# Machine Learning for IoT

Homework 2

Matteo Donadio Student ID: s303903 Politecnico di Torino Alessandro Maria Feri' Student ID: s305949 Politecnico di Torino Mattia Sabato Student ID: s305849 Politecnico di Torino

### I. EVOLUTIONARY NEURAL ARCHITECTURE SEARCH

## A. Encoding and Mapping

The task of identifying a lightweight and reliable architecture, along with determining the optimal hyper-parameters for a fast pre-processing and effective training, does not come with a simple solution. Being these aspects greatly intertwined, we opted to conduct a Neural Architecture Search (NAS) utilizing a Genetic Algorithm (GA) to efficiently optimize them in parallel.

| Operation   | Shape                                     | #Parameters                             |
|---|---|---|
| Input   | (61, 8, 1)                                | 0                                       |
| Conv [2, 2, valid, L2]  BatchNorm  RElU                                 | (30, 4, 32)<br>(30, 4, 32)<br>(30, 4, 32) | 128<br>128<br>0                         |
| Conv [2, 2, same, L2]  BatchNorm  RE1U                                  | (15, 2, 32)<br>(15, 2, 32)<br>(15, 2, 32) | 4,096<br>128<br>0                       |
| DW-Conv [3, 1, same, -]  BatchNorm  RE1U                                | (15, 2, 32)<br>(15, 2, 32)<br>(15, 2, 32) | 288<br>128<br>0                         |
| Global AvgPool<br>Dense<br>SoftMax                                      | (32)<br>(2)<br>(2)                        | 0<br>66<br>0                            |
| Total Parameters:<br>Trainable Parameters:<br>Non-trainable Parameters: | 4,962<br>4,770<br>192                     | (19.38KB)<br>(18.63KB)<br>(768.00 Byte) |

TABLE I: Best Performing Architecture

We define a triple x = (Preprocessing, Architecture, Training)as a single vector  $x \in \mathbb{N}^n$ , whose components consist of integer values carrying encoded information, and we denote it as an individual. Each value serves as a gene, representing an encoded tunable hyper-parameter, whereas the ordered sequence of all genes corresponding to an individual forms the genome. For instance, the gene at position 1 may be an integer value representing a multiple of 4, indicating the batch\_size, while the gene at position 12 may be a binary value indicating the presence or absence of a kernel regularization for a convolution. We then define a mapping  $f(x): \mathbb{N}^n \to \mathcal{H}$  which decodes an individual x into a general function h = f(x), wrapping preprocessing, training and architecture into a unified object. The best performing architecture we have identified is presented in Table I, where we report the kernel size, stride, padding and kernel regularizer for convolutions. Table II shows the hyperparameters associated with the architecture's remaining genome, along with their corresponding performances after training for 60 epochs.

#### B. Strategy

Having established the encoding, mapping, and the supported values for the genome, we adopted a  $(\mu, \lambda)$  strategy with  $\mu = 6$  and  $\lambda = 10$ . In each *generation*, 3 individuals are chosen using a roulette-wheel approach, where the probability is proportional to their relative *fitness*. Subsequently, a total of  $\lambda$  offsprings are generated by applying randomly chosen uniform/one-cut/two-cuts crossovers and one-bit mutations with a probability of 0.4. In each generation, the fitness of each individual is lazily computed as the maximum accuracy on the test set over a total of 15 epochs. Notice how a multi-objective approach, where the fitness is alternatively computed as the reciprocal of the size with a certain probability, could yield a *Pareto front* where both the accuracy and the model size are extremely optimized.

|                | Hyper-parameter/ <b>KPI</b>   | Value  |
|----------------|---|--|
| Pre-processing | sampling_rate frame_length frame_step mel_bins lower_freq upper_freq MFCC | 16,000<br>0.032<br>0.016<br>14<br>0<br>4,000                   |
| Training       | batch_size epochs optimizer initial_lr final_lr schedule power_decay      | 16<br>60<br>RMSprop<br>10e-3<br>10e-5<br>PolynomicalDecay<br>2 |
| Performance    | Accuracy (%) TFLite Size (KB) Savings                                     | 99.00<br>21.57<br>57.04  |

TABLE II: Best Hyper-parameters and Performances

## II. RESULTS

Notably, no additional optimization was required, thanks to the thorough self-optimization carried out from the beginning and the careful design of the search spaces. The incorporation of depthwise convolutions, crucial for reducing the number of weights, was restricted to the third layer alone, since their application to other layers proved to be detrimental to accuracy. Moreover, the adoption of Mel-Frequency Cepstral Coefficients, widely utilized in speech processing, has shown to be very effective and able to enhance the accuracy despite a slowdown in the pre-processing. Finally, the use of a GA has incredibly reduced the amount of time needed to find an acceptable solution, avoiding the burden of conducting a grid search and allowing solutions to self-optimize themselves.