# Lossless Compression of Deep Neural Networks (results-presentation)

*Vilhelm Toivonen*

*01.12.2023*

Instructor: Nikita *Belyak*
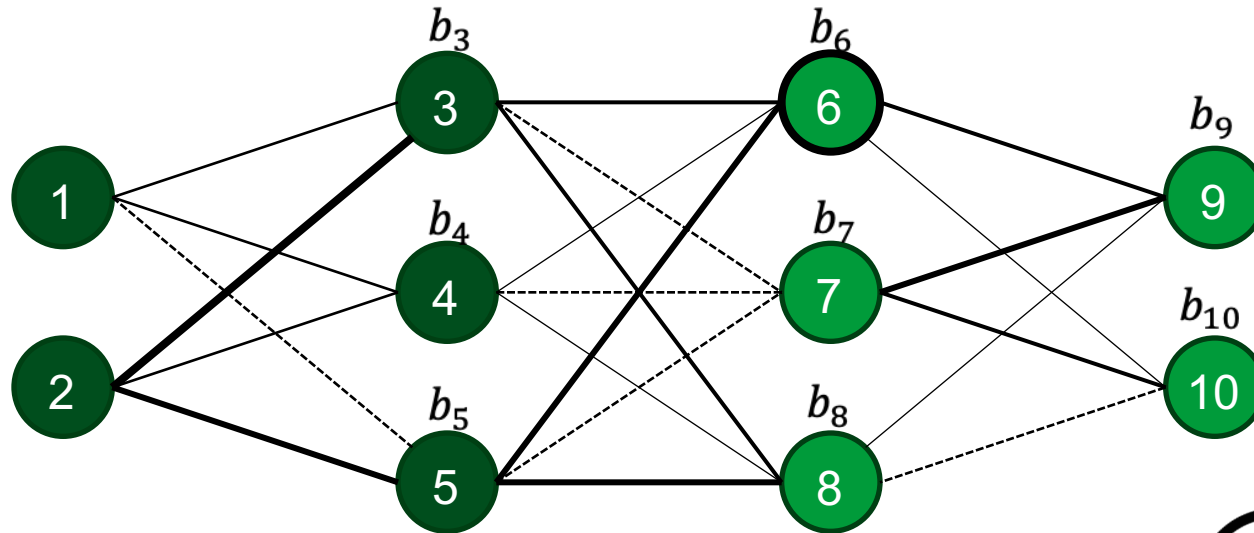
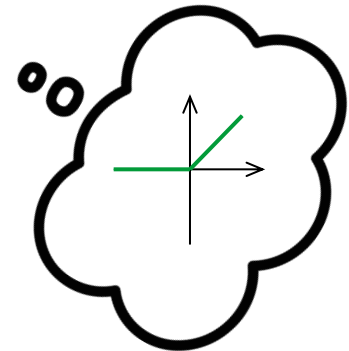Supervisor: Fabricio *Oliveira*

# Background – DNNs



Larger DNNs make calculations more intensive
- Slow forward passes
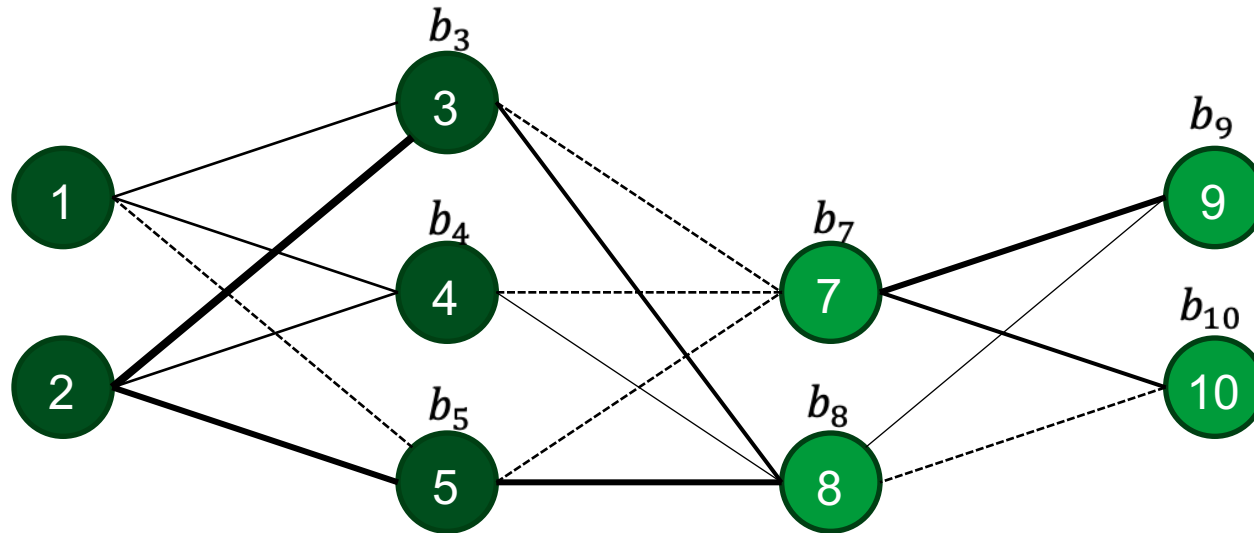- Computationally expensive to create mathematical programming problems
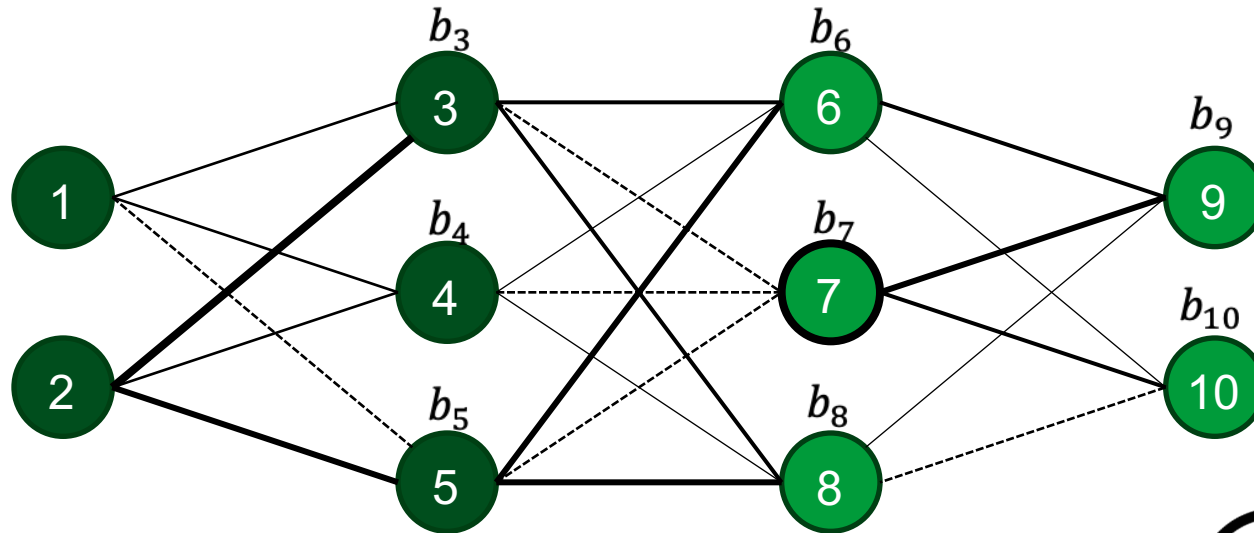
# Example – Upper bound
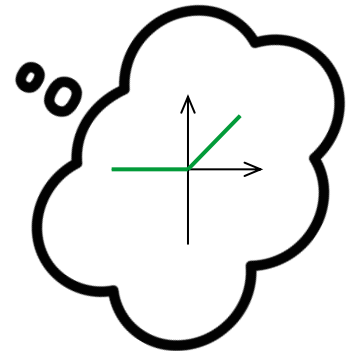


$$l = 2$$
$$i = 1$$
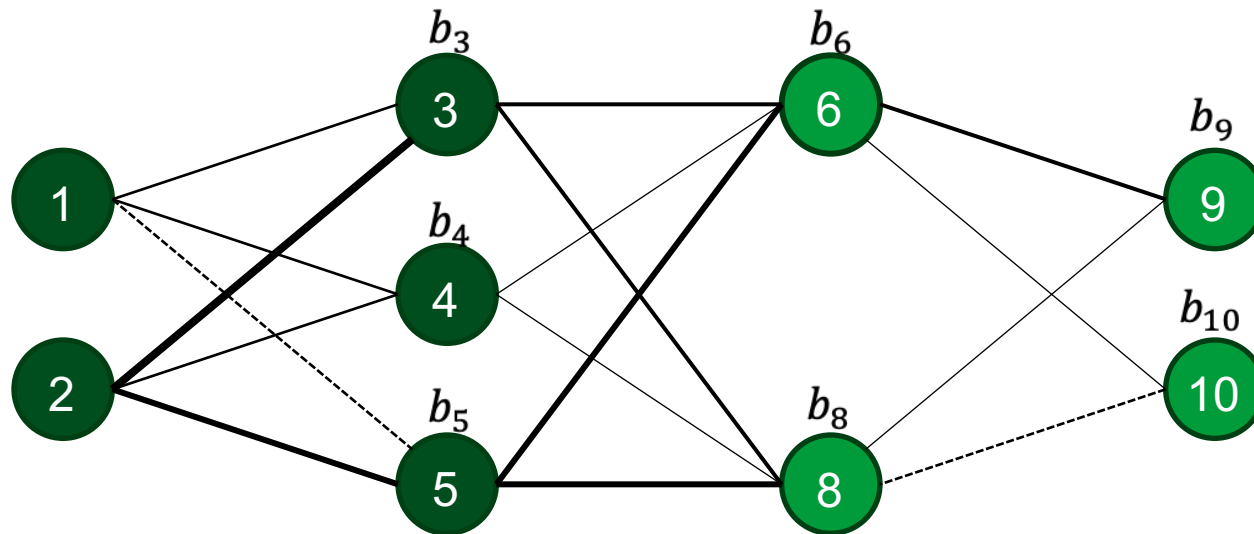$$G_i^l < 0$$

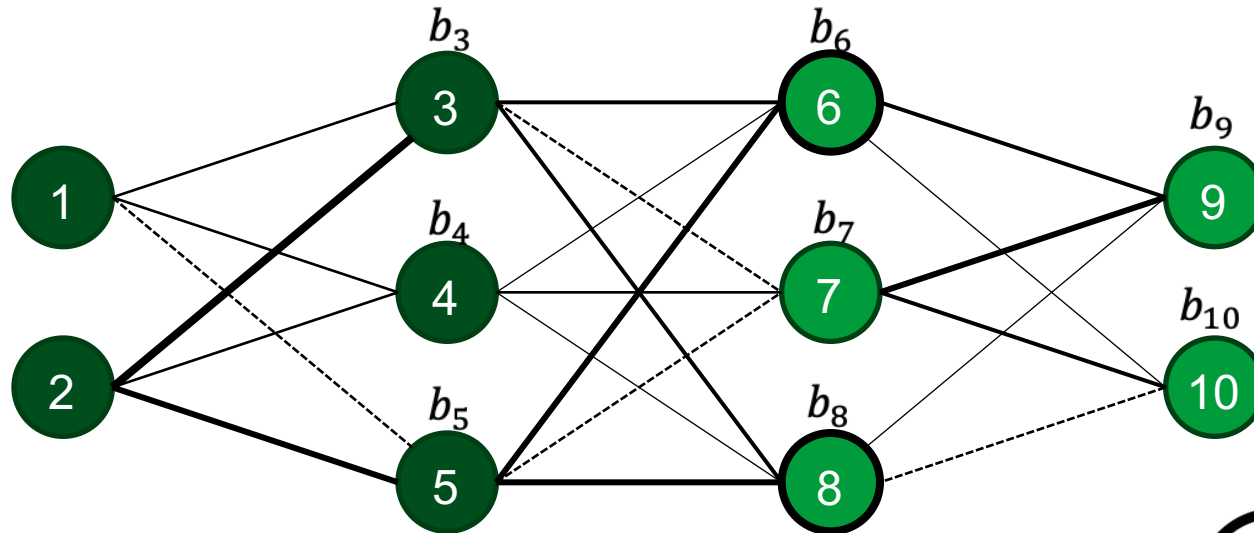# Example – Upper bound

# Example – Zero weights



$$l = 2$$
$$i = 2$$
$$W_i^l = 0$$

# Example – Zero weights

# Example – Linear dependence
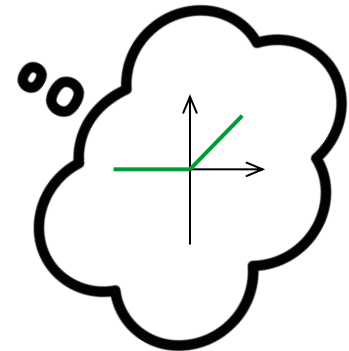


$$l = 2$$
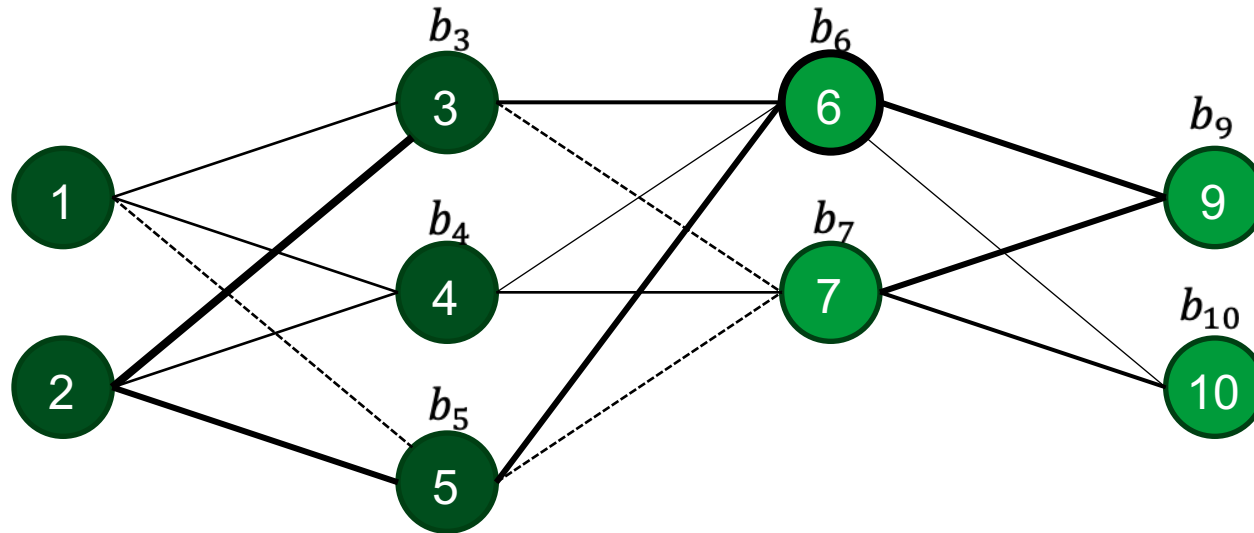$$i = 1$$
$$\bar{G}_i^l > 0$$

$$l = 2$$
$$i = 3$$
$$\bar{G}_i^l > 0$$

# Example – Linear dependence

# Methods

- Train the models using different optimizers and regularizations
- Find the upper bounds $G$ and the lower bounds $\overline{G}$ for each neuron
- Prune the network with the illustrated algorithm
- Check that the output of the pruned network matches the output of the original network, and log results

```
 1: for l ← 1, . . . , L do
 2:     S ← {}                                          ▷ Set of stable units left in layer l
 3:     Unstable ← False                                ▷ If there are unstable units in layer l
 4:     for i ← 1, . . . , n_l do
 5:         if G_i^l < 0 for x ∈ D or W_i^l = 0 then    ▷ Stably inactive, constant output
 6:             if i < n_l or |S| > 0 or Unstable then
 7:                 if W_i^l = 0 and b_i^l > 0 then
 8:                     for j ← 1, . . . , n_{l+1} do
 9:                         b_j^{l+1} ← b_j^{l+1} + w_{ji}^{l+1} b_i^l
10:                     end for
11:                 end if
12:                 Remove unit i from layer l            ▷ Unit i is not necessary
13:             end if
14:         else if Ḡ_i^l > 0 for x ∈ D then            ▷ Stably active
15:             if rank(W_{S∪{i}}^l) > |S| then
16:                 S ← S ∪ {i}                          ▷ Keep unit in the network
17:             else
18:                 Find {α_k}_{k∈S} such that w_i^l = Σ_{k∈S} α_k w_k^l
19:                 for j ← 1, . . . , n_{l+1} do
20:                     for k ∈ S do
21:                         w_{jk}^{l+1} ← w_{jk}^{l+1} + α_k w_{ji}^{l+1}
22:                     end for
23:                     b_j^{l+1} ← b_j^{l+1} + w_{ji}^{l+1} (b_i^l − Σ_{k∈S} α_k b_k^l)
24:                 end for
25:                 Remove unit i from layer l            ▷ Unit i is no longer necessary
26:             end if
27:         else
28:             Unstable ← True
29:         end if
30:     end for
31:     if not Unstable then                             ▷ All units left in layer l are stable
32:         if |S| > 0 then                              ▷ The units left have varying outputs
33:             Create matrix W̄ ∈ ℝ^{n_l×n_{l+1}} and vector b̄ ∈ ℝ^{n_{l+1}}
34:             for i ← 1, . . . , n_{l+1} do
35:                 b̄_i ← b_i^{l+1} + Σ_{k∈S} w_{ik}^{l+1} b_k^l
36:                 for j ← 1, . . . , n_{l−1} do
37:                     w̄_{ij} ← Σ_{k∈S} w_{kj}^l w_{ik}^{l+1}
38:                 end for
39:             end for
40:             Remove layer l; replace parameters in next layer with W̄ and b̄
41:         else                                          ▷ Only unit left in layer l has constant output
42:             Compute output Υ for any input χ ∈ D
43:             (W^{L+1}, b^{L+1}) ← (0, Υ)               ▷ Set constant values in output layer
44:             Remove layers 1 to L and break           ▷ Remove all hidden layers and leave
45:         end if
46:     end if
47: end for
```

# Training the models – variations

Trained with Pytorch, and used the MSE loss function.

Optimizers:

1. Adam
2. AdaDelta
3. SGD
4. SGD with momentum

Model architecture
$[2,1024,512,512,256,1]$
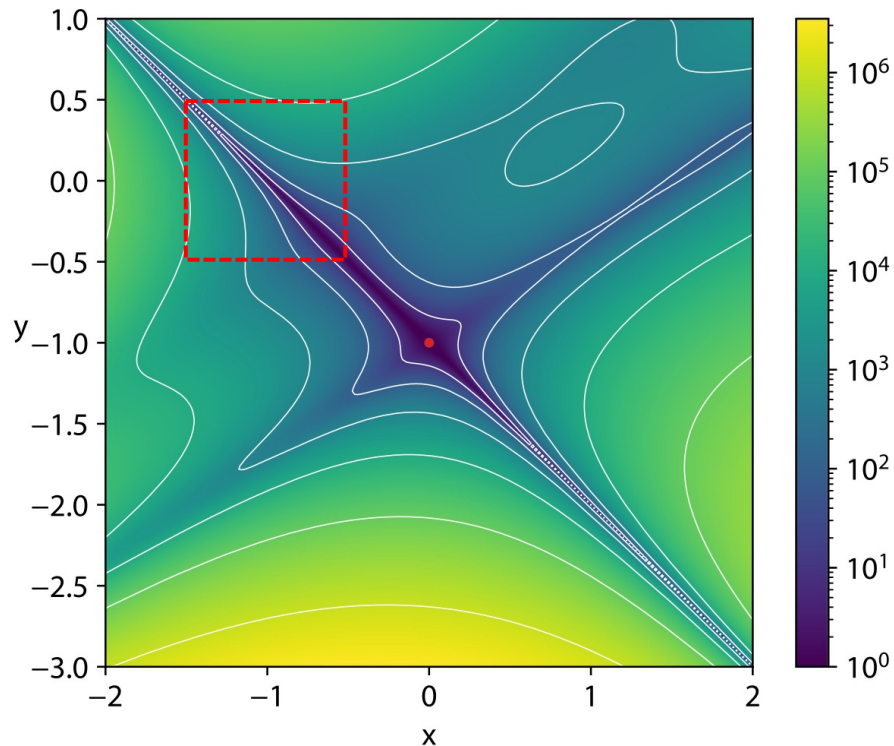
Regularization parameters:

1. L1 (Lasso) regularization, $\lambda_{l_1} \in$ $[0,0.001,0.01,0.1,0.5]$
2. L2 (Ridge) regularization, $\lambda_{l_2} \in$ $[0,0.001,0.01,0.1,0.5]$
3. L1+L2 (Elastic net) regularization, $\left(\lambda_{l_1}, \lambda_{l_2}\right) \in [(\lambda, \lambda), \lambda$ $\in [0.001,0.01,0.1,0.5]$

$\rightarrow 13 \times 4 = 52$ models

Aalto-yliopisto
Perustieteiden
korkeakoulu

Systeemianalyysin
laboratorio

# Data - Goldstein price

$$x \in [-1.5, -0.5]$$
$$y \in [-0.5, 0.5]$$



Training samples: 40,000
Testing samples: 8,000

Aalto-yliopisto
Perustieteiden
korkeakoulu

Systeemianalyysin
laboratorio
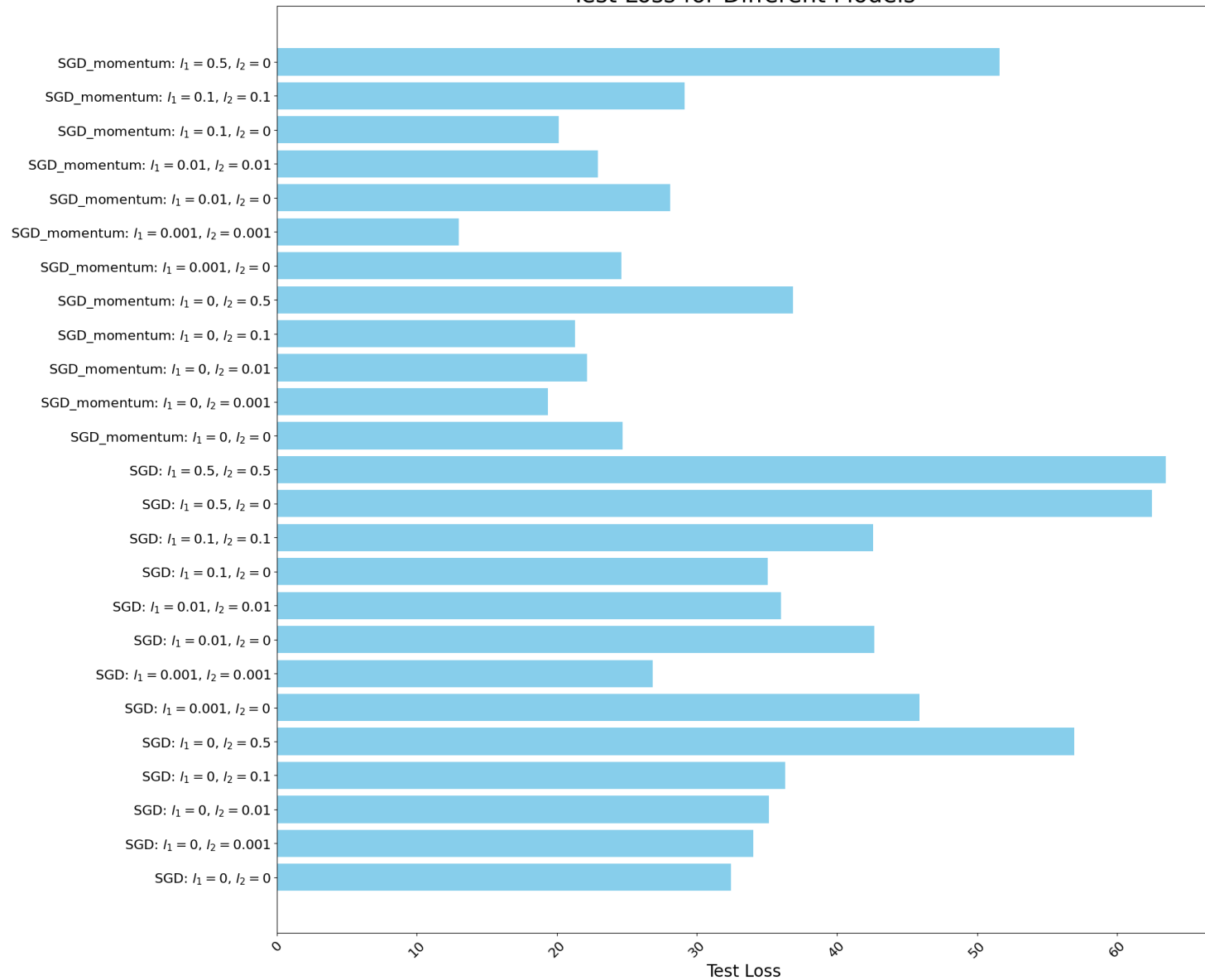
# Training the models – hyperparameters

| Optimizer | Epochs | Training time min/model on T4 | Training time min/model on M2 Max | Learning rate | Learning rate decay $\gamma$ | Gradient norm clipped to |
|---|---|---|---|---|---|---|
| Adam | 600 | 10 | 14 | 0.001 | 0.95 | 5 |
| AdaDelta | 2000 | 31 | 45 | 1.5 | 0.9985 | False |
| SGD | 2000 | 30 | 43 | 0.025 | 0.998 | 5 |
| SGD with momentum | 800 | 13 | 19 | 0.005 | 0.995 | 5 |

# Training the models – losses

# Testing the models



Test Loss for Different Models

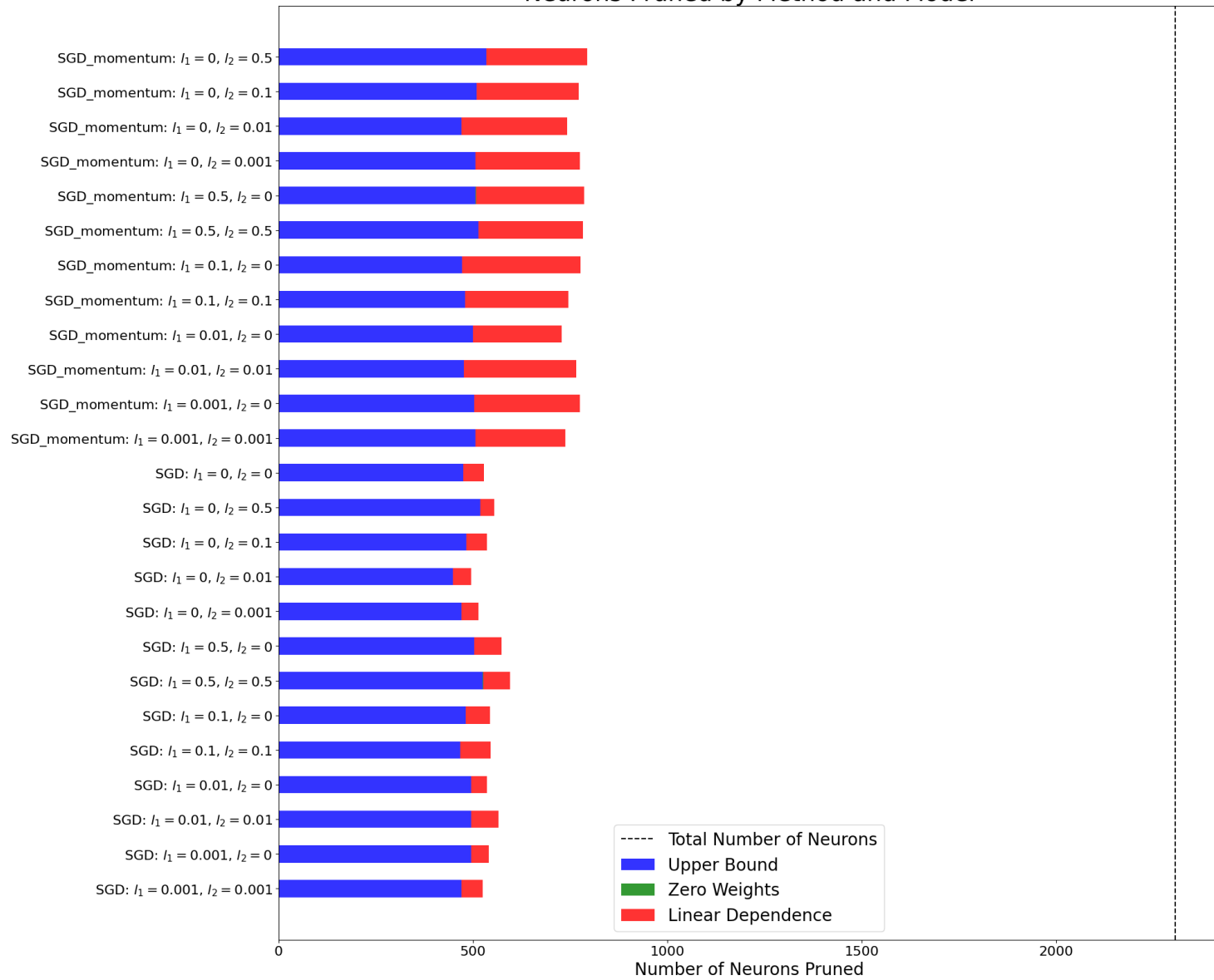# Testing the models



Test Loss for Different Models

# Bounding and pruning the network

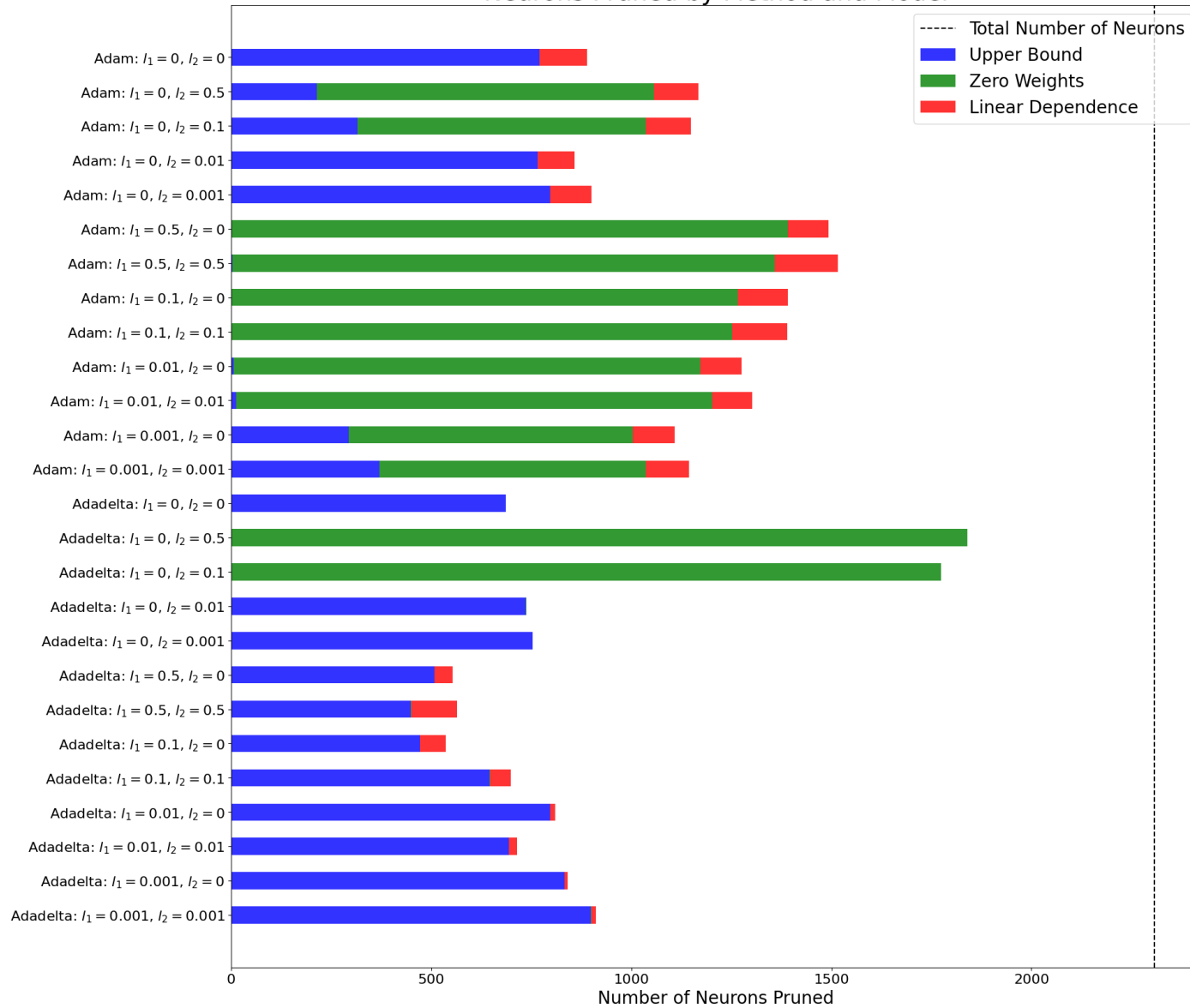Bounds were calculated by computing the LP relaxations for the MIP problems for each neuron.

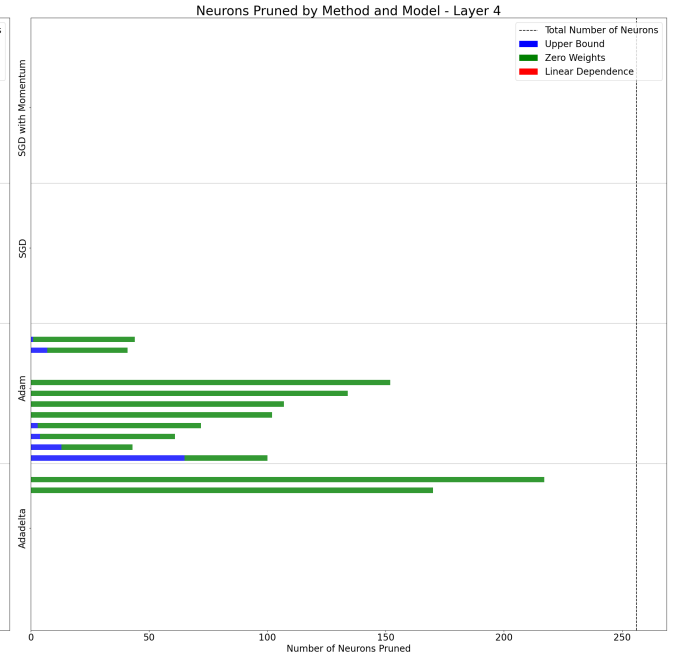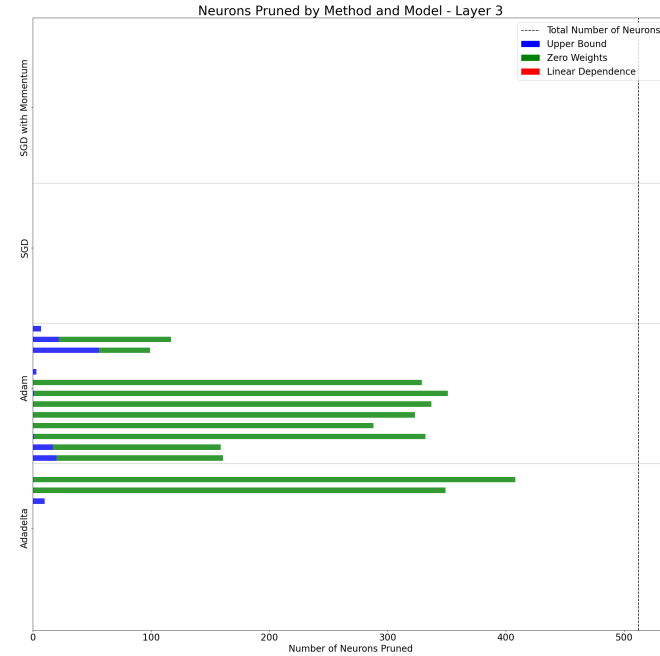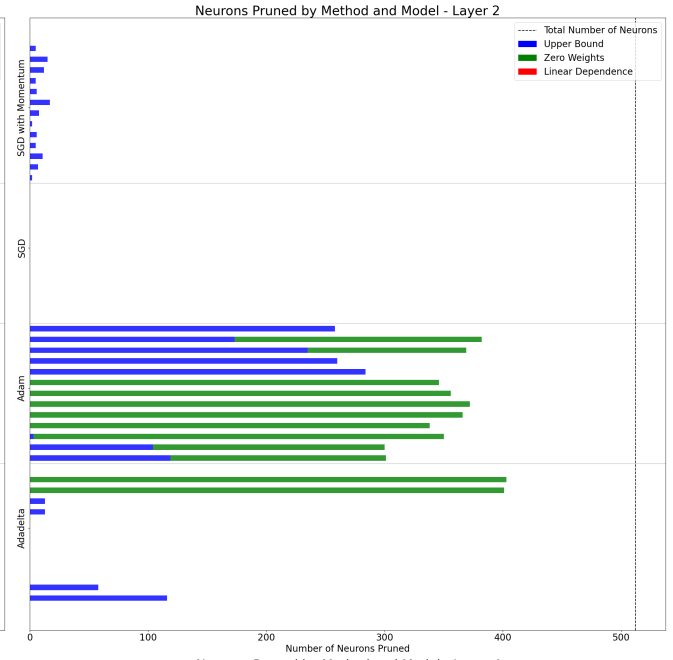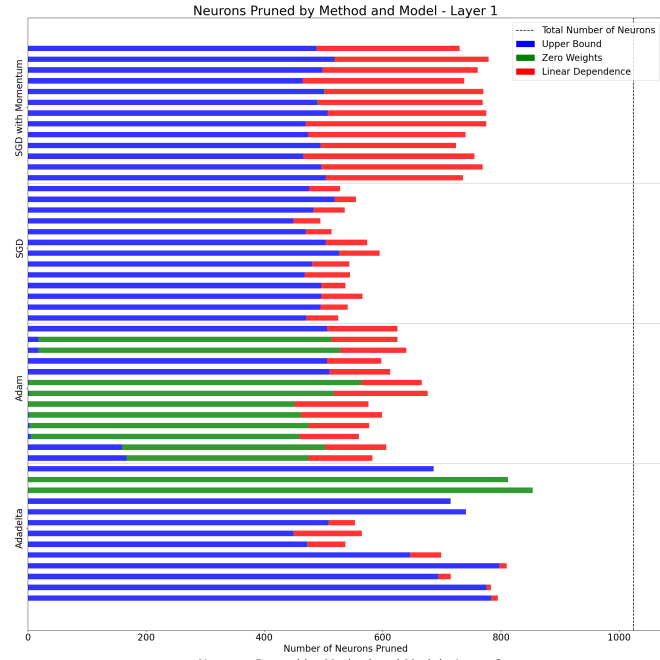For each model, the bounding took $\approx 8$ min.

# Results



Neurons Pruned by Method and Model

# Results



Neurons Pruned by Method and Model

# Results by layer

# Conclusions



The AdaDelta optimizer with $\lambda_{l_2} = 0.5$ or $\lambda_{l_2} = 0.1$ yield the best results.

The Adam optimizer yield consistently good results.

# Limitations

Only trained one model for each combination.

Only tried LP relaxed bounds.

Only considered strict inequalities in upper and lower bound clauses.

Only considered DNNs with the ReLU activation.

Did not consider other architectures, like CNNs, RNNs or attention models.

# Sources and materials

- Thiago Serra, Abhinav Kumar, and Srikumar Ramalingam 2020. Lossless Compression of Deep Neural Networks. Bucknell University and the University of Utah. Usa

- Linkola, J. 2023. Reformulating deep neural networks as mathematical programming problems. Bachelor thesis. Aalto-University. School of Science. Espoo.

- ML_as_MO package (https://github.com/gamma-opt/ML_as_MO). Includes implementations for calculating bounds