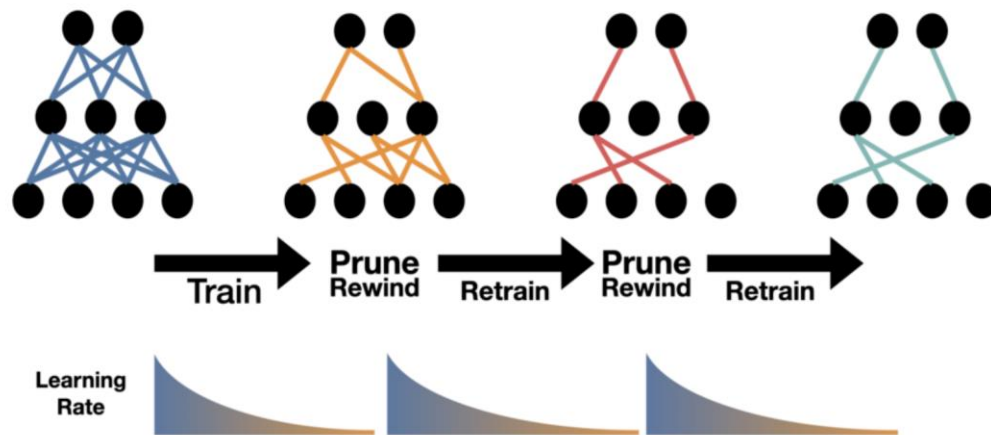Patricia Gschoßmann

# Model Pruning

# **Standard Pruning**

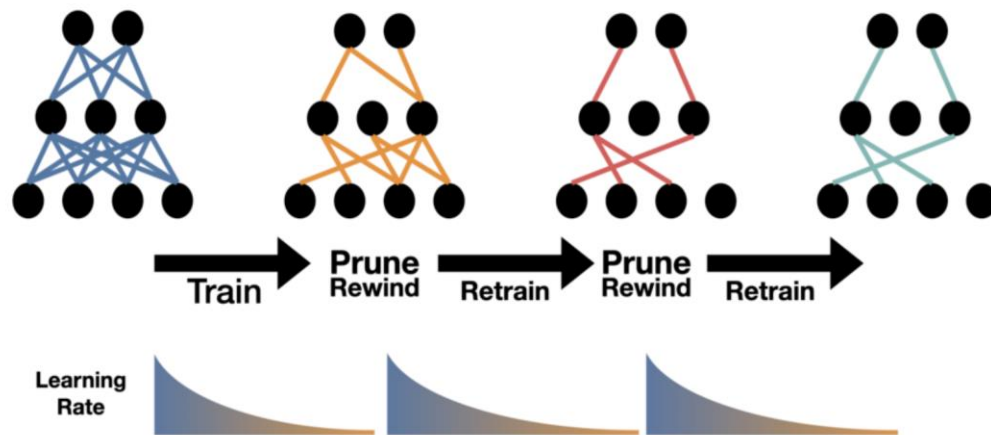How to?

# **Standard Pruning**
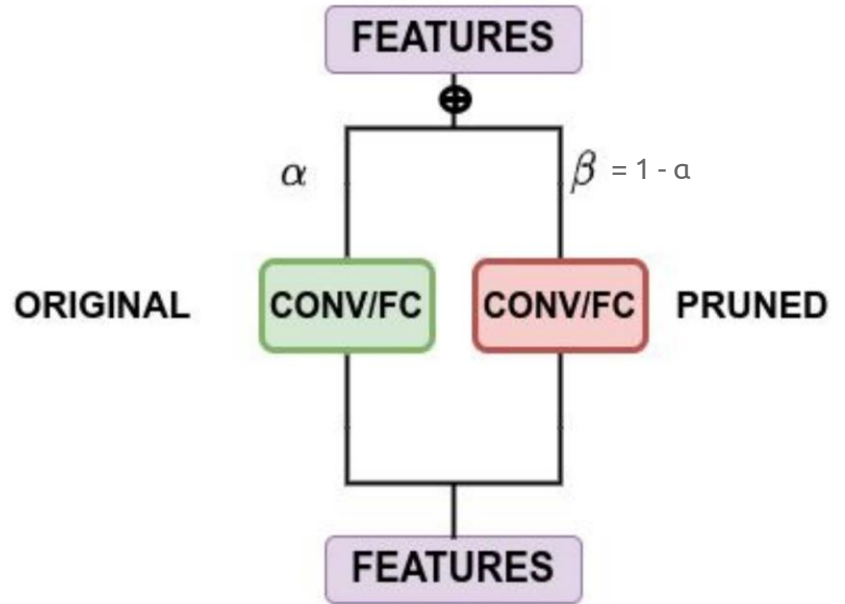
How to?

How can you speed up
the process?



https://news.mit.edu/2020/foolproof-way-shrink-deep-learning-models-0430

# New Approach

How to?

# New Approach

How to?

Advantage: Start at e.g. 60%



FEATURES

$\alpha$         $\beta$ = 1 - α

ORIGINAL   CONV/FC   CONV/FC   PRUNED

FEATURES

Learning for Self Driving Cars - Tutorial Week 7

# Original model

# Original model

## Original model

64

Conv1
$128 \times 64 \times 3 \times 3$

128

ReLU

128

Conv2
$128 \times 128 \times 3 \times 3$

128

## With parallel branches

64

Conv1
$128 \times 64 \times 3 \times 3$

128

128

ReLU

128

ConvA
$64 \times 128 \times 1 \times 1$

64

ReLU

64

ConvB
$128 \times 64 \times 1 \times 1$

128

$\alpha$

$1 - \alpha$

$+$

128

Conv2
$128 \times 128 \times 3 \times 3$

128

**Original model**

Conv1
$128 \times 64 \times 3 \times 3$

128

ReLU

128

Conv2
$128 \times 128 \times 3 \times 3$

128

**With parallel branches**

64

Conv1
$128 \times 64 \times 3 \times 3$

128

torch.einsum("**a**bcd, e**a**fg -> ebcd", conv1_weights, convA_weights)

128

ReLU

128

128

ConvA
$64 \times 128 \times 1 \times 1$

64

ReLU

64

ConvB
$128 \times 64 \times 1 \times 1$

128

$\alpha$

$1 - \alpha$

$+$

128

torch.einsum("**a**bcd, e**a**fg -> ebfg", convB_weights, conv2_weights)

Conv2
$128 \times 128 \times 3 \times 3$

128

# Original model

Conv1
$128 \times 64 \times 3 \times 3$

64

128

ReLU

128

Conv2
$128 \times 128 \times 3 \times 3$

128

# With parallel branches

64

Conv1
$128 \times 64 \times 3 \times 3$

128

128

ReLU

128

$\alpha$

128

Conv2
$128 \times 128 \times 3 \times 3$

128

ConvA
$64 \times 128 \times 1 \times 1$

64

ReLU

64

ConvB
$128 \times 64 \times 1 \times 1$

128

$1 - \alpha$

$+$

# Pruned model

64

Conv1A
$64 \times 64 \times 3 \times 3$

64

ReLU

64

Conv2B
$128 \times 64 \times 3 \times 3$

64

# Challenges

## α-schedule

- Type of decay (step, exponential, multiplicative, …)
- Decay rate

## Pretrained model

Should the weights of the pretrained model be freezed or updated?

## Learning rate

- How big/small?
- Static or dynamic?
- Schedulers?

# VGG16

Original performance:
- **Training accuracy: 100%**
- **Validation accuracy: 92%, Test accuracy: 91.3%**

# VGG16

Original performance:
- ○ **Training accuracy: 100%**
- ○ **Validation accuracy: 92%, Test accuracy: 91.3%**

Standard pruning results:
- ○ **Validation accuracy: 90%**
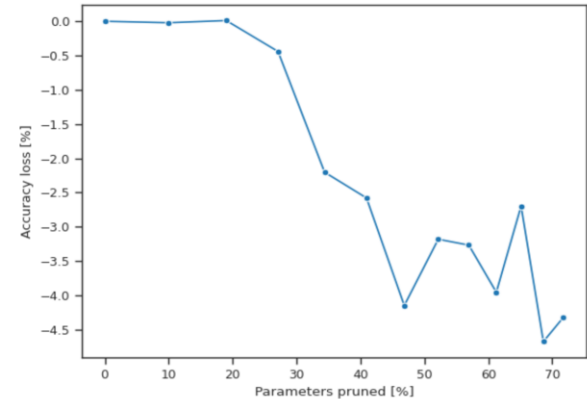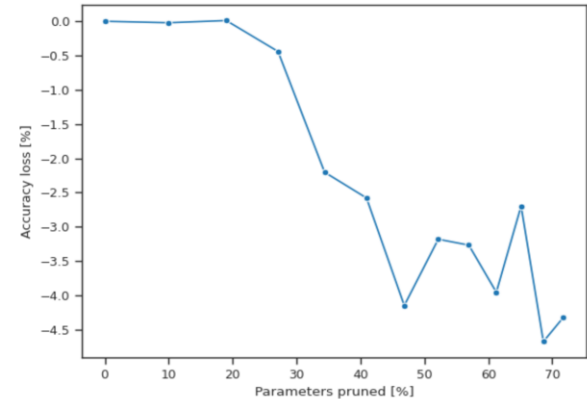- ○ **Pruned model size: ≈ 35% of original**

# VGG16

Original performance:
- **Training accuracy: 100%**
- **Validation accuracy: 92%, Test accuracy: 91.3%**

Standard pruning results:
- **Validation accuracy: 90%**
- **Pruned model size: ≈ 35% of original**

➡ Goal: **Prune 65% at once**

# VGG16: Results

# VGG16: Results

Failed attempts/mistakes:
- **Freeze pretrained + reduceLrOnPlateau-scheduler + constant initial LR (0.001)**
- **Step decay/Exp. decay**

# VGG16: Results

Failed attempts/mistakes:
- **Freeze pretrained + reduceLrOnPlateau-scheduler + constant initial LR (0.001)**
- **Step decay/Exp. decay**

Final setup:
- **Update pretrained + step decay (-0.1) + reduceLrOnPlateau-scheduler**
- **At start: Initial LR of 0.001, at end: Repeated increase of initial LR x10**

# VGG16: Results

Failed attempts/mistakes:
- **Freeze pretrained + reduceLrOnPlateau-scheduler + constant initial LR (0.001)**
- **Step decay/Exp. decay**

Final setup:
- **Update pretrained + step decay (-0.1) + reduceLrOnPlateau-scheduler**
- **At start: Initial LR of 0.001, at end: Repeated increase of initial LR x10**

Results:
- **Validation accuracy: 88.5%**
- **Parameters remaining: ≈ 35%**

| | Original | Pruned |
|---|---|---|
| **Time/img (GPU)** | ≈0.0531s | ≈0.0529s |
| **GPU VRAM** | 1.53GB | 0.53GB |
| **Time/img (CPU)** | ≈0.0409s | ≈0.0257s |
| **RAM consumption** | 3.64GB | 2.32GB |

# RGB Autoencoder

Original performance (input of size 3x128x128, bottleneck of size 128x16x16):
- **Training loss: 1.14**
- **Validation loss: 1.40, Test loss: 1.41**



| Original | Reconstructed | Original | Reconstructed |

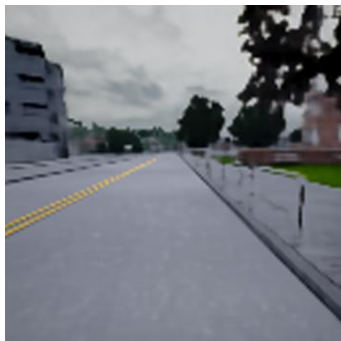# RGB Autoencoder

Original performance (input of size 3x128x128, bottleneck of size 128x16x16):
- **Training loss: 1.14**
- **Validation loss: 1.40, Test loss: 1.41**



Original            Reconstructed                Original            Reconstructed

→ Goal: **Prune 65% at once – bottleneck of size 44x16x16**

# RGB Autoencoder: Approach and Results

# RGB Autoencoder: Approach and Results

Possible approaches:
- **Prune encoder – then decoder**
- **Prune en- and decoder simultaneously**

# RGB Autoencoder: Approach and Results

Possible approaches:
- **Prune encoder – then decoder**
- **Prune en- and decoder simultaneously**

Setup: **Same setup as with VGG**

# RGB Autoencoder: Approach and Results

Possible approaches:
- ○ **Prune encoder – then decoder**
- ○ **Prune en- and decoder simultaneously**
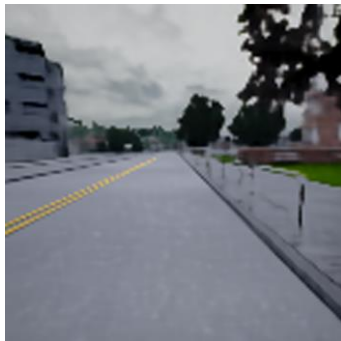
Setup: **Same setup as with VGG**

Results – Prune encoder:
- ● **Training loss: 2.31**
- ● **Validation loss: 2.59**
- ● **Test loss: 2.65**

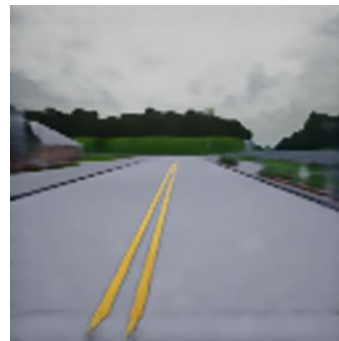|  | Original | Pruned |
|---|---|---|
| **Time/img (GPU)** | ≈0.0573s | ≈0.0582s |
| **GPU VRAM** | 6.83GB | 5.23GB |
| **Time/img (CPU)** | ≈0.0891s | ≈0.0836s |
| **RAM consumption** | 7.97GB | 7.89GB |

# RGB Autoencoder: Results – Prune Encoder
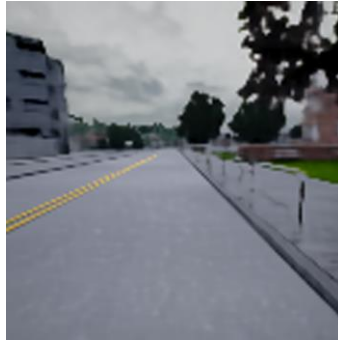


Img. 1

Img. 2

Original

Reconstructed
(before pruning)

Reconstructed
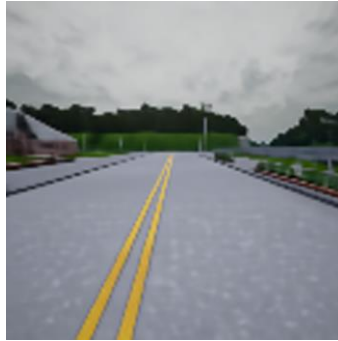(after pruning,
before multiplication)

Reconstructed
(after pruning,
after multiplication)

# RGB Autoencoder: Intermediate Results – Prune Both

Img. 1



Img. 2



Original

Reconstructed
(before pruning)

Reconstructed
(after pruning,
before multiplication)

# Conclusion

## α-schedule

- Type of decay (step, exponential, multiplicative, ...)
- Decay rate

## Pretrained model

Should the weights of the pretrained model be freezed or updated?

## Learning rate

- How big/small?
- Static or dynamic?
- Schedulers?

# Conclusion

## α-schedule

- Type of decay (step, exponential, multiplicative, …)
- Decay rate

**Open question**

## Pretrained model

Should the weights of the pretrained model be freezed or updated?

## Learning rate

- How big/small?
- Static or dynamic?
- Schedulers?

# Conclusion

## α-schedule

- Type of decay (step, exponential, multiplicative, ...)
- Decay rate

**Open question**

## Pretrained model

Should the weights of the pretrained model be freezed or updated?

**Update weights**

## Learning rate

- How big/small?
- Static or dynamic?
- Schedulers?

# Conclusion

## α-schedule

- Type of decay (step, exponential, multiplicative, …)
- Decay rate

**Open question**

## Pretrained model

Should the weights of the pretrained model be freezed or updated?

**Update weights**

## Learning rate

- How big/small?
- Static or dynamic?
- Schedulers?

**Initial LR of 0.001
Increase at local minima**

# Future Work

- Prune autoencoder decoder
- Compare results of both autoencoder approaches

- …

- Prune autoencoder with skip connections
- In general: Test new pruning approach with a larger α-decay

- …

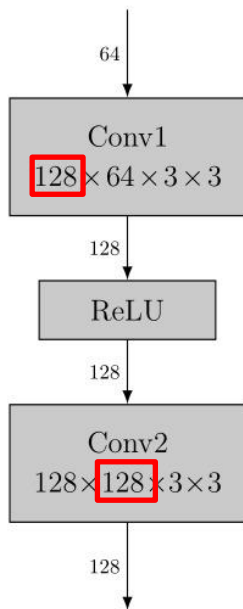- Goal: Prune DCGAN

# Thanks!

Does anyone have any questions?

# Sources

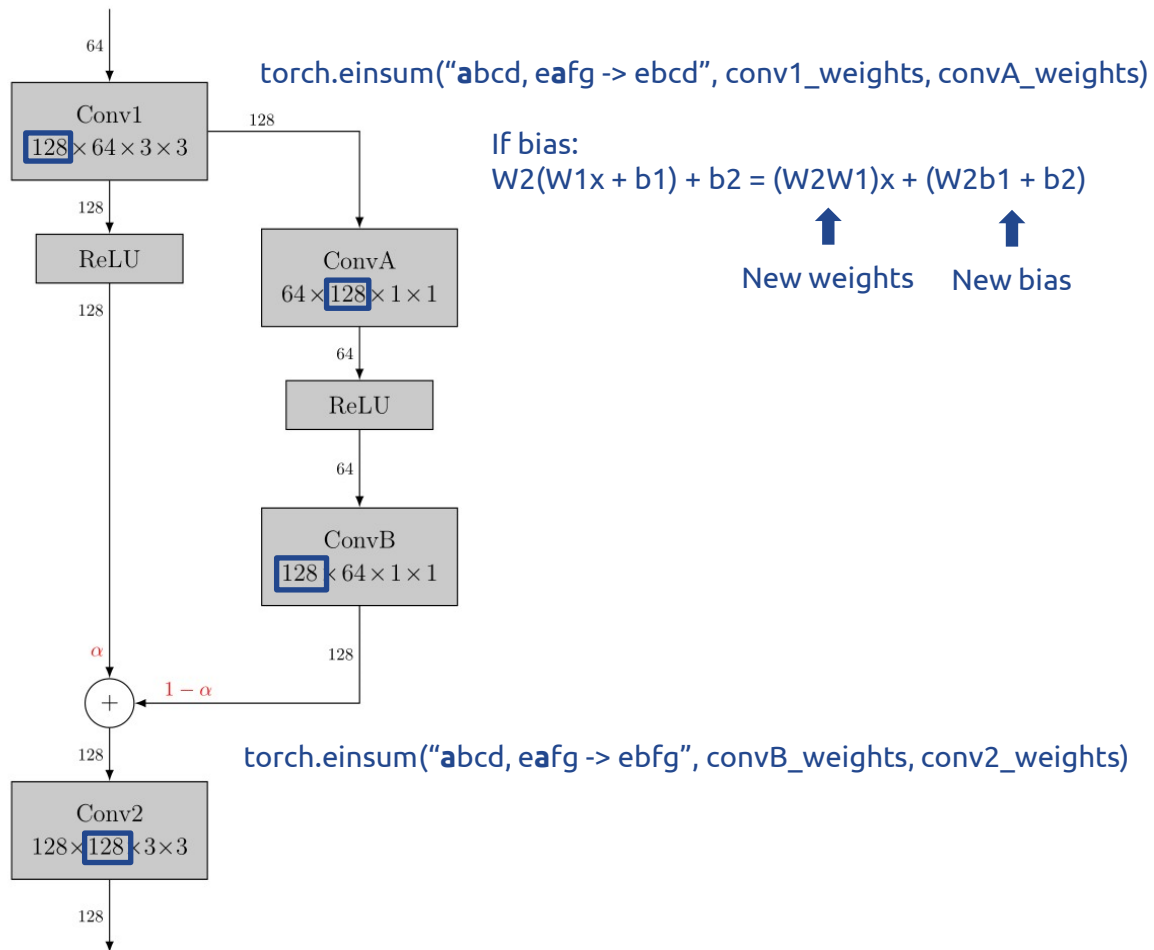Presentation template by Slidesgo

# Original model
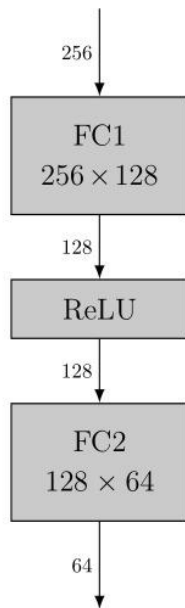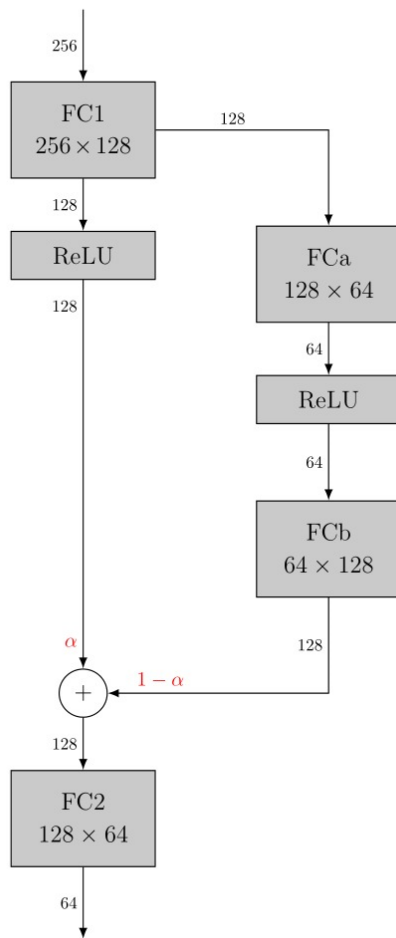
# With parallel branches



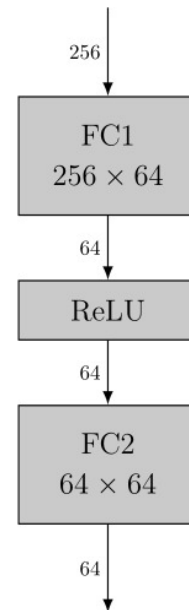torch.einsum("**a**bcd, e**a**fg -> ebcd", conv1_weights, convA_weights)

If bias:
$W2(W1x + b1) + b2 = (W2W1)x + (W2b1 + b2)$

New weights        New bias

torch.einsum("**a**bcd, e**a**fg -> ebfg", convB_weights, conv2_weights)

**Original model**

256

FC1
$256 \times 128$

128

ReLU

128

FC2
$128 \times 64$

64

**With parallel branches**

256

FC1
$256 \times 128$

128

128

ReLU

128

FCa
$128 \times 64$

64

ReLU

64

FCb
$64 \times 128$

128

$\alpha$

$1 - \alpha$

+

128

FC2
$128 \times 64$

64

**Pruned model**

256

FC1
$256 \times 64$

64

ReLU

64

FC2
$64 \times 64$

64

# Original model



# With parallel branches

# Pruned model

64 × 1 × 1

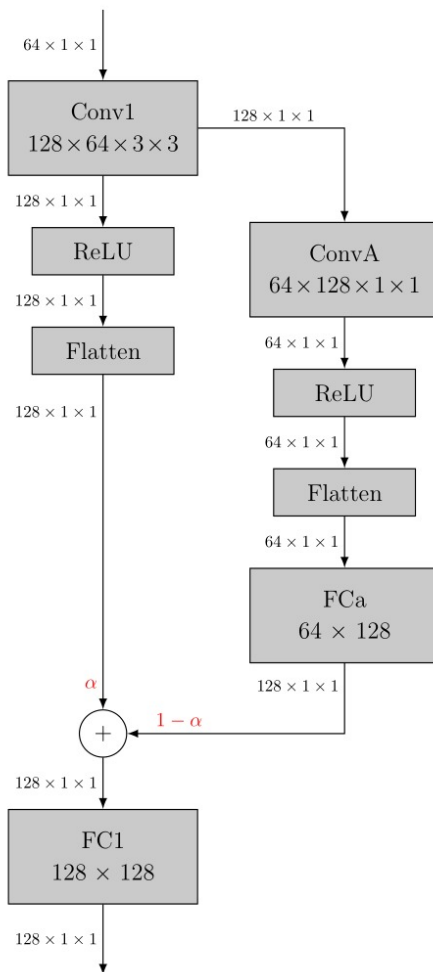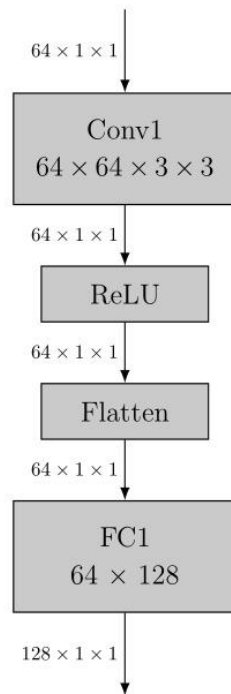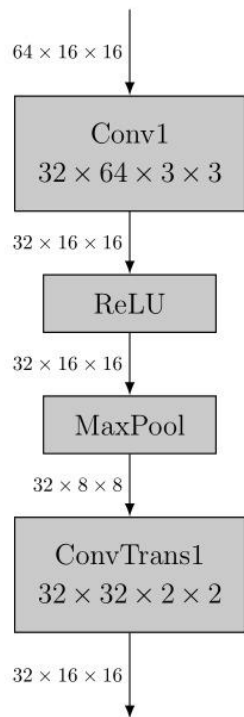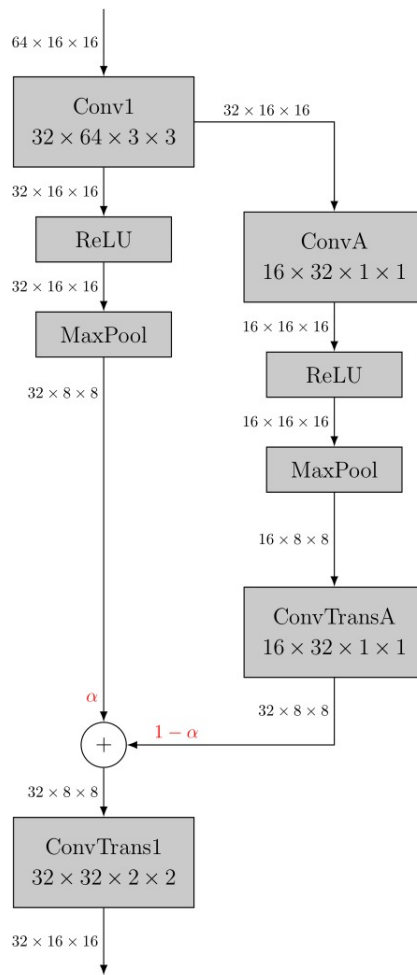Conv1
128 × 64 × 3 × 3

128 × 1 × 1

128 × 1 × 1

ReLU

128 × 1 × 1

Flatten

128 × 1 × 1

ConvA
64 × 128 × 1 × 1

64 × 1 × 1

ReLU

64 × 1 × 1

Flatten

64 × 1 × 1

FCa
64 × 128

128 × 1 × 1

α

1 − α

+

128 × 1 × 1

FC1
128 × 128

128 × 1 × 1

64 × 1 × 1

Conv1
128 × 64 × 3 × 3

128 × 1 × 1

ReLU

128 × 1 × 1

Flatten

128 × 1 × 1

FC1
128 × 128

128 × 1 × 1

64 × 1 × 1

Conv1
64 × 64 × 3 × 3

64 × 1 × 1

ReLU

64 × 1 × 1

Flatten

64 × 1 × 1

FC1
64 × 128

128 × 1 × 1