) 5×5 convolution with 64 filters and ReLU activations, (2) 3×3 max-pooling

278 with a stride of 2, (3) local response normalization [14] with $k = 1, \alpha = 0.001/9, \beta = 0.75$,
 279 (4) 5×5 convolution with 64 filters and ReLU activations, (5) local response normalization with
 280 $k = 1, \alpha = 0.001/9, \beta = 0.75$, (6) 3×3 max-pooling with a stride of 2, (7) 384-unit fully-connected
 281 layer with ReLU activations, (8) 192-unit fully-connected, linear layer, and (9) a softmax layer.

282 Inputs to the network are sized $24 \times 24 \times 3$, rather than $32 \times 32 \times 32$ as provided in the dataset;
 283 this enables more sophisticated data augmentation. To force the network to better learn spatial
 284 invariance, random 24×24 croppings are selected from each full-size image, which are randomly
 285 flipped longitudinally, randomly lightened or darkened, and their contrast is randomly perturbed.
 286 Furthermore, to attain quicker convergence, an image’s mean pixel value and variance are subtracted
 287 and divided, respectively, from the whole image during training and evaluation. CIFAR-10 networks
 288 were trained with SGD, L^2 regularization with a weight decay of 0.004, a batch size of 1024, and an
 289 initial learning rate of 0.05 that decays by a factor of 0.1 every 350 epochs.

290 Implementation details

291 Due to the large number of partial training sessions that are needed for both the discovery and
 292 optimization phases, training is distributed across the network to a cluster of dedicated machines that
 293 use Condor [27] for scheduling. Each machine in this cluster has one NVIDIA GeForce GTX Titan
 294 Black GPU and two Intel Xeon E5-2603 (4 core) CPUs running at 1.80GHz with 8GB of memory.
 295 Training itself is implemented with TensorFlow [1] in Python. The primary components of GLO (i.e.,
 296 the genetic algorithm and CMA-ES) are implemented in Swift. These components run centrally on
 297 one machine and asynchronously dispatch work to the Condor cluster over SSH. Code is available at:
 298 http://bit.ly/neurips2019_GLO_code

299 Analysis

300 This section presents a symbolic analysis of the Baikal loss function, followed by experiments that
 301 attempt to elucidate why Baikal works better than the cross-entropy loss. A likely explanation is that
 302 Baikal results in implicit regularization.

303 Binary classification

304 Loss functions used on the MNIST dataset, being a 10-dimensional classification problem, are
 305 difficult to plot and visualize graphically. In this section, loss functions are analyzed in the context of
 306 binary classification; where $n = 2$, the Baikal loss expands to:

$$\mathcal{L}_{\text{Baikal2D}} = -\frac{1}{2} \left(\log(y_0) - \frac{x_0}{y_0} + \log(y_1) - \frac{x_1}{y_1} \right) \propto -\log(y_0) + \frac{x_0}{y_0} - \log(1 - y_0) + \frac{1 - x_0}{1 - y_0}, \quad (1)$$

307 since vectors x and y sum to 1, by consequence of being passed through a softmax function, for
 308 binary classification $x = [x_0, 1 - x_0]$ and $y = [y_0, 1 - y_0]$. This constraint simplifies the binary
 309 Baikal loss to a function of two variables (x_0 and y_0). This same methodology can be applied to the
 310 cross-entropy loss and BaikalCMA.

311 In practice, true labels are assumed to be correct
 312 with certainty, thus, x_0 is equal to either 0 or
 313 1. The specific case where $x_0 = 1$ is plotted in
 314 Figure 5 for the cross-entropy loss, Baikal, and
 315 BaikalCMA. The cross-entropy loss is shown to
 316 be monotonically decreasing, while Baikal and
 317 BaikalCMA counterintuitively show an increase
 318 in the loss value as the predicted label, y_0 , ap-
 319 proaches the true label x_0 . Section 5 provides
 320 reasoning for this unusual phenomenon.

321 As also seen in Figure 5, the minimum for the
 322 Baikal loss where $x = 1$ lies around 0.71, while
 323 the minimum for the BaikalCMA loss where
 324 $x = 1$ lies around 0.77. This, along with the
 325 more pronounced slope around $y = 0.5$ is likely

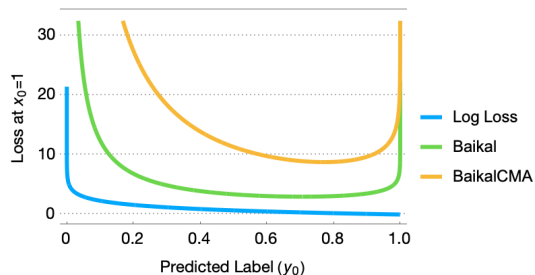


Figure 5: Binary classification loss functions at $x_0 = 1$. Correct predictions lie on the right side of the graph, and vice versa. The log loss is shown to be monotonically decreasing, while Baikal and BaikalCMA present counterintuitive, sharp increases in loss as predictions, approach the true label.

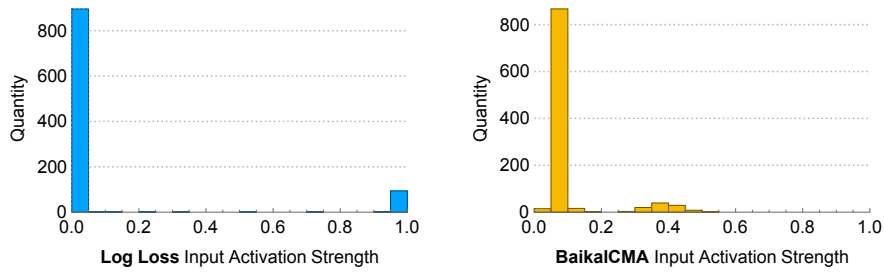


Figure 6: Loss function input activation strength histograms for cross-entropy loss and BaikalCMA. The peaks are likely shifted with BaikalCMA due to implicit regularization. These histograms match those from a network trained with a confidence regularizer [21].

326 a reason why BaikalCMA performs better than
 327 Baikal.

328 **Implicit regularization**

329 The Baikal and BaikalCMA loss functions are unusual in that they are not monotonically decreasing
 330 (see the previous section for more details). At first glance, this behavior may seem undesirable;
 331 however, this may be an advantageous trait that implicitly provides a form of regularization (enabling
 332 better generalization). This is strongly supported by [21], where researchers built a confidence
 333 regularizer, on top of cross-entropy loss, that penalizes low entropy prediction distributions. The
 334 bimodal distribution of output probabilities that the researchers found on MNIST is nearly identical
 335 to that which can be found on a network trained with Baikal or BaikalCMA.

336 Histograms of the output probability distributions of network trained with the cross-entropy loss and
 337 BaikalCMA on the test dataset, after 15,000 steps of training on MNIST, are shown in Figure 6. Note
 338 that the abscissae in Figures 5 and 6 correspond with each other, thus one can qualitatively see how
 339 the channel-shaped curves for BaikalCMA may contribute to the shift in histogram peaks.

340 Furthermore, the improved behavior under small-dataset conditions described in Section 4.1 backs
 341 this theory of implicit regularization, since less overfitting was observed when using Baikal and
 342 BaikalCMA.