

Band-limited Training and Inference for Convolutional Neural Networks

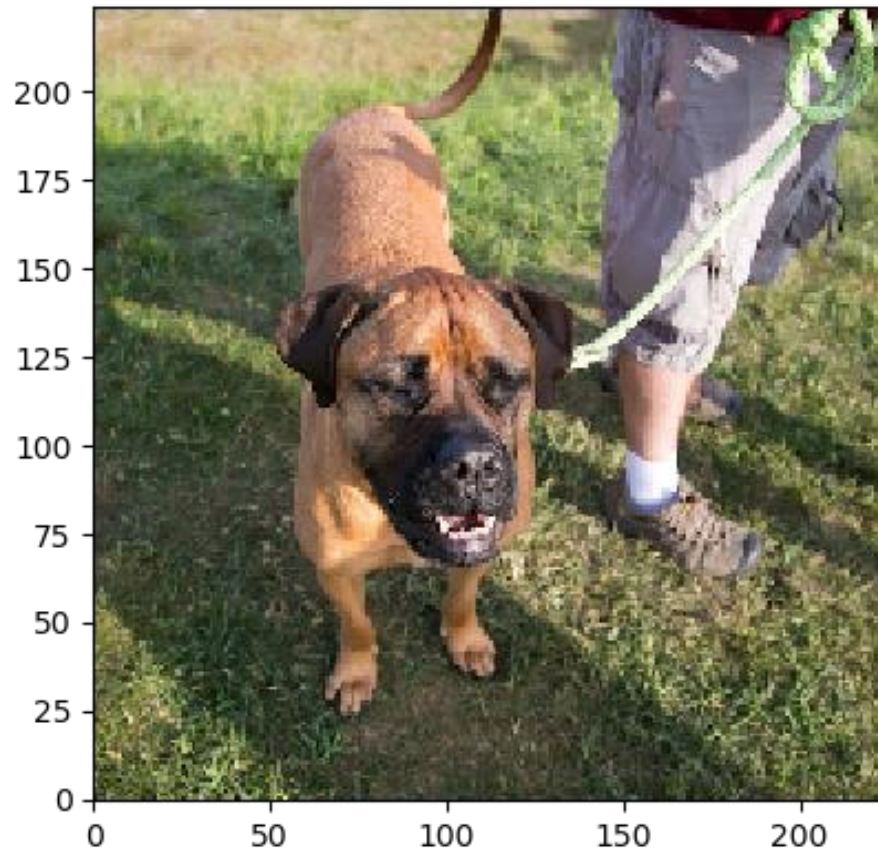
Adam Dziedzic*, John Paparrizos*, Sanjay Krishnan,
Aaron Elmore, Michael Franklin



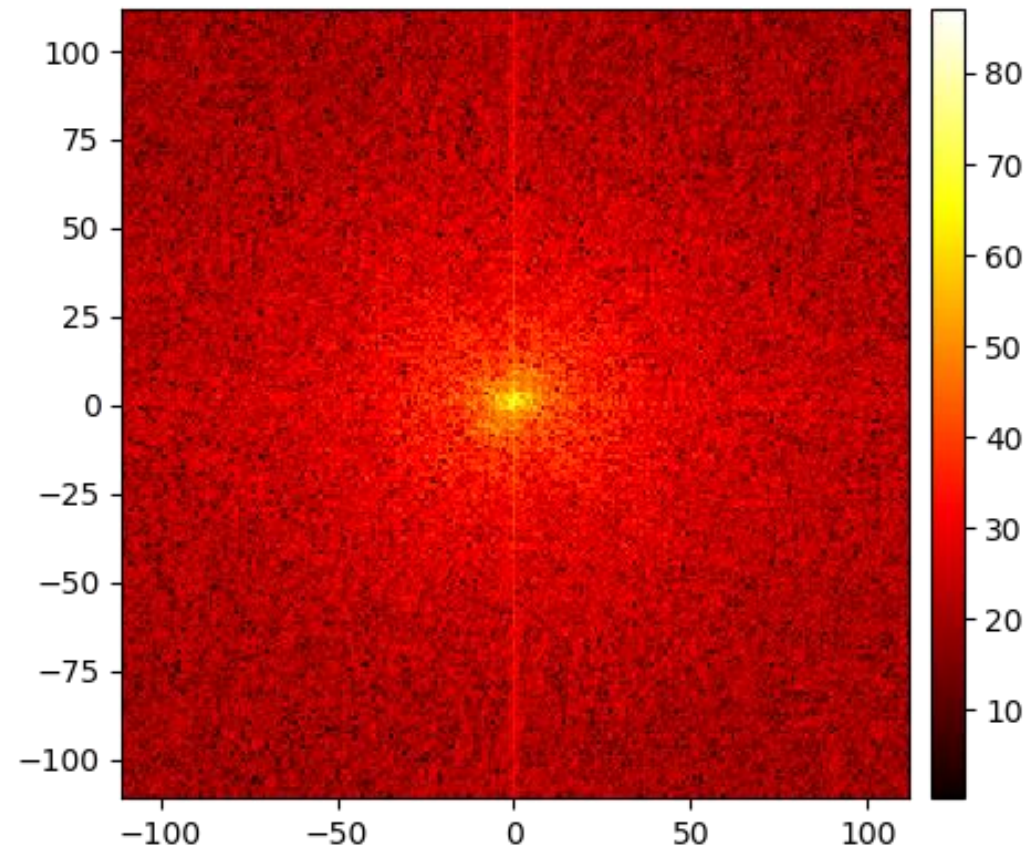
Natural images

More information put in lower frequencies

Original image



Spatial domain

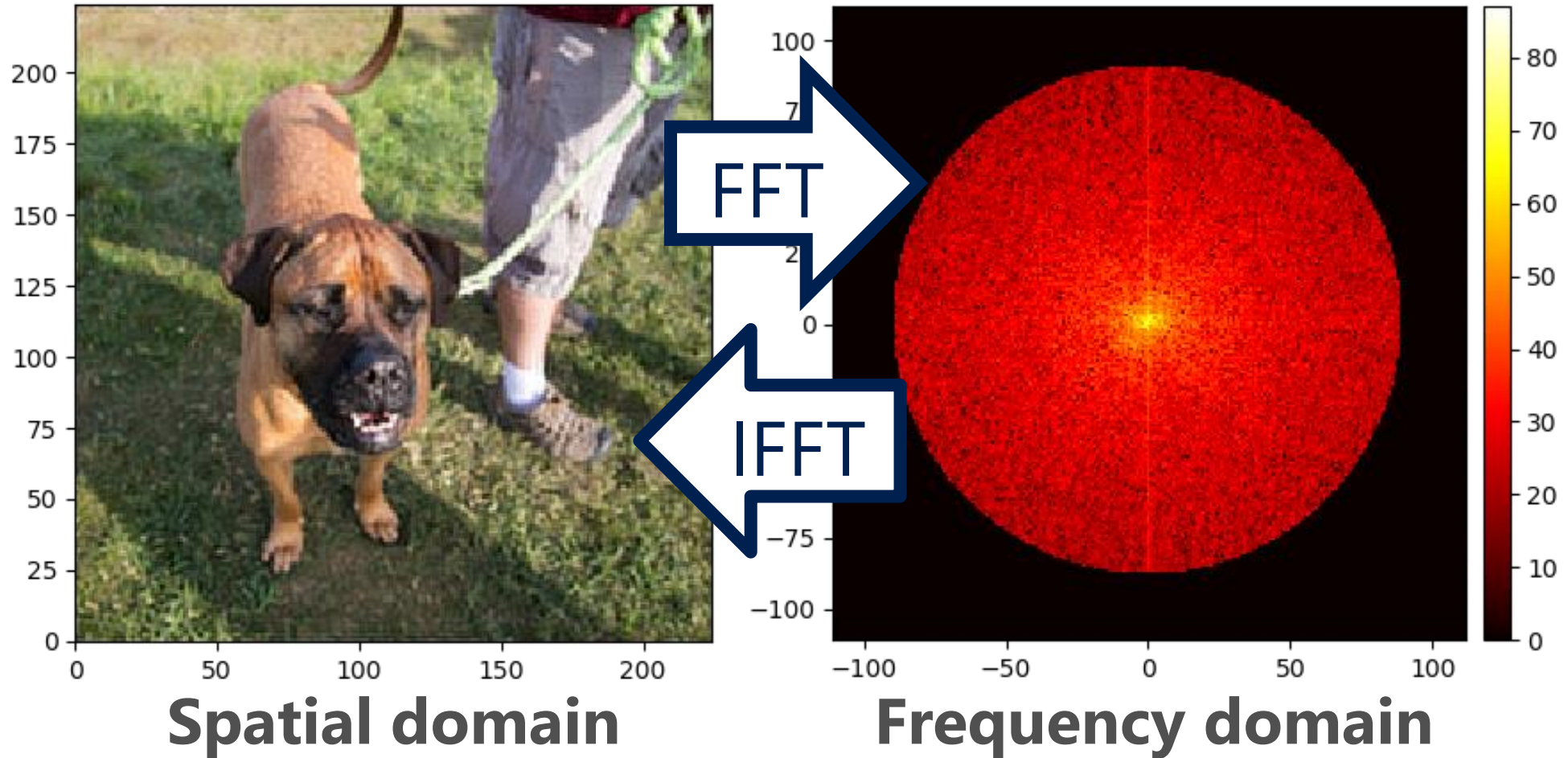


Frequency domain

Natural images

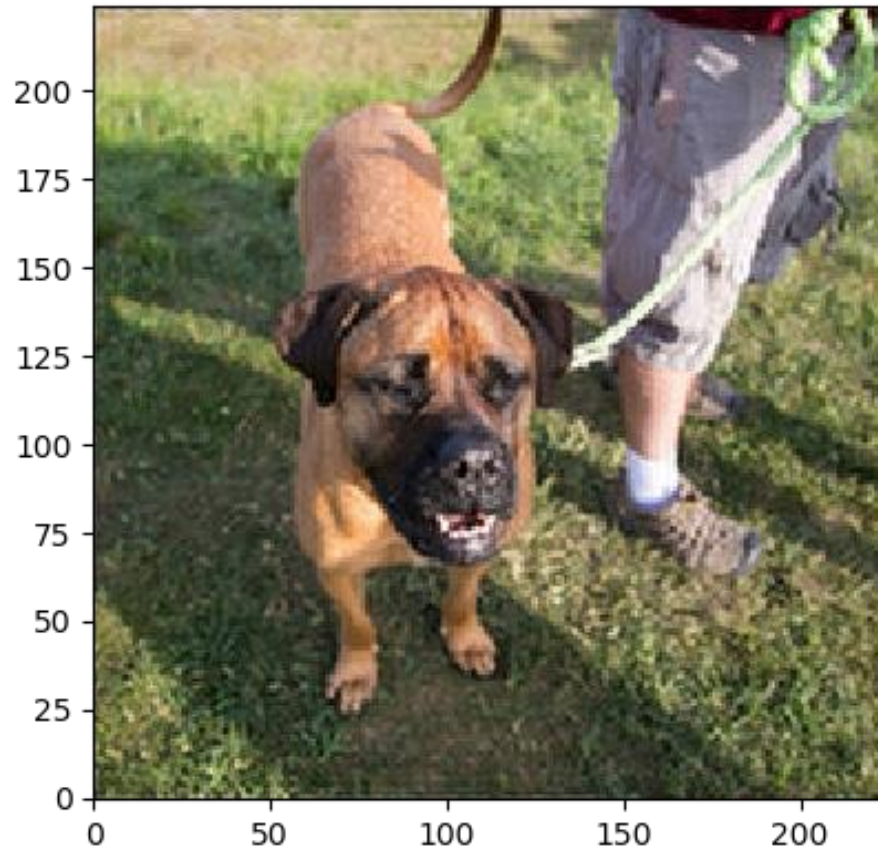
Transformations between the domains

Compression 50%

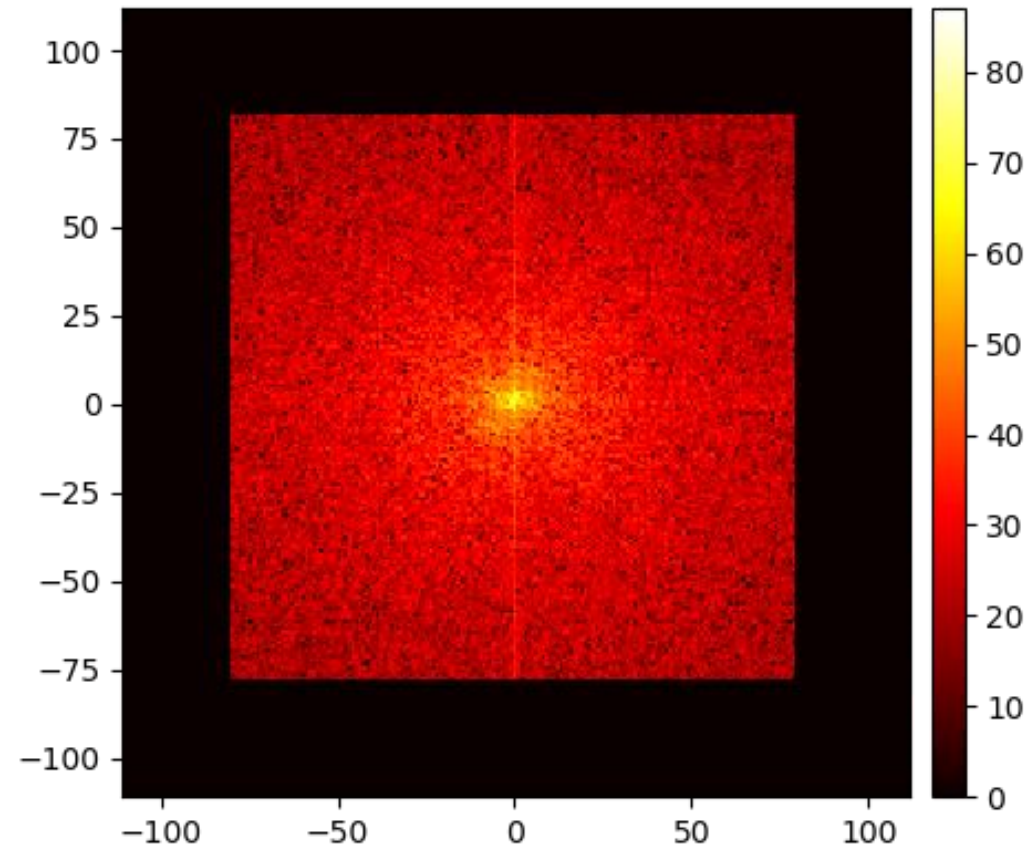


Method for ConvNets to constrain the frequency band in convolution operation for efficiency

Compression 50% in practice



Spatial domain

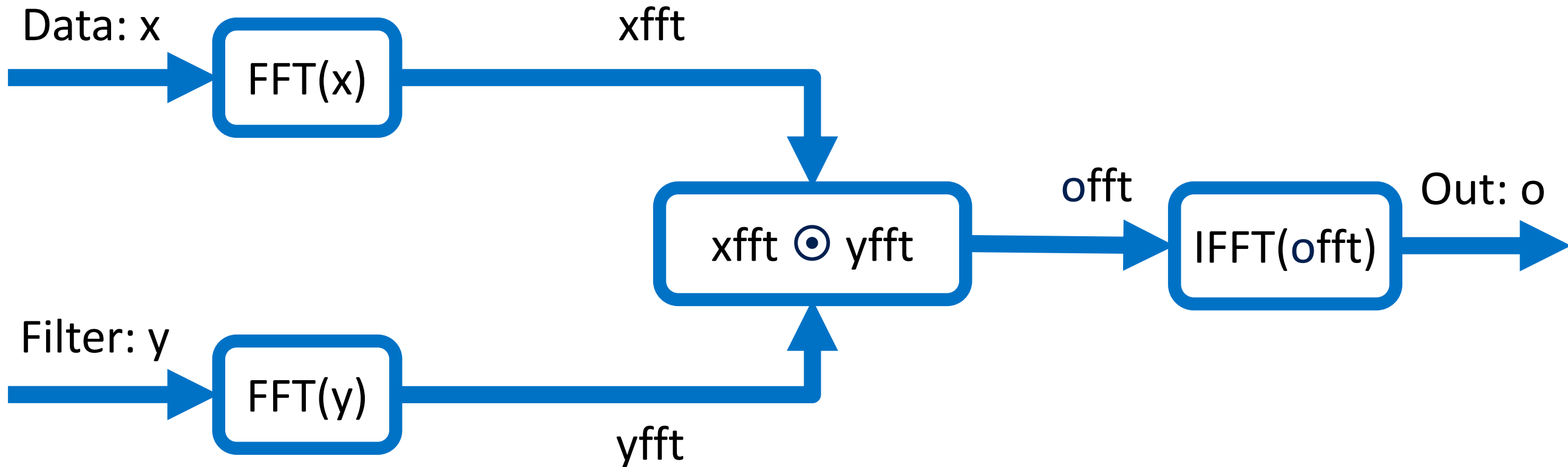


Frequency domain

FFT based convolution

Mathieu et al.: "Fast Training of Convolutional Networks through FFTs"

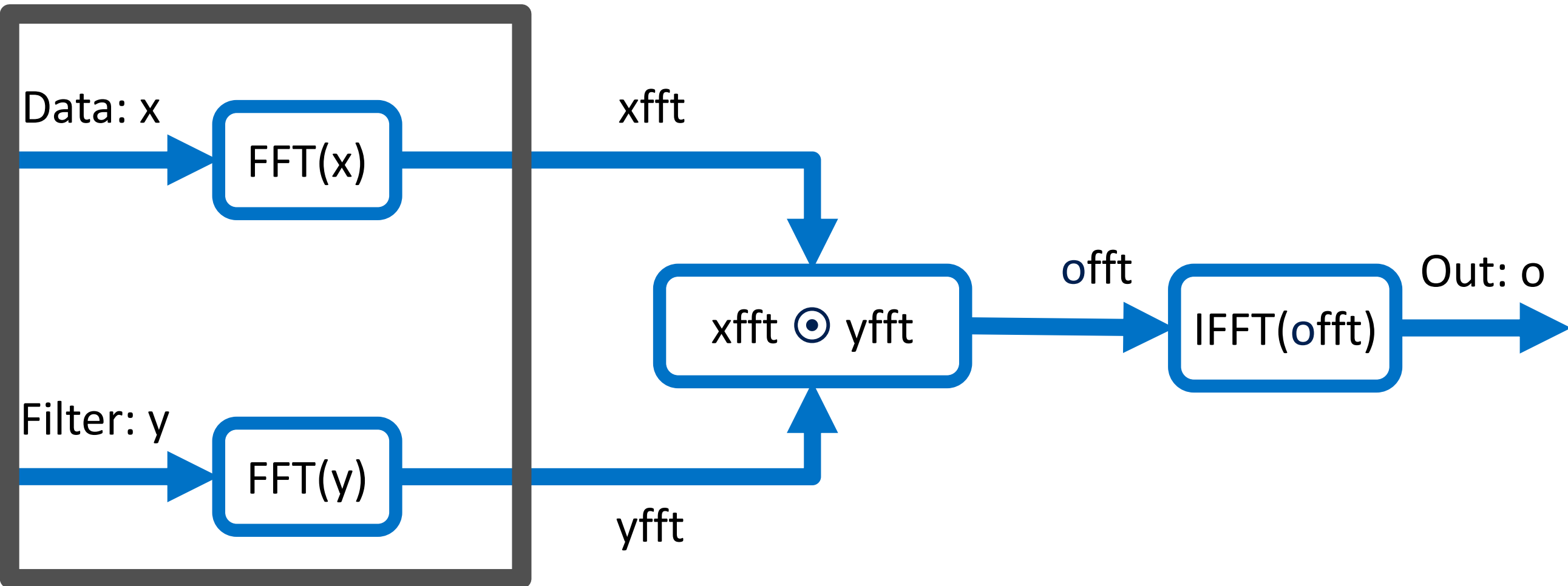
Vasilache et al.: "Fast Convolutional Nets With fbfft: A GPU Performance Evaluation"



FFT based convolution

Mathieu et al.: "Fast Training of Convolutional Networks through FFTs"

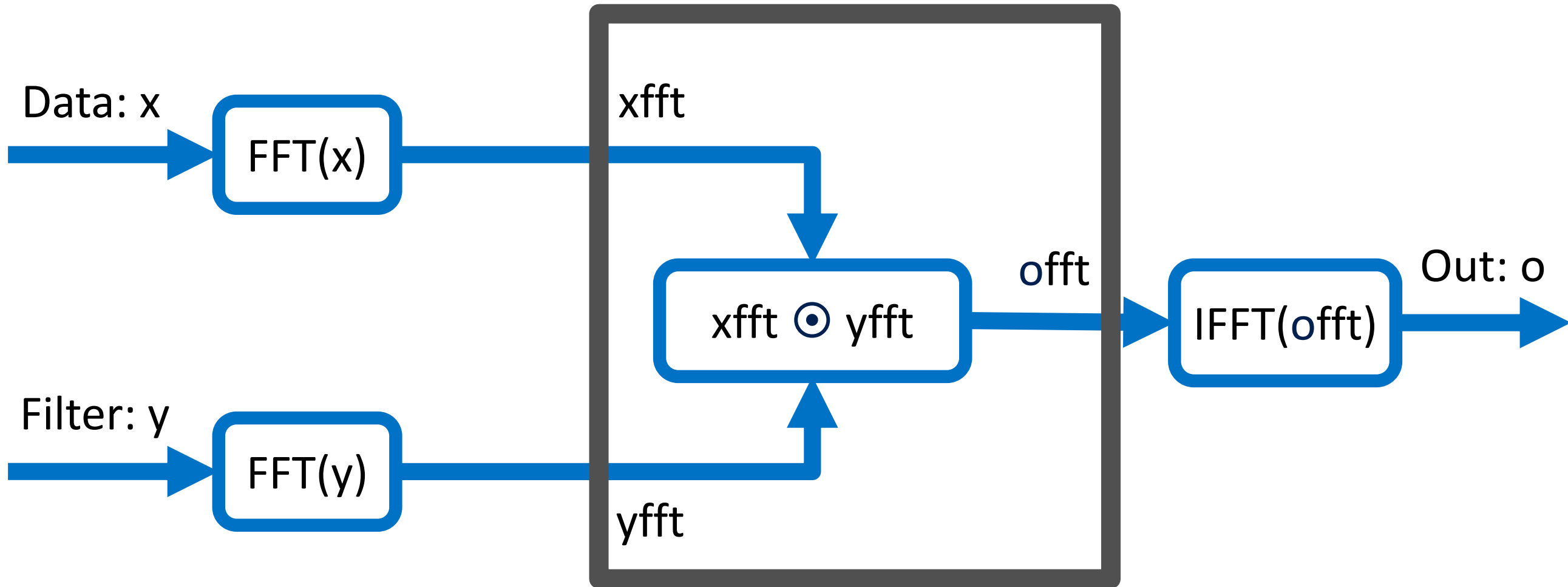
Vasilache et al.: "Fast Convolutional Nets With fbfft: A GPU Performance Evaluation"



FFT based convolution

Mathieu et al.: "Fast Training of Convolutional Networks through FFTs"

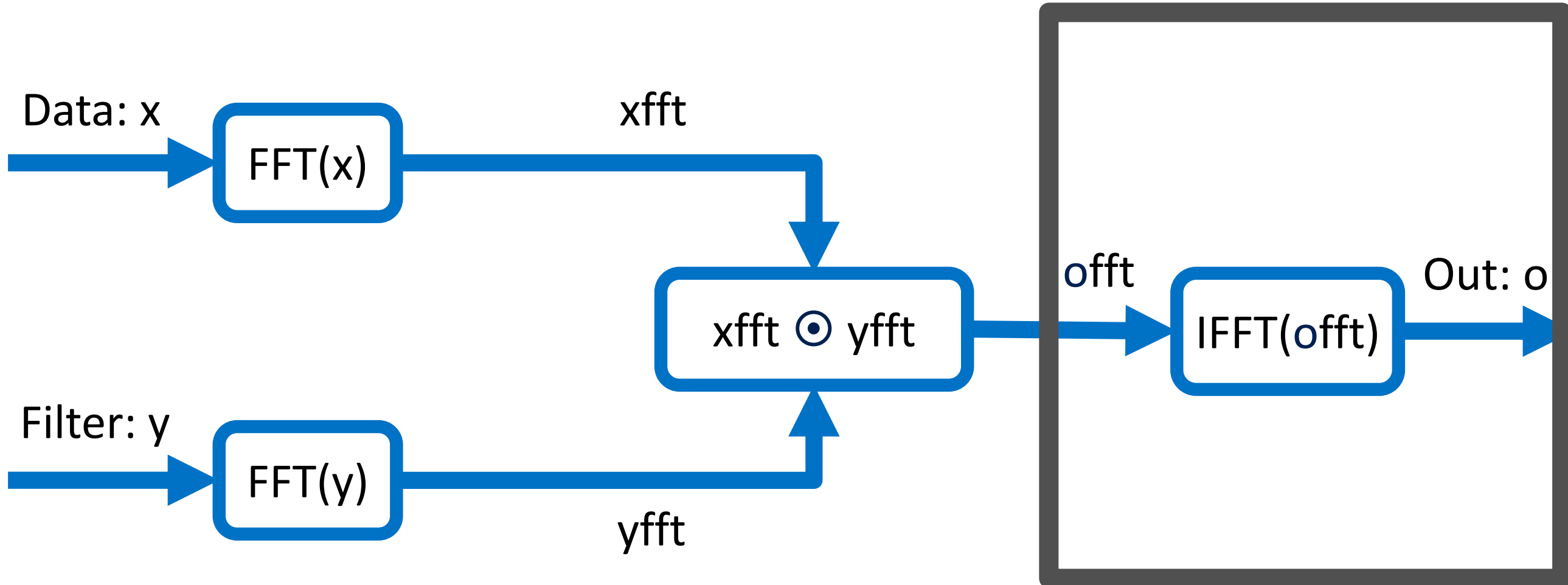
Vasilache et al.: "Fast Convolutional Nets With fbfft: A GPU Performance Evaluation"



FFT based convolution

Mathieu et al.: "Fast Training of Convolutional Networks through FFTs"

Vasilache et al.: "Fast Convolutional Nets With fbfft: A GPU Performance Evaluation"

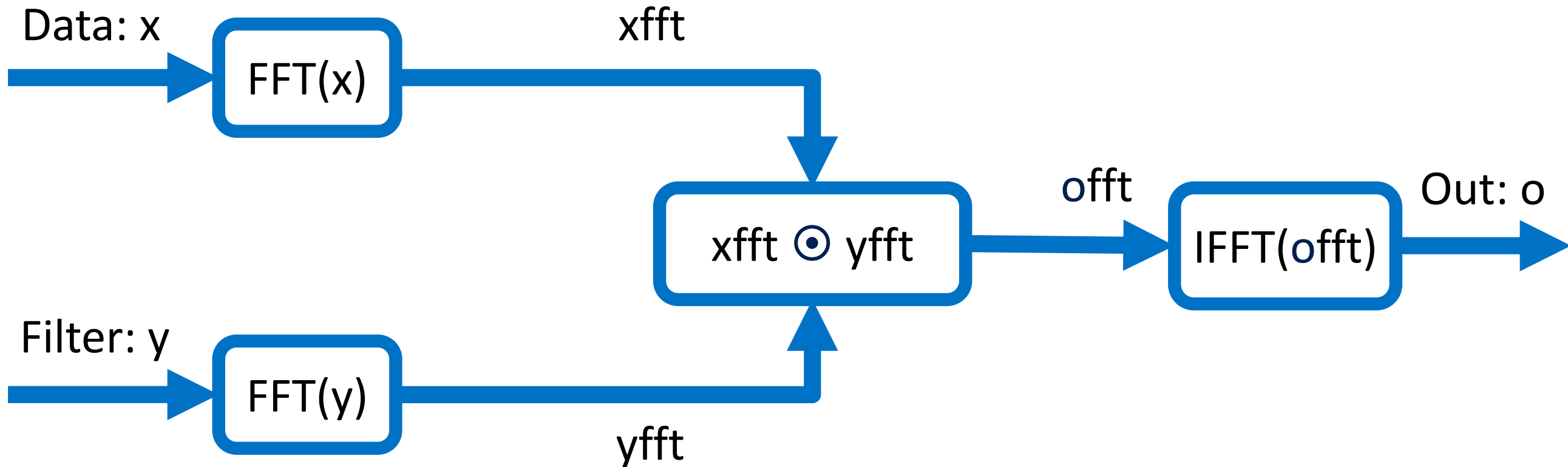


FFT based convolution

Mathieu et al.: "Fast Training of Convolutional Networks through FFTs"

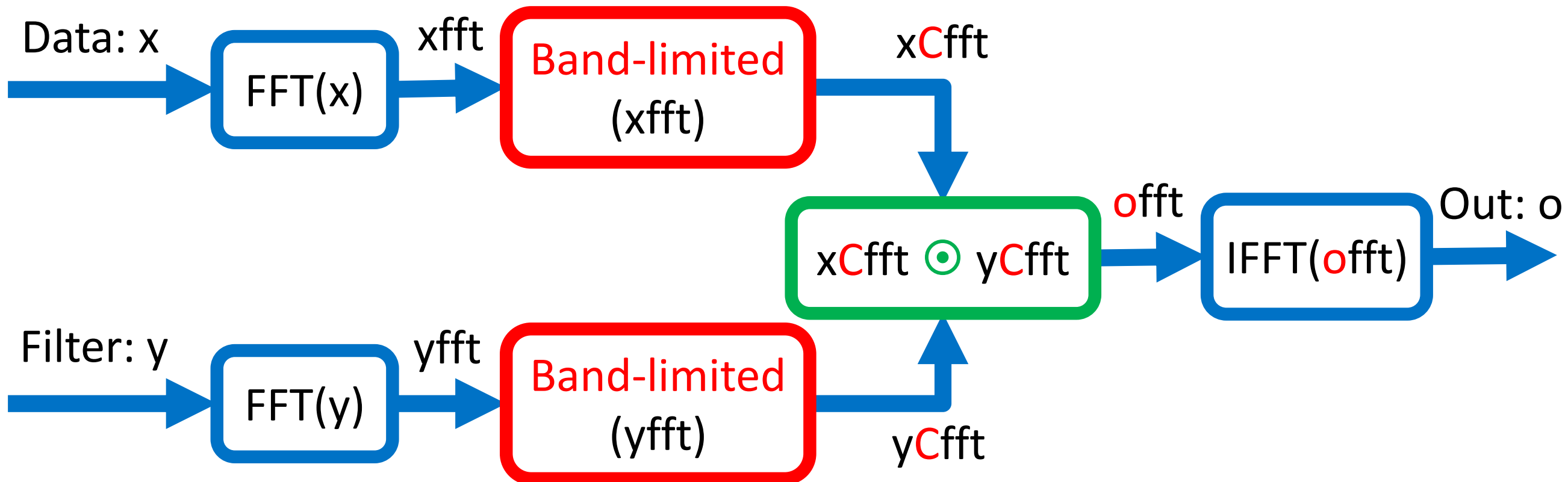
Vasilache et al.: "Fast Convolutional Nets With fbfft: A GPU Performance Evaluation"

cuDNN: Substantial memory workspace needed for intermediate results.

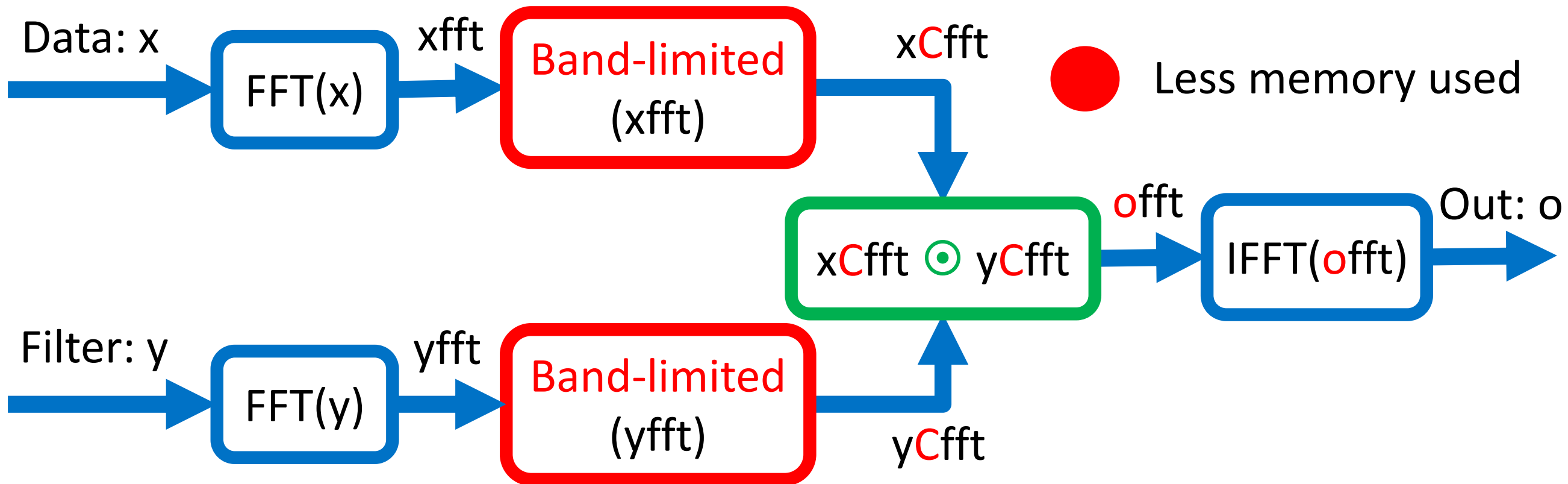


Band-limited FFT based convolution

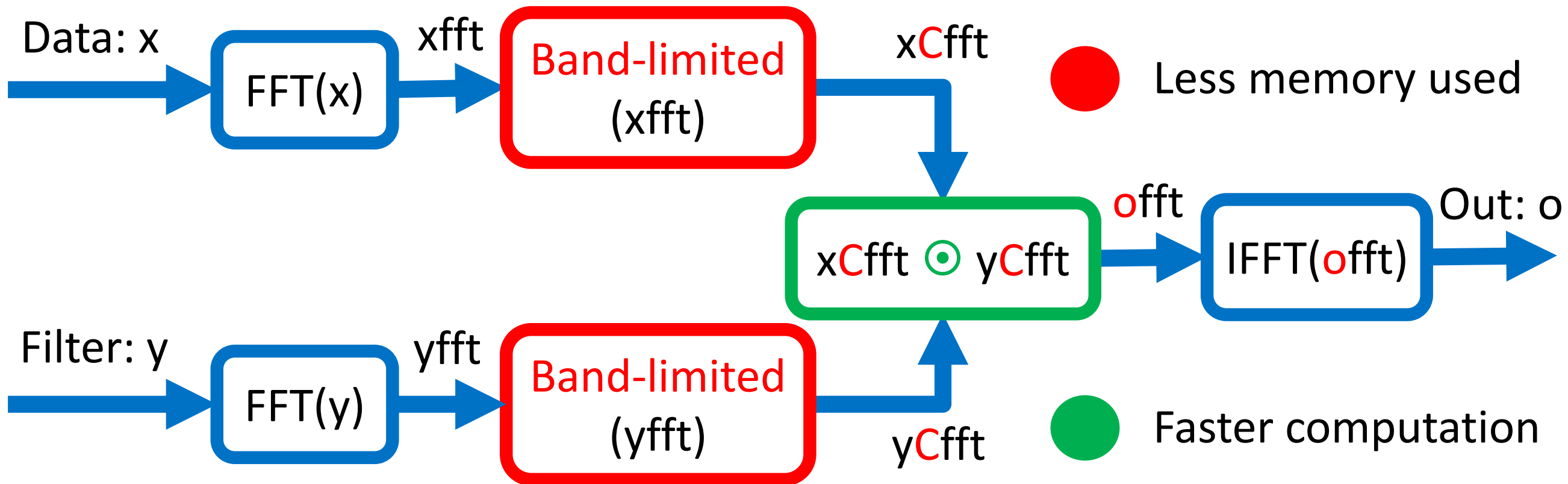
Band-limiting = masking out high frequencies



Band-limited FFT based convolution

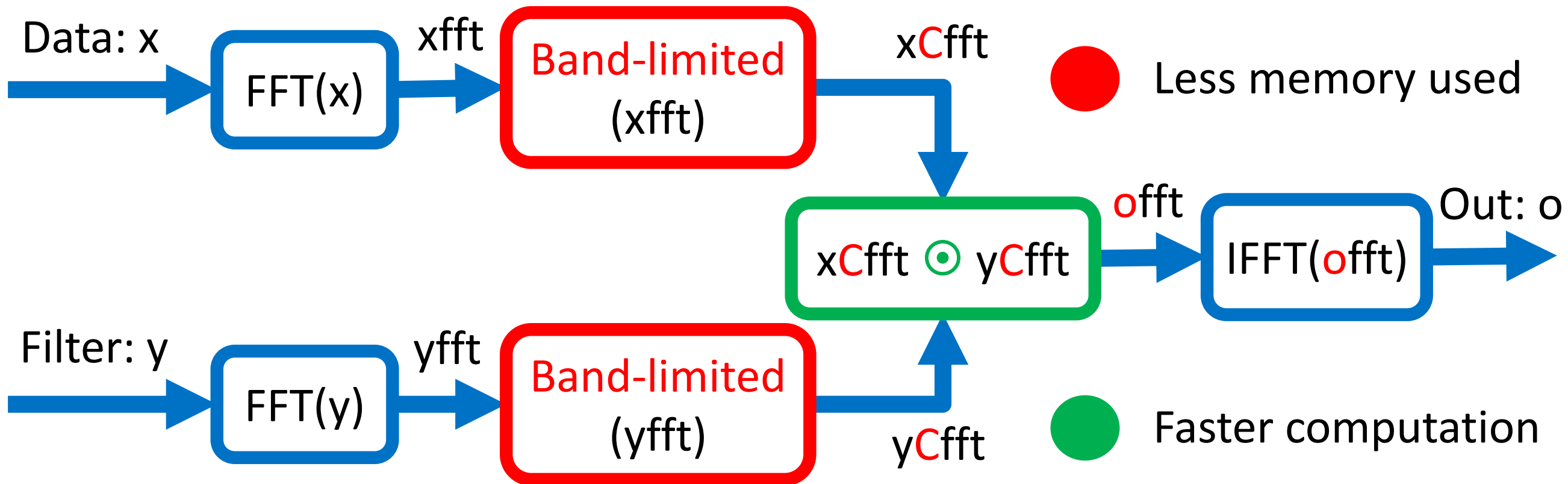


Band-limited FFT based convolution



Band-limited FFT based convolution

Preserve enough of the spectrum to retain high accuracy of models.

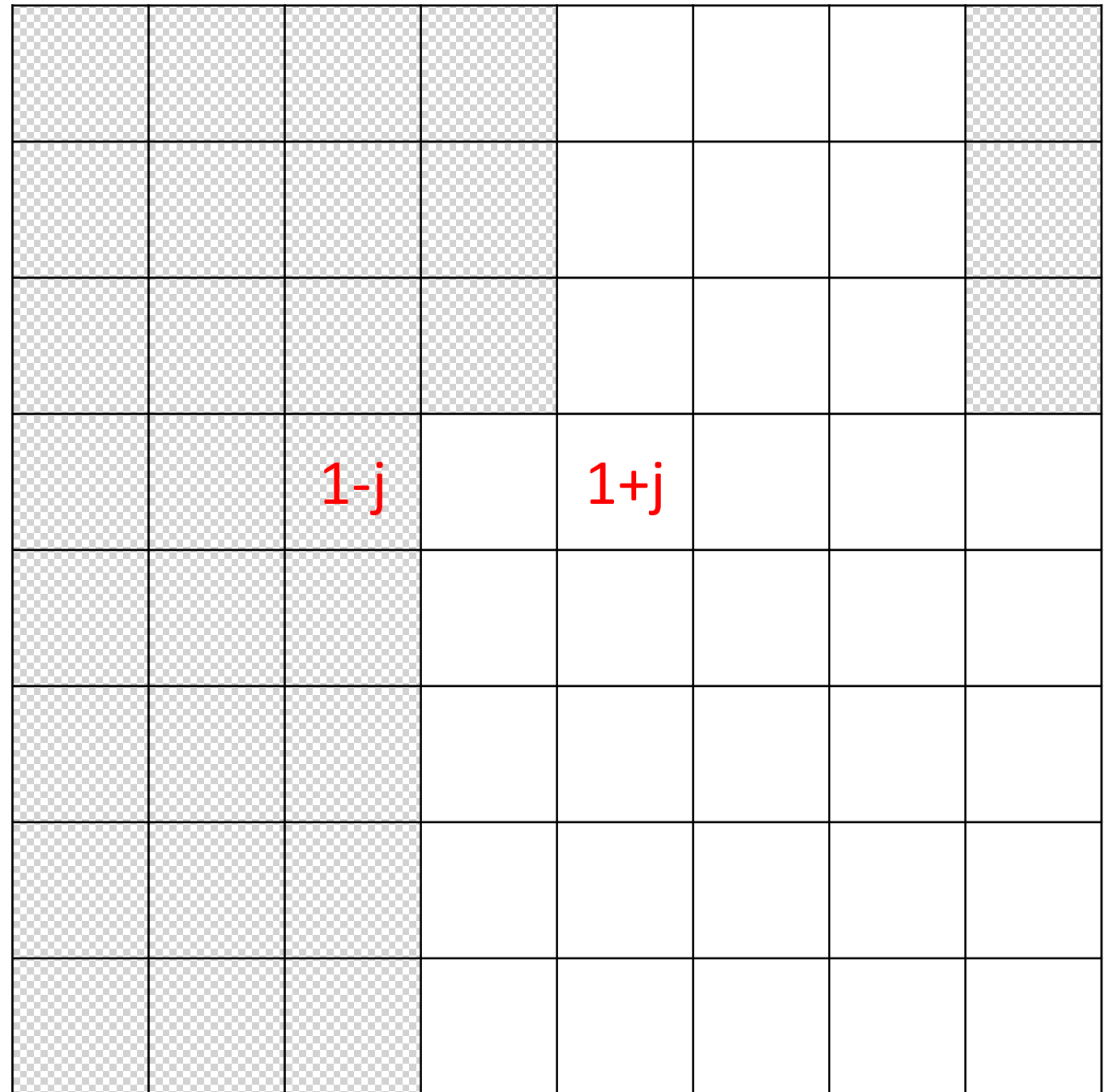


Band-limiting Technique

1. FFT of an input data

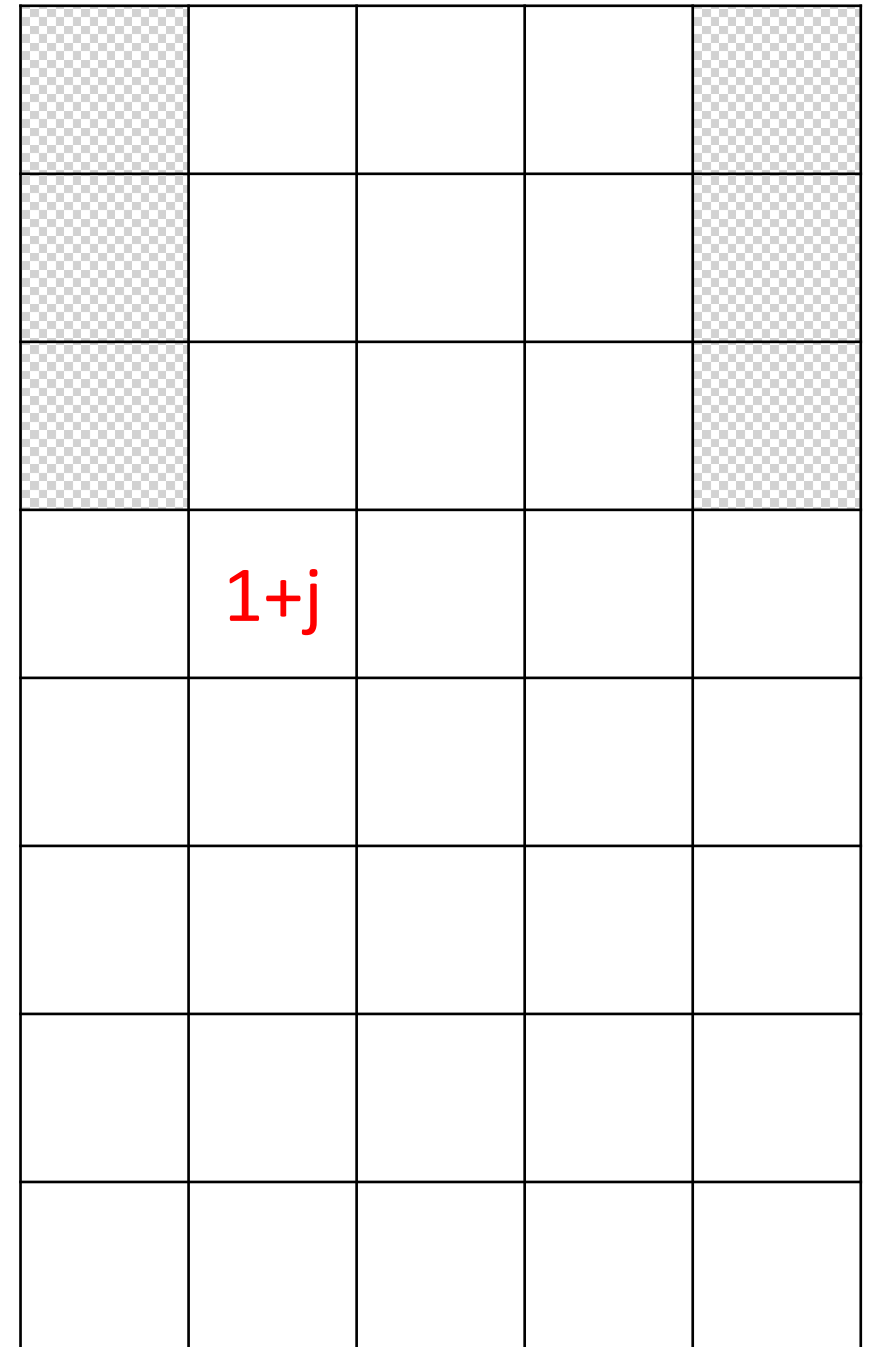
Band-limiting Technique

1. FFT of an input data
2. Conjugate symmetry



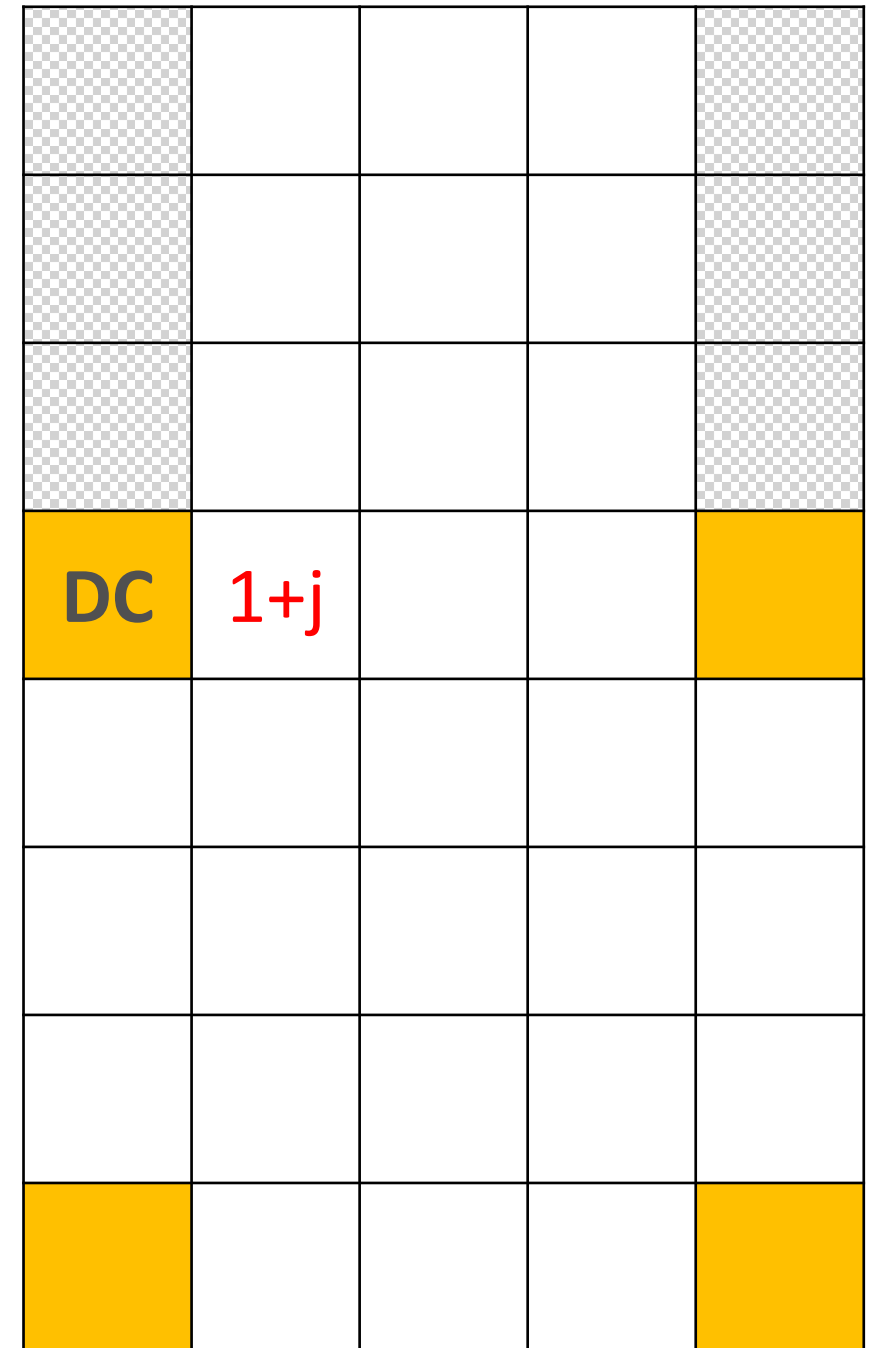
Band-limiting Technique

1. FFT of an input data
2. Conjugate symmetry



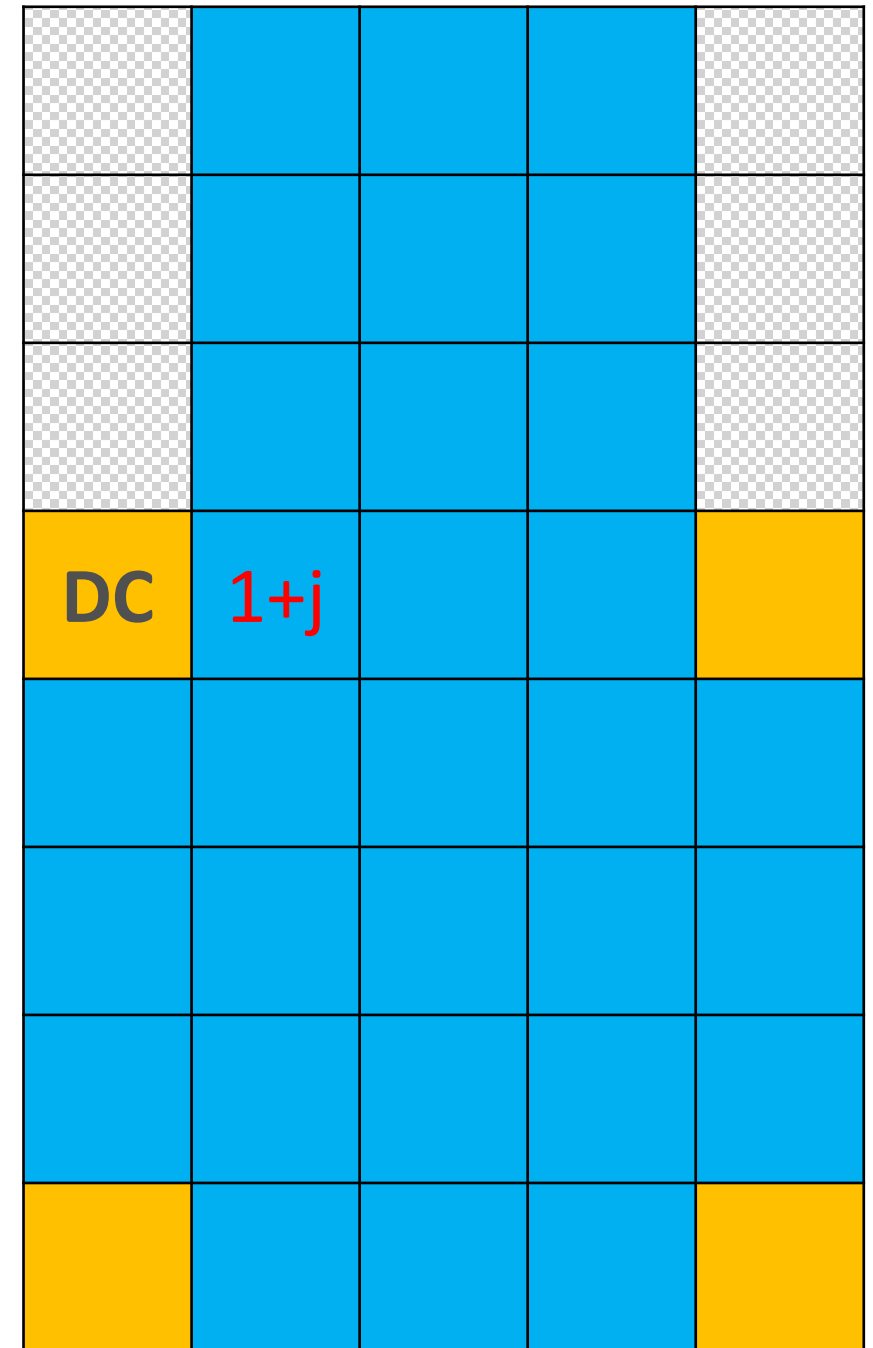
Band-limiting Technique

1. FFT of an input data
2. Conjugate symmetry
3. Real values



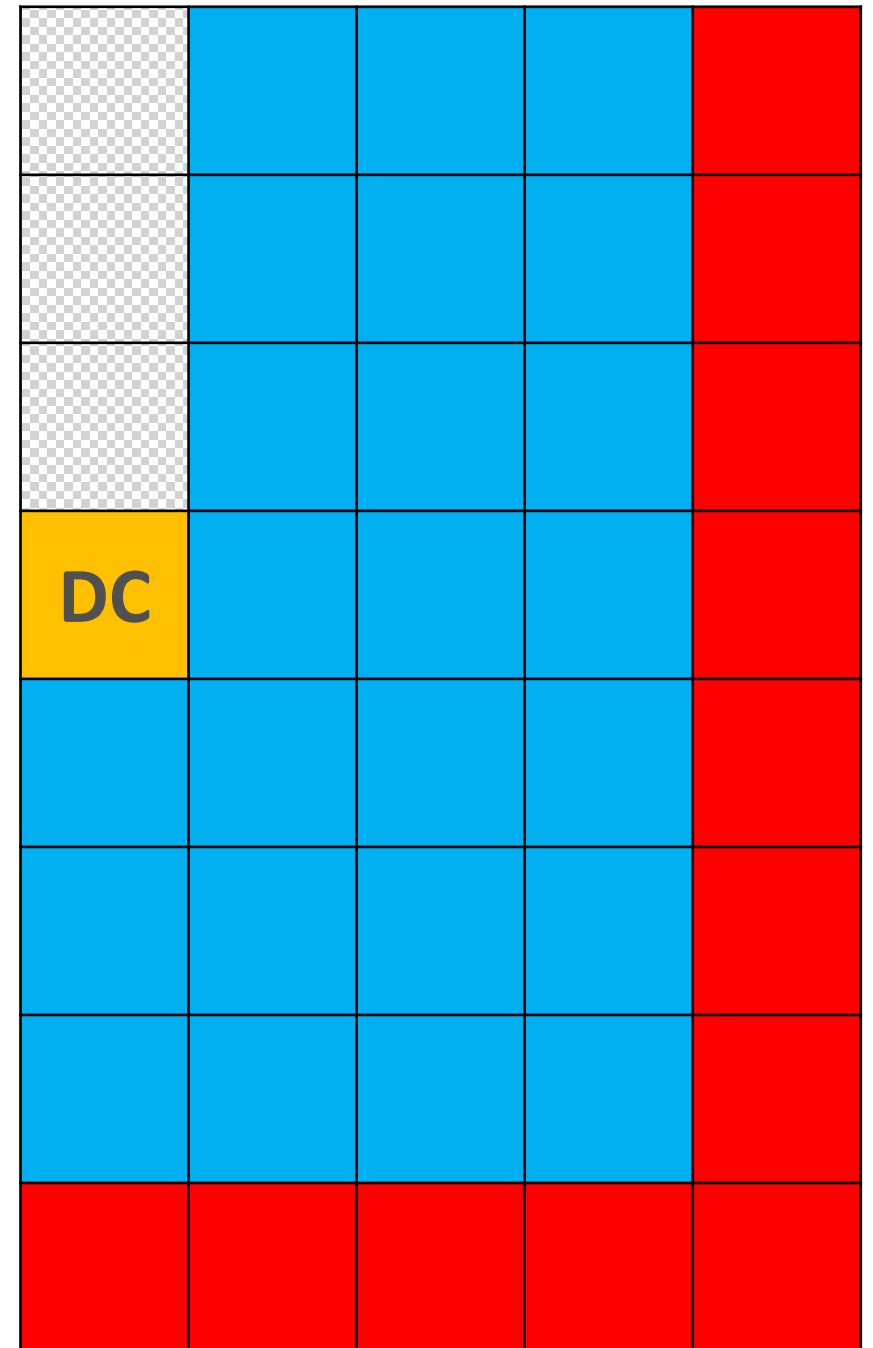
Band-limiting Technique

1. FFT of an input data
2. Conjugate symmetry
3. Real values
4. No constraints



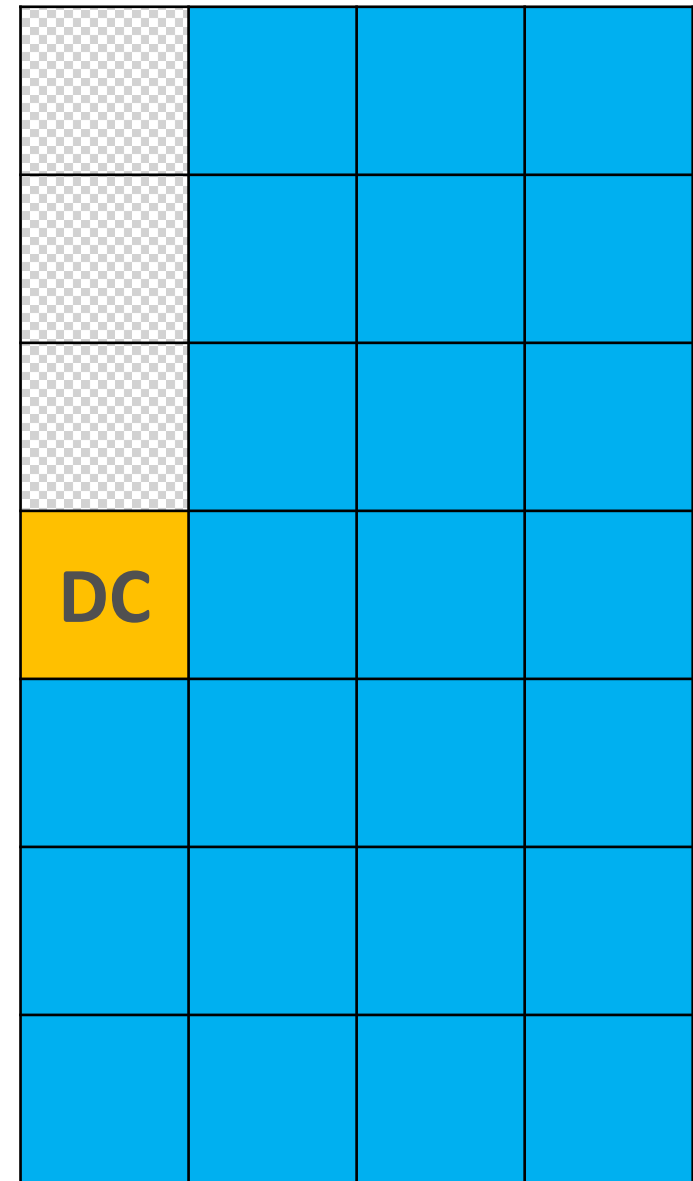
Band-limiting Technique

1. FFT of an input data
2. Conjugate symmetry
3. Real values
4. No constraints
5. 1st compression



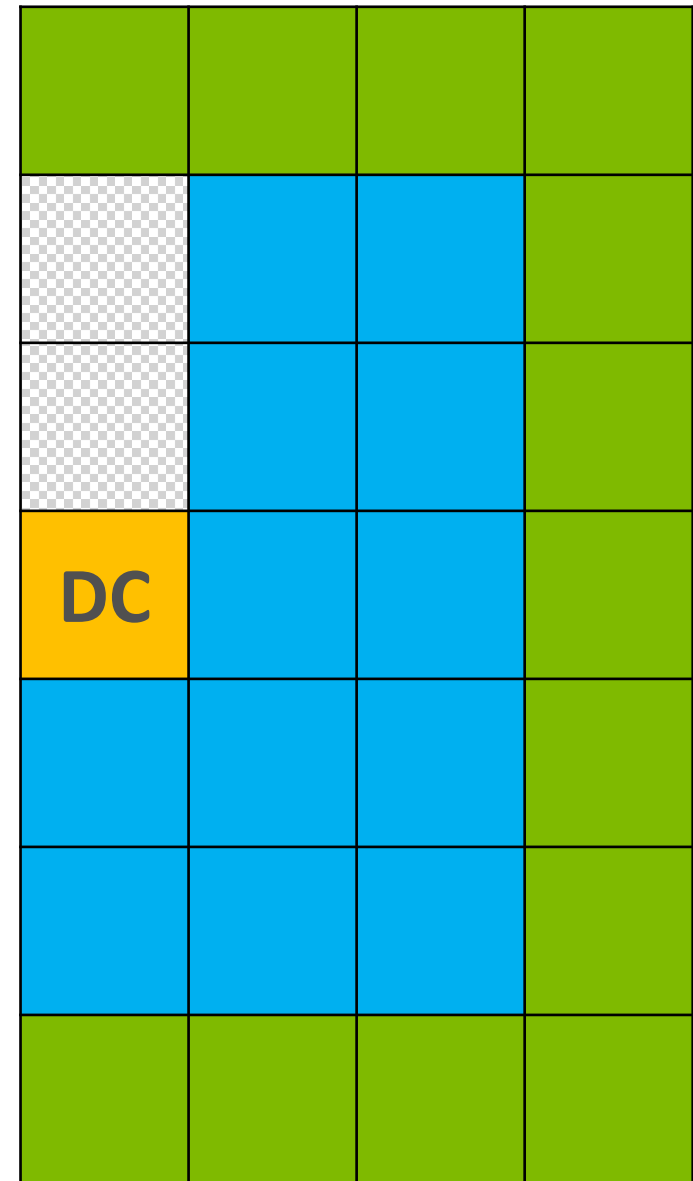
Band-limiting Technique

1. FFT of an input data
2. Conjugate symmetry
3. Real values
4. No constraints
5. 1st compression



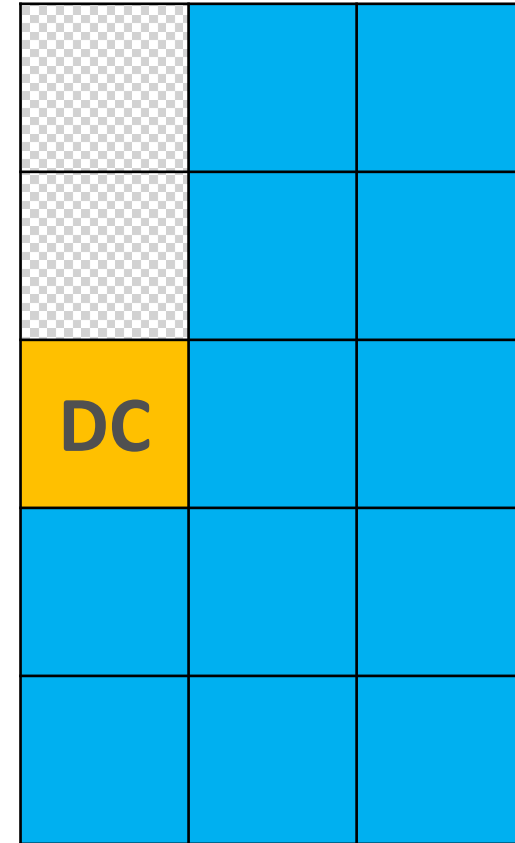
Band-limiting Technique

1. FFT of an input data
2. Conjugate symmetry
3. Real values
4. No constraints
5. 1st compression
6. 2nd compression

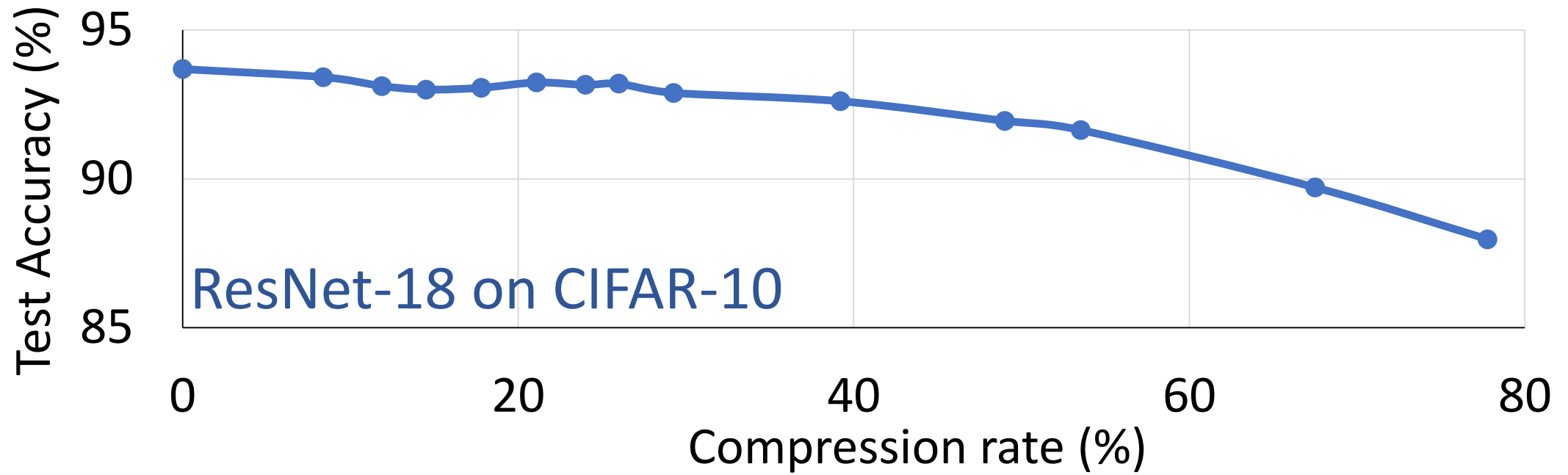


Band-limiting Technique

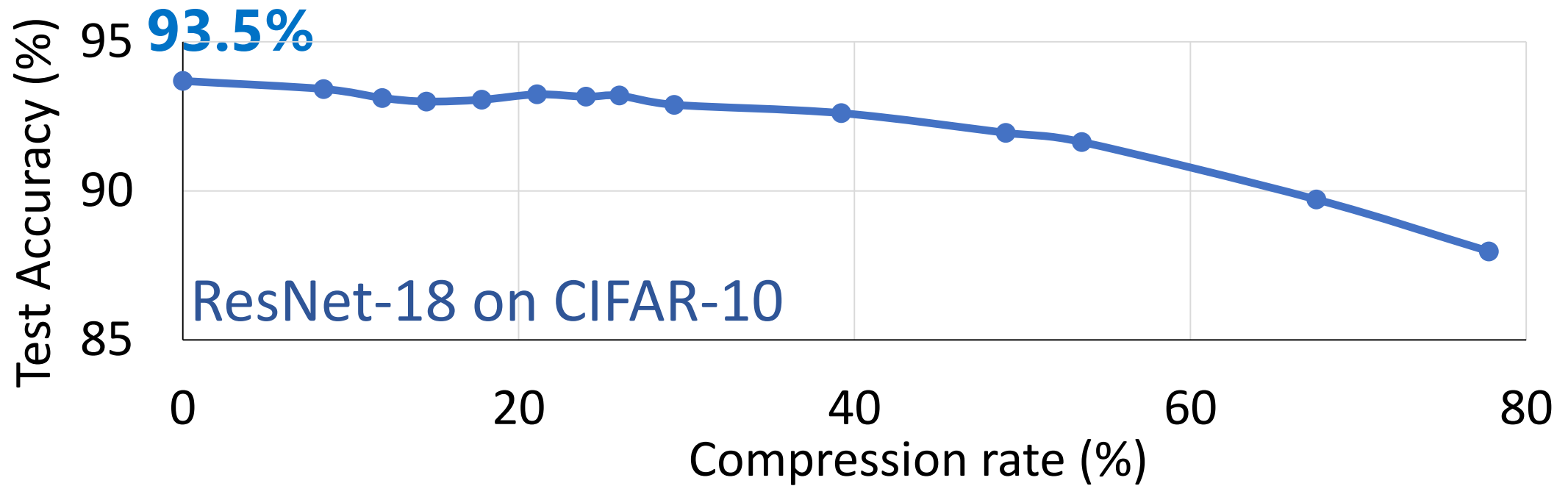
1. FFT of an input data
2. Conjugate symmetry
3. Real values
4. No constraints
5. 1st compression
6. 2nd compression



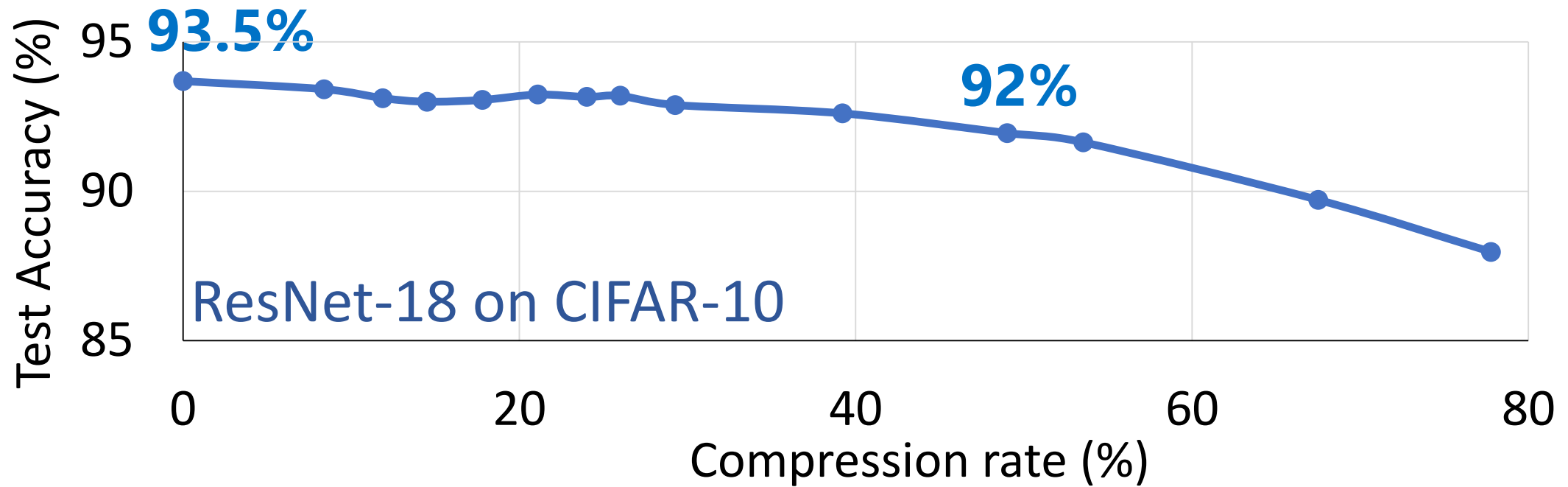
Effects of band-limiting on accuracy



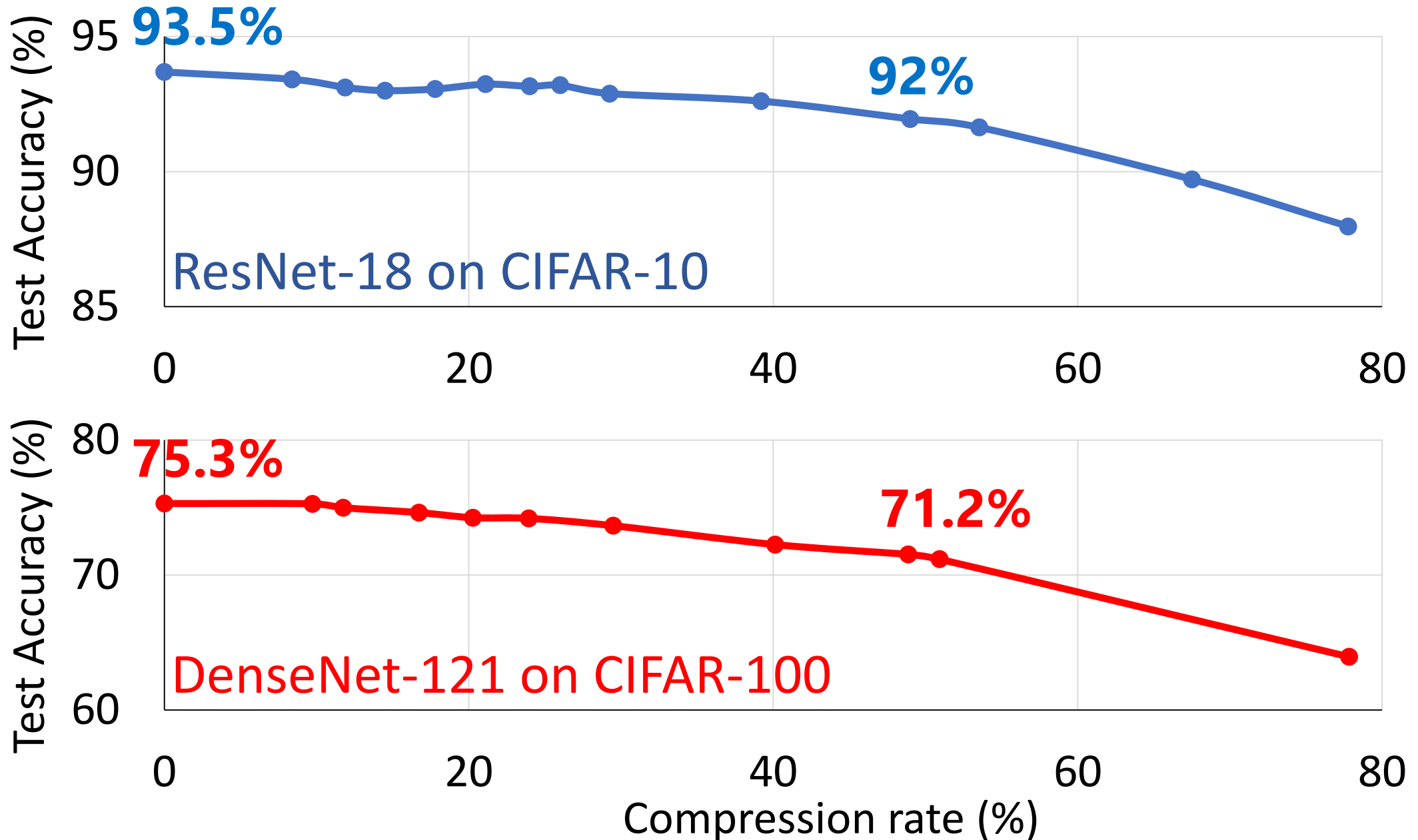
Effects of band-limiting on accuracy



Effects of band-limiting on accuracy



Effects of band-limiting on accuracy



Main take-aways from Band-limited CNNs

- Method to constrain the frequency band in convolution.
- Models trained with band-limiting **gracefully degrade** the accuracy as the function of the compression rate.
- Effectively control resource usage (GPU/CPU and memory).
- The low frequency coefficients learned first during training.
- The same compression rate applied to training and inference.
- The more band-limited model, the more **robust to attacks**.
- Applicable to **other domains**: time-series & speech data.

Thank you

Poster: 6:30-9:00 PM @ Pacific Ballroom #132

github.com/adam-dziedzic/bandlimited-cnns

ady@uchicago.edu

Backup



Why is FFT based convolution important?

- The theoretical properties of the Fourier domain are well understood. No such properties in other domains (Winograd).
- ResNet and DenseNet architectures use 7x7 filters in first layers.
- FFT based convolution can be combined with spectral pooling.
- Band-limiting in the first layer serves as a simple defense.
- A standard algorithm included in popular frameworks (cuDNN).
- Gradient acts as a large filter in the backward pass.
- Zlateski et al. suggest using FFT based convolution on CPUs.
- The 1D FFT convolution for DSP where large filters are used.

Band-limited FFT based convolution *formally*

Cross-correlate input data and filter: $x *_c y$

$$F_x[\omega] = F(x[n]) \quad F_y[\omega] = F(y[n])$$

$$x *_c y = F^{-1}(F_x[\omega] \odot F_y[\omega])$$

Spectrum of convolution: $S[\omega] = F_x[\omega] \odot F_y[\omega]$

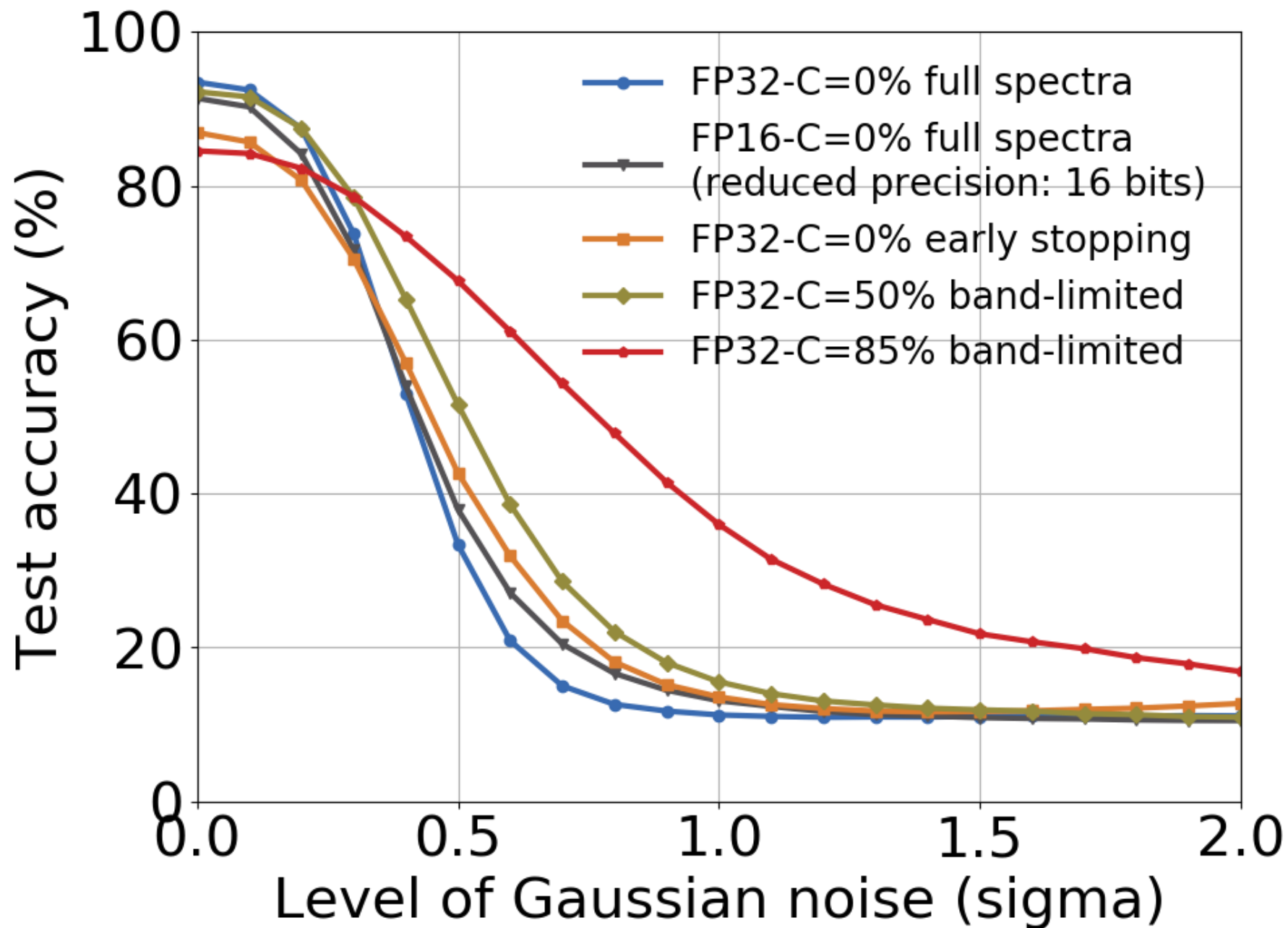
$$M_c[\omega] = \begin{cases} \mathbf{1}, & \omega \leq c \\ \mathbf{0}, & \omega > c \end{cases}$$

$$x *_c y = F^{-1}[(F_x[\omega] \odot M_c[\omega]) \odot (F_y[\omega] \odot M_c[\omega])]$$

$$x *_c y = F^{-1}(S[\omega] \odot M_c[\omega])$$

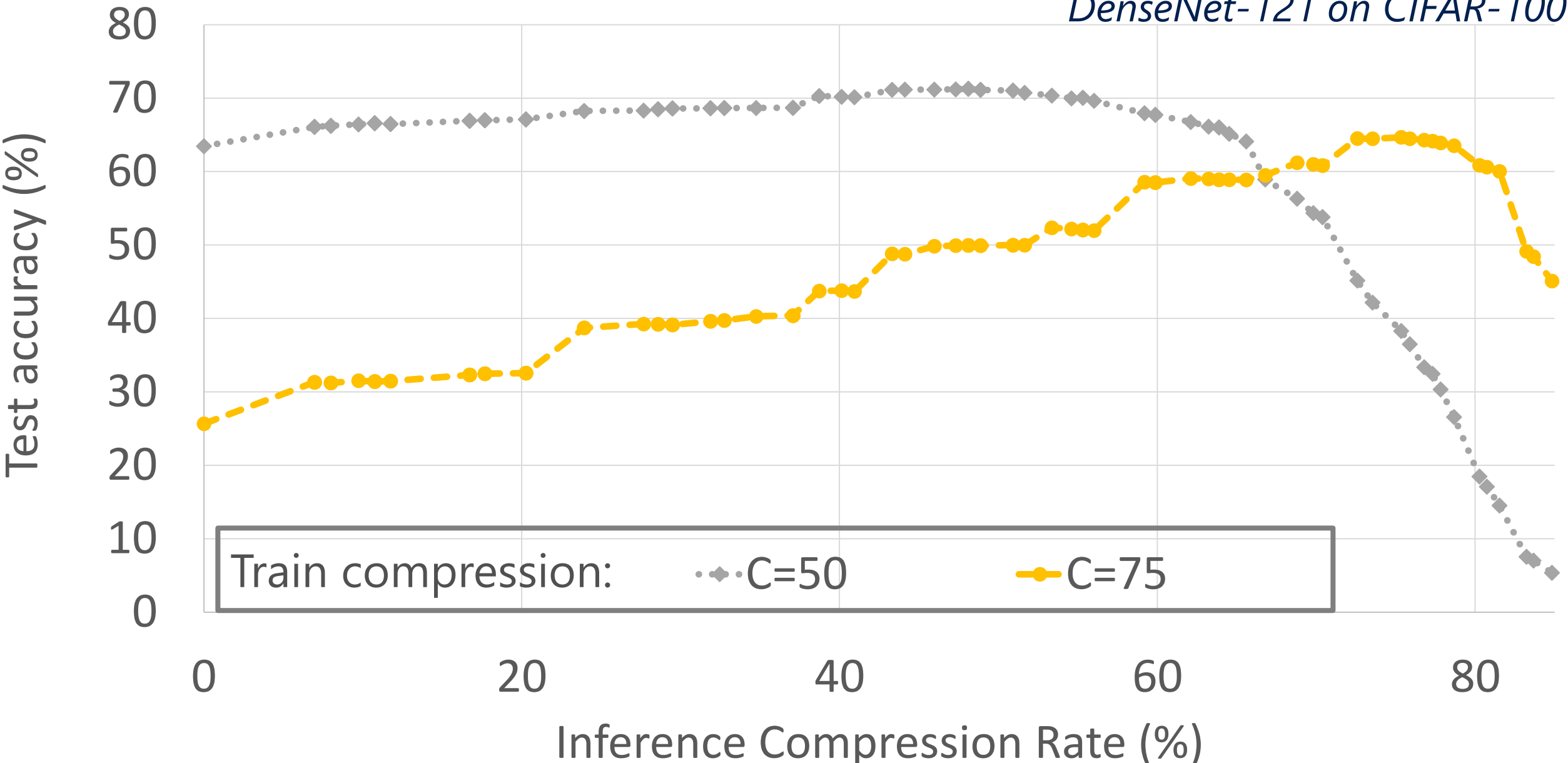
Energy (Parseval's theorem): $\sum_{n=0}^{N-1} |x[n]|^2 = \sum_{\omega=0}^{2\pi} |F_x(\omega)|^2$ ₃₁

Robustness to noise



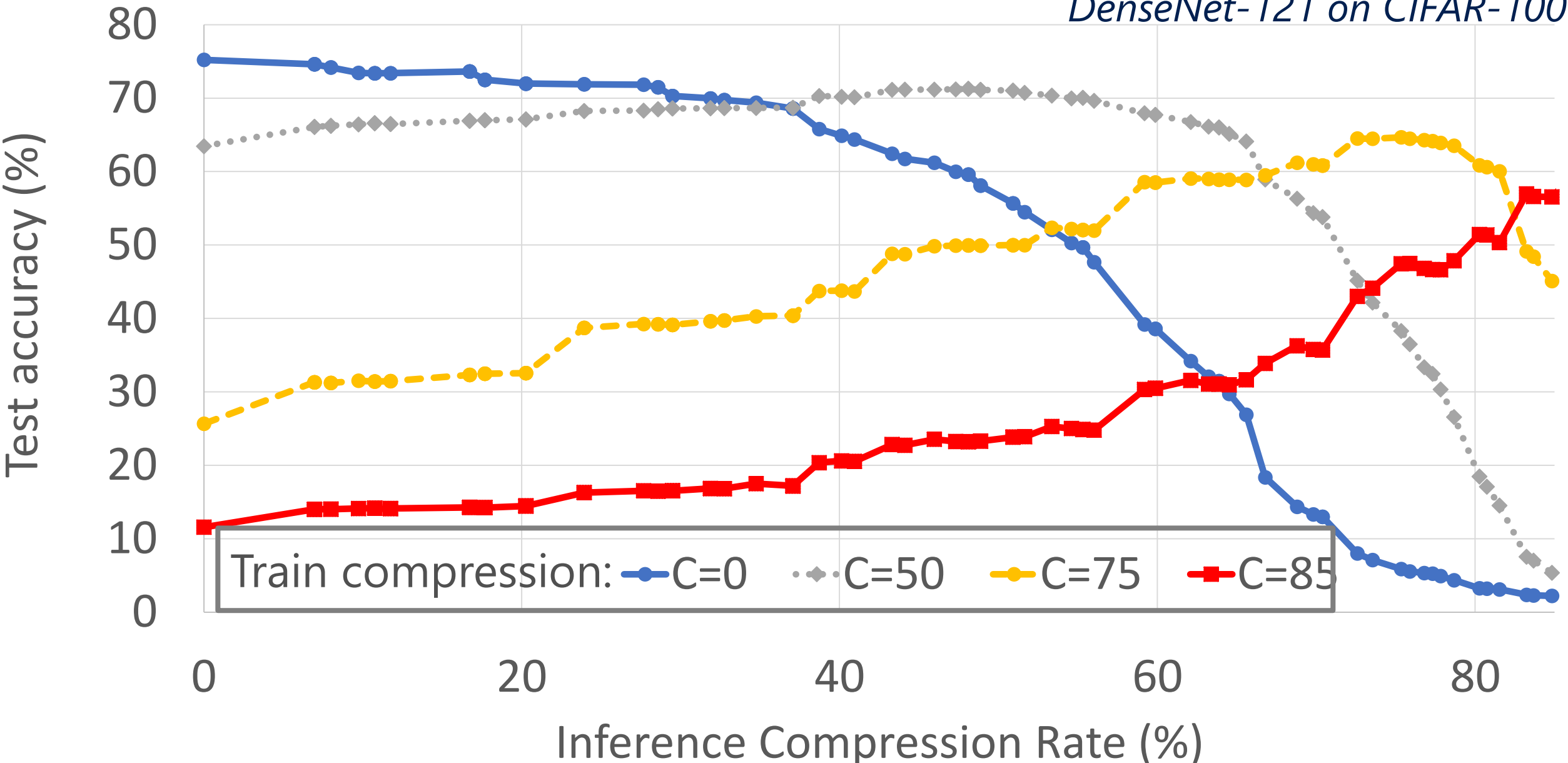
Compression Rate for Training vs Inference

DenseNet-121 on CIFAR-100



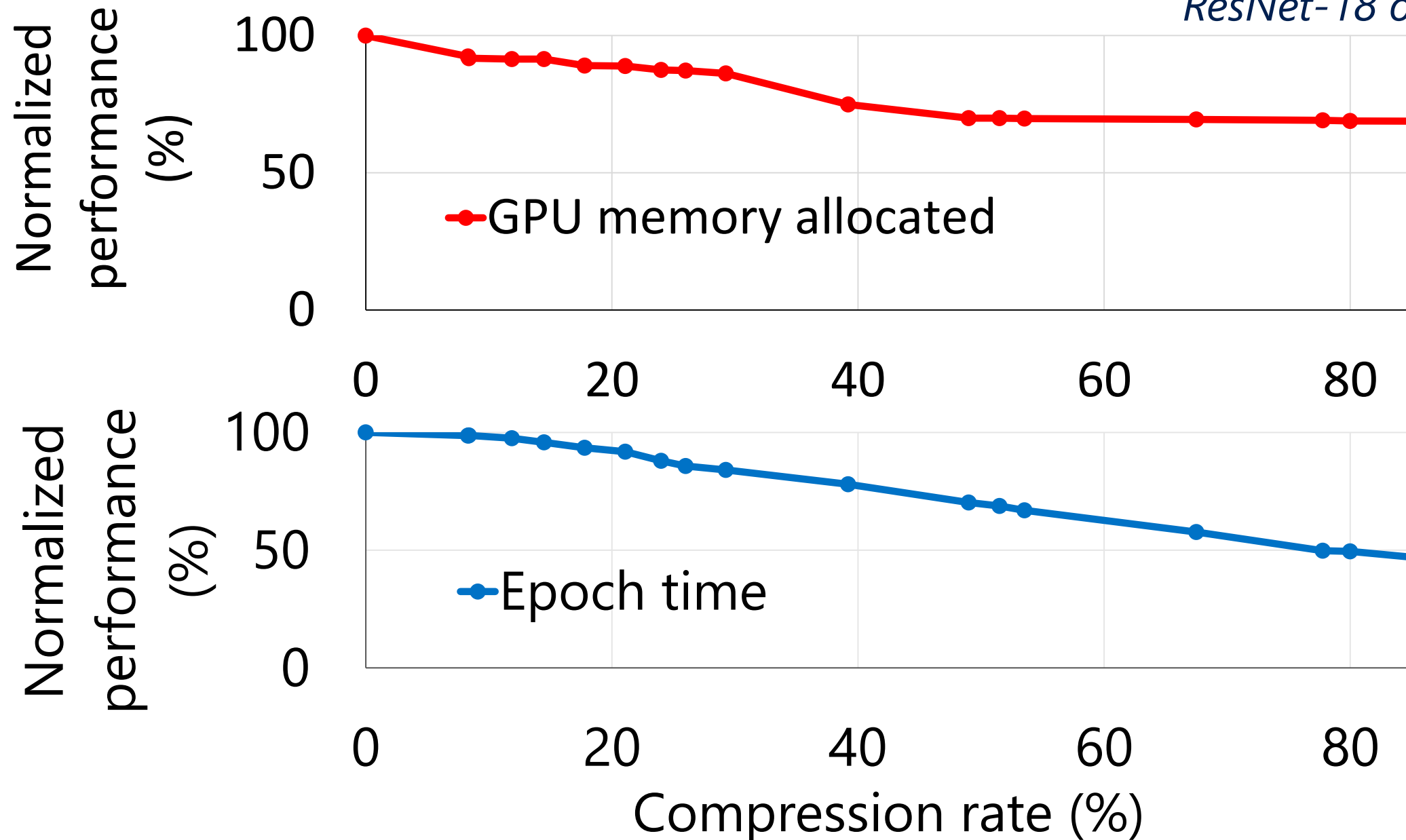
Compression Rate for Training vs Inference

DenseNet-121 on CIFAR-100

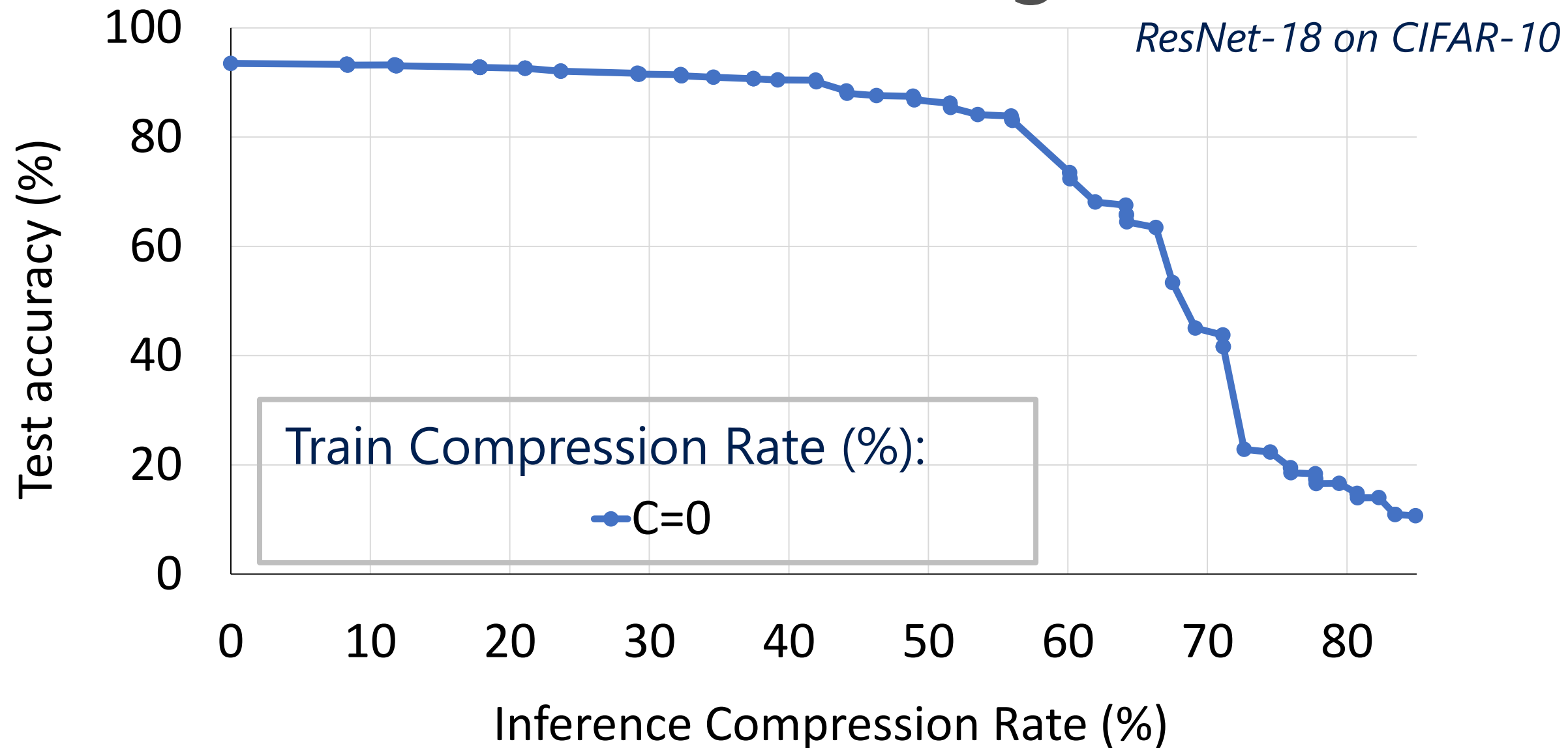


Effectively control resource usage

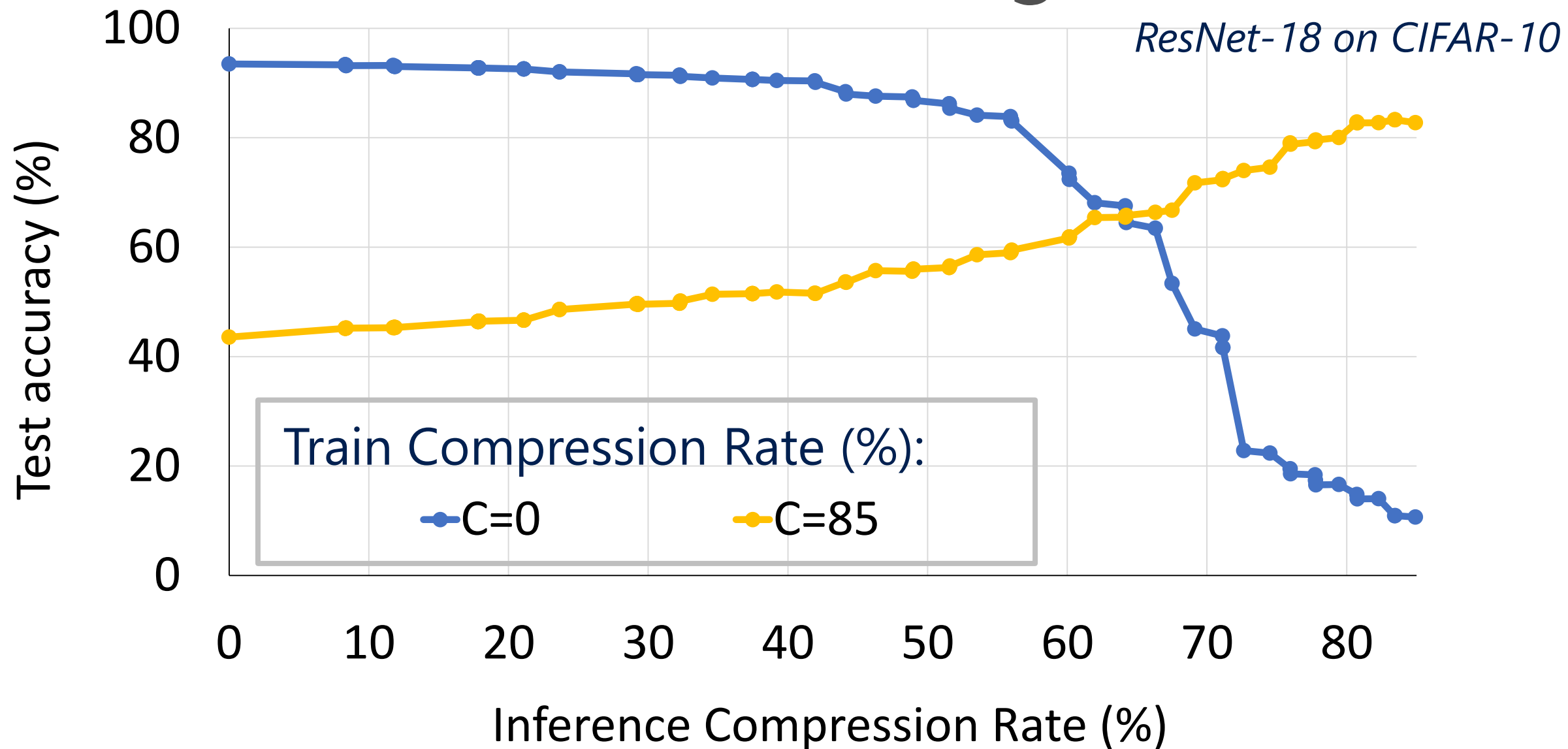
ResNet-18 on CIFAR-10



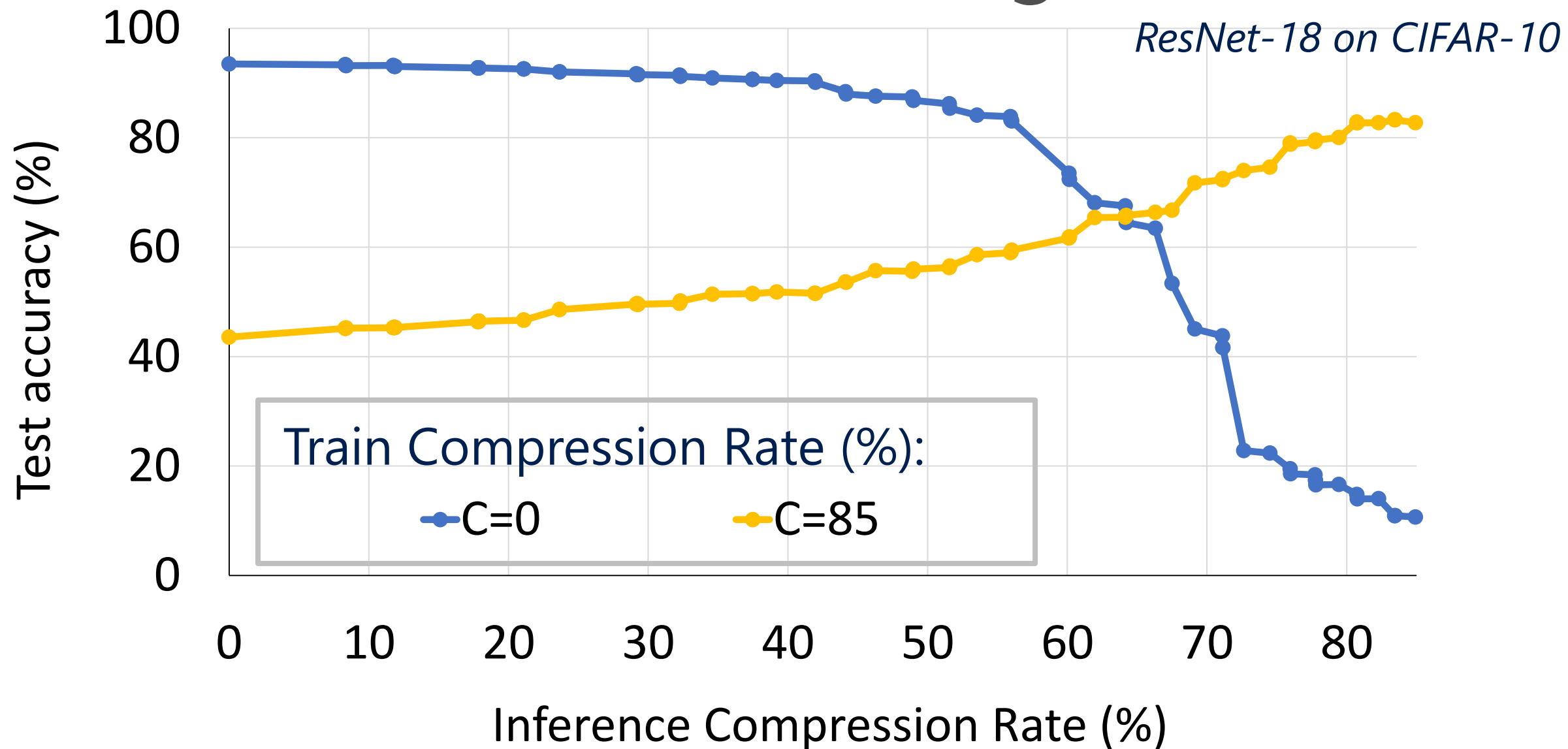
Compression Rate for Training vs Inference



Compression Rate for Training vs Inference

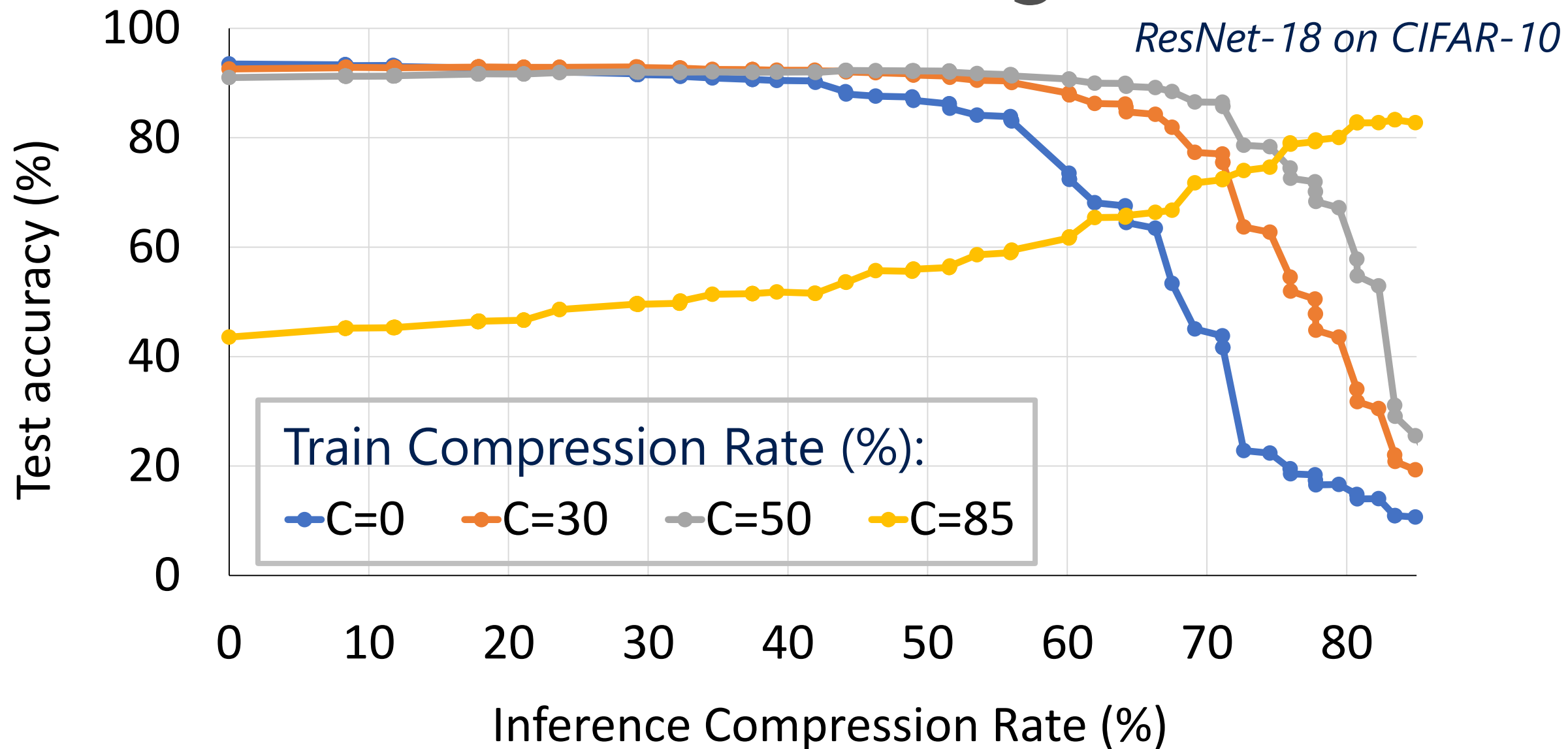


Compression Rate for Training vs Inference



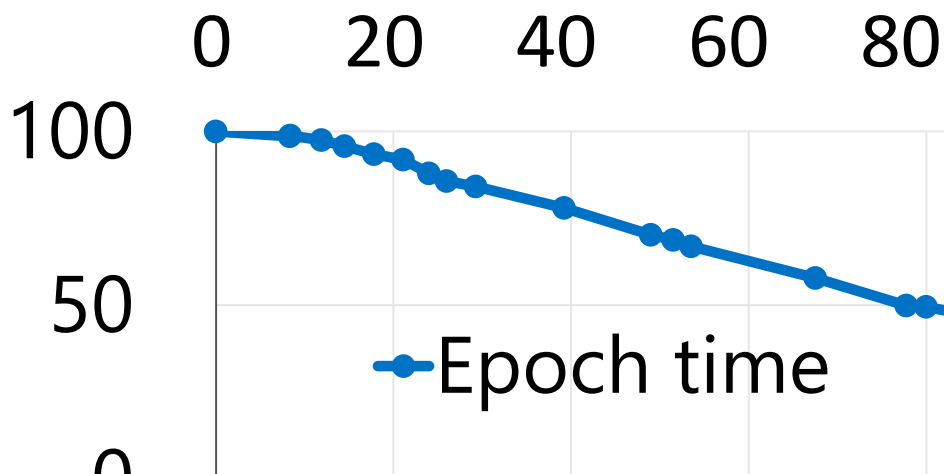
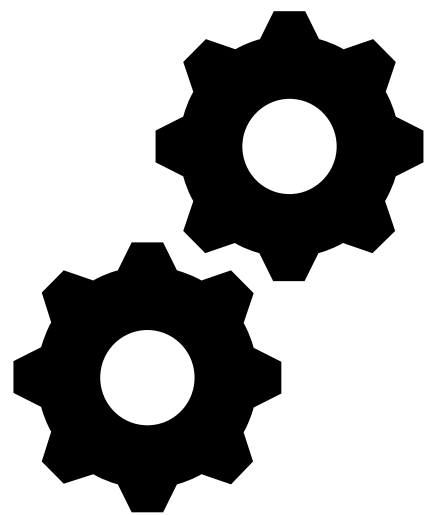
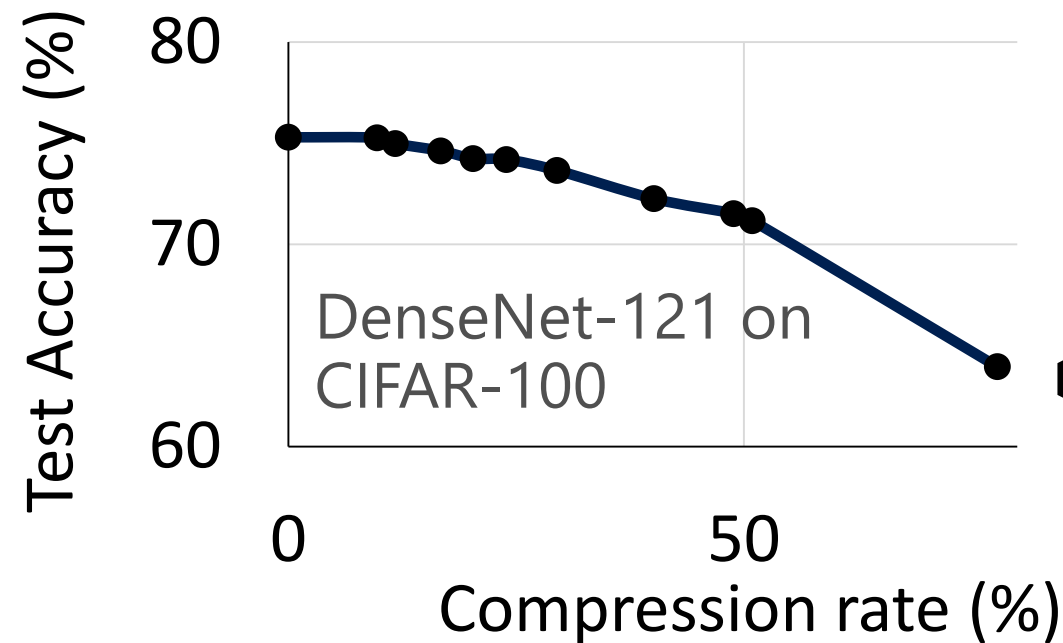
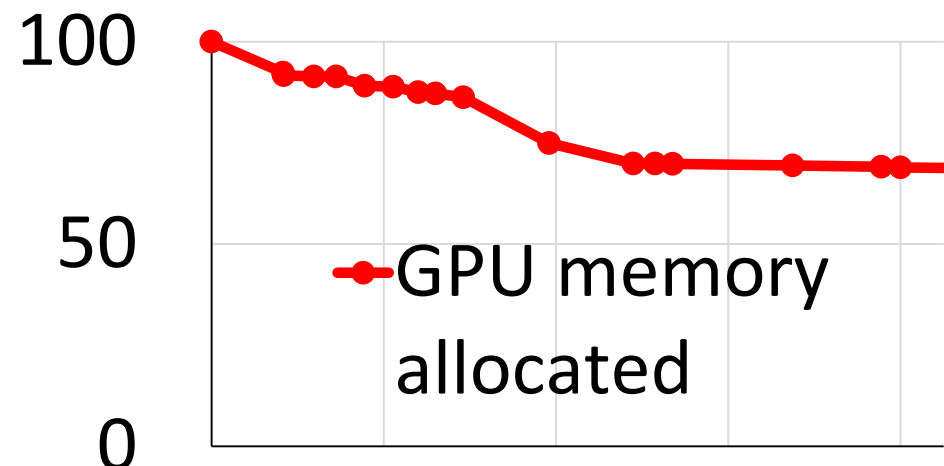
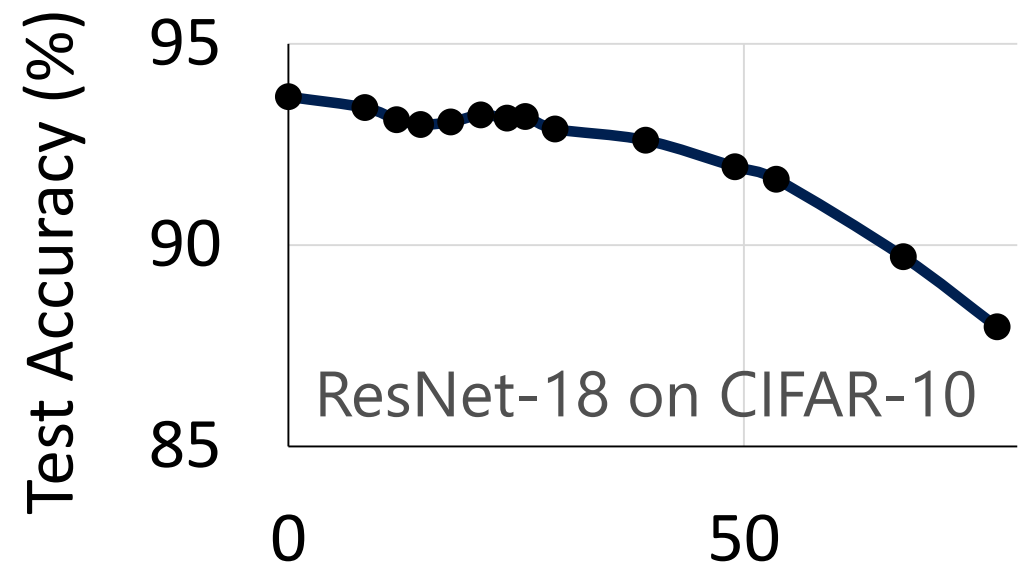
Smooth degradation of accuracy during inference

Compression Rate for Training vs Inference



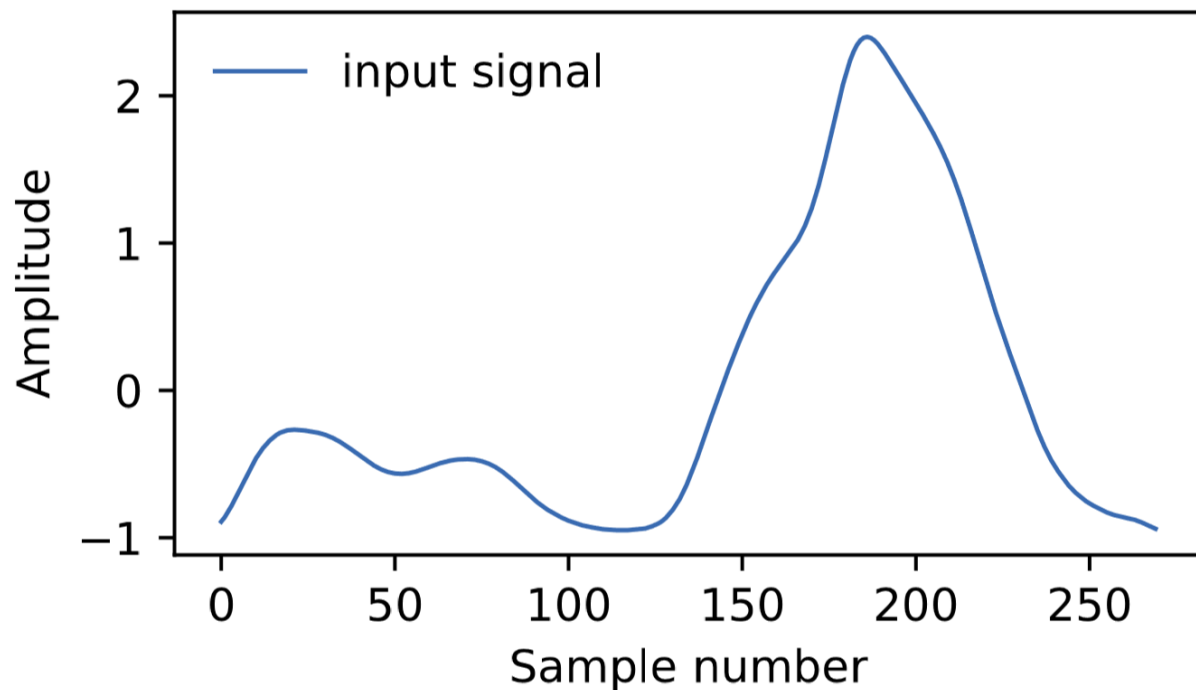
Apply the same compression rate to training and inference

Tuning: Accuracy vs Higher Performance

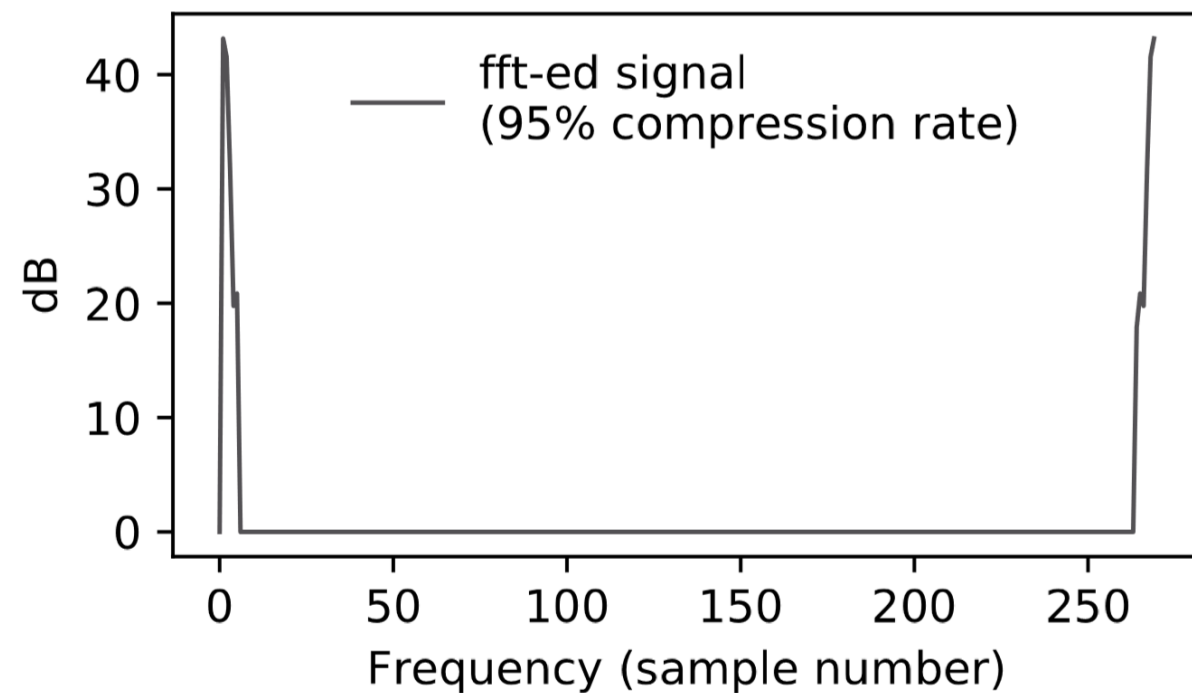
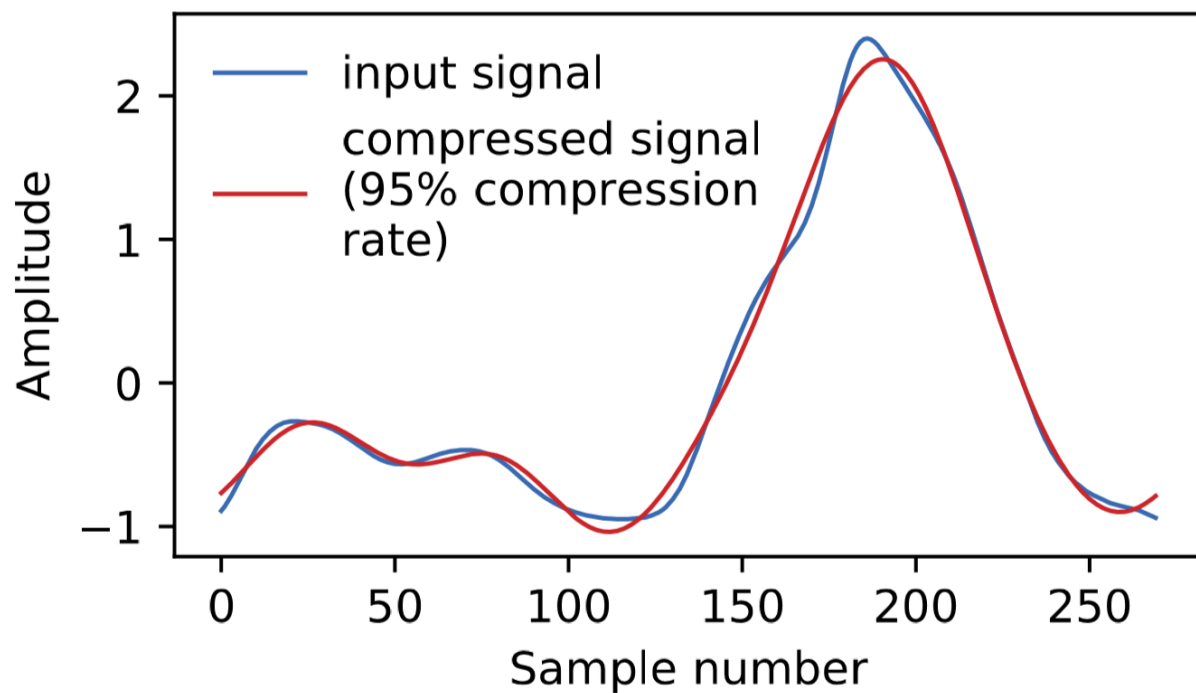
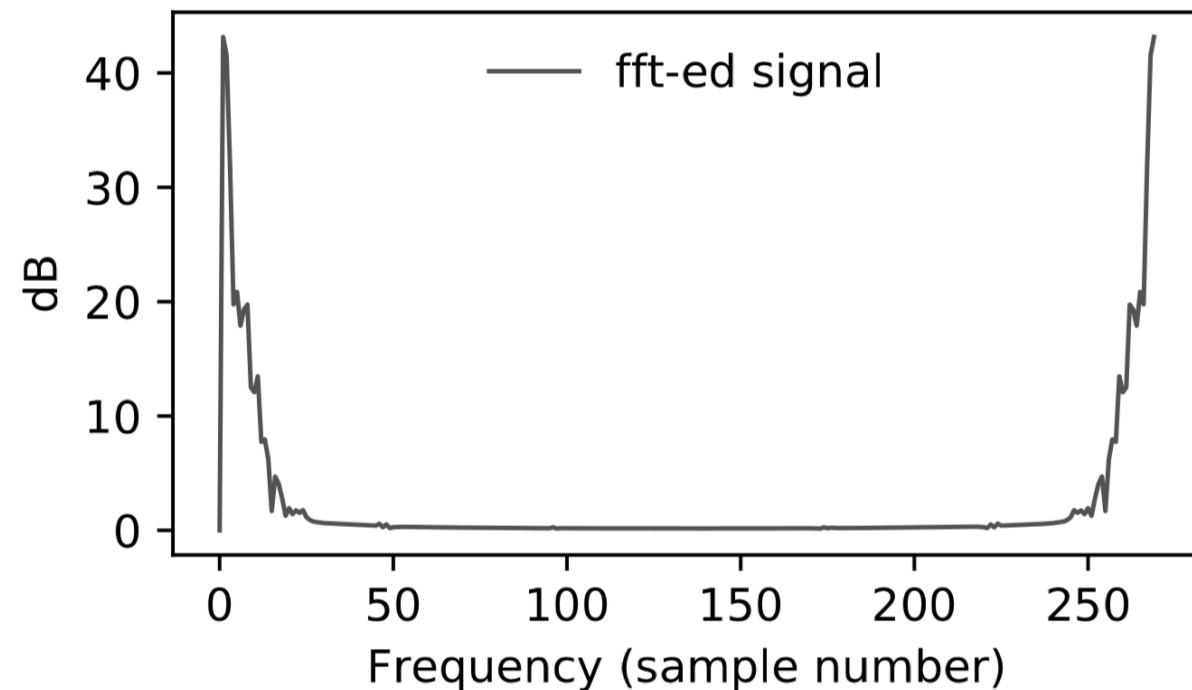


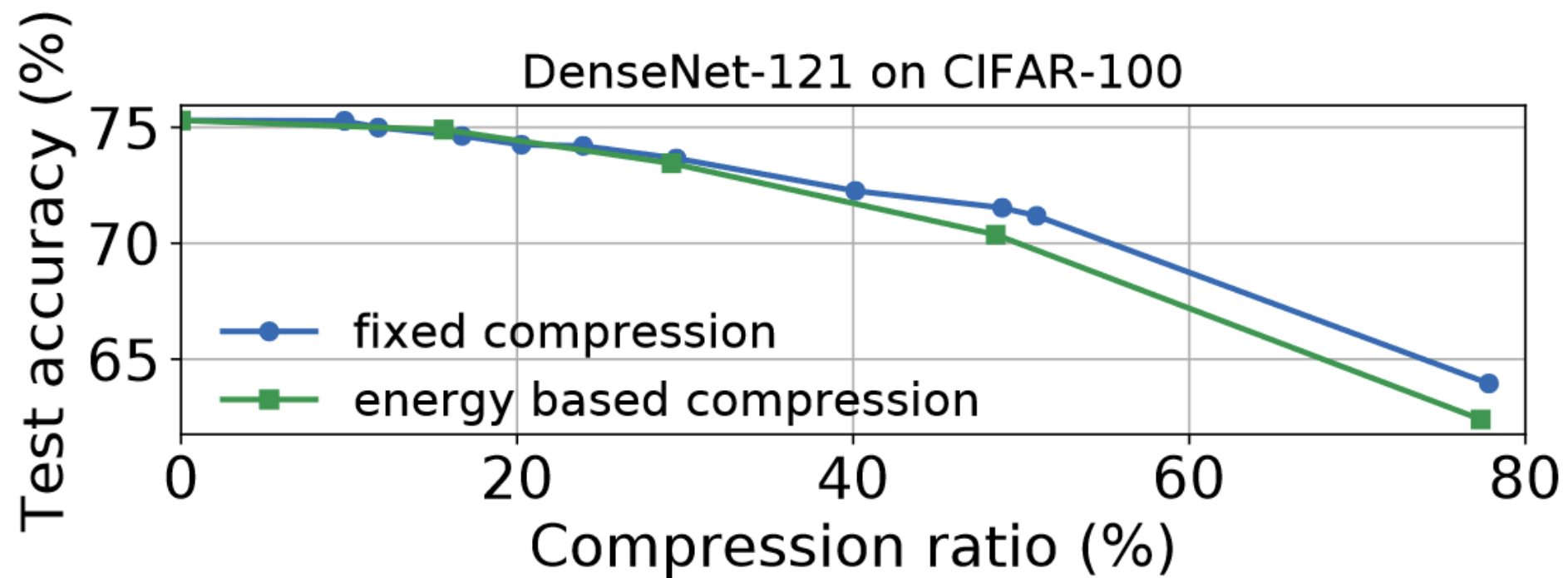
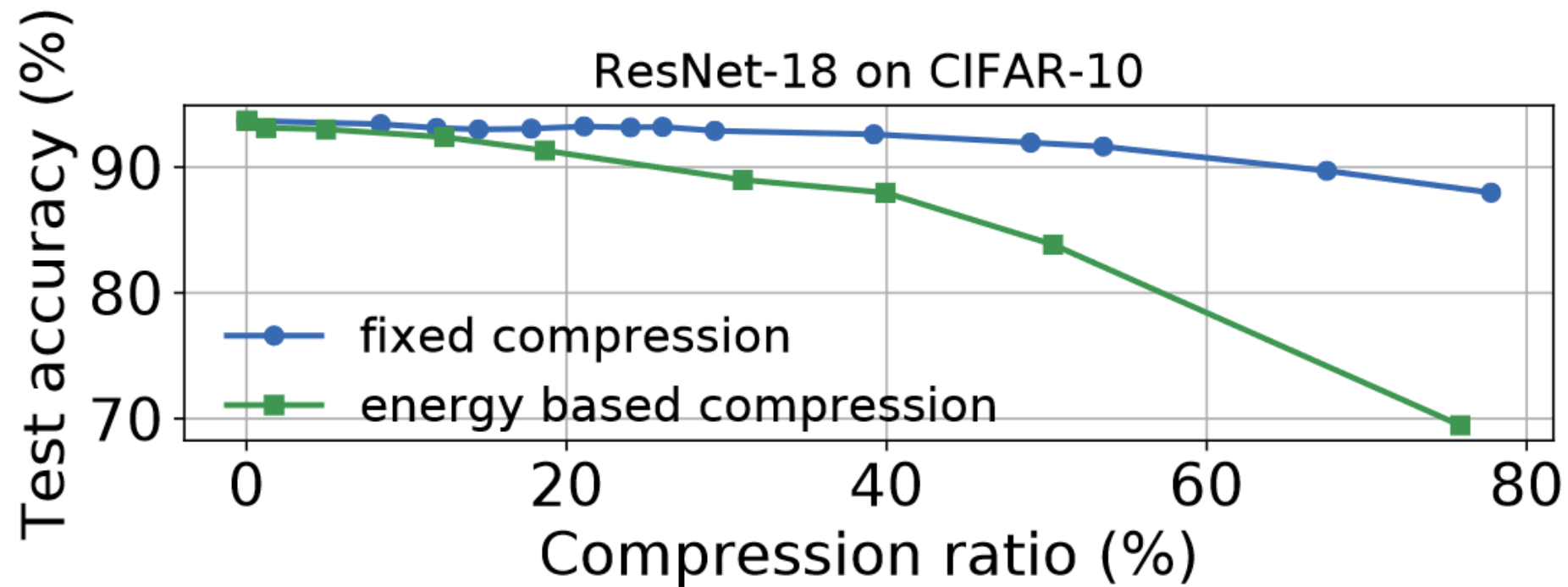
Compression rate (%) 40

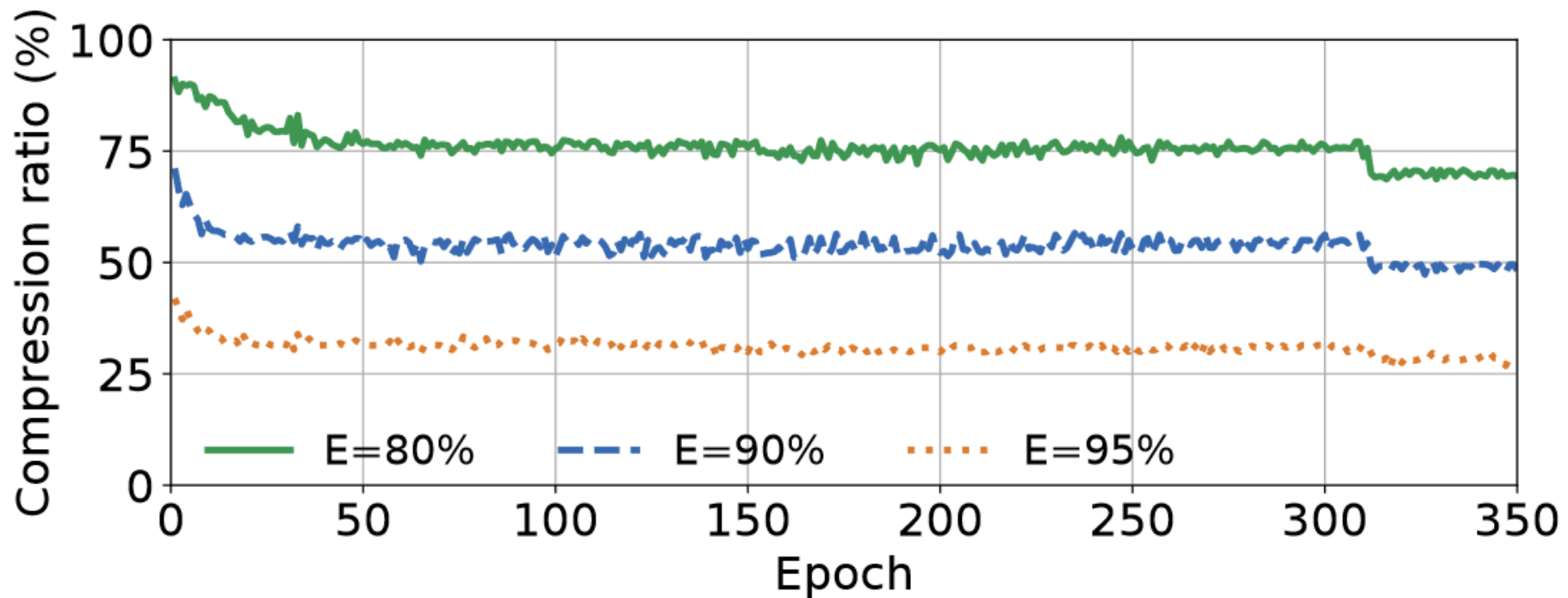
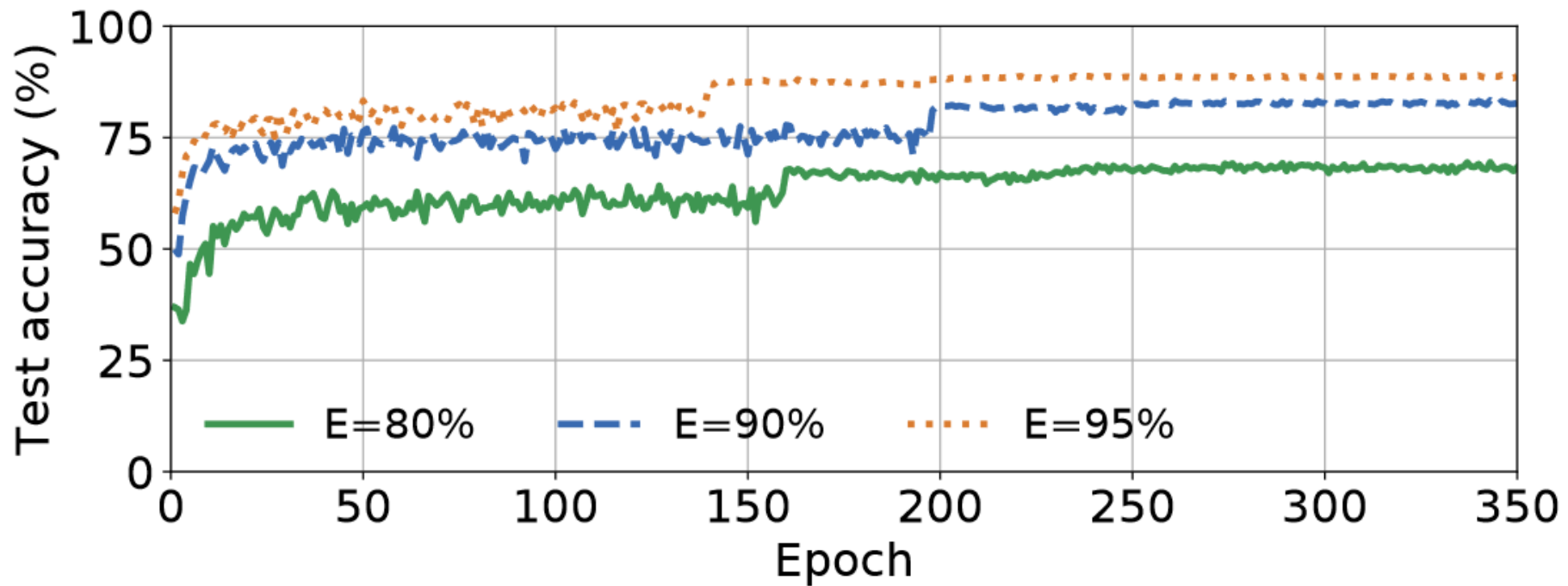
Time domain

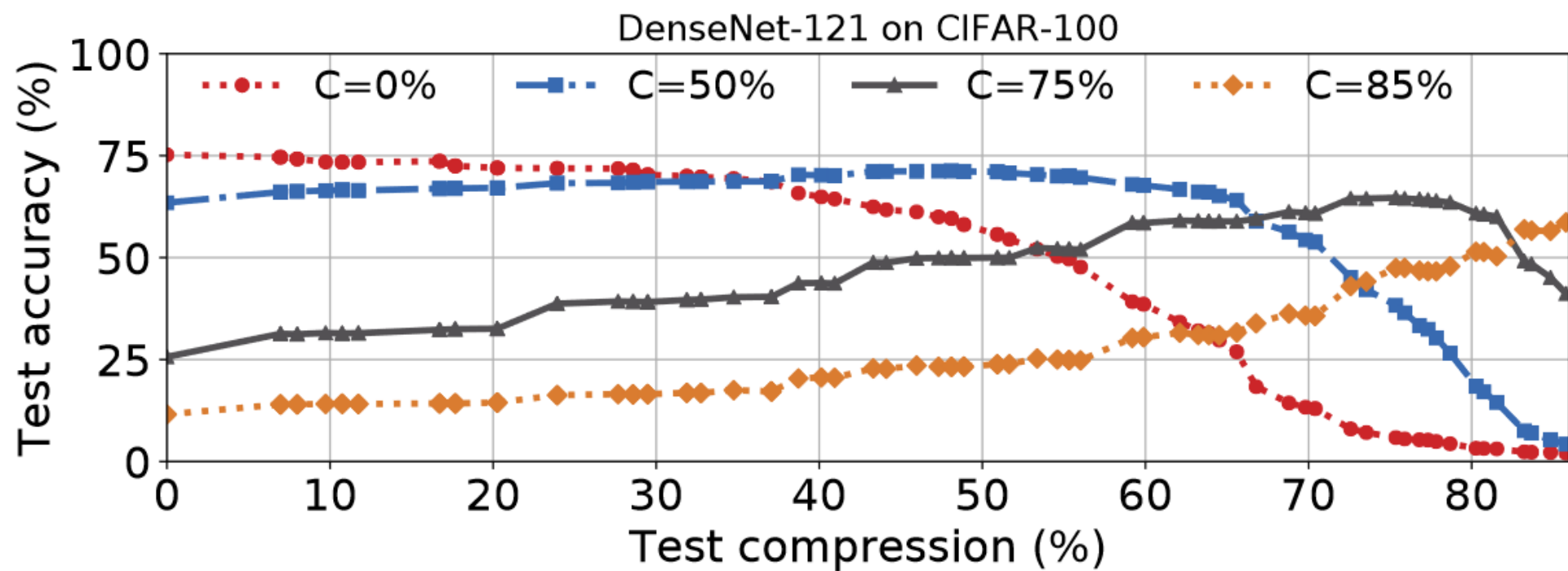
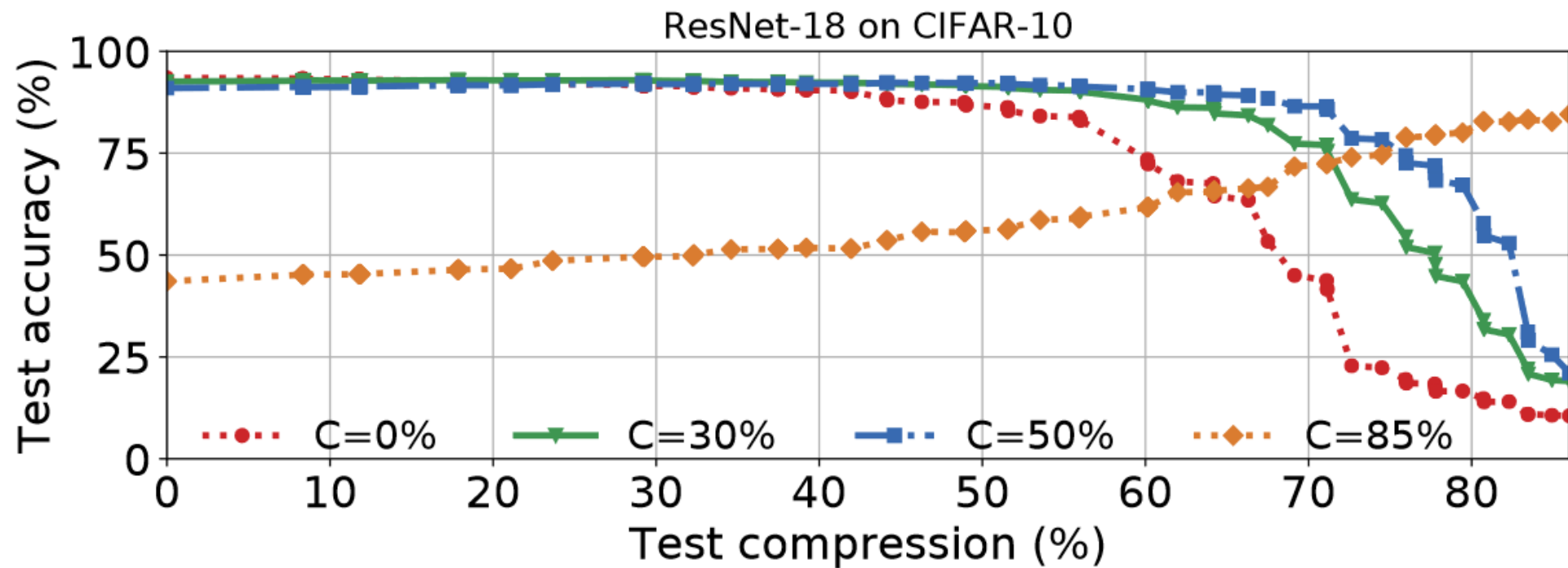


Frequency domain

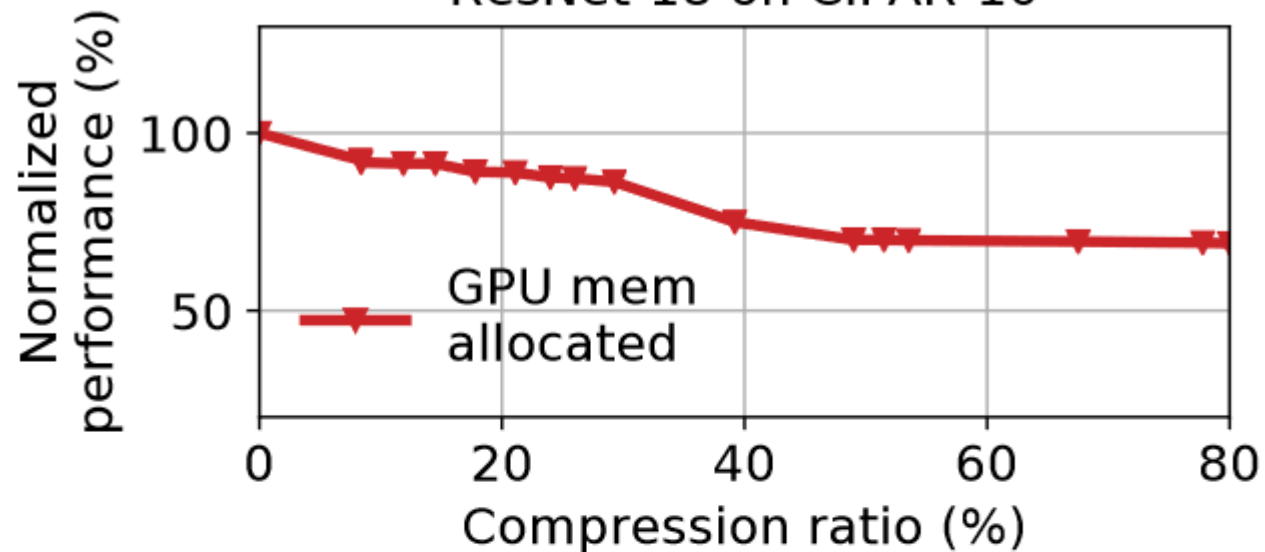




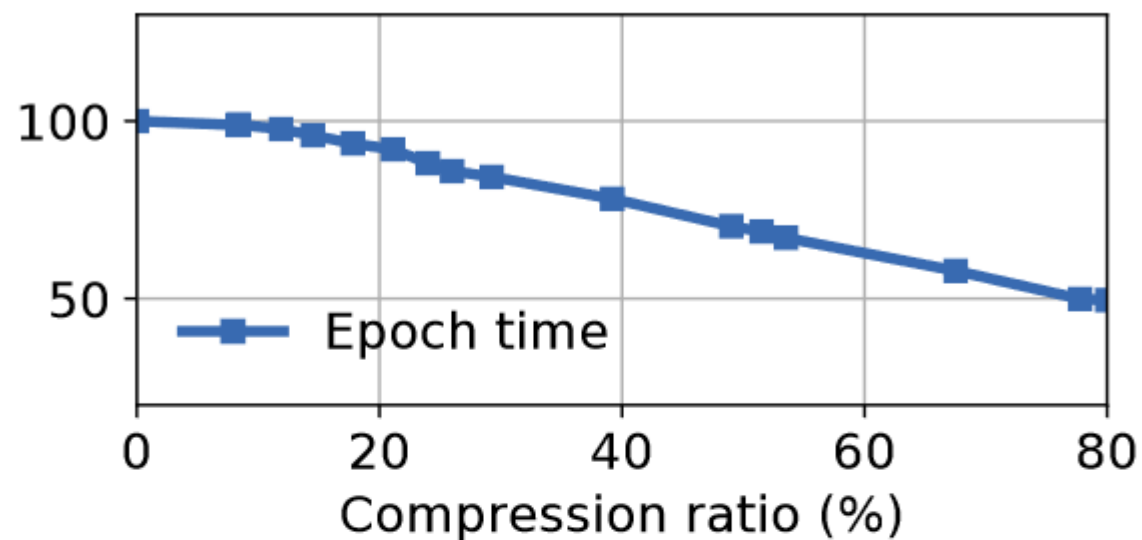




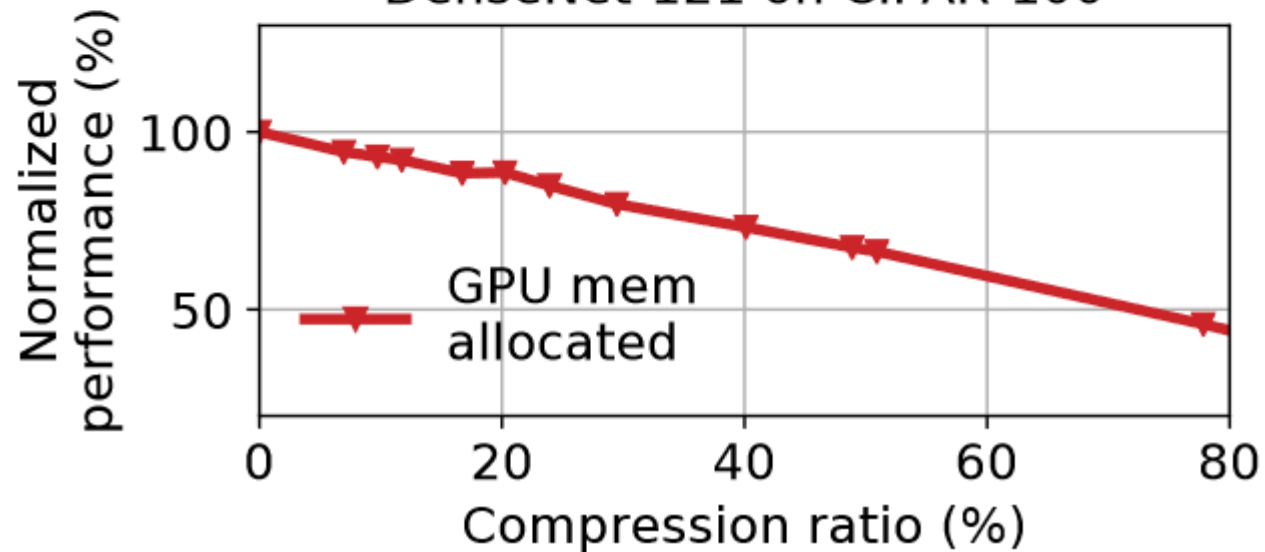
ResNet-18 on CIFAR-10



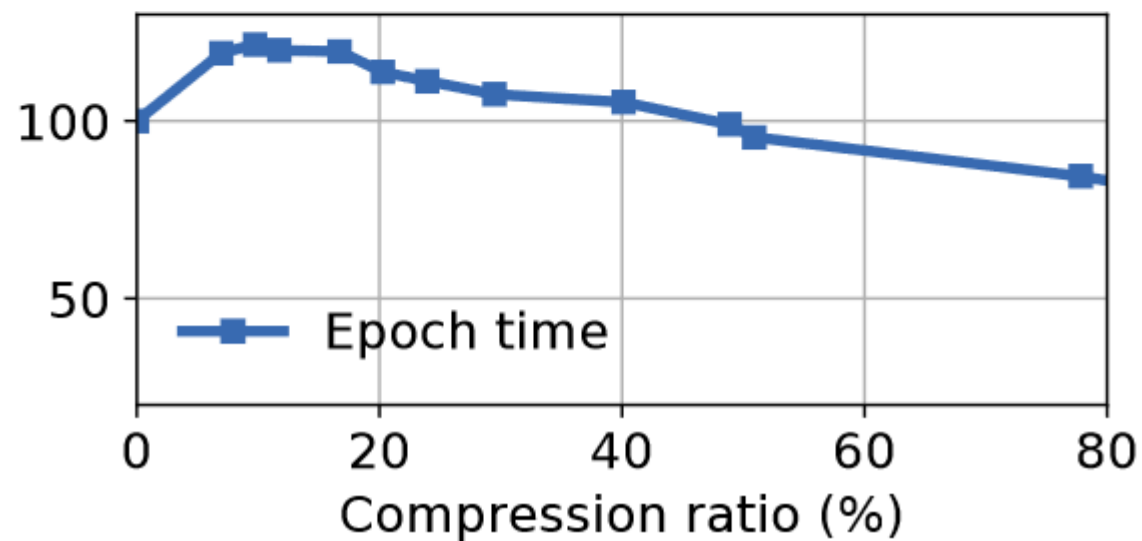
ResNet-18 on CIFAR-10



DenseNet-121 on CIFAR-100



DenseNet-121 on CIFAR-100



“Speaking of longer term, it would be nice if the community migrated to a fully open sourced implementation for all of this [convolution operations, etc.]. This stuff is just too important to the progress of the field for it to be **locked away in proprietary implementations**. The more people working together on this the better for everyone. There's plenty of room to compete on the hardware implementation side.”

Scott Gray

<https://github.com/soumith/convnet-benchmarks/issues/93>