**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: for layer in MODEL do
 2:     if layer is Conv2D or SeparableConv2D then
 3:         filter, bias ← GetWeights(layer)
 4:         for x in filter and bias do
 5:             sign ← GetSign(x)
 6:             exp ← GetExponent(x)
 7:             fullexp ← 2^(E_size−1) − 1                    ▷ Get full range value
 8:             cman ← GetCustomMantissa(x, M_size)
 9:             leftman ← GetLeftoverMantissa(x, M_size)
10:             if exp < −fullexp then
11:                 x ← 0
12:             else if exp > fullexp then
13:                 x ← (−1)^sign · 2^fullexp · (1 + (1 − 2^(−M_size)))
14:             else
15:                 if 2^(STDM_size−M_size−1) − 1 < leftman then
16:                     cman ← cman + 1                       ▷ Above halfway
17:                     if 2^(M_size) − 1 < cman then
18:                         cman ← 0                          ▷ Correct mantissa overflow
19:                         exp ← exp + 1
20:                     end if
21:                 end if
22:                 x ← (−1)^sign · 2^exp · (1 + cman · 2^(−M_size))
23:             end if
24:         end for
25:         SetWeights(layer, filter, bias)
26:     end if
27: end for
```