

Accelerating Video Super Resolution for Mobile Device

CSC 591/791-025 Final Project

Brian Park, Oliver Fowler

Outline

Video Super Resolution Task

Super Resolution Model Choice

Optimizations

Experimental Results

Conclusion

Video Super Resolution Task

- SOTA Super Resolution models are very large
- Tradeoff between quality and speed
- Strict latency requirements for real time video

| FPS | Latency (ms) |
|-----|--------------|
| 24 | 41.67 |
| 30 | 33.33 |
| 60 | 16.67 |

DNN Model Choice

- Find a good baseline model
 - A model that achieves good baseline performance on cloud and compatible with XGen
- Reported any XGen issues we came across
- NTIRE 2022 Challenge on Efficient Super Resolution (CVPR 2022)

| Team | Main Track | Sub-Track 1 | Sub-Track 2 | PSNR [Val.] | PSNR [Test] | Ave. Time [ms] | #Params [M] | FLOPs [G] | #Acts [M] | GPU Mem. [M] | #Conv |
|-------------------|------------|--------------------|---------------------|-------------|-------------|------------------------|-----------------------|------------------------|------------------------|-------------------------|-------|
| ByteESR | 1 | 22 ₍₁₁₎ | 33 ₍₂₎ | 29.00 | 28.72 | 27.11 ₍₁₎ | 0.317 ₍₁₁₎ | 19.70 ₍₁₁₎ | 80.05 ₍₆₎ | 377.91 ₍₄₎ | 39 |
| NJU_Jet | 2 | 37 ₍₁₈₎ | 44 ₍₆₎ | 29.00 | 28.69 | 28.07 ₍₂₎ | 0.341 ₍₁₈₎ | 22.28 ₍₁₉₎ | 72.09 ₍₄₎ | 204.60 ₍₁₎ | 34 |
| NEESR | 3 | 10 ₍₄₎ | 27 ₍₁₎ | 29.01 | 28.71 | 29.97 ₍₃₎ | 0.272 ₍₄₎ | 16.86 ₍₆₎ | 79.59 ₍₅₎ | 575.99 ₍₉₎ | 59 |
| Super | 4 | 26 ₍₁₂₎ | 55 ₍₁₀₎ | 29.00 | 28.71 | 32.09 ₍₄₎ | 0.326 ₍₁₄₎ | 20.06 ₍₁₂₎ | 93.82 ₍₁₀₎ | 663.07 ₍₁₅₎ | 59 |
| MegSR | 5 | 18 ₍₉₎ | 43 ₍₅₎ | 29.00 | 28.68 | 32.59 ₍₅₎ | 0.290 ₍₉₎ | 17.70 ₍₉₎ | 91.72 ₍₈₎ | 640.63 ₍₁₂₎ | 64 |
| rainbow | 6 | 16 ₍₈₎ | 34 ₍₃₎ | 29.01 | 28.74 | 34.10 ₍₆₎ | 0.276 ₍₆₎ | 17.98 ₍₁₀₎ | 92.80 ₍₉₎ | 309.23 ₍₃₎ | 59 |
| VMCL.Taobao | 7 | 29 ₍₁₄₎ | 57 ₍₁₁₎ | 29.01 | 28.68 | 34.24 ₍₇₎ | 0.323 ₍₁₃₎ | 20.97 ₍₁₆₎ | 98.67 ₍₁₁₎ | 633.00 ₍₁₀₎ | 40 |
| Bilibili AI | 8 | 15 ₍₇₎ | 41 ₍₄₎ | 29.00 | 28.70 | 34.67 ₍₈₎ | 0.283 ₍₈₎ | 17.61 ₍₇₎ | 90.50 ₍₇₎ | 633.74 ₍₁₁₎ | 64 |
| NKU_ESR | 9 | 12 ₍₅₎ | 48 ₍₇₎ | 29.00 | 28.66 | 34.81 ₍₉₎ | 0.276 ₍₇₎ | 16.73 ₍₅₎ | 111.12 ₍₁₃₎ | 662.51 ₍₁₄₎ | 65 |
| NJUST_RESTORARION | 10 | 54 ₍₂₇₎ | 89 ₍₁₅₎ | 28.99 | 28.68 | 35.76 ₍₁₀₎ | 0.421 ₍₂₈₎ | 27.67 ₍₂₆₎ | 108.66 ₍₁₂₎ | 643.95 ₍₁₃₎ | 52 |
| TOVBU | 11 | 43 ₍₂₁₎ | 96 ₍₁₉₎ | 29.00 | 28.71 | 38.32 ₍₁₁₎ | 0.376 ₍₂₃₎ | 22.38 ₍₂₀₎ | 113.55 ₍₁₅₎ | 867.17 ₍₂₇₎ | 64 |
| Alpan Team | 12 | 18 ₍₁₀₎ | 51 ₍₉₎ | 29.01 | 28.75 | 39.63 ₍₁₂₎ | 0.326 ₍₁₅₎ | 12.31 ₍₃₎ | 115.52 ₍₁₆₎ | 439.37 ₍₅₎ | 132 |
| Dragon | 13 | 38 ₍₁₉₎ | 70 ₍₁₃₎ | 29.01 | 28.69 | 41.80 ₍₁₃₎ | 0.358 ₍₂₀₎ | 21.11 ₍₁₈₎ | 120.15 ₍₁₇₎ | 260.00 ₍₂₎ | 131 |
| TieGuoDun Team | 14 | 54 ₍₂₇₎ | 104 ₍₂₁₎ | 28.95 | 28.65 | 42.35 ₍₁₄₎ | 0.433 ₍₂₉₎ | 27.10 ₍₂₅₎ | 112.03 ₍₁₄₎ | 788.13 ₍₂₂₎ | 64 |
| HiImageTeam | 15 | 7 ₍₃₎ | 70 ₍₁₃₎ | 29.00 | 28.72 | 47.75 ₍₁₅₎ | 0.242 ₍₃₎ | 14.51 ₍₄₎ | 151.36 ₍₂₃₎ | 861.84 ₍₂₅₎ | 100 |
| xilinxSR | 16 | 66 ₍₃₄₎ | 107 ₍₂₂₎ | 29.05 | 28.75 | 48.20 ₍₁₆₎ | 0.790 ₍₃₄₎ | 51.76 ₍₃₂₎ | 136.31 ₍₁₈₎ | 471.37 ₍₇₎ | 38 |
| cipher | 17 | 50 ₍₂₄₎ | 111 ₍₂₃₎ | 29.00 | 28.72 | 51.42 ₍₁₇₎ | 0.407 ₍₂₆₎ | 25.25 ₍₂₄₎ | 155.35 ₍₂₄₎ | 770.82 ₍₂₀₎ | 67 |
| NJU_MCG | 18 | 13 ₍₆₎ | 66 ₍₁₂₎ | 28.99 | 28.71 | 52.02 ₍₁₈₎ | 0.275 ₍₅₎ | 17.65 ₍₈₎ | 212.35 ₍₂₇₎ | 511.08 ₍₈₎ | 84 |
| IMGWLH | 19 | 34 ₍₁₇₎ | 91 ₍₁₇₎ | 29.01 | 28.72 | 56.34 ₍₁₉₎ | 0.362 ₍₂₁₎ | 20.10 ₍₁₃₎ | 136.35 ₍₁₉₎ | 753.02 ₍₁₉₎ | 113 |
| imglhl | 20 | 45 ₍₂₂₎ | 92 ₍₁₈₎ | 29.03 | 28.75 | 56.88 ₍₂₀₎ | 0.381 ₍₂₄₎ | 23.26 ₍₂₁₎ | 144.05 ₍₂₁₎ | 451.21 ₍₆₎ | 127 |
| whu_sigma | 21 | 63 ₍₃₂₎ | 132 ₍₃₀₎ | 29.02 | 28.73 | 61.04 ₍₂₁₎ | 0.705 ₍₃₃₎ | 43.88 ₍₃₀₎ | 142.91 ₍₂₀₎ | 1011.54 ₍₂₈₎ | 64 |
| Aselsan Research | 22 | 27 ₍₁₃₎ | 98 ₍₂₀₎ | 29.02 | 28.73 | 63.18 ₍₂₂₎ | 0.317 ₍₁₂₎ | 20.71 ₍₁₅₎ | 206.05 ₍₂₆₎ | 799.52 ₍₂₃₎ | 134 |
| Drintea | 23 | 59 ₍₃₁₎ | 121 ₍₂₇₎ | 29.00 | 28.70 | 75.52 ₍₂₃₎ | 0.589 ₍₃₁₎ | 36.92 ₍₂₈₎ | 148.05 ₍₂₂₎ | 734.54 ₍₁₇₎ | 67 |
| GDUT_SR | 24 | 50 ₍₂₄₎ | 136 ₍₃₁₎ | 29.05 | 28.75 | 75.70 ₍₂₄₎ | 0.414 ₍₂₇₎ | 24.80 ₍₂₃₎ | 260.05 ₍₂₈₎ | 1457.98 ₍₃₄₎ | 195 |
| Giantpandacy | 25 | 63 ₍₃₂₎ | 150 ₍₃₄₎ | 29.07 | 28.76 | 87.87 ₍₂₅₎ | 0.683 ₍₃₂₎ | 45.07 ₍₃₁₎ | 361.23 ₍₃₁₎ | 1272.95 ₍₃₁₎ | 122 |
| neptune | 26 | 39 ₍₂₀₎ | 123 ₍₂₉₎ | 28.99 | 28.69 | 101.69 ₍₂₆₎ | 0.316 ₍₁₀₎ | 38.03 ₍₂₉₎ | 269.48 ₍₂₉₎ | 1179.05 ₍₄₅₎ | 45 |
| XPixel | 27 | 3 ₍₁₎ | 49 ₍₈₎ | 29.01 | 28.69 | 140.47 ₍₂₇₎ | 0.156 ₍₁₎ | 9.50 ₍₂₎ | 65.76 ₍₃₎ | 729.94 ₍₁₆₎ | 43 |
| NJUST_ESR | 28 | 3 ₍₁₎ | 89 ₍₁₅₎ | 28.96 | 28.68 | 164.80 ₍₂₈₎ | 0.176 ₍₂₎ | 8.73 ₍₁₎ | 160.43 ₍₂₅₎ | 1346.74 ₍₃₃₎ | 25 |
| TeamInception | 29 | 57 ₍₃₀₎ | 146 ₍₃₃₎ | 29.12 | 28.82 | 171.56 ₍₂₉₎ | 0.505 ₍₃₀₎ | 32.42 ₍₂₇₎ | 502.27 ₍₃₄₎ | 866.16 ₍₂₆₎ | 74 |
| cceNBgdd | 30 | 33 ₍₁₆₎ | 114 ₍₂₄₎ | 28.97 | 28.67 | 180.60 ₍₃₀₎ | 0.339 ₍₁₆₎ | 21.11 ₍₁₇₎ | 404.16 ₍₃₃₎ | 739.65 ₍₁₈₎ | 197 |
| ZLZ | 31 | 55 ₍₂₉₎ | 118 ₍₂₆₎ | 29.00 | 28.72 | 183.43 ₍₃₁₎ | 0.372 ₍₂₂₎ | 64.45 ₍₃₃₎ | 57.51 ₍₂₎ | 1244.23 ₍₃₀₎ | 16 |
| Express | 32 | 31 ₍₁₅₎ | 117 ₍₂₅₎ | 29.04 | 28.77 | 203.16 ₍₃₂₎ | 0.339 ₍₁₇₎ | 20.41 ₍₁₄₎ | 325.53 ₍₃₀₎ | 853.27 ₍₂₄₎ | 148 |
| Just Try | 33 | 70 ₍₃₅₎ | 170 ₍₃₅₎ | 29.12 | 28.81 | 247.90 ₍₃₃₎ | 0.832 ₍₃₅₎ | 135.30 ₍₃₅₎ | 392.43 ₍₃₂₎ | 2387.93 ₍₃₅₎ | 207 |
| ncepuxplorers | 34 | 47 ₍₂₃₎ | 137 ₍₃₂₎ | 29.09 | 28.79 | 317.66 ₍₃₄₎ | 0.390 ₍₂₅₎ | 23.73 ₍₂₂₎ | 994.25 ₍₃₅₎ | 771.54 ₍₂₁₎ | 374 |
| mju_mnu | 35 | 53 ₍₂₆₎ | 121 ₍₂₇₎ | 29.06 | 28.79 | 332.28 ₍₃₅₎ | 0.345 ₍₁₉₎ | 78.81 ₍₃₄₎ | 46.76 ₍₁₎ | 1310.72 ₍₃₂₎ | 40 |

The following methods are not ranked since their validation/testing PSNR are not on par with the baseline.

| | | | | | | | | |
|---------------------|-------|-------|---------|-------|--------|---------|----------|-----|
| VirtualReality Team | 27.35 | 27.26 | 2231.32 | 0.423 | 423.16 | 2731.08 | 3336.88* | 82 |
| NTU607QCO-ESR | 27.79 | 27.61 | 38.85 | 0.433 | 27.06 | 108.89 | 776.38 | 60 |
| Strong Tiger | 29.00 | 28.61 | 34.92 | 0.560 | 36.64 | 78.91 | 641.13 | 23 |
| VAP | 29.01 | 28.47 | 23.96 | 0.175 | 10.83 | 70.93 | 507.64 | 63 |
| Multicog | 29.20 | 28.66 | 20.97 | 0.243 | 2.25 | 142.22 | 1422.22 | 100 |
| Set5Baby Team | 29.20 | 28.66 | 20.97 | 0.243 | 2.25 | 142.22 | 1422.22 | 100 |
| NWPU_SweetDreamLab | 29.20 | 28.66 | 20.97 | 0.243 | 2.25 | 142.22 | 1422.22 | 100 |
| SSL | 29.20 | 28.66 | 20.97 | 0.243 | 2.25 | 142.22 | 1422.22 | 100 |
| RFDN AIM2020 Winner | | | | | | | | |
| IMDN_baseline | | | | | | | | |

* This solution uses too much GPU memory. Image

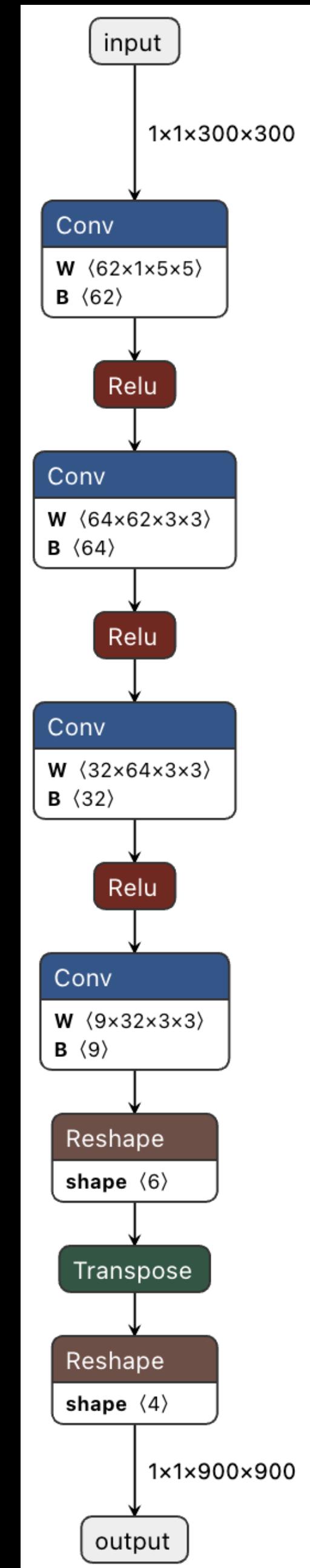


Compatibility Matrix of SOTA SR Models

| SOTA Model | XGen Baseline Compatibility | NNI Pruning Compatibility | Baseline Inference on V100 (ms) | Baseline Inference on Galaxy S10e (ms) |
|------------------|-----------------------------|---------------------------|---------------------------------|--|
| Twitter | ✓ | ✓ | 2.39 | 72.46 |
| ByteDance | ✓ | ✗ | 20.97 | 166.00 |
| WDSR (XGen Demo) | ✓ | ✗ | 54.96 | 424.70 |
| RDN | ✗ | ✗ | 67.84 | ✗ |
| IMDN | ✗ | ✗ | 27.17 | ✗ |
| Real-ERSGAN | ✗ | ✗ | 200.12 | ✗ |

SuperResolutionTwitter

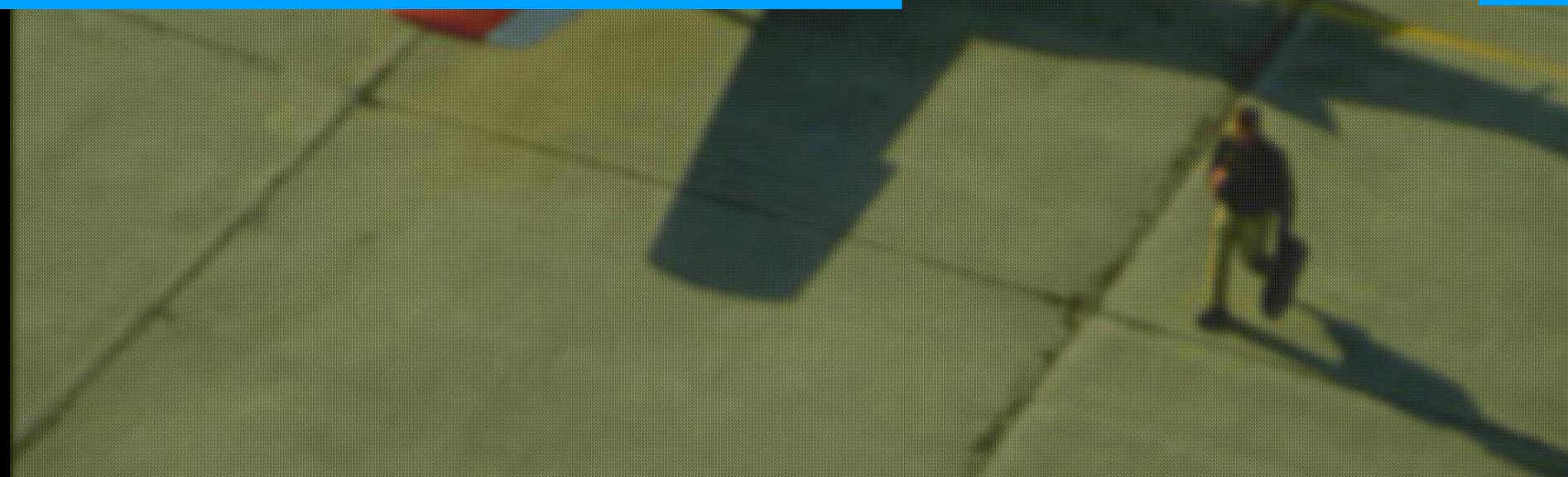
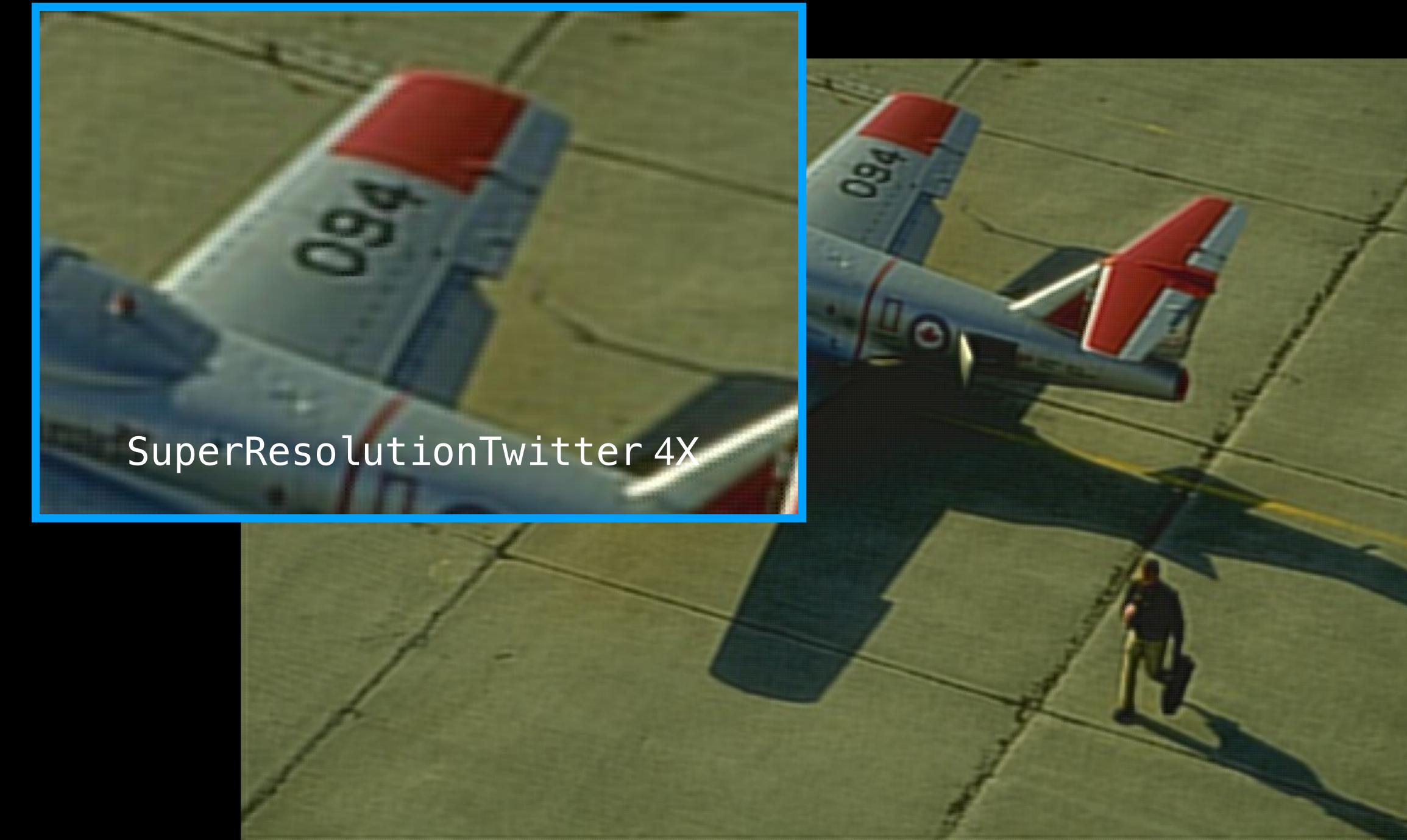
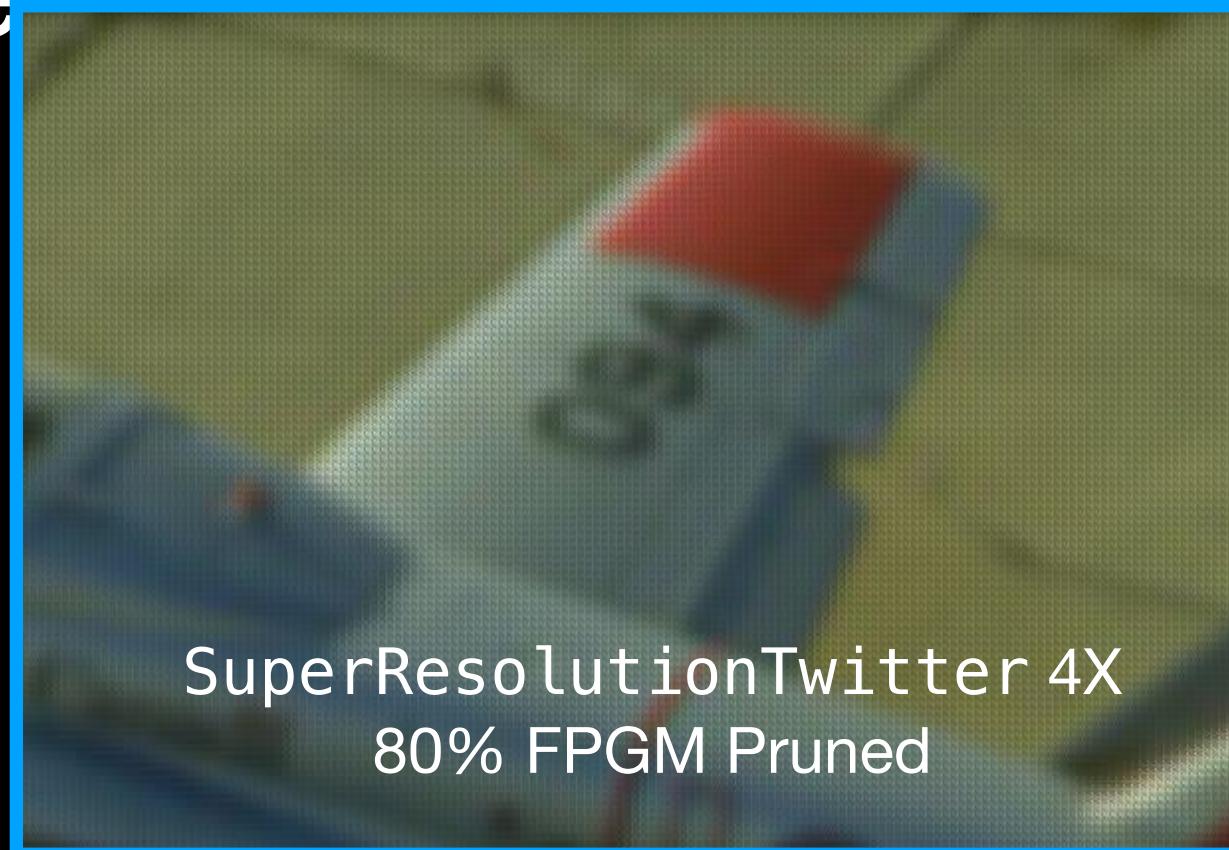
- ***Real-Time Single Image and Video Super-Resolution Using an Efficient Sub-Pixel Convolutional Neural Network***, Twitter
- CVPR 2016
- <https://arxiv.org/abs/1609.05158>
- Trained on BSD300 dataset



Optimizations

Pruning

- Pruning degrades quality of model
- Pruning aggressively over 80% will produce ~~even worse results~~



Optimizations

Quantization

- Tried modifying the ONNX file to FP16 weights
- FP16 supported by both mobile and cloud hardware
- XGen does this for us, so no manual optimizations required

Qualcomm
snapdragon



855 mobile platform

Adreno 640

50% More ALUs*
FP32 & FP16

Hexagon 690

New Tensor Accelerator

- QTI designed
- Dedicated to AI
- Multidimensional math and integrated nonlinear functions

4x Vector eXtensions*

Optimized scalar

Voice Assistant

INT16, INT8 & Mixed

Kryo 485

New dot product instructions

FP32 & INT8

*Compared to Snapdragon 845
Qualcomm Adreno, Qualcomm Spectra, Qualcomm Hexagon, Qualcomm Processor Security and Qualcomm Kryo are products of Qualcomm Technologies, Inc. and/or its subsidiaries.

Optimizations

Color Optimization

- Used YCbCr instead of RGB channels
- Training and inference on only the luma (Y) channel
- Reduces tensor size by 3X from $(N, 3, H, W)$ to $(N, 1, H, W)$
- Overhead of deconstruction and reconstruction at every frame

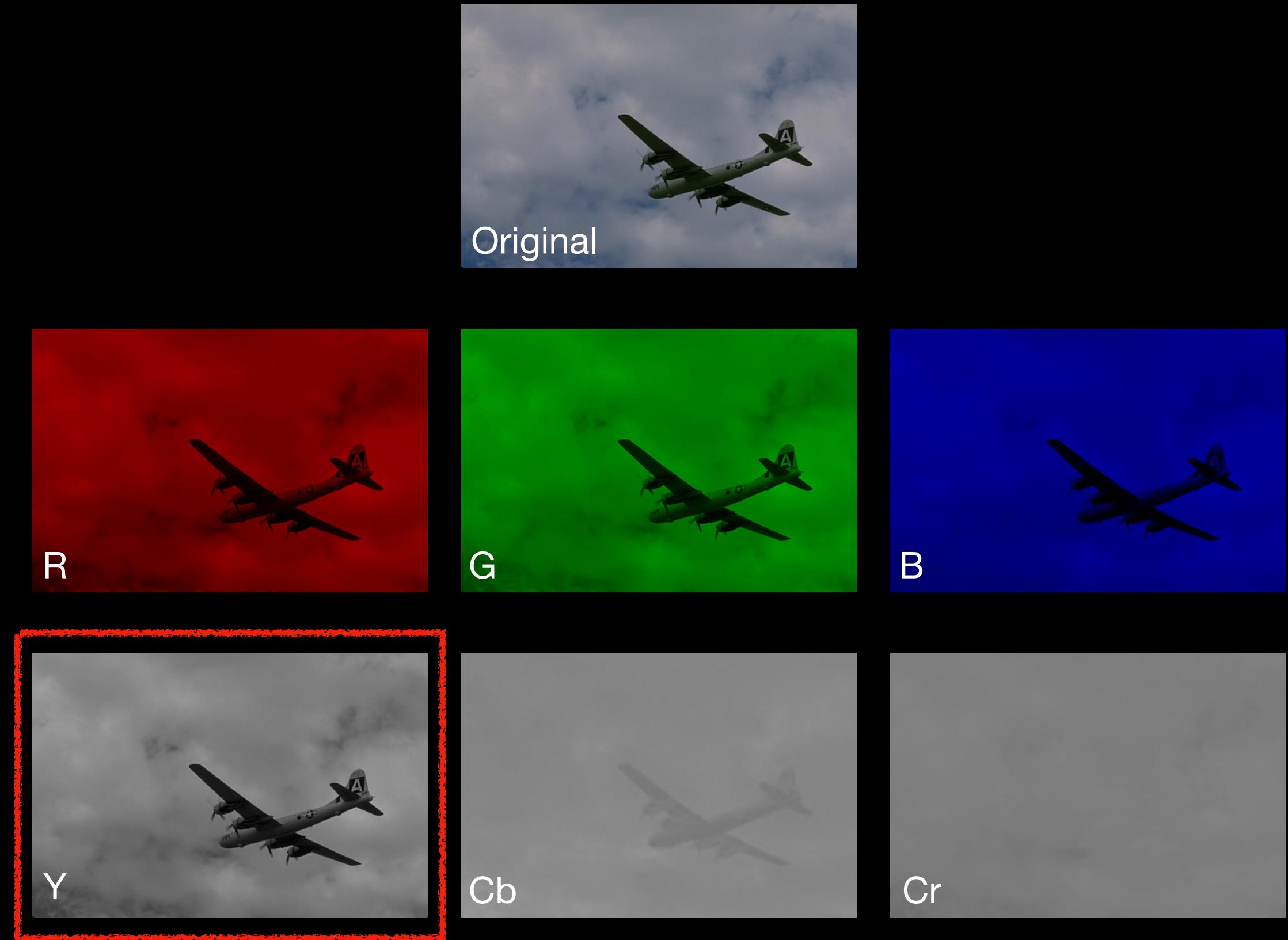


Inputs to Super Resolution Model

Optimizations

Color Optimization

- Used YCbCr instead of RGB channels
- Training and inference on only the luma (Y) channel
- Reduces tensor size by 3X from $(N, 3, H, W)$ to $(N, 1, H, W)$
- Overhead of deconstruction and reconstruction at every frame

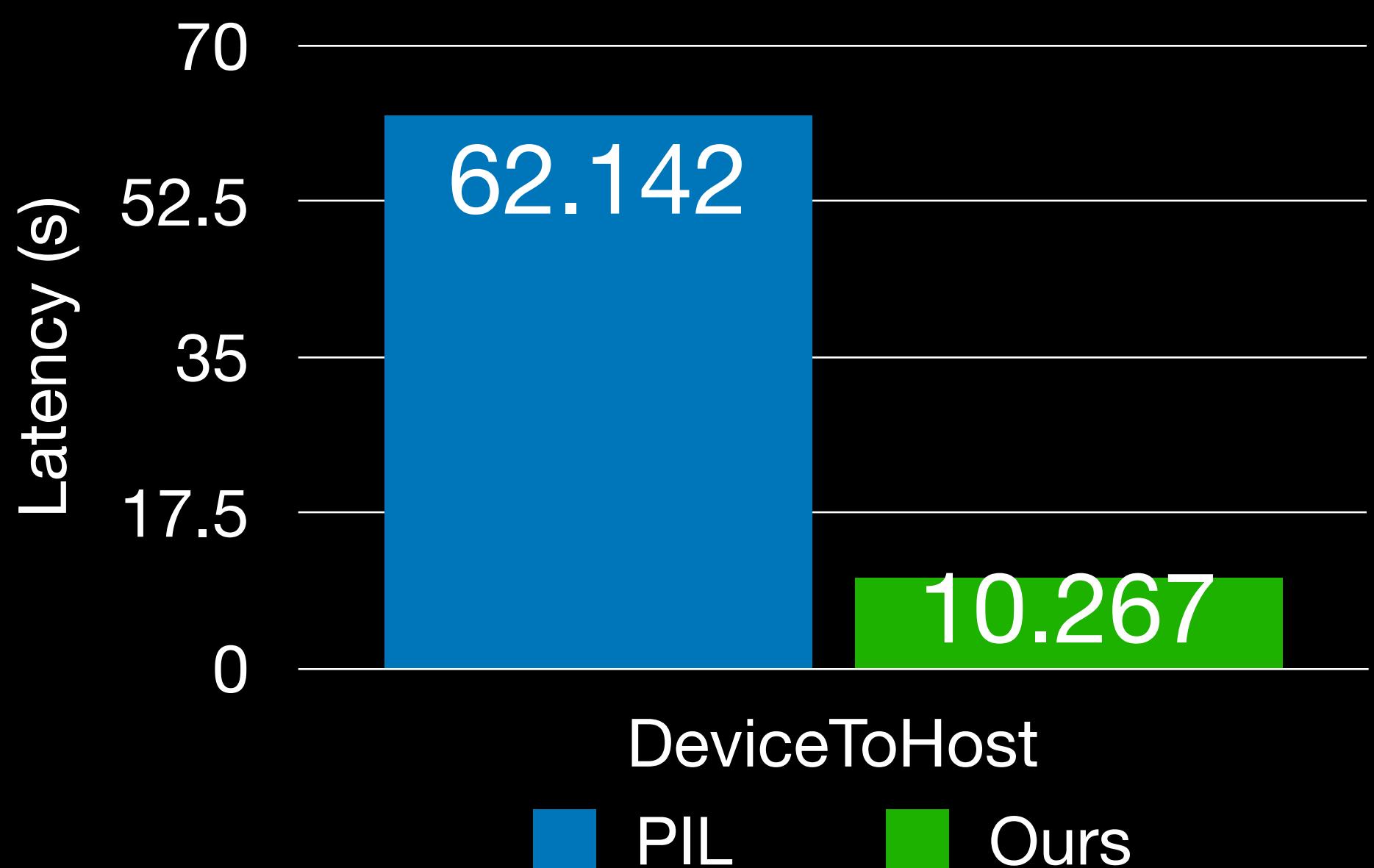
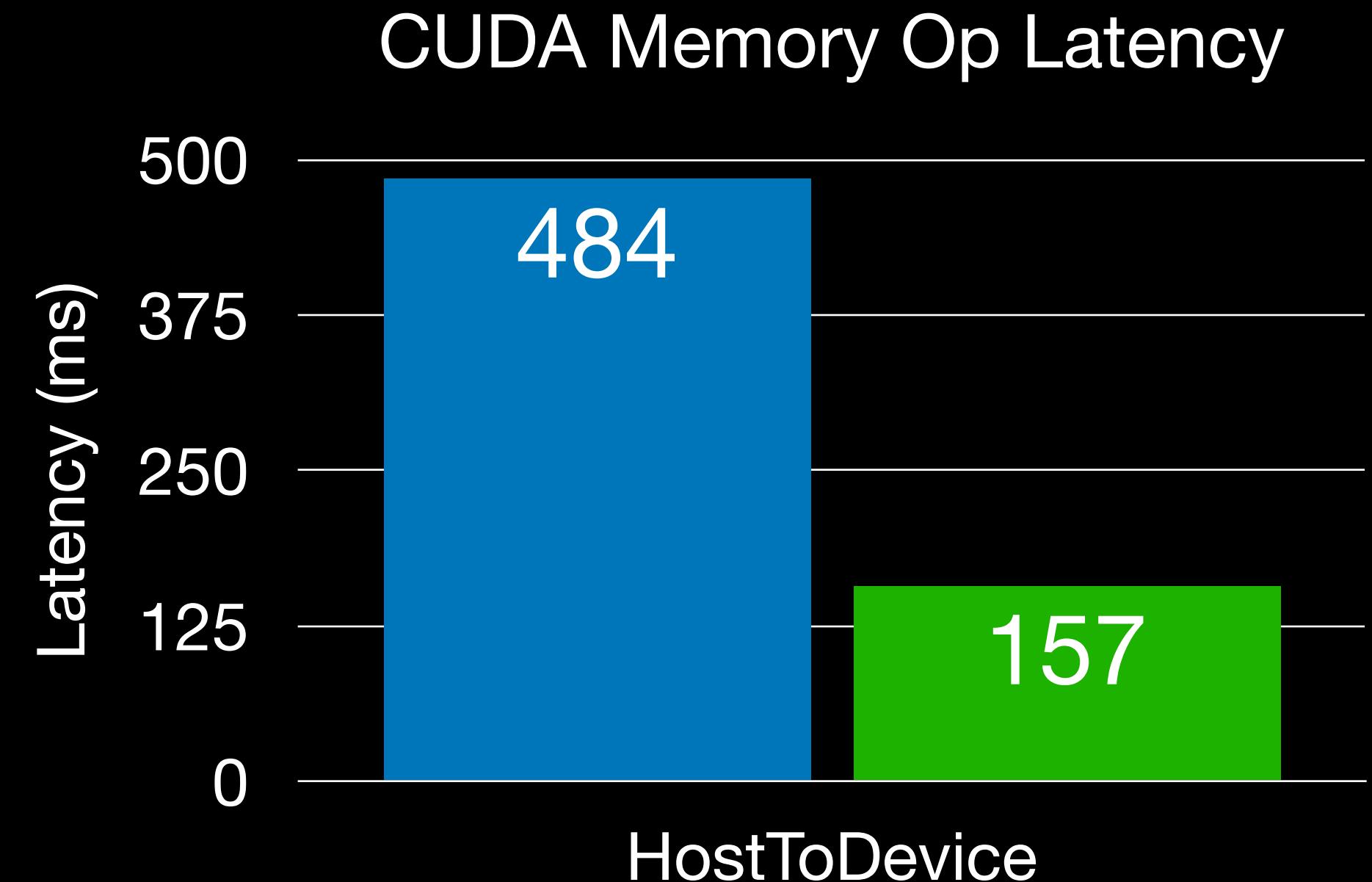


Inputs to Super Resolution Model

Optimizations

Color Optimization

- Originally used Python Imaging Library (PIL) to do the color conversion
 - PIL is written only for the CPU
 - Overhead of GPU/CPU memory transfers (`cudaMemcpyDeviceToHost/`
`cudaMemcpyHostToDevice`)
 - Optimized code by converting CPU operations into CUDA operations via PyTorch



Experimental Results

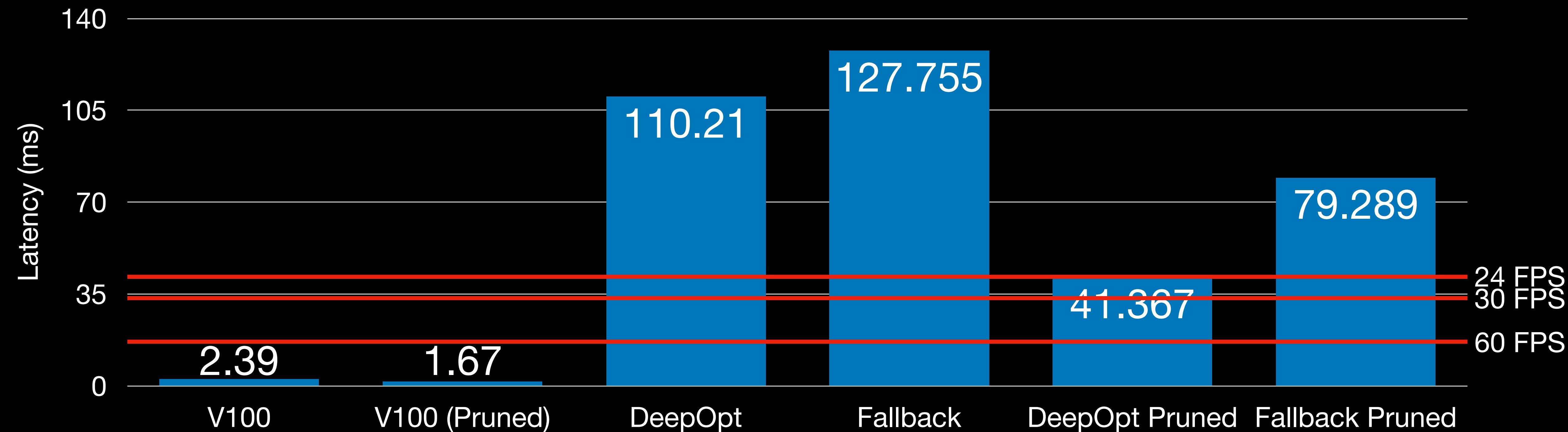
Benchmark Configuration

- Benchmarked video frames upscaled from 360p to 2K resolution
 - Upscale factor set to 4X
 - From input of 640×360 to output 2560×1440

Experimental Results

Latency

SuperResolutionTwitter 360p to 2K Inference



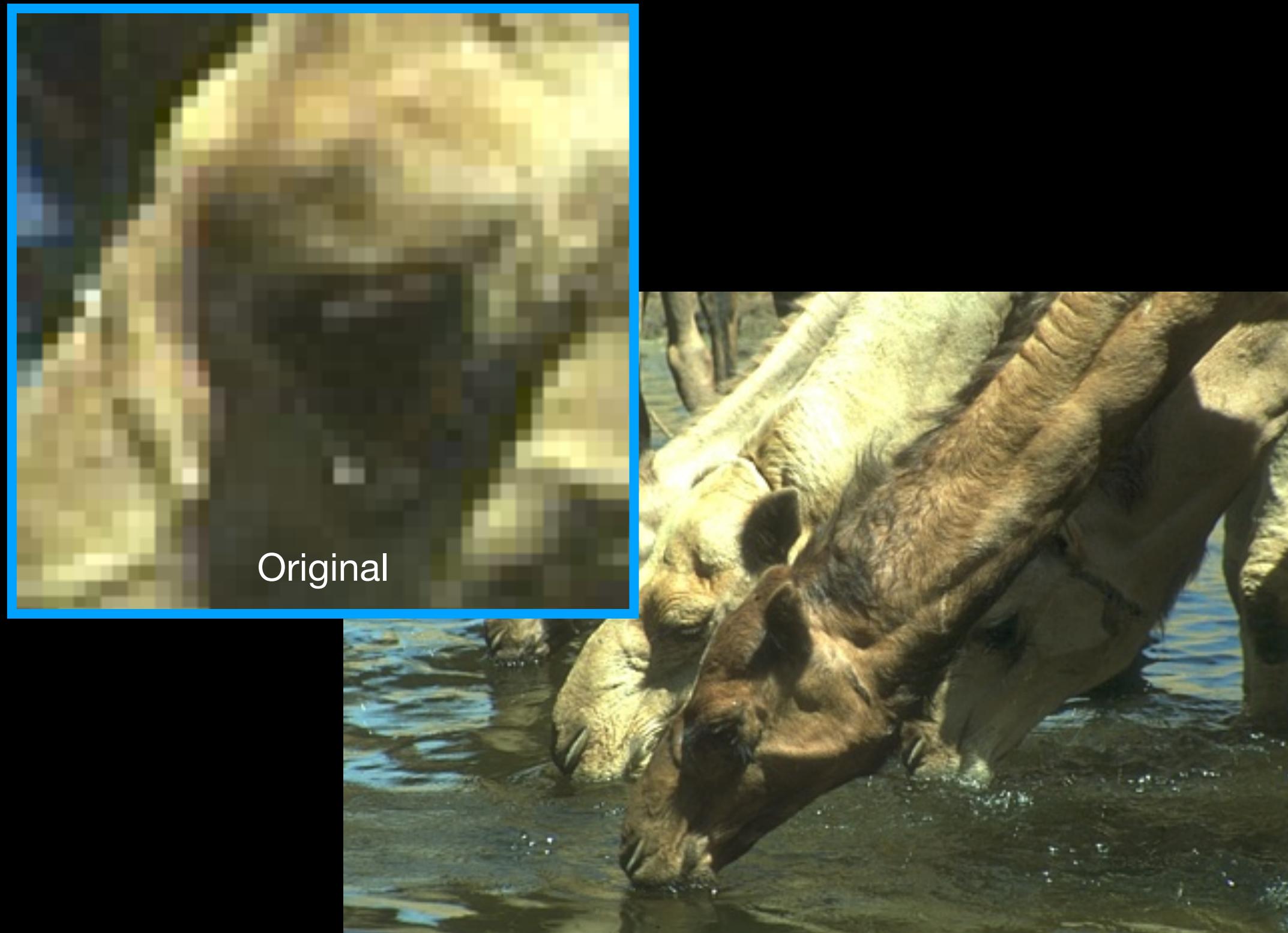
Original: 244KB

Pruned: 100KB

*L1NormPruner 60% Pruned

Experimental Results

Quality



PSNR: 6.47



PSNR: 20.69

Experimental Results

Other Frameworks

- Successfully tried compilation on other frameworks
 - TVM (Apache)
 - TensorRT (NVIDIA)
 - ONNX Runtime (Microsoft)
 - CoreML (Apple)

Challenges and Lessons Learned

- Hard to achieve a workable model on XGen due to limited functionality of opset
- Limited XGen resources with GPU and smartphone
- Easy to achieve and collect results on HPC GPU, but not for mobile
- Learned more about computer vision and super resolution

Next Steps

- Try more Super Resolution architectures
 - RNN could be used for video super resolution, as video frames hold redundant data between frames
- Aggressively apply other optimizations without impacting quality
- Try different datasets

Thank You

- ARC (NCSU Systems Lab)
- Bridges-2 (PSC)
- Lab XGen Server and Samsung Galaxy S10e
(Dr. Shen)

**NC STATE
UNIVERSITY**



