

OCR实战

第二课 单字符识别

第二课 单字符识别

- 单字符识别
 - 模版匹配
 - 特征提取 + 传统分类器
 - CNN模型
- 过切分的路径选择
 - VITERBI + BEAM SEARCH
- 语言模型
 - N-GRAM
 - LSTM
 - 2D-PCFG

单字符识别

$$D(i, j) = \sum_{m=1}^M \sum_{n=1}^N [S^{ij}(m, n) - T(m, n)]^2$$

欧式距离相似性函数

- 模版匹配

$$\sum_{m=1}^M \sum_{n=1}^N [S^{ij}(m, n)]^2 - 2 \sum_{m=1}^M \sum_{n=1}^N S^{ij}(m, n) \cdot T(m, n) + \sum_{m=1}^M \sum_{n=1}^N [T(m, n)]^2$$

展开后相似性函数



$$R(i, j) = \frac{\sum_{m=1}^M \sum_{n=1}^N S^{ij}(m, n) \cdot T(m, n)}{\sqrt{\sum_{m=1}^M \sum_{n=1}^N [S^{ij}(m, n)]^2} \sqrt{\sum_{m=1}^M \sum_{n=1}^N [T(m, n)]^2}}$$

归一化相似性函数

- method=CV_TM_SQDIFF

$$R(x, y) = \sum_{x', y'} (T(x', y') - I(x + x', y + y'))^2$$

- method=CV_TM_SQDIFF_NORMED

- method=CV_TM_CCORR

$$R(x, y) = \frac{\sum_{x', y'} (T(x', y') - I(x + x', y + y'))^2}{\sqrt{\sum_{x', y'} T(x', y')^2 \cdot \sum_{x', y'} I(x + x', y + y')^2}}$$

$$R(x, y) = \sum_{x', y'} (T(x', y') \cdot I(x + x', y + y'))$$

- method=CV_TM_CCORR_NORMED

$$R(x, y) = \frac{\sum_{x', y'} (T(x', y') \cdot I(x + x', y + y'))}{\sqrt{\sum_{x', y'} T(x', y')^2 \cdot \sum_{x', y'} I(x + x', y + y')^2}}$$

- method=CV_TM_CCOEFF

$$R(x, y) = \sum_{x', y'} (T'(x', y') \cdot I'(x + x', y + y'))$$

where

$$\begin{aligned} T'(x', y') &= T(x', y') - 1/(w \cdot h) \cdot \sum_{x'', y''} T(x'', y'') \\ I'(x + x', y + y') &= I(x + x', y + y') - 1/(w \cdot h) \cdot \sum_{x'', y''} I(x + x'', y + y'') \end{aligned}$$

- method=CV_TM_CCOEFF_NORMED

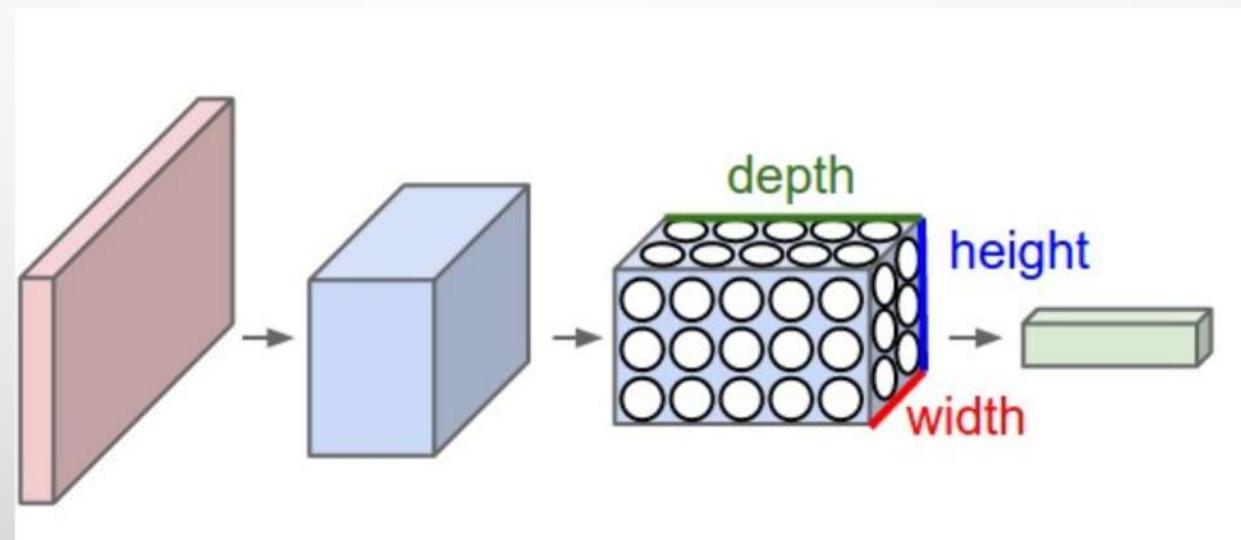
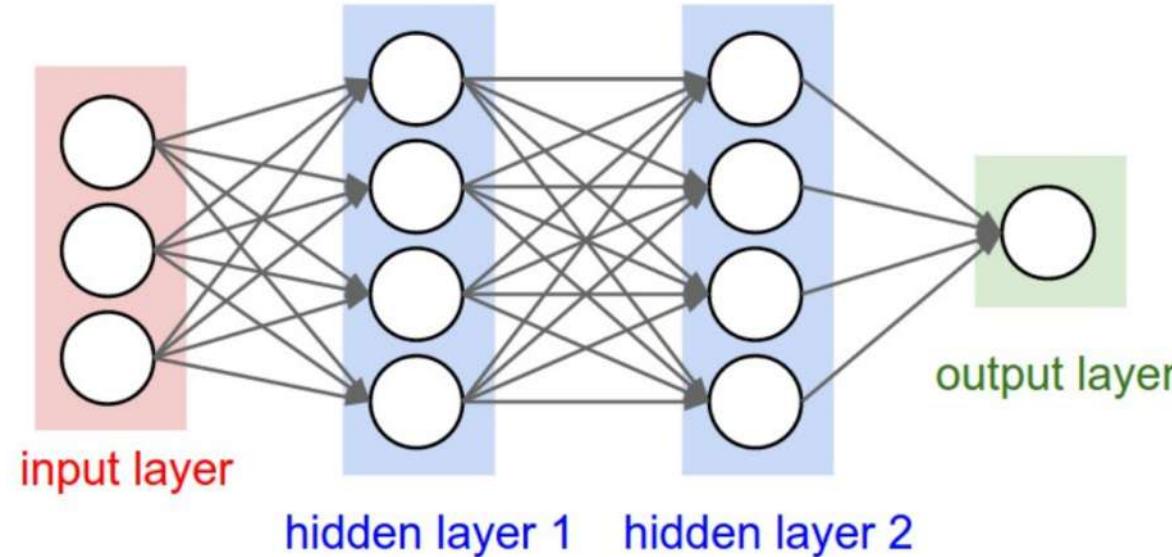
$$R(x, y) = \frac{\sum_{x', y'} (T'(x', y') \cdot I'(x + x', y + y'))}{\sqrt{\sum_{x', y'} T'(x', y')^2 \cdot \sum_{x', y'} I'(x + x', y + y')^2}}$$

单字符识别

- 特征提取 + 传统分类器
 - 位图
 - 垂直/水平投影
 - 占空比 / 孔洞 / 穿刺
 - 方向 / 弧度
 - ...
- +
 - SVM
 - MLP

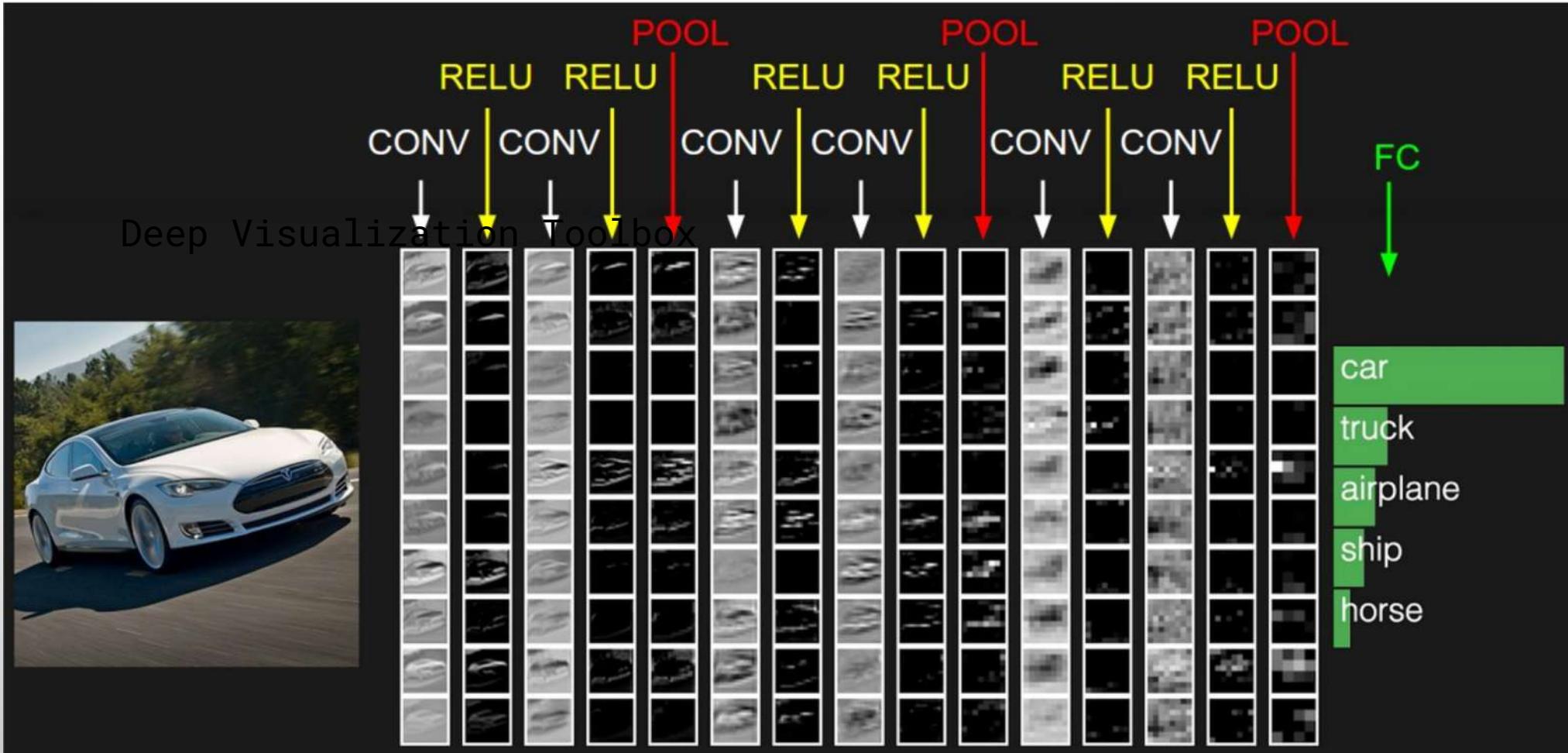
单字符识别

- CNN模型



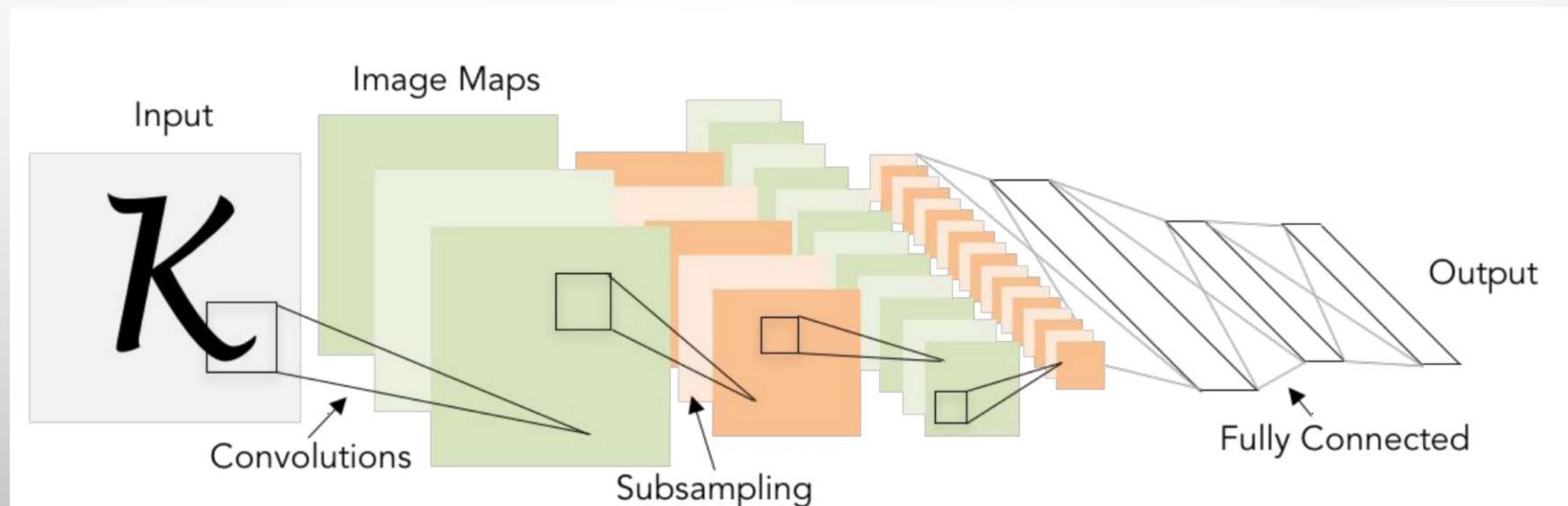
单字符识别

- CNN模型



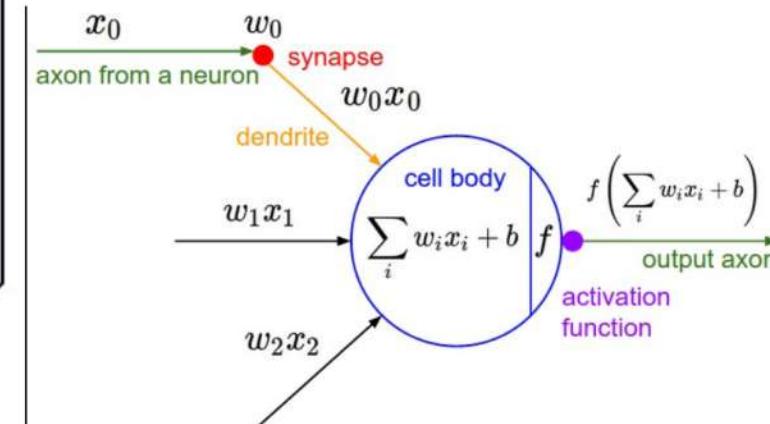
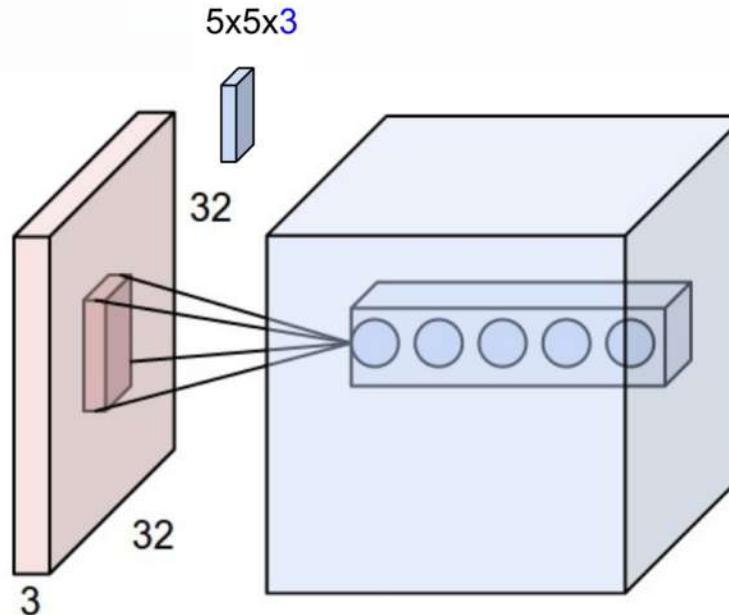
单字符识别

- CNN模型

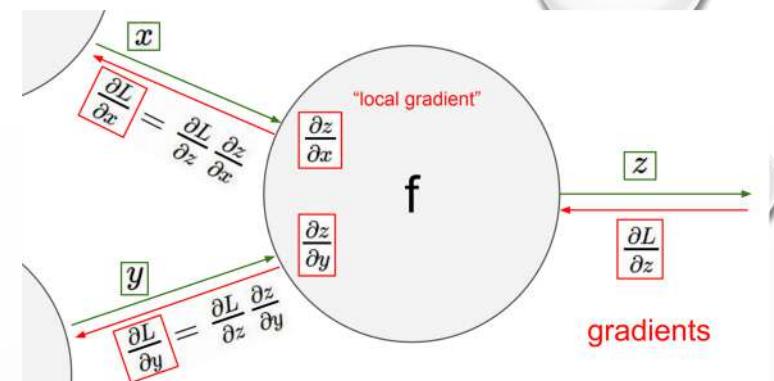


单字符识别

- CNN模型



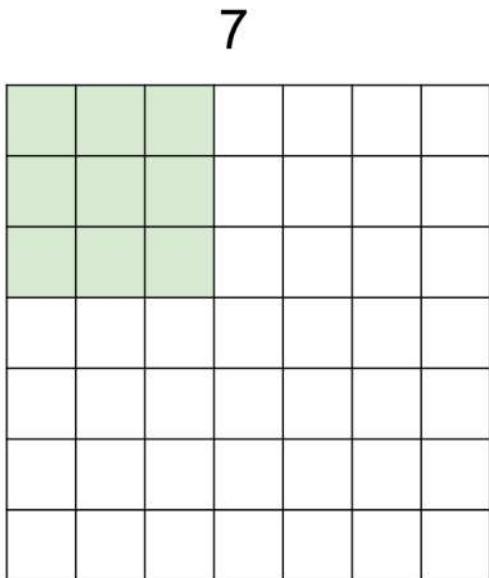
Left: An example input volume in red (e.g. a 32x32x3 CIFAR-10 image), and an example volume of neurons in the first Convolutional layer. Each neuron in the convolutional layer is connected only to a local region in the input volume spatially, but to the full depth (i.e. all color channels). Note, there are multiple neurons (5 in this example) along the depth, all looking at the same region in the input - see discussion of depth columns in text below. Right: The neurons from the Neural Network chapter remain unchanged: They still compute a dot product of their weights with the input followed by a non-linearity, but their connectivity is now restricted to be local spatially.



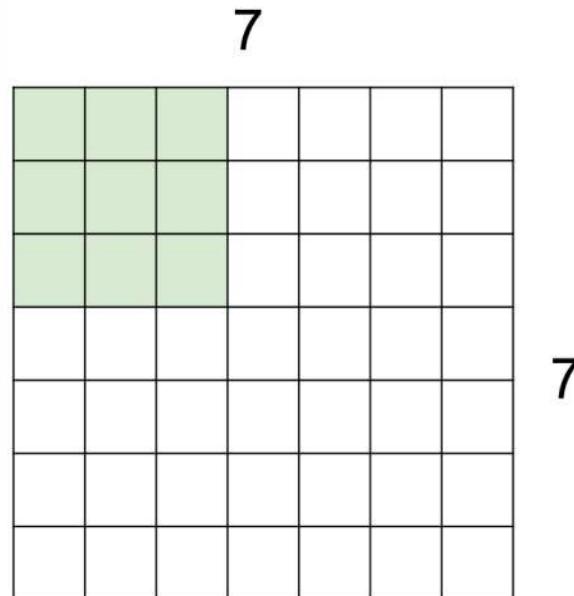
单字符识别

- CNN模型

A closer look at spatial dimensions:

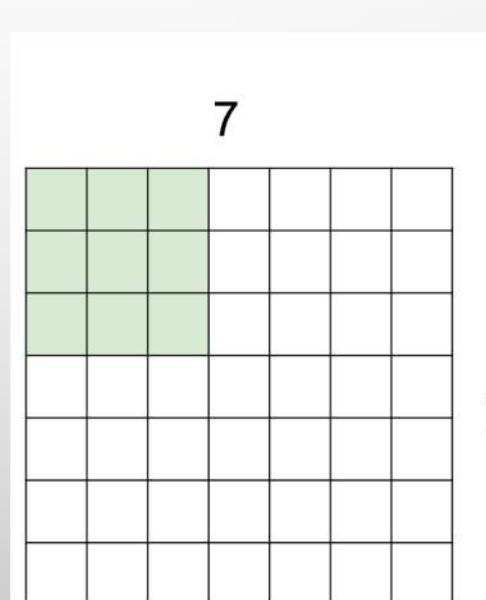


7x7 input (spatially)
assume 3x3 filter



7

7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**



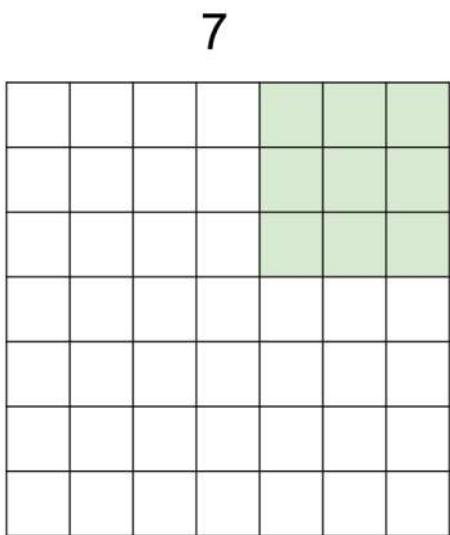
7

7x7 input (spatially)
assume 3x3 filter
applied **with stride 3?**

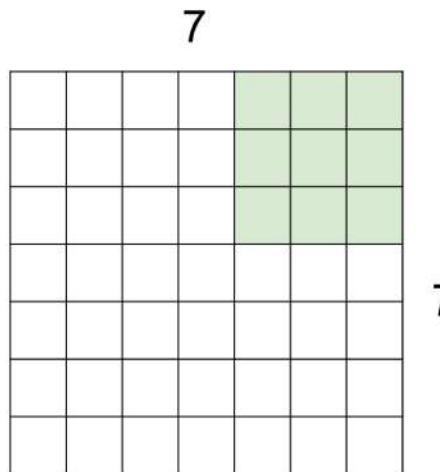
单字符识别

- CNN模型

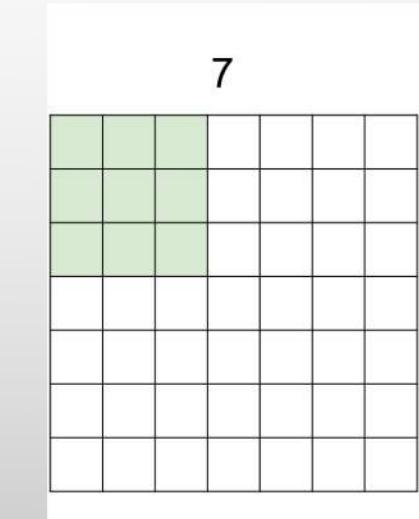
A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter
=> 5x5 output



7



7

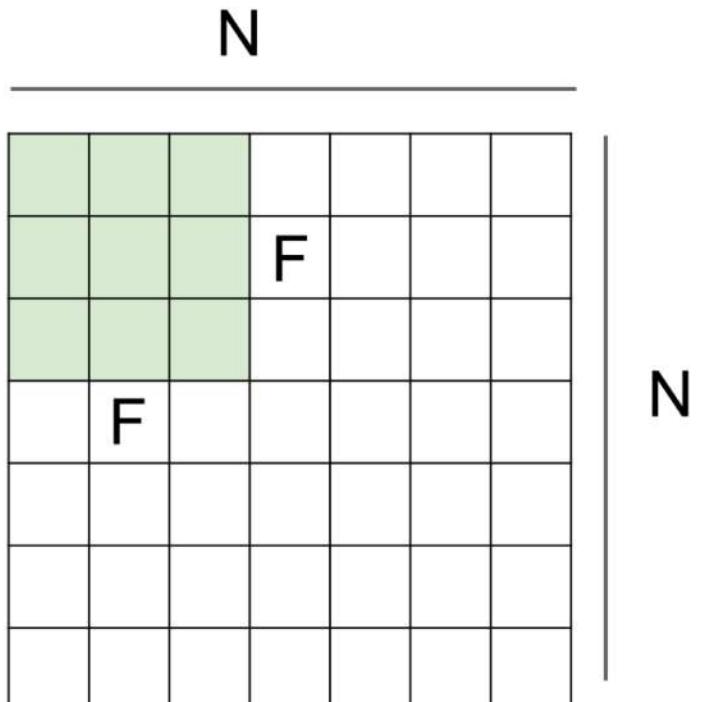
7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**
=> 3x3 output!

7x7 input (spatially)
assume 3x3 filter
applied **with stride 3?**

doesn't fit!
cannot apply 3x3 filter on
7x7 input with stride 3.

单字符识别

- CNN模型



Output size:
(N - F) / stride + 1

e.g. $N = 7$, $F = 3$:
stride 1 => $(7 - 3)/1 + 1 = 5$
stride 2 => $(7 - 3)/2 + 1 = 3$
stride 3 => $(7 - 3)/3 + 1 = 2.33$:\

单字符识别

- CNN模型

0	0	0	0	0	0		
0							
0							
0							
0							

e.g. input 7x7

3x3 filter, applied with **stride 1**

pad with 1 pixel border => what is the output?

(recall:)

$$(N - F) / \text{stride} + 1$$

单字符识别

In practice: Common to zero pad the border

- CNN模型

0	0	0	0	0	0		
0							
0							
0							
0							

e.g. input 7x7

3x3 filter, applied with **stride 1**

pad with 1 pixel border => what is the output?

7x7 output!

in general, common to see CONV layers with stride 1, filters of size FxF, and zero-padding with $(F-1)/2$. (will preserve size spatially)

e.g. $F = 3 \Rightarrow$ zero pad with 1

$F = 5 \Rightarrow$ zero pad with 2

$F = 7 \Rightarrow$ zero pad with 3

单字符识别

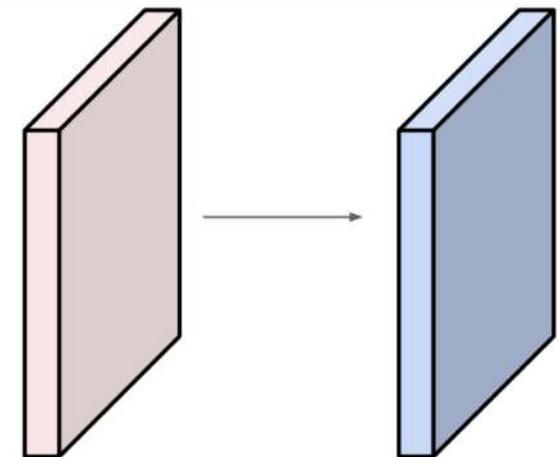
- CNN模型

Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

Output volume size: ?

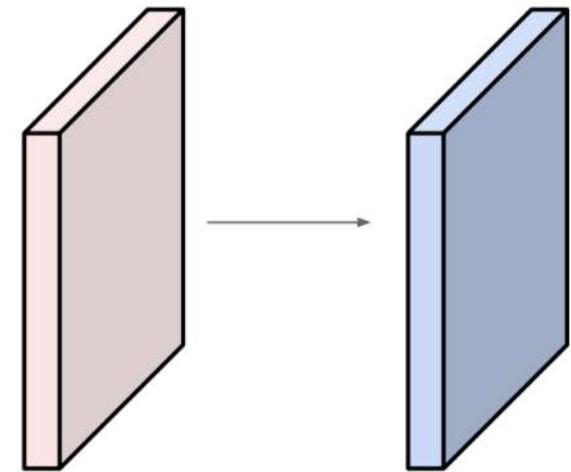


单字符识别

- CNN模型

Examples time:

Input volume: **32x32x3**
10 5x5 filters with stride **1**, pad **2**



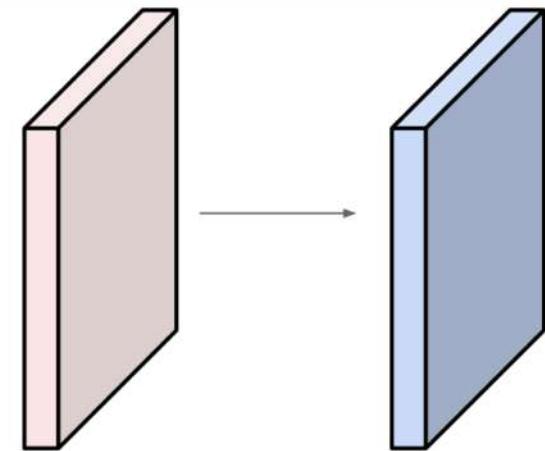
Output volume size:
 $(32+2*2-5)/1+1 = 32$ spatially, so
32x32x10

单字符识别

- CNN模型

Examples time:

Input volume: **32x32x3**
10 5x5 filters with stride 1, pad 2



Number of parameters in this layer?

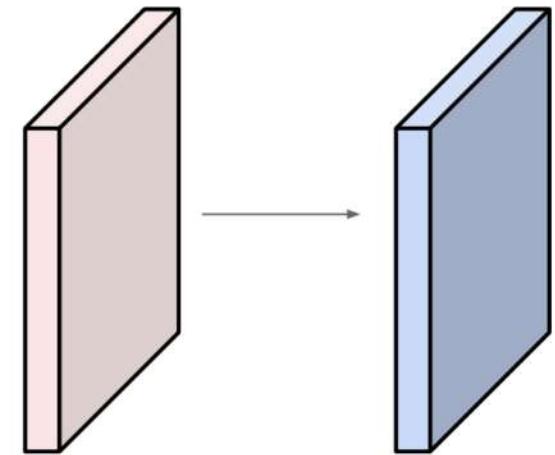
单字符识别

- CNN模型

Examples time:

Input volume: **32x32x3**

10 **5x5** filters with stride 1, pad 2



Number of parameters in this layer?

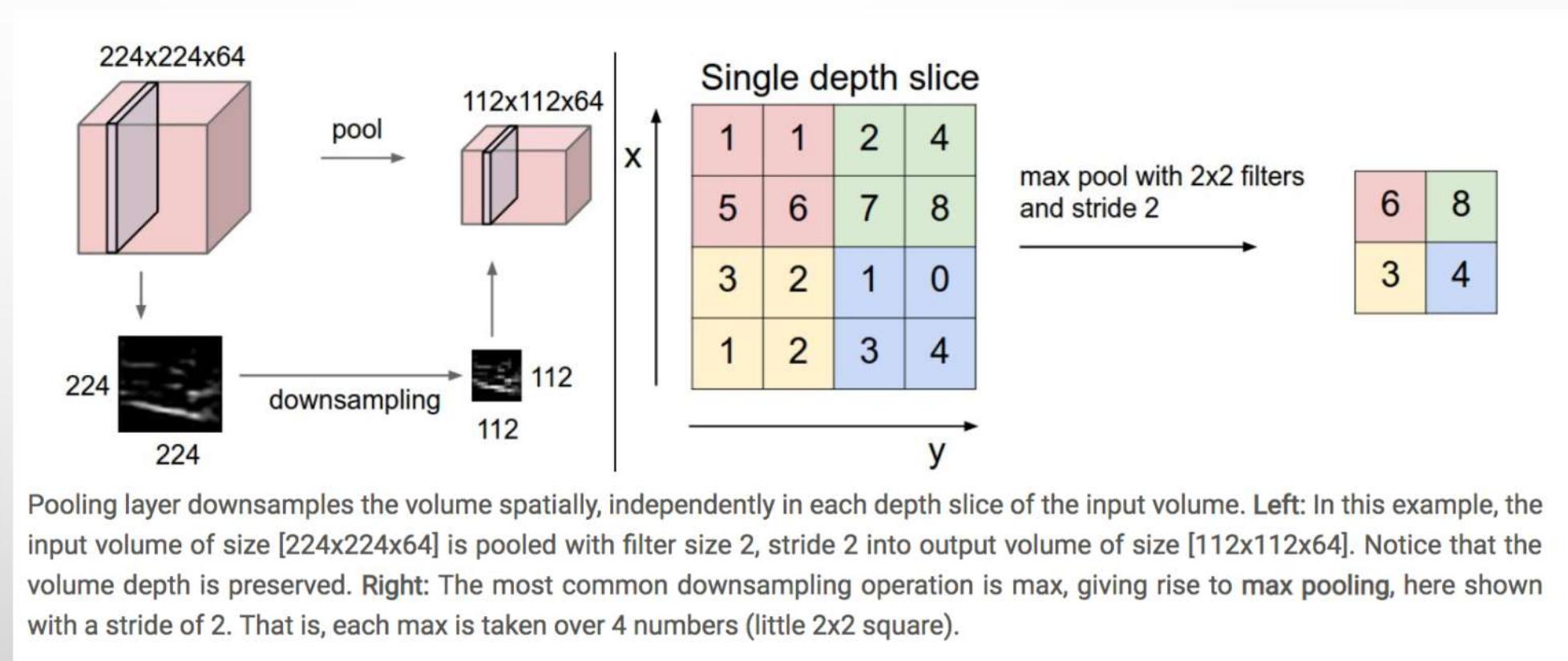
each filter has **5*5*3 + 1 = 76** params

(+1 for bias)

$$\Rightarrow 76 * 10 = 760$$

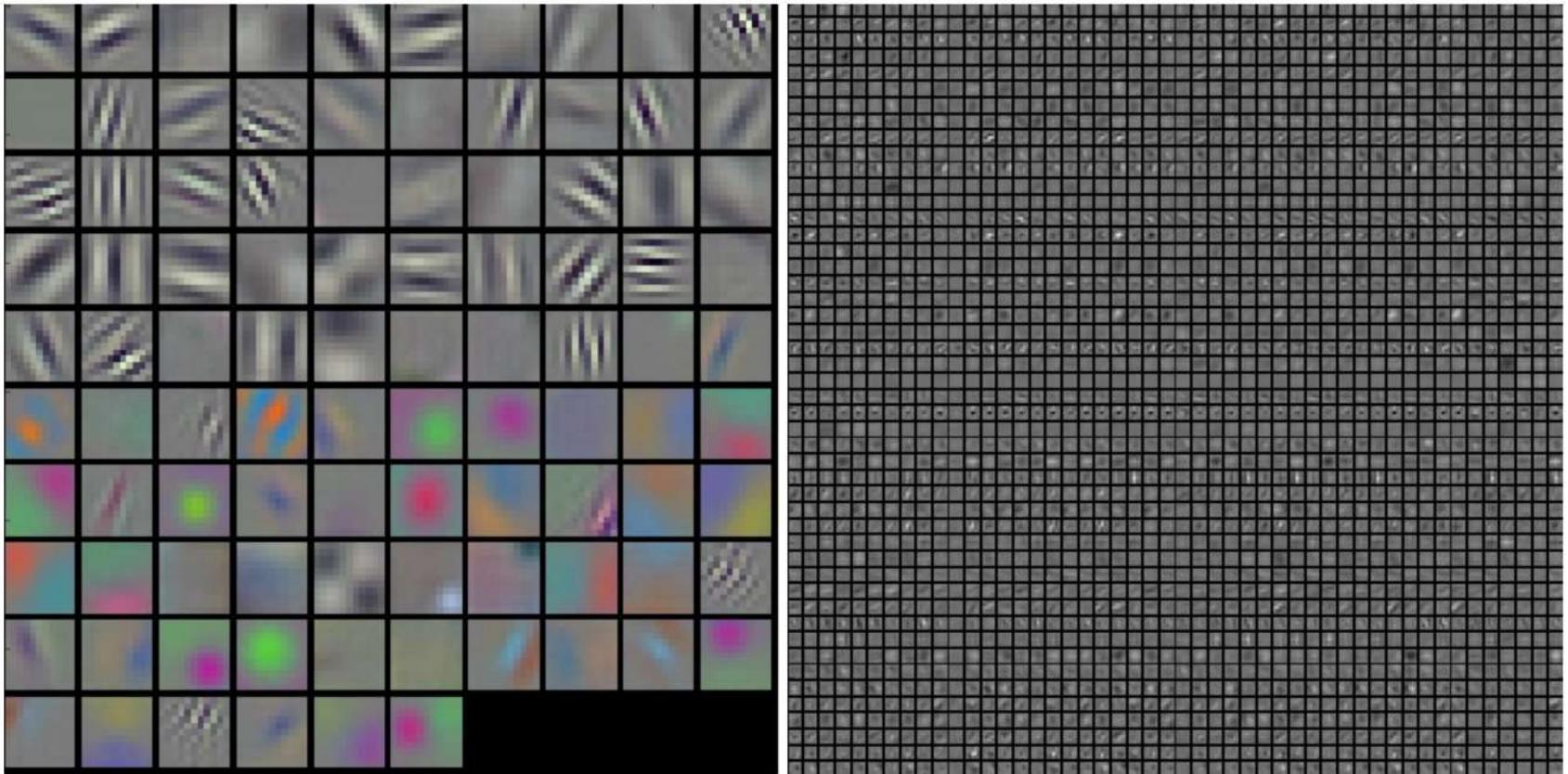
单字符识别

- CNN模型



单字符识别

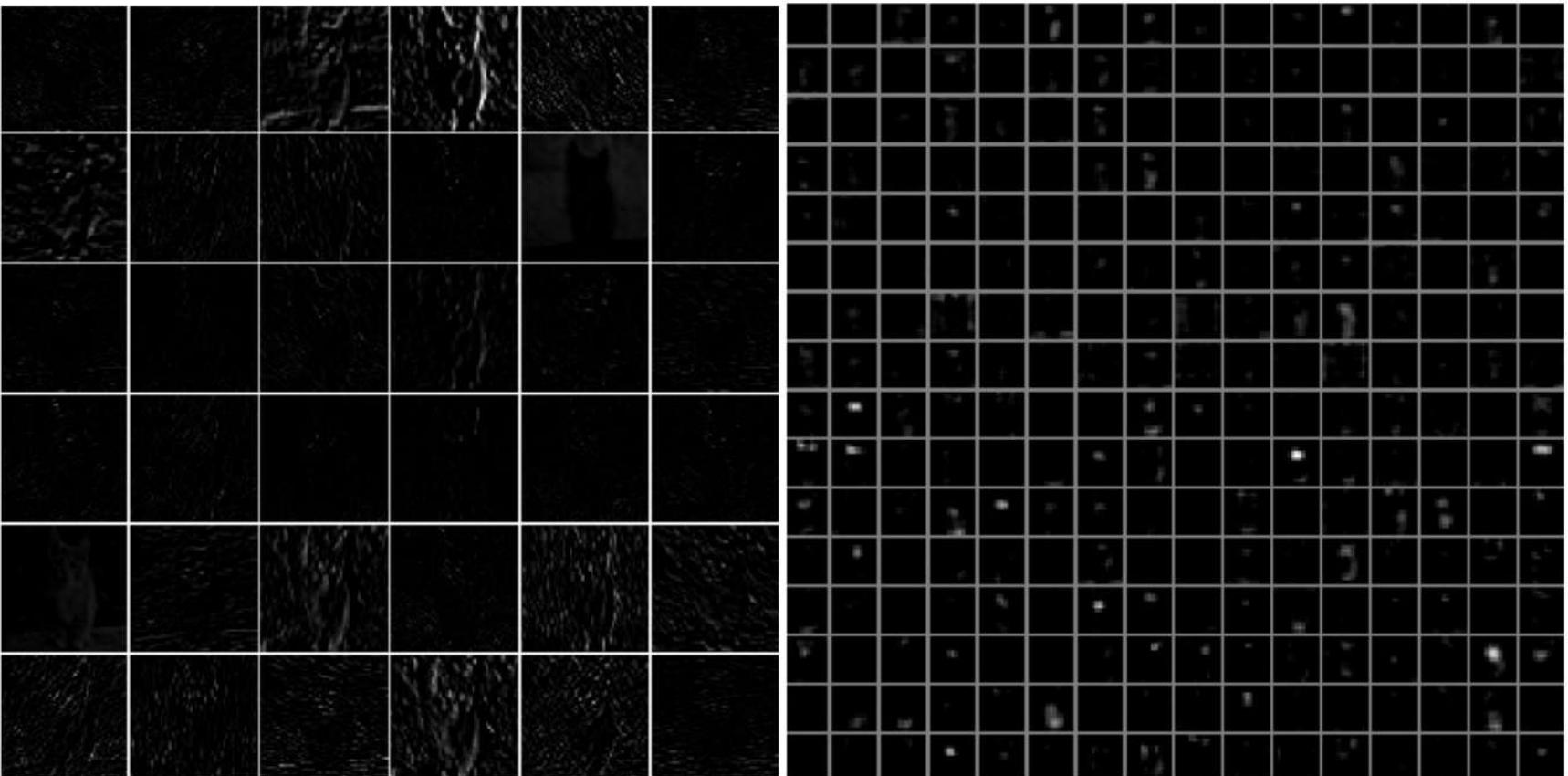
- CNN模型



Typical-looking filters on the first CONV layer (left), and the 2nd CONV layer (right) of a trained AlexNet. Notice that the first-layer weights are very nice and smooth, indicating nicely converged network. The color/grayscale features are clustered because the AlexNet contains two separate streams of processing, and an apparent consequence of this architecture is that one stream develops high-frequency grayscale features and the other low-frequency color features. The 2nd CONV layer weights are not as interpretable, but it is apparent that they are still smooth, well-formed, and absent of noisy patterns.

单字符识别

- CNN模型



Typical-looking activations on the first CONV layer (left), and the 5th CONV layer (right) of a trained AlexNet looking at a picture of a cat. Every box shows an activation map corresponding to some filter. Notice that the activations are sparse (most values are zero, in this visualization shown in black) and mostly local.

单字符识别

- CNN模型



Maximally activating images for some POOL5 (5th pool layer) neurons of an AlexNet. The activation values and the receptive field of the particular neuron are shown in white. (In particular, note that the POOL5 neurons are a function of a relatively large portion of the input image!) It can be seen that some neurons are responsive to upper bodies, text, or specular highlights.

单字符识别

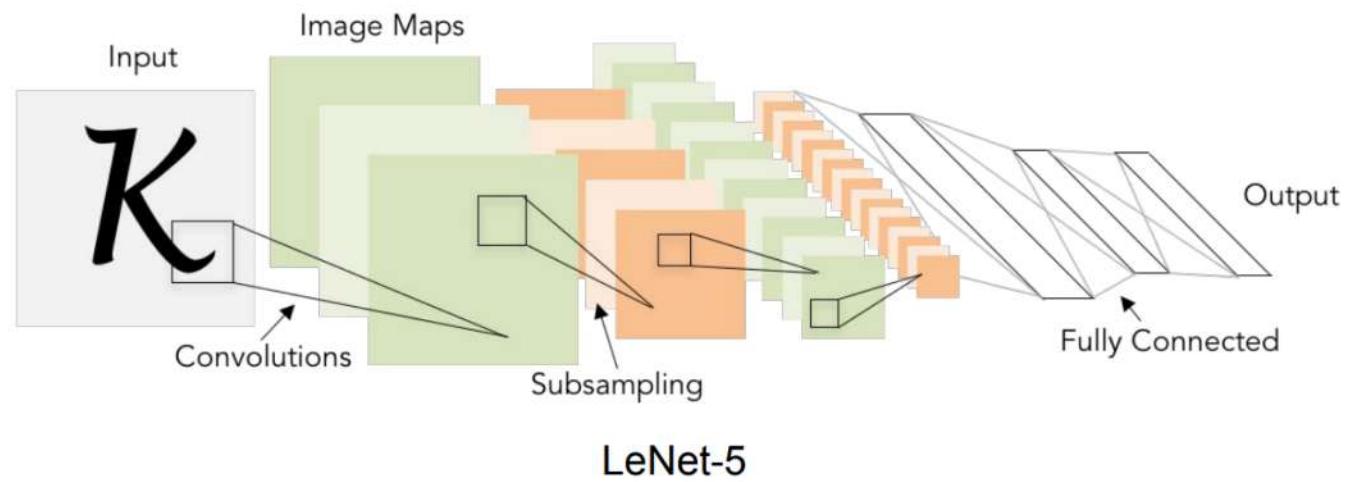
- CNN模型
 - DEEP VISUALIZATION TOOLBOX
 - [HTTPS://WWW.YOUTUBE.COM/WATCH?V=AGKFIQ4IGAM&T=7S](https://www.youtube.com/watch?v=AGKFIQ4IGAM&t=7s)
 - CONVNETJS
 - [HTTP://CS.STANFORD.EDU/PEOPLE/KARPATHY/CONVNETJS/](http://cs.stanford.edu/people/karpathy/convnetjs/)

单字符识别

- CNN模型

Conv filters were 5x5, applied at stride 1
Subsampling (Pooling) layers were 2x2 applied at stride 2
i.e. architecture is [CONV-POOL-CONV-POOL-FC-FC]

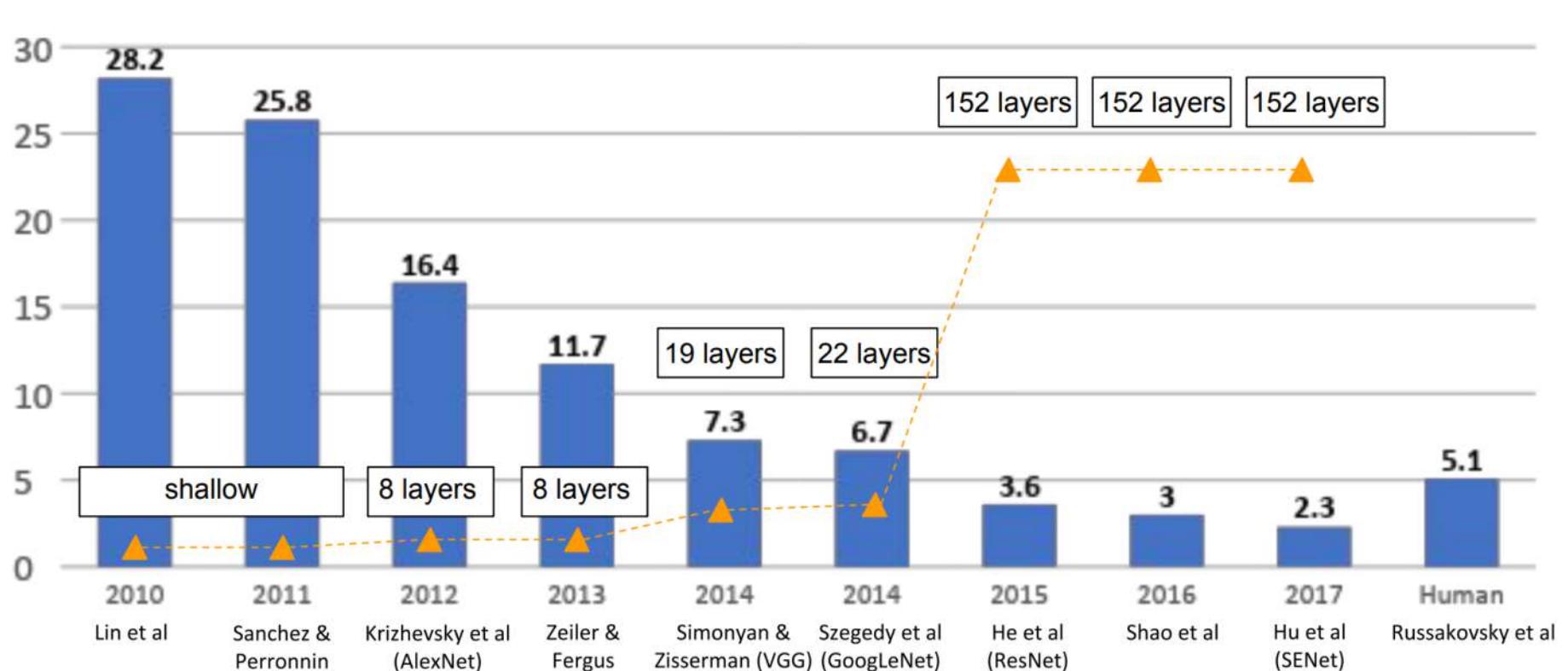
Gradient-based learning applied to document recognition *[LeCun, Bottou, Bengio, Haffner 1998]*



单字符识别

- CNN模型

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



单字符识别

Case Study: VGGNet

[Simonyan and Zisserman, 2014]

- CNN模型

Small filters, Deeper networks

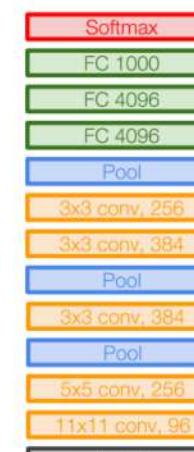
8 layers (AlexNet)

-> 16 - 19 layers (VGG16Net)

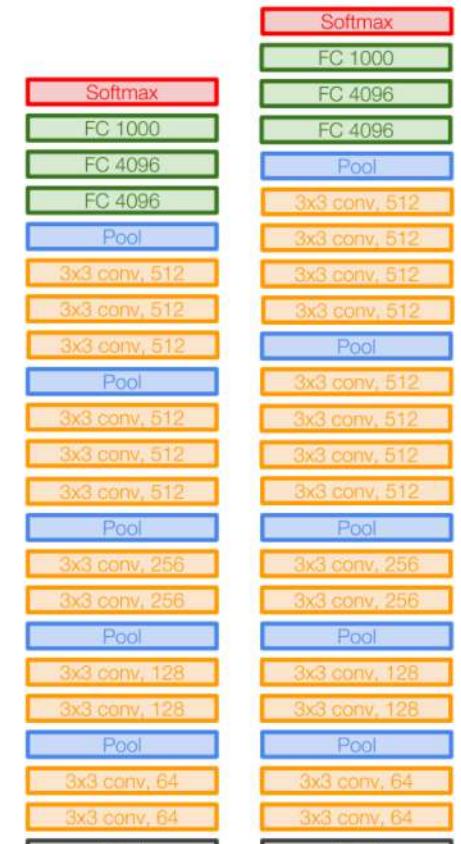
Only 3x3 CONV stride 1, pad 1
and 2x2 MAX POOL stride 2

11.7% top 5 error in ILSVRC'13
(ZFNet)

-> 7.3% top 5 error in ILSVRC'14



AlexNet



VGG16

VGG19

单字符识别

Case Study: VGGNet

[Simonyan and Zisserman, 2014]

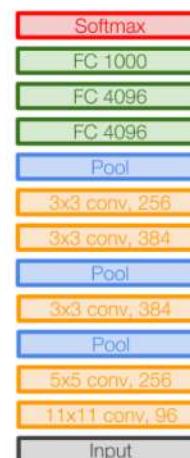
- CNN模型

Q: Why use smaller filters? (3x3 conv)

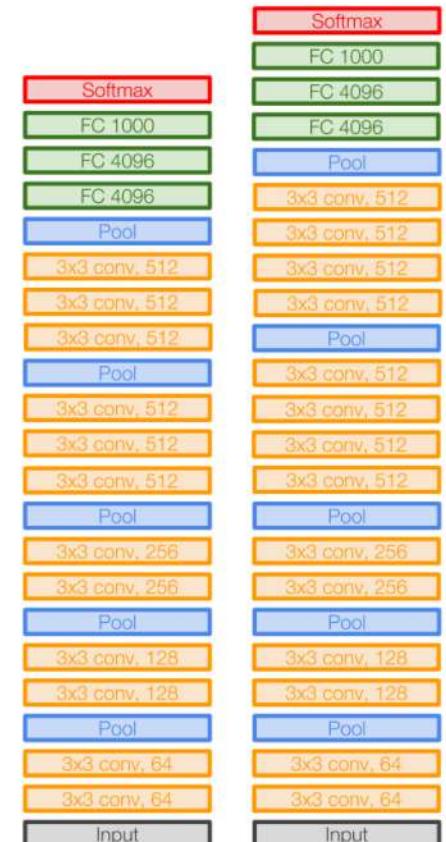
Stack of three 3x3 conv (stride 1) layers has same **effective receptive field** as one 7x7 conv layer

But deeper, more non-linearities

And fewer parameters: $3 * (3^2 C^2)$ vs. $7^2 C^2$ for C channels per layer



AlexNet



VGG16

VGG19

单字符识别

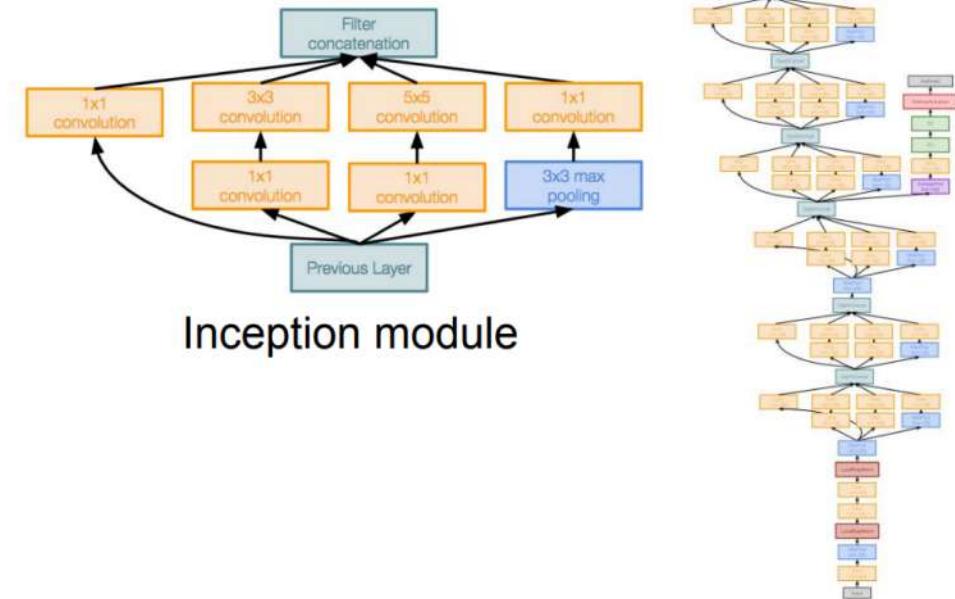
- CNN模型

Case Study: GoogLeNet

[Szegedy et al., 2014]

Deeper networks, with computational efficiency

- 22 layers
- Efficient “Inception” module
- No FC layers
- Only 5 million parameters!
12x less than AlexNet
- ILSVRC’14 classification winner
(6.7% top 5 error)

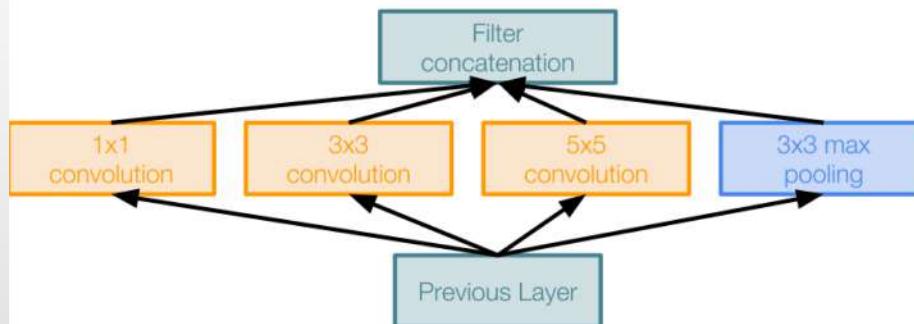


单字符识别

Case Study: GoogLeNet

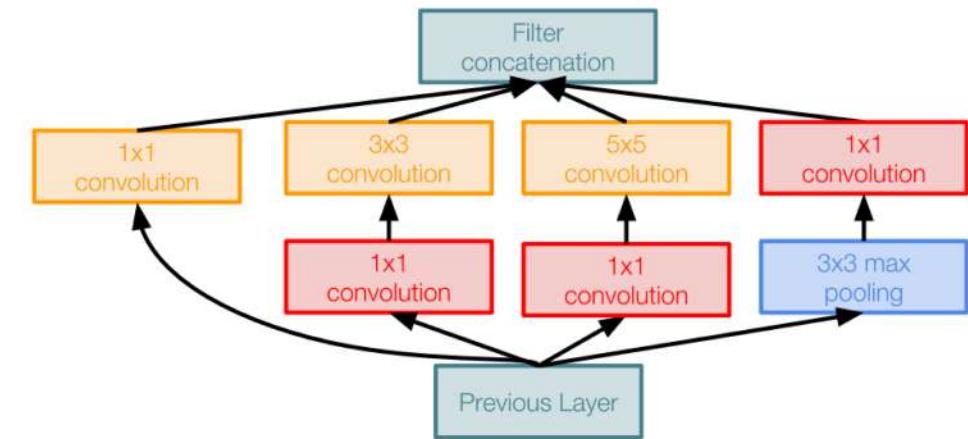
- CNN模型

[Szegedy et al., 2014]



Naive Inception module

1x1 conv “bottleneck”
layers

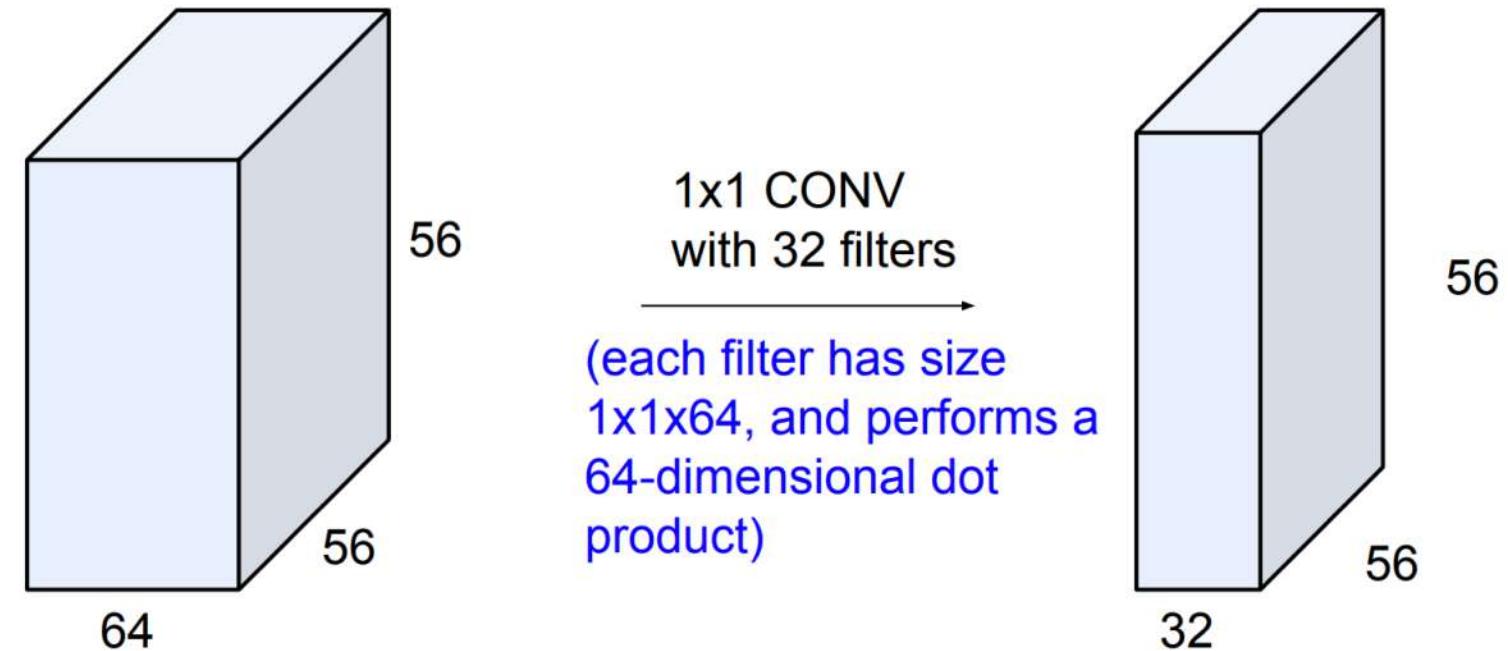


Inception module with dimension reduction

单字符识别

(btw, 1x1 convolution layers make perfect sense)

- CNN模型

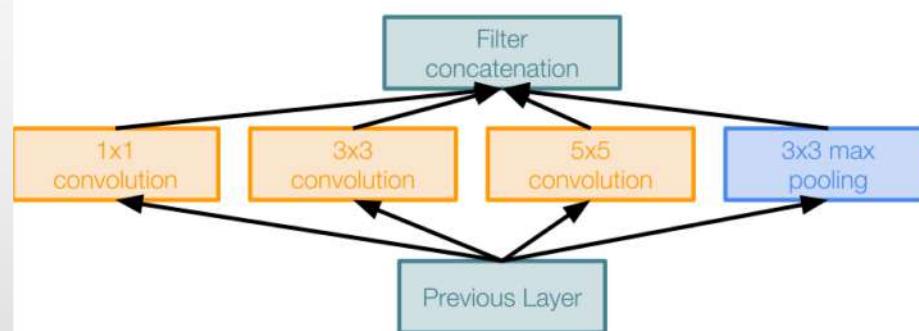


单字符识别

Case Study: GoogLeNet

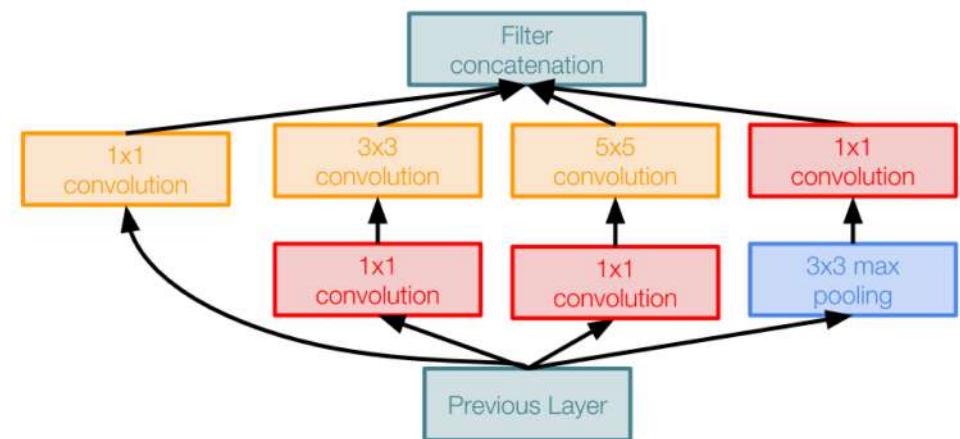
- CNN模型

[Szegedy et al., 2014]



Naive Inception module

1x1 conv “bottleneck”
layers



Inception module with dimension reduction

单字符识别

- CNN模型

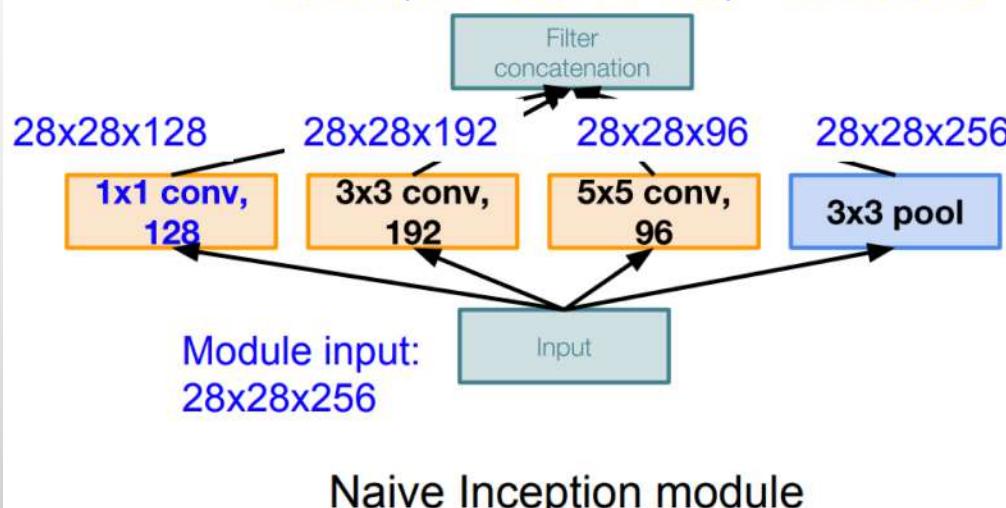
Case Study: GoogLeNet

[Szegedy et al., 2014]

Example:

Q3: What is output size after filter concatenation?

$$28 \times 28 \times (128 + 192 + 96 + 256) = 28 \times 28 \times 672$$



Q: What is the problem with this?
[Hint: Computational complexity]

Conv Ops:

[1x1 conv, 128] $28 \times 28 \times 128 \times 1 \times 1 \times 256$

[3x3 conv, 192] $28 \times 28 \times 192 \times 3 \times 3 \times 256$

[5x5 conv, 96] $28 \times 28 \times 96 \times 5 \times 5 \times 256$

Total: 854M ops

Very expensive compute

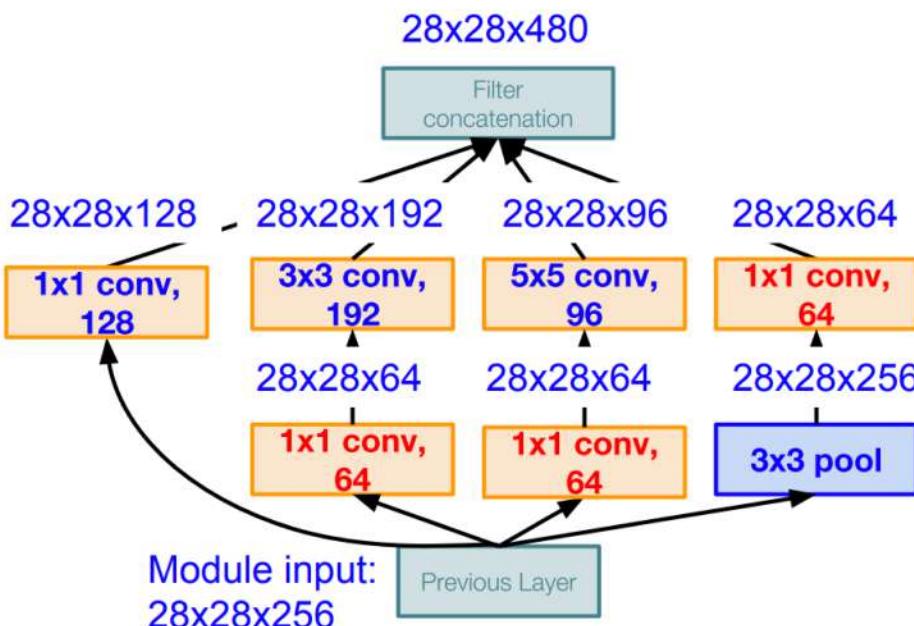
Pooling layer also preserves feature depth, which means total depth after concatenation can only grow at every layer!

单字符识别

- CNN模型

Case Study: GoogLeNet

[Szegedy et al., 2014]



Inception module with dimension reduction

Using same parallel layers as naive example, and adding “1x1 conv, 64 filter” bottlenecks:

Conv Ops:

[1x1 conv, 64] 28x28x64x1x1x256
[1x1 conv, 64] 28x28x64x1x1x256
[1x1 conv, 128] 28x28x128x1x1x256
[3x3 conv, 192] 28x28x192x3x3x64
[5x5 conv, 96] 28x28x96x5x5x64
[1x1 conv, 64] 28x28x64x1x1x256

Total: 358M ops

Compared to 854M ops for naive version
Bottleneck can also reduce depth after pooling layer

单字符识别

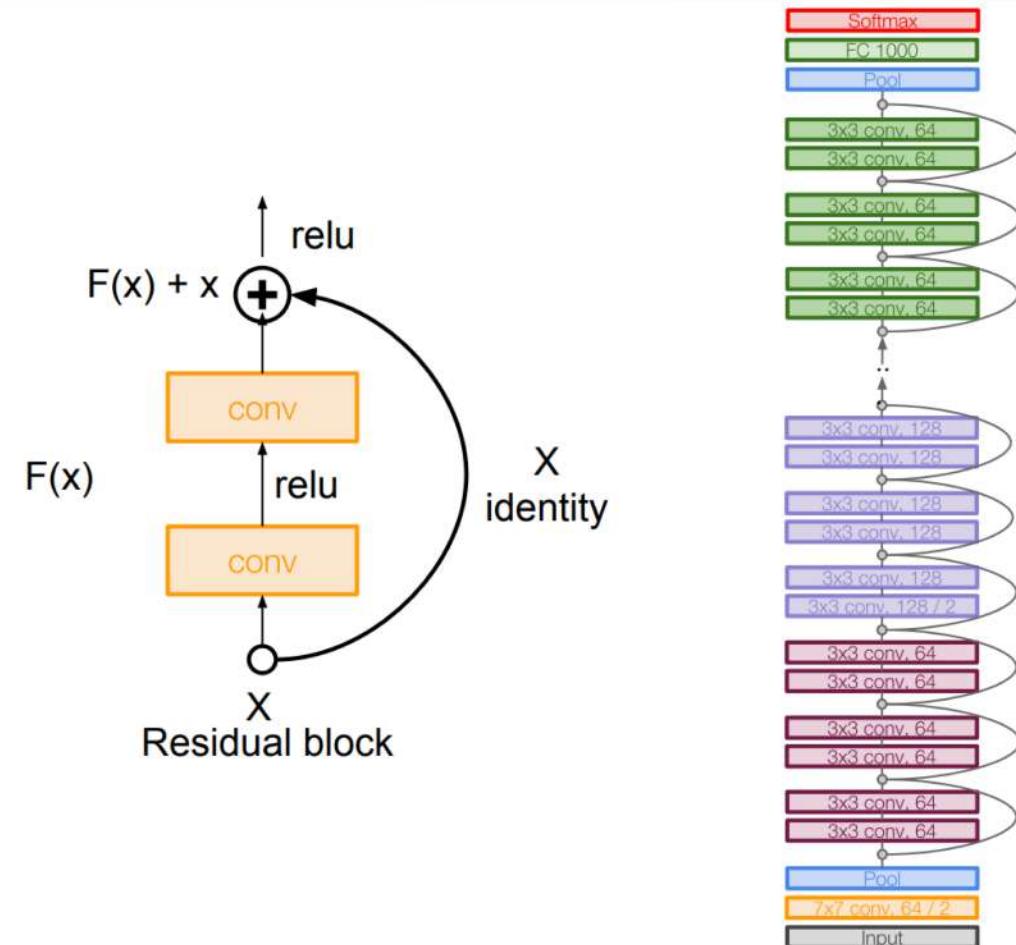
- CNN模型

Case Study: ResNet

[He et al., 2015]

Very deep networks using residual connections

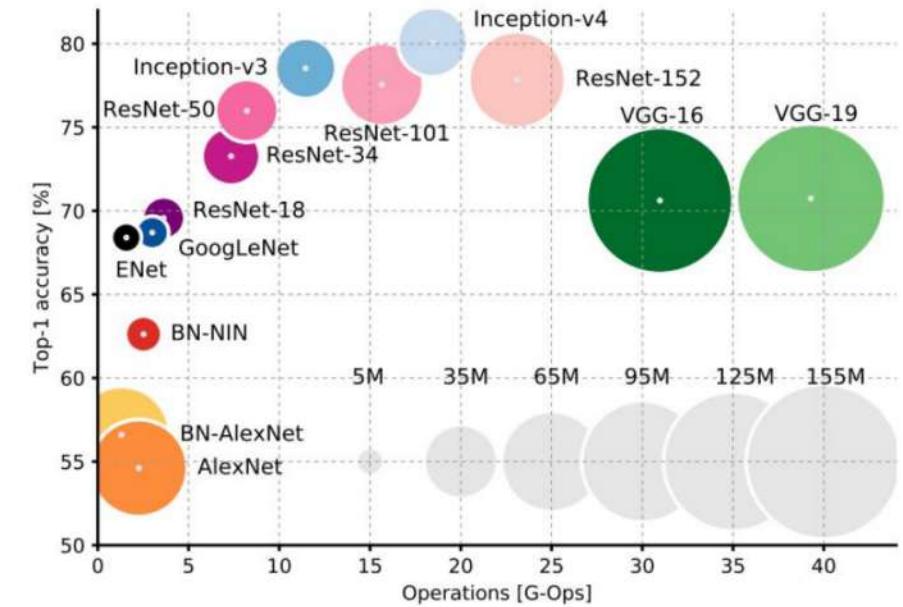
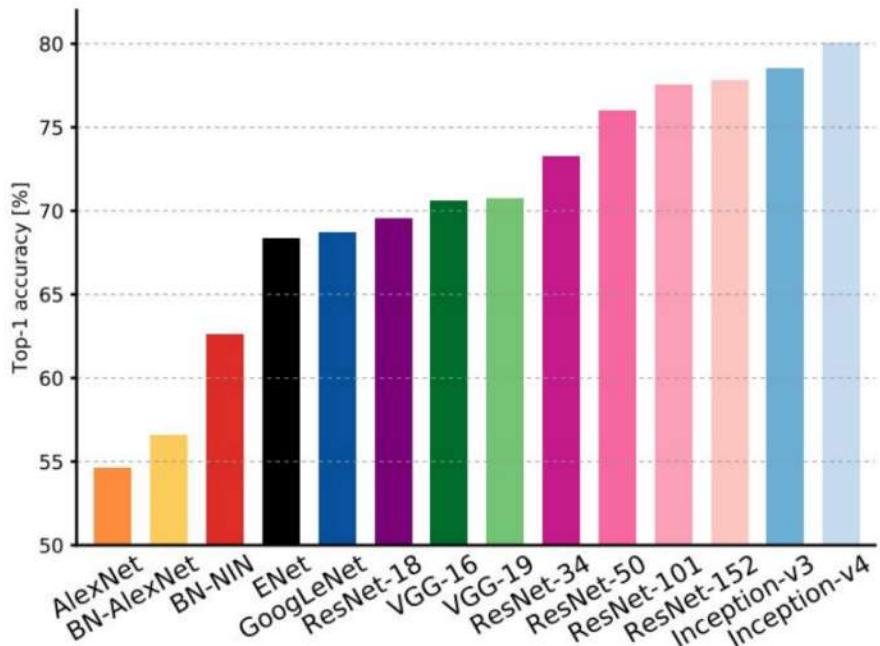
- 152-layer model for ImageNet
- ILSVRC'15 classification winner (3.57% top 5 error)
- Swept all classification and detection competitions in ILSVRC'15 and COCO'15!



单字符识别

Comparing complexity...

- CNN模型

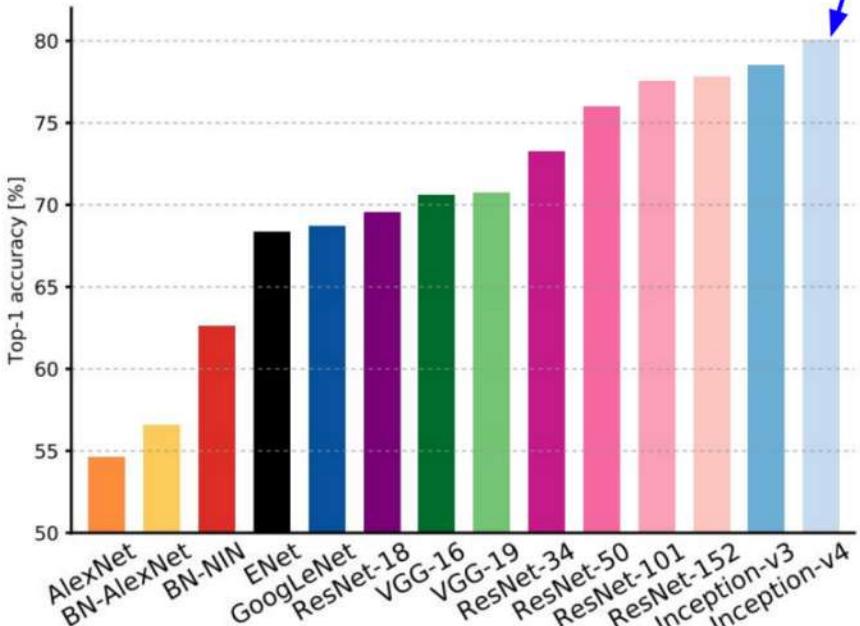


An Analysis of Deep Neural Network Models for Practical Applications, 2017.

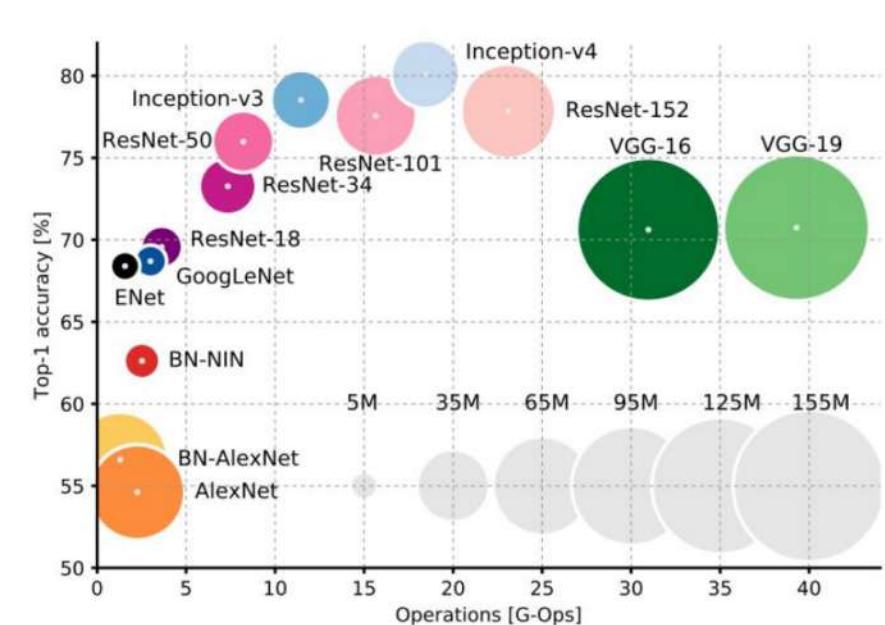
单字符识别

Comparing complexity...

- CNN模型



Inception-v4: Resnet + Inception!

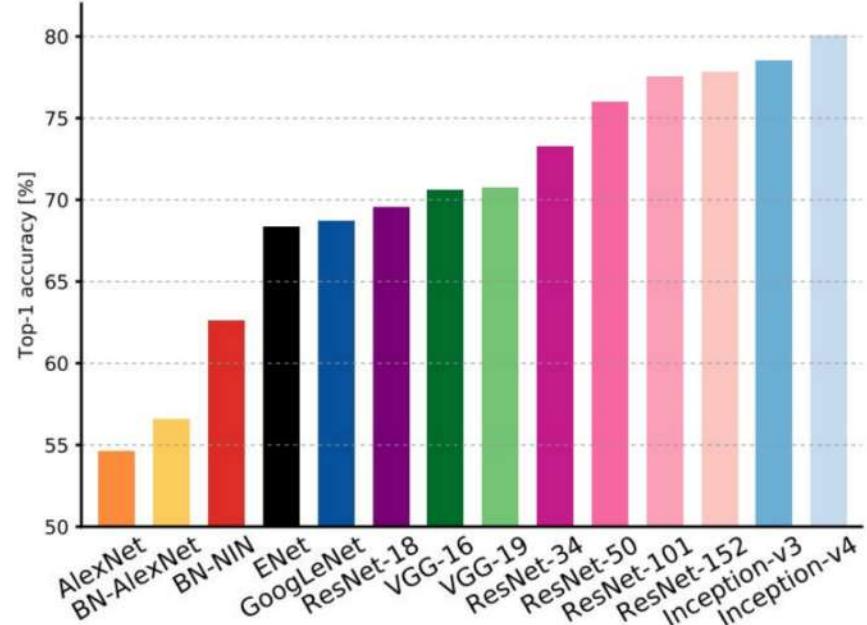


An Analysis of Deep Neural Network Models for Practical Applications, 2017.

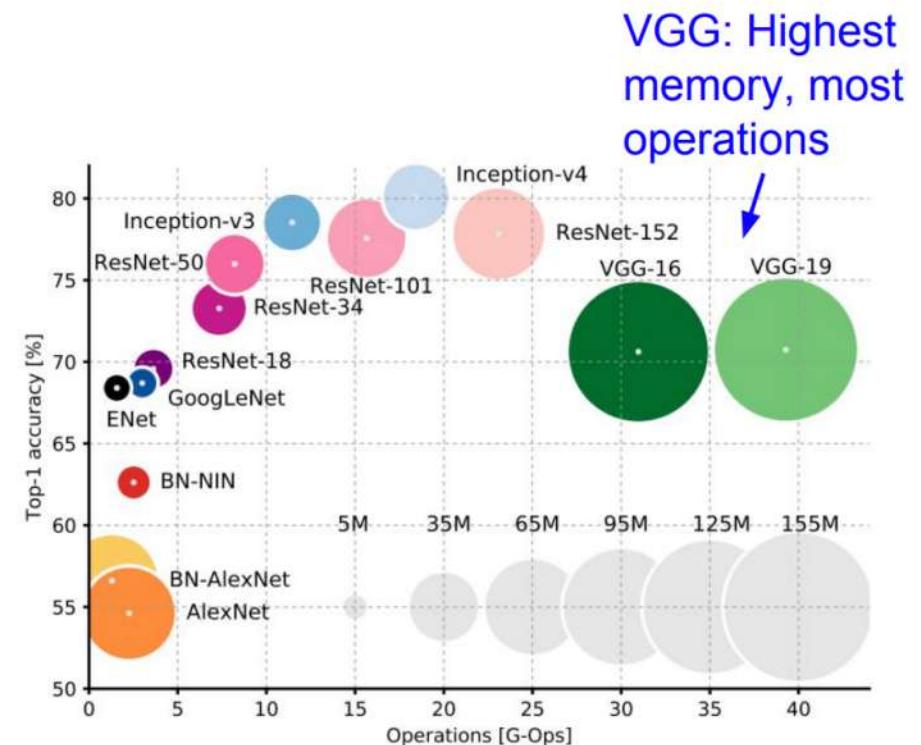
单字符识别

- CNN模型

Comparing complexity...



An Analysis of Deep Neural Network Models for Practical Applications, 2017.

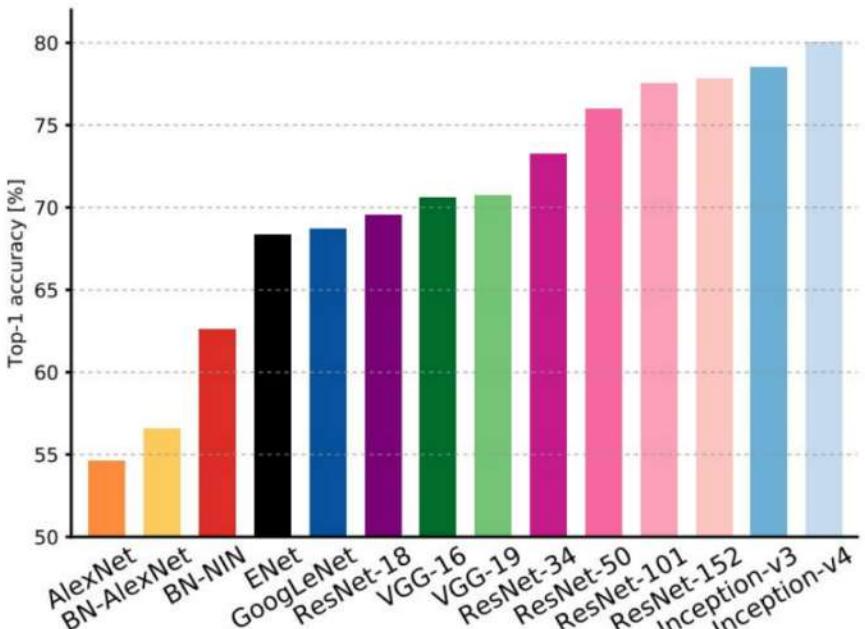


VGG: Highest
memory, most
operations

单字符识别

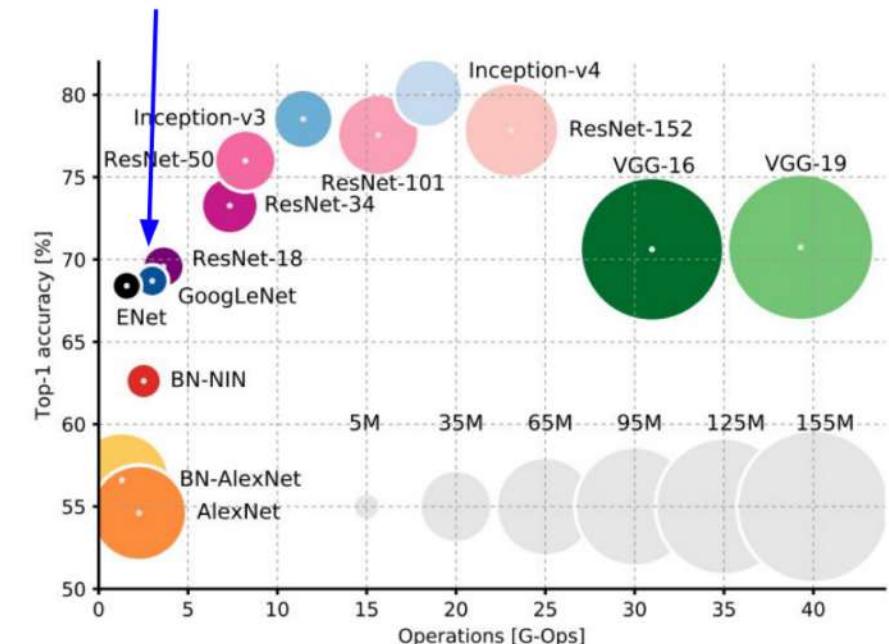
Comparing complexity...

- CNN模型



An Analysis of Deep Neural Network Models for Practical Applications, 2017.

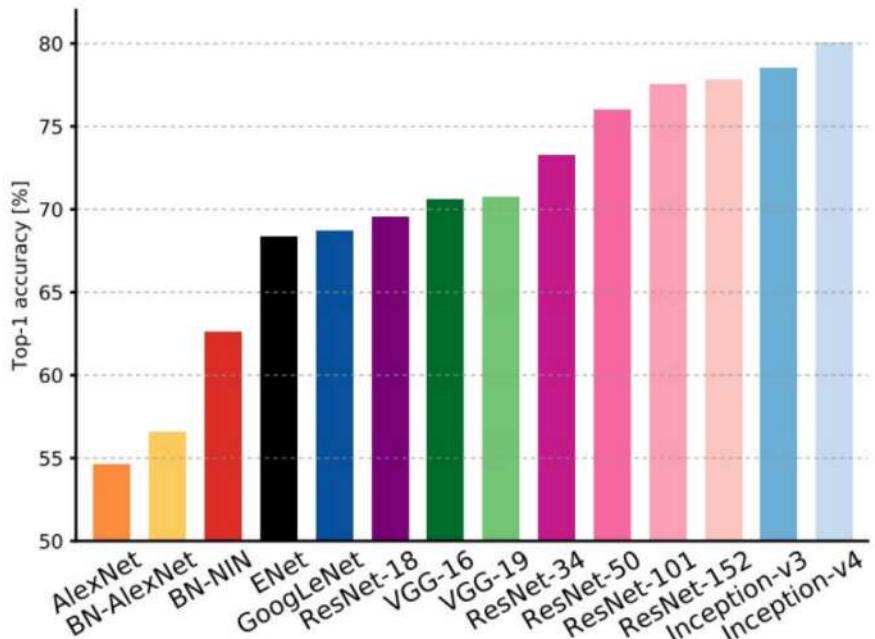
GoogLeNet:
most efficient



单字符识别

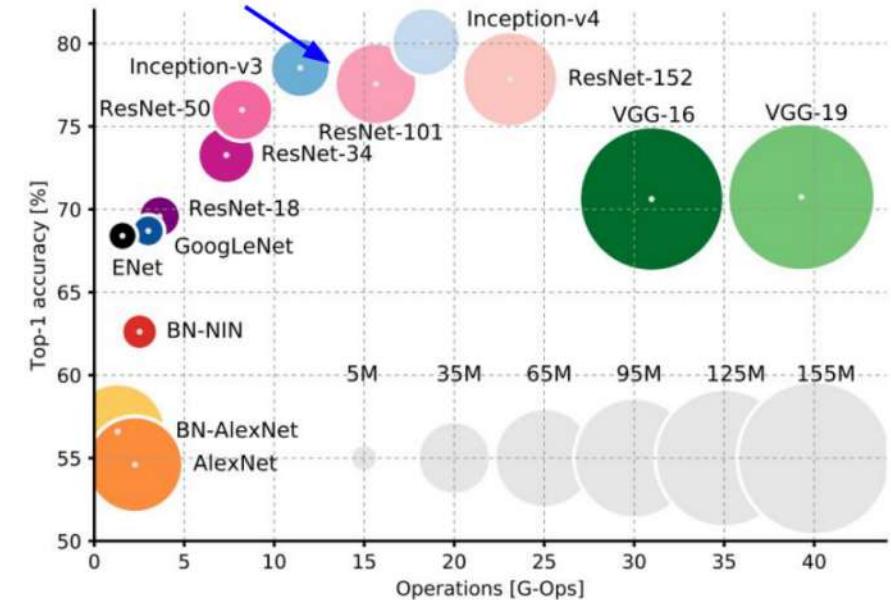
- CNN模型

Comparing complexity...



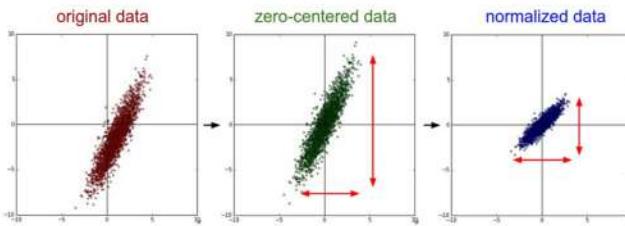
An Analysis of Deep Neural Network Models for Practical Applications, 2017.

ResNet:
Moderate efficiency depending on
model, highest accuracy



单字符识别

- CNN模型
 - 输入图片归一化
 - 数据增强
 - PADDING
 - 灰度图还是二值图？
 - 英文单词可以用CNN学？
宽高比变化范围大怎么办？



休息一下

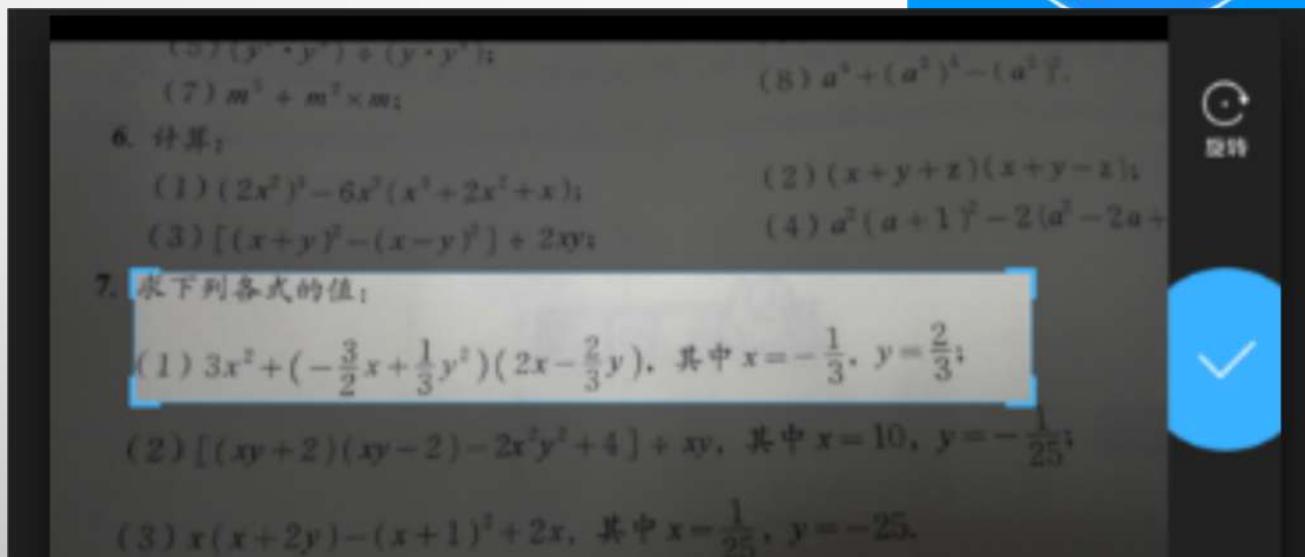
- DEEP DREAM
 - [HTTPS://WWW.YOUTUBE.COM/WATCH?V=BSSMBPMPEYQ&T=683S](https://www.youtube.com/watch?v=BSSMBPMPEYQ&t=683s)

第二课 单字符识别

- 单字符识别
 - 模版匹配
 - 特征提取 + 传统分类器
 - CNN模型
- 过切分的路径选择
 - VITERBI + BEAM SEARCH
- 语言模型
 - N-GRAM
 - LSTM
 - 2D-PCFG

实例一

• 通用文本识别



Wang Lin is a good 1. She is twenty-one years old. She drives a car in a 2. She 3 from Sunday to Friday. Her home isn't 4 the factory. She gets up 5 at six. She goes to work at 6:30, 6 7:30 she must get there. She has 7 in the factory. She has lunch there, too. She 8 the factory at 5:00 in the afternoon. She cooks and then does housework in the evening. She likes 9 very much. She 10 to bed at about 10:30 pm.

2015年02月04日 00:36

④ 搜索结果

完形填空。

Wang Lin is a good 1. She is twenty-one years old. She drives a car in a 2. She

3 from Sunday to Friday. Her home isn't 4 the factory. She gets up 5 at six

2. 下列句中的成语使用不正确的一项是 ()
- A. 任何不称职的或者愚蠢得不可救药的人，都看不见这衣服。
 - B. 这可骇人听闻了。难道我是一个愚蠢的人吗？
 - C. 不知道什么时候，出现了一个神通广大的女神，叫做女娲。
 - D. 接着一阵手舞足蹈的跳跃和欢呼，表示他获得生命的欢乐。

2015年02月04日 00:33

④ 搜索结果

题干

下列词语中加红词语运用不恰当的是

- A: 不知道什么时候，出现了一个神通广大^①的女神，叫做女娲。

(Sub) Text Line Extraction

八年级数学上(苏三册)

26. 已知, 点 C 是 $\angle MAN$ 平分线上一点, $\square BCD$ 的两边 CB 、 CD 分别与射线 AM 、 AN 于 B、D 两点, 且 $\angle BCD + \angle MAN = 180^\circ$. 过点 C 作 $CE \perp AB$, 垂足为 E.
- (1) 当点 E 在线段 AB 上时(如图 1), 求证: $\square NDX \Leftrightarrow \square BCE = 90^\circ$;
- (2) 当点 E 在线段 AB 的延长线上时(如图 2), 求证: $AB - AD = 2BE$;
- (3) 在(2)的条件下, 若 $\angle MAN = 60^\circ$, 连接 BD , 作 $\square ABD$ 的平分线 BF 交 AD 于点 O, 连接 DO 并延长交 AB 于点 G. $BC = 2$, $DF = 3$, 求线段 BD .

N

Character Segmentation

基于行片段拟合边界，再次对连通域进行划分



Character Segmentation

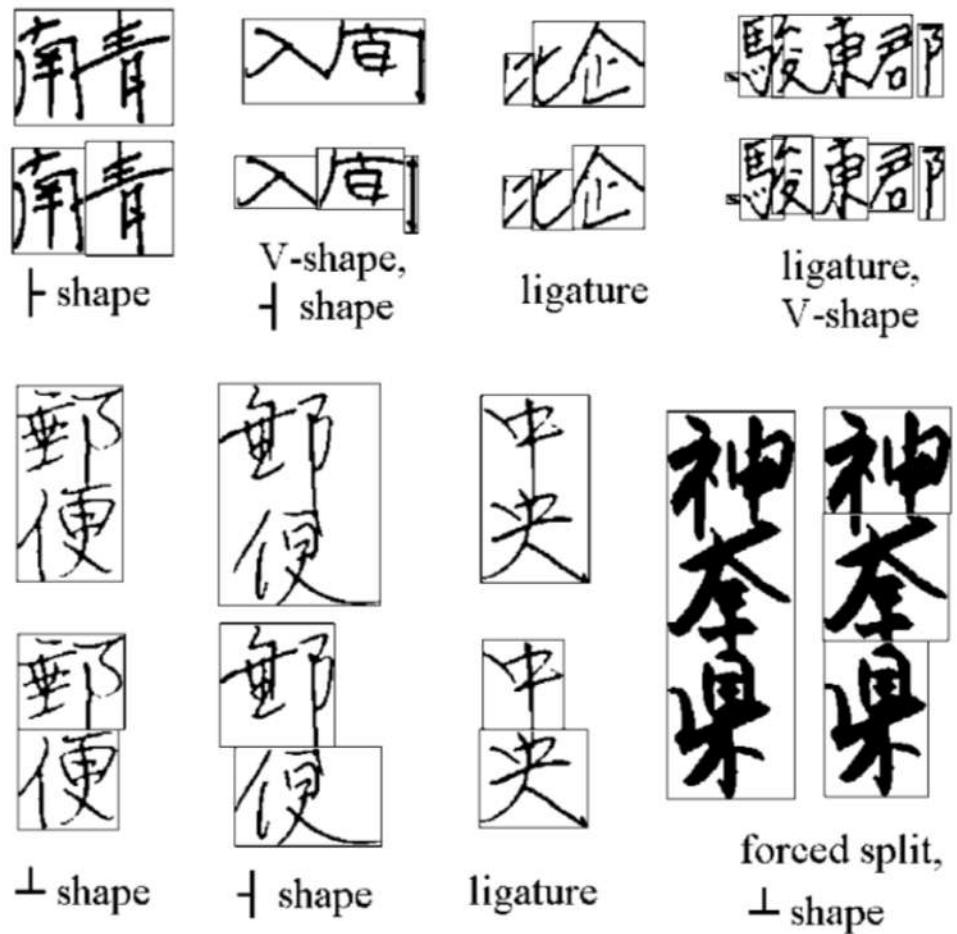
计算分割点置信度，在粗分割后，递归地以置信度最大点进行分割



投影法

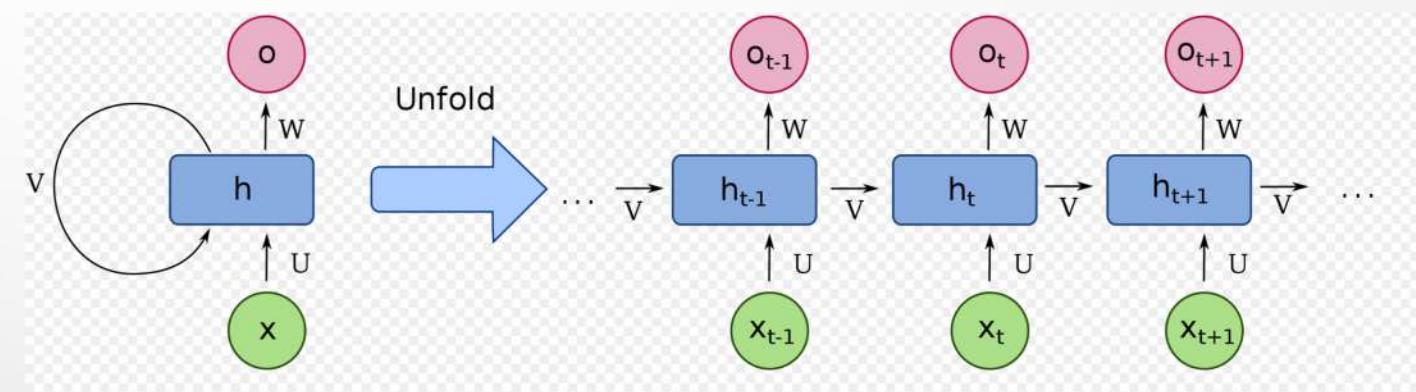
Character Segmentation

计算分割点置信度，在粗分割后，递归地以置信度最大点进行分割



Character Segmentation

计算分割点置信度，在粗分割后，递归地以置信度最大点进行分割



遗留问题

- 字符分割

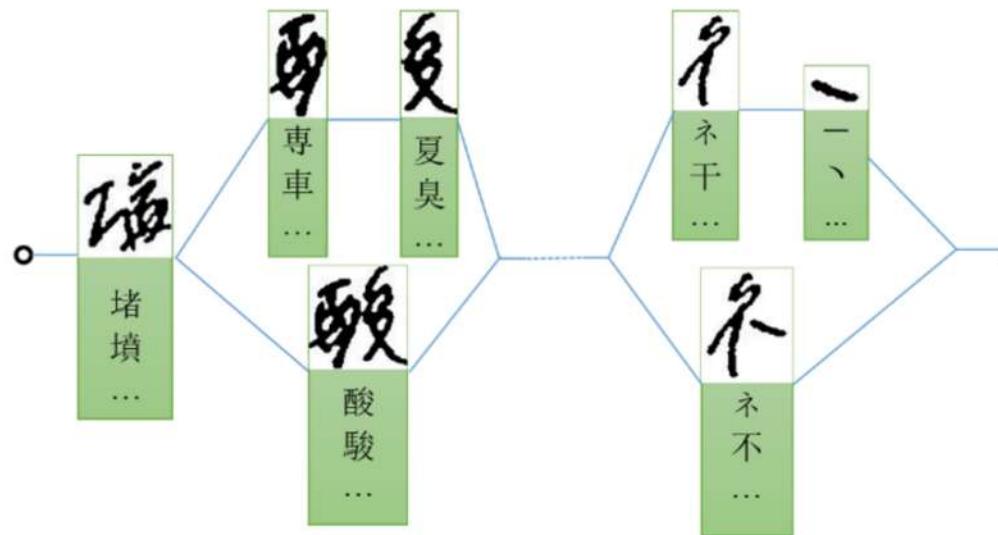
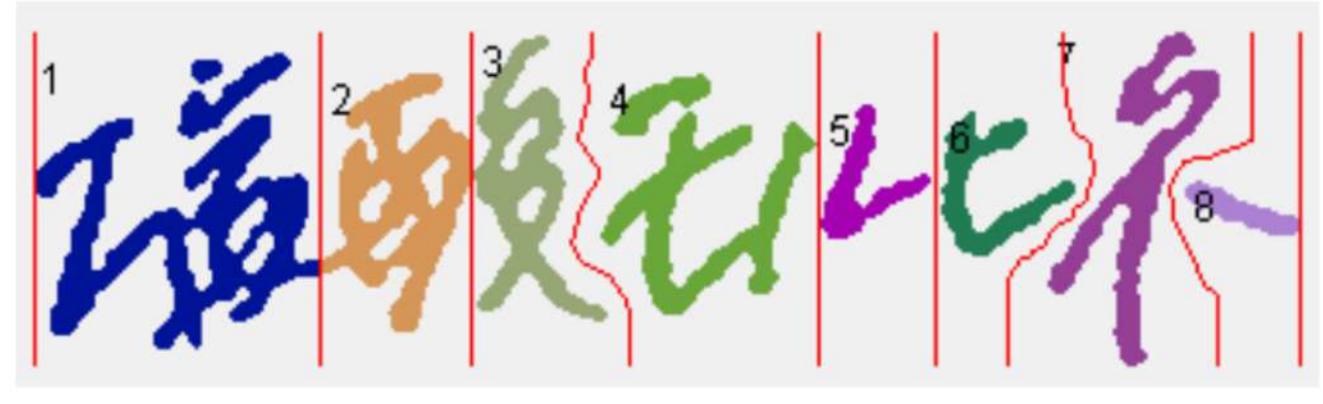


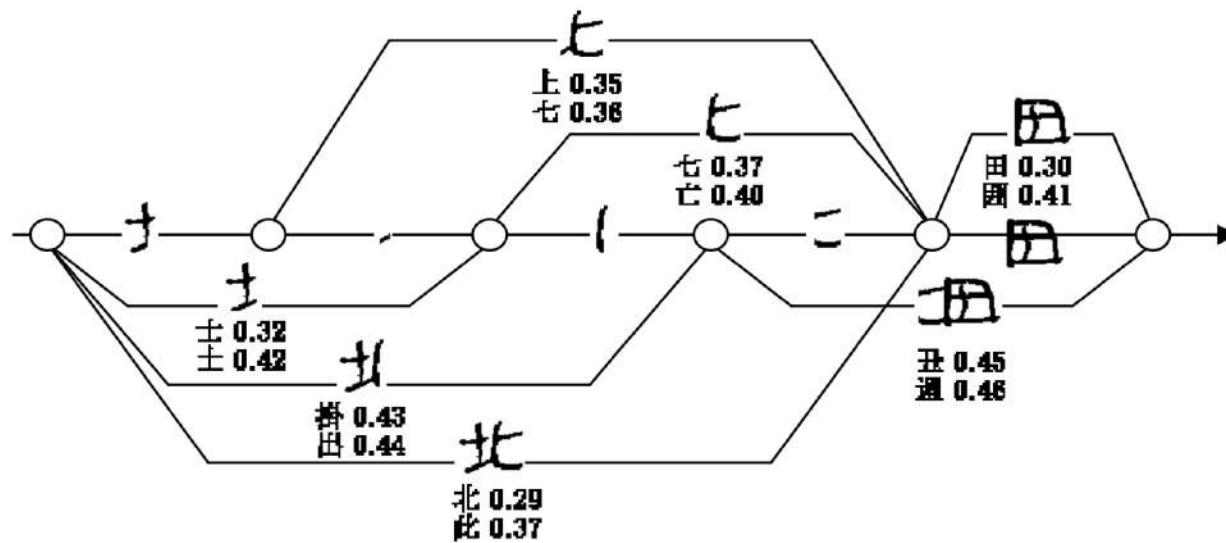
Fig. 11: Candidate lattice diagram.

路径选择

- VITERBI算法 + BEAM SEARCH

多步骤每步多选择模型的最优选择问题，其在每一步的所有选择都保存了前续所有步骤到当前步骤当前选择的最小总代价（或者最大价值）以及当前代价的情况下前继步骤的选择。

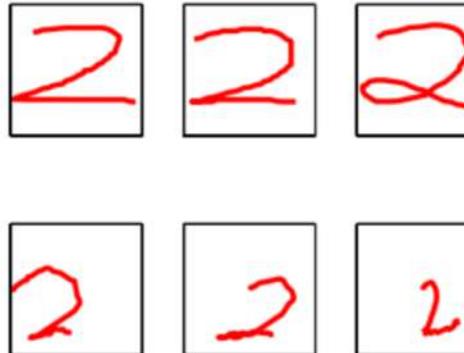
依次计算完所有步骤后，通过回溯的方法找到最优选择路径。



路径选择

- HMM + VITERBI算法

Top row: examples of on-line handwritten digits. Bottom row: synthetic digits sampled generatively from a left-to-right hidden Markov model that has been trained on a data set of 45 handwritten digits.



The most natural unit of handwriting is a letter. A letter is represented by a 7-state left-to-right HMM. The HMM model is illustrated in Figure 2-3.

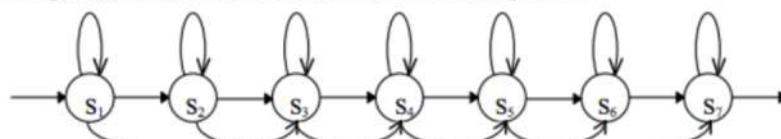


Figure 2-3: A 7-state HMM for a letter.

The left-to-right type of HMM, a special class of HMMs, have an additional property that the state index is non-decreasing as the time increases, i.e.

$$a_{ij} = P(Q_n=S_j|Q_{n-1}=S_i) = 0, \quad i > j. \quad (2.31)$$

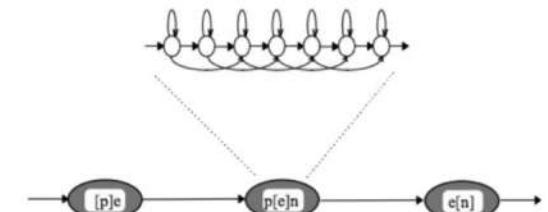


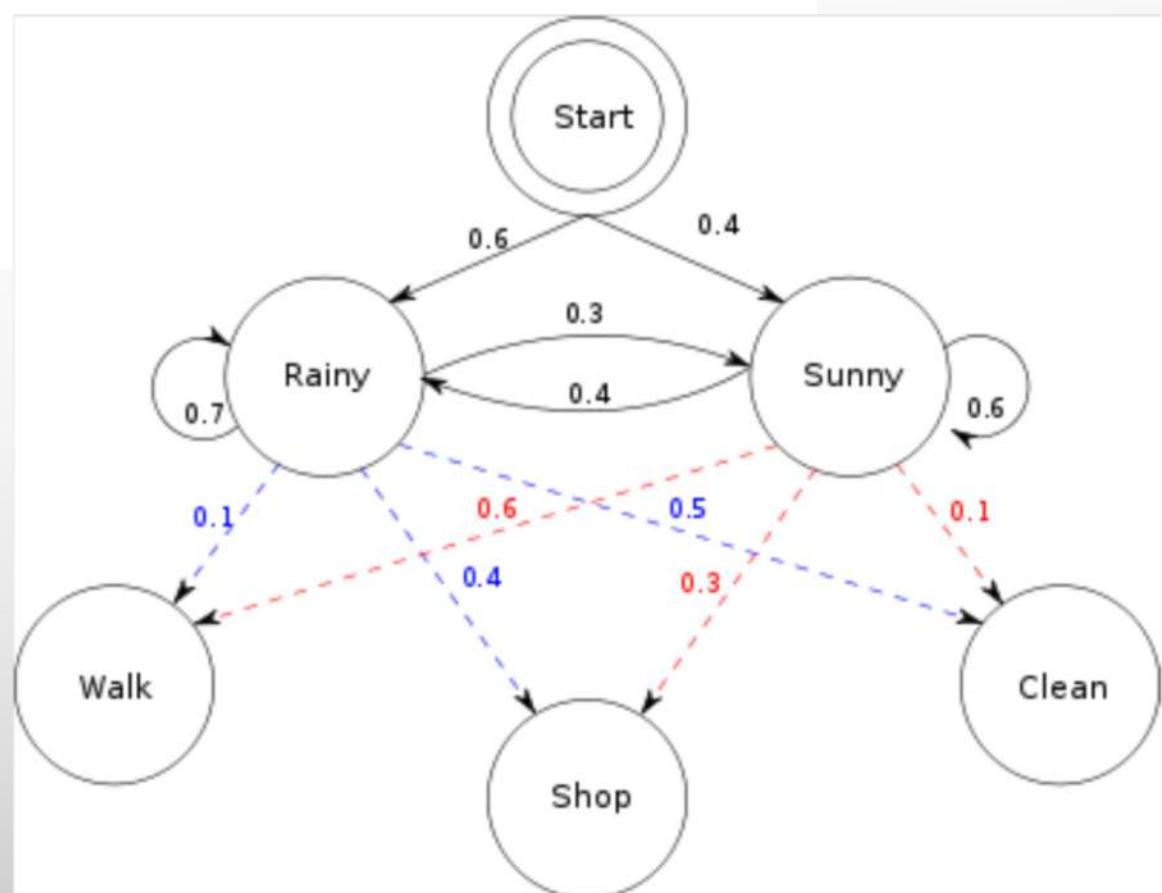
Figure 2-5: Illustration of an HMM modeling of the word "pen", consists of three 7-state HMMs for the trigraph "/p/e", "p/e/n", and "e/n".

HMM 模型

隐马尔可夫模型由初始状态概率向量 π 、状态转移概率矩阵 A 和观测概率矩阵 B 决定。 π 和 A 决定状态序列， B 决定观测序列。因此，隐马尔可夫模型 λ 可以用三元符号表示，即

$$\lambda = (A, B, \pi)$$

A, B, π 称为隐马尔可夫模型的三要素。

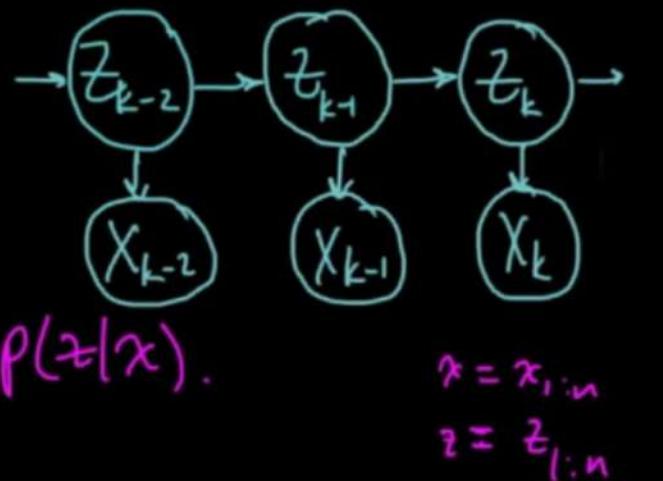


HMM + VITERBI

Viterbi Algorithm

Given: x_1, \dots, x_n
Assume dists.

Goal: Compute $z^* = \arg \max_z p(z|x)$.



Rk: If $f(a) \geq 0 \forall a$ and $g(ab) \geq 0 \forall a, b$,

$$\text{then } \max_{a,b} f(a) g(ab) = \max_a \left(f(a) \max_b g(a, b) \right).$$

HMM + VITERBI

Rk: If $f(a) \geq 0 \ \forall a$ and $g(a,b) \geq 0 \ \forall a,b$,

$$\text{then } \max_{a,b} f(a)g(a,b) = \max_a \left(f(a) \max_b g(a,b) \right).$$

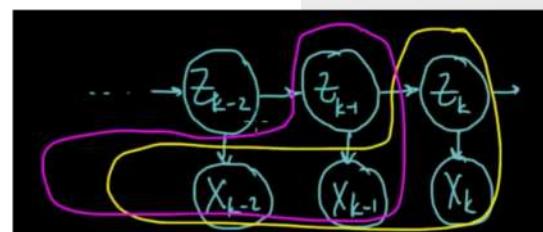
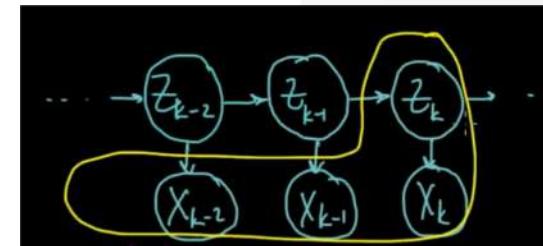
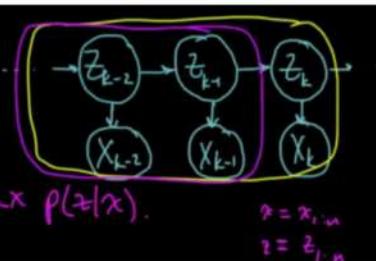
Note:
 $\arg\max_{z_{1:n}} p(z|x) = \arg\max_{z_{1:n}} p(z,x)$

$z = z_{1:n}$

Viterbi Algorithm

Given: x_1, \dots, x_n
 Assume dists.

Goal: Compute $z^* = \arg\max_z p(z|x)$.



$$\begin{aligned}
 \mu_k(z_k) &= \max_{z_{1:k-1}} p(z_{1:k}, x_{1:k}) \\
 &= \max_{z_{1:k-1}} \underbrace{p(x_k|z_k)}_{f(a)} \underbrace{p(z_k|z_{k-1})}_{g(a|b)} p(z_{1:k-1}, x_{1:k-1}) \\
 &= \max_{z_{k-1}} \left(\underbrace{p(x_k|z_k)}_{\text{known}} \underbrace{p(z_k|z_{k-1})}_{\text{known}} \max_{z_{1:k-2}} \underbrace{p(z_{1:k-1}, x_{1:k-1})}_{\mu_{k-1}(z_{k-1})} \right)
 \end{aligned}$$

HMM + VITERBI

$$= \max_{z_{k-1}} \left[p(x_k | z_k) p(z_k | z_{k-1}) \max_{\tilde{z}_{1:k-2}} p(\tilde{z}_{1:k-1} | x_{1:k-1}) \right] \\ \mu_{k-1}(z_{k-1})$$

$$\mu_k(z_k) = \max_{z_{k-1}} p(x_k | z_k) p(z_k | z_{k-1}) / \mu_{k-1}(z_{k-1}) \text{ for } k=2, \dots, n$$

$$\mu_1(z_1) = p(z_1, x_1) = p(z_1) p(x_1 | z_1)$$

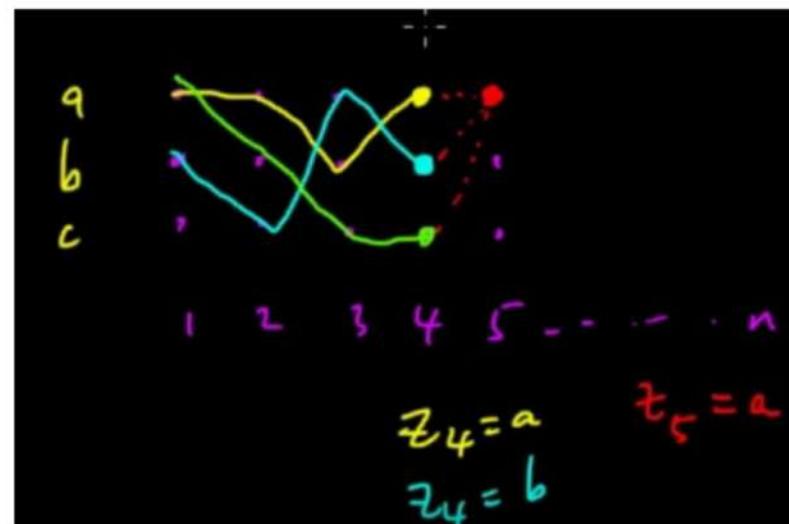
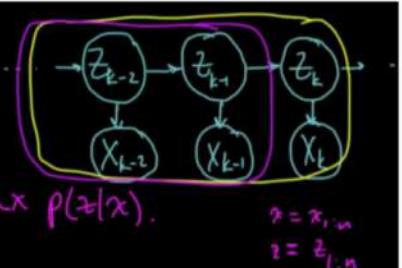
$$\mu_n(z_n) = \max_{z_{1:n-1}} p(x_{1:n}, z_{1:n})$$

$$\max_{z_n} \mu_n(z_n) = \max_{z_{1:n}} p(x_{1:n}, z_{1:n})$$

Viterbi Algorithm

Given: x_1, \dots, x_n
Assume dists.

Goal: Compute $z^* = \arg \max_{z} p(z|x)$.



第二课 单字符识别

- 单字符识别
 - 模版匹配
 - 特征提取 + 传统分类器
 - CNN模型
- 过切分的路径选择
 - VITERBI + BEAM SEARCH
- 语言模型
 - N-GRAM
 - LSTM
 - 2D-PCFG

当 $n=1$, 一个一元模型 (unigram model)即为 :

$$P(w_1, w_2, \dots, w_m) = \prod_{i=1}^m P(w_i)$$

- N-GRAM

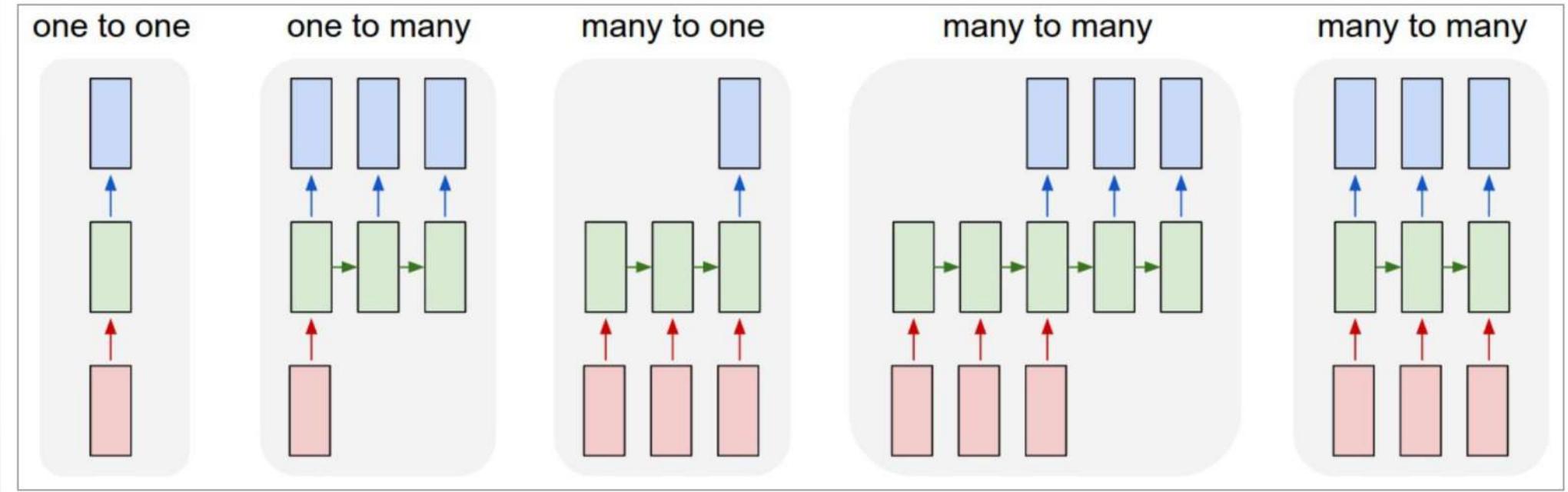
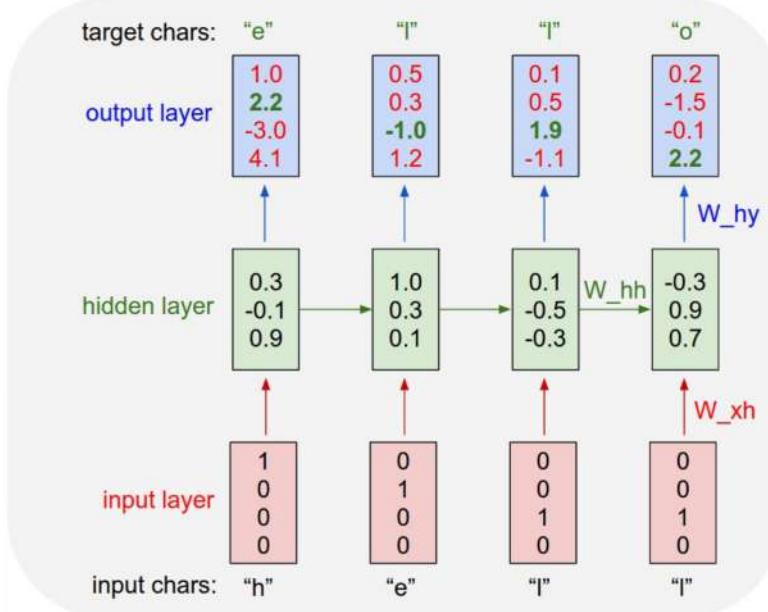
当 $n=2$, 一个二元模型 (bigram model)即为 :

$$P(w_1, w_2, \dots, w_m) = \prod_{i=1}^m P(w_i | w_{i-1})$$

当 $n=3$, 一个三元模型 (trigram model)即为

$$P(w_1, w_2, \dots, w_m) = \prod_{i=1}^m P(w_i | w_{i-2} w_{i-1})$$

- LSTM



- CFG (CONTEXT FREE GRAMMAR)

上下文无关语法 (Context-Free Grammars) :

一个上下文无关语法是一个四元组2:

$$G = (N, \Sigma, R, S)$$

其中N是一组非终结符(non-terminal symbols)

Σ 是一组终结符(terminal symbols)

R是一系列的变换规则, 即是:

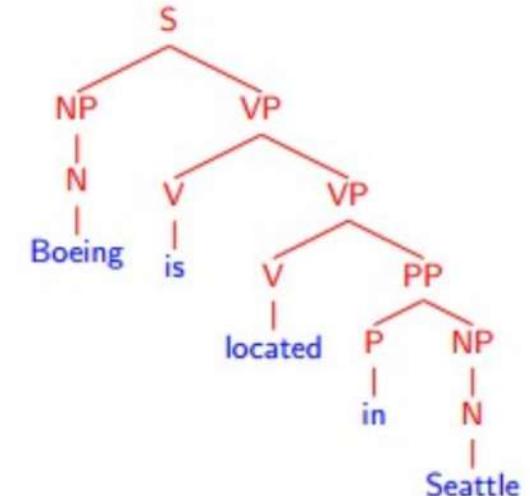
$$X \rightarrow Y_1 Y_2 \dots Y_n \quad \text{for } n \geq 0, X \in N, Y_i \in (N \cup \Sigma)$$

S是一个特殊的开始符号

INPUT:

Boeing is located in Seattle.

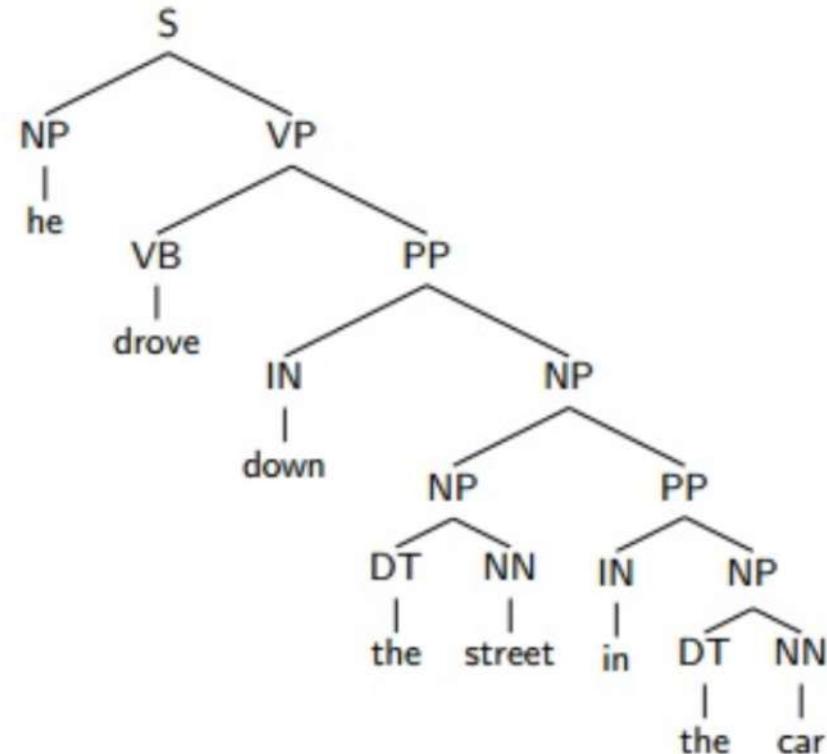
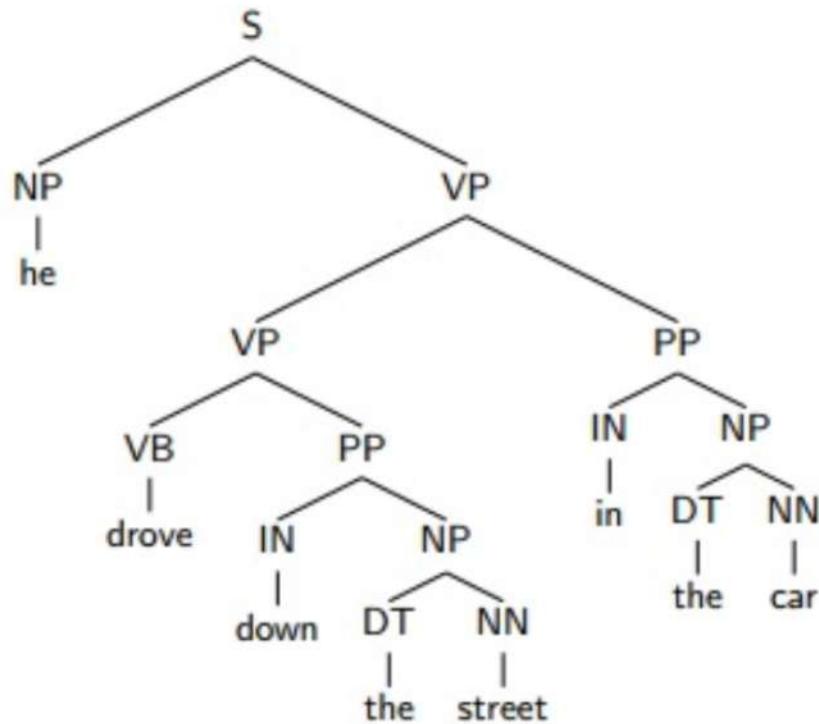
OUTPUT:



Noun Phrases (NP) Verb Phrases (VP) Sentences (S)
 (N = noun, V = verb, DT = determiner)

CFG (context free grammar)

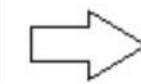
Ambiguity



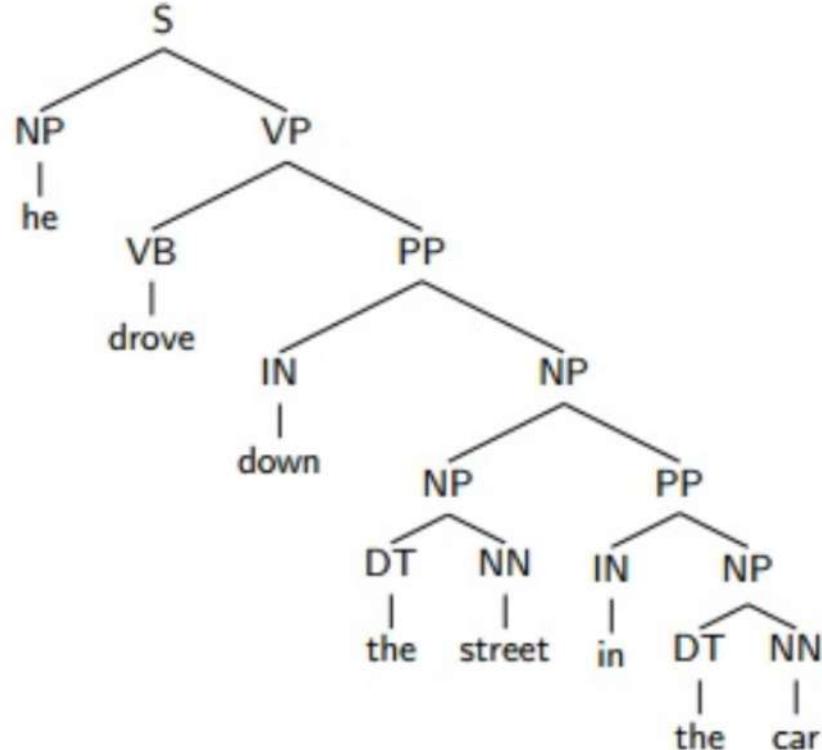
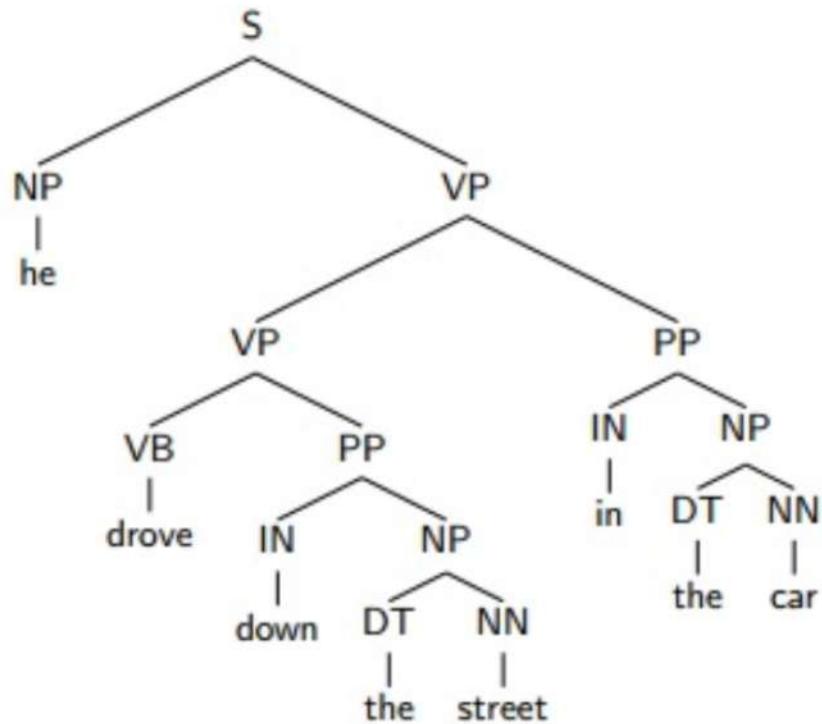
PCFG (probabilistic context free grammar)

Ambiguity

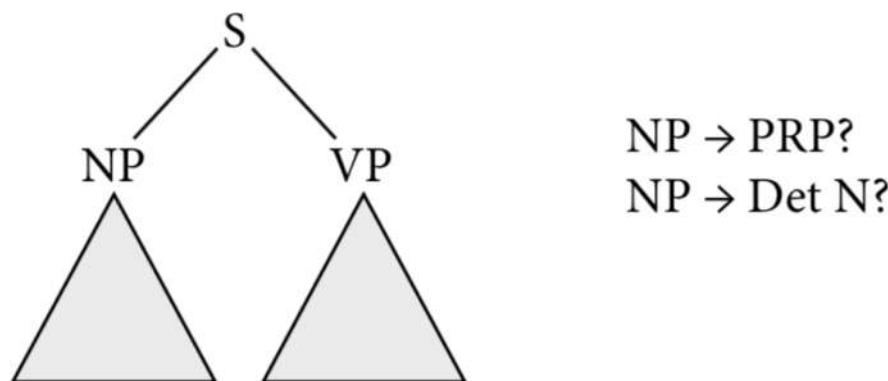
$S \Rightarrow NP VP$	
$VP \Rightarrow Vi$	
$VP \Rightarrow Vt NP$	
$VP \Rightarrow VP PP$	
$NP \Rightarrow DT NN$	
$NP \Rightarrow NP PP$	
$PP \Rightarrow P NP$	



$S \Rightarrow NP VP$	1.0
$VP \Rightarrow Vi$	0.4
$VP \Rightarrow Vt NP$	0.4
$VP \Rightarrow VP PP$	0.2
$NP \Rightarrow DT NN$	0.3
$NP \Rightarrow NP PP$	0.7
$PP \Rightarrow P NP$	1.0



- Context-free grammar: One rule can only “see” parent and its children, not anything above or below.
- PCFG: Assumes statistical independence of all rewrite events.



$$a_0 = x^2 + \frac{\sqrt{\pi}}{3}$$

- 2D-PCFG

$$a_0 = x^2 + \frac{\sqrt{\pi}}{3}$$

$$(o_1 X^{o_2} + o_3 Y^{o_4})^{o_6} {}^{o_7}$$

$$\|f\|^2_\alpha=\frac{\Gamma(n+\alpha)}{2^n\pi^n\Gamma(\alpha)}\int_D(-\rho)^\alpha|f|^2\,\mathrm{d}V$$

$$k^*(w') \geq k^*(w_0)-\varepsilon\}\geq \frac{3\pi^n\delta_1^{2n-1}}{2\sqrt{n}4^{2n-1}}$$

$$\sum_{i=0}^{k-1} (\alpha_i-\beta_i)\bar{\delta}_i=(\beta_k-\alpha_k)\bar{\delta}_k.\prod_{p\in P}C(p),$$

$$\sigma(T)(\subset \sigma(A)) \hspace{1cm} (\sigma=i_{k-1}, 2 \leqslant$$

$$\int e^{x^2}x^3dx \quad \int e^{x^2}x^3dx \quad \int e^{x^2}x^3dx$$

$$\sum_{n=1}^N\left(-1\right) ^n\sin (nx)\qquad\qquad\sum_{n=1}^N\left(-1\right) ^n\sin \left(nx\right)$$

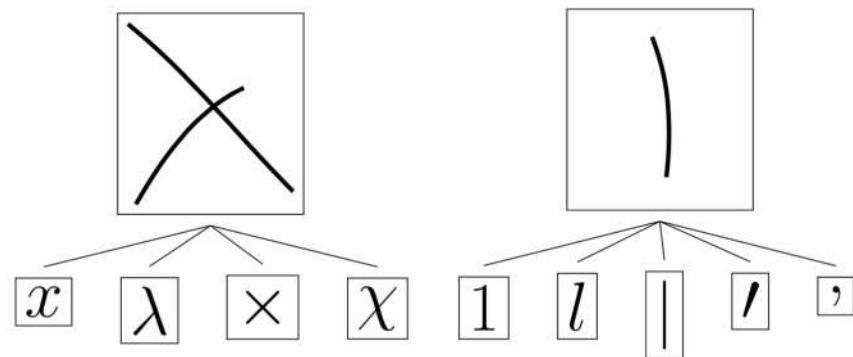
$$\frac{\sin{(x)}}{\cos{(x)}} = \tan{(x)}$$

$$\sqrt{q+j}-\left[z^{\omega \infty }\right]$$

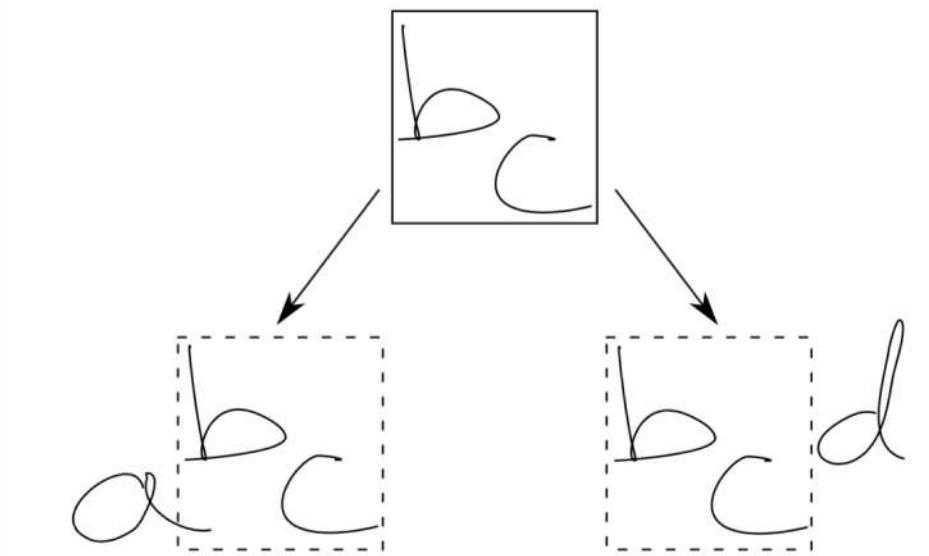
$$\sqrt{\frac{2z^3}{\sqrt{4z}}}$$

$$\begin{array}{l} \text{(a)} \quad \pm \frac{2}{5} \\ \text{(b)} \quad |-| < x \end{array}$$

$$\begin{array}{cccc} i & x^2 & Y & d \\ \text{(a)} & \text{(b)} & \text{(c)} & \text{(d)} \end{array}$$



$$A \times B = \emptyset$$



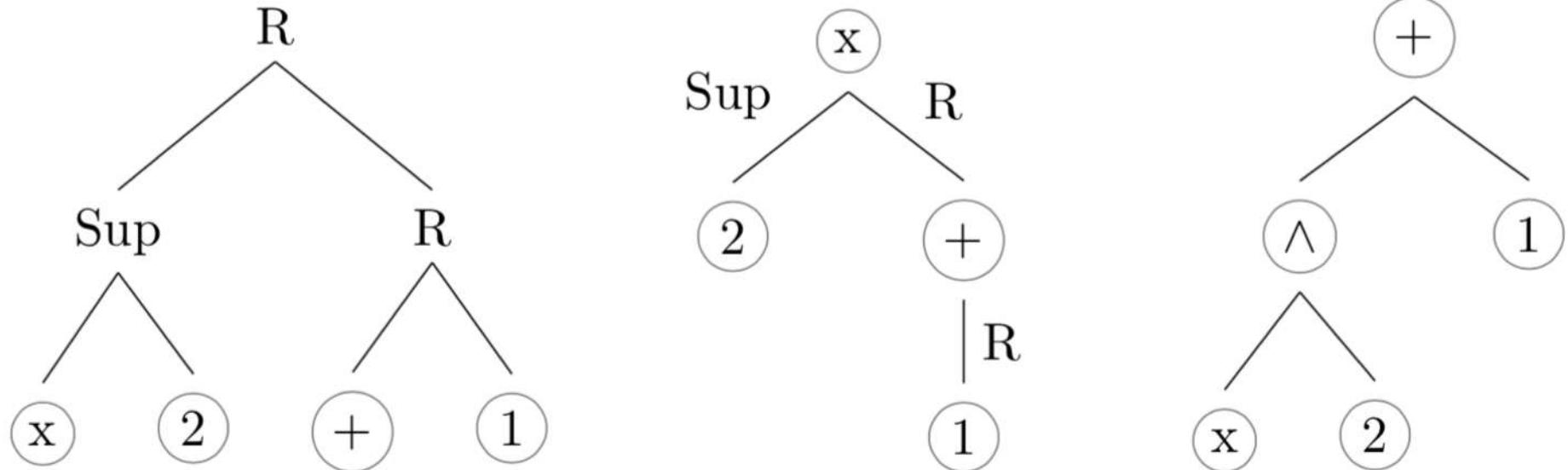


Figure 1.8: Example of tree representation of the expression $x^2 + 1$.
Left to right: relational tree, symbol layout tree, and operator tree.

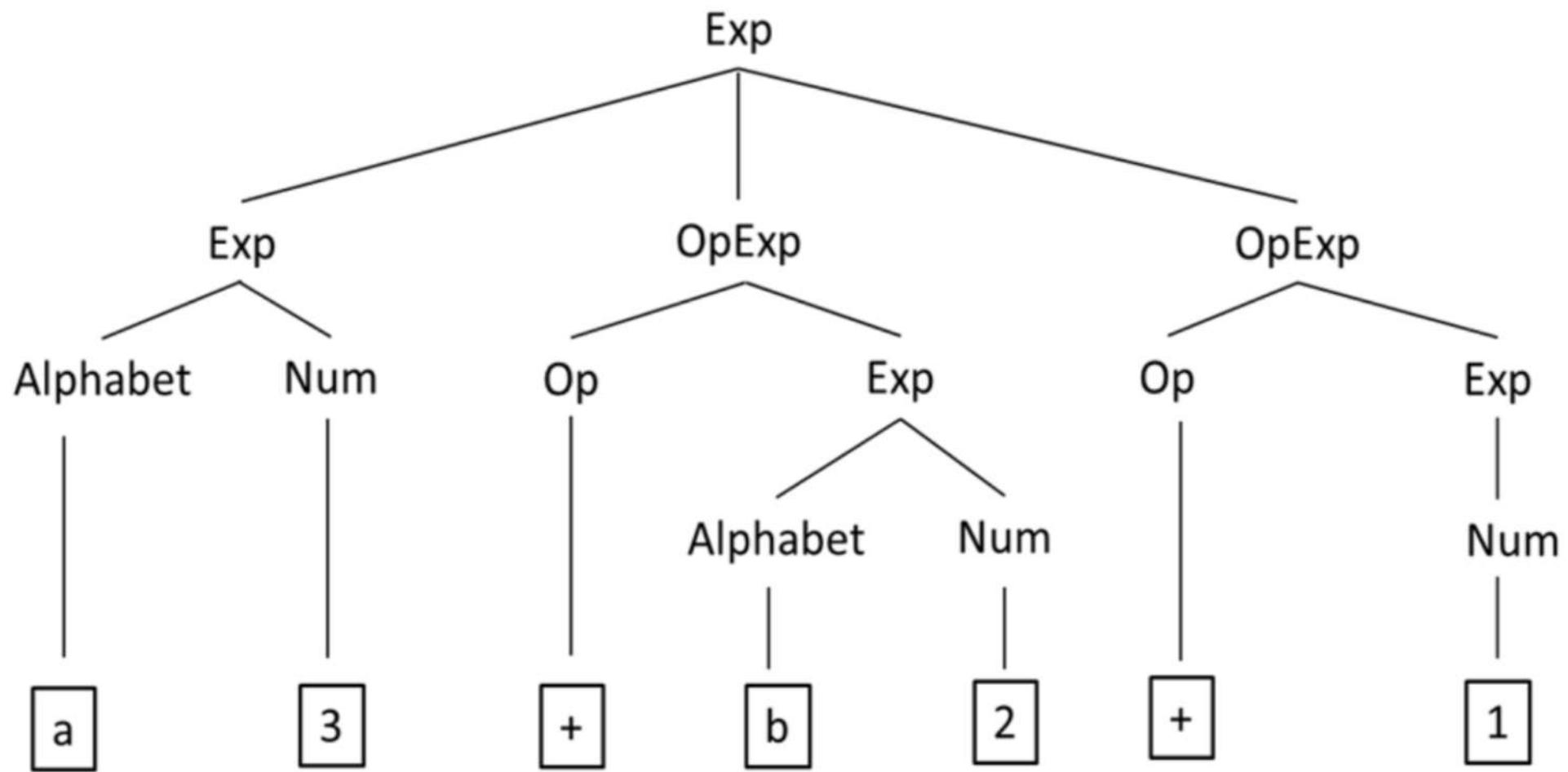
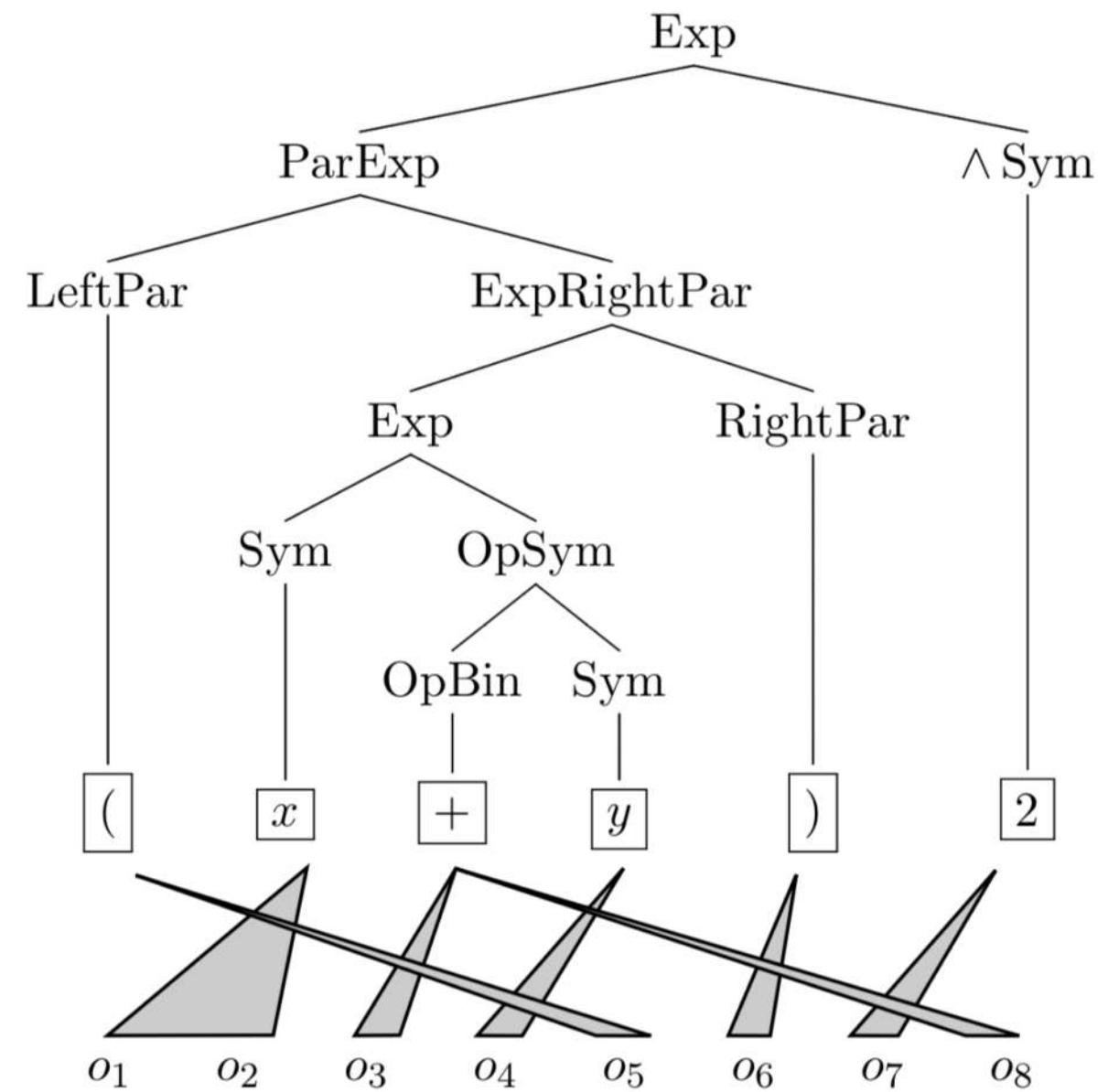


Figure 2.4 Parsing tree for $a^3 + b^2 + 1$

- 2D-PCFG

$$(o_1 \times o_2 + o_3 \times o_4)^2$$

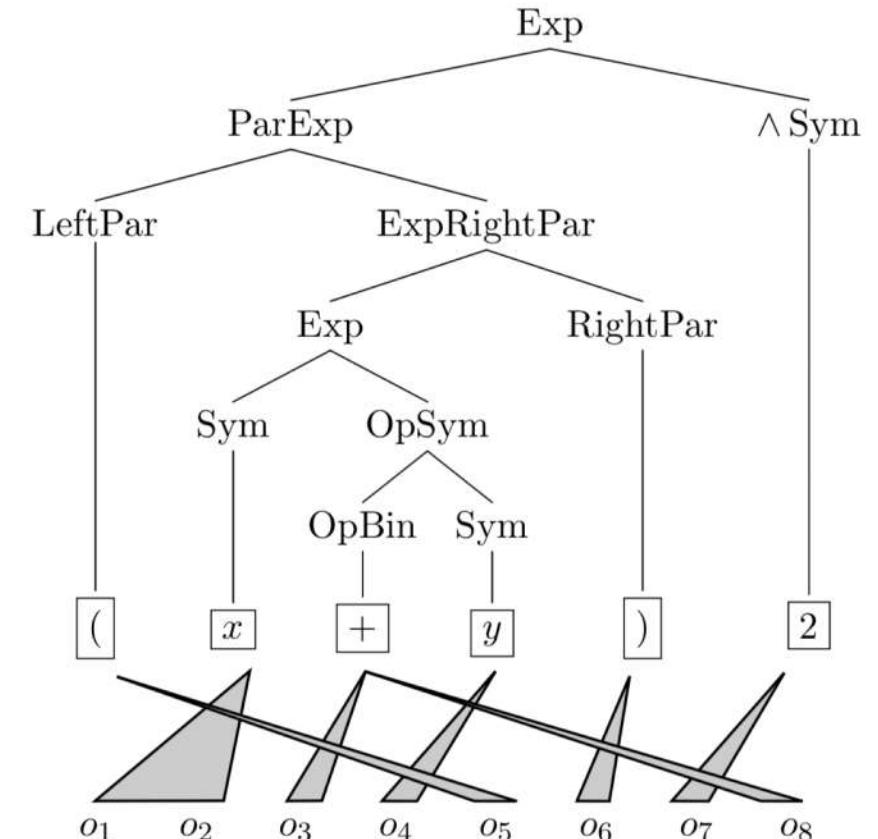
(o₅ o₂ o₈ o₃ + o₄)² o₇



Definition 2.4.1. A *Context-Free Grammar* (CFG) G is a 4-tuple $(\mathcal{N}, \Sigma, S, \mathcal{P})$ where \mathcal{N} is a finite set of nonterminal symbols, Σ is a finite set of terminal symbols ($\mathcal{N} \cap \Sigma = \emptyset$), $S \in \mathcal{N}$ is the start symbol of the grammar, and \mathcal{P} is a finite set of rules: $A \rightarrow \alpha$, $A \in \mathcal{N}$, $\alpha \in (\mathcal{N} \cup \Sigma)^+$.

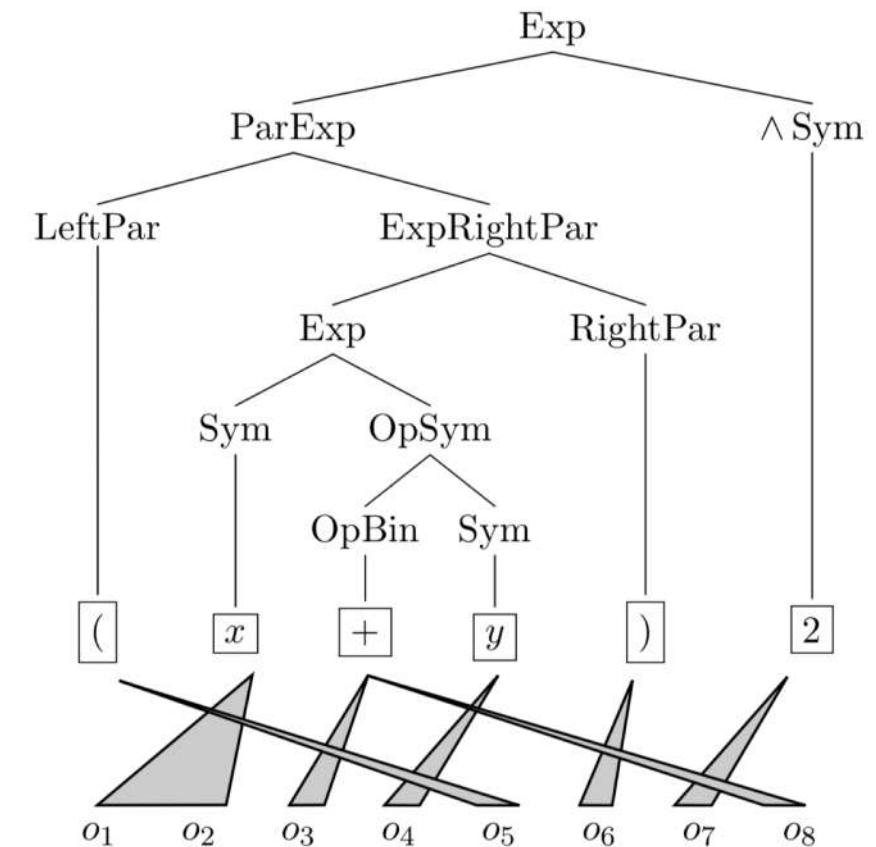
A CFG in Chomsky Normal Form (CNF) is a CFG in which the rules are of the form $A \rightarrow BC$ or $A \rightarrow a$ (where $A, B, C \in \mathcal{N}$ and $a \in \Sigma$). □

$$(o_5 X^{o_2} + o_3 Y^{o_4})^{o_6} {}^{o_7} 2^{o_8}$$



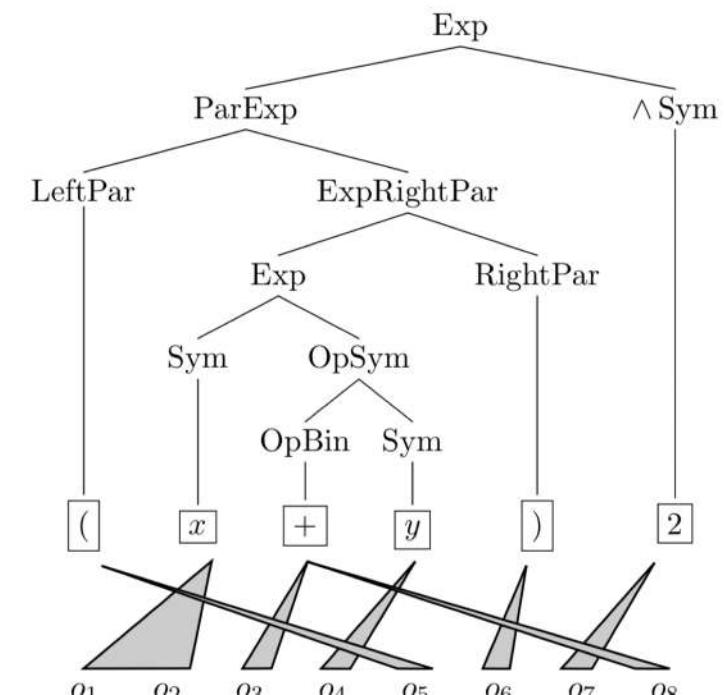
Definition 2.4.2. A *Probabilistic Context-Free Grammar* (PCFG) \mathcal{G} is defined as a pair (G, p) , where G is a CFG and $p : P \rightarrow]0, 1]$ is a probability function of rule application such that $\forall A \in \mathcal{N} : \sum_{i=1}^{n_A} p(A \rightarrow \alpha_i) = 1$, where n_A is the number of rules associated with nonterminal symbol A .

$$\left(\begin{matrix} o_5 \\ o_1 \times o_2 \end{matrix} + \begin{matrix} o_3 \\ o_8 \end{matrix} + \begin{matrix} o_4 \\ y \end{matrix} \right)^2 o_6 o_7$$



Definition 2.4.3. A *Two-Dimensional Probabilistic Context-Free Grammar* (2D-PCFG) is a generalization of a PCFG, where terminal and nonterminal symbols describe two-dimensional regions. This grammar in Chomsky Normal Form results in two types of rules: terminal rules and binary rules. First, the terminal rules $A \rightarrow a$ represent the mathematical symbols which are ultimately the terminal symbols of 2D-PCFG. Second, the binary rules $A \xrightarrow{r} BC$ have an additional parameter r that represents a given spatial relationship, and its interpretation is that regions B and C must be spatially arranged according to the spatial relationship r .

$$\left(\begin{matrix} o_5 & o_2 \\ o_1 & o_3 \end{matrix} + \begin{matrix} o_8 \\ o_4 \end{matrix} \right)^2 o_6$$



```

0
7 #Non-terminals
3 Exp
9 Sym
0 ExpOp
1 OpUn
2 ROpUn
3 OBExp
4 OpBin
5 OverExp
5 Over
7 OverSym
3 LeftPar
9 RightPar
0 RPExp
1 SSExp
2 BigOpExp
3 BigOp
4 Sqrt
5 Func
5 2Let
7 Let
3 SupSym
9 Num
0 Point
1 Dec

```

#Initial symbol of the grammar

START

Exp
Sym

#Terminal productions

PTERM

#S	->	[t1, t2, ... tn]
Sym		symbols.term
Let		letters.term
Num		numbers.term
Point		point.term
Over		over.term
BigOp		bigop.term
OpUn		opun.term
ROpUn		ropun.term
OpBin		opbin.term
OverSym		oversym.term
SupSym		supsym.s term
LeftPar		leftpar.term
RightPar		rightpar.term
Sqrt		sqrt.term

#Binary productions

PBIN

#pr	t	S	->	A	B	TeX-out
1.0	Sup	Exp		ExpOp	SupSym	"{\$1}^{\$2}"
1.0	Sup	Exp		Sym	SupSym	"{\$1}^{\$2}"
1.0	Sub	Exp		BigOp	Exp	"\$1_{-}{\$2}"
1.0	Sub	Exp		BigOp	Sym	"\$1_{-}{\$2}"
1.0	H	Exp		Exp	Exp	"\$1 \$2"
1.0	H	Exp		Exp	Sym	"\$1 \$2"
1.0	H	Exp		Sym	Exp	"\$1 \$2"
1.0	H	Exp		Sym	Sym	"\$1 \$2"
1.0	H	Exp		OpUn	Exp	"\$1 \$2"
1.0	H	Exp		OpUn	Sym	"\$1 \$2"
1.0	H	Exp		Exp	ROpUn	"\$1 \$2"
1.0	H	Exp		Sym	ROpUn	"\$1 \$2"
1.0	H	Exp		Exp	OBExp	"\$1 \$2"
1.0	H	Exp		Sym	OBExp	"\$1 \$2"
1.0	H	OBExp		OpBin	Exp	"\$1 \$2"
1.0	H	OBExp		OpBin	Sym	"\$1 \$2"
1.0	Sup	Exp		ExpOp	Exp	"{\$1}^{\$2}"
1.0	Sup	Exp		ExpOp	Sym	"{\$1}^{\$2}"
1.0	Sup	Exp		Sym	Exp	"{\$1}^{\$2}"
1.0	Sup	Exp		Sym	Sym	"{\$1}^{\$2}"
1.0	Sub	Exp		ExpOp	Exp	"{\$1}_{-}{\$2}"
1.0	Sub	Exp		Sym	Exp	"{\$1}_{-}{\$2}"
1.0	Sub	Exp		ExpOp	Sym	"{\$1}_{-}{\$2}"
1.0	Sub	Exp		Sym	Sym	"{\$1}_{-}{\$2}"
1.0	V	Exp		Exp	OverExp	"\frac{\$1}{\$2}"
1.0	V	Exp		Sym	OverExp	"\frac{\$1}{\$2}"
1.0	V	OverExp		Over	Exp	"\$2"
1.0	V	OverExp		Over	Sym	"\$2"
1.0	Vs	Exp		OverSym	Exp	"\$1{\$2}"
1.0	Vs	ExpOp		OverSym	Exp	"\$1{\$2}"
1.0	Vs	Exp		OverSym	Sym	"\$1{\$2}"
1.0	Vs	ExpOp		OverSym	Sym	"\$1{\$2}"
1.0	Vs	Exp		OverSym	Exp	"\$1{\$2}"
1.0	Vs	Exp		OverSym	Sym	"\$1{\$2}"

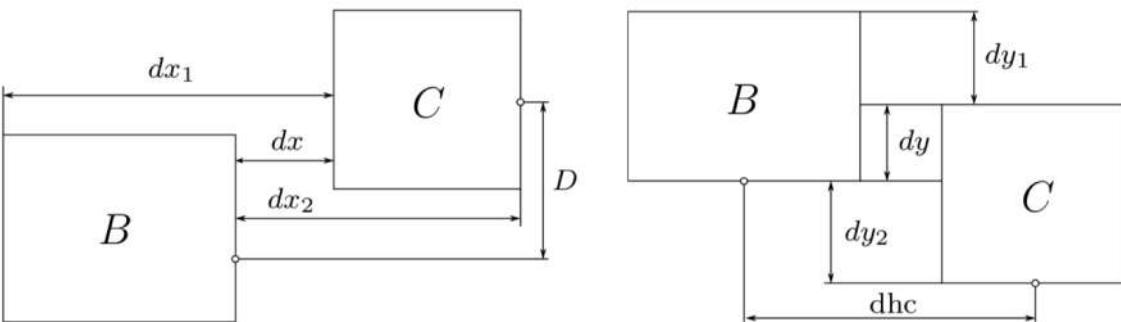
$$\frac{2}{5}$$

Fraction $\xrightarrow{\text{below}}$ Expression OverExp
OverExp $\xrightarrow{\text{below}}$ HorzLine Expression

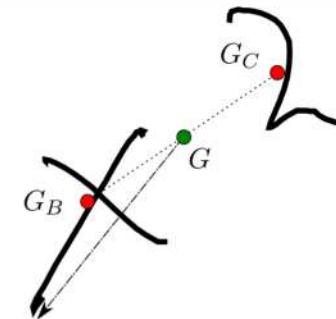
Matrix(1,6)					
[EXP] x^y+2 : 0.0000001					
[SYM] x^{y-12} : 0.00000005					
[EXP] $xy+2$: 0.00000001					
Matrix(1,5)			Matrix(2,6)		
[XOP] x^y+ : 0.000001		[EXP] C^y+2 : 0.000001		[EXP] C^yT2 : 0.000005	
[XOP] $xy+$: 0.0000005		[SYM] C^{y-12} : 0.0000001			
[EXP] $xy-1$: 0.0000001					
Matrix(1,4)		Matrix(2,5)		Matrix(3,6)	
[XOP] x^y- : 0.00003		[XOP] C^y+ : 0.00001		[EXP] $y+2$: 0.00001	
[XOP] $xy-$: 0.00001		[EXP] C^yT : 0.000005		[SYM] y_{Tz} : 0.000005	
[EXP] C^{y-1} : 0.000001					
Matrix(1,3)		Matrix(2,4)		Matrix(3,5)	
[SYM] x^y : 0.0003		[XOP] C^y- : 0.0005		[XOP] $y+$: 0.0001	
[EXP] xy : 0.0001		[XOP] $Cy-$: 0.0003		[SYM] y_T : 0.0001	
[EXP] $y-1$: 0.0001					
Matrix(1,2)		Matrix(2,3)		Matrix(3,4)	
[SYM] x : 0.005		[SYM] C^y : 0.005		[XOP] $y-$: 0.01	
[SYM] Cy : 0.003		[XOP] Cy : 0.003		[OP] $+$: 0.005	
[XPAR] $(y : 0.003)$				[SYM] T : 0.002	
				[NUM] -12 : 0.0001	
				[NUM] 12 : 0.01	
Matrix(1,1)		Matrix(2,2)		Matrix(3,3)	
[RPAR]) : 0.2		[RPAR] (: 0.1		[SYM] y : 0.1	
x_1 : 0.1		x_2 : 0.1		[OP] - : 0.2	
[SYM] C : 0.1				$+_1$: 0.1	
				T_1 : 0.1	
				[NUM] 1 : 0.2	
				$+_2$: 0.1	
				T_2 : 0.1	
				[NUM] 2 : 0.2	
				[SYM] z : 0.1	

$x^y + 2$

Figure 2.5 Example of a search for most likely expression candidate using the CYK algorithm



$$h(B, C) = [H, D, \text{dhc}, dx, dx_1, dx_2, dy, dy_1, dy_2]$$



Symbol pair (centers for x , 2 and midpoint shown).

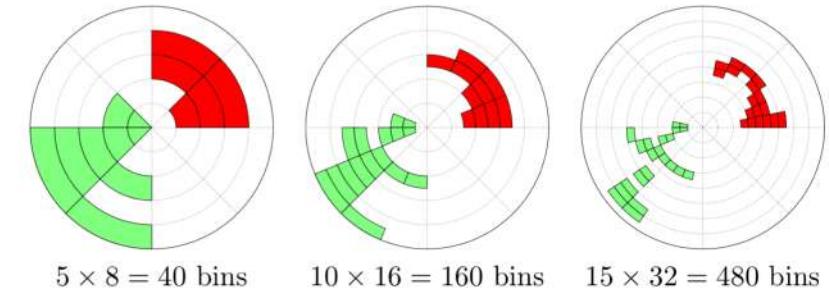
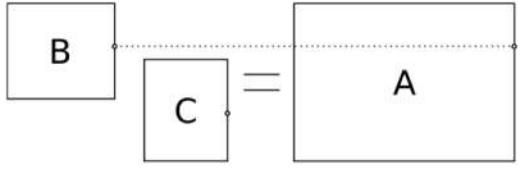
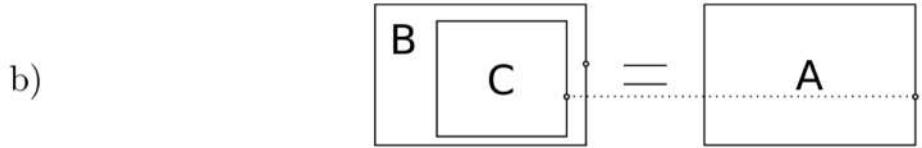
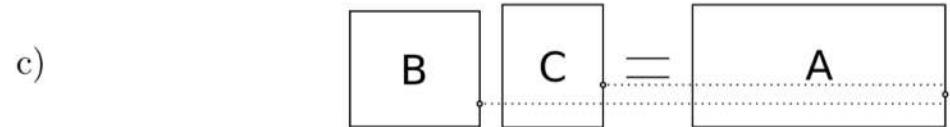
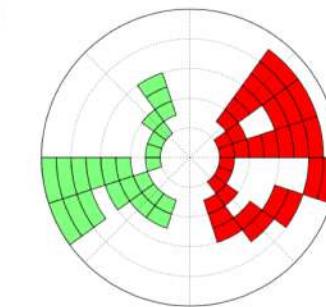


Figure 4.4: Varying distance (n) \times angle (m) resolution in a polar histogram layout descriptor. Values shown using green (-1), red (+1), and white (0).

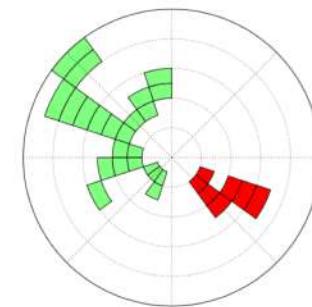
a) 

b) 

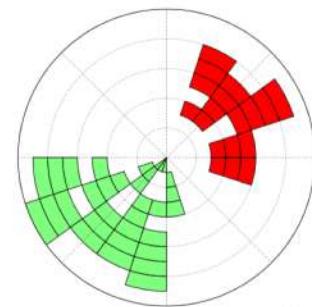
c) 



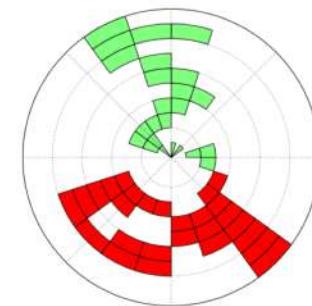
Right: dB



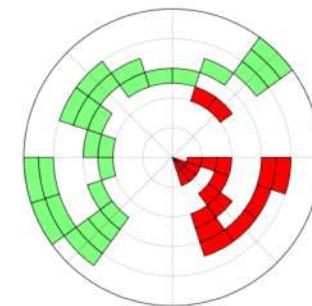
Subscript: X_u



Superscript: $\exp^{\frac{h}{d}}$



Below: $\sum_{C>h}$



Inside: \sqrt{i}

下周预告

第三课 字符序列识别

知识点1：概览RNN以及LSTM + CTC、CRNN、RARE

知识点2：概览注意力机制Attention，以及DRAW / DRAW、Transformer

实战项目： 实战CNN和RNN的综合体CRNN

THANKS