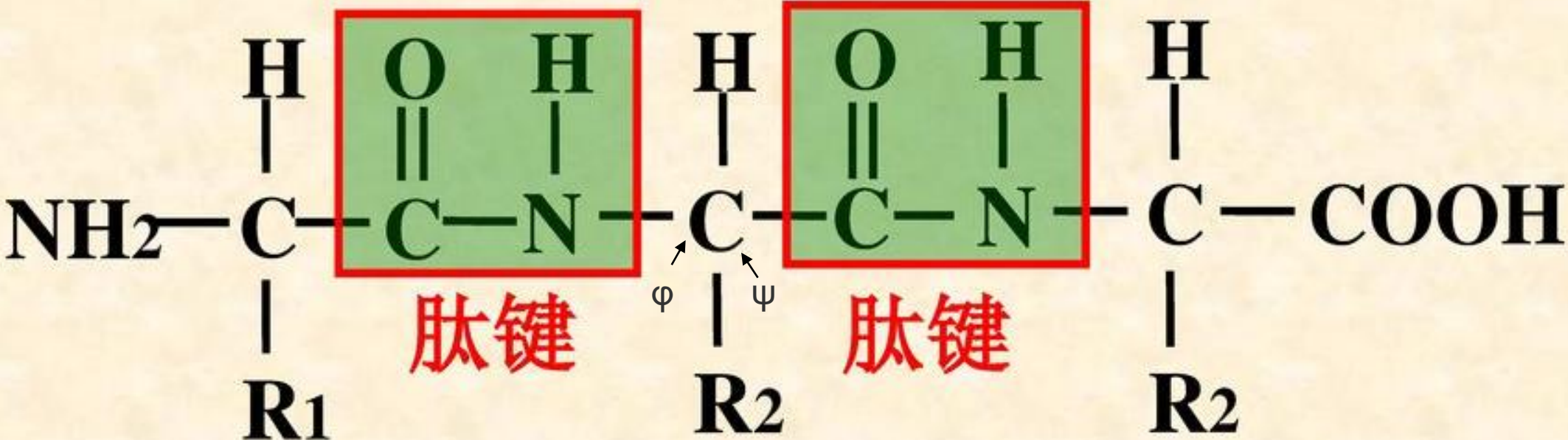


蛋白质二面角预测

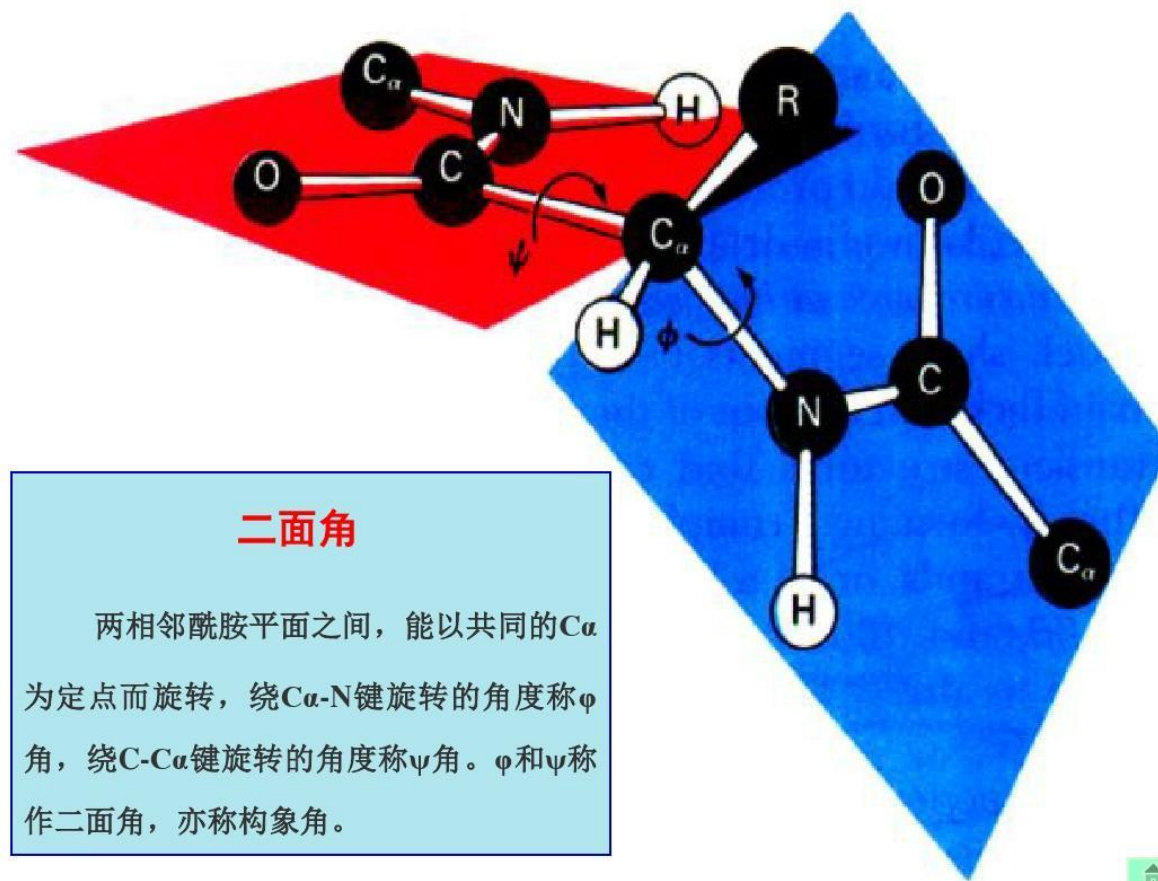
二面角 φ 和 ψ



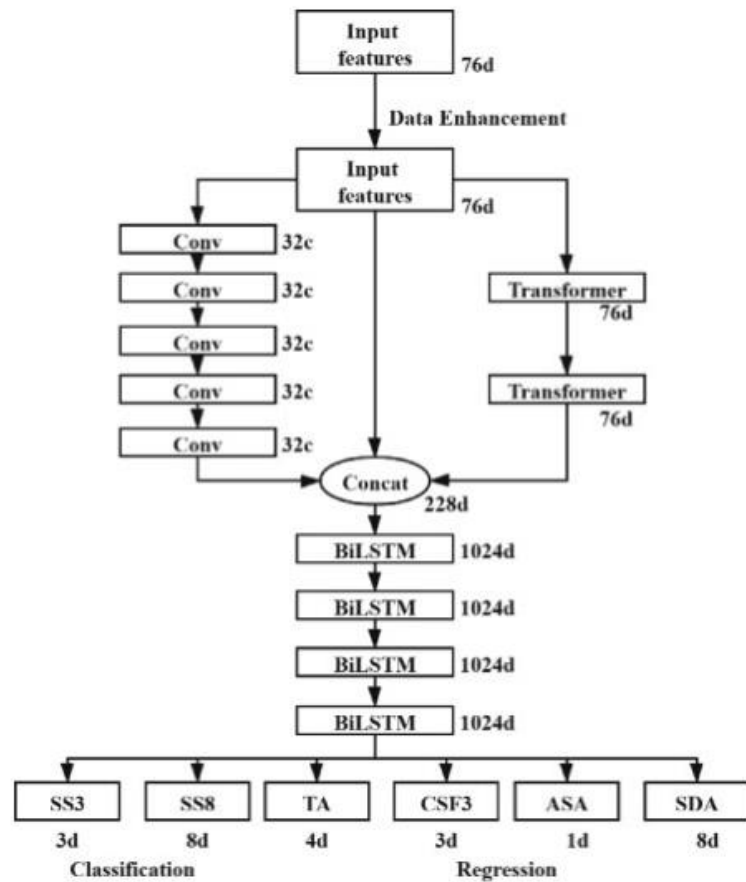
三肽

二肽

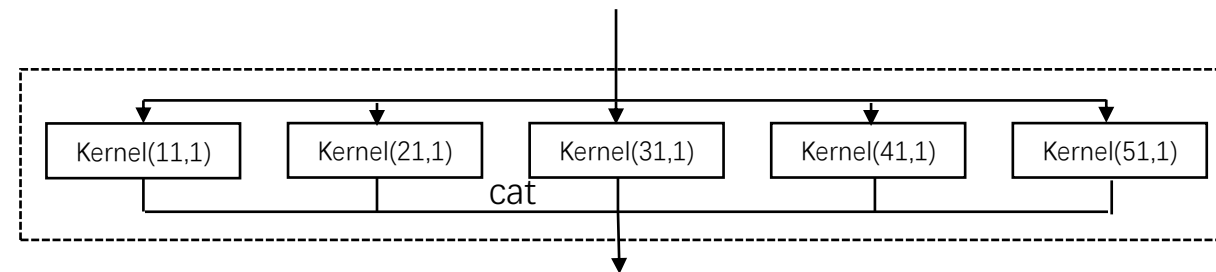
H_2O H_2O



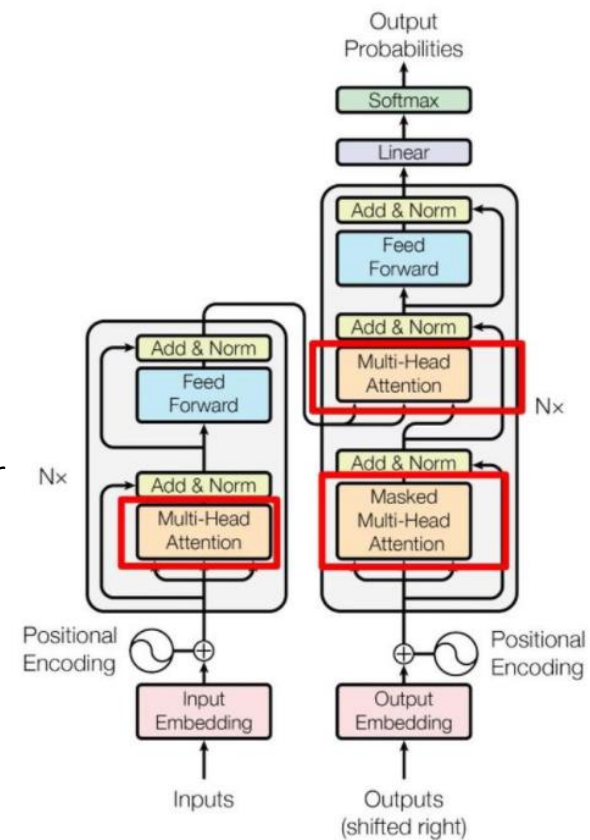
OPUS-TASS



Conv



Transformer



OPUS-TASS

Inputs: pssm (20dim) + hhm (30dim) + pc (7dim) + psp (19dim)

Outputs: SS3 + SS8 + TA + CSF3 + ASA + SDA

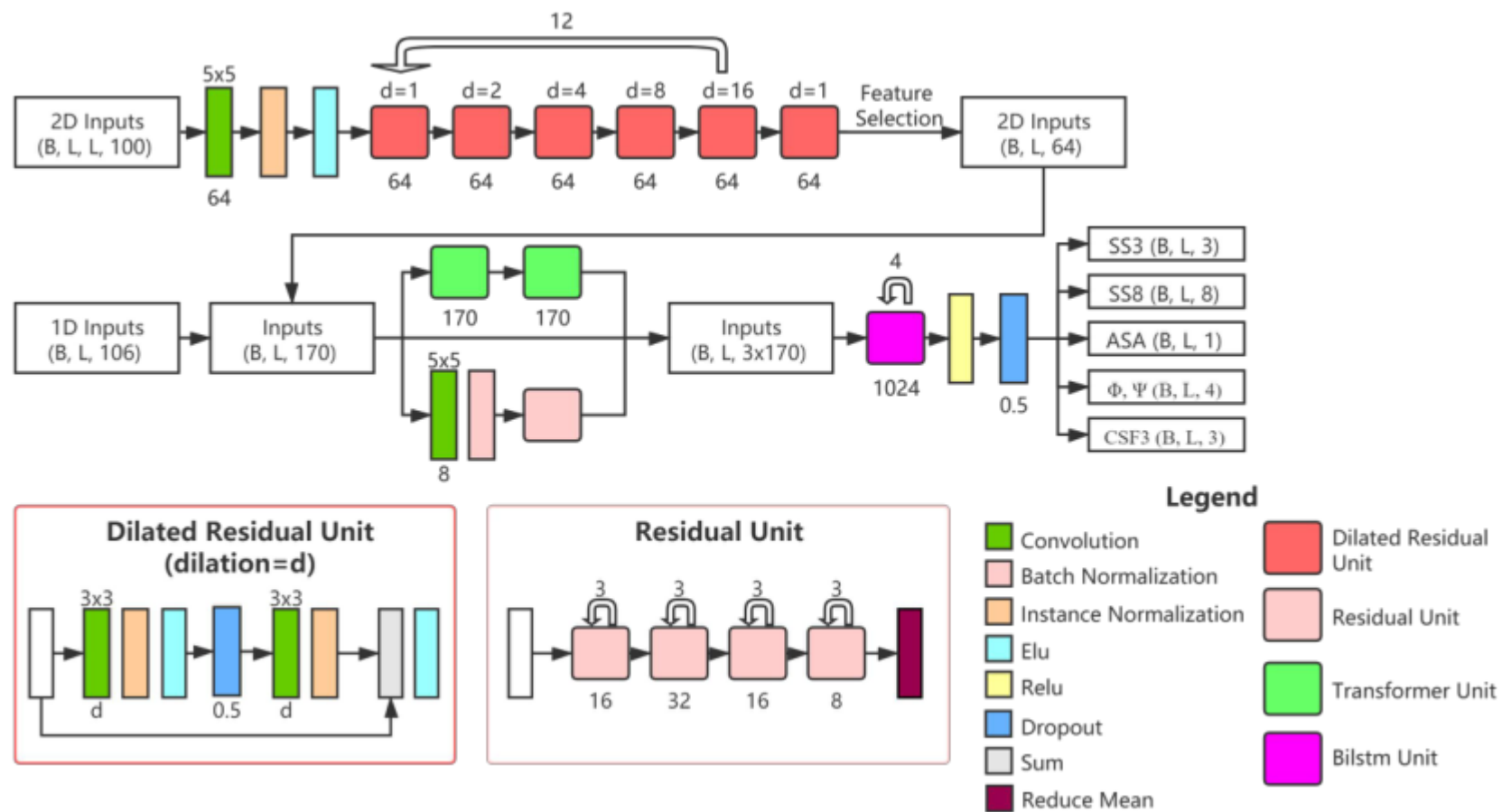
Table 3. Performance of different predictors on TEST2016 and TEST2018

Predictors	SS3	SS8	MAE(ϕ)	MAE(ψ)
TEST2016				
SPOT-1D ^a	87.16	77.10	16.27	23.26
OPUS-TASS	87.79	78.01	15.78	22.46
TEST2018				
MUFOLD ^b	84.78	73.66	17.78	27.24
NetSurfP-2.0 ^b	85.31	73.81	17.90	26.63
SPOT-1D ^a	86.18	75.41	16.89	24.87
OPUS-TASS	86.84	76.59	16.40	24.06

Table 4. Performance of different predictors on CAMEO93

Predictors	SS3	SS8	MAE(ϕ)	MAE(ψ)
SPOT-1D ^a	87.72	77.15	16.89	23.02
w/ OPUS-Refine	–	–	16.65	22.51
OPUS-TASS	89.06	78.87	16.56	22.56
w/ OPUS-Refine	–	–	16.28	21.98

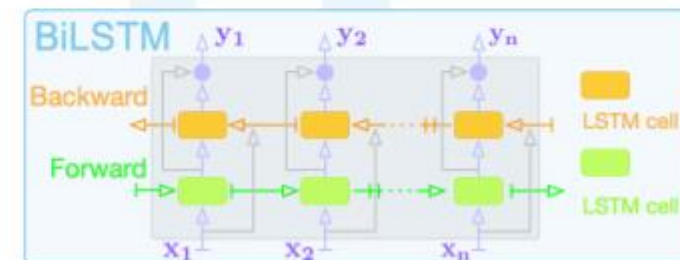
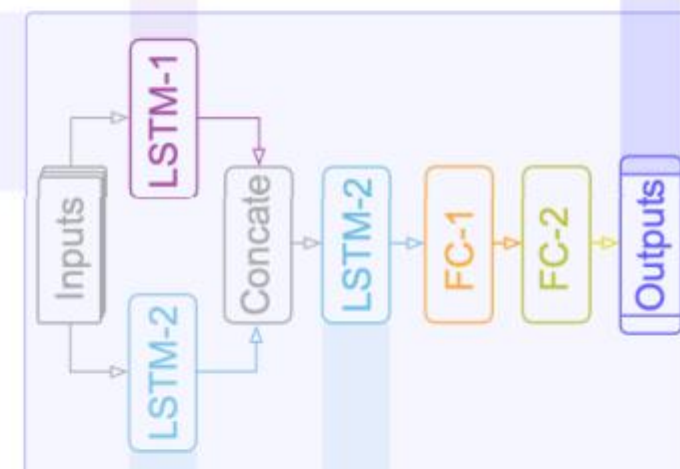
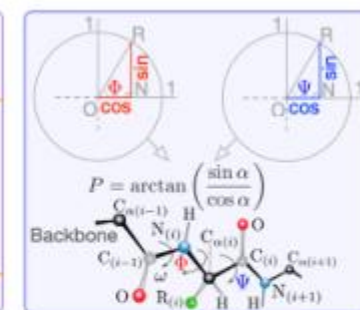
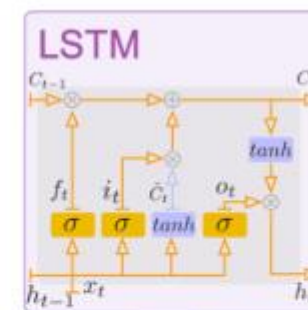
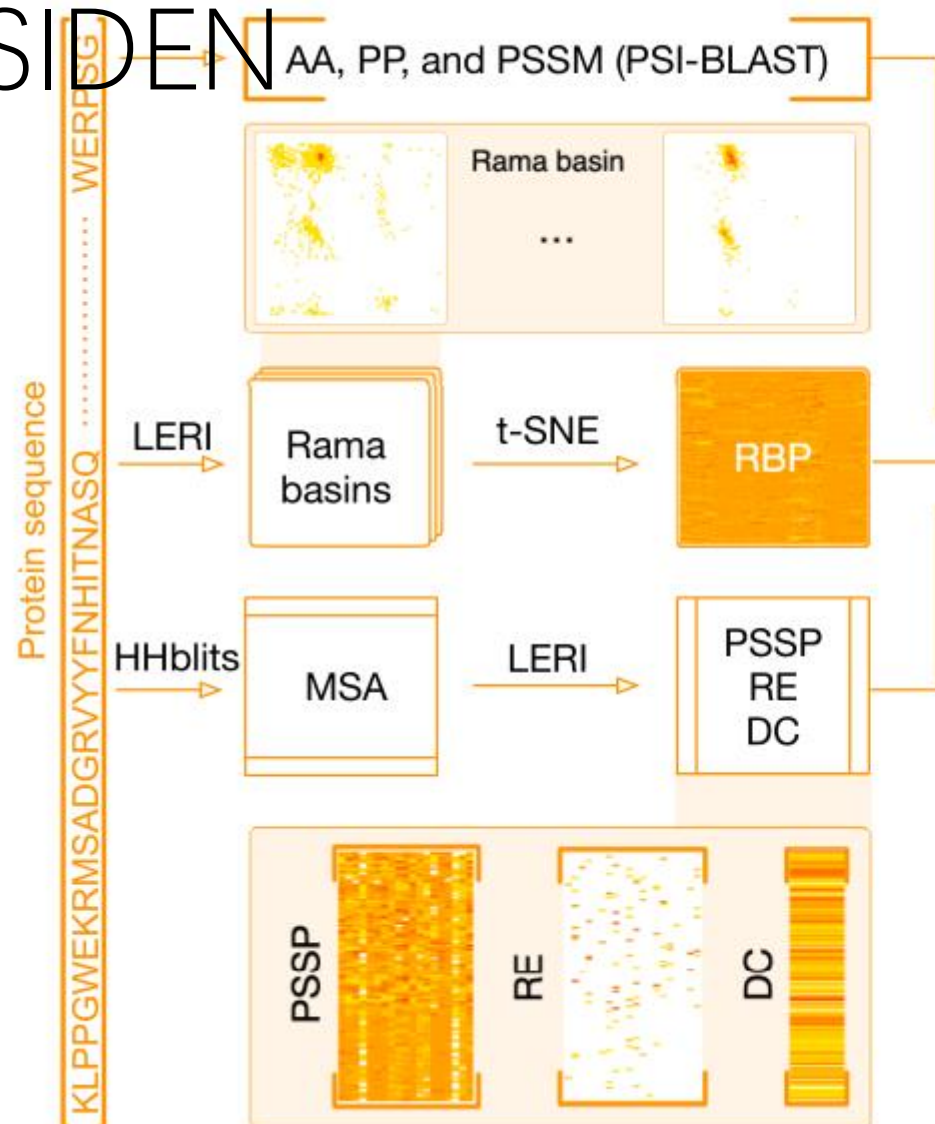
OPUS-TASS2



OPUS-TASS2

	SS3	SS8	MAE(Φ)	MAE(Ψ)	ASA
CAMEO-Hard61 (60)					
NetSurfP-2.0	83.78	70.38	20.1	29.99	0.779
SPOT-1D	83.69	70.72	19.55	29.97	0.775
OPUS-TASS	84.15	72.12	19.26	29.47	-
OPUS-TASS2	84.55	72.5	19.07	28.79	0.797
CASP-FM (56)					
NetSurfP-2.0	80.68	69.14	19.94	31.43	0.749
SPOT-1D	82.37	71.11	19.39	30.1	0.744
OPUS-TASS	83.4	73.27	18.85	28	-
OPUS-TASS2	85.96	76.28	17.94	25.17	0.804
CASP14 (15)					
NetSurfP-2.0	75.39	61.87	22.62	40.54	0.68
SPOT-1D	75.19	61.41	23.19	43.98	0.663
OPUS-TASS	77.3	63.53	21.91	38.93	-
OPUS-TASS2	80.87	68.26	20.53	33.48	0.735

ESIDEN



ESIDEN

Table 3. Performance of different methods on the TEST2016 and TEST2018.

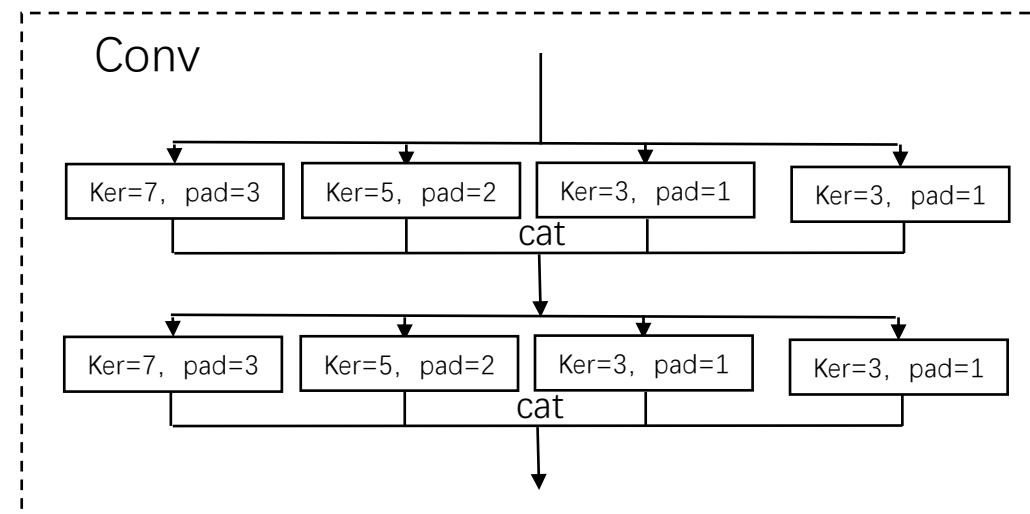
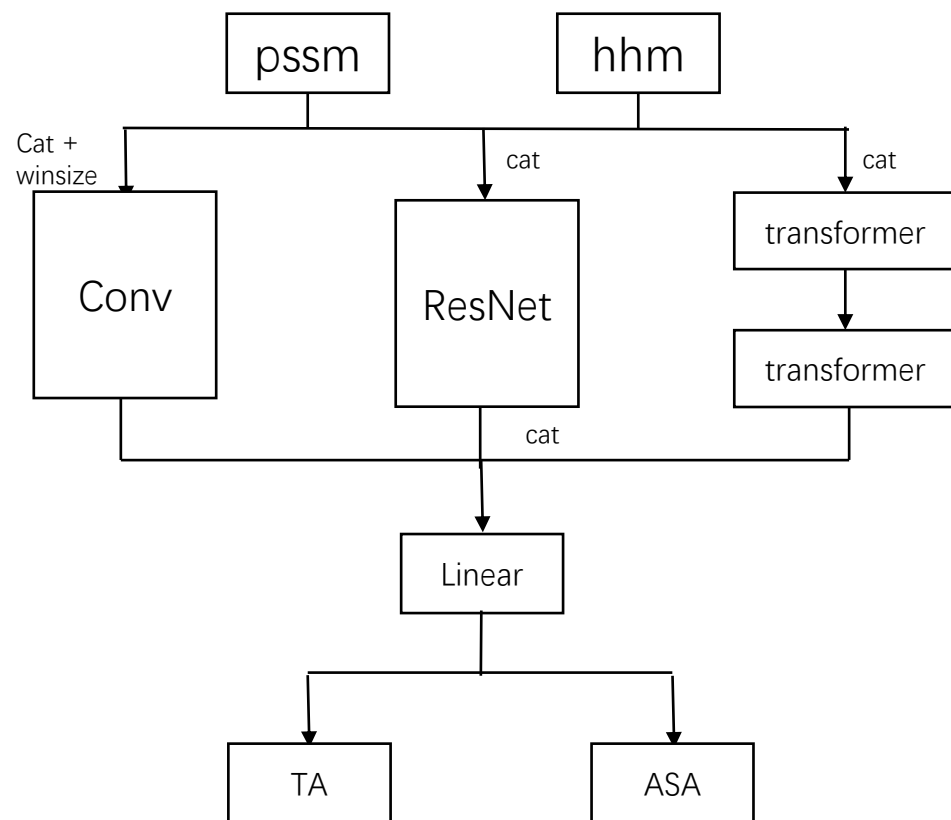
Method	Test2016		Test2018	
	MAE (ϕ)	MAE (ψ)	MAE (ϕ)	MAE (ψ)
Spider3*	17.88	26.66	18.38	28.10
RaptorX-Angle*	18.08	26.68	21.01	35.95
SPOT-1D*	16.27	23.26	16.89	24.87
OPUS-TASS [†]	15.78	24.46	16.40	24.06
ESIDEN	15.61	19.30	17.09	21.70

Table 4. Comparison among different methods on the TFM targets of the CASP11, CASP12, CASP13, and CASP14

Method	CASP11 (27)		CASP12 (11)		CASP13 (13)		CASP14 (8)	
	MAE (ϕ)	MAE (ψ)	MAE (ϕ)	MAE (ψ)	MAE (ϕ)	MAE (ψ)	MAE (ϕ)	MAE (ψ)
Spider3*	19.19	34.63	21.14	34.92	22.48	38.46	23.41	38.79
RaptorX-Angle*	20.33	40.05	21.71	38.22	22.73	41.18	24.95	48.06
SPOT-1D*	18.54	25.77	20.21	31.71	22.60	34.28	23.42	33.44
ESIDEN	17.93	23.87	19.68	28.63	22.51	32.68	23.16	29.89

*The results are obtained locally using the Spider3, RaptorX-Angle, and SPOT-1D standalone packages, respectively

以往用过的模型模型



TA: $\sin(\varphi)$ 、 $\cos(\varphi)$ 、 $\sin(\psi)$ 、 $\cos(\psi)$

训练小批量选择模型

Epoch = 20

Batchsize = 8(蛋白质)

Batchnum = 250

Samplenum = 2000

将batch中的每个蛋白质
都扩充到和该batch中最
长的蛋白质的残基数一样
的长度

进入model进行训练

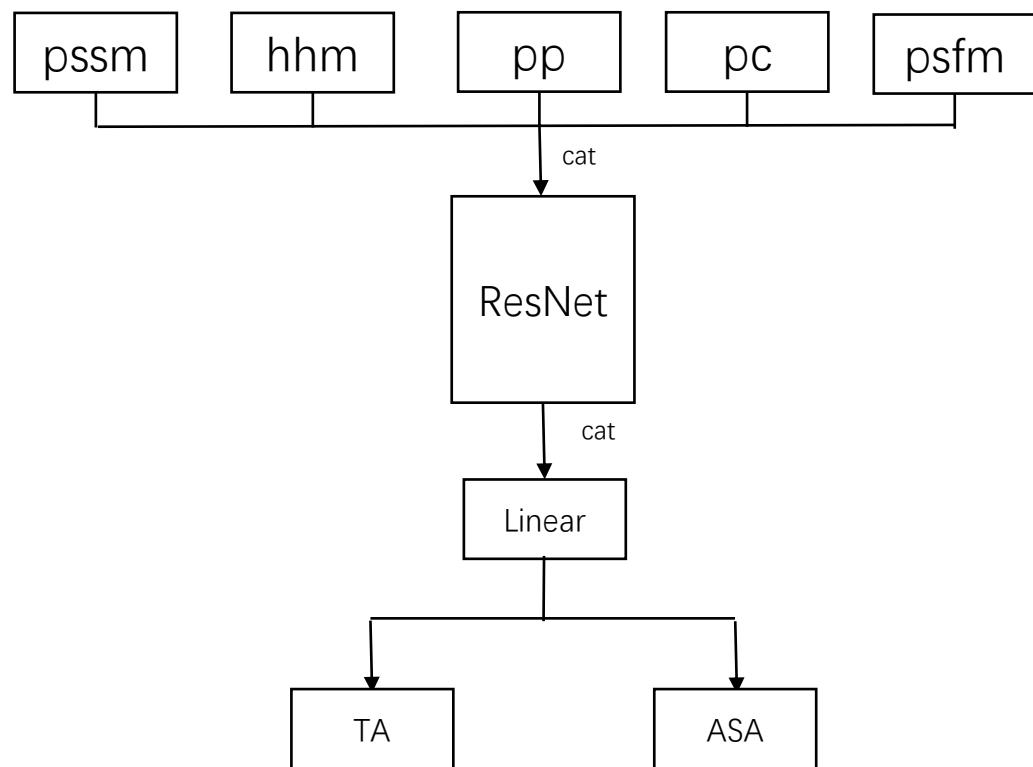
训练后再截取回原本每个
蛋白质的原始长度，进行
mae计算

实验结果

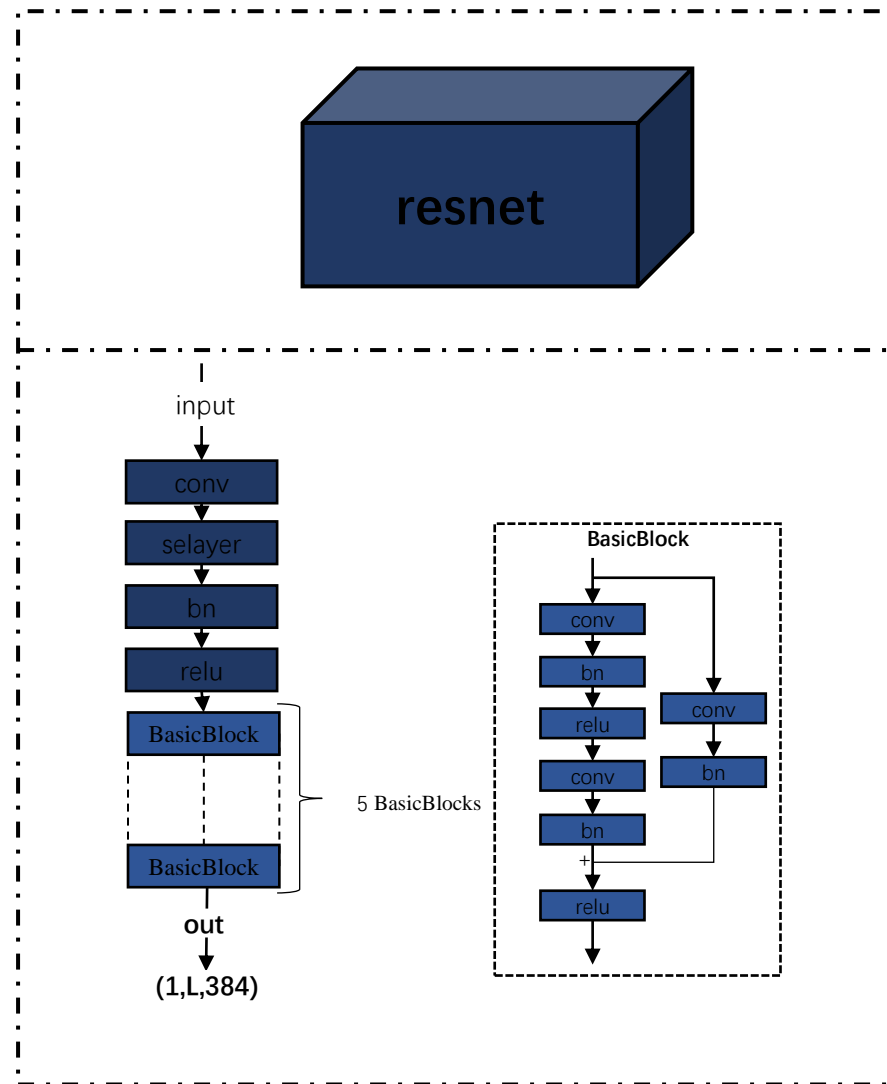
PSSM+HHM + pp + pc
Samplenum = 2000

实验序号	模型	MAE (Phi)	MAE (Psi)
1	OPUS-TASS	19.62	29.99
2	Transformer+BiLSTM	19.94	30.31
3	ResNet+BiLSTM	19.63	30.28
4	Transformer+ResNet+BiLSTM	19.88	29.86
5	CNN+Transformer+ResNet+SElayer+BiLSTM	19.75	30.24
6	(Transformer * ResNet) + BiLSTM	19.71	30.03
7	Transformer+CNN+BiLSTM	19.87	30.43

以往用过的模型模型



TA: $\sin(\varphi)$ 、 $\cos(\varphi)$ 、 $\sin(\psi)$ 、 $\cos(\psi)$



实验结果

ResNet+BiLSTM

Epoch = 20

Batchsize = 1000(残基数)

Samplenum = 2000

实验序号	winsize	MAE (Phi)	MAE (Psi)
1	9	19.49	29.71
2	15	19.55	30.27
3	5	19.15	29.45
4	1	18.85	29.13

实验结果

ResNet(PLUS)+BiLSTM

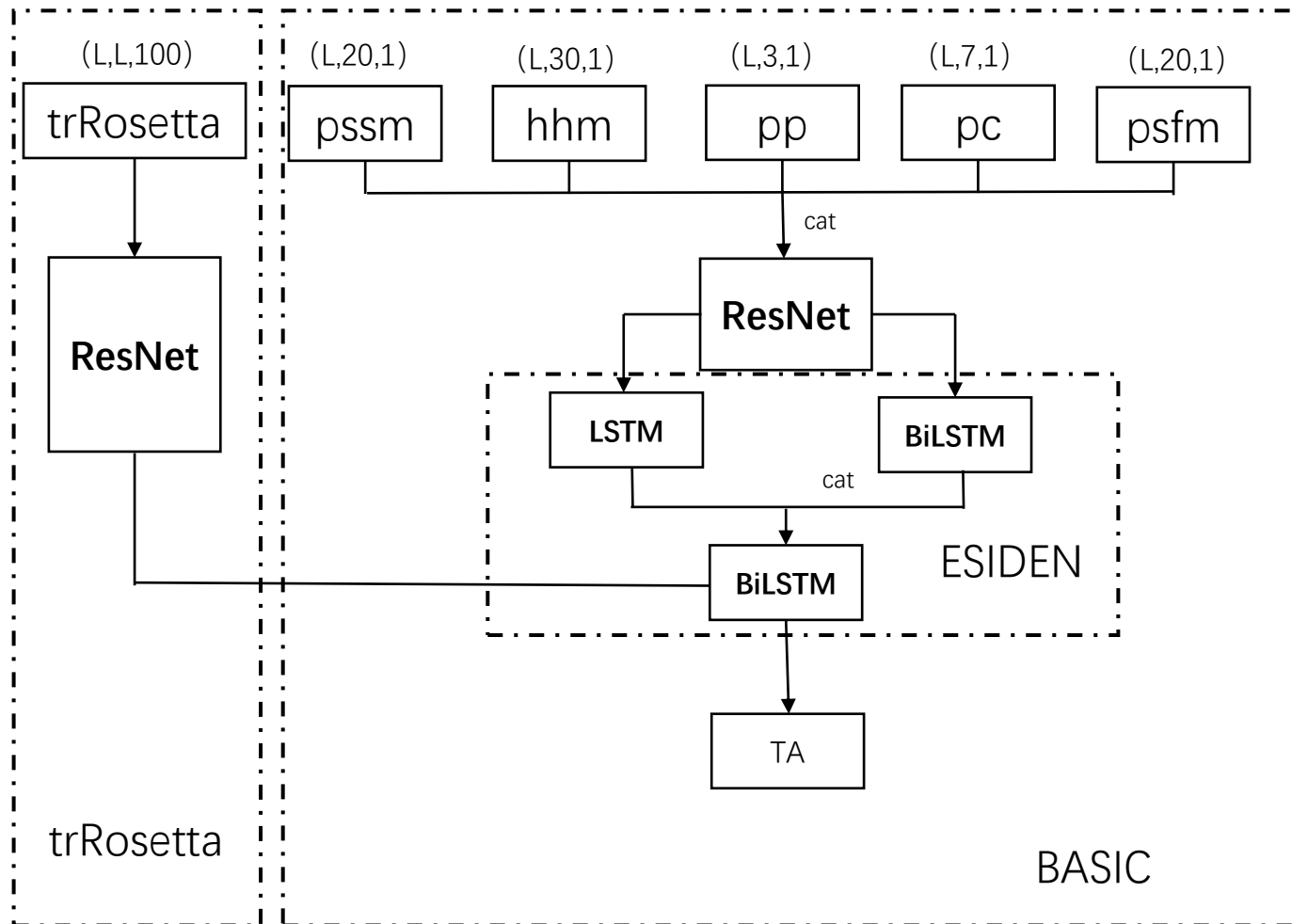
Epoch = 20

Batchsize = 1000(残基数)

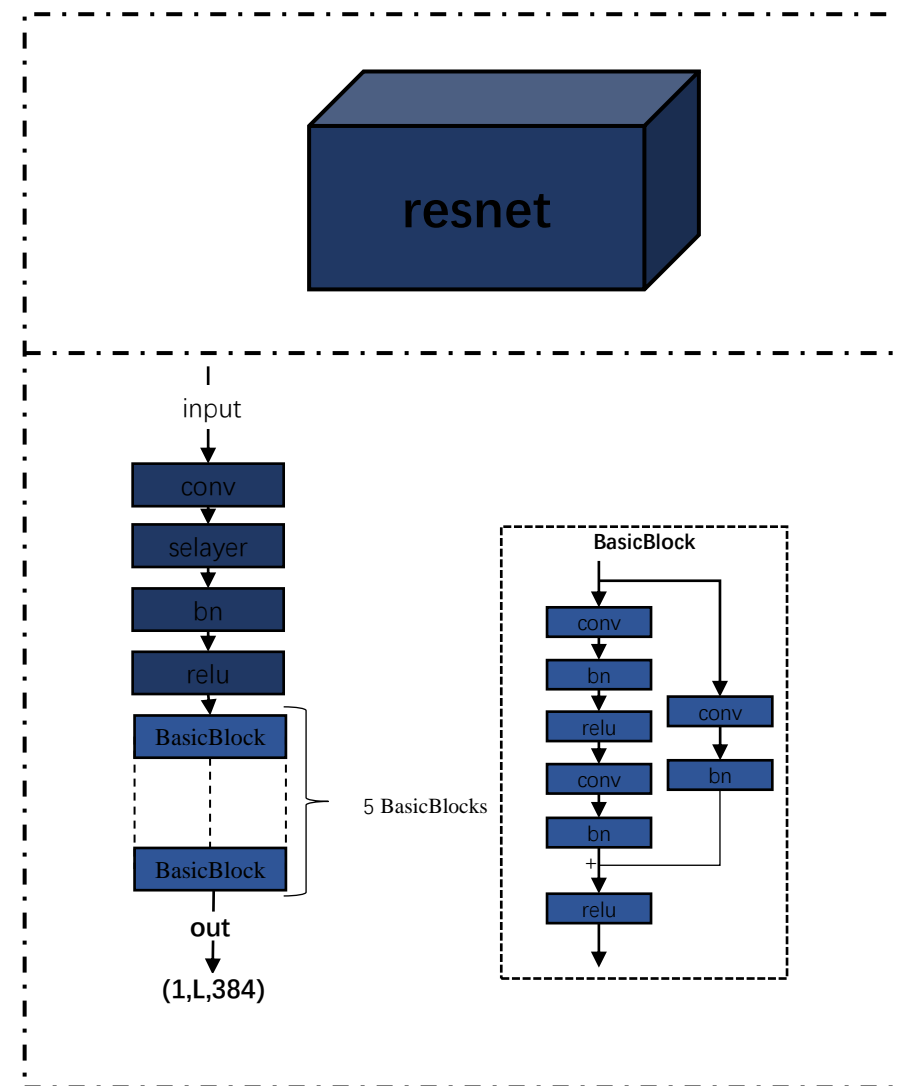
Samplenum = 10024 (all)

实验序号	winsize	MAE (Phi)	MAE (Psi)
1	1	17.50	25.49
2			
3			
4			

现在模型（加入trRosetta）



TA: $\sin(\varphi)$ 、 $\cos(\varphi)$ 、 $\sin(\psi)$ 、 $\cos(\psi)$



训练小批量选择模型

因为想要加入trRosetta的特征 (L,L,100)，其中包含着残基间距离和位置信息，所以改为单个蛋白质为单位进行训练

Epoch = 20

Batchsize = 1(蛋白质计数)

Samplenum = 2000

验证单个特征有效性

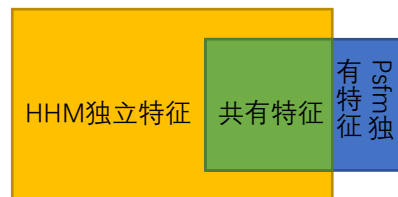
实验序号	Feature(基于pssm)	MAE (Phi)	MAE (Psi)
1	pssm=P	19.93+0.09	30.50+0.42
2	P+RE	19.69	30.00
3	P+DC	19.85	30.35
4	P+RE+DC	19.70	29.86
5	P+pi(position)	20.12	31.12
6	P+hyd	19.21	29.89
7	P+onehot	19.00	29.29
8	P+pKx	19.85	30.54
9	P+HHM	18.62	27.33
10	P+psfm	18.81	27.86
11	P+pp	19.31	30.24
12	P+pc	19.30	30.36

验证特征有效性补充

实验序号	Feature(基于pssm)	MAE (Phi)	MAE (Psi)
1	pssm=P	19.93+0.09	30.50+0.42
2	P+HHM	18.62	27.33
3	P+psfm	18.81	27.86
4	P+HHM+psfm	18.60	27.28
5	P+HHM+psfm+onehot	17.73	26.17
6	P+HHM+onehot	17.79	26.36
7	P+psfm+onehot	18.00	26.86

结论：

psfm中的大部分有效特征HHM中都有， 但还含有少量有效特征， HHM特征中除了psfm有效特征， 还有比psfm更多的独立有效特征



验证特征有效性补充

实验序号	Feature(基于pssm)	MAE (Phi)	MAE (Psi)
1	pssm=P	19.93+0.09	30.50+0.42
2	P+HHM+psfm	18.60	27.28
3	P+hhm+psfm+hyd+pp+pc	17.62	26.06
4	P+hhm+psfm+hyd+pp+pc+re+dc	17.88	26.41
5	P+HHM+psfm+onehot	17.73	26.17
6	P+HHM+onehot	17.79	26.36
7	P+psfm+onehot	18.00	26.86
8	P+hhmdim+psfm+hyd	17.74	26.39
9	P+hhmdim+pp+psp	17.78	26.25
10	P+hhmdim+pp	17.84	26.41
11	P+hhm+psfm+onehot+hyd+pp+pc+psp	17.68	26.19
12	P+hhm+psfm+onehot+hyd+pp+pc	17.65	26.09

Small batch train (2000个) 改版后

(C): Epoch变大还可以再收敛

实验序号	feature	MAE (Phi)	MAE (Psi)
1	Basic+tr101	17.98	25.30
2	Basic+tr151	17.71(C)	25.13(C)
3	Basic+opus-x	17.96	25.37
4	Basic+tr21	17.83	25.31
5	Basic+tr11	17.75	25.15
6	Basic+tr251	18.24	25.74
7	Basic+tr1	20.05	31.47

Small batch train (2000↑) Inception
151trsize

实验序号	feature	MAE (Phi)	MAE (Psi)
1	res1	19.54	30.07
2	res2	19.77	31.2

model select(没什么用)

实验序号	Model selcet	MAE (Phi)	MAE (Psi)
1	model1	17.61	25.02
2	model2	17.71	25.04
3	model3	17.72	25.00
4	model4	18.07	25.42
5	model5	17.74	25.02
6	model6	17.89	25.37
7	model7	17.69	25.11
8	model8	17.68	25.05
9	model9	17.75	25.24
10	model10	17.74	25.12
	model11	18.55	25.59
12	model12	17.82	24.93

将二级结构预测结果作为特征

Sample_num: 10024

实验序号	feature	MAE (Phi)	MAE (Psi)
1	Pssm+hhm+pp+pc+hyd+psfm=basic	17.50	25.49
3	Basic + SS8	15.17	17.83

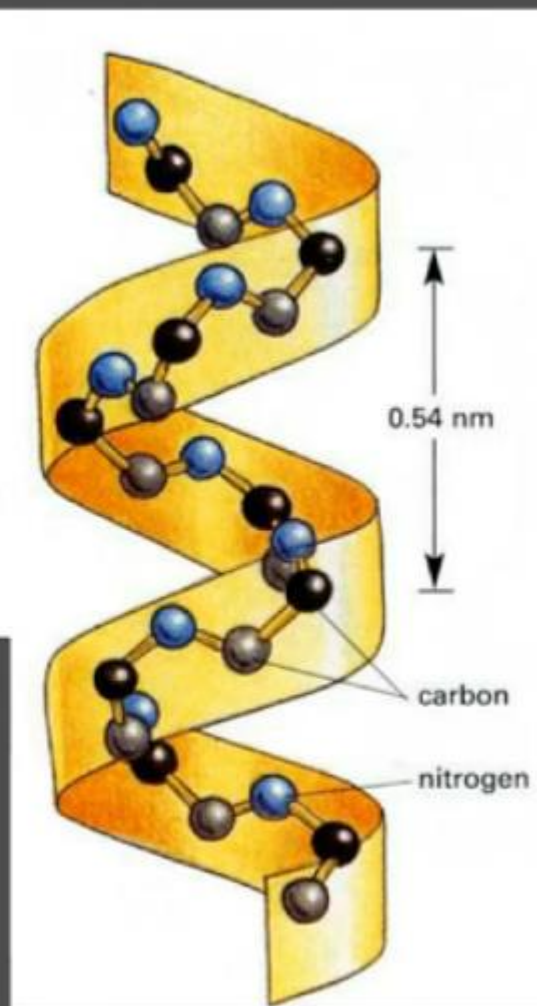
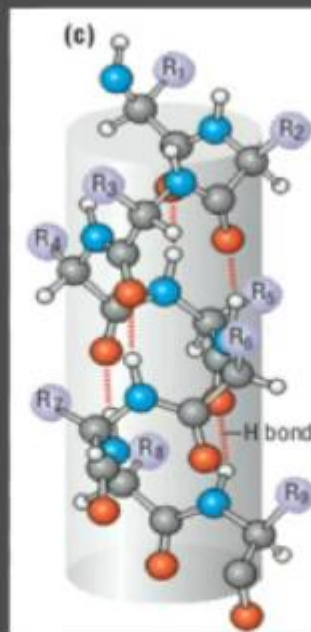
二级结构与二面角关系

■ α 螺旋

● 相关参数

- 3.6氨基酸残基/圈；
- 100° /氨基酸残基；
- 螺距0.54nm；
- 0.15nm/残基
- 螺旋半径0.23nm；
- Φ 和 ψ 分别为 -57° 和 -47°

● R基团位于螺旋外侧，不参与螺旋形成，大小形状和带电荷性质影响螺旋形成和稳定。



二级结构与二面角关系

■ α 螺旋

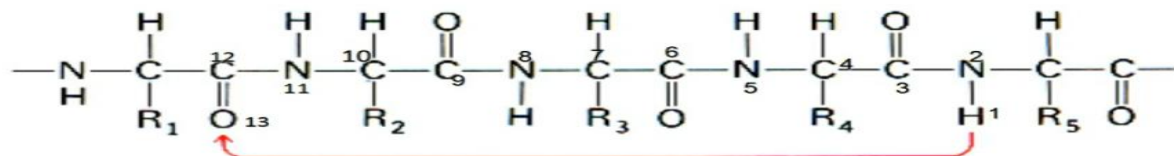
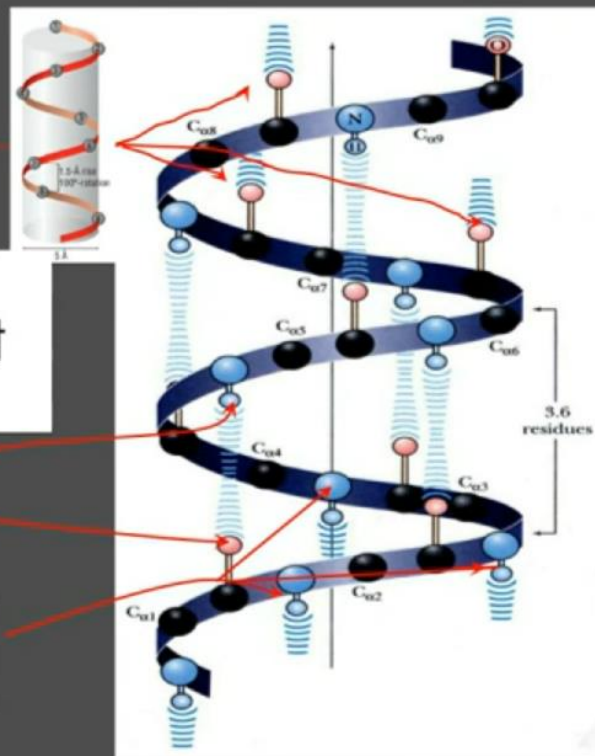
羧基端后3个肽键内C=O不能形成 α -螺旋的氢键

α -螺旋的氢键连接发生在C=O和NH之间，被氢键封闭的环含有13个原子；

第n个残基

第n+4个残基

氨基端前3个肽键内NH不能形成 α -螺旋的氢键



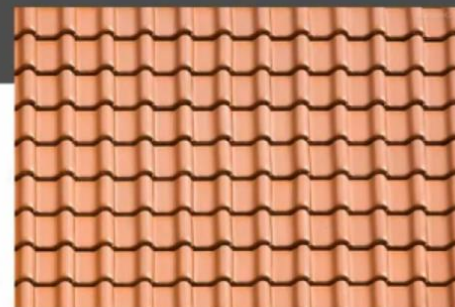
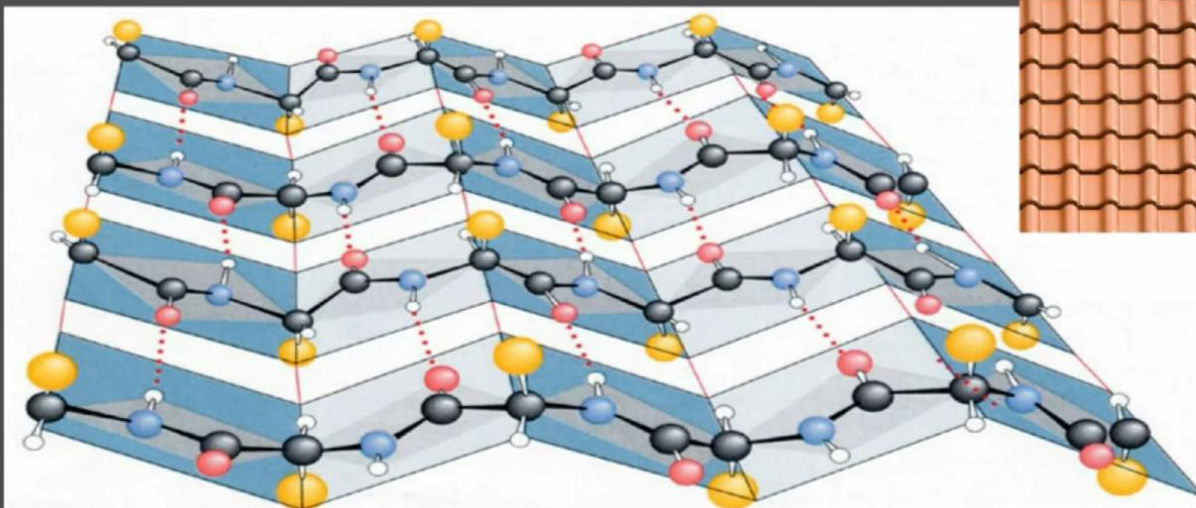
Figure

In the α helix, the NH group of residue n is hydrogen bonded to the CO group of residue $(n - 4)$.

二级结构与二面角关系

■ β 折叠片-----一种相当伸展的结构

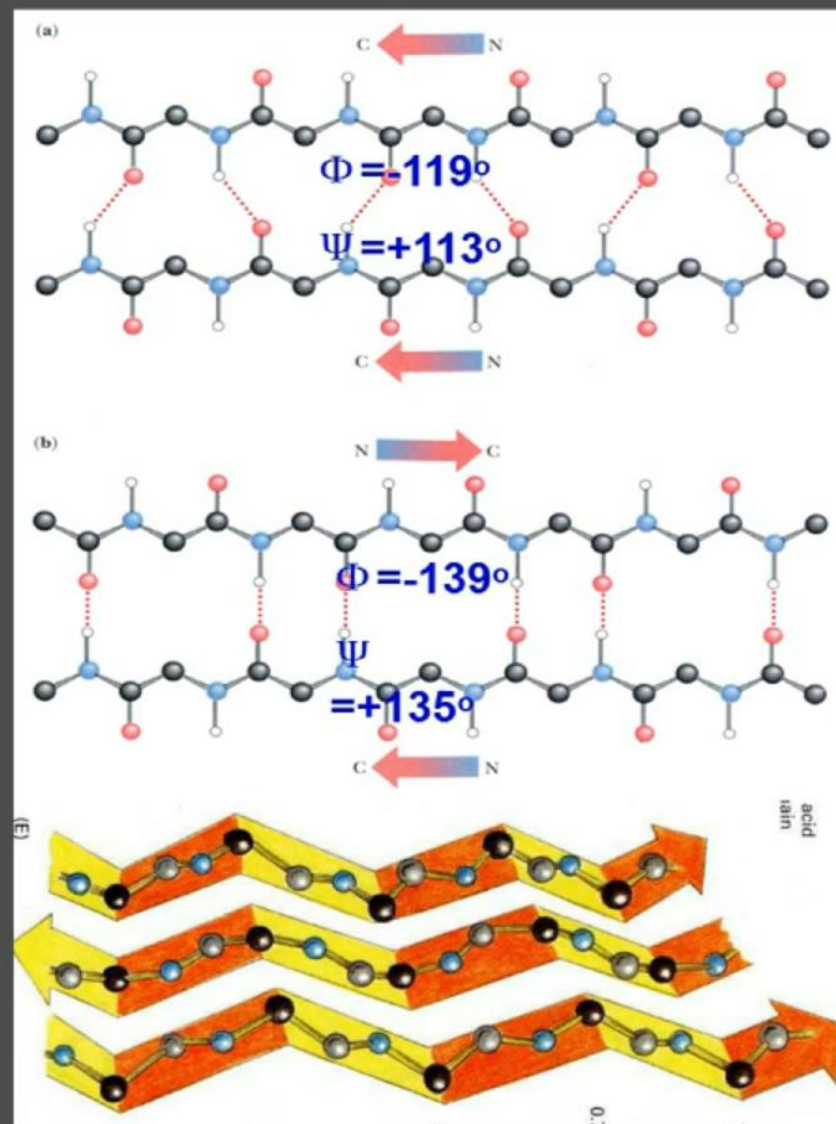
- 定义：几乎完全伸展的肽段侧向聚集在一起，相邻主链间氨基和羰基形成氢键，这样的多肽构型就是 β -折叠（ β -sheet），或 β -折叠片层（ β -plated sheet），其中每一股肽段成为 β -股（ β -strand）；一条多肽链不同肽段，两条或多条肽链间肽段，或不同蛋白质各提供一个 β -股。



二级结构与二面角关系

● 结构特点

- 肽段几乎完全伸展，肽平面间成锯齿状；
- 肽段平行排列，肽段间形成氢键
- 侧链基团垂直于相邻两个肽平面的交线，交替分布在折叠片层的两侧；
- 平行和反平行两类，反平行折叠形成的氢键N-H-O三个原子位于同一直线，较稳定；
- 反平行折叠的每一个氨基酸残基上升0.347nm，平行折叠上升0.325nm

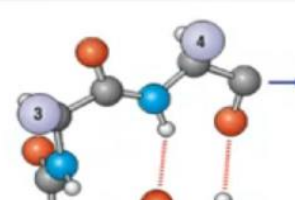
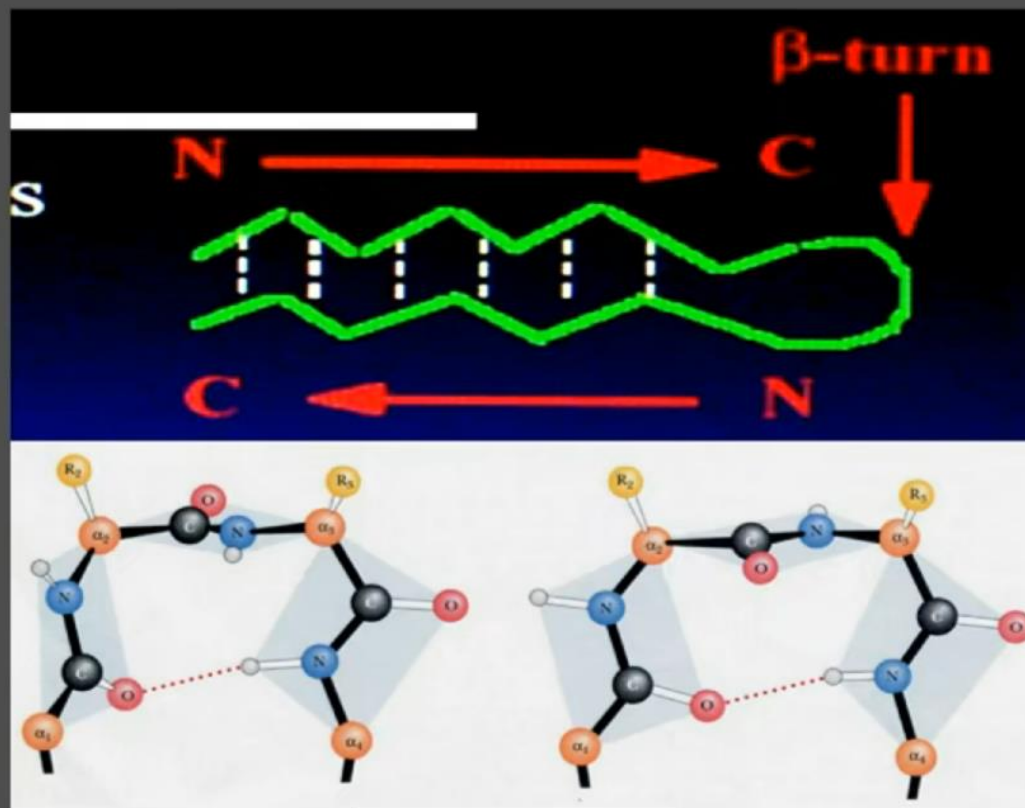


二级结构与二面角关系

■ β -转角和 β 凸起

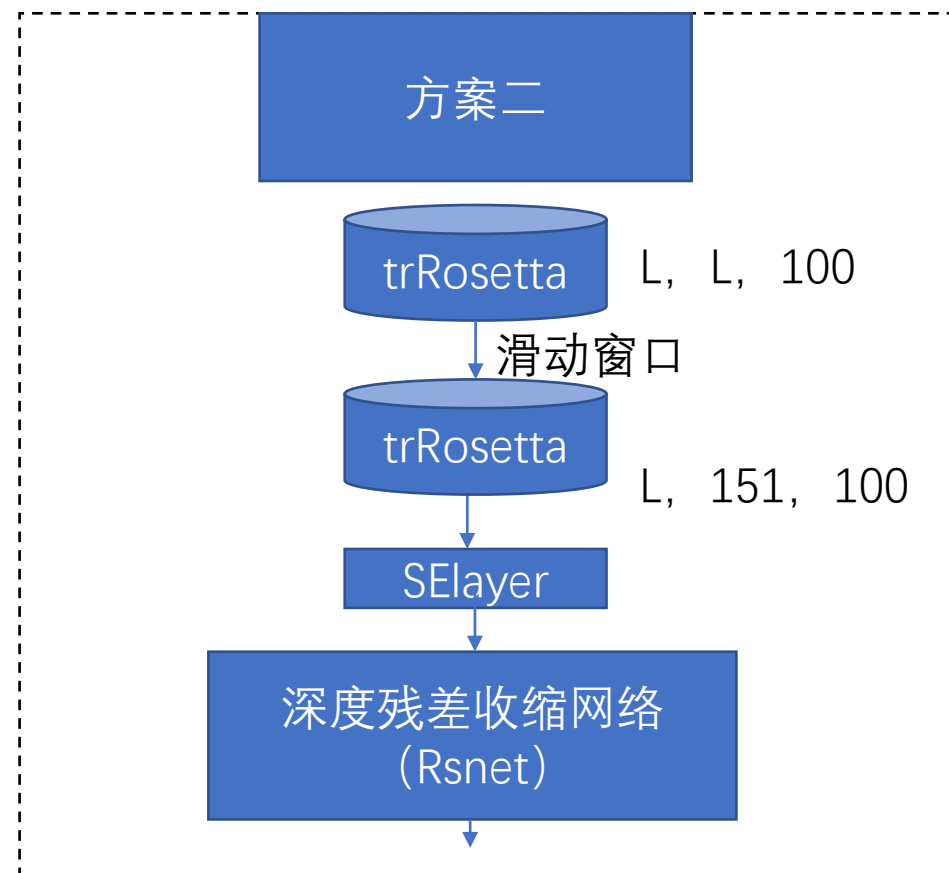
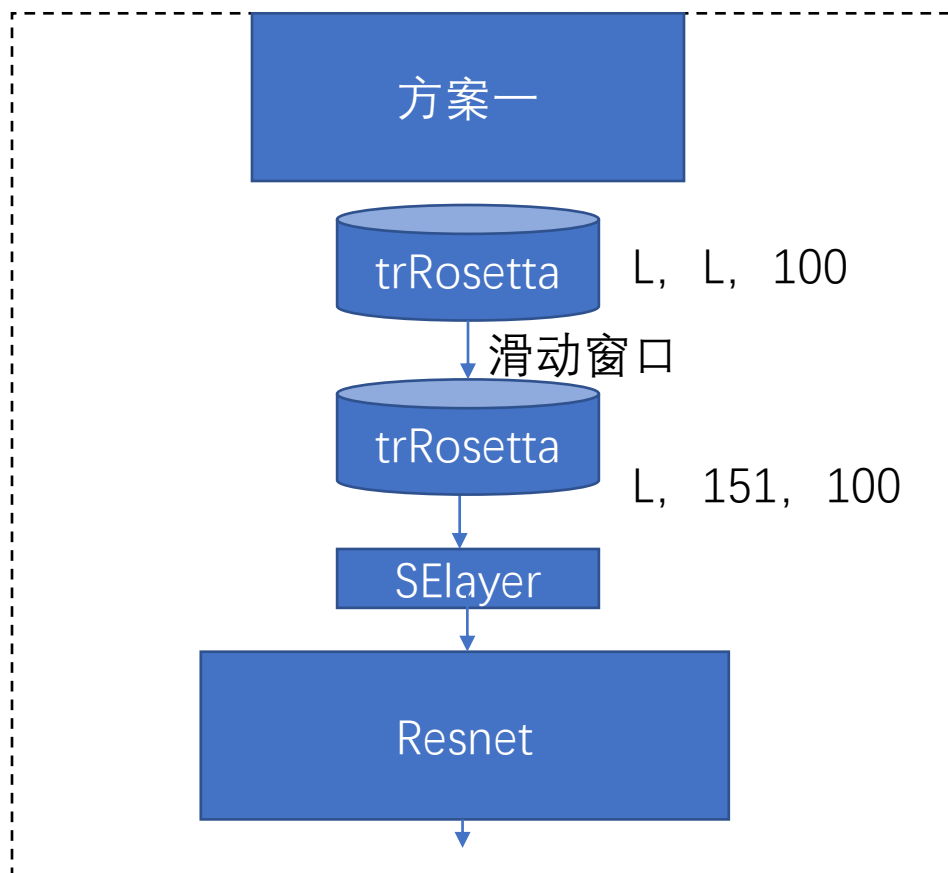
● β -转角

- 一种非重复性结构
- 4个连续的氨基酸残基组成
- 主链骨架以180°返回折叠
- C1羰基氧与C4亚氨基氢形成氢键；C1与C4之间距离小于0.7nm
- Gly、Pro频率高



加入trRosetta预测残基间距离和方位信息 (10024)

- 将taRosetta中的预测预测残基间距离和方位信息 (L,L,100)，经过如下处理后加入新特征



深度残差收缩网络

- 分为2种：DRSN-CS和DRSN-CW
- 面向的是带有“噪声”的信号，将“软阈值化”作为“收缩层”引入残差模块之中，并且提出了自适应设置阈值的方法。实际上，这里的“噪声”可以宽泛地理解为“与当前任务无关的特征信息”
- 核心算法的部分（软阈值化和其导数公式）

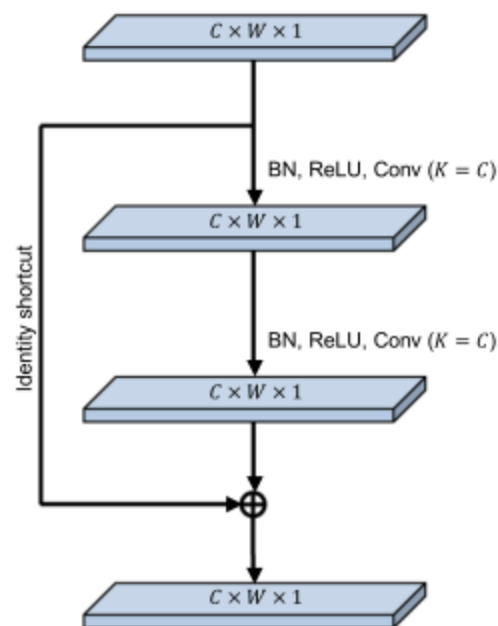
$$y = \begin{cases} x - \tau & x > \tau \\ 0 & -\tau \leq x \leq \tau \\ x + \tau & x < -\tau \end{cases}$$

$$\frac{\partial y}{\partial x} = \begin{cases} 1 & x > \tau \\ 0 & -\tau \leq x \leq \tau \\ 1 & x < -\tau \end{cases} .$$

DRSN-CS中的模块RSBU-CS

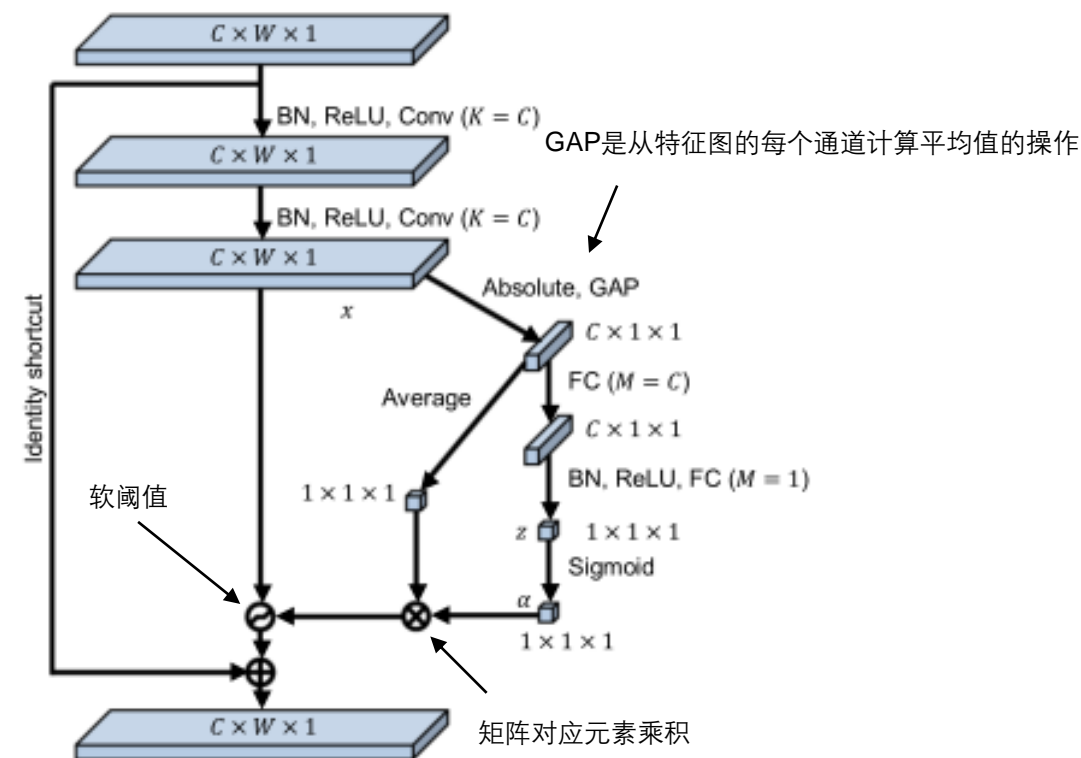
C 、 W 和 1 分别是特征图的通道数、宽度和高度， K 是卷积层中卷积核的数量

ResNet中的残差块



Residual building units (RBUs)

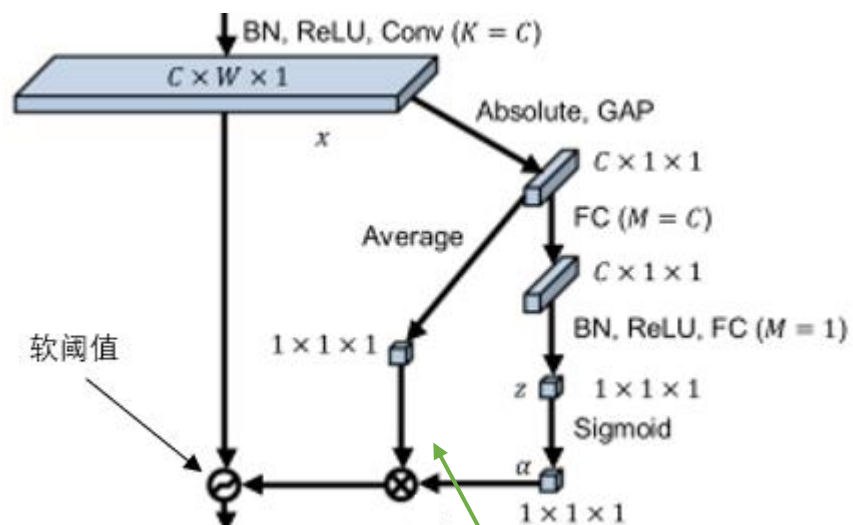
DRSN-CS中的残差块



RSBU-CS

优势：阈值由反向传播自动学习调整，类似于加权处理

DRSN-CS



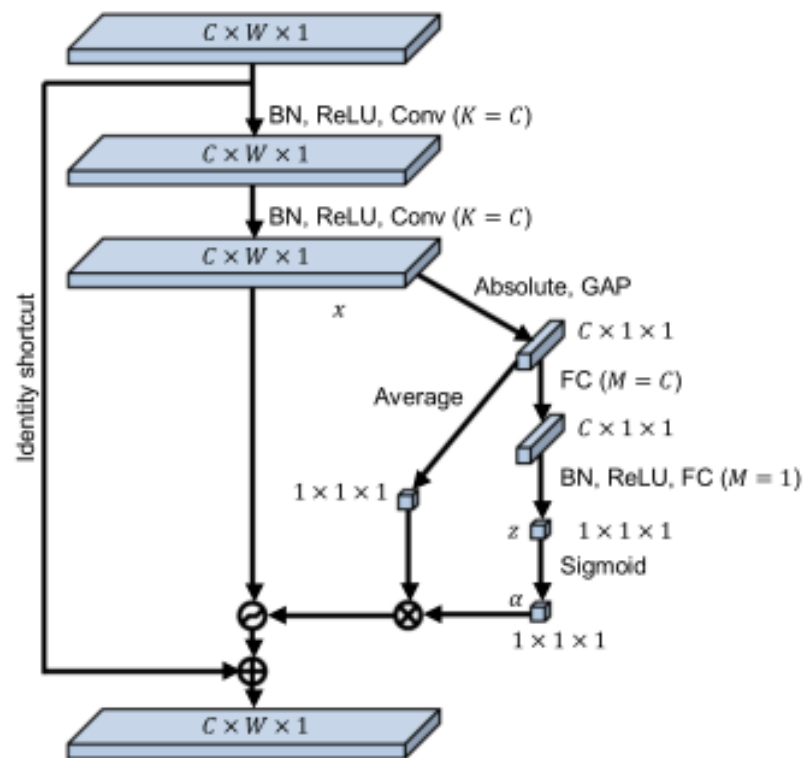
$$\tau = \alpha \cdot \text{average}_{i,j,c} |x_{i,j,c}|$$

软阈值 经sigmoid后的输出 输入特征 宽 高 通道数

DRSN-CW和DRSN-CS的对比

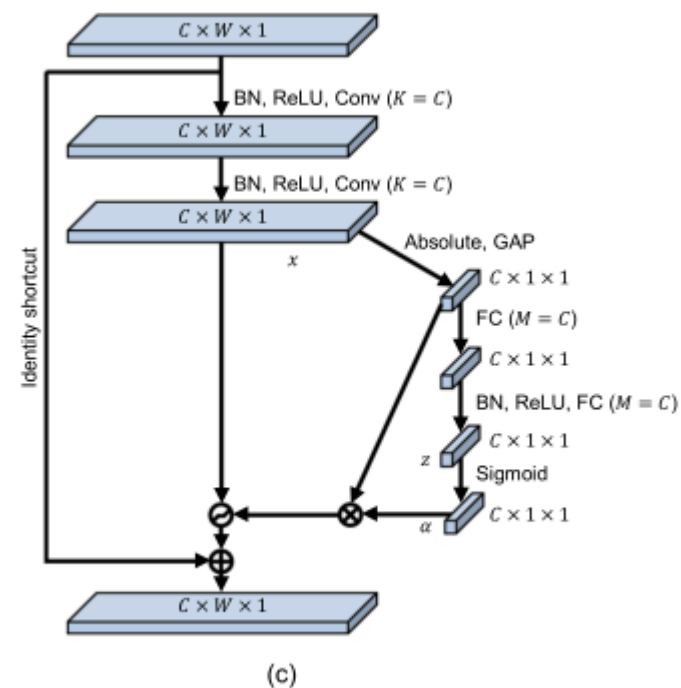
C、W和1分别是特征图的通道数、宽度和高度，K是卷积层中卷积核的数量

DRSN-CS中的残差块



RSBU-CS

DRSN-CW中的残差块



RSBU-CW

加入trRosetta残基间距离和方位信息 (10024)

实验序号	feature	MAE (Phi)	MAE (Psi)
1	Pssm+hhm+psfm+hyd+pp+pc+ss8=basic	15.17	17.80
2	Basic+RsNet151	14.75	17.12
3	Basic+ResNet151	待验证	待验证
4			

同源性比对

二面角测试集： 1212

二级结构训练集： 31338

因此要做共 $1212 \times 31338 = 37,981,656$ 次同源性比对

同源性大于0.25个数： 8257306

同源性大于0.25个数： 401673

PHI结果

平均数比较结果:

都好 = 330

Q8好phi不好 = 307

Q8不好phi好 = 211

都不好 = 197

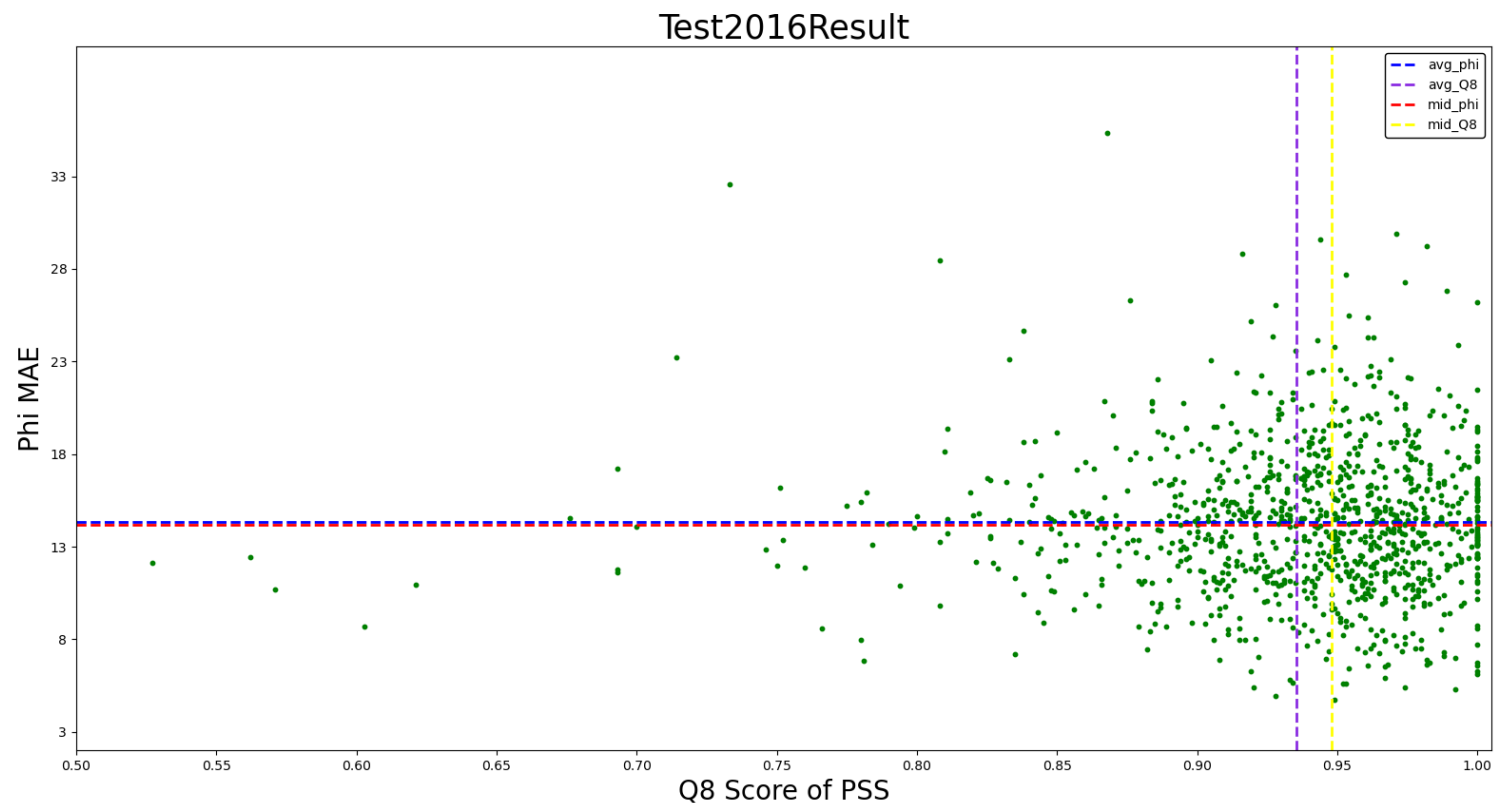
中位数比较结果:

都好 = 267

Q8好phi不好 = 251

Q8不好phi好 = 255

都不好 = 272



PSI结果

平均数比较结果:

都好 = 352

Q8好phi不好 = 286

Q8不好phi好 = 228

都不好 = 179

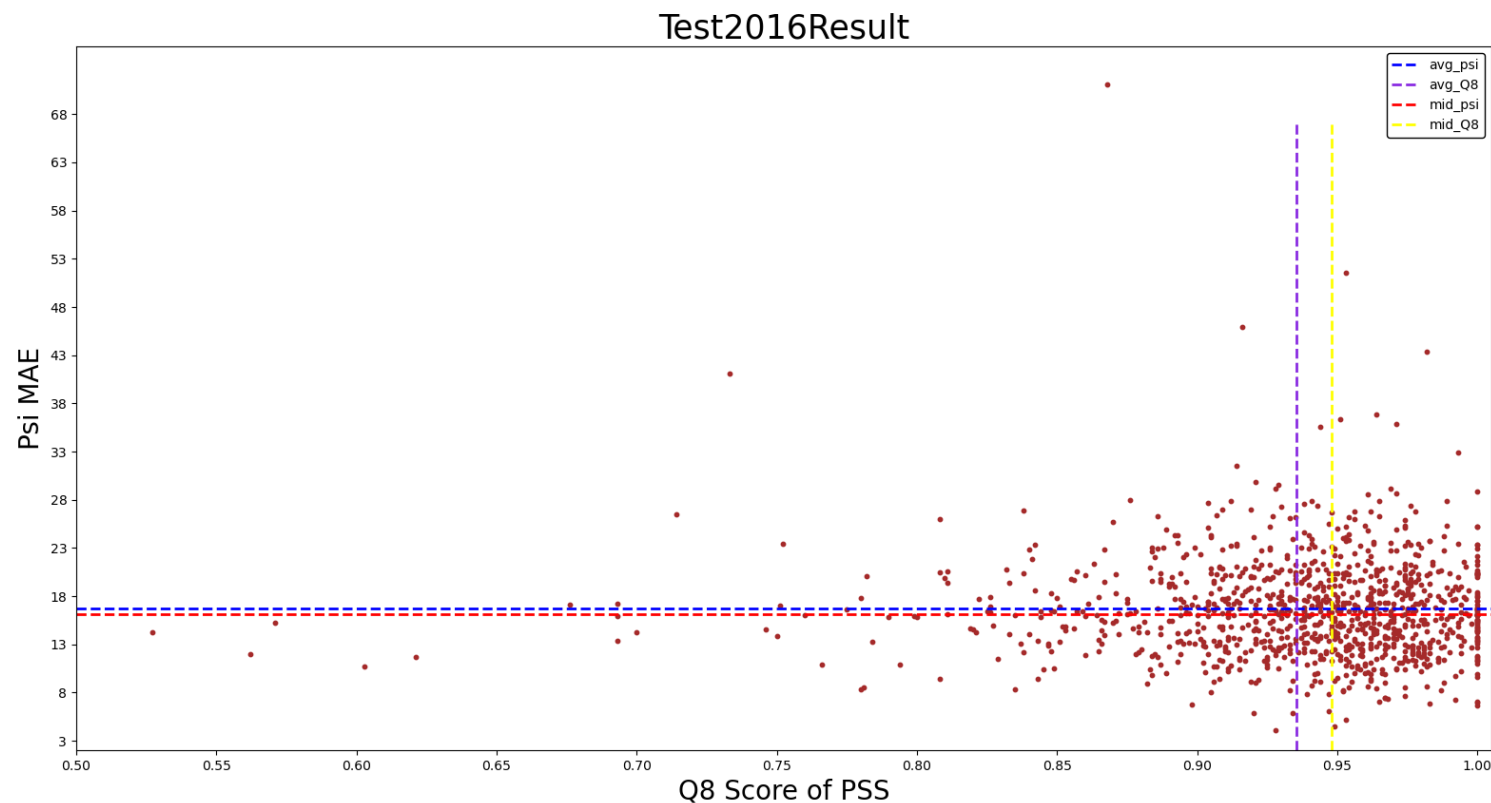
中位数比较结果:

都好 = 267

Q8好phi不好 = 252

Q8不好phi好 = 254

都不好 = 272



Q8好phi不好

phi: 该蛋白质中所有Q8预测不正确的残基的phi值的平均值

psi: 该蛋白质中所有Q8预测不正确的残基的psi值的平均值

phi_avg: 该蛋白质所有残基的phi值的平均值

psi_avg: 该蛋白质所有残基的psi值的平均值

共307个蛋白质

225个phi 大于 phi_avg

240个psi大于 psi_avg

实验序号	protein	phi	phi_avg
1	5wqwA	21.07	13.26
2	5a48A	18.25	14.00
3	5a6sB	5.37	14.59
4	5m3iB	2.17	11.15
5	5kx4B	22.84	9.18

Q8好psi不好

phi: 该蛋白质中所有Q8预测不正确的残基的phi值的平均值

psi: 该蛋白质中所有Q8预测不正确的残基的psi值的平均值

phi_avg: 该蛋白质所有残基的phi值的平均值

psi_avg: 该蛋白质所有残基的psi值的平均值

共286个蛋白质

206个phi 大于 phi_avg

223个psi大于 psi_avg

实验序号	protein	phi	phi_avg
1	5a4oB		
2			
3			
4			
5			

Q8不好phi好

phi: 该蛋白质中所有Q8预测不正确的残基的phi值的平均值
psi: 该蛋白质中所有Q8预测不正确的残基的psi值的平均值
phi_avg: 该蛋白质所有残基的phi值的平均值
psi_avg: 该蛋白质所有残基的psi值的平均值

共211个蛋白质
192个phi 大于 phi_avg
197个psi大于 psi_avg

实验序号	protein	phi	psi	Phi_avg	Psi_avg
1	5a4oB	12.82	22.55	11.22	17.02
2	5ucvB	26.26	28.96	13.09	14.50
3	5fvdA	31.25	50.45	12.92	19.27

4

Q8不好psi好

phi: 该蛋白质中所有Q8预测不正确的残基的phi值的平均值
psi: 该蛋白质中所有Q8预测不正确的残基的psi值的平均值
phi_avg: 该蛋白质所有残基的phi值的平均值
psi_avg: 该蛋白质所有残基的psi值的平均值

共228个蛋白质
210个phi 大于 phi_avg
212个psi大于 psi_avg

实验序号	protein	phi	psi	Phi_avg	Psi_avg
1	5a4oB	12.82	22.55	11.22	17.02
2	5ucvB	26.26	28.96	13.09	14.50
3	5fvdA	31.25	50.45	12.92	19.27
4					

Val (983)

实验序号	Model	MAE (Phi)	MAE (Psi)
1	SPOT-1D	16.27	22.53
2	OPUS-TASS	16.50	23.53
3	SFIPTA	15.39	18.25
4	SFIPTA2	15.42	18.25

Test2016 (1212)

实验序号	Model	MAE (Phi)	MAE (Psi)
1	spider3	17.88	26.66
2	SPOT-1D	16.27	23.26
3	OPUS-TASS	15.78	24.46
4	ESIDEN	15.61	19.30
5	SFIPTA	15.18	18.14
6	SFIPTA2	15.18	18.08

Test2018 (250)

实验序号	Model	MAE (Phi)	MAE (Psi)
1	spider3	18.38	28.10
2	SPOT-1D	16.89	24.87
3	OPUS-TASS	16.40	24.06
4	ESIDEN	17.09	21.70
5	Mine	15.65	18.62
5	Mine2	15.74	18.64

CASP12 (55)

实验序号	Model	MAE (Phi)	MAE (Psi)
1	spider3	21.14	34.92
2	SPOT-1D	20.21	31.71
3	ESIDEN	19.68	28.63
4	Mine	17.63	23.29
5	mine2	17.67	23.70

CASP13 (32)

实验序号	Model	MAE (Phi)	MAE (Psi)
1	spider3	22.48	38.46
2	SPOT-1D	22.60	34.28
3	ESIDEN	22.51	32.68
4	Mine	16.69	22.18
4	Mine2	16.81	22.46

CASP-FM (56)

实验序号	Model	MAE (Phi)	MAE (Psi)
1	NetSurfP-2.0	19.94	31.43
2	SPOT-1D	19.39	30.1
3	OPUS-TASS	18.85	28.00
4	Mine	18.33	25.17
4	Mine	18.30	25.53

验证单个特征及其组合的有效性

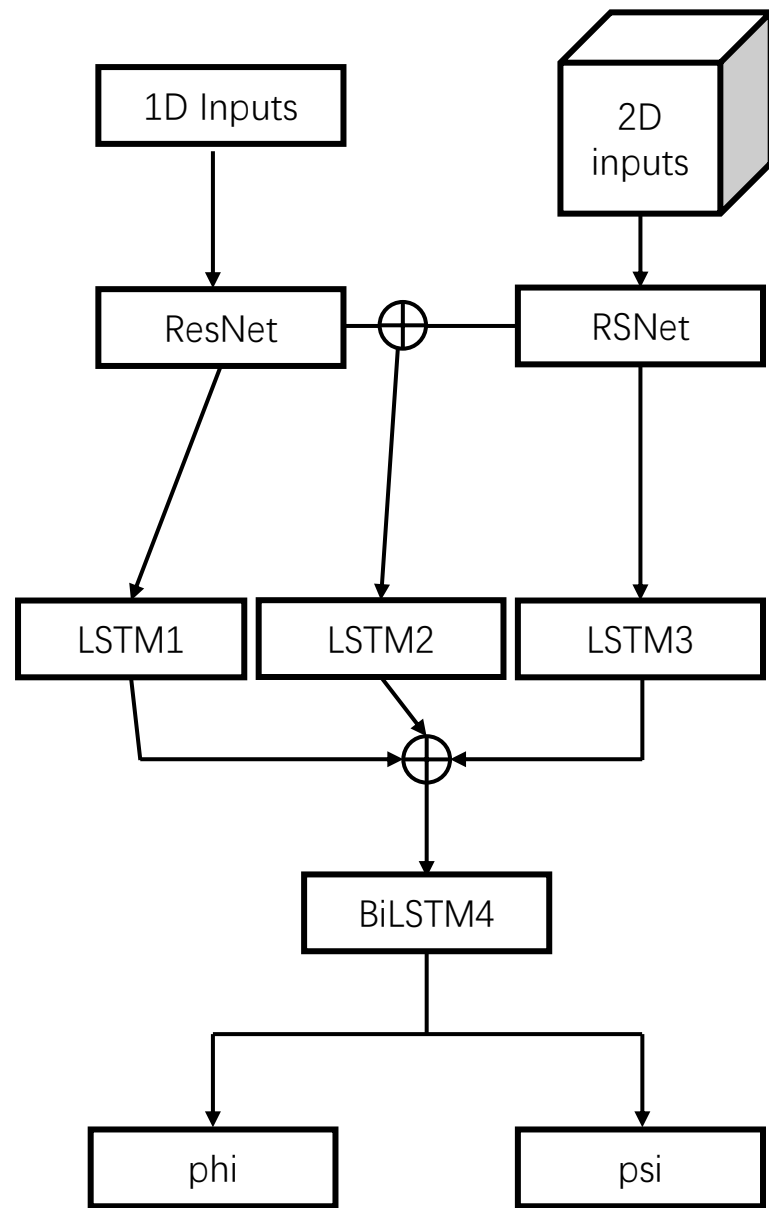
Model	Val_Phi	Val_Psi	Test2016_Phi	Test2016_Psi	epoch
Basic = pssm+hhm+pp	20.94	32.86	18.71	31.43	14
Basic+psfm	20.28	33.16	17.80	27.48	21
Basic+onehot	19.28	30.55	18.16	27.41	22
Basic+pc	21.19	32.51	18.87	27.71	24
Basic+hyd	20.62	32.55	18.31	27.36	16
Basic+psfm+pc	20.24	31.56	18.49	27.21	14
Basic+onehot+pc	21.57	34.56	18.53	27.50	23
Basic+psfm+pc+hyd	19.65	32.31	18.05	28.04	25
Basic+onehot+pc+hyd	19.78	29.35	18.67	27.41	12
Basic+psfm+pc+hyd+ss8	16.34	19.66	15.62	18.89	14
Basic+onehot+pc+hyd+ss8	15.95	19.66	15.09	18.72	16
Basic+psfm+onehot+pc+hyd+ss8	16.09	19.11	15.57	18.56	10

验证单个特征及其组合的有效性（去掉增强）

Model	Val_Phi	Val_Psi	Test2016_Phi	Test2016_Psi	epoch
Basic = pssm+hhm+pp	20.13	32.70	18.19	27.34	30
Basic+psfm	20.86	34.87	18.84	28.31	24
Basic+onehot	22.68	39.71	17.98	27.19	11
Basic+pc	24.75	42.94	20.70	30.59	
Basic+hyd	21.72	34.60	18.48	27.41	11
Basic+psfm+pc（集群）	20.55	34.36	18.01	27.64	11
Basic+onehot+pc	19.38	29.55	17.77	25.97	17
Basic+psfm+pc+hyd（D123）	22.05	38.76	18.19	27.29	13
Basic+onehot+pc+hyd(集群)	20.92	36.49	17.83	27.36	16
Basic+psfm+pc+hyd+ss8（集群）	16.36	19.23	15.84	18.81	11
Basic+onehot+pc+hyd+ss8	16.63	19.35	15.60	18.38	6
Basic+psfm+onehot+pc+hyd+ss8 （集群）	16.05	18.99	15.54	18.53	10

验证窗口大小

winsize	Val_Phi	Val_Psi	Test2016_Phi	Test2016_Psi	epoch
3 (jiqun)	15.72	17.93	15.82	17.90	
5 (jiqun)	15.41	17.62	15.32	17.10	
7 (D123)					
9 (jiqun)	15.36	17.47	15.36	17.71	7
11					
13					
15					
17					
19					
21					



目前存在的问题

- 模型在所有测试集上的预测结果MAE值都比其他预测器要好，但在验证集上 ϕ 角的MAE值不理想,并且在验证集上收敛非常快。

实验序号	Model	MAE (Phi)	MAE (Psi)
1	SPOT-1D	16.27	22.53
2	OPUS-TASS	16.50	23.04
3	SFIPTA	19左右	22.5左右

- 造成这个结果的原因可能是模型泛化能力不够，于是采用数据增强，改变模型框架（加入dropout等）和参数的方式，尝试增强泛化能力

数据增强

- 什么是数据增强？作用是什么？
- 数据增强实在不实质性增加数据的情况下,从原始数据中加工出更多的数据表示，提高原数据的数量及质量，以接近更多于更多数据量产生的价值
- 我采用的是随机裁剪的方式进行数据增强，其原理是：

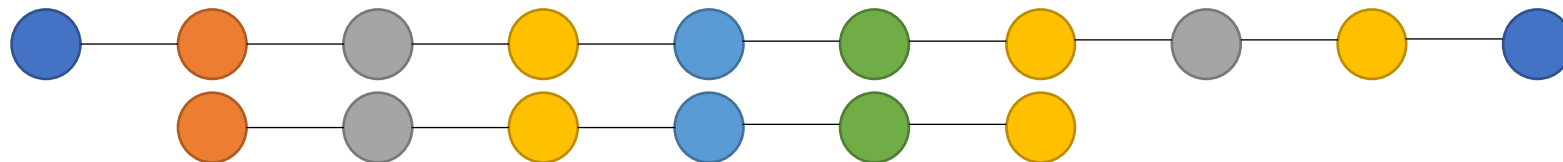
在每一次迭代（epoch）中，将每一条蛋白质进行随机裁剪，使得对于同一条蛋白质，在不同epoch中，有着不同的氨基酸序列（等于一条蛋白质衍生出N个epoch个数的子蛋白质）。

随机裁剪

原序列



epoch1



epoch2



•
•
•

•
•
•

epochN

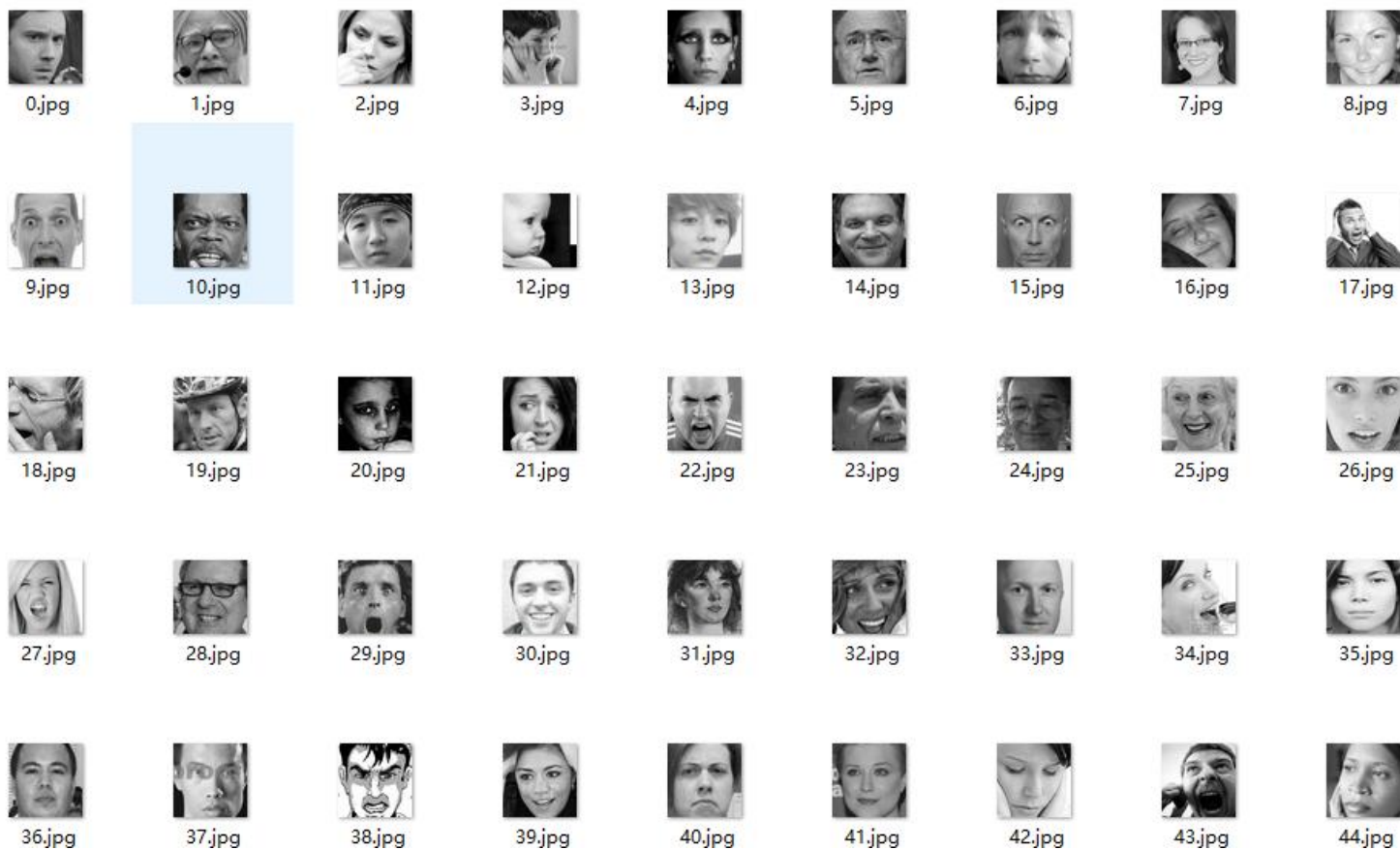


基于深度人脸语义识别的 学评教自动评价研究

- 任务：通过学生上课时的人脸表情状态，判断学生是否有认真听课。
- 做法：用fer2013数据集预训练一个7分类的模型（生气、厌恶、害怕、开心、伤心、惊讶、中性），将其预测结果以onehot的形式作为学评教模型的输入，经过网络得到二分类输出，判断该生是否认真听课。

Fer2013数据集 (28079train, 3589val, 3589test)

(生气、厌恶、害怕、开心、伤心、惊讶、中性)



自制数据集 (100张)



29.jpg



30.jpg



31.jpg



32.jpg



33.jpg



34.jpg



35.jpg



36.jpg



37.jpg



38.jpg



39.jpg



40.jpg



41.jpg



42.jpg



43.jpg



44.jpg

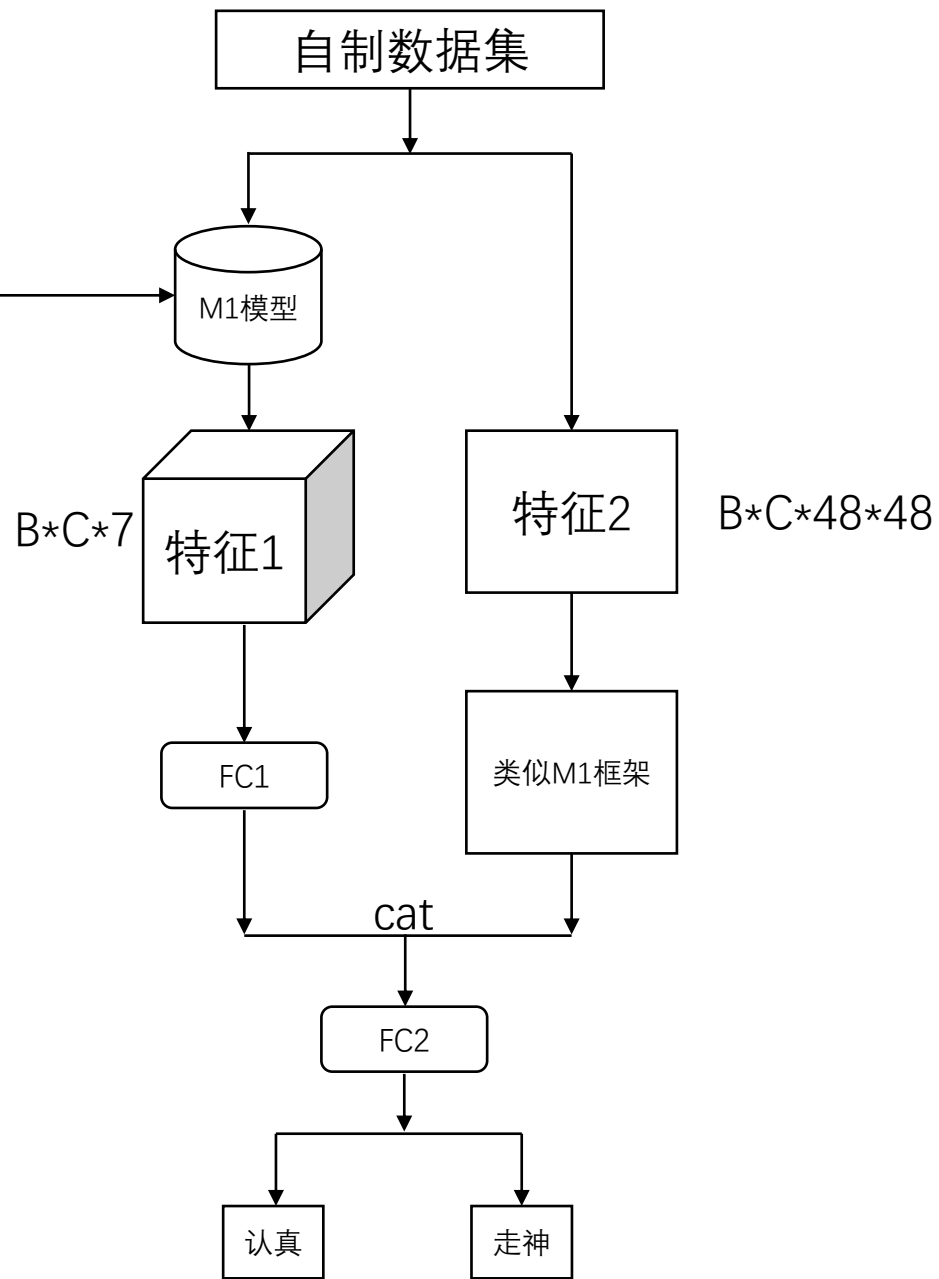
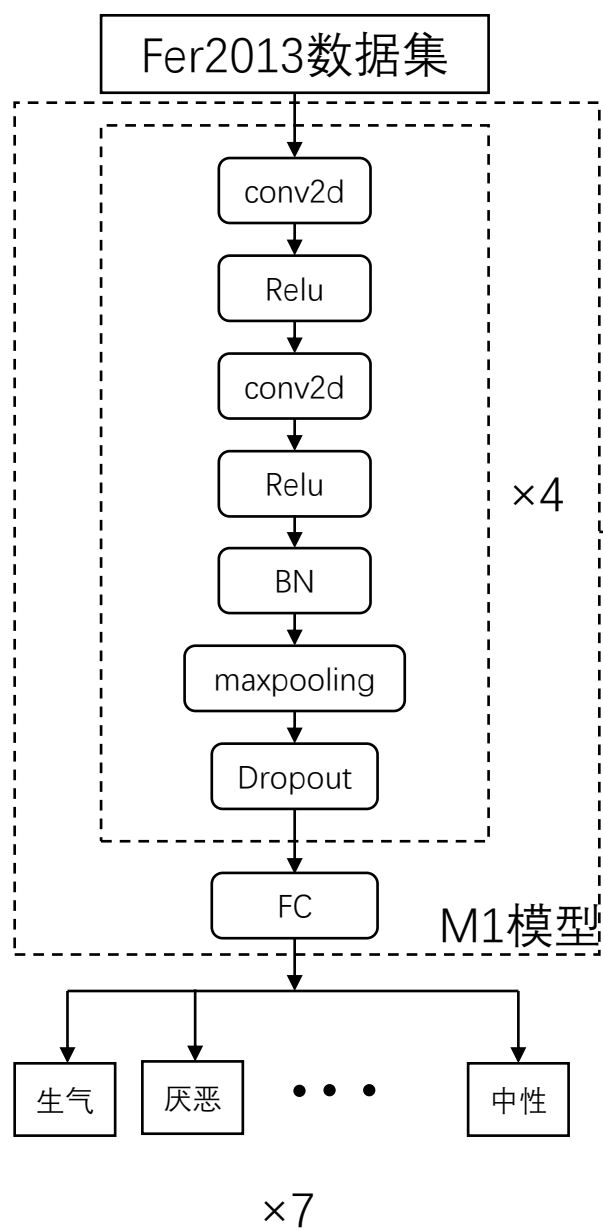


45.jpg



46.jpg





打算

根据