香港中文大學
The Chinese University of Hong Kong

# Hardware Friendly Computer Vision

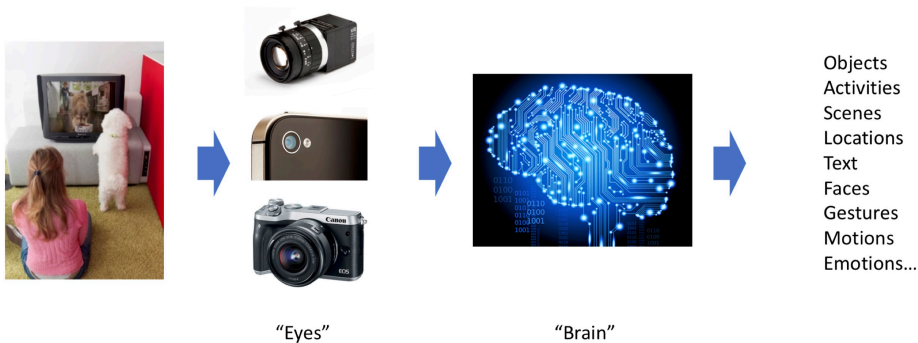Bei Yu
Department of Computer Science & Engineering
Chinese University of Hong Kong
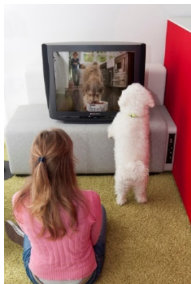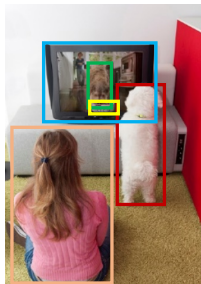byu@cse.cuhk.edu.hk

November 15, 2022

# Computer Vision

- Humans use their eyes and their brains to visually sense the world.
- Computers user their cameras and computation to visually sense the world



"Eyes"                    "Brain"

Objects
Activities
Scenes
Locations
Text
Faces
Gestures
Motions
Emotions...

Classification     Detection     Segmentation     Sequence

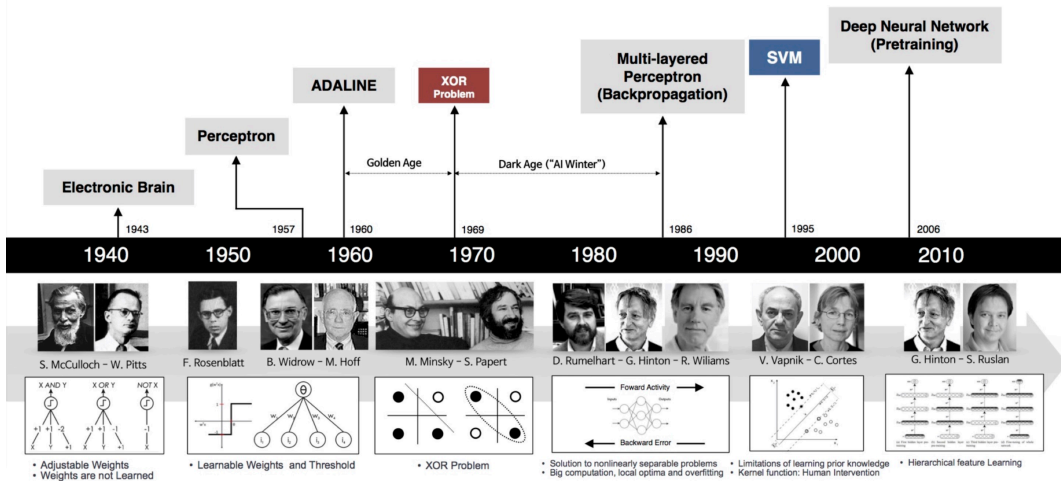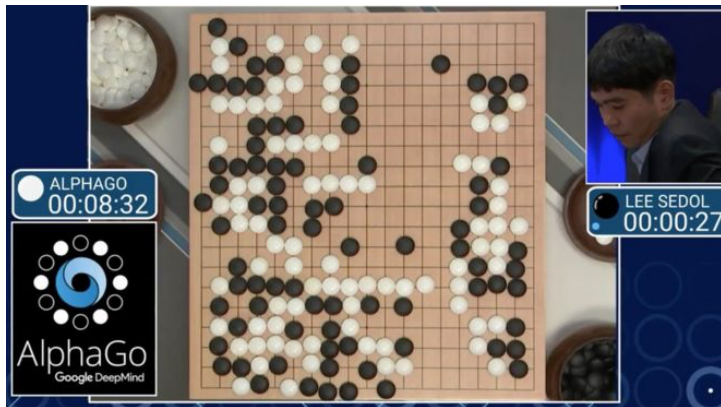Image   →   Region   →   Pixel     Video

- The rises of SVM, Random forest
- No theory to play
- Lack of training data
- Benchmark is insensitive
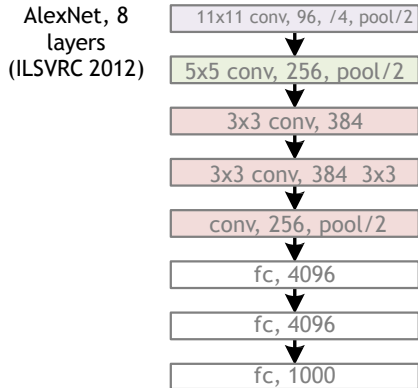- Difficulties in optimization
- Hard to reproduce results

## Curse

"Deep neural networks are no good and could never be trained."

- A fast learning algorithm for deep belief nets. [Hinton et.al 1996]
- Data + Computing + Industry Competition
- NVidia's GPU, Google Brain (16,000 CPUs)
- Speech: Microsoft [2010], Google [2011], IBM
- Image: AlexNet, 8 layers [Krizhevsky et.al 2012] (26.2% -> 15.3%)
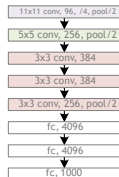
# Revolution of Depth

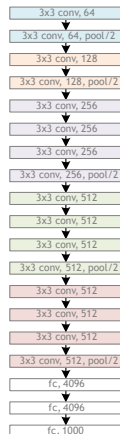AlexNet, 8
layers
(ILSVRC 2012)

| 11x11 conv, 96, /4, pool/2 |
| 5x5 conv, 256, pool/2 |
| 3x3 conv, 384 |
| 3x3 conv, 384  3x3 |
| conv, 256, pool/2 |
| fc, 4096 |
| fc, 4096 |
| fc, 1000 |

# Revolution of Depth

AlexNet, 8 layers (ILSVRC 2012)

```
11x11 conv, 96, /4, pool/2
5x5 conv, 256, pool/2
3x3 conv, 384
3x3 conv, 384
3x3 conv, 256, pool/2
fc, 4096
fc, 4096
fc, 1000
```

VGG, 19 layers (ILSVRC 2014)

```
3x3 conv, 64
3x3 conv, 64, pool/2
3x3 conv, 128
3x3 conv, 128, pool/2
3x3 conv, 256
3x3 conv, 256
3x3 conv, 256
3x3 conv, 256, pool/2
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512, pool/2
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512, pool/2
fc, 4096
fc, 4096
fc, 1000
```

GoogleNet, 22 layers (ILSVRC 2014)



Slide Credit: He et al. (MSRA)

# Revolution of Depth

AlexNet, 8 layers (ILSVRC 2012)

VGG, 19 layers (ILSVRC 2014)
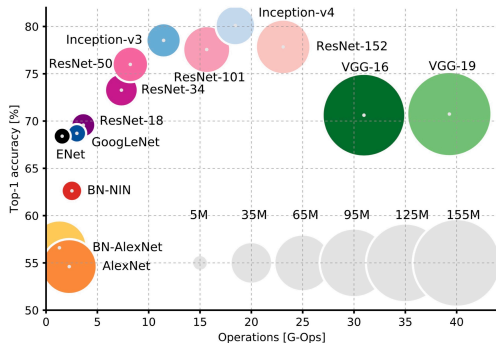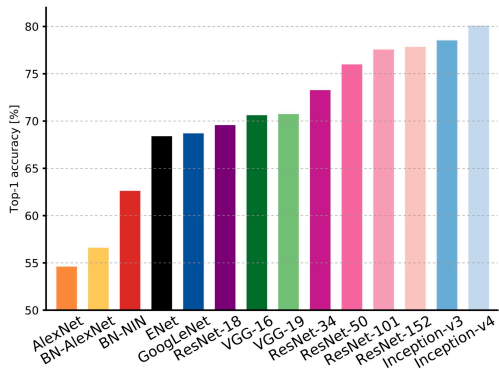
ResNet, 152 layers (ILSVRC 2015)

- AlexNet (Krizhevsky, Sutskever, and E. Hinton 2012) 233MB
- Network in Network (Lin, Q. Chen, and Yan 2013) 29MB
- VGG (Simonyan and Zisserman 2015) 549MB
- GoogleNet (Szegedy, Liu, et al. 2015) 51MB
- ResNet (K. He et al. 2016) 215MB
- Inception-ResNet (Szegedy, Vanhoucke, et al. 2016)
- DenseNet (Huang et al. 2017)
- Xception (Chollet 2017)
- MobileNetV2 (Sandler et al. 2018)
- ShuffleNet (Zhang, Zhou, et al. 2018)

- AlexNet (Krizhevsky, Sutskever, and E. Hinton 2012) 233MB

- Network in Network (Lin, Q. Chen, and Yan 2013) 29MB

- VGG (Simonyan and Zisserman 2015) 549MB

- GoogleNet (Szegedy, Liu, et al. 2015) 51MB

- ResNet (K. He et al. 2016) 215MB

- Inception-ResNet (Szegedy, Vanhoucke, et al. 2016) 23MB

- DenseNet (Huang et al. 2017) 80MB

- Xception (Chollet 2017) 22MB

- MobileNetV2 (Sandler et al. 2018) 14MB

- ShuffleNet (Zhang, Zhou, et al. 2018) 22MB

---

[1] Alfredo Canziani, Adam Paszke, and Eugenio Culurciello (2017). "An analysis of deep neural network models for practical applications". In: *arXiv preprint*.
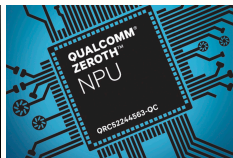
Convolution layer is one of the most expensive layers

- Computation pattern
- Emerging challenges

## More and more end-point devices with limited memory

- Cameras
- Smartphone
- Autonomous driving

Hard to distribute large models through over-the-air update



**This item is over 100MB.**
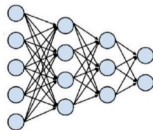Microsoft Excel will not download
until you connect to Wi-Fi.

Cancel          OK

2

[2]Song Han and William J. Dally (2018). "Bandwidth-efficient Deep Learning". In: *Proc. DAC*, 147:1–147:6.

AlphaGo: 1920 CPUs and 280 GPUs,
**$3000 electric bill** per game

on mobile: drains battery
on data-center: increases TCO

[3]

---

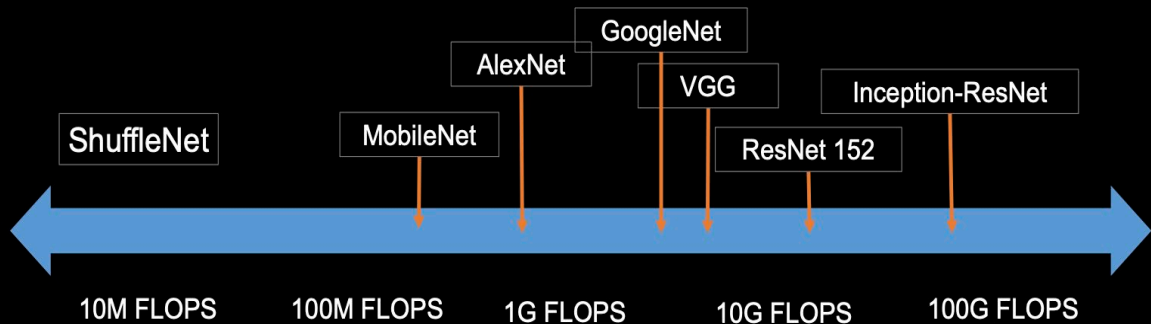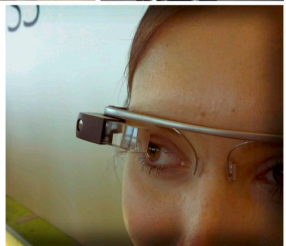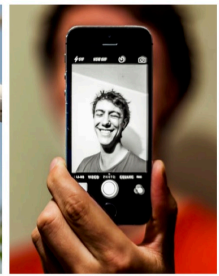[3]Song Han and William J. Dally (2018). "Bandwidth-efficient Deep Learning". In: *Proc. DAC*, 147:1–147:6.

# Application Category

| Both | Datacenter | Edge |
|------|-----------|------|
| Intel, Nvidia, IBM, Xilinx, HiSilicon, Google, Baidu, Alibaba Group, Cambricon, DeePhi, Bitmain, Wave Computing | AMD, Microsoft, Apple, Tencent Cloud,Aliyun, Baidu Cloud, HUAWEI Cloud, Fujitsu, Nokia, Facebook, HPE, Thinkforce, Cerebras, Graphcore, Groq, SambaNova Systems, Adapteva, PEZY | Qualcomm, Samsung, STMicroelectronics, NXP, MediaTek, Rockchip, Amazon_AWS, ARM, Synopsys, Imagination, CEVA, Cadence, VeriSilicon, Videantis, Horizon Robotics, Chipintelli, Unisound, AISpeech, Rokid, KnuEdge, Tenstorrent, ThinCI, Koniku, Knowm, Mythic, Kalray, BrainChip, AImotive, DeepScale, Leepmind, Krtkl, NovuMind, REM, TERADEEP, DEEP VISION, KAIST DNPU, Kneron, Esperanto Technologies, Gyrfalcon Technology, GreenWaves Technology, Lightelligence, Lightmatter, ThinkSilicon, Innogrit, Kortiq, Hailo,Tachyum |

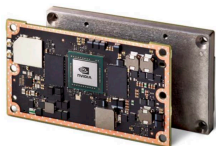Source: https://basicmi.github.io/Deep-Learning-Processor-List/

# Computing Spectrum



ShuffleNet    MobileNet    AlexNet    GoogleNet    VGG    ResNet 152    Inception-ResNet

10M FLOPS    100M FLOPS    1G FLOPS    10G FLOPS    100G FLOPS

CPU
(Raspberry Pi3)

GPU
(Jetson TX2)

FPGA
(UltraZed)

ASIC
(Movidius)

Flexibility

Power/Performance
Efficiency

# Convolution Basis

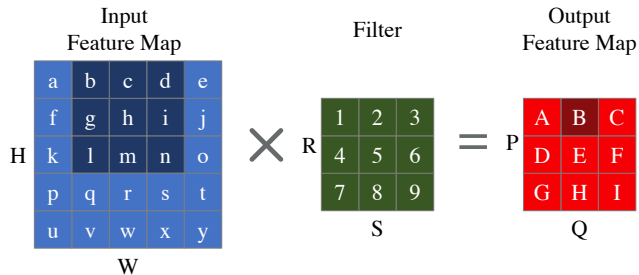| Input Feature Map | Filter | Output Feature Map |
|---|---|---|

**Input Feature Map**

| a | b | c | d | e |
|---|---|---|---|---|
| f | g | h | i | j |
| k | l | m | n | o |
| p | q | r | s | t |
| u | v | w | x | y |

H, W

**Filter**

| 1 | 2 | 3 |
|---|---|---|
| 4 | 5 | 6 |
| 7 | 8 | 9 |

R, S

**Output Feature Map**

| A | B | C |
|---|---|---|
| D | E | F |
| G | H | I |

P, Q

- H: Height of input feature map
- W: Width of input feature map
- R: Height of filter
- S: Width of filter
- P: Height of output feature map
- Q: Width of output feature map

$$A = a \cdot 1 + b \cdot 2 + c \cdot 3$$
$$+ f \cdot 4 + g \cdot 5 + h \cdot 6$$
$$+ k \cdot 7 + l \cdot 8 + m \cdot 9$$

Input Feature Map

Filter

Output Feature Map

| a | b | c | d | e |
|---|---|---|---|---|
| f | g | h | i | j |
| k | l | m | n | o |
| p | q | r | s | t |
| u | v | w | x | y |

H

W

R

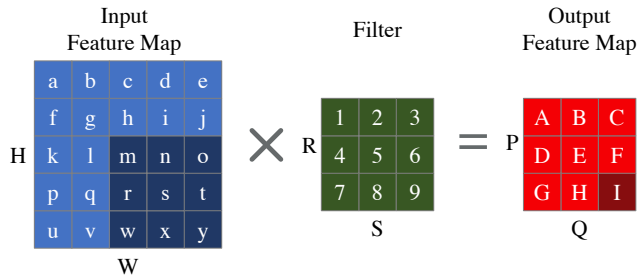| 1 | 2 | 3 |
|---|---|---|
| 4 | 5 | 6 |
| 7 | 8 | 9 |

S

P

| A | B | C |
|---|---|---|
| D | E | F |
| G | H | I |

Q

- H: Height of input feature map
- W: Width of input feature map
- R: Height of filter
- S: Width of filter
- P: Height of output feature map
- Q: Width of output feature map
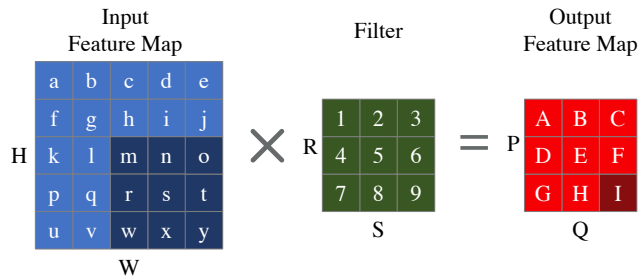- stride: # of rows/columns traversed per step

Input Feature Map × Filter = Output Feature Map

| | | | | |
|---|---|---|---|---|
| a | b | c | d | e |
| f | g | h | i | j |
| k | l | m | n | o |
| p | q | r | s | t |
| u | v | w | x | y |

Filter:

| | | |
|---|---|---|
| 1 | 2 | 3 |
| 4 | 5 | 6 |
| 7 | 8 | 9 |

Output Feature Map:

| | | |
|---|---|---|
| A | B | C |
| D | E | F |
| G | H | I |

- H: Height of input feature map
- W: Width of input feature map
- R: Height of filter
- S: Width of filter
- P: Height of output feature map
- Q: Width of output feature map
- stride: # of rows/columns traversed per step

**Input Feature Map**
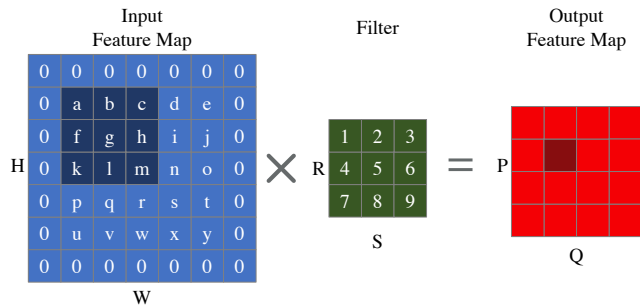
| a | b | c | d | e |
|---|---|---|---|---|
| f | g | h | i | j |
| k | l | m | n | o |
| p | q | r | s | t |
| u | v | w | x | y |

H

W

**Filter**

| 1 | 2 | 3 |
|---|---|---|
| 4 | 5 | 6 |
| 7 | 8 | 9 |

R

S

**Output Feature Map**

| A | B | C |
|---|---|---|
| D | E | F |
| G | H | I |

P

Q

- H: Height of input feature map
- W: Width of input feature map
- R: Height of filter
- S: Width of filter
- P: Height of output feature map
- Q: Width of output feature map
- stride: # of rows/columns traversed per step

Input Feature Map
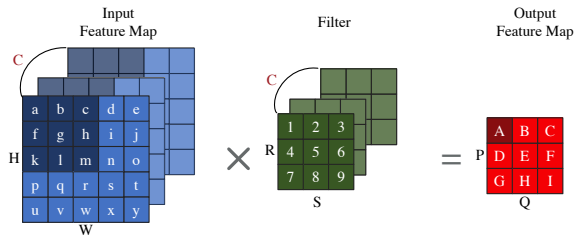
Filter

Output Feature Map

- H: Height of input feature map
- W: Width of input feature map
- R: Height of filter
- S: Width of filter
- P: Height of output feature map
- Q: Width of output feature map
- stride: # of rows/columns traversed per step

$$P = \frac{(H - R)}{\text{stride}} + 1;$$
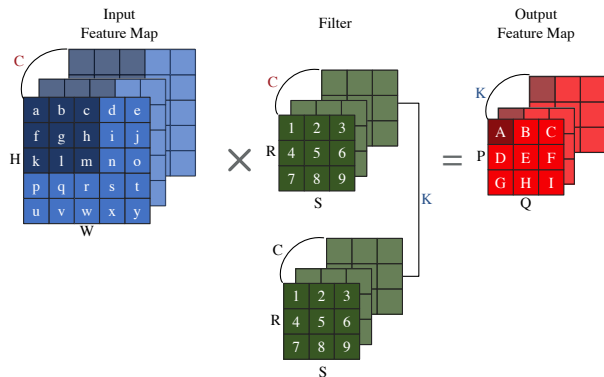$$Q = \frac{(W - S)}{\text{stride}} + 1.$$

Input Feature Map

Filter

Output Feature Map

- H: Height of input feature map
- W: Width of input feature map
- R: Height of filter
- S: Width of filter
- P: Height of output feature map
- Q: Width of output feature map
- stride: # of rows/columns traversed per step
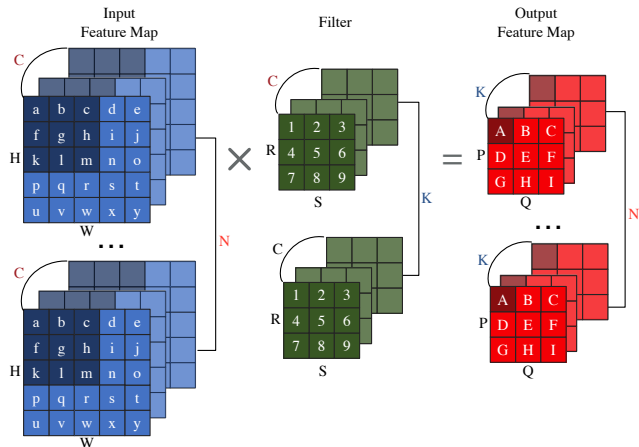- padding: # of zero rows/columns added

$$P = \frac{(H - R + 2 \cdot \text{pad})}{\text{stride}} + 1;$$
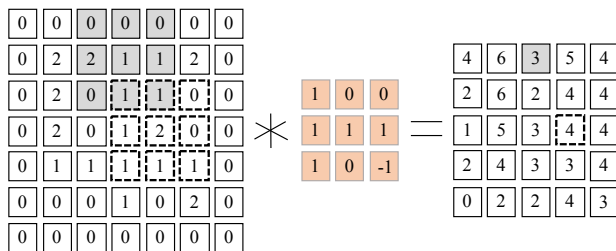$$Q = \frac{(W - S + 2 \cdot \text{pad})}{\text{stride}} + 1.$$

- H: Height of input feature map
- W: Width of input feature map
- R: Height of filter
- S: Width of filter
- P: Height of output feature map
- Q: Width of output feature map
- stride: # of rows/columns traversed per step
- padding: # of zero rows/columns added
- C: # of input channels

# 3D-Convolution



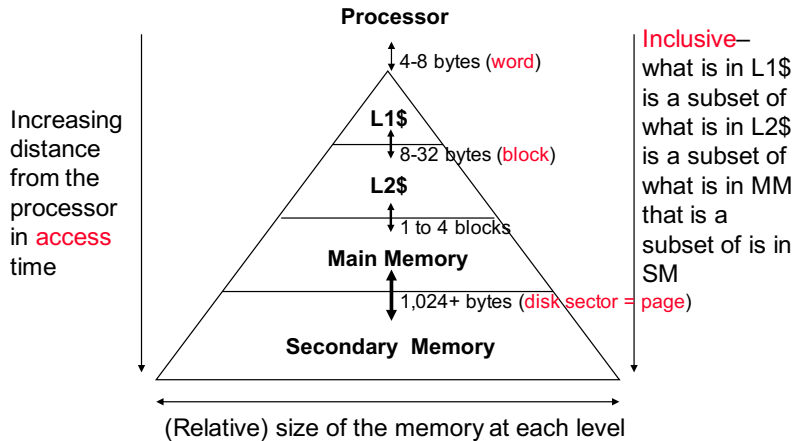- H: Height of input feature map
- W: Width of input feature map
- R: Height of filter
- S: Width of filter
- P: Height of output feature map
- Q: Width of output feature map
- stride: # of rows/columns traversed per step
- padding: # of zero rows/columns added
- C: # of input channels
- K: # of output channels

Input Feature Map

Filter

Output Feature Map

- H: Height of input feature map
- W: Width of input feature map
- R: Height of filter
- S: Width of filter
- P: Height of output feature map
- Q: Width of output feature map
- stride: # of rows/columns traversed per step
- padding: # of zero rows/columns added
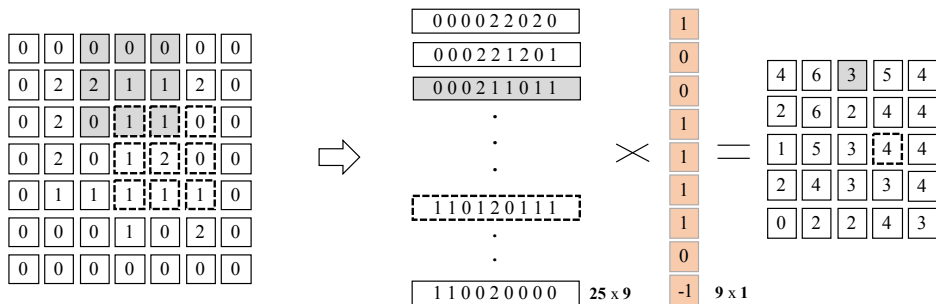- C: # of input channels
- K: # of output channels
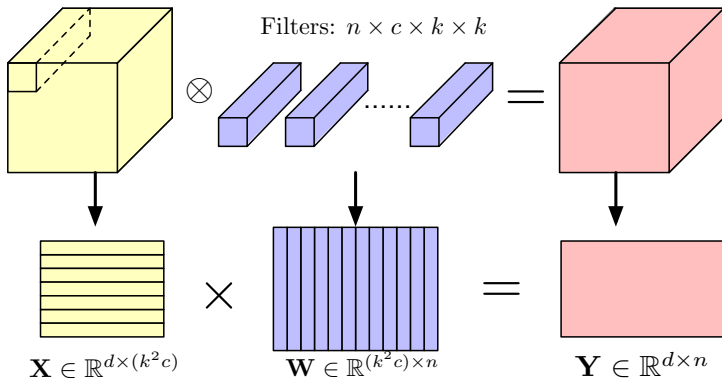- N: Batch size

Direct convolution: No extra memory overhead

- Low performance
- Poor memory access pattern due to geometry-specific constraint
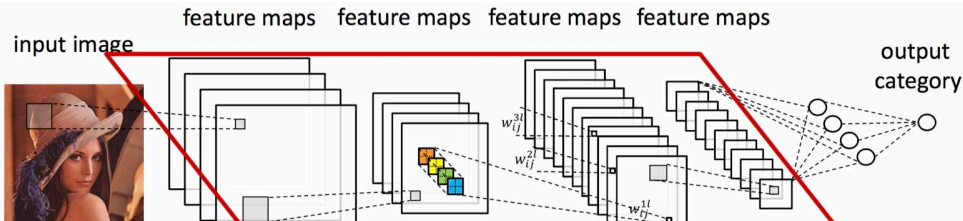- Relatively short dot product

Processor

4-8 bytes (word)

**L1$**

8-32 bytes (block)

**L2$**

1 to 4 blocks

**Main Memory**

1,024+ bytes (disk sector = page)

**Secondary  Memory**

Increasing distance from the processor in access time

(Relative) size of the memory at each level

Inclusive– what is in L1$ is a subset of what is in L2$ is a subset of what is in MM that is a subset of is in SM

- Spatial locality
- Temporal Locality

- Large extra memory overhead
- Good performance
- BLAS-friendly memory layout to enjoy SIMD/locality/parallelism
- Applicable for any convolution configuration on any platform

Filters: $n \times c \times k \times k$

$\mathbf{X} \in \mathbb{R}^{d \times (k^2 c)}$     $\mathbf{W} \in \mathbb{R}^{(k^2 c) \times n}$     $\mathbf{Y} \in \mathbb{R}^{d \times n}$

- Transform convolution to matrix multiplication
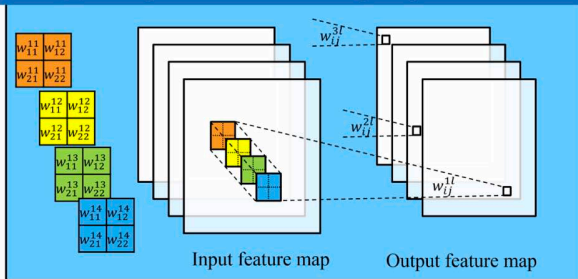- Unified calculation for both convolution and fully-connected layers
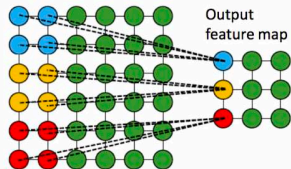
# Algorithm Design Level

feature maps    feature maps    feature maps    feature maps

input image

output category

$w_{ij}^{3l}$

$w_{ij}^{2l}$

$w_{ij}^{1l}$
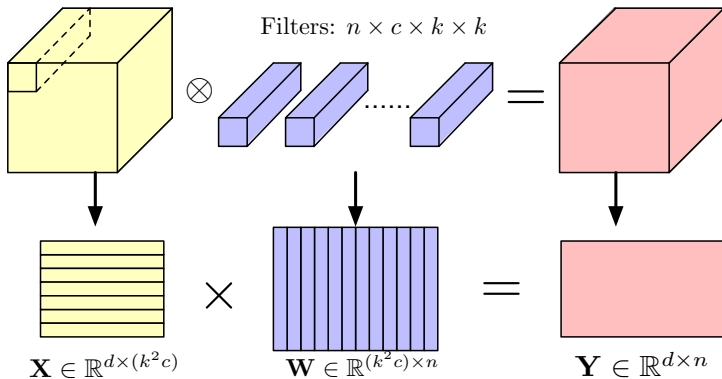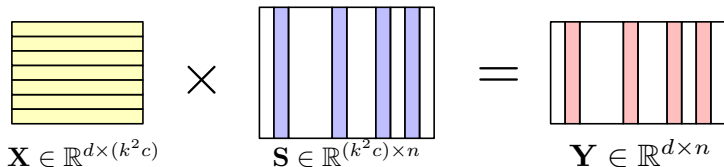
**Convolutional layers account for over 90% computation**

[1] A. Krizhevsky, etc. Imagenet classification with deep convolutional neural networks. NIPS 2012.
[2] J. Cong and B. Xiao. Minimizing computation in convolutional neural networks. ICANN 2014

$w_{11}^{11}$ $w_{12}^{11}$
$w_{21}^{11}$ $w_{22}^{11}$

$w_{11}^{12}$ $w_{12}^{12}$
$w_{21}^{12}$ $w_{22}^{12}$

$w_{11}^{13}$ $w_{12}^{13}$
$w_{21}^{13}$ $w_{22}^{13}$

$w_{11}^{14}$ $w_{12}^{14}$
$w_{21}^{14}$ $w_{22}^{14}$

$w_{ij}^{3l}$

$w_{ij}^{2l}$

$w_{ij}^{1l}$

Input feature map          Output feature map

Max-pooling is optional

Input feature map

Output feature map

Filters: $n \times c \times k \times k$

$$\mathbf{X} \in \mathbb{R}^{d \times (k^2 c)} \quad \times \quad \mathbf{W} \in \mathbb{R}^{(k^2 c) \times n} \quad = \quad \mathbf{Y} \in \mathbb{R}^{d \times n}$$

- Transform convolution to matrix multiplication
- Unified calculation for both convolution and fully-connected layers

$$\mathbf{X} \in \mathbb{R}^{d \times (k^2 c)} \qquad \mathbf{S} \in \mathbb{R}^{(k^2 c) \times n} \qquad \mathbf{Y} \in \mathbb{R}^{d \times n}$$

## Sparse DNN

- *Sparsification*: weight pruning;

- *Compression*: compressed sparse format for storage;

- *Potential acceleration*: sparse matrix multiplication algorithm.

[4]Wei Wen et al. (2016). "Learning structured sparsity in deep neural networks". In: *Proc. NIPS*, pp. 2074–2082.

[5]Yihui He, Xiangyu Zhang, and Jian Sun (2017). "Channel Pruning for Accelerating Very Deep Neural Networks". In: *Proc. ICCV*.

$$\mathbf{X} \in \mathbb{R}^{d \times (k^2 c)} \quad \mathbf{U} \in \mathbb{R}^{(k^2 c) \times r} \quad \mathbf{V} \in \mathbb{R}^{1^2 r \times n} \quad \mathbf{Y} \in \mathbb{R}^{d \times n}$$

## Low-rank DNN

- *Low-rank approximation*: matrix decomposition or tensor decomposition.
- *Compression and acceleration*: less storage required and less FLOP in computation.

---

[6]Xiangyu Zhang, Jianhua Zou, et al. (2015). "Efficient and accurate approximations of nonlinear convolutional networks". In: *Proc. CVPR*, pp. 1984–1992.

[7]Xiyu Yu et al. (2017). "On compressing deep models by low rank and sparse decomposition". In: *Proc. CVPR*, pp. 7370–7379.

ReLU

- Activation unit: `ReLU`
- Error more sensitive to positive response;
- Enlarge the solution space.

$$\min_{\boldsymbol{W}} \sum_{i=1}^{N} \|\boldsymbol{W}\boldsymbol{X}_i - \boldsymbol{Y}_i\|_F \rightarrow \min_{\boldsymbol{W}} \sum_{i=1}^{N} \|r(\boldsymbol{W}\boldsymbol{X}_i) - \boldsymbol{Y}_i\|_F$$

- $\boldsymbol{X}$: input feature map
- $\boldsymbol{Y}$: output feature map

[8]Xiangyu Zhang, Jianhua Zou, et al. (2015). "Efficient and accurate approximations of nonlinear convolutional networks". In: *Proc. CVPR*, pp. 1984–1992.

- Simultaneous low-rank approximation and network sparsification;
- Non-linearity is taken into account.
- Acceleration is achieved with structured sparsity.

---

[9]Yuzhe Ma et al. (2019). "A Unified Approximation Framework for Non-Linear Deep Neural Networks". In: *Proc. ICTAI*.

Tensor Reconstruction Module (TRM). The pipeline of TRM consists of two main steps, sub-attention map generation and global context reconstruction. The processing from top to bottom (see ↓) indicates the sub-attention map generation from three dimensions (channel / height / width). The processing from left to right (see $A_1 + A_2 + \cdots + A_r = A$) denotes the global context reconstruction from low-rank to high-rank.

[10]Wanli Chen et al. (2020). "Tensor low-rank reconstruction for semantic segmentation". In: *Proc. ECCV*, pp. 52–69.

# Compilation Level

Left column:

```
 1  #include <stdio.h>
 2  #include <memory.h>
 3  #include <time.h>
 4  #include <stdlib.h>
 5  #include <sys/time.h>
 6
 7  double a_arr[1024][1024];
 8  double b_arr[1024][1024];
 9
10  int main()
11  {
12      int N = 1024;
13      for(int i=0;i<1024;i++) for(int j=0;j<1024;j++)
14      {
15          a_arr[i][j] = 1; b_arr[i][j] = 2;
16      }
17      double sum;
18
19      struct timeval startTime,endTime;
20      float Timeuse;
21      gettimeofday(&startTime,NULL);
22
23      // ========================================
24      //    following are the key operations
25      for(int i=0; i<1024; i++){
26          for(int j=0; j<1024; j++){
27              sum += a_arr[j][i] * b_arr[j][i];
28          }
29      }
30      // ========================================
31
32      gettimeofday(&endTime,NULL);
33      Timeuse = 1000000 * (endTime.tv_sec-startTime.tv_sec) + (endTime.tv_usec-startTime.tv_us
34      printf("========== Cache Optimization Demo ========== \n");
35      printf("total timeuse = %.2f us \n",Timeuse);
36
37      return 0;
38  }
39
```
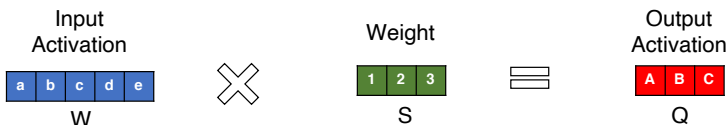
Right column:

```
 1  #include <stdio.h>
 2  #include <memory.h>
 3  #include <time.h>
 4  #include <stdlib.h>
 5  #include <sys/time.h>
 6
 7  double a_arr[1024][1024];
 8  double b_arr[1024][1024];
 9
10  int main()
11  {
12      int N = 1024;
13      for(int i=0;i<1024;i++) for(int j=0;j<1024;j++)
14      {
15          a_arr[i][j] = 1; b_arr[i][j] = 2;
16      }
17      double sum;
18
19      struct timeval startTime,endTime;
20      float Timeuse;
21      gettimeofday(&startTime,NULL);
22
23      // ========================================
24      //    following are the key operations
25      for(int i=0; i<1024; i++){
26          for(int j=0; j<1024; j++){
27              sum += a_arr[i][j] * b_arr[i][j];
28          }
29      }
30      // ========================================
31
32      gettimeofday(&endTime,NULL);
33      Timeuse = 1000000 * (endTime.tv_sec-startTime.tv_sec) + (endTime.tv_usec-startTime.tv_u
34      printf("========== Cache Optimization Demo ========== \n");
35      printf("total timeuse = %.2f us \n",Timeuse);
36
37      return 0;
38  }
39
```

Same complexity; same real runtime?

Input Activation

| a | b | c | d | e |

W

Weight

| 1 | 2 | 3 |

S

Output Activation
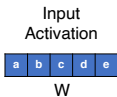
| A | B | C |

Q

```
for(q=0; q<Q; q++){
  for (s=0; s<S; s++){
    OA[q] += IA[q+s] * W[s];
  }
}
```
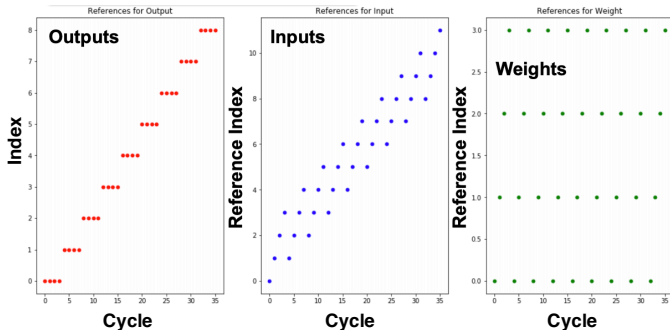
**Output Stationary (OS)**
**Dataflow**

```
for (s=0; s<S; s++){
  for(q=0; q<Q; q++){
    OA[q] += IA[q+s] * W[s];
  }
}
```
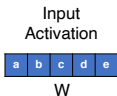
**Weight Stationary (WS)**
**Dataflow**

```
for(q=0; q<Q; q++){ // Q =9
  for (s=0; s<S; s++){ // S=4
    OA[q] += IA[q+s] * W[s];
  }
}
```

**Input Activation**

| a | b | c | d | e |
|---|---|---|---|---|

W

**Weight**

| 1 | 2 | 3 |
|---|---|---|

S

**Output Activation**

| A | B | C |
|---|---|---|

Q

```
for (s=0; s<S; s++){// S=4
  for(q=0; q<Q; q++){// Q =9
    OA[q] += IA[q+s] * W[s];
  }
```



References for Output — Index / Cycle

References for Input — Index Value - Input[] / Cycle

References for Weight — Index Value - Weight[] / Cycle

```
1    for (n=0; n<N; n++) {
2    for (k=0; k<K; k++) {
3    for (p=0; p<P; p++) {
4    for (q=0; q<Q; q++) {
5        OA[n][k][p][q]= 0;
6        for (r=0; r<R; r++) {
7        for (s=0; s<S; s++) {
8        for (c=0; c<C; c++) {
9            h = p * stride - pad + r;
10           w = q * stride - pad + s;
11           OA[n][k][p][q] += IA[n][c][h][w] * W[k][c][r][s];
12   } } } } } } }
```
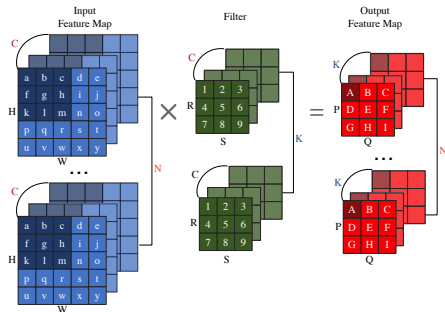
```
1     for (n=0; n<N; n++) {
2     for (r=0; r<R; r++) {
3     for (s=0; s<S; s++) {
4     for (c=0; c<C; c++) {
5     for (k=0; k<K; k++) {
6         float curr_w = W[r][s][c][k];
7         for (p=0; p<P; p++) {
8         for (q=0; q<Q; q++) {
9             h = p * stride - pad + r;
10            w = q * stride - pad + s;
11            OA[n][k][p][q] += IA[n][c][h][w] * curr_w;
12    } } } } } } }
```

```
1    for (n=0; n<N; n++) {
2    for (r=0; r<R; r++) {
3    for (s=0; s<S; s++) {
4        spatial_for (c=0; c<C; c++) {
5        spatial_for (k=0; k<K; k++) {
6        float curr_w = W[r][s][c][k];
7        for (p=0; p<P; p++) {
8        for (q=0; q<Q; q++) {
9            h = p * stride - pad + r;
10           w = q * stride - pad + s;
11           OA[n][k][p][q] += IA[n][c][h][w] * curr_w;
12   } } } } } } }
```
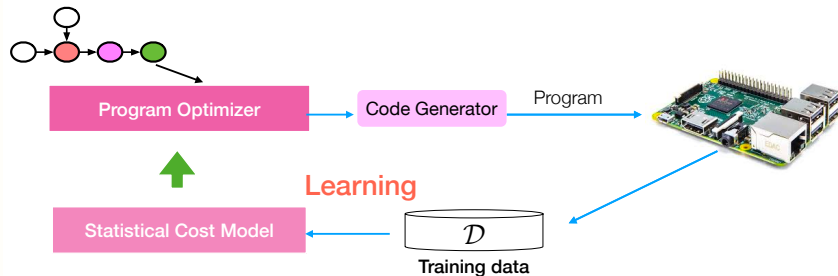
```
1        for (n=0; n<N; n++) {
2        for (r=0; r<R; r++) {
3        for (s=0; s<S; s++) {
4        for (c_t=0; c_t<C/16; c_t++) {
5        for (k_t=0; k_t<K/64; k_t++) {
6        spatial_for (c_s=0; c_s<16; c_s++) {
7        spatial_for (k_s=0; k_s<64; k_s++) {
8            int curr_c = c_t * 16 + c_s;
9            int curr_k = k_t * 64 + k_s;
10           float curr_w = W[r][s][curr_c][curr_k];
11           for (p=0; p<P; p++) for (q=0; q<Q; q++) {
12               h = p * stride - pad + r; w = q * stride - pad + s;
13               OA[n][curr_k][p][q] += IA[n][curr_c][h][w] * curr_w;
14       } } } } } }
```

```
1      for (n=0; n<N; n++) {
2      for (r=0; r<R; r++) {
3      for (s=0; s<S; s++) {
4      for (c_t=0; c_t<C/16; c_t++) {
5      for (k_t=0; k_t<K/64; k_t++) {
6      spatial_for (c_s=0; c_s<16; c_s++) {
7      spatial_for (k_s=0; k_s<64; k_s++) {
8          int curr_c = c_t * 16 + c_s;
9          int curr_k = k_t * 64 + k_s;
10         float curr_w = W[r][s][curr_c][curr_k];
11         for (p=0; p<P; p++) for (q=0; q<Q; q++) {
12             h = p * stride - pad + r; w = q * stride - pad + s;
13             OA[n][curr_k][p][q] += IA[n][curr_c][h][w] * curr_w;
14     } } } } } }
```

### Questions:

- How many configurations we have?

- How to search for the BEST configuration?

- How to handle different backend devices?

## Layer-wise Optimization: Autotuning



Tuning algorithms:

- Active learning.
- Transfer learning.
- Reinforcement learning.

- Batch transductive experimental design
- Bootstrap-guided adaptive optimization



---
[12]Qi Sun, Chen Bai, Hao Geng, et al. (2021). "Deep neural network hardware deployment optimization via advanced active learning". In: *Proc. DATE*, pp. 1510–1515.
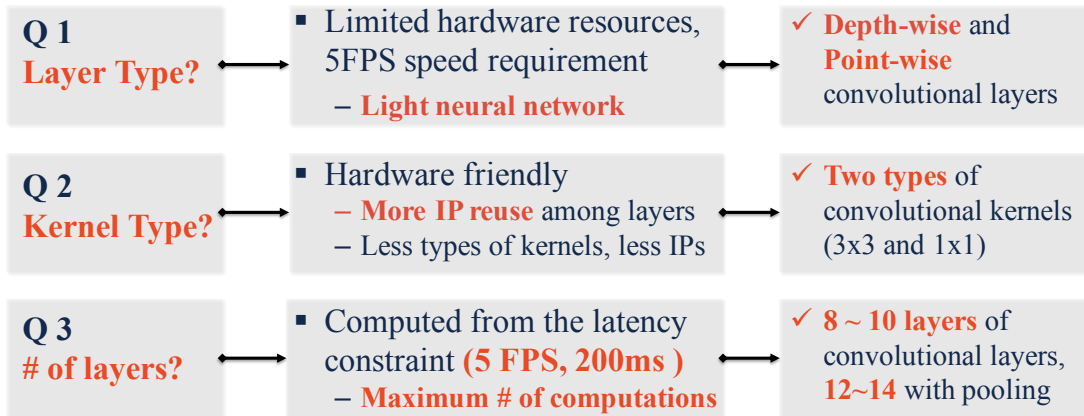
1. **preparation**: learn a deep Gaussian process model from historical data

2. **transfer**: transfer knowledge of the DGP model to new tasks

3. **optimal searching**: guide the optimization of new tasks with the tuned DGP model



---

[13]Qi Sun, Chen Bai, Tinghuan Chen, et al. (2021). "Fast and Efficient DNN Deployment via Deep Gaussian Transfer Learning". In: *Proc. ICCV*.

# Hardware Implementation

| | | |
|---|---|---|
| **Q 1**<br>**Layer Type?** | ■ Limited hardware resources, 5FPS speed requirement<br>– **Light neural network** | ✓ **Depth-wise** and **Point-wise** convolutional layers |
| **Q 2**<br>**Kernel Type?** | ■ Hardware friendly<br>– **More IP reuse** among layers<br>– Less types of kernels, less IPs | ✓ **Two types** of convolutional kernels (3x3 and 1x1) |
| **Q 3**<br>**# of layers?** | ■ Computed from the latency constraint **(5 FPS, 200ms )**<br>– **Maximum # of computations** | ✓ **8 ~ 10 layers** of convolutional layers, **12~14** with pooling |

1. (optional) DNN Design in C/C++ (`example`)
2. Generate RTL (`example`) by tool Vivado HLS



3. Generate bitstream by tool Vivado IDE
4. Load bitstream to FPGA board

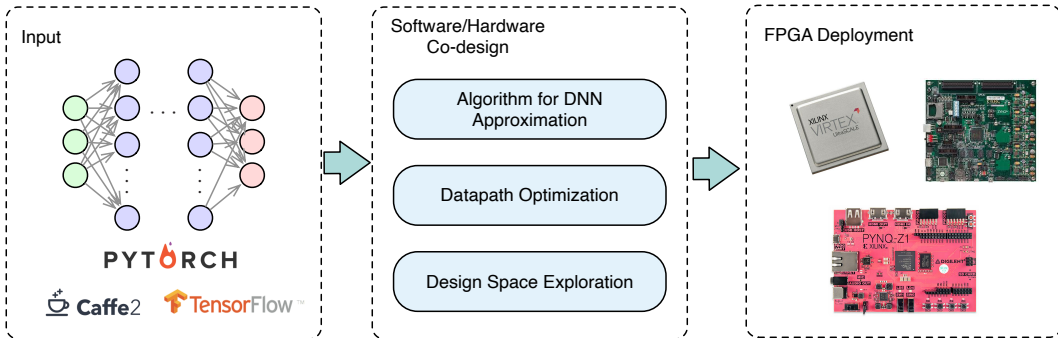(a)                    (b)                    (c)

(a) 540p $\rightarrow$ 1080p SR running on KV260; (b) Physical implementation on KV260; (c) Physical implementation on ZCU104.

- Comparison with SOTA FPGA solutions and commercial GPU and ASIC
- 540p→1080p SR is based on KV260
- 720p→1440p is based on ZCU104

| Accelerator | 540p→1080p | | | 720p→1440p | | |
|---|---|---|---|---|---|---|
| | Latency (ms) | FPS | Price ($) | Latency (ms) | FPS | Price ($) |
| NX GPU | 35.36 | 28.28 | 399 | 61.83 | 16.17 | **399** |
| Ascend 310 | 48.55 | 20.60 | 999 | 87.57 | 11.42 | 999 |
| Xilinx DPU | >37.31 | < 26.80 | 199 | >56.14 | < 17.81 | 1554 |
| DNNBuilder | >186.57 | < 5.36 | 199 | >82.92 | < 12.06 | 1554 |
| Ours | **27.94** | **35.79** | **199** | **39.74** | **25.16** | 1554 |

Input

Software/Hardware
Co-design

Algorithm for DNN
Approximation

Datapath Optimization

Design Space Exploration

FPGA Deployment

THANK YOU!