

Algorithm 1: Custom floating-point quantization

Input: $MODEL$ as the CNN

Input: E_{size} as the target exponent bit size

Input: M_{size} as the target mantissa bits size

Input: $STDM_{size}$ as the IEEE 754 mantissa bit size

Output: $MODEL$ as the quantized CNN

```
1 foreach layer in  $MODEL$  do
2   if layer is Conv2D or SeparableConv2D then
3     filter, bias  $\leftarrow$  GetWeights(layer)
4     foreach x in filter and bias do
5       sign  $\leftarrow$  GetSign(x)
6       exp  $\leftarrow$  GetExponent(x)
7       fullexp  $\leftarrow$   $2^{E_{size}-1} - 1$  // Get full range value
8       cman  $\leftarrow$  GetCustomMantissa(x,  $M_{size}$ )
9       leftman  $\leftarrow$  GetLeftoverMantissa(x,  $M_{size}$ )
10      if exp < -fullexp then
11        | x  $\leftarrow$  0
12      else
13        if exp > fullexp then
14          | x  $\leftarrow$   $(-1)^{sign} \cdot 2^{fullexp} \cdot (1 + (1 - 2^{-M_{size}}))$ 
15        else
16          if  $2^{STDM_{size}-M_{size}-1} - 1 < leftman$  then
17            | cman  $\leftarrow$  cman + 1 // Above halfway
18            if  $2^{M_{size}} - 1 < cman$  then
19              | cman  $\leftarrow$  0 // Correct mantissa overflow
20              | exp  $\leftarrow$  exp + 1
21            | x  $\leftarrow$   $(-1)^{sign} \cdot 2^{exp} \cdot (1 + cman \cdot 2^{-M_{size}})$ 
22      SetWeights(layer, filter, bias)
```