

AI VISUAL ALGORITHMS AND THEIR COMPUTING SYSTEMS

Kuan-Hung Chen/陳冠宏

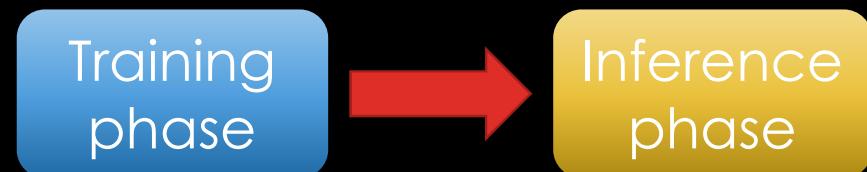
OUTLINE

- How to obtain an accurate AI algorithm?
- What are key computations in AI visual algorithms?
- What are major concerns for an AI visual computing system?

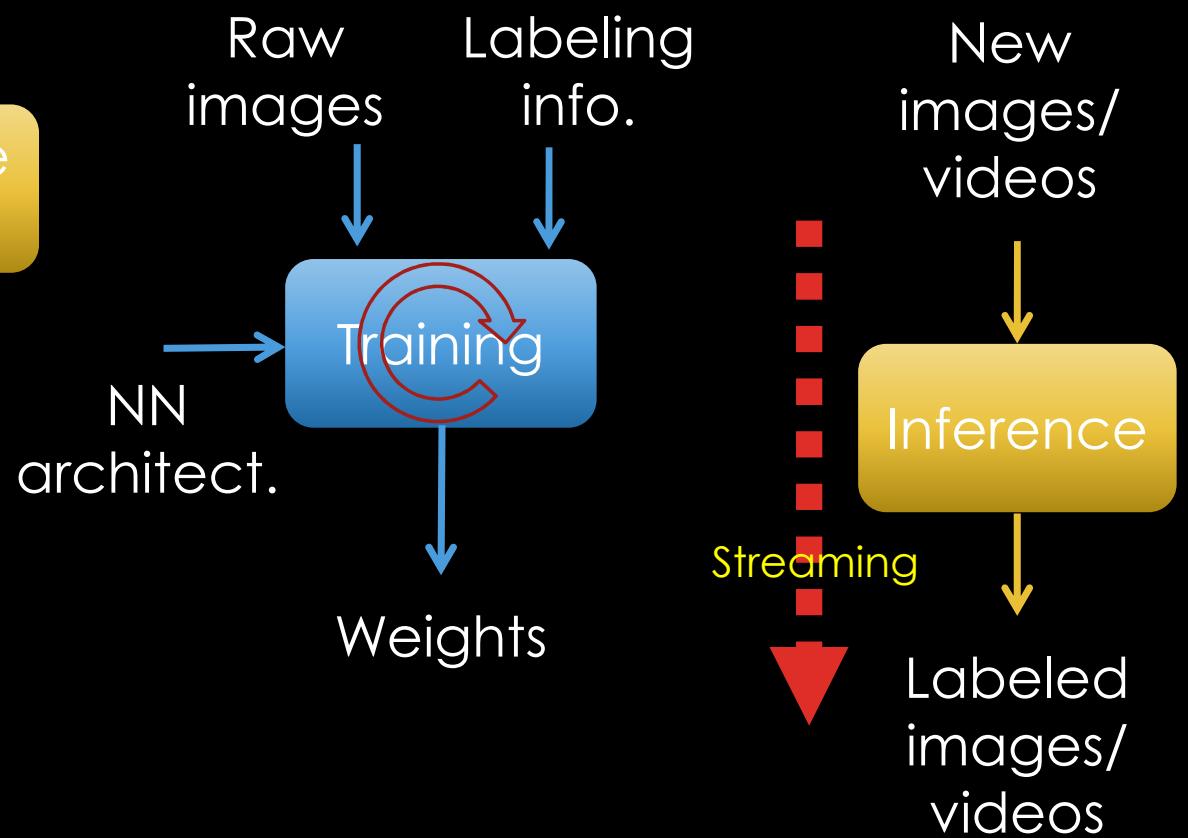
OUTLINE

- How to obtain an accurate AI algorithm?
- What are key computations in AI visual algorithms?
- What are major concerns for a AI visual computing system?

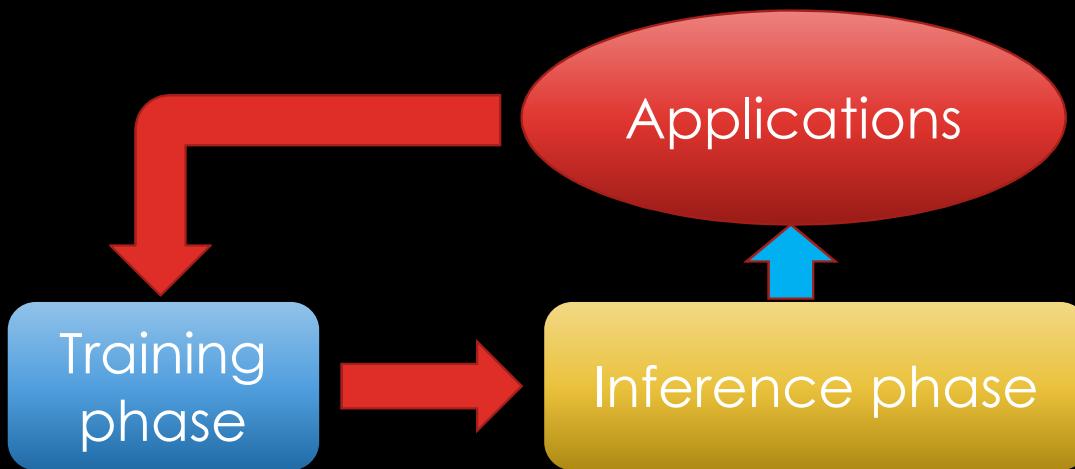
TRAINING VS. INFERENCE



Procedure	Sample amount	Algo. Type	Latency tolerance
Training	Finite	Close loop	High
Inference	Infinite	Open loop	Low



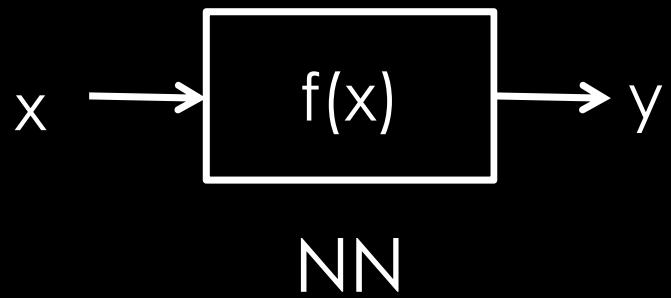
TRAINING VS. INFERENCE



- Data labeling services
- Autonomous vehicle
- Security
- Biomedical imaging
- Smart robots
- Industry 4.0

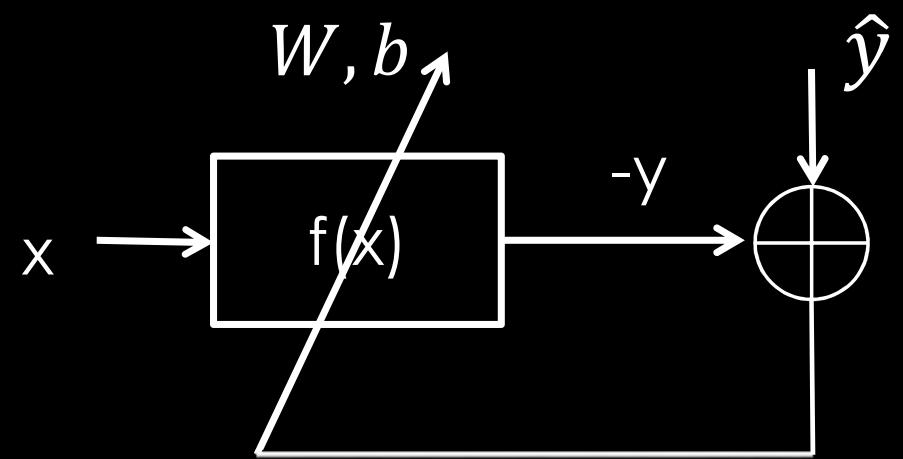
Procedure	Sample amount	Algo. Type	Latency tolerance	Device quantity	Device unit price	Existing type	Operating period
Training	Finite	Close loop	High	Few	High	Cloud	Intermittent
Inference	Infinite	Open loop	Low	Many	Low	Edge	Continued

TRAINING VS. INFERENCE



$$y = f(x) = W \cdot X + b$$

(Inference side)

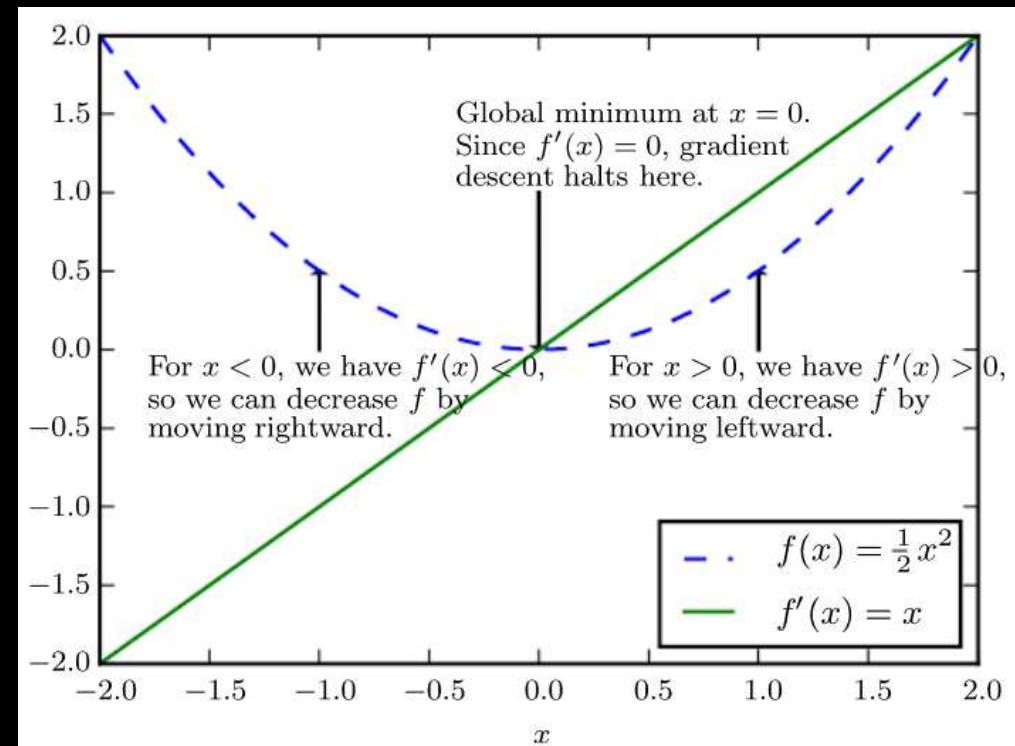


$$|\hat{y} - y| \rightarrow 0$$

(Training side)

GRADIENT DESCENT

- For $x < 0$, $f'(x) < 0$, so we can decrease f by moving rightward
- For $x > 0$, $f'(x) > 0$, so we can decrease f by moving leftward
- The step size is correlated to the value of the gradient
 - Large gradient \rightarrow large step size
 - Small gradient \rightarrow small step size



GRADIENT DESCENT

$$f(x + \underline{\Delta x}) \approx f(x) + \underline{\Delta x \cdot f'(x)}$$

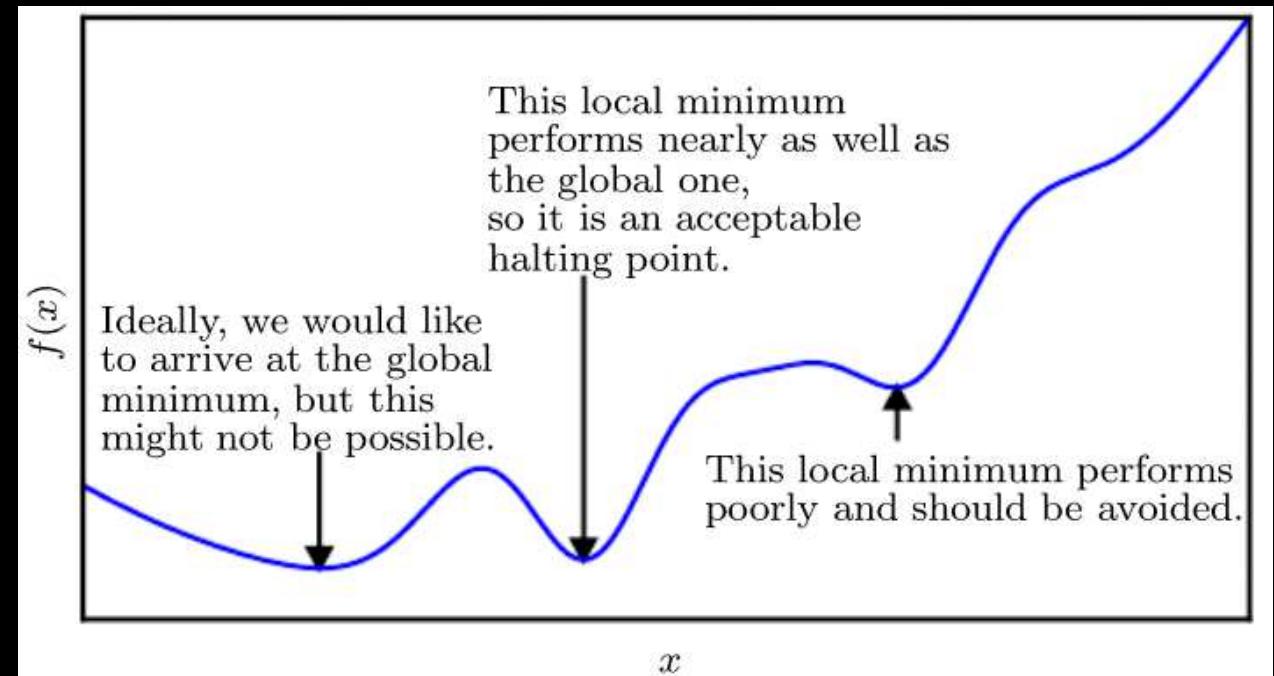
$$x' = x + \Delta x$$

$$= x - (\epsilon \nabla_x f(x))$$

↑ ↑
learning Gradient
rate of $f(x)$

GRADIENT DESCENT

- Local minimum issue
 - The training did not converge to the global minimum
 - The step size may need to become large again



GENERALIZATION

- Generalization- the ability to perform well on previously unseen inputs
 - Training error $\left\| \bar{W} \cdot \bar{x}^{(\text{train})} - y^{(\text{train})} \right\|_2^2$
 - Test error (generalization error) $\left\| \bar{W} \cdot \bar{x}^{(\text{test})} - y^{(\text{test})} \right\|_2^2$
- How well an AI algorithm will perform depends on
 - Make the training error small
 - Make the gap between training error and test error small

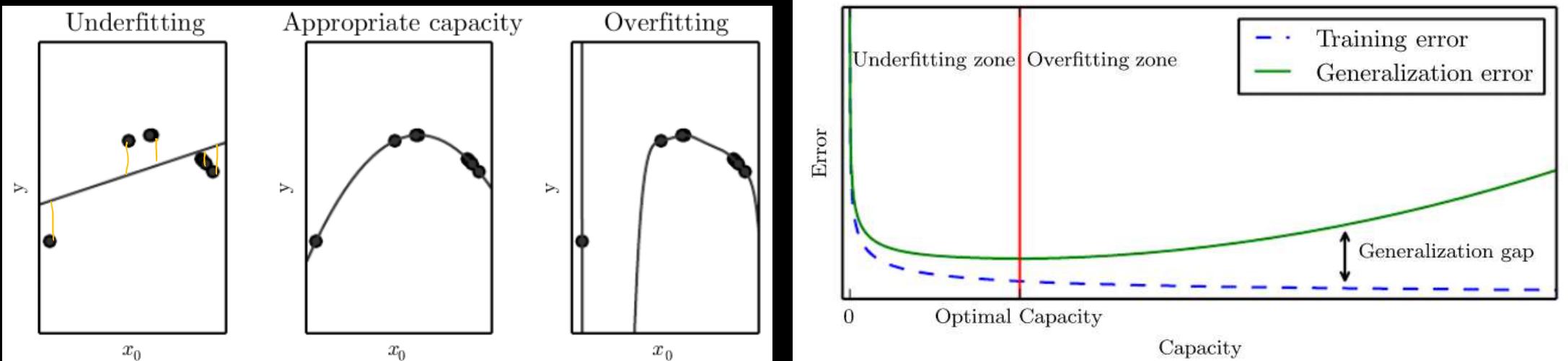
CAPACITY, UNDERFITTING, OVERFITTING

- Capacity-the ability of a model to fit a wide variety of functions
 - Models with low capacity may struggle to fit the training set
 - Models with high capacity can overfit by memorizing properties of the training set that do not serve well on the test set

CAPACITY, UNDERFITTING, OVERFITTING

$$Y_1 = b + \omega X$$

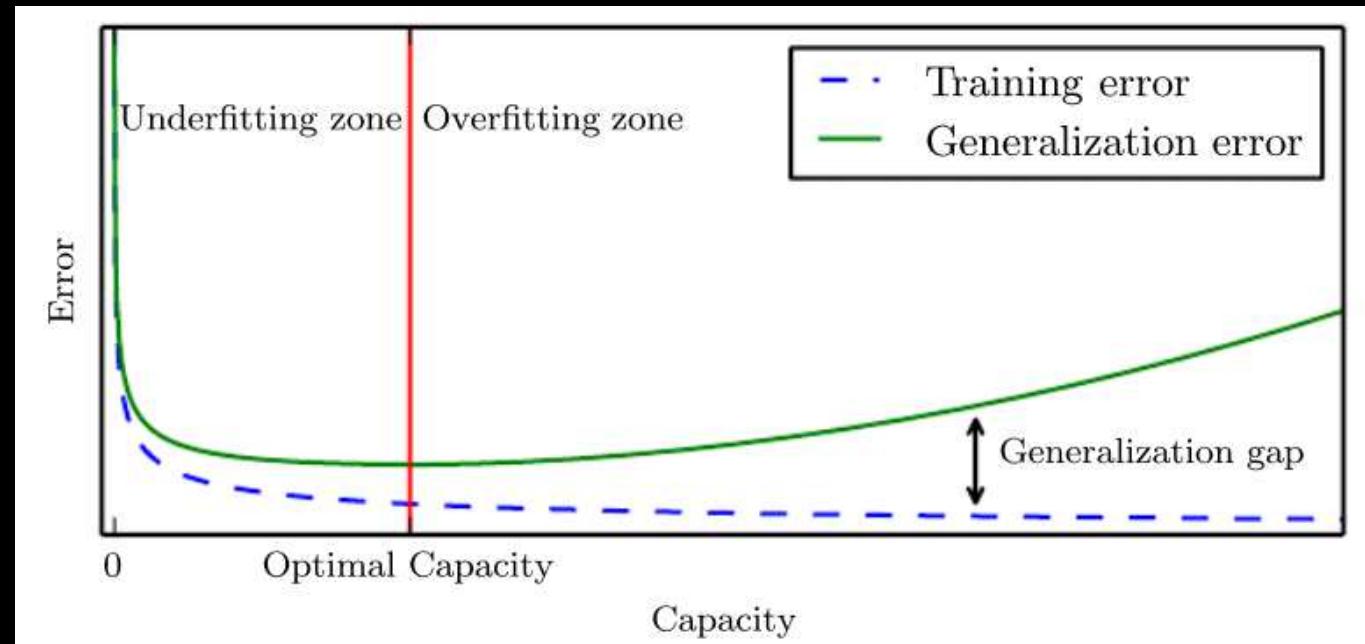
$$Y_2 = b + \omega_1 X + \omega_2 X^2$$



$$Y_q = b + \sum_{i=1}^q \omega_i X^i$$

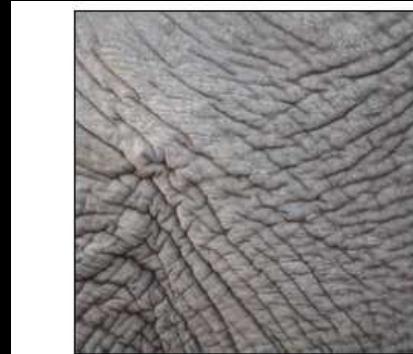
HOW TO OBTAIN AN ACCURATE AI ALGORITHM?

- Increase depth of NN
- Increase the variety of dataset



THE VARIETY OF DATASET (1/2)

- ImageNet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness



(a) Texture image
81.4% **Indian elephant**
10.3% indri
8.2% black swan



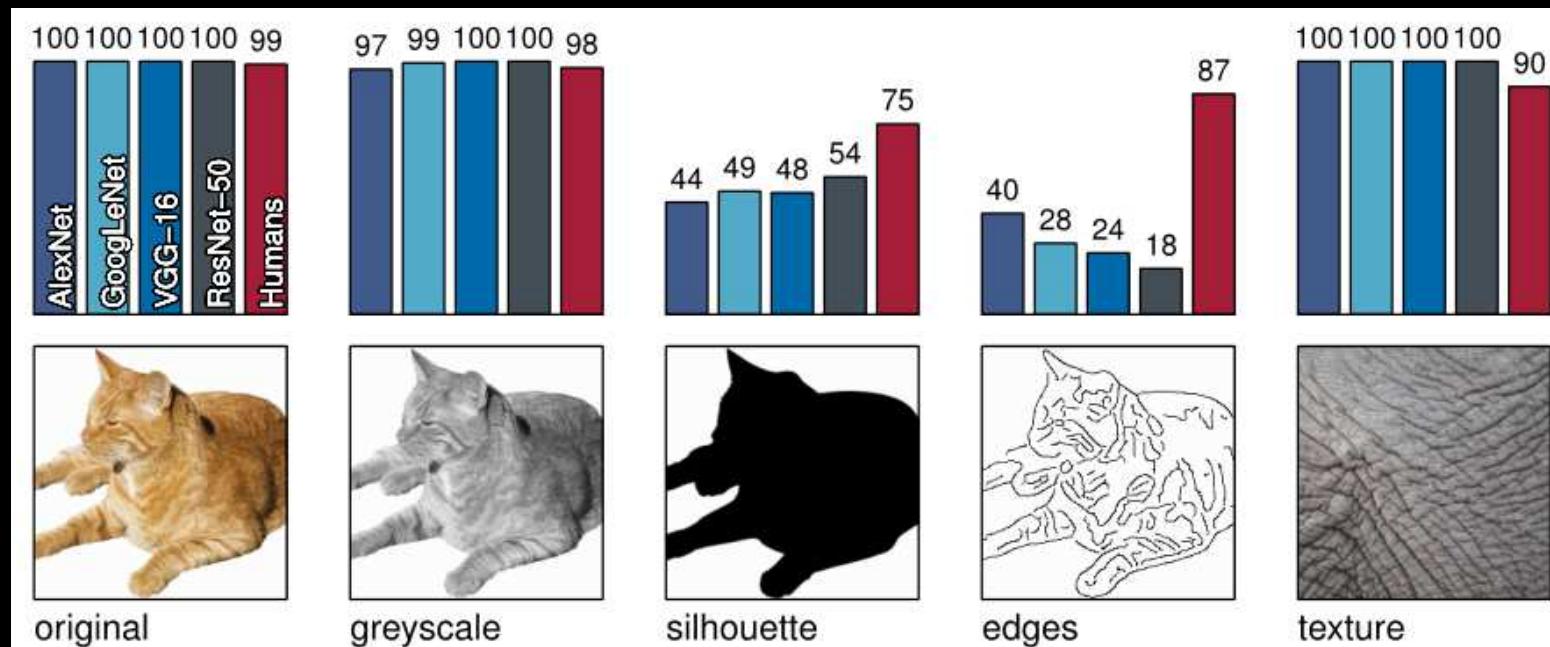
(b) Content image
71.1% **tabby cat**
17.3% grey fox
3.3% Siamese cat



(c) Texture-shape cue conflict
63.9% **Indian elephant**
26.4% indri
9.6% black swan

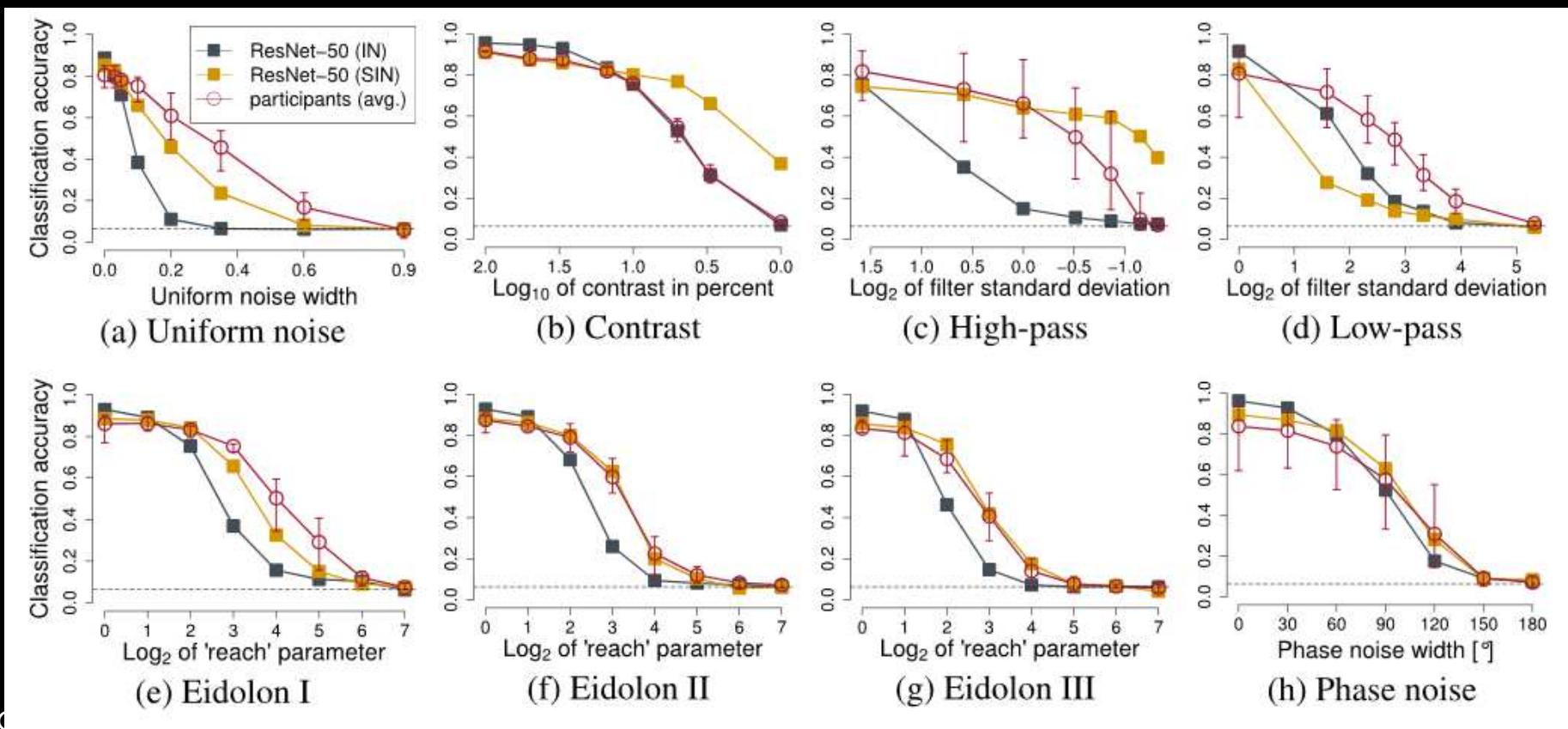
From a conference paper at ICLR 2019

THE VARIETY OF DATASET (2/2)



From a conference paper at ICLR 2019

INCREASING SHAPE BIAS IMPROVES ACCURACY AND ROBUSTNESS



From

OUTLINE

- How to obtain an accurate AI algorithm?
- What are key computations in AI visual algorithms?
- What are major concerns for an AI visual computing system?

KEY COMPUTATIONS IN AI VISUAL ALGORITHMS

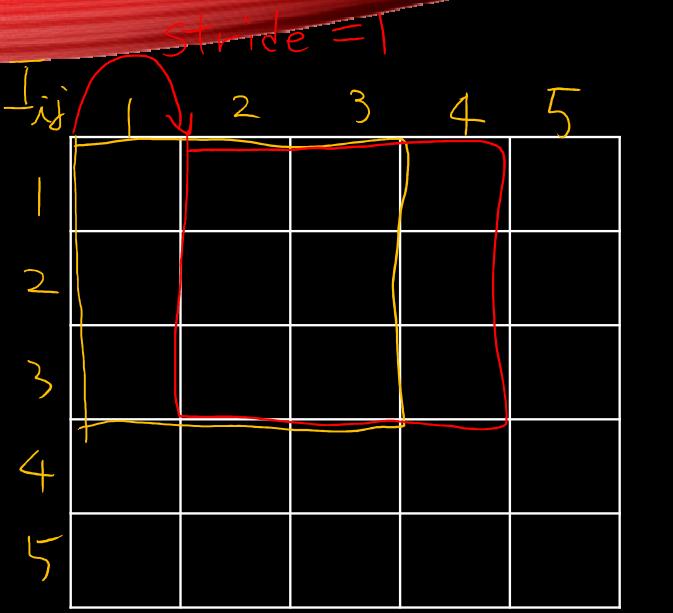
- Convolution: to extract features
- Pooling: down sampling with key info kept
- Activation (ReLU): to add non-linearity
- Residual layer: make deep NN convergent

CONVOLUTION

$$\begin{aligned} & \overline{I_{11}} \cdot \overline{W_{11}} + \overline{I_{12}} \overline{W_{12}} + \overline{I_{13}} \overline{W_{13}} + \\ & \overline{I_{21}} \overline{W_{21}} + \overline{I_{22}} \overline{W_{22}} + \overline{I_{23}} \overline{W_{23}} + \\ & \overline{I_{31}} \overline{W_{31}} + \overline{I_{32}} \overline{W_{32}} + \overline{I_{33}} \overline{W_{33}} = 0_{11} \end{aligned}$$

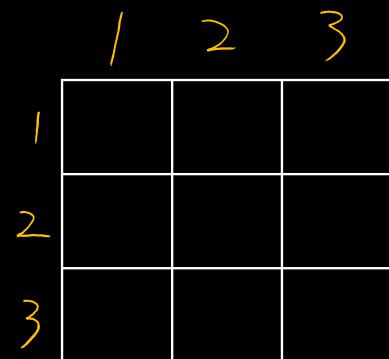
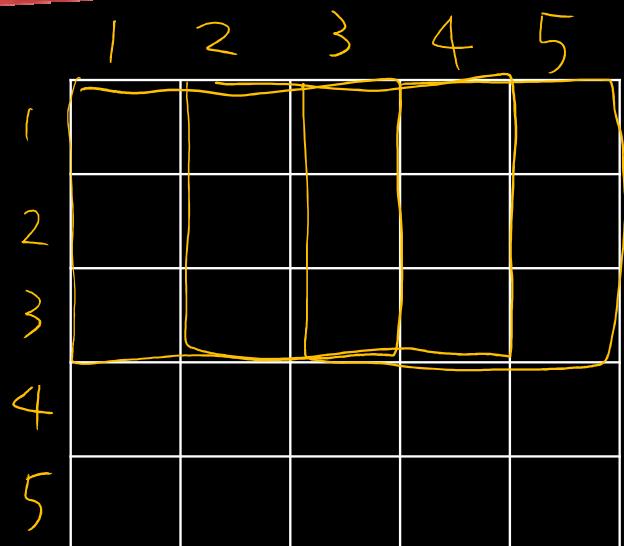
\searrow stride = 1

$$I_{12} \overline{W}_{11} + I_{13} \overline{W}_{12} + I_{14} \overline{W}_{13} + \\ I_{22} \overline{W}_{21} + I_{23} \overline{W}_{22} + I_{24} \overline{W}_{23} + \\ \underline{\hspace{10cm}} = O_{12}$$



A handwritten diagram illustrating a neural network architecture. At the top, the words "Feature map" are written in yellow. Below them, a neural network layer is depicted as a rectangle containing several nodes, with arrows pointing from the "Input" layer down into it. The word "Input" is written in yellow at the bottom of the diagram.

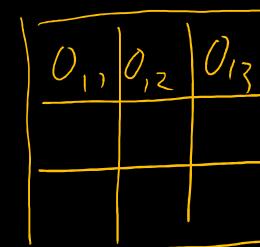
CONVOLUTION



Stride = 1



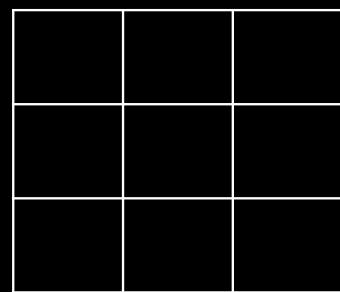
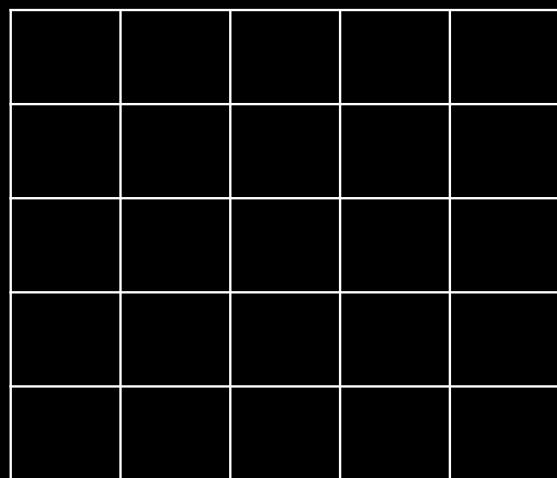
conv.



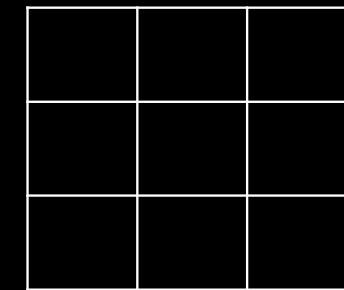
3×3

Output fmap

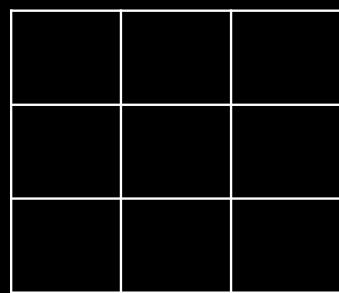
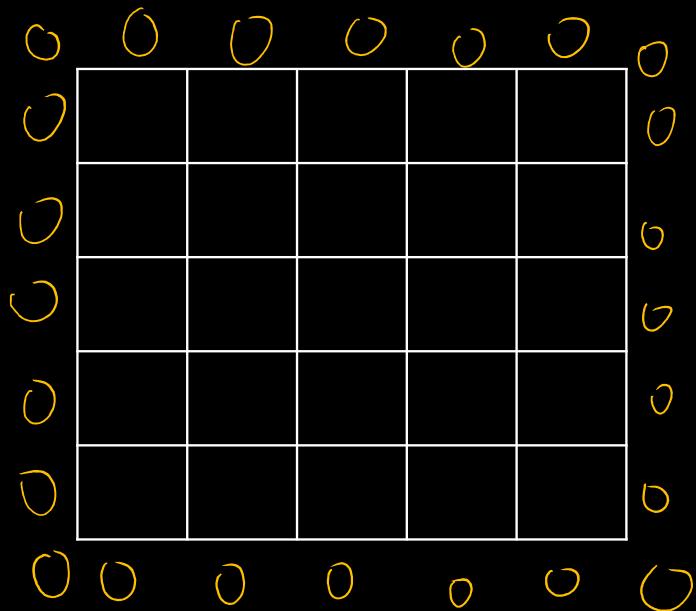
CONVOLUTION



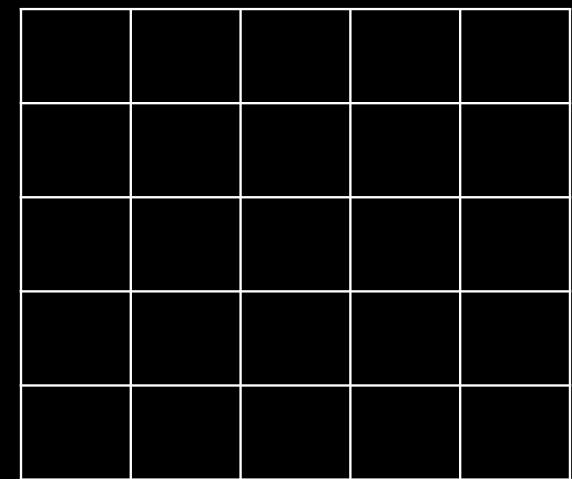
Convolution,
stride=1



ZERO PADDING



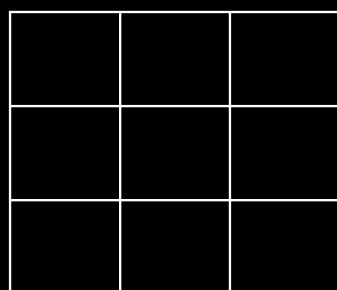
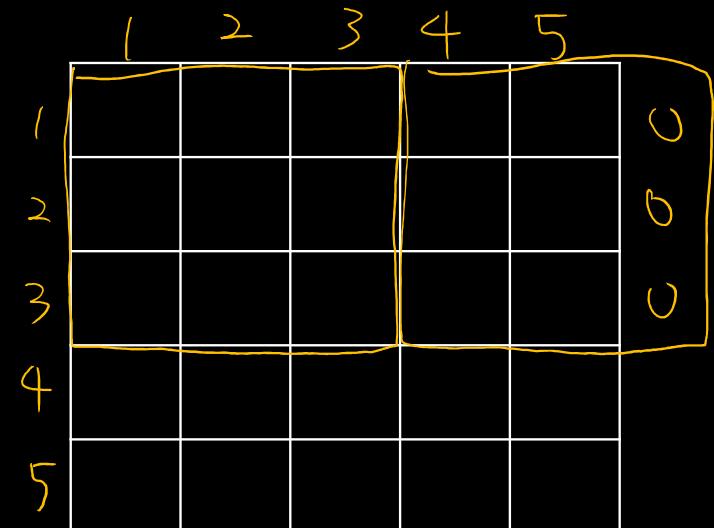
Convolution,
stride=1



5×5

Out Fmap

STRIDE MORE THAN ONE



Stride = 3

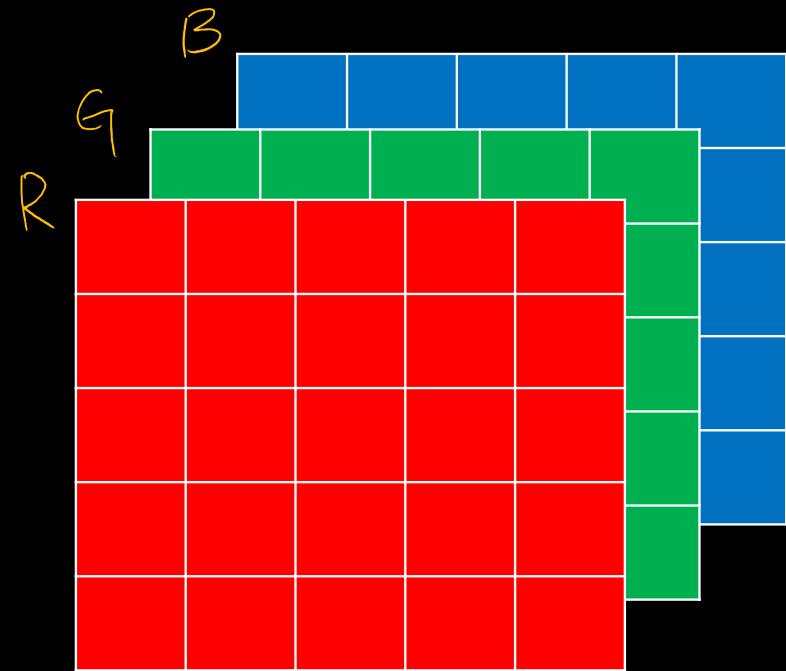
O₁₁ the same as stride = 1

$$O_{12} = I_{14} \cdot \bar{W}_{11} + I_{15} \cdot \bar{W}_{12} + I_{16} \cdot \bar{W}_{13} +$$

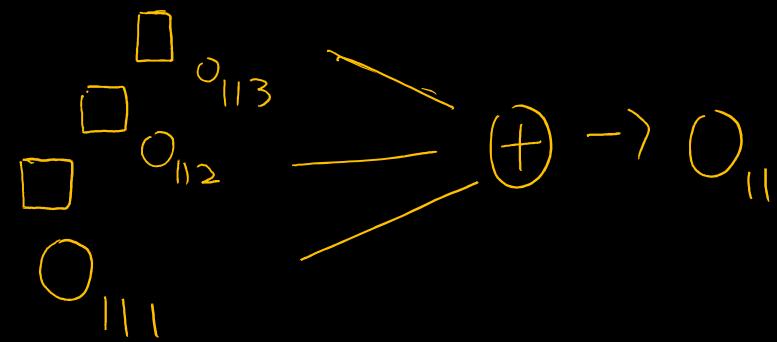
-

-

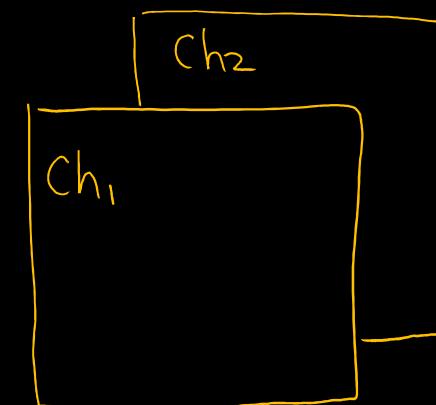
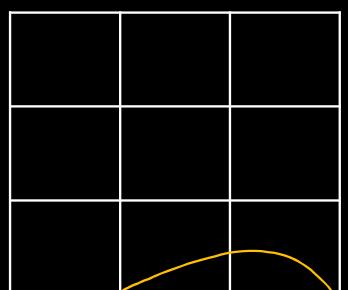
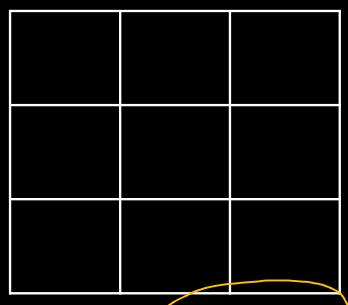
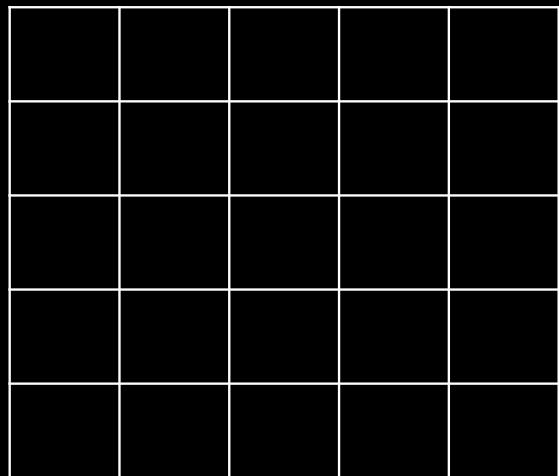
MULTI-CHANNELS



Color image



MULTI-FILTERS



KEY COMPUTATIONS IN AI VISUAL ALGORITHMS

- Convolution: to extract features
- Pooling: down sampling with key info kept
- Activation (ReLU): to add non-linearity
- Residual layer: make deep NN convergent

池化

POOLING (POOL) LAYER

2 x 2 pooling, stride 2

35	10	33	15
25	12	5	21
22	18	19	25
7	20	22	13

Max pooling

35	33
22	25

Average pooling

20	19
17	20

POOL LAYER IMPLEMENTATION

```
for (n=0;n<N;n++) {  
    for (m=0;m<M;m++) {  
        for (x=0;x<F;x++) {  
            for (y=0;y<E;y++) {  
                max = -Inf;  
                for (i=0;i<R;i++) {  
                    for (j=0;j<S;j++) {  
                        if(I[n][m][Ux+i][Uy+j]>max){  
                            max =I[n][m][Ux+i][Uy+j];  
                        }  
                    }  
                }  
                O[n][m][x][y] = max;  
            }  
        }  
    }  
}
```

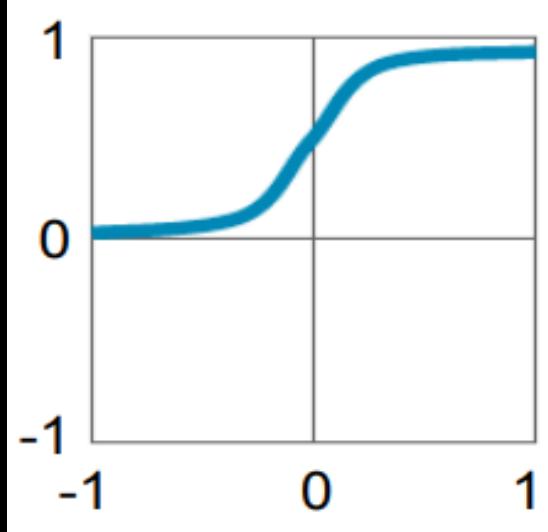
搜尋max pool值

KEY COMPUTATIONS IN AI VISUAL ALGORITHMS

- Convolution: to extract features
- Pooling: down sampling with key info kept
- Activation (ReLU): to add non-linearity
- Residual layer: make deep NN convergent

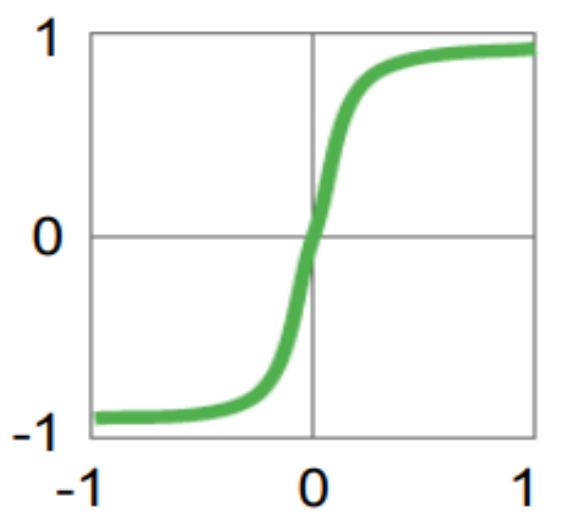
TRADITIONAL ACTIVATION FUNCTIONS

Sigmoid



$$y = 1/(1 + e^{-x})$$

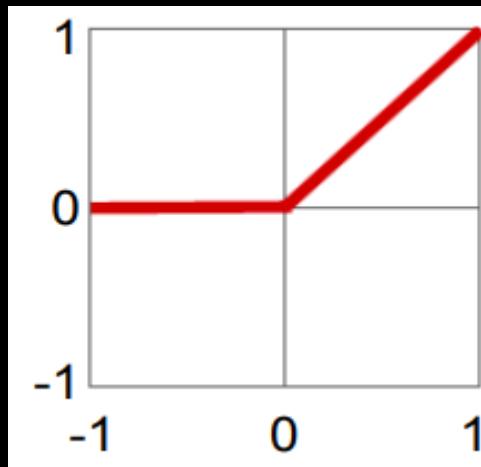
Hyperbolic Tangent



$$y = (e^x - e^{-x})/(e^x + e^{-x})$$

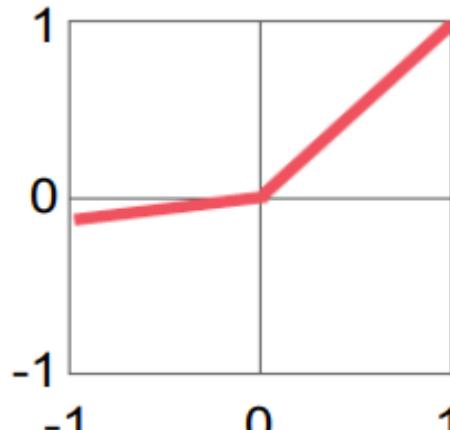
MODERN ACTIVATION FUNCTIONS

Rectified Linear Unit
(ReLU)



$$y = \max(0, x)$$

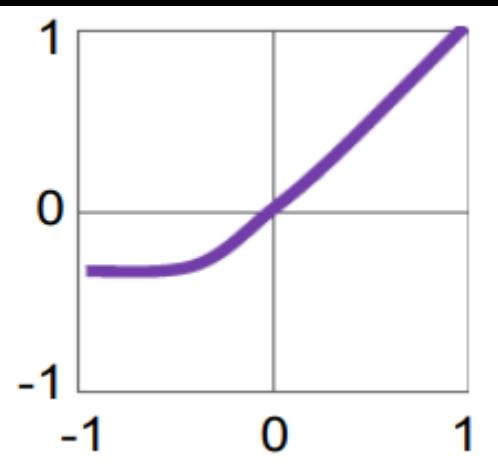
Leaky ReLU



$$y = \max(ax, x)$$

$a = \text{small const}$

Exponential LU



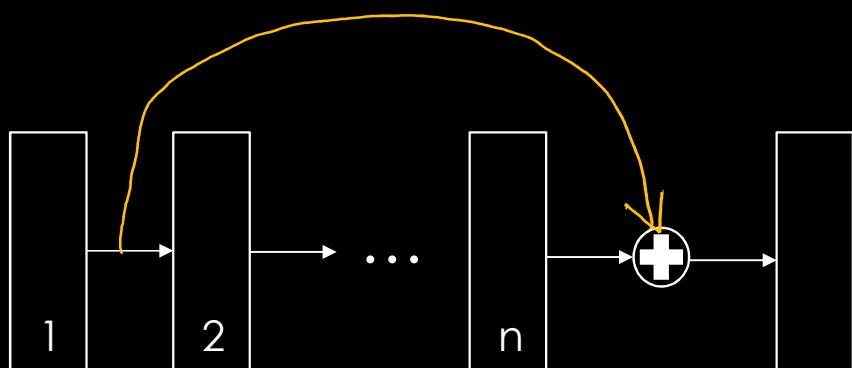
$$y = \text{exp}(x^2 - 1)$$

KEY COMPUTATIONS IN AI VISUAL ALGORITHMS

- Convolution: to extract features
- Pooling: down sampling with key info kept
- Activation (ReLU): to add non-linearity
- Residual layer: make deep NN convergent

RESIDUAL LAYER

- Also named shortcut or skip connection
- To solve the gradient vanishing problem



⊕ : element-wise addition or concatenation

$$\begin{matrix} I_{ij} \\ \oplus \\ F_{ij} \end{matrix} = O_{ij}$$

The same dimension and channel numbers



The same dimension number, but channel numbers are not necessary the same

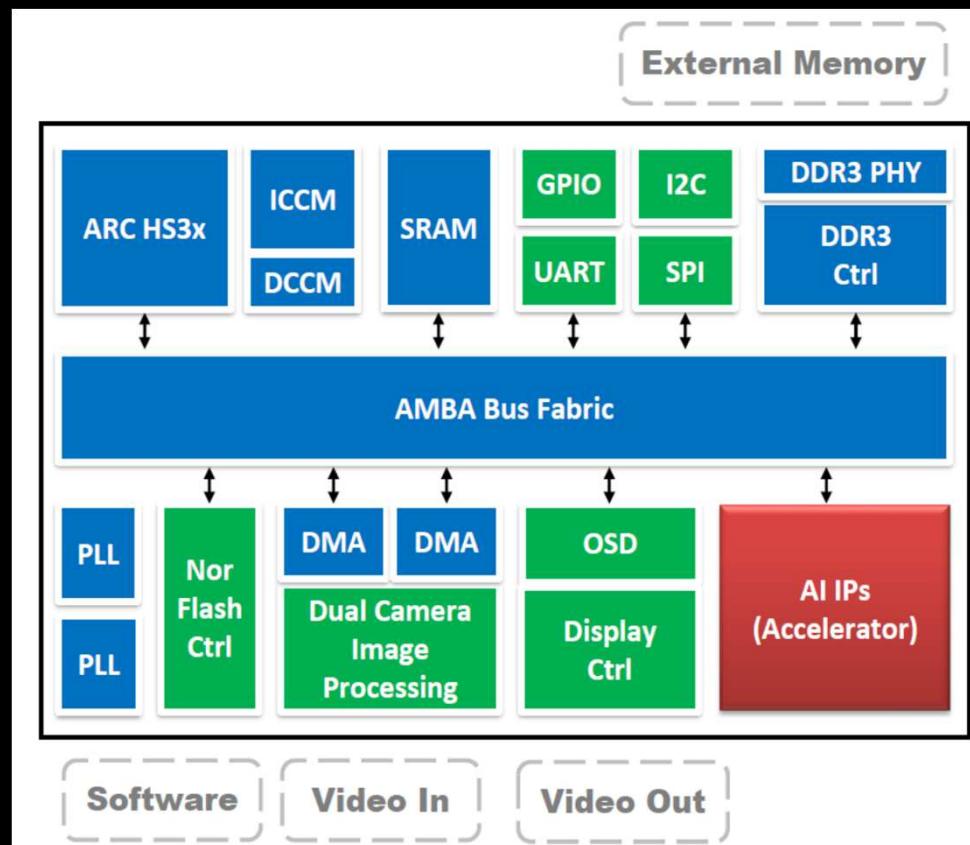
OUTLINE

- How does object classification work on computers?
- How does object detection work on computers?
- What are key computations in AI visual algorithms?
- What are major concerns for a AI visual computing system?

A COMPUTING SYSTEM

- Mainly composed of processing element (PE), cache/SRAM/on-chip buffer, DRAM, and peripherals.
- A PE is composed of one multiplier and one adder, or sometimes one multiplier and accumulator (MAC), several registers, and a local controller.

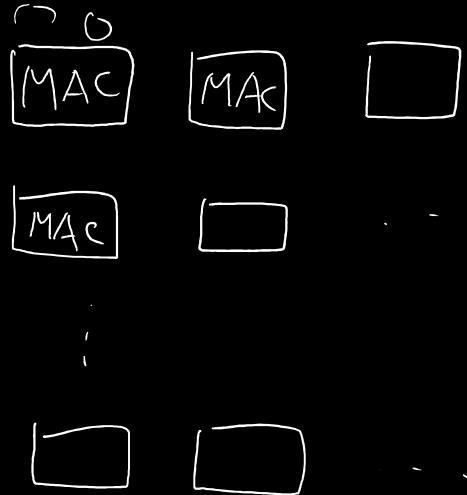
AI SOC DESIGN PLATFORM OF CIC



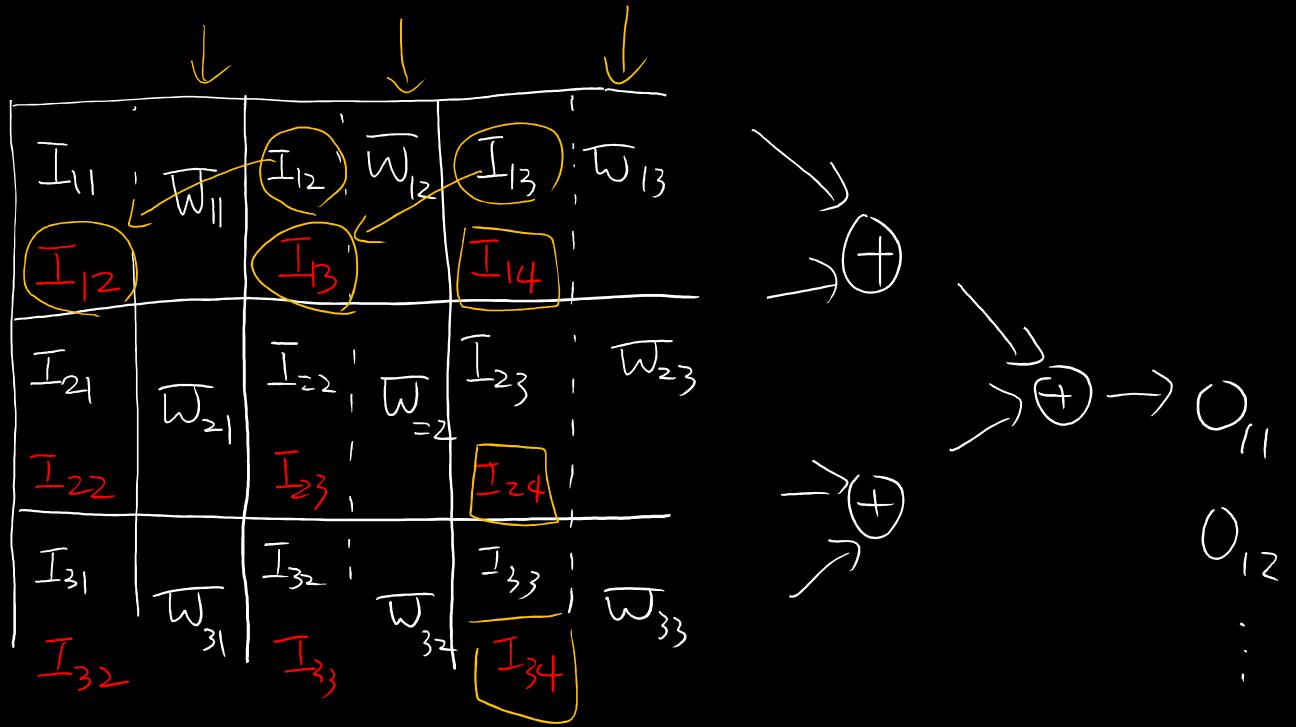
PARALLEL COMPUTATION

- Utilize data independency
 - Convolution operation in AI visual networks



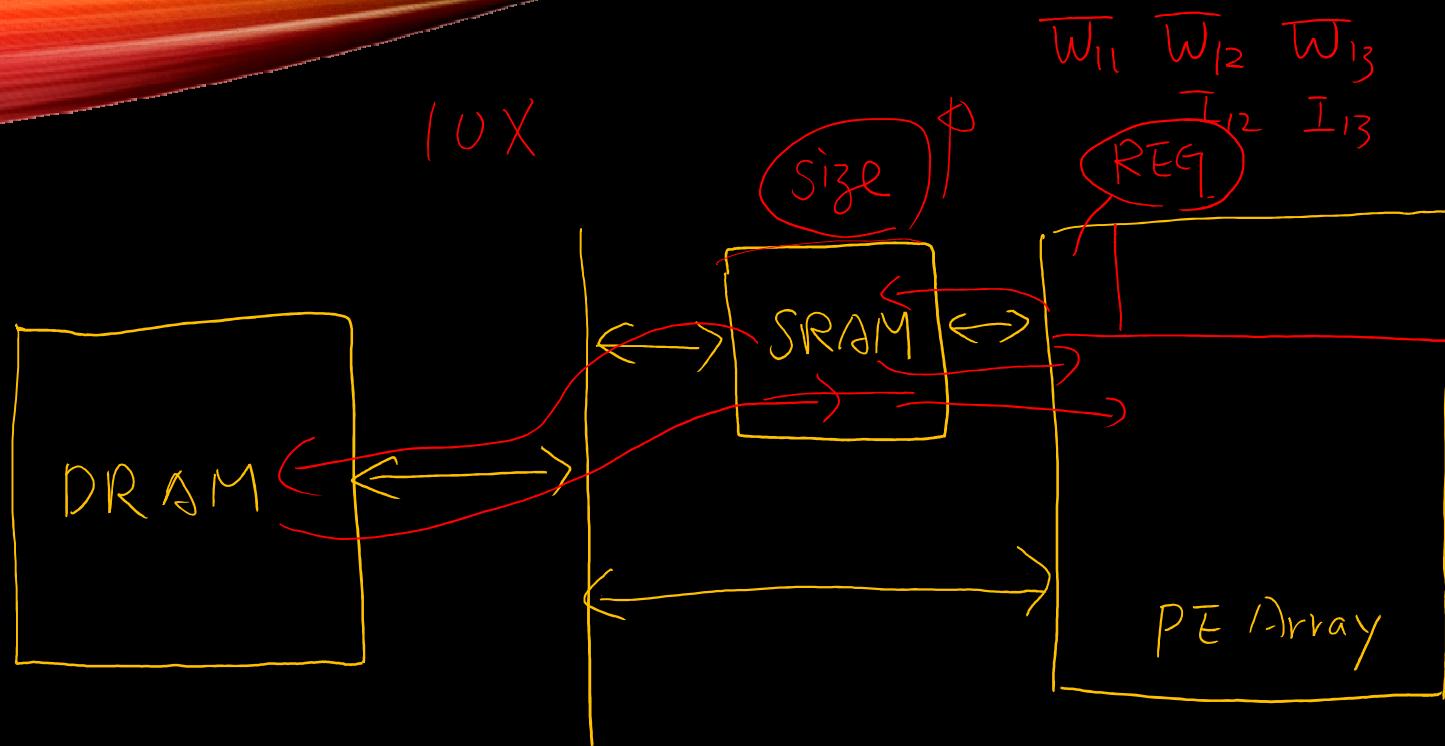


PE Array



DATA REUSE

- To lower down the DRAM transfer burden, and so that to accelerate the AI computation and save energy consumption.
- Need local storage. The size of local buffer versus the data reuse ratio is an important design trade-off.



DRAM BANDWIDTH (DBW)

- DRAM: Dynamic Random Access Memory, be used as a massive data storage in computing systems due to its low-cost merit.

BASIS FOR COMPARISON	SRAM	DRAM
Speed	Faster	Slower
Size	Small	Large
Cost	Expensive	Cheap
Used in	Cache memory	Main memory
Density	Less dense	Highly dense
Construction	Complex and uses transistors and latches.	Simple and uses capacitors and very few transistors.
Single block of memory requires	6 transistors	Only one transistor.
Charge leakage property	Not present	Present hence require power refresh circuitry
Power consumption	Low	High

DBW REQUIREMENTS FOR COMPUTING CNNS

- Because of the limited on-chip cache size, data such as IFMs, weights, and OFMs are necessary to be moved between DRAM and SRAM. This forms the DRAM bandwidth requirement.
- Taking Agilev3 for example, the data transferred between DRAM and SRAM for 30 fps of 416*416 input image resolution may be as high as 3.02 GB/s based on that 72kB kernel SRAM and 169kB IFM SRAM are equipped on-chip.

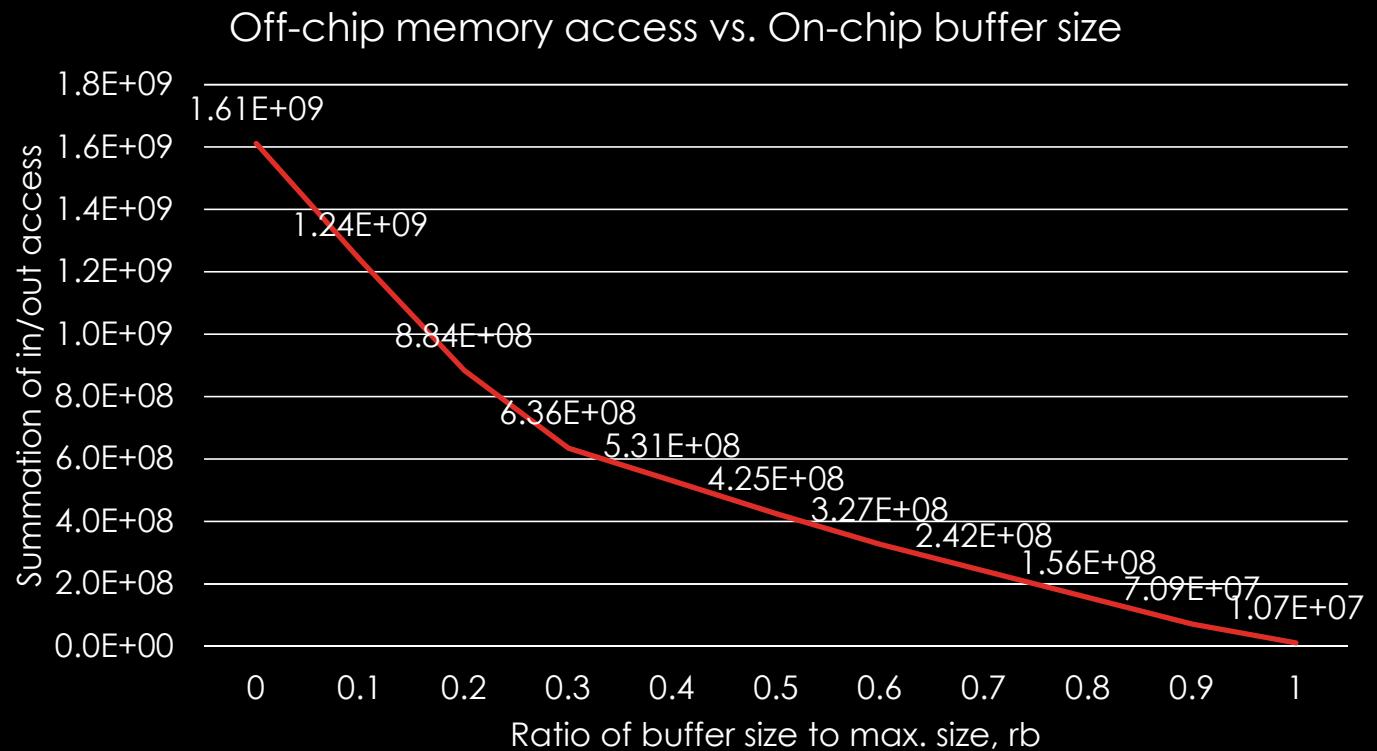
TOTAL AVAILABLE DBW

- The total bandwidth is the product of
 - **Base DRAM clock frequency**
 - **Number of data transfers per clock:** Two, in the case of double data rate (DDR, DDR2, DDR3, DDR4) memory.
 - **Memory bus (interface) width:** Each DDR, DDR2, or DDR3 memory interface is 64 bits wide. Those 64 bits are sometimes referred to as a line.
 - **Number of interfaces:** Modern personal computers typically use two memory interfaces (dual-channel mode) for an effective 128-bit bus width.
- For example, a computer with dual-channel memory and one DDR2-800 module per channel running at 400 MHz would have a theoretical maximum memory bandwidth of:

400,000,000 clocks per second \times 2 lines per clock \times 64 bits per line \times 2 interfaces =
102,400,000,000 (102.4 billion) bits per second (in bytes, 12,800 MB/s or 12.8 GB/s)

LIMITED ON-CHIP BUFFER SIZE EFFECT

- NN: YOLO v2
- Input image resolution: 224*224
- Note
 - The data is for only a single frame



SUGGESTED DRAM TYPE

Width	Height	fps	DBW needed (GB/s) Original/optimized	Suggested DRAM type
416	416	30	3.02/1.93*	DDR-400, PC-3200
1080	720	30	13.57/8.67*	DDR4-2400, PC4-19200/DR3-1600, PC3-12800
1920	1080	30	36.19/23.13*	No available DRAM/DDR4-3200, PC4-25600

* DBW of **integer** AgileNet, a light-weight CNN for mobile use.

Names	Memory clock	I/O bus clock	Transfer rate	Theoretical bandwidth
DDR-200, PC-1600	100 MHz	100 MHz	200 MT/s	1.6 GB/s
DDR-400, PC-3200	200 MHz	200 MHz	400 MT/s	3.2 GB/s
DDR2-800, PC2-6400	200 MHz	400 MHz	800 MT/s	6.4 GB/s
DDR3-1600, PC3-12800	200 MHz	800 MHz	1600 MT/s	12.8 GB/s
DDR4-2400, PC4-19200	300 MHz	1200 MHz	2400 MT/s	19.2 GB/s
DDR4-3200, PC4-25600	400 MHz	1600 MHz	3200 MT/s	25.6 GB/s

MOBILE AI PLATFORMS

Platform	CPU	GPU	Performance	DBW (GB/s)
Jetson Nano	4 cortex A57	128 CUDA cores	472 GOPs	25.6
Jetson TX2	2 Denver cores and 4 cortex A57	256 pascal gpu cores	1.33 TOPs	59.7
Jetson AGX Xavier	8 Carmel cores and ARM 8.2 64b CPU	512 volta gpu cores with 64 tensor cores	32 TOPs	136.5

SUGGESTED PLATFORM

Width	Height	fps	Operation required (GOPS)	Network	Suggested platform
416	416	30	981	AgileV3	Jetson TX2 or Jetson Nano for 14 fps
1080	720	30	4405		Jetson TX2
1920	1080	30	11752		Jetson AGX Xavier

Width	Height	fps	Operation required (GOPS)	Network	Suggested platform
416	416	30	1560	YOLOv3	Jetson TX2/Jetson AGX Xavier
416	416	30	1803	YOLOv4	Jetson TX2/Jetson AGX Xavier

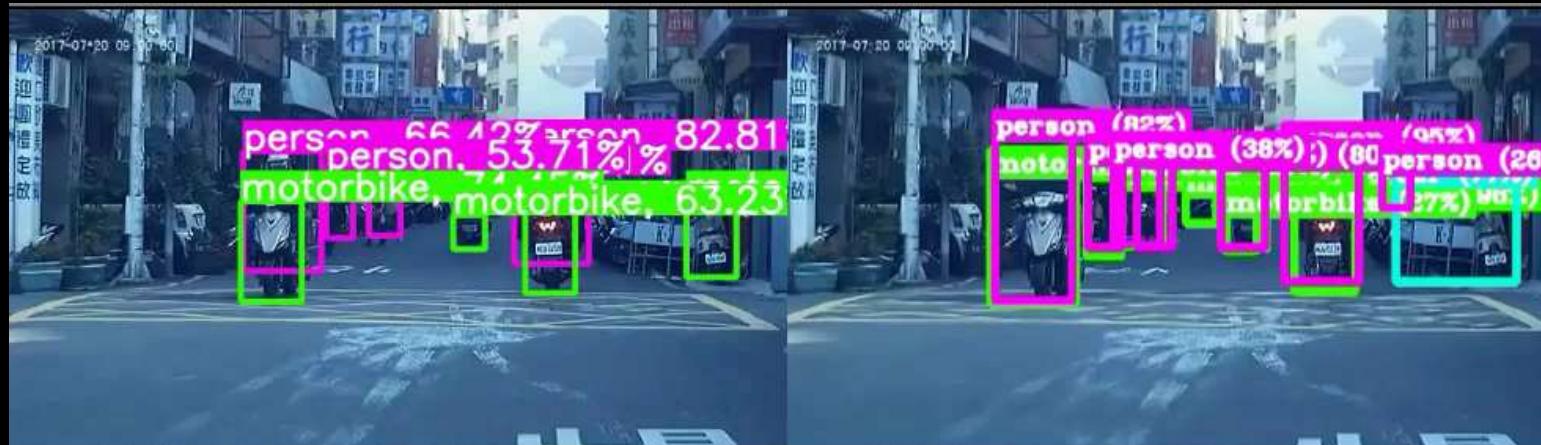
COMPARISON OF DETECTION NNS

Network	Word length	No. of conv. layers	Model size (MB)	Conv. IO* (Mega)	Required GOPS*	Year
Agilev3	FP32	43	65.39	2023.8	480	2019
YOLOv3	FP32	74	241.78	3352.5	980	2018
YOLOv4	FP32	109	251.15	4795.8	900	2020
YOLOv4-tiny	FP32	21	23.10	707.1	100	2020
HarDNet+SSD	FP32	88	98.04	2361.3	770	2019
YOLOv5 - s	FP16	70	14.2	1075.8	110	2020

* for 416x416 @ 30 fps

QUALITATIVE COMPARISON

agilev3



yolov3

yolov4



yolov5s

The raw video is provided by mPower Microelectronics Inc.

ADVANCED OPTIMIZATION TOPICS

- Auto evaluation of customized testing samples
- Auto labeling of sample classes that are not available yet
 - Labelimg, Ezlabel, llabeler

PROS AND CONS OF GROUNDDRUTH LABELING TOOLS

Labeling tool	Pros	Cons
Labelimg	<ul style="list-style-type: none"> 1. Easy to use 2. User-friendly Interface 	<ul style="list-style-type: none"> 1. Manual labeling, 2. Time-consuming for multiple images labeling 3. No video labeling.
EZLabel	<ul style="list-style-type: none"> 1. User-friendly Interface 2. Auto image labeling 3. A cloud-based operation platform 	<ul style="list-style-type: none"> 1. No video labeling
llabeler-v1	<ul style="list-style-type: none"> 1. Auto image labeling 2. Auto video labeling 3. Allow to interactively update image labeling result. 	<ul style="list-style-type: none"> 1. On-site computer operation only 2. Need sufficient computing resource

THANK YOU FOR LISTENING