

# 딥러닝 몸풀기

누구나 이해할 수 있는 딥러닝

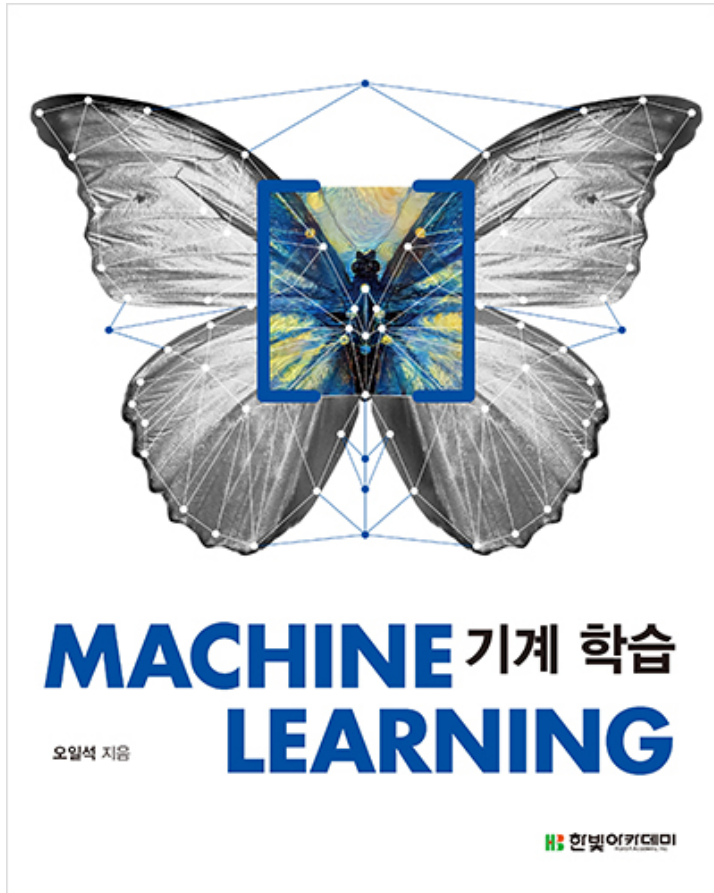
Local Laboratory

# 딥러닝 몸풀기

Local Laboratory

1. 인공지능과 기계학습 그리고 딥러닝
2. 다층 퍼셉트론
3. 딥러닝의 기초

# Reference

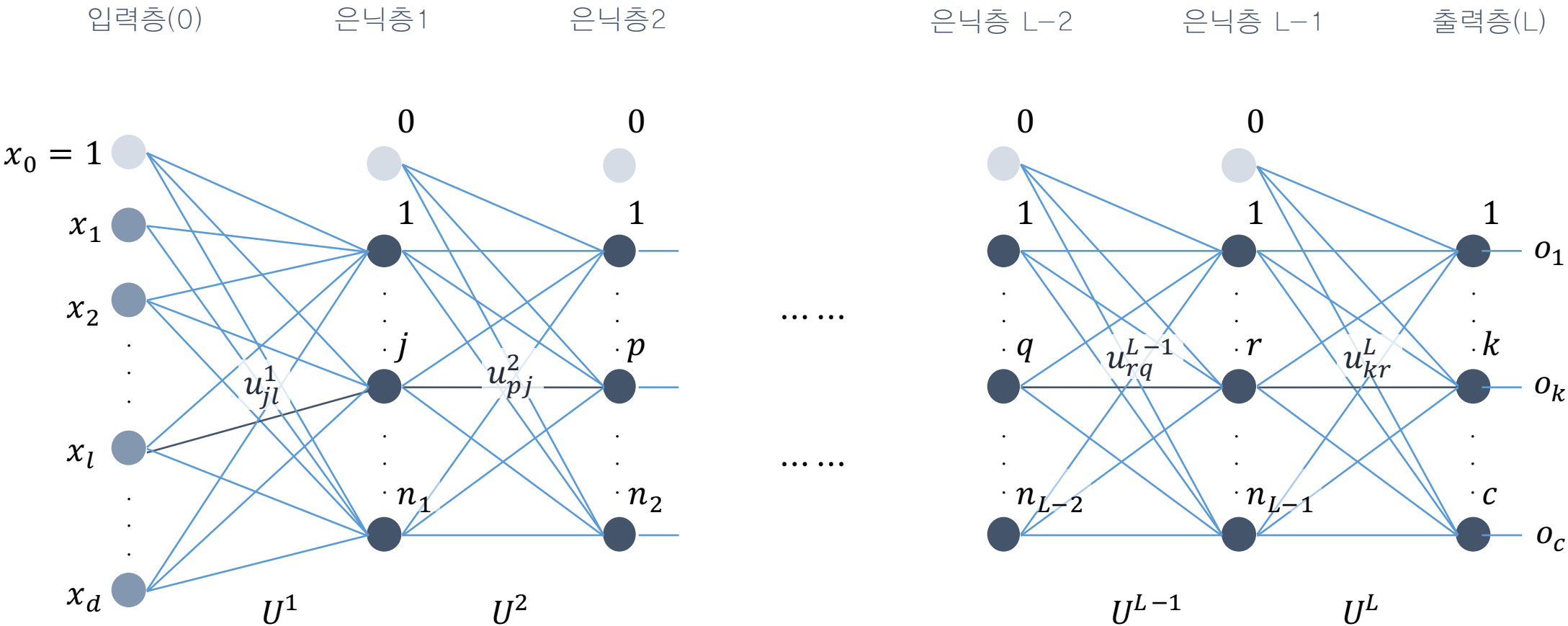


기계학습, 오일석



Self-Driving Car Nano degree, Udacity

# 딥러닝 Deep MLP / Deep Neural Network



입력층 d+1개, 출력층 c개, L-1개의 은닉층  
l번째 은닉층의 노드 수를  $n_l$ 로 표기

$$U^l = \begin{pmatrix} \begin{bmatrix} u_{10}^l & u_{11}^l & \cdots & u_{1n_{l-1}}^l \\ u_{20}^l & u_{21}^l & & u_{2n_{l-1}}^l \\ \vdots & & \ddots & \vdots \\ u_{n_l 0}^l & u_{n_l 1}^l & \cdots & u_{n_l n_{l-1}}^l \end{bmatrix} \end{pmatrix}, \quad l = 1, 2, 3, \dots, L$$

DMLP의 가중치 행렬

$l-1$ 번째 층과  $l$ 번째 층을 연결하는 가중치는 총  $(n_{l-1} + 1)n_l$

$$\mathbf{o} = \mathbf{f}(\mathbf{x}) = \mathbf{f}_L \left( \cdots \mathbf{f}_2 \left( \mathbf{f}_1(\mathbf{x}) \right) \right)$$

입력층 특징  $\mathbf{x}$ 를 내부 표현으로 바꿔쓰면,

$$\mathbf{z}^0 = (z_0, z_1, z_2, \dots, z_{n_0})^T = (1, x_1, x_2, \dots, x_d)^T$$

$l$ 번째 층의  $j$ 번째 노드가 수행하는 연산

$$z_j^l = \tau_l(s_j^l) \quad s_j^l = \mathbf{u}_j^l \mathbf{z}^{l-1}$$

$$\mathbf{z}^{l-1} = (1, z_1^{l-1}, z_2^{l-1}, \dots, z_{n_{l-1}}^{l-1})^T, \quad \mathbf{u}_j^l = (u_{j0}^l, u_{j1}^l, \dots, u_{jn_{l-1}}^l)$$

행렬표기를 이용하여  $l$ 번째 층의 연산 전체를 쓰면,

$$\mathbf{z}^l = \tau_l(\mathbf{U}^l \mathbf{z}^{l-1}), \quad 1 \leq l \leq L$$

$L$ 번째 층(출력층)의 그레이디언트 계산

$$\delta_k^L = \tau'_L(s_k^L)(y_k - o_k), \quad 1 \leq k \leq c$$

$$\frac{\partial J}{\partial u_{kr}^L} = -\delta_k^L z_r^{L-1}, \quad 0 \leq r \leq n_{L-1}, 1 \leq k \leq c$$

$l + 1$ 번째 층의 정보를 이용하여  $l$ 번째 층의 그레이디언트 계산( $l = L - 1, L - 2, \dots, 1$ )

$$\delta_j^l = \tau'_L(s_j^l) \sum_{p=1}^{n_{l+1}} \delta_p^{l+1} u_{pj}^{l+1}, \quad 1 \leq j \leq n_l$$

$$\frac{\partial J}{\partial u_{ji}^l} = -\delta_j^l z_i^{l-1}, \quad 0 \leq i \leq n_{l-1}, 1 \leq j \leq n_l$$

## 미니배치 스토캐스틱 경사 하강법

### 알고리즘 4-1 DMLP를 위한 미니배치 스토캐스틱 경사 하강법

입력: 훈련집합  $\mathbb{X}$ 와  $\mathbb{Y}$ , 학습률  $\rho$ , 미니배치 크기  $t$

출력: 가중치 행렬  $\mathbf{U}^l, l = 1, 2, \dots, L$

```

1   $\mathbf{U}^l, l = 1, 2, \dots, L$ 을 초기화한다.
2  repeat
3     $\mathbb{X}$ 와  $\mathbb{Y}$ 에서  $t$ 개의 샘플을 무작위로 뽑아 미니배치  $\mathbb{X}'$ 와  $\mathbb{Y}'$ 를 만든다.
4    for ( $l=1$  to  $L$ )  $\Delta \mathbf{U}^l = \mathbf{0}$ 
5      for ( $\mathbb{X}'$ 의 샘플 각각에 대해)
6        현재 처리하는 샘플을  $\mathbf{x} = (x_0, x_1, x_2, \dots, x_d)^T, \mathbf{y} = (y_1, y_2, \dots, y_c)^T$ 라 표기한다.
7         $x_0, z_0^1, z_0^2, \dots, z_0^L$ 을 1로 설정한다.
8        // 전방 계산
9         $\mathbf{x}$ 를  $\mathbf{z}^0$ 에 대입한다. // 식 (4.3)
10       for ( $l=1$  to  $L$ ) // 왼쪽 층에서 오른쪽 층으로 진행하면서 전방 계산
11         for ( $j=1$  to  $n_l$ ) // 각 노드에 대해
12            $s_j^l = \mathbf{u}_j^l \mathbf{z}^{l-1}$  // 식 (4.4)
13            $z_j^l = \tau_l(s_j^l)$  // 식 (4.4)
14           // 오류 역전파의 단계 1: 그레이디언트 계산
15         for ( $k=1$  to  $c$ )  $\delta_k^L = \tau'_L(s_k^L)(y_k - o_k)$  // 식 (4.6)
16         for ( $k=1$  to  $c$ ) for ( $r=0$  to  $n_{l+1}$ )  $\Delta u_{kr}^L = \Delta u_{kr}^L + (-\delta_k^L z_r^{L-1})$ 
17         for ( $l=L-1$  to 1) // 오른쪽 층에서 왼쪽 층으로 진행하면서 오류 역전파
18           for ( $j=1$  to  $n_l$ )  $\delta_j^l = \tau'_l(s_j^l) \sum_{p=1}^{n_{l+1}} \delta_p^{l+1} u_{pj}^{l+1}$  // 식 (4.8)
19           for ( $j=1$  to  $n_l$ ) for ( $i=0$  to  $n_{l-1}$ )  $\Delta u_{ji}^l = \Delta u_{ji}^l + (-\delta_j^l z_i^{l-1})$ 
20           // 오류 역전파의 단계 2: 가중치 갱신
21         for ( $l=L$  to 1)
22           for ( $j=1$  to  $n_l$ ) for ( $i=0$  to  $n_{l-1}$ )  $u_{ji}^l = u_{ji}^l - \rho \left(\frac{1}{t}\right) \Delta u_{ji}^l$ 
23 until (멈춤 조건)
    
```



퍼셉트론  
Perception

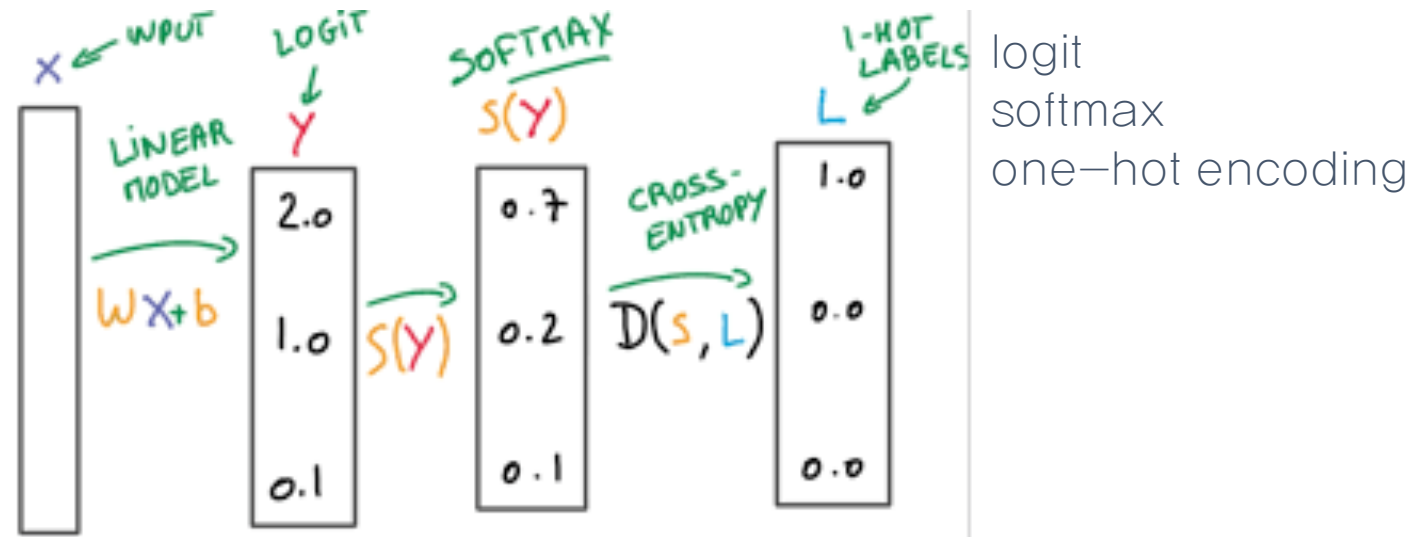
다층 퍼셉트론  
MLP

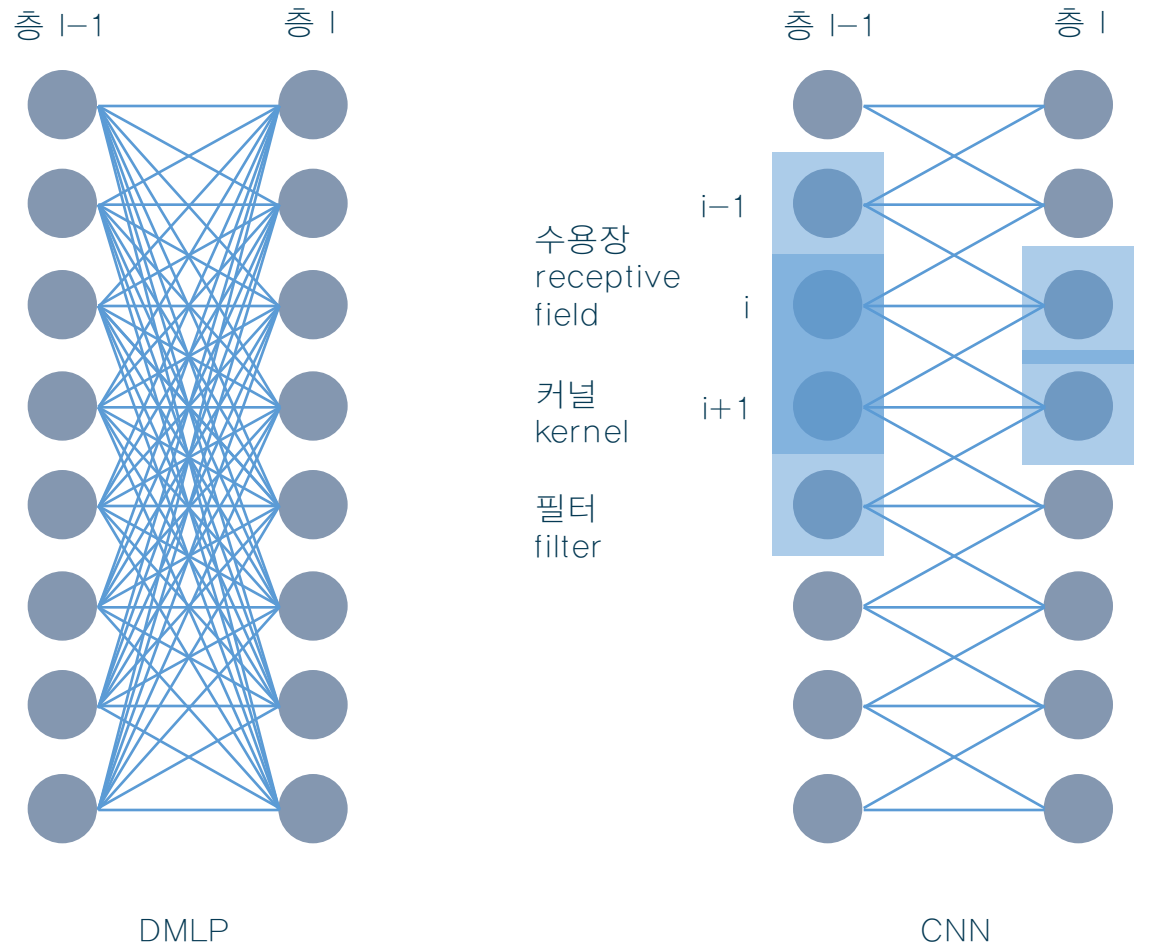
깊은 다층 퍼셉트론  
DMLP or DNN

활성함수: 계단함수  
비용함수: 평균제곱 오차

시그모이드함수  
평균제곱 오차

ReLU와 변형들  
교차 엔트로피 또는 로그우도



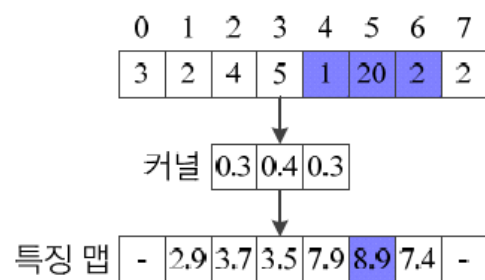


## DMLP와 CNN의 비교

- DMLP
  - 완전연결 구조로 높은 복잡도
  - 학습이 매우 느리고 과잉적합 우려
- CNN
  - 컨볼루션 연산을 이용한 부분연결 (희소 연결) 구조로 복잡도 크게 낮춤
  - 컨볼루션 연산은 좋은 특징 추출

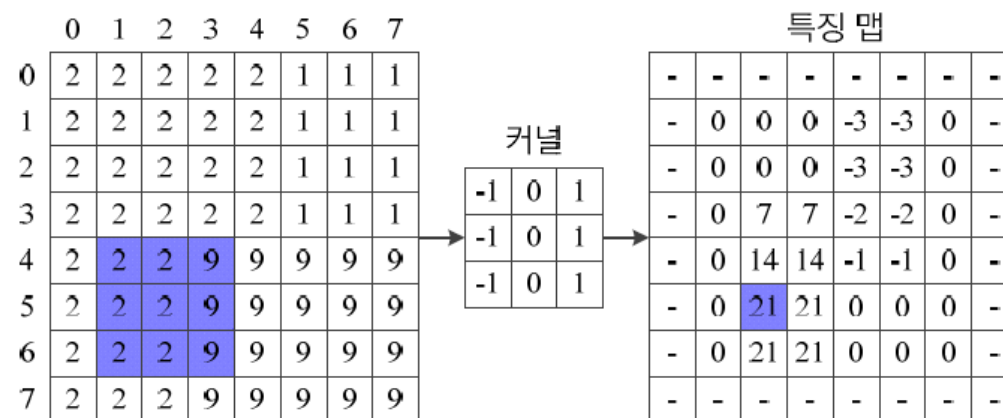
## CNN

- 격자 구조(영상, 음성 등)를 갖는 데이터에 적합
- 수용장은 *receptive field* 인간시각과 유사
- 가변 크기의 입력 처리 가능

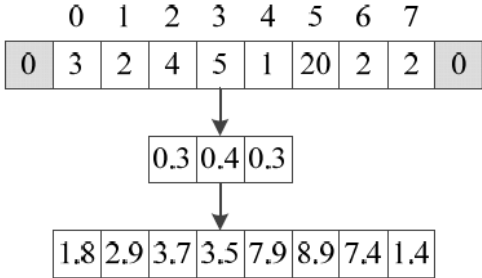
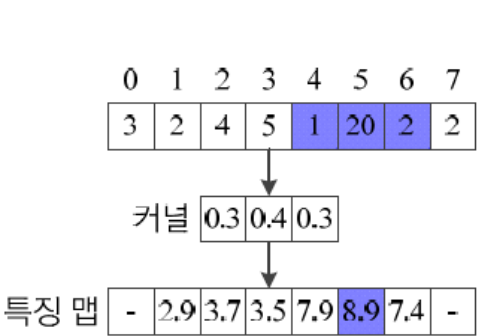


특징맵 feature map

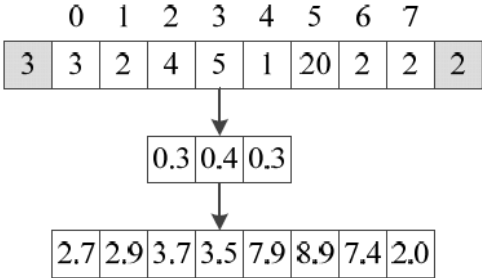
$$s(i) = z \circledast u = \sum_{x=(h-1)/2}^{(h-1)/2} z(i+x)u(x)$$



$$s(j, i) = z \circledast u = \sum_{y=(h-1)/2}^{(h-1)/2} \sum_{x=(h-1)/2}^{(h-1)/2} z(j+y, i+x)u(y, x)$$

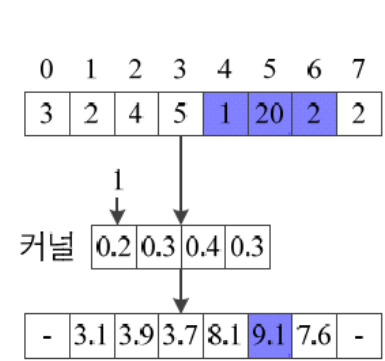


(a) 0 덧대기



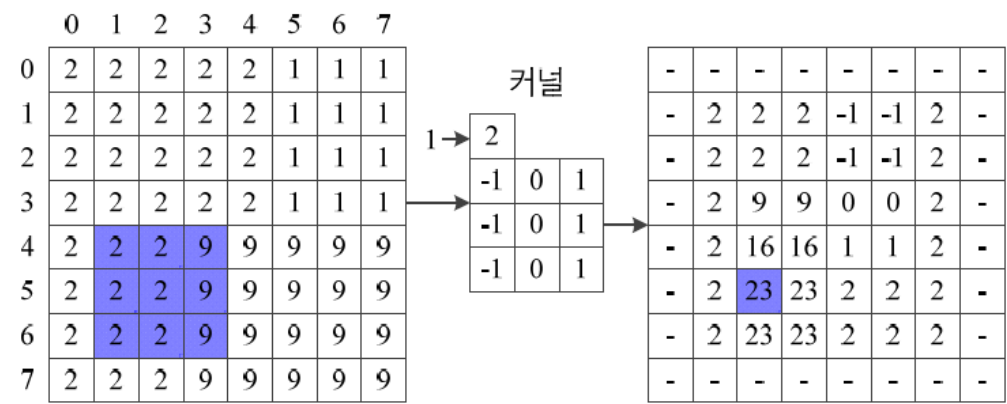
(b) 복사 덧대기

그림 4-7 덧대기(회색 노드가 덧댄 노드)



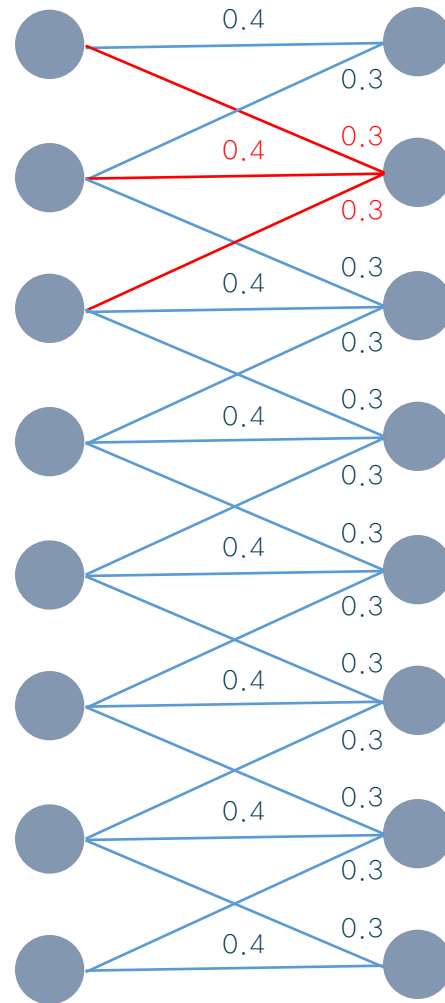
(a) 1차원 컨볼루션

그림 4-8 바이어스



(b) 2차원 컨볼루션

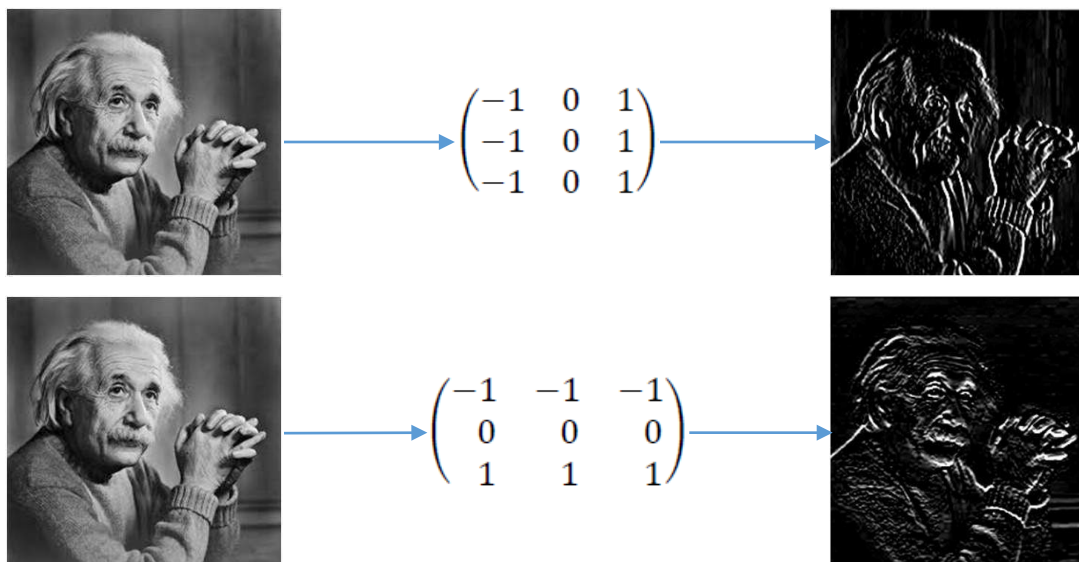
## 가중치 공유



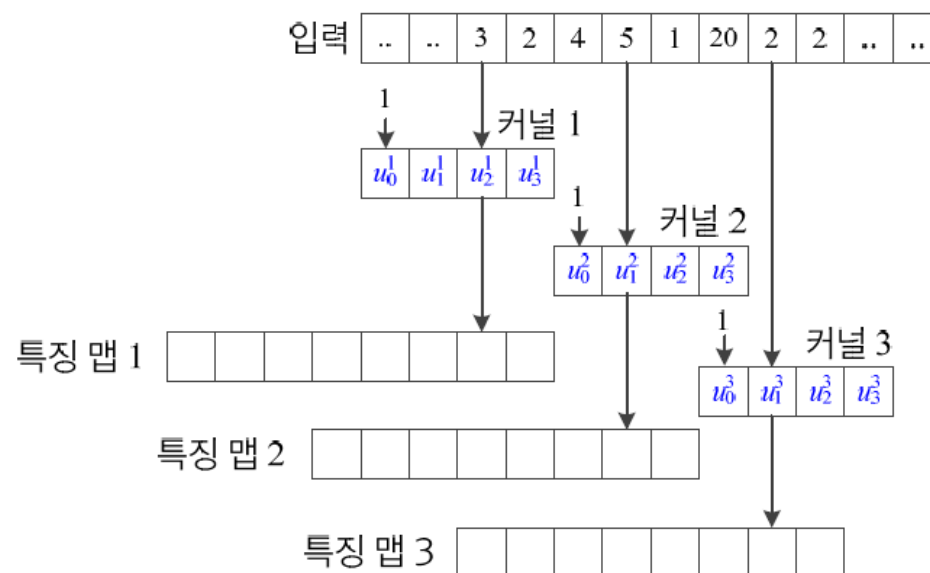
- 모든 노드가 동일한 커널을 사용 (즉 가중치를 공유)하므로 매개변수는 3개에 불과
- 모델의 복잡도가 크게 낮아짐

커널의 값에 따라 커널이 추출하는 특징이 달라짐

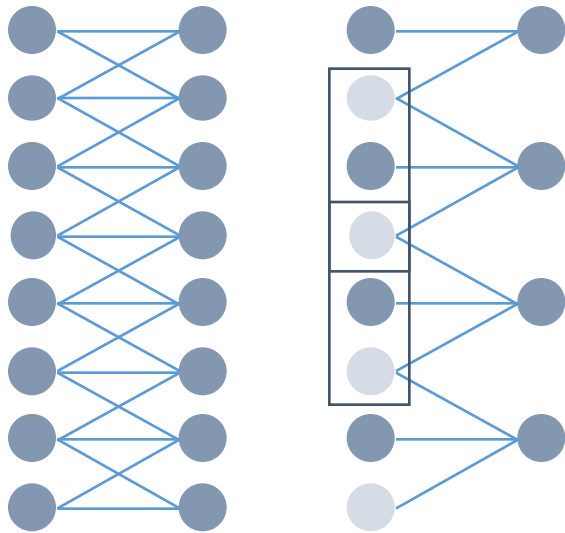
하나의 커널만 사용하는 것보다 여러개의 커널을 이용하는 것이 특징 추출에 도움을 줌



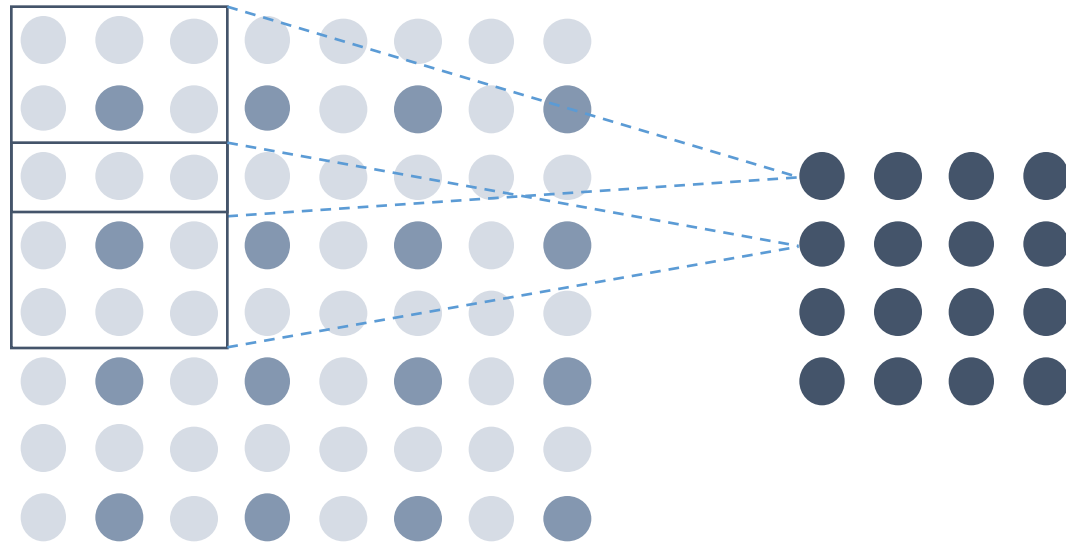




- 커널을 사람이 설계하지 않고, 학습으로 알아냄
- 예) 2차원 영상이  $7 \times 7$  커널을 64개 사용한다면,  
학습은  $(7 \times 7 + 1) \times 64 = 3200$ 개의 매개변수를 찾아내야 함
- DMLP와 마찬가지로 오류 역전파로 커널을 학습



1차원 데이터



2차원 데이터(예: 영상, 이미지)

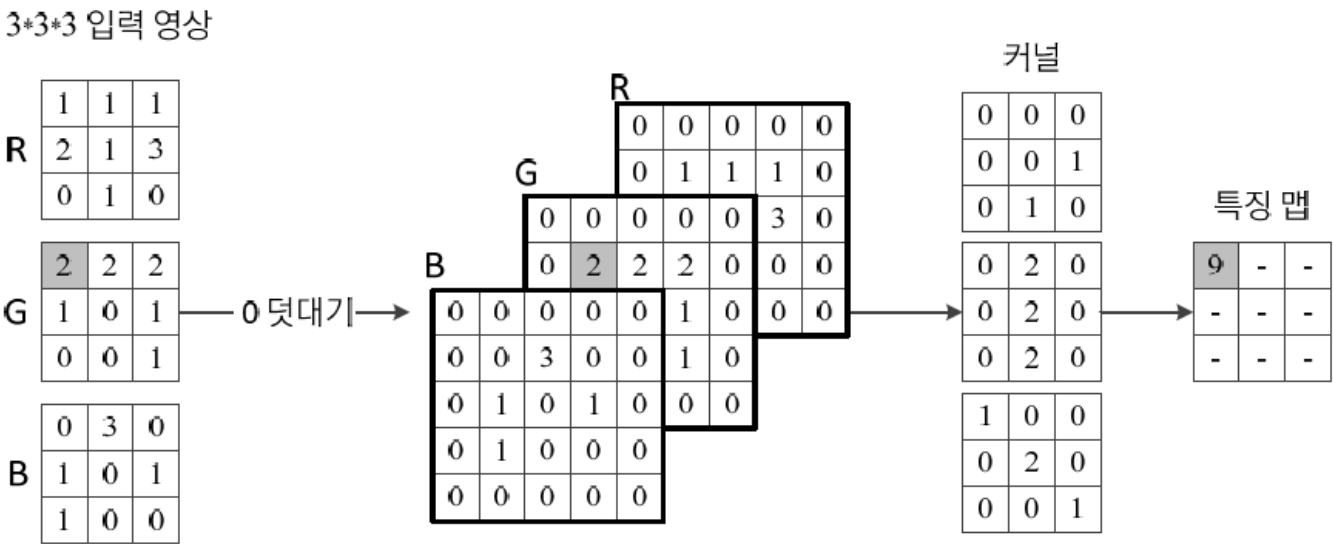
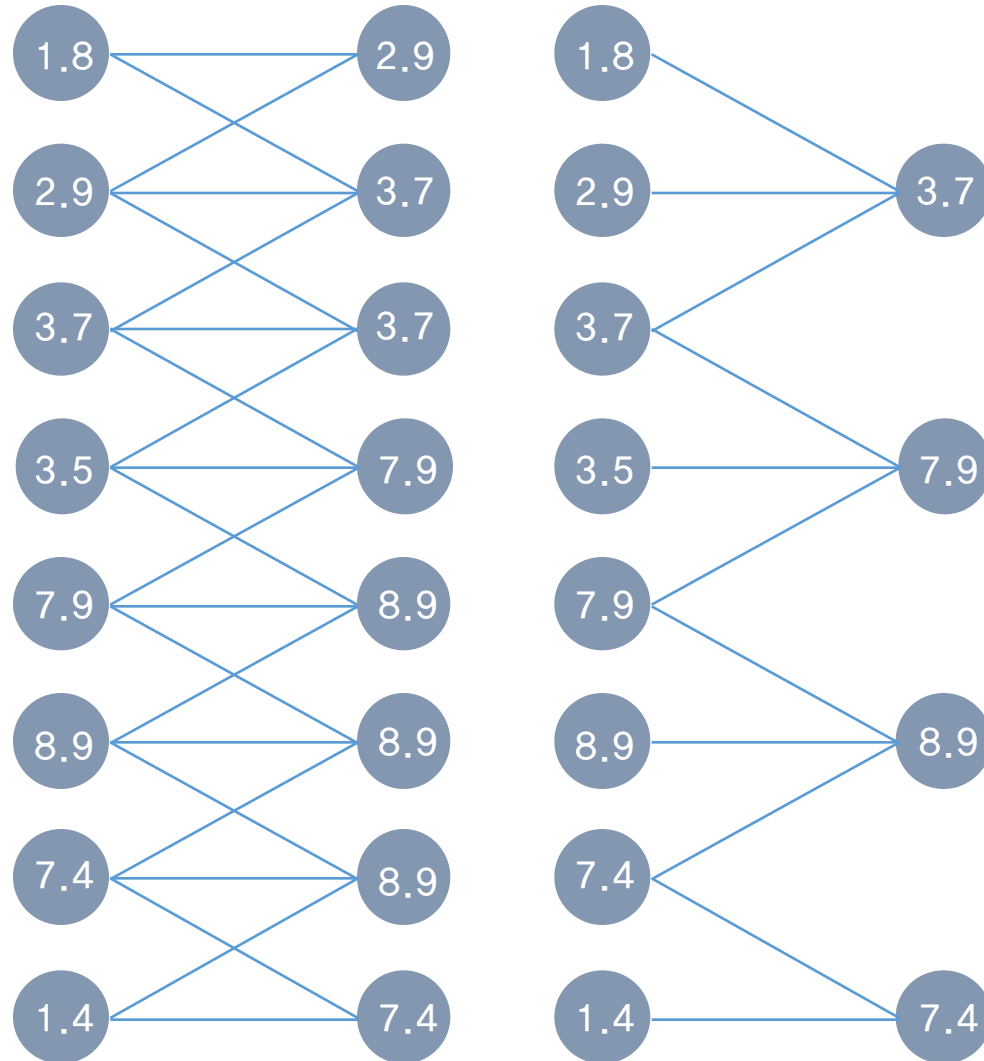


그림 4-14 텐서의 컨볼루션 연산(0 덧대기 적용)

$$\underbrace{\begin{pmatrix} 0 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 2 & 1 \end{pmatrix}}_R \underbrace{\begin{pmatrix} 0 & 0 & 0 \\ 0 & 2 & 2 \\ 0 & 1 & 0 \end{pmatrix}}_G \underbrace{\begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 3 \\ 0 & 1 & 0 \end{pmatrix}}_B \otimes \underbrace{\begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix}}_{c_1} \underbrace{\begin{pmatrix} 0 & 2 & 0 \\ 0 & 2 & 0 \\ 0 & 2 & 0 \end{pmatrix}}_{c_2} \underbrace{\begin{pmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 1 \end{pmatrix}}_{c_3} = 9$$



최대풀링, 평균 풀링, 가중치 평균 풀링

보폭(stride)을 크게 하면 다운샘플링 효과

풀링 연산의 특성

- 풀링은 상세 내용에서 요약통계를 추출함
- 매개변수가 없음
- 특징 맵의 수를 그대로 유지함
- 작은 이동에 둔감 → 물체인식이나 영상 검색 등에 효과적임

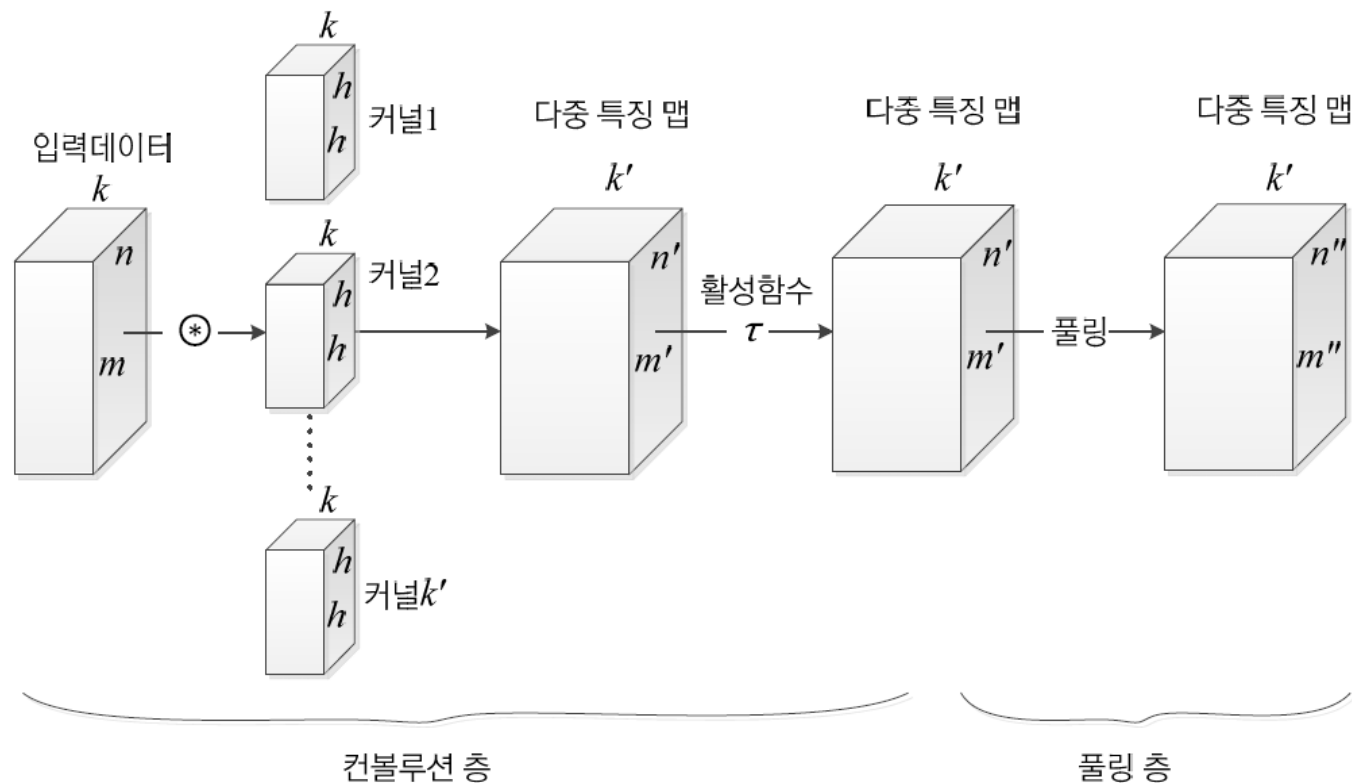
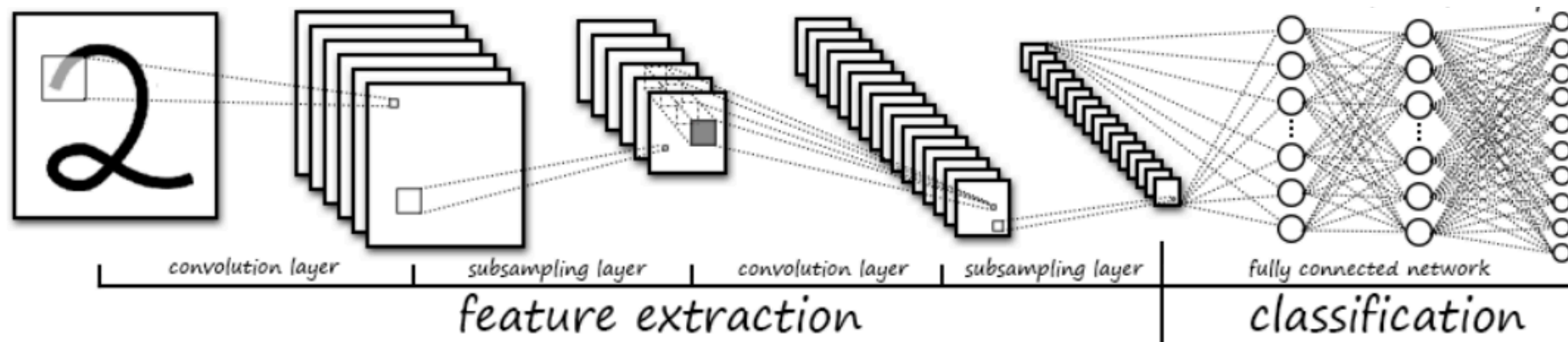
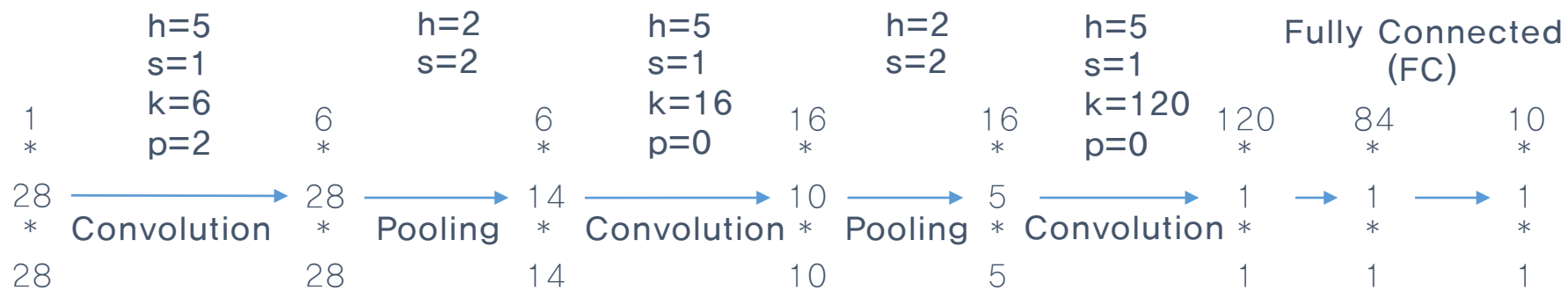


그림 4-18 CNN의 빌딩블록

- CNN은 빌딩블록을 이어 붙여 깊은 구조를 만듦
- 전형적인 빌딩블록 : 컨볼루션층  $\rightarrow$  활성화함수(주로 ReLU)  $\rightarrow$  풀링층
- 다중커널을 사용하여 다중 특징 맵을 추출함

# Convolution Neural Network 응용사례 : LeNet-5

로컬랩서울 locallab-seoul.com



DMLP는 특징 벡터의 크기가 달라지면 연산 불가능

CNN은 가변 크기를 다룰 수 있는 강점

컨볼루션층에서 보폭을 조정한다거나, 풀링층에서 커널이  
나 보폭을 조정하여 특징 맵 크기를 조절



컨볼루션층 5개와 완전연결(FC)층 3개  
8개 층에 290400-186624-64896-43264-4096-4096-1000개의 노드 배치  
컨볼루션층은 200만개, FC층은 6500만개 가량의 매개변수  
FC층에 30배 많은 매개변수 ☾ 향후 CNN은 FC층의 매개변수를 줄이는 방향으로 발전

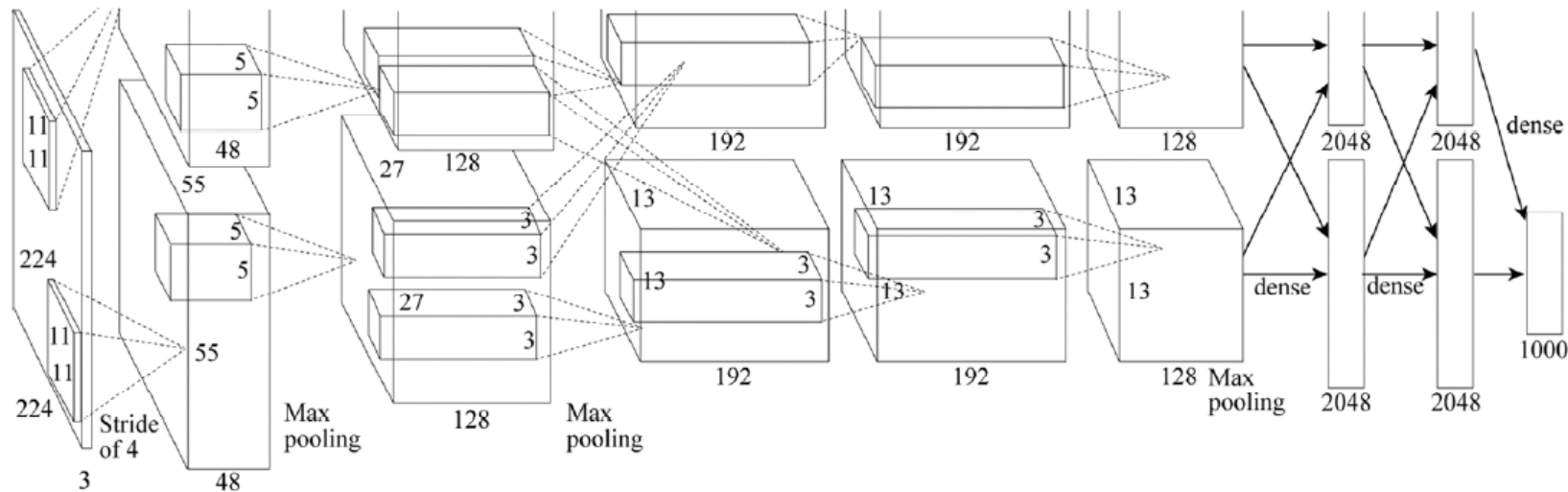


그림 4-21 AlexNet 구조[Krizhevsky2012]

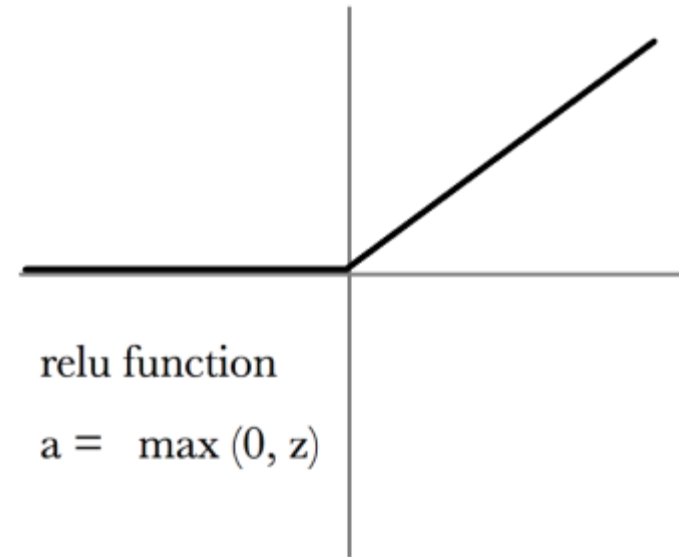
## AlexNet이 학습에 성공한 요인

### 외부 요인

- ImageNet이라는 대용량 데이터베이스
- GPU를 사용한 병렬처리

### 내부 요인

- 활성화함수로 ReLU 사용
- 지역 반응 정규화 기법 적용
- 과잉적합 방지하는 여러 규제 기법 적용
- 데이터 확대(잘라내기과 반전으로 2048배로 확대)
- 드롭아웃 등



## VGGNet의 핵심 아이디어

3\*3의 작은 커널을 사용하여 신경망을 더욱 깊게 만듦

컨볼루션층 8~16개를 두어 AlexNet의 5개에 비해 2~3배 깊어짐

16층짜리 VGG-16(컨볼루션 13층+FC 3층)

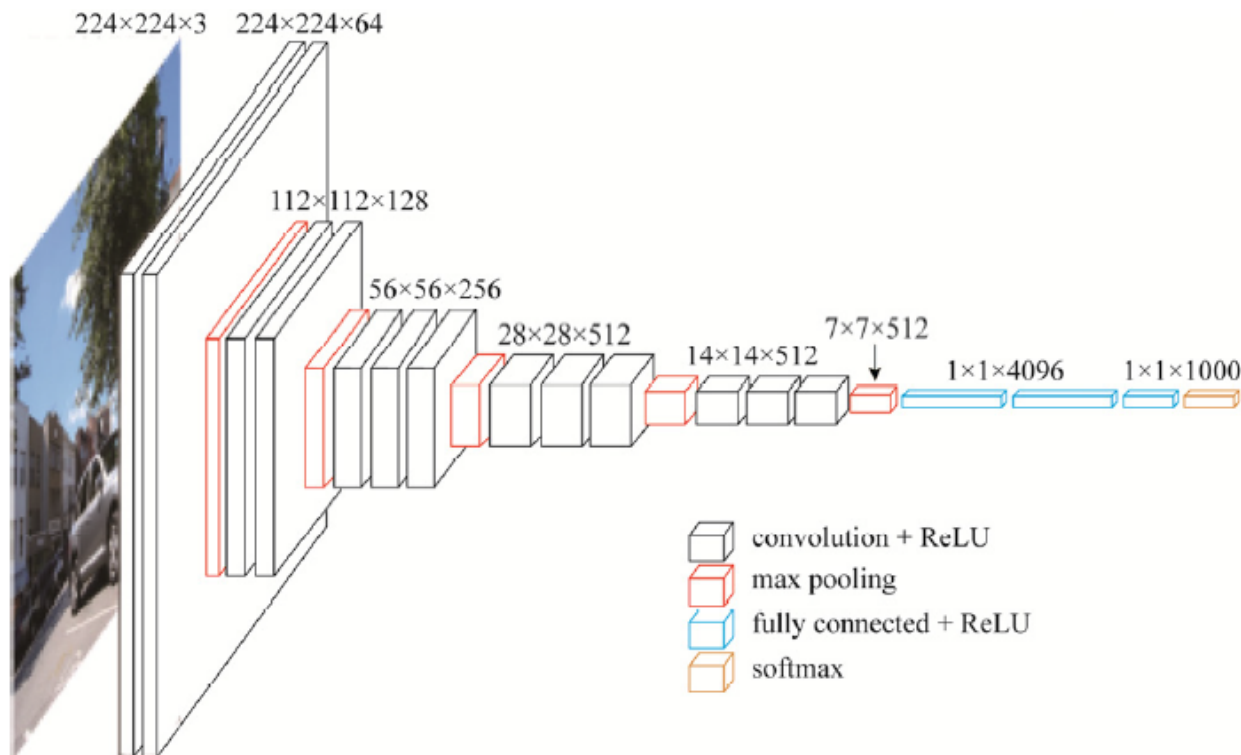


그림 4-22 VGGNet 구조[Simonyan2015]

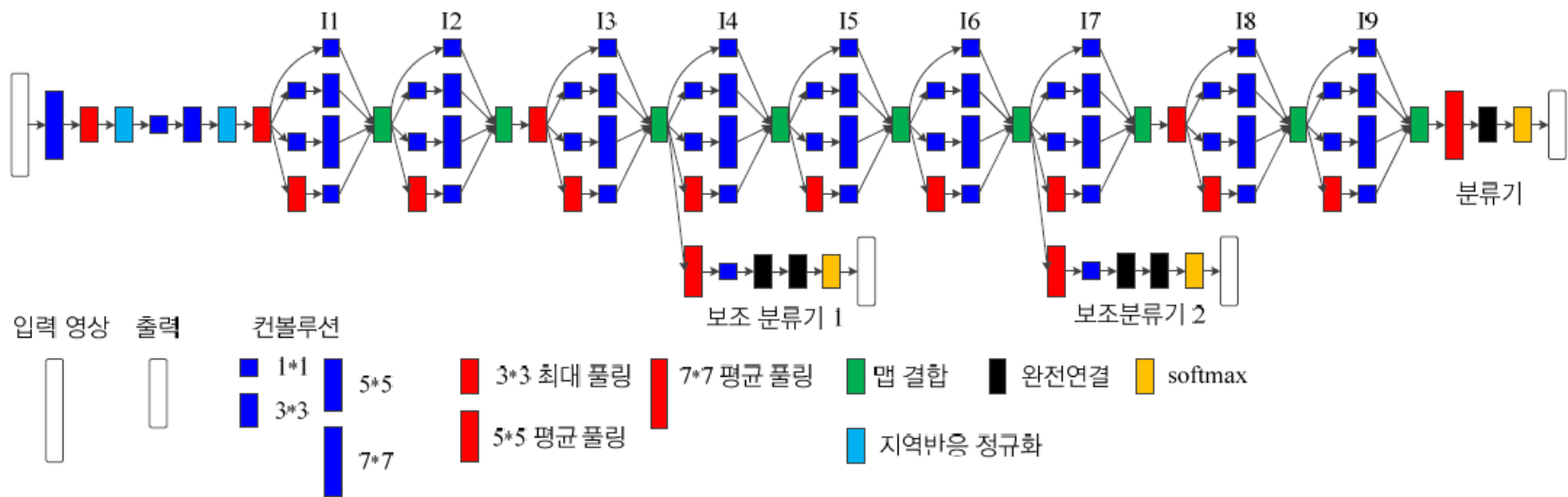


그림 4-27 GoogLeNet의 구조

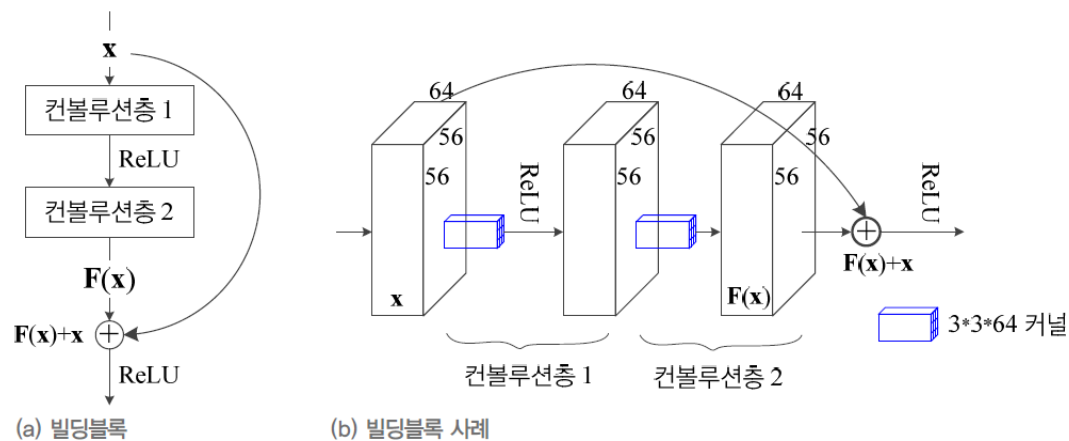


그림 4-28 잔류 학습의 구조와 동작

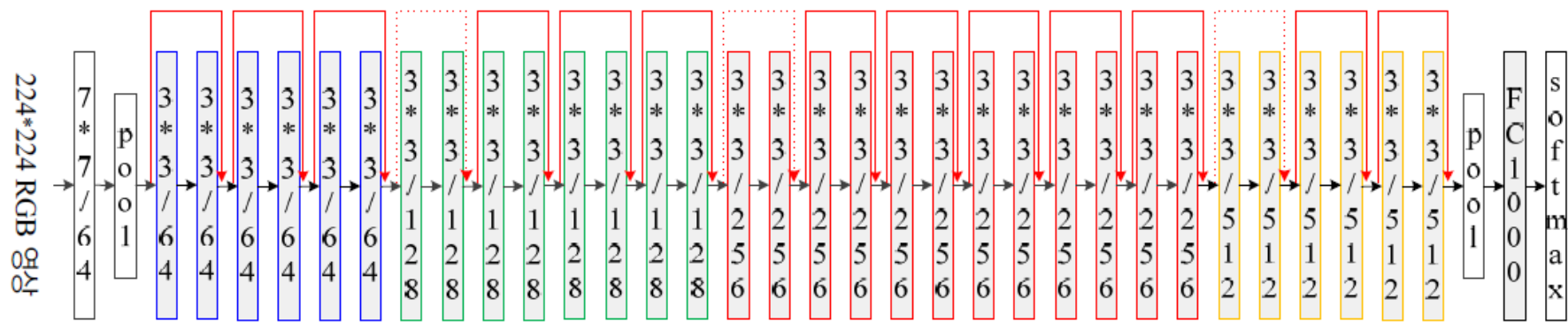


그림 4-29 ResNet 예제(34층)

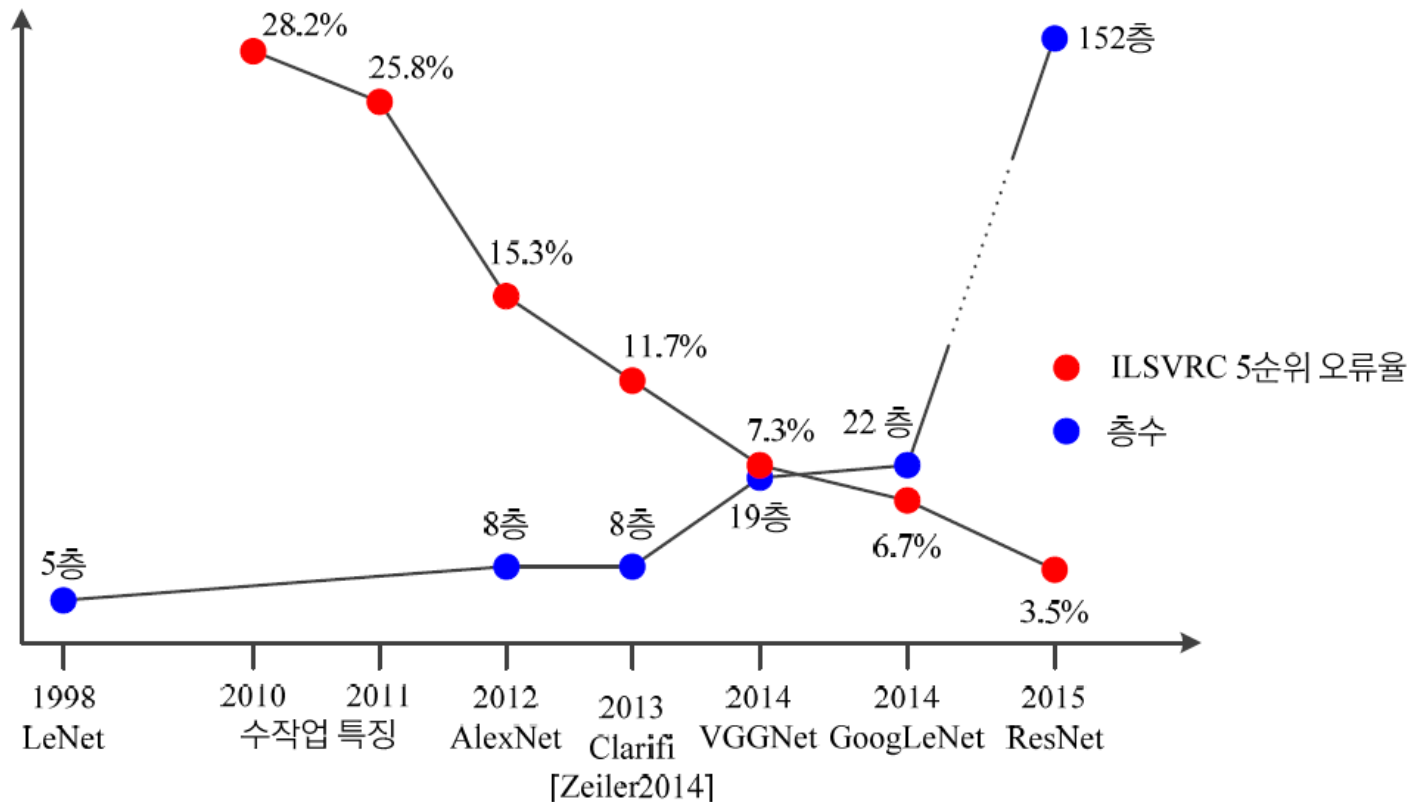
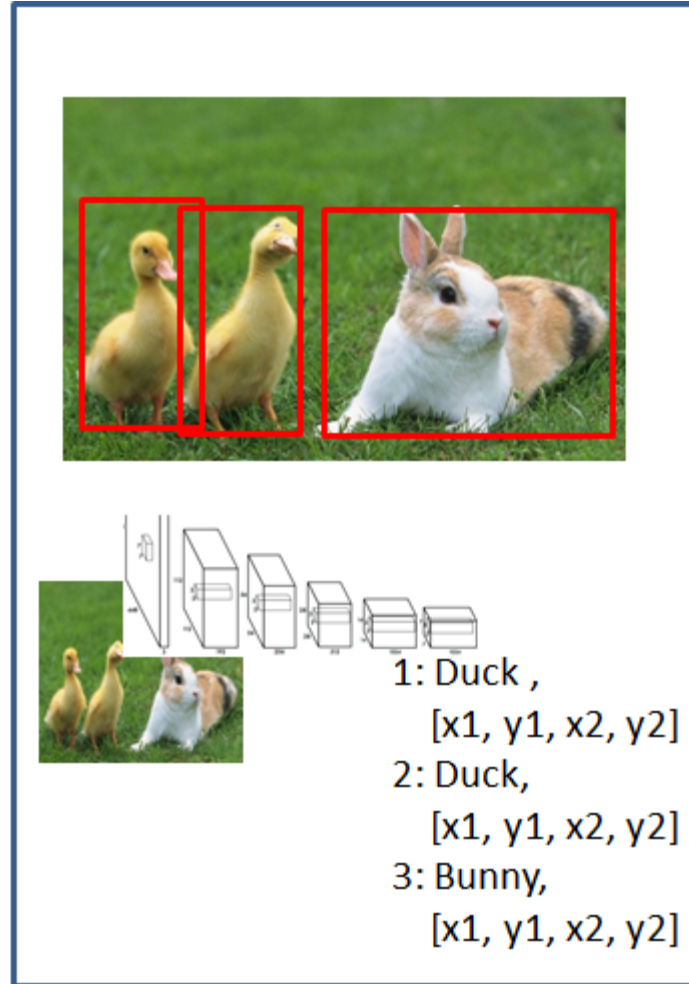


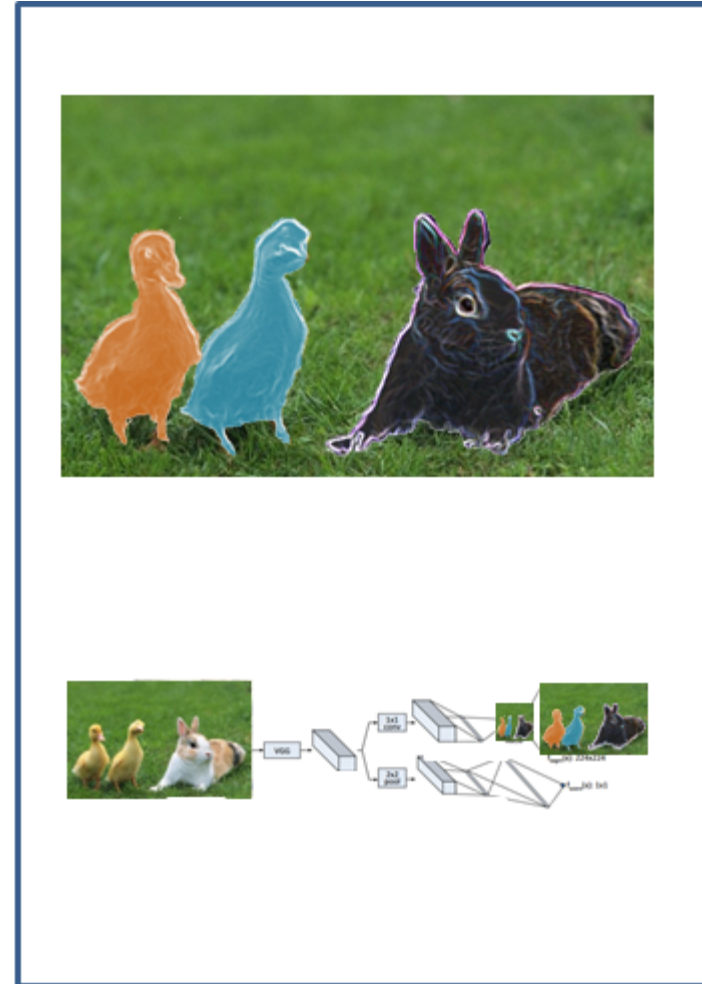
그림 4-30 CNN의 발전 추세



Classification



Object Detection



Segmentation

