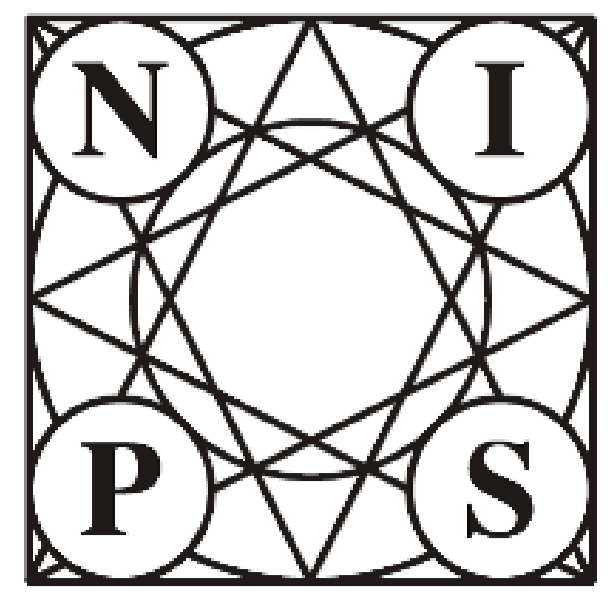




# Model compression as constrained optimization, with application to neural nets

Miguel Á. Carreira-Perpiñán and Yerlan Idelbayev, UC MERCED



## 1 Abstract

Compressing neural nets is an active research problem, given the large size of state-of-the-art nets for tasks such as object recognition, and the computational limits imposed by mobile devices. Firstly, we give a general formulation of model compression as constrained optimization. This makes the problem of model compression well defined and amenable to the use of modern numerical optimization methods. Our formulation includes many types of compression: quantization, low-rank decomposition, pruning, lossless compression and others. Then, we give a general algorithm to optimize this nonconvex problem based on a penalty function (quadratic penalty or augmented Lagrangian) and alternating optimization. This results in a “learning-compression” algorithm, which alternates a learning step of the uncompressed model, independent of the compression type, with a compression step of the model parameters, independent of the learning task. This algorithm is guaranteed to find the best compressed model for the task under standard assumptions. It is simple to implement in existing deep learning toolboxes and efficient, with a runtime comparable to that of training a reference model in the first place.

Supported by NSF IIS-1423515 and NVIDIA GPU donation

## 2 Motivation

- Deep learning very successful in some practical applications: vision, speech, NLP
- Very good classification accuracy if training large, deep neural nets on large datasets (with several GPUs over several days...).
- It is of interest to deploy high-accuracy deep nets on limited-computation devices (mobile phones, IoT, etc.), but this requires compressing the deep net to meet the device’s limitations in memory, runtime, energy, bandwidth, etc.
- Surprisingly high compression rates often possible because deep nets are vastly over-parameterized in practice. It seems this makes it easier to learn a high-accuracy net.
- Lossy compression will degrade accuracy so we want to compress optimally.

## 3 Related work

Many papers in the last 2–3 years, although neural net compression was already studied in the 1980–90s. Some common approaches:

- Direct compression: train a reference net, compress its weights using standard compression techniques. quantization (binarization, low precision), pruning weights/neurons, low-rank, etc. Simple and fast, but it ignores the loss function, so the accuracy deteriorates a lot for high compression rates.
- Embedding the compression mechanism in the backprop training. remove neuron/weight values with low magnitude on the fly binarize or round weight values in the backprop forward pass Often no convergence guarantees.
- Other ad-hoc algorithms proposed for specific compression techniques.
- Multiple compression techniques can be combined.

## 4 Our general formulation for model compression

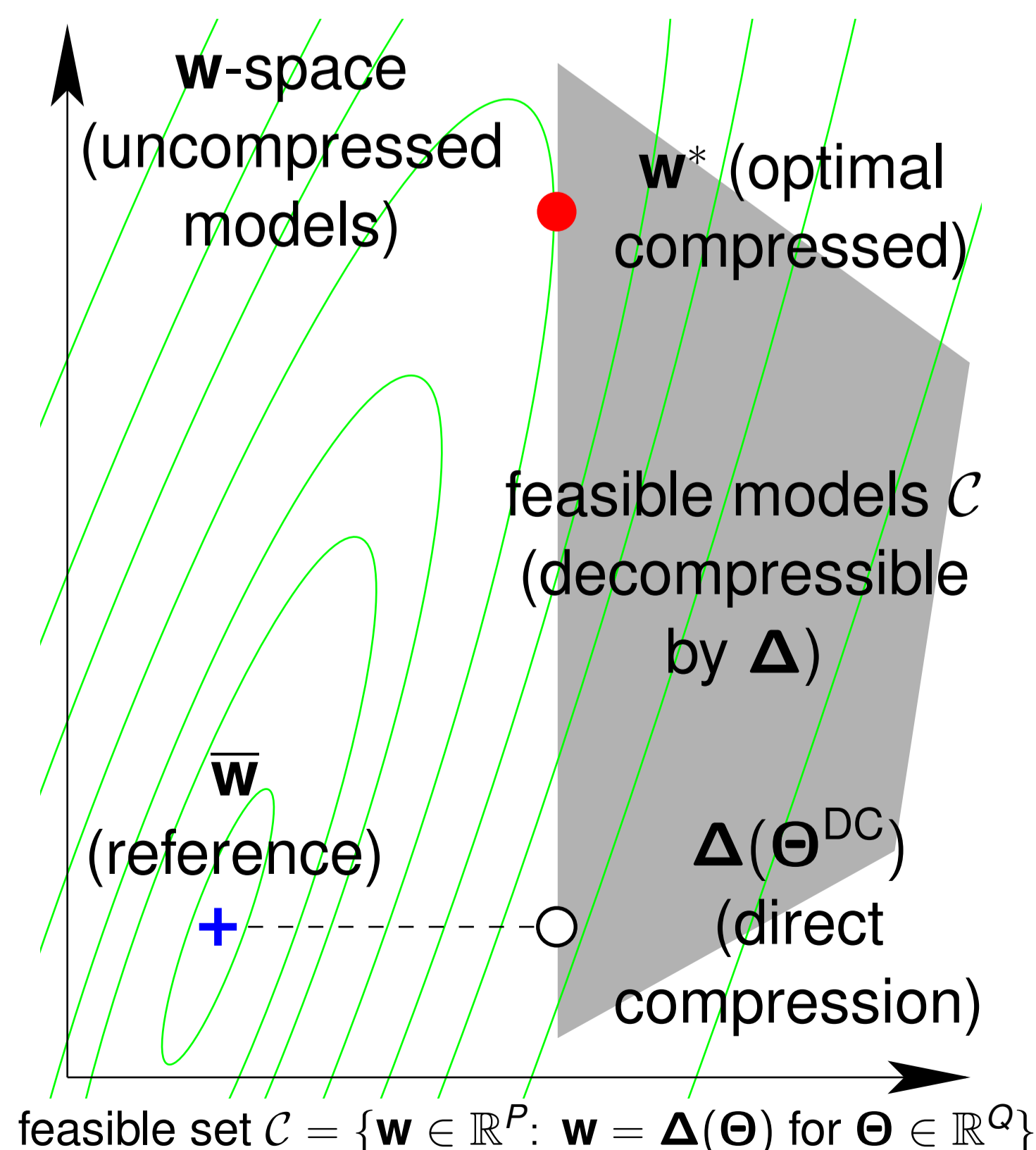
Desiderata for a good solution to the problem:

- A precise, general mathematical definition of model compression problem that is amenable to numerical optimization tools. We want to achieve the **lowest loss possible for a given compression technique and compression rate**
- We would like an algorithm that is generic, applicable to many compression techniques, has optimality guarantees, easy to integrate in existing deep learning toolboxes and efficient in training

Compression and decompression are usually seen as algorithms, but here we regard them as mathematical mappings in parameter space.

Our framework includes well-known types of compression which are abstracted in  $\Delta(\Theta)$ :

- **Low-rank compression** defines  $\Delta(\mathbf{U}, \mathbf{V}) = \mathbf{UV}^T$ . The compression mapping is given by the singular value decomposition (SVD) of  $\mathbf{W}$ .
- **Quantization** uses a discrete mapping  $\Delta$  given by assigning each weight to one of  $K$  codebook values. Mapping is given by  $k$ -means or by a form of rounding.
- **Pruning** defines  $\mathbf{w} = \Delta(\theta) = \theta$  where  $\mathbf{w}$  is real and  $\theta$  is constrained to have few nonzero values. The compression mapping involves some kind of thresholding.



## 5 The “learning-compression” (LC) algorithm

Use a penalty method (e.g. quadratic penalty): as  $\mu \rightarrow \infty$ , minimize:

$$Q(\mathbf{w}, \Theta; \mu) = L(\mathbf{w}) + \frac{\mu}{2} \|\mathbf{w} - \Delta(\Theta)\|^2$$

using alternating optimization over  $\mathbf{w}$  and  $\Theta$ :

**L step:**  $\min_{\mathbf{w}} L(\mathbf{w}) + \frac{\mu}{2} \|\mathbf{w} - \Delta(\Theta)\|^2$  – independent of the compression

**C step:**  $\min_{\Theta} \|\mathbf{w} - \Delta(\Theta)\|^2$  – independent of the loss and dataset.

This algorithm converges to a local stationary point of the constrained problem under standard assumptions as  $\mu \rightarrow \infty$ :

smooth loss  $L(\mathbf{w})$  and decompression mapping  $\Delta(\Theta)$  sufficiently accurate optimizations on L and C steps.

**L step always takes the same form** regardless of the compression technique:

$$\min_{\mathbf{w}} L(\mathbf{w}) + \frac{\mu}{2} \|\mathbf{w} - \Delta(\Theta)\|^2$$

- This is the original loss on the uncompressed weights  $\mathbf{w}$ , regularized with a quadratic term on the weights (which “decay” towards a validly compressed model  $\Delta(\Theta)$ ).
- It can be optimized in the same way as the reference net. Simply add  $\mu(\mathbf{w} - \Delta(\Theta))$  to the gradient.
- With large datasets we typically use SGD and clip the learning rates so they never exceed  $\frac{1}{\mu}$  to avoid oscillations as  $\mu \rightarrow \infty$ .

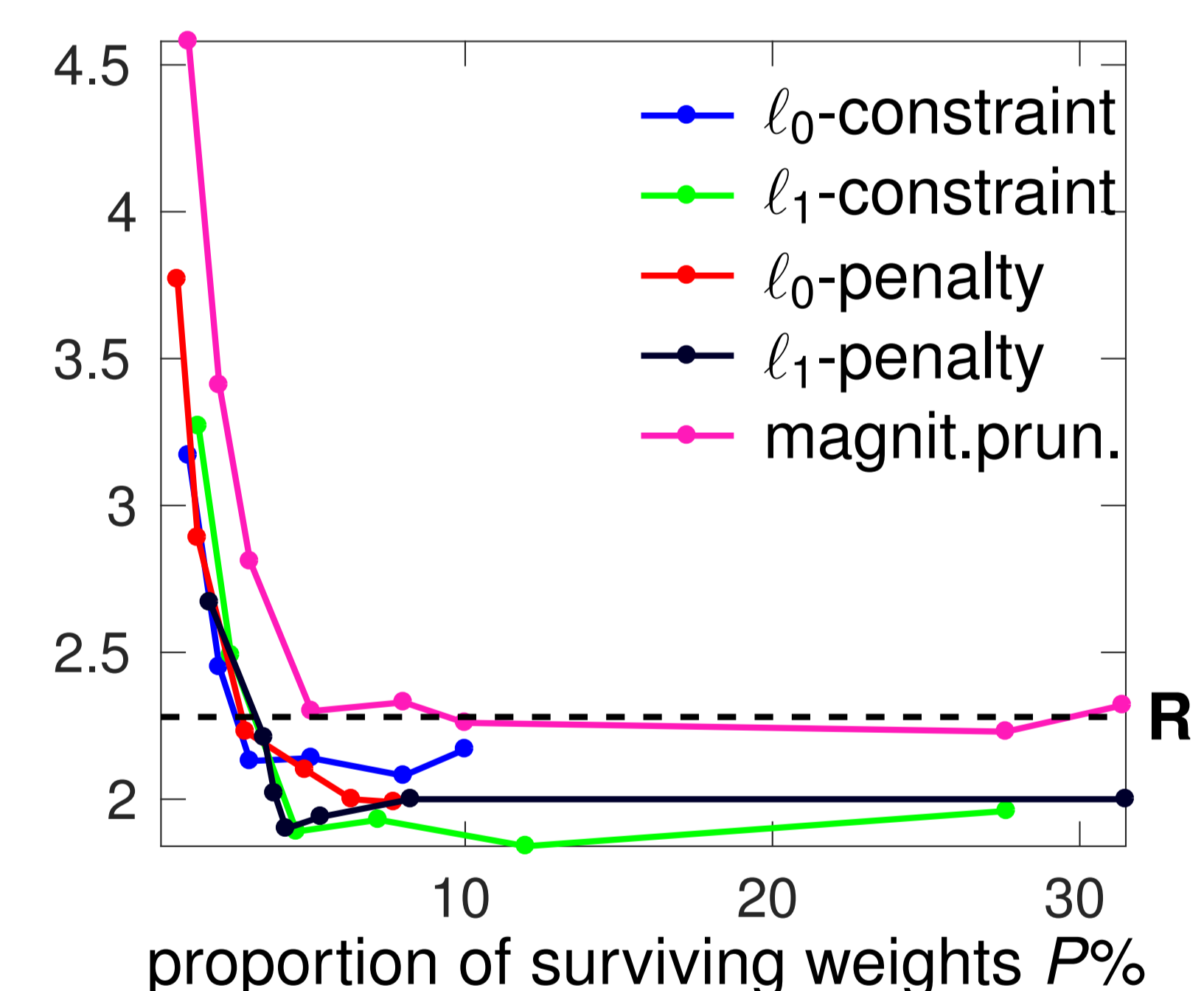
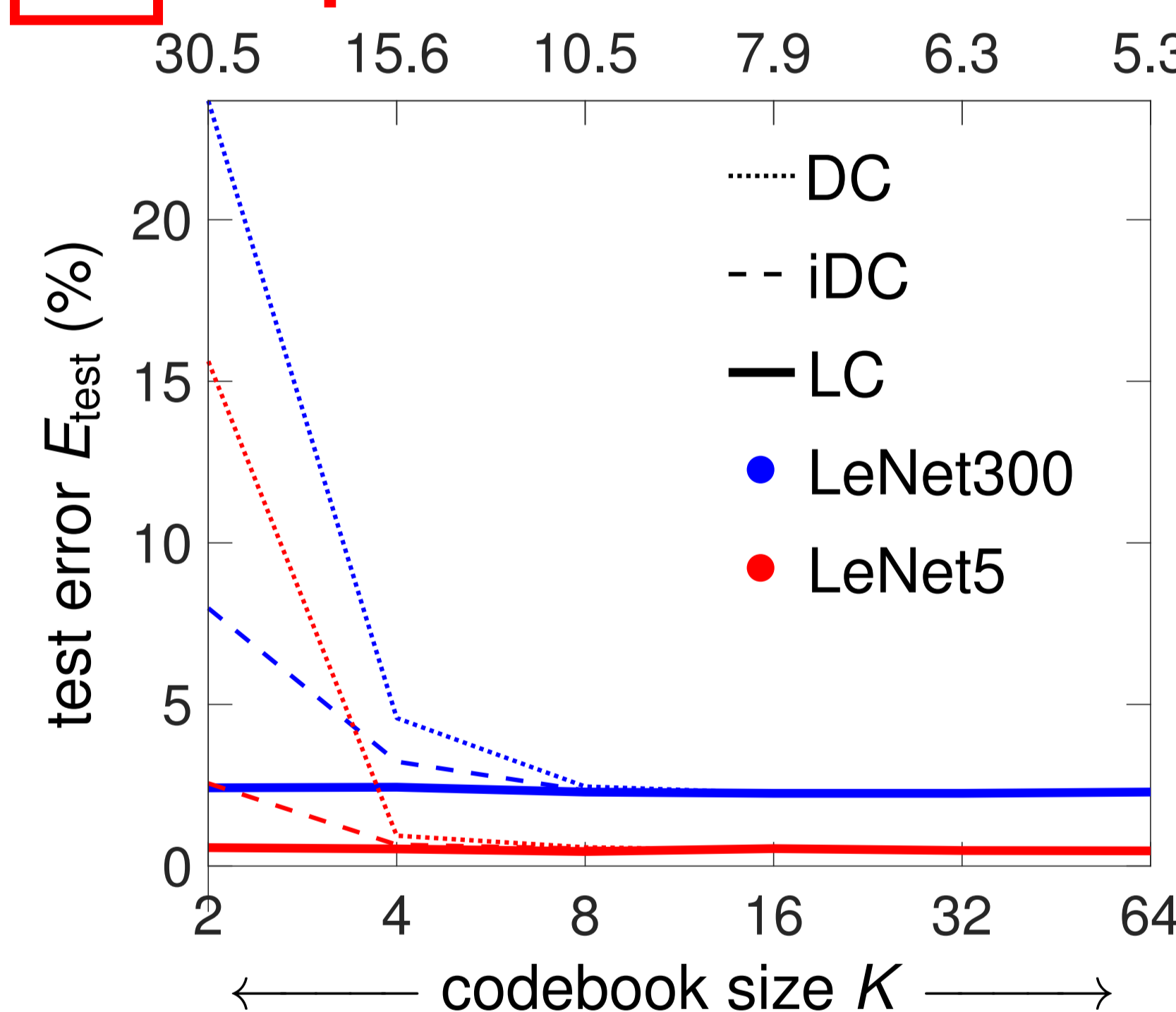
Solving the **C-step** means optimally compressing the current weights  $\mathbf{w}$ :

$$\min_{\Theta} \|\mathbf{w} - \Delta(\Theta)\|^2$$

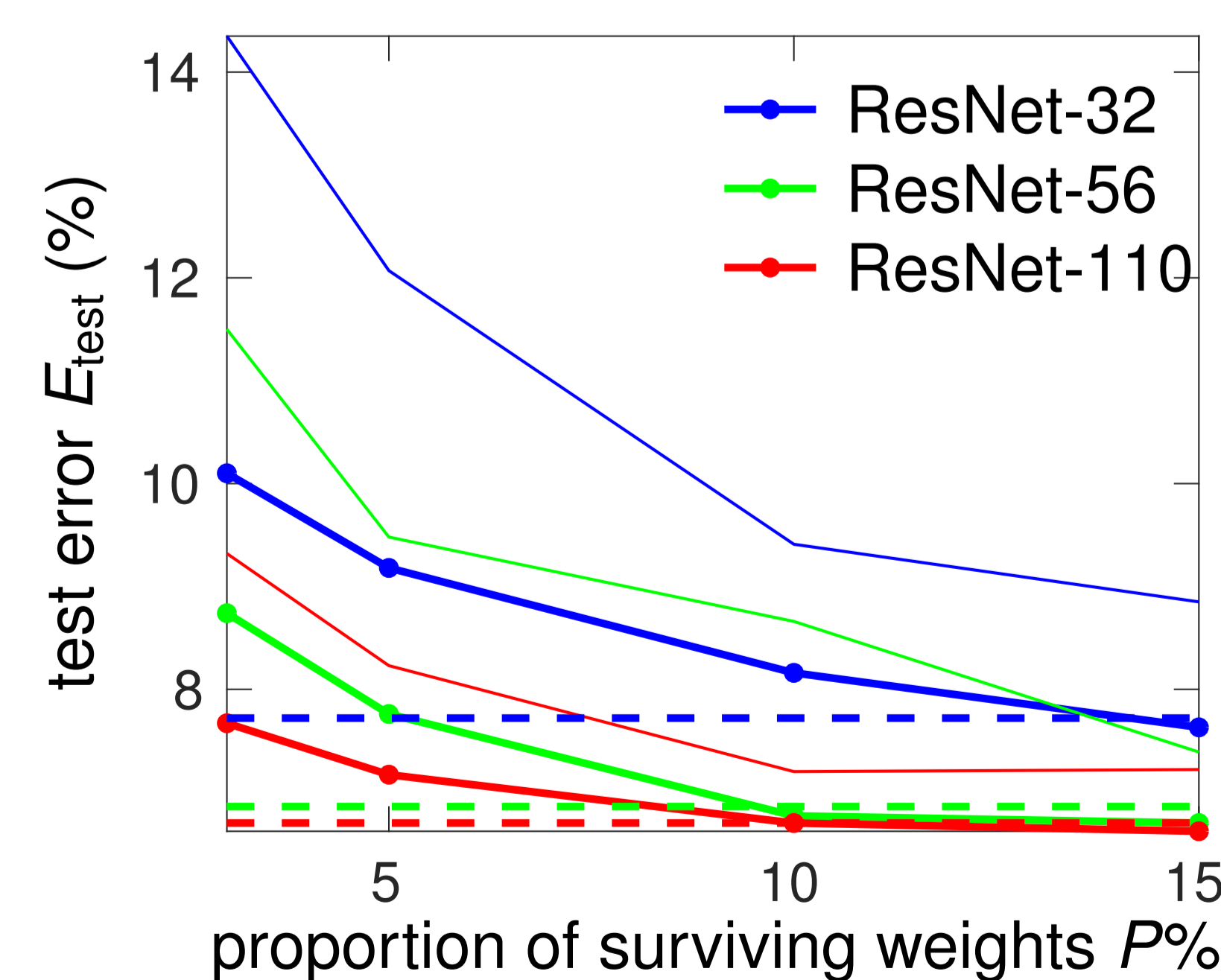
- Fast: it does not require the dataset (which appears in  $L(\mathbf{w})$ ).
- Equivalent to orthogonal projection  $\Theta = \Pi(\mathbf{w})$  on the feasible set and finds the closest compressed model to  $\mathbf{w}$ .
- The solution depends on the choice of the decompression mapping  $\Delta$ , and is known for many common compression techniques.
  - low rank: SVD
  - binarization: binarize weights
  - quantization in general:  $k$ -means
  - pruning: zero all but top weights

So the **C step** simply requires calling the corresponding compression subroutine as a black box!

## 6 Experiments



LC algorithm applied to compression of benchmark nets on MNIST. (Left) quantization using adaptive codebooks for **LeNet300** and **LeNet5**; **no degradation in error for LC**. (Right) Pruning using different norms ( $\ell_0, \ell_1$ ) and forms (penalty, constraint).



Results of pruning **ResNets** of different depth using LC, trained on **CIFAR10** (60K images, 10 classes) and having 0.4M, 0.8M and 1.7M parameters. Thick lines – LC, thin lines – magnitude based pruning, dashed horizontal lines – reference models. **LC achieves marginally better errors in high pruning regions,  $P < 5\%$  and always better than pruning-retraining strategy.**

We compress **AlexNet** (ILSVRC2012, 1.2M images, 1000 classes) having 61M parameters and trained to accuracy 79.5%, using  $k = 2$  codebook ( $\times 32$  compression) and achieve error of 77.1%.

### References

- “Model compression as constrained optimization, with application to neural nets”:
- 1) “Part I: general framework”, arxiv 2017 (1707.01209)
  - 2) “Part II: quantization”, arxiv 2017 (1707.04319)
  - 3) “Part III: pruning”, arxiv 2017