

2020  
embedded  
**VISION**  
summit®

# Designing Bespoke CNNs for Target Hardware

Woonhyun Nam  
Director, Algorithms  
StradVision  
September, 2020



**STRADVISION**

- We develop deep learning-based perception software for Advanced Driver Assistance Systems (ADAS) and Autonomous Vehicles

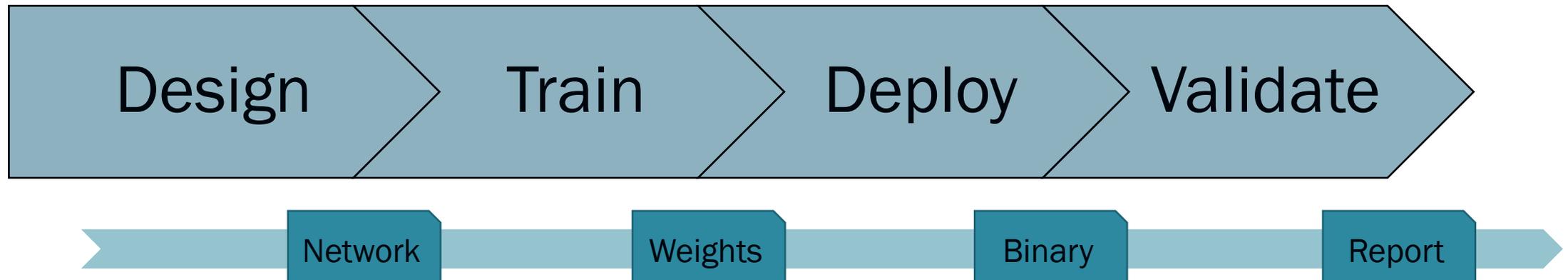


For more detailed information, please refer to <https://youtu.be/tK7F8yvpiGA> and <https://stradvision.com/>

- Development Process of Deep Learning-based Software
- Scalability Challenge
- Layer Transformations
- Case Studies
- Conclusions

# Development Process of Deep Learning-based Software

Repeated for each task (e.g., object detection) and target platform



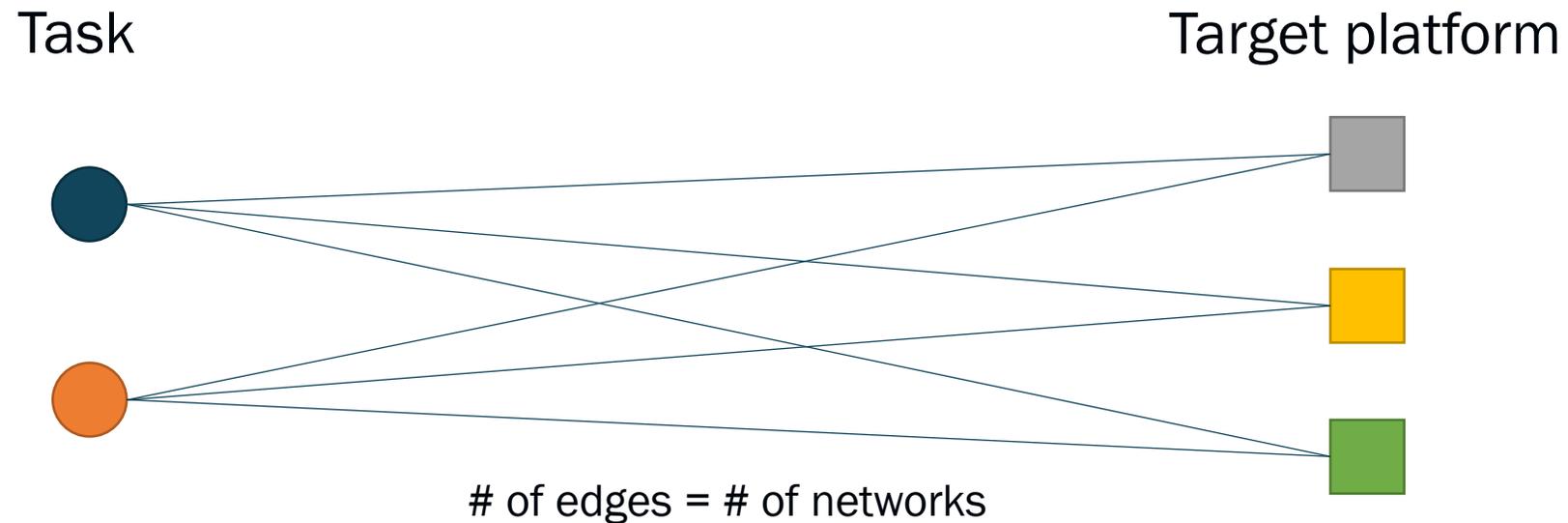
- (Pros)
- **Reuse** existing CNN designs
  - **Tweak** some layers

- (Cons)
- **Retrain from scratch** even for small changes
  - **Revalidate from scratch** if output changes

# Scalability Challenge

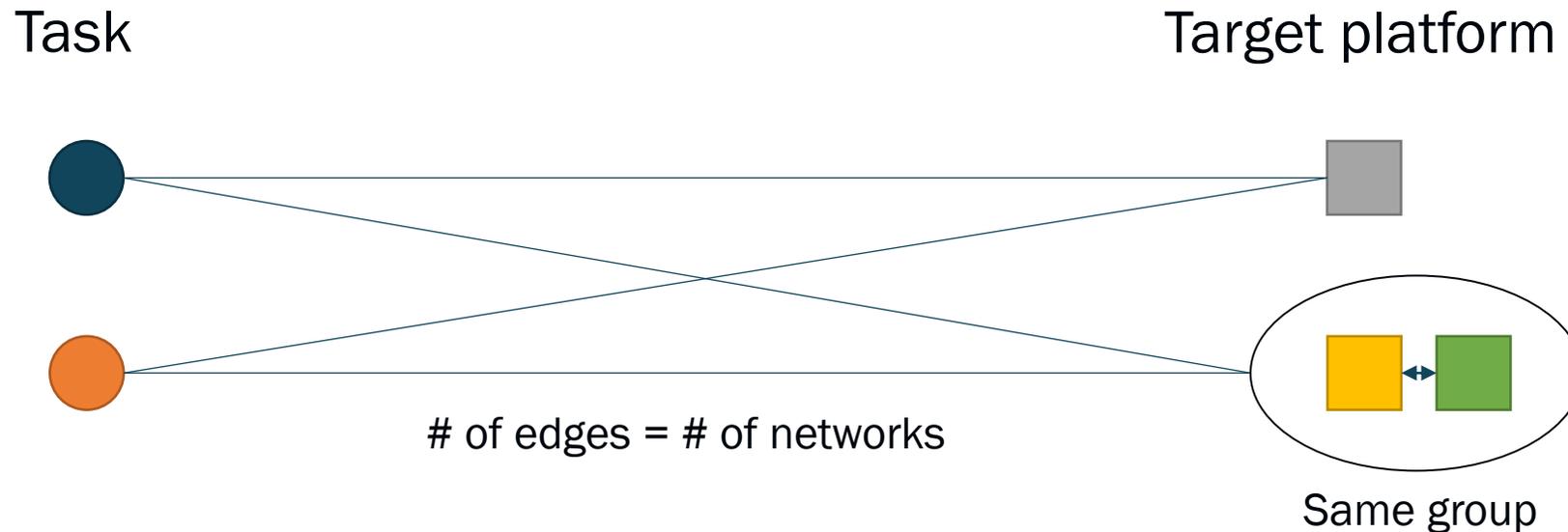
## Scalability challenge

- Number of networks to maintain grows rapidly
- More than 100 networks become difficult to manage



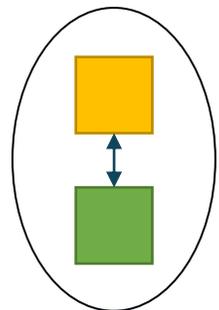
## Main idea

- Grouping target platforms based on similar computational characteristics (or capabilities) with respect to a given CNN transformation

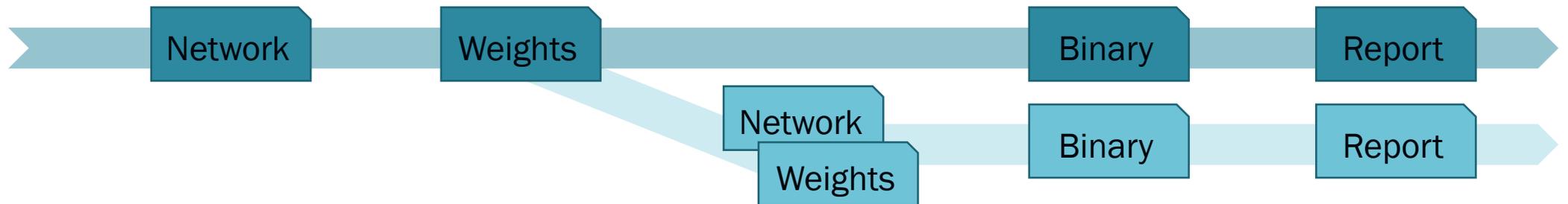


# Scalable Development Process

- Design/Train CNNs for each group
- Transform CNNs without retraining within same group



Same group



## Cost-effective method

- Increasing computational efficiency
  - E.g., lower latency, higher throughput, cheaper hardware
- Increasing engineering efficiency
  - E.g., shorter development time, less engineers, more projects
- Increasing economic efficiency
  - E.g., lower usage of cloud services for training, lower production-cost

# Layer Transformations

- Basic operations
- Case studies
  - YUV Input
  - Fixed-size Kernel Acceleration in Convolution
  - No FC Support
  - No Custom Layer Support in Quantization
  - Block-level Accumulation in Convolution

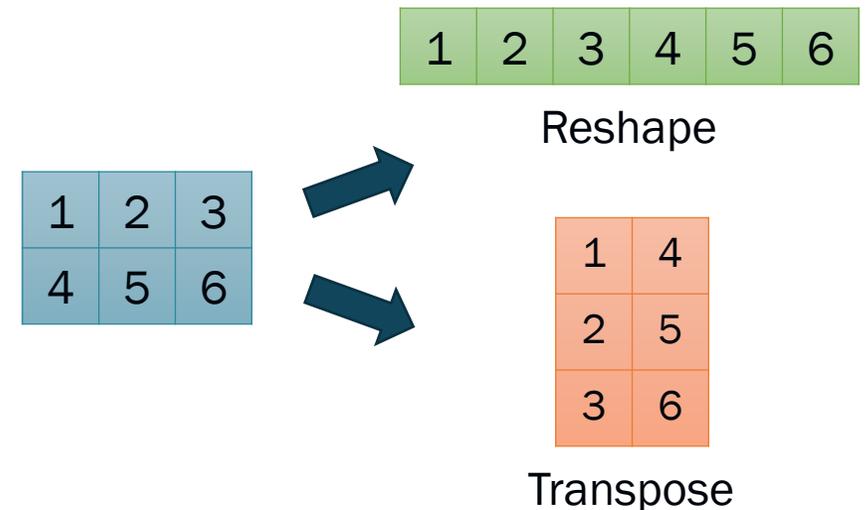
## Basic operations

- (Linear) layer fusion
  - Conv & Conv  $\rightarrow$  Conv
- Memory layout change
  - For Input, Weight, and Output tensors
  - Reshape, Transpose
- Zero-filling

$$y = ax + b \quad \text{and} \quad z = cy + d$$

$$\rightarrow z = (ac)x + (bc + d)$$

New weight      New bias



# Case Studies

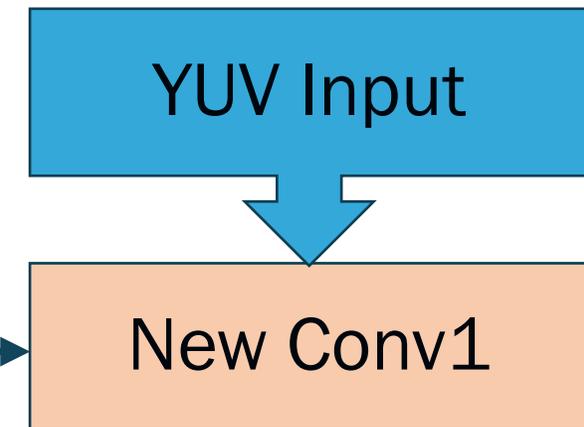
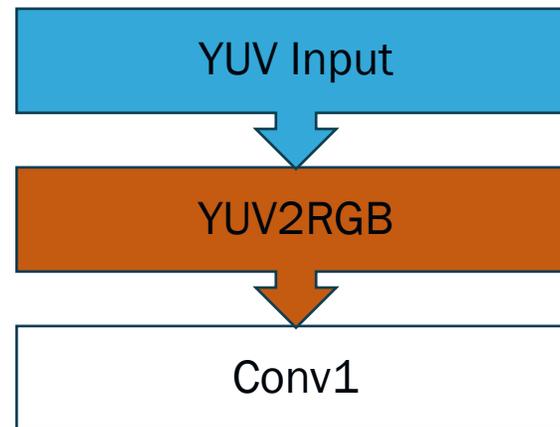
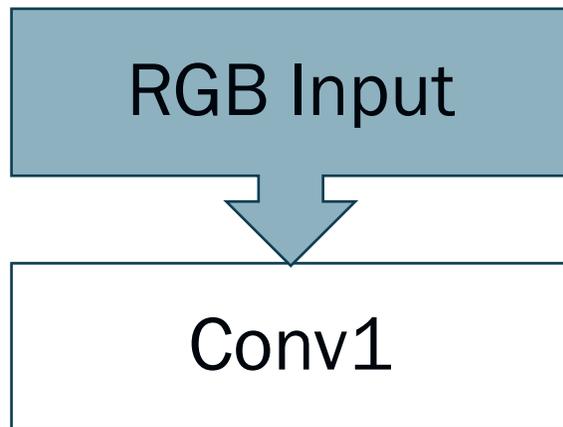
## Hypothetical target hardware

- YUV input, instead of RGB
- YUV2RGB equation is given

$$\begin{bmatrix} R \\ G \\ B \end{bmatrix} = \begin{bmatrix} W_{RY} & W_{RU} & W_{RV} \\ W_{GY} & W_{GU} & W_{GV} \\ W_{BY} & W_{BU} & W_{BV} \end{bmatrix} \begin{bmatrix} Y + b_Y \\ U + b_U \\ V + b_V \end{bmatrix}$$

YUV2RGB equation

→ Fusion of YUV2RGB and the first Conv



Fusion

## Equations for new weights and new bias

$$Z = [w_R \quad w_G \quad w_B] \begin{bmatrix} R \\ G \\ B \end{bmatrix} + b$$

Conv1 equation

$$\begin{bmatrix} R \\ G \\ B \end{bmatrix} = \begin{bmatrix} w_{RY} & w_{RU} & w_{RV} \\ w_{GY} & w_{GU} & w_{GV} \\ w_{BY} & w_{BU} & w_{BV} \end{bmatrix} \begin{bmatrix} Y + b_Y \\ U + b_U \\ V + b_V \end{bmatrix}$$

YUV2RGB equation



Fusion

$$Z = [w_R \quad w_G \quad w_B] \begin{bmatrix} w_{RY} & w_{RU} & w_{RV} \\ w_{GY} & w_{GU} & w_{GV} \\ w_{BY} & w_{BU} & w_{BV} \end{bmatrix} \begin{bmatrix} Y \\ U \\ V \end{bmatrix} + b + [w_R \quad w_G \quad w_B] \begin{bmatrix} w_{RY} & w_{RU} & w_{RV} \\ w_{GY} & w_{GU} & w_{GV} \\ w_{BY} & w_{BU} & w_{BV} \end{bmatrix} \begin{bmatrix} b_Y \\ b_U \\ b_V \end{bmatrix}$$

New weight

New Conv1 equation

New bias

## Hypothetical target hardware

- Only 5x5 kernel acceleration supported
- Applying zero-filling for smaller kernels
- Computational utilization
  - 1x1 Conv ~ 4%, 3x3 Conv ~ 36%

$W_{11}$	$W_{12}$	$W_{13}$
$W_{21}$	$W_{22}$	$W_{23}$
$W_{31}$	$W_{32}$	$W_{33}$

➔  
Zero filling

0	0	0	0	0
0	$W_{11}$	$W_{12}$	$W_{13}$	0
0	$W_{21}$	$W_{22}$	$W_{23}$	0
0	$W_{31}$	$W_{32}$	$W_{33}$	0
0	0	0	0	0

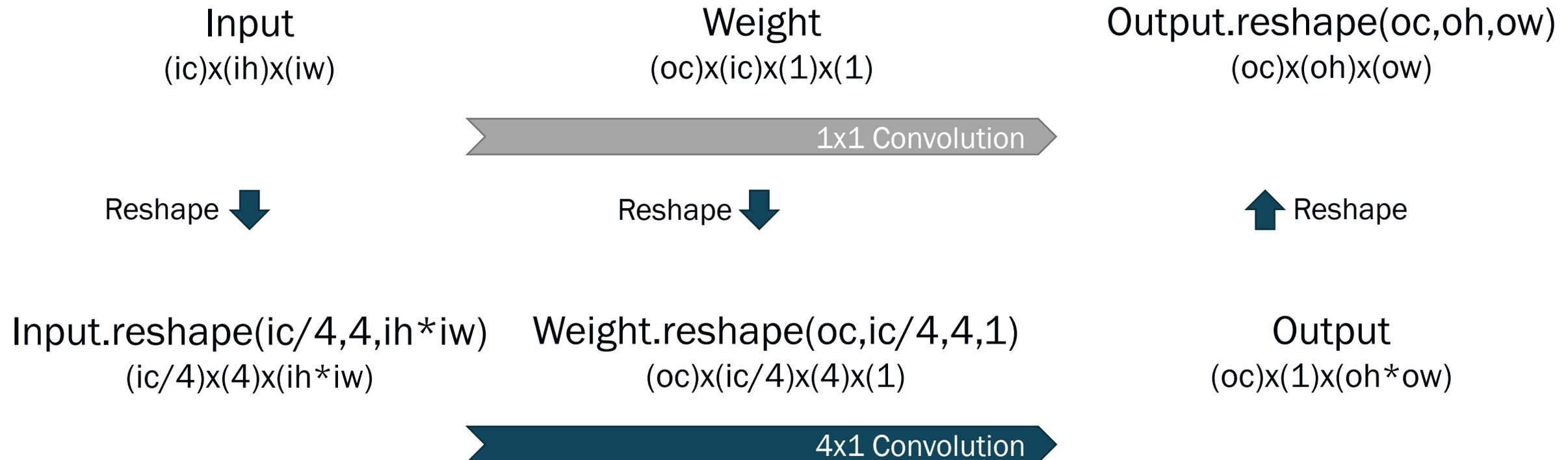
# Case Study – Fixed-size Kernel Acceleration in Convolution

→ Convert 1x1 Conv into 4(kh)x1(kw) Conv

iw=input width  
ih=input height  
ic=input channel

kw=kernel width  
kh=kernel height

ow=output width  
oh=output height  
oc=output channel



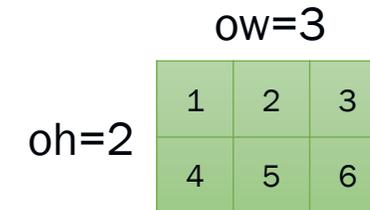
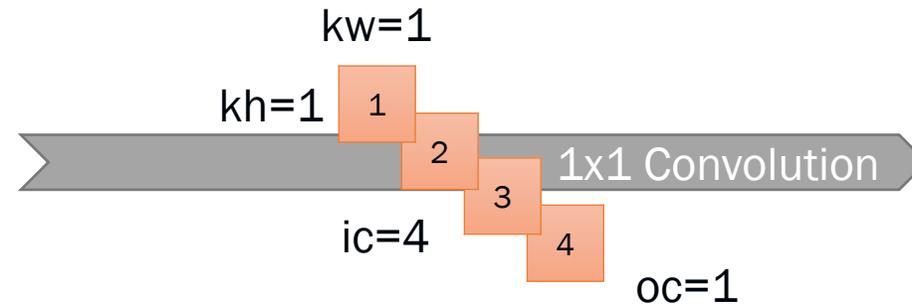
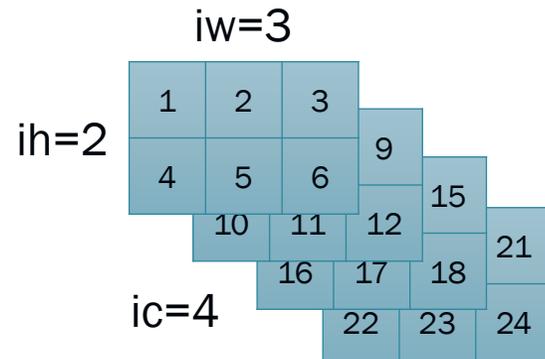
# Case Study – Fixed-size Kernel Acceleration in Convolution

→ Convert 1x1 Conv into 4(kh)x1(kw) Conv

iw=input width  
ih=input height  
ic=input channel

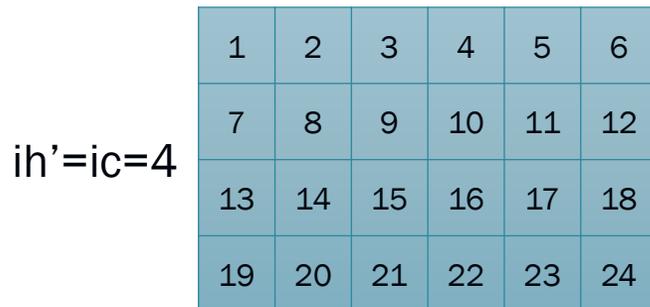
kw=kernel width  
kh=kernel height

ow=output width  
oh=output height  
oc=output channel



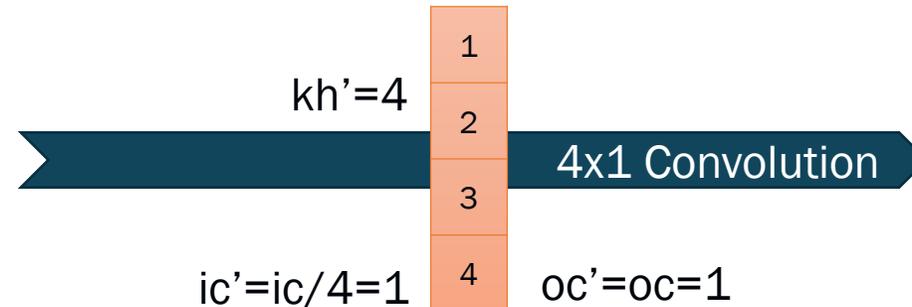
Reshape ↓

iw'=ih\*iw=6



Reshape ↓

kw'=1



↑ Reshape

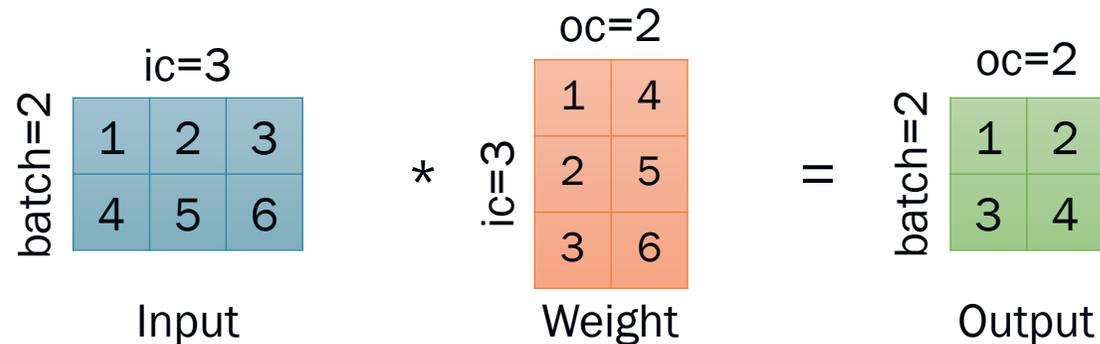
ow'=iw'=6



# Case Study – No FC Support

## Hypothetical target hardware

- No support for FC layers, or
- Only single-batched FC layers supported

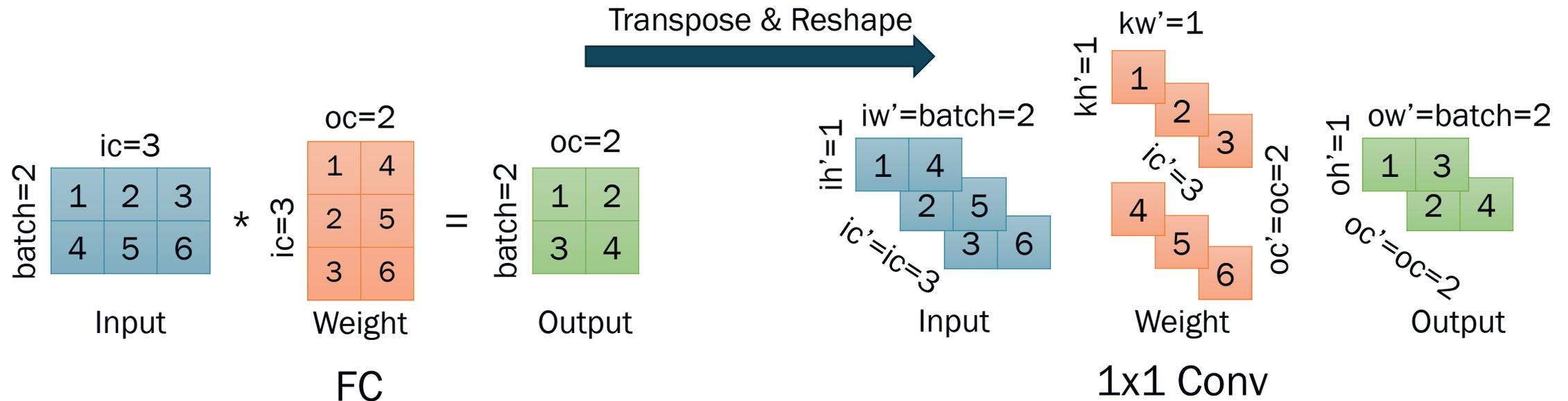


Matrix Multiplication in FC layer

# Case Study – No FC Support

## → Convert FC into 1x1 Conv

$$\begin{array}{l}
 \text{Input} \quad * \quad \text{Weight} \quad = \quad \text{Output} \\
 (\text{batch}) \times (\text{ic}) \quad (\text{ic}) \times (\text{oc}) \quad (\text{batch}) \times (\text{oc})
 \end{array}
 \quad
 \begin{array}{l}
 \text{transpose(Input)} \\
 \text{.reshape(ic,1,batch)} \\
 (\text{ic}) \times (\text{1}) \times (\text{batch})
 \end{array}
 *
 \begin{array}{l}
 \text{transpose(Weight)} \\
 \text{.reshape(oc,ic,1,1)} \\
 (\text{oc}) \times (\text{ic}) \times (\text{1}) \times (\text{1})
 \end{array}
 =
 \begin{array}{l}
 \text{transpose(Output)} \\
 \text{.reshape(oc,1,batch)} \\
 (\text{oc}) \times (\text{1}) \times (\text{batch})
 \end{array}$$



## Hypothetical target hardware

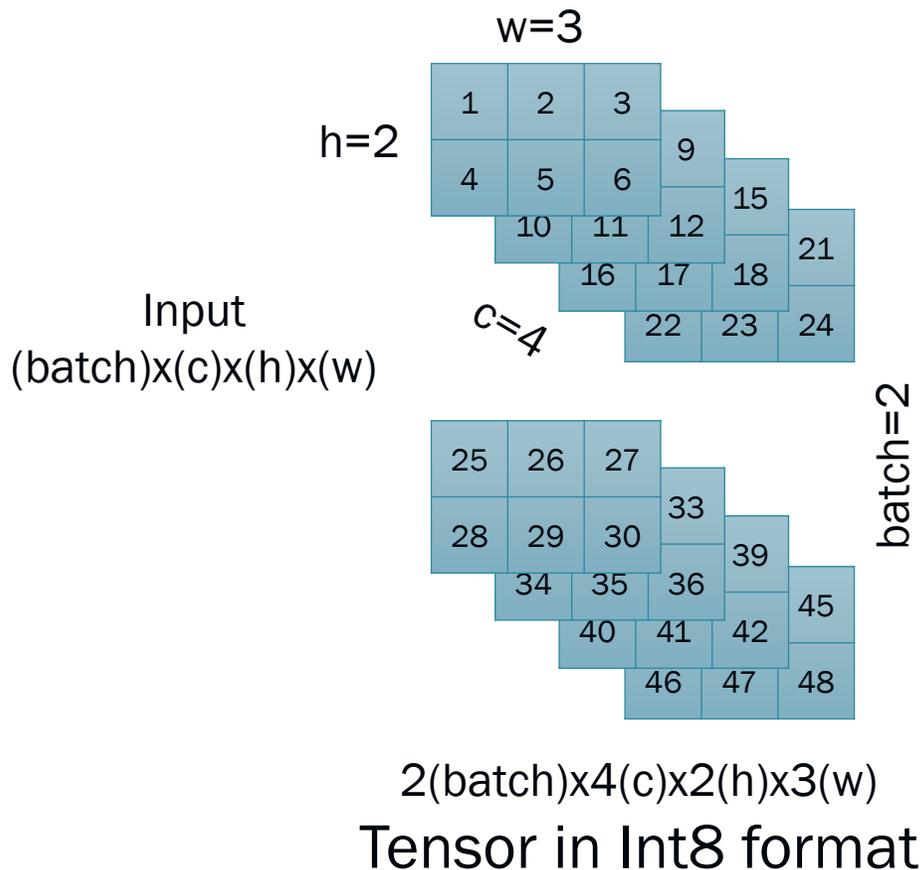
- Quantization should be done by a tool provided by hardware vendor
- No support for custom layers
- The tool accepts only image files as quantization inputs

→ Convert input tensors of custom layers into image files

- Useful for two-stage detection networks (E.g., Faster RCNN)

# Case Study – No Custom Layer Support in Quantization

## Convert 4D tensors into image files



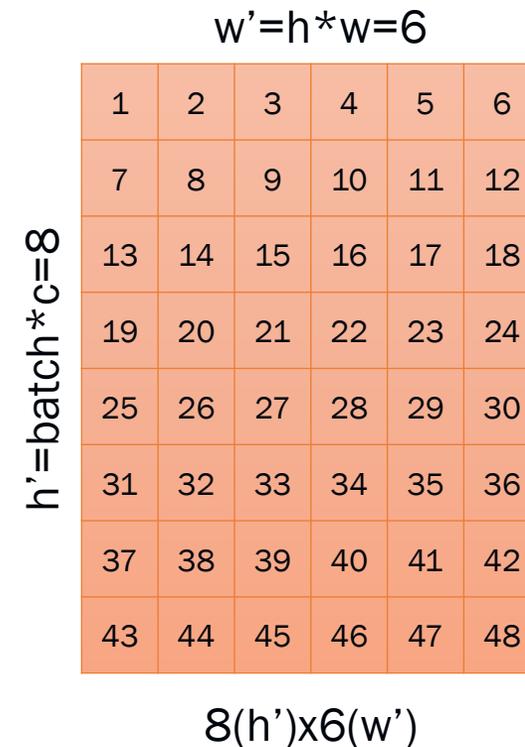
Reshape & Add 128



Subtract 128 & Reshape



Input.reshape(batch\*c,h\*w) + 128  
(batch\*c)x(h\*w)



Gray image in Uint8 format

# Case Study – Block-level Accumulation in Convolution

## Hypothetical target hardware

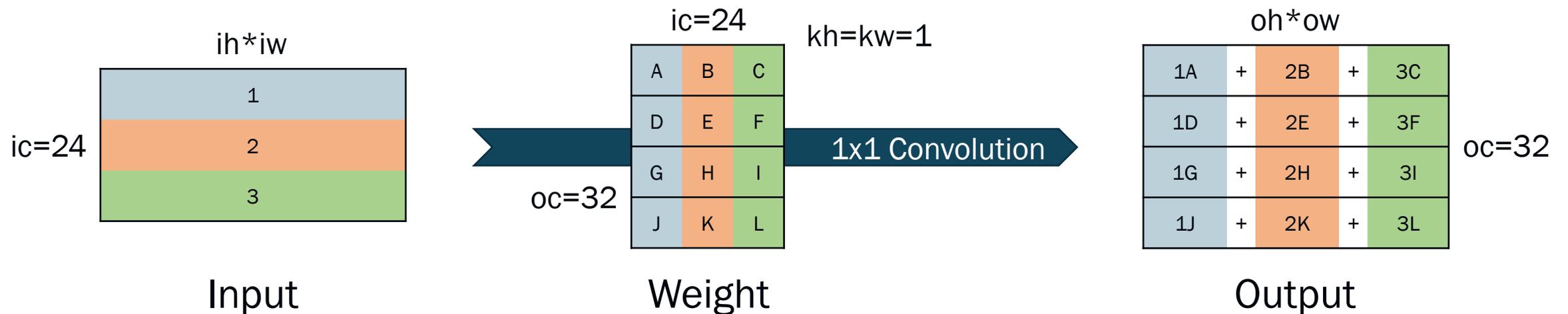
- Accumulating intermediate results of Conv in  $8(\text{oc}) \times 8(\text{ic})$  channels
- For example,

$$\begin{aligned} & \text{Conv} \left( \begin{array}{c} \text{Input} \\ 24(\text{ic}) \times 32(\text{ih}) \times 32(\text{iw}) \end{array}, \begin{array}{c} \text{Weight} \\ 8(\text{oc}) \times 24(\text{ic}) \times 1(\text{kh}) \times 1(\text{kW}) \end{array} \right) \\ &= \text{Conv}_{8(\text{oc}) \times 8(\text{ic})} \left( \begin{array}{c} \text{Input}[0:8,:::] \\ 8(\text{ic}) \times 32(\text{ih}) \times 32(\text{iw}) \end{array}, \begin{array}{c} \text{Weight}[:,0:8,:::] \\ 8(\text{oc}) \times 8(\text{ic}) \times 1(\text{kh}) \times 1(\text{kW}) \end{array} \right) \\ & \quad + \text{Conv}_{8(\text{oc}) \times 8(\text{ic})} \left( \begin{array}{c} \text{Input}[8:16,:::] \\ 8(\text{ic}) \times 32(\text{ih}) \times 32(\text{iw}) \end{array}, \begin{array}{c} \text{Weight}[:,8:16,:::] \\ 8(\text{oc}) \times 8(\text{ic}) \times 1(\text{kh}) \times 1(\text{kW}) \end{array} \right) \\ & \quad + \text{Conv}_{8(\text{oc}) \times 8(\text{ic})} \left( \begin{array}{c} \text{Input}[16:24,:::] \\ 8(\text{ic}) \times 32(\text{ih}) \times 32(\text{iw}) \end{array}, \begin{array}{c} \text{Weight}[:,16:24,:::] \\ 8(\text{oc}) \times 8(\text{ic}) \times 1(\text{kh}) \times 1(\text{kW}) \end{array} \right) \end{aligned}$$

# Case Study – Block-level Accumulation in Convolution

## Hypothetical target hardware

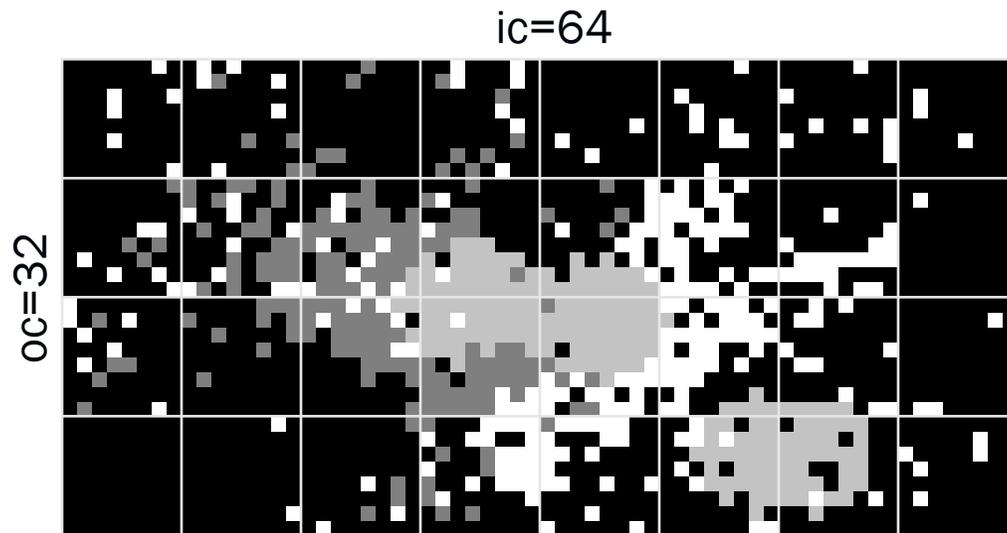
- Accumulating intermediate results of Conv in  $8(oc) \times 8(ic)$  channels
- For example,



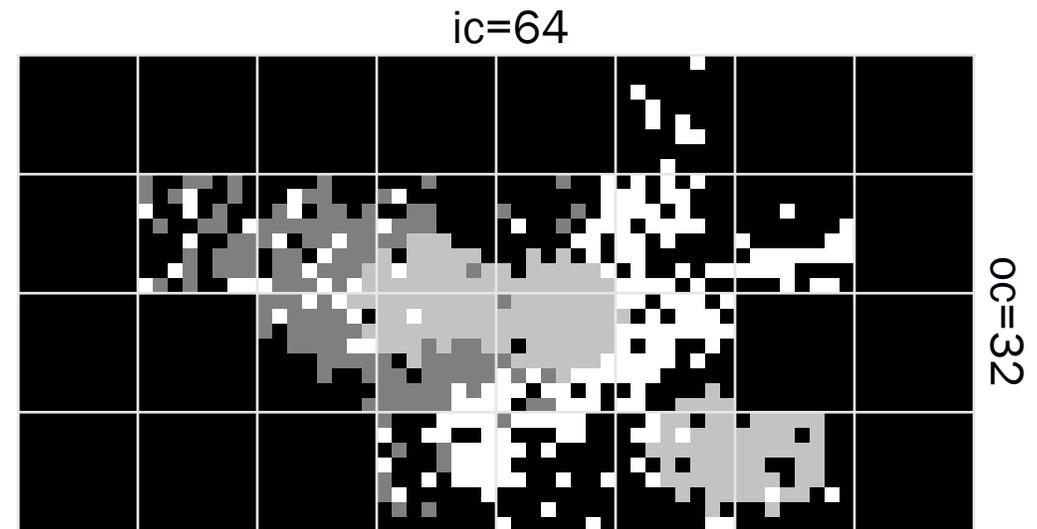
# Case Study – Block-level Accumulation in Convolution

→ 8(oc)x8(ic) block-level sparsification

$$\text{Weight}_{\text{sparse}}(\text{i-th block}) = \begin{cases} 0 & \text{if } \sum |\text{Weight}_{\text{dense}}(\text{i-th block})| < \theta \\ \text{Weight}_{\text{dense}}(\text{i-th block}) & \text{otherwise} \end{cases}$$



1x1 Conv Weight (dense)



1x1 Conv Weight (53% sparse)

# Conclusions

- Layer transformation techniques can reduce production-cost especially for retraining of networks while supporting a diverse array of target platforms
- It may also be beneficial for deep learning hardware manufacturers to support a wide variety of networks not yet natively supported by their hardware

- Sparsification
  - “Exploring the Granularity of Sparsity in Convolution Neural Networks”, CVPR2017  
[https://openaccess.thecvf.com/content\\_cvpr\\_2017\\_workshops/w29/html/Mao\\_Exploring\\_the\\_Granularity\\_CVPR\\_2017\\_paper.html](https://openaccess.thecvf.com/content_cvpr_2017_workshops/w29/html/Mao_Exploring_the_Granularity_CVPR_2017_paper.html)