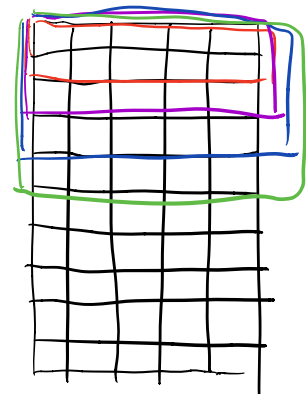


G output

16

20



不同 filter size

概率

seq-size

Vocab-size

a	0.1	0.2	0.3	0.12	...	...
b						
c						
d						
e						

简单示意图

seq-size

Vocab-size

one-hot

0	0	1	0	0

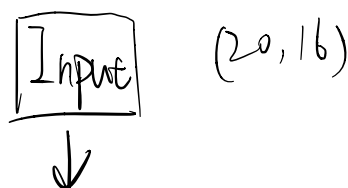
这里不考虑 batch, 假设矩阵 shape 为 (20, 16)

seq-size

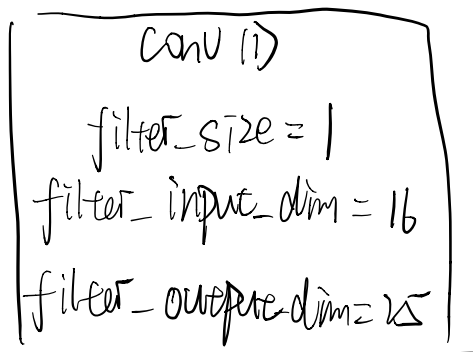
Vocab-size

DIM = 25 为 feature 数

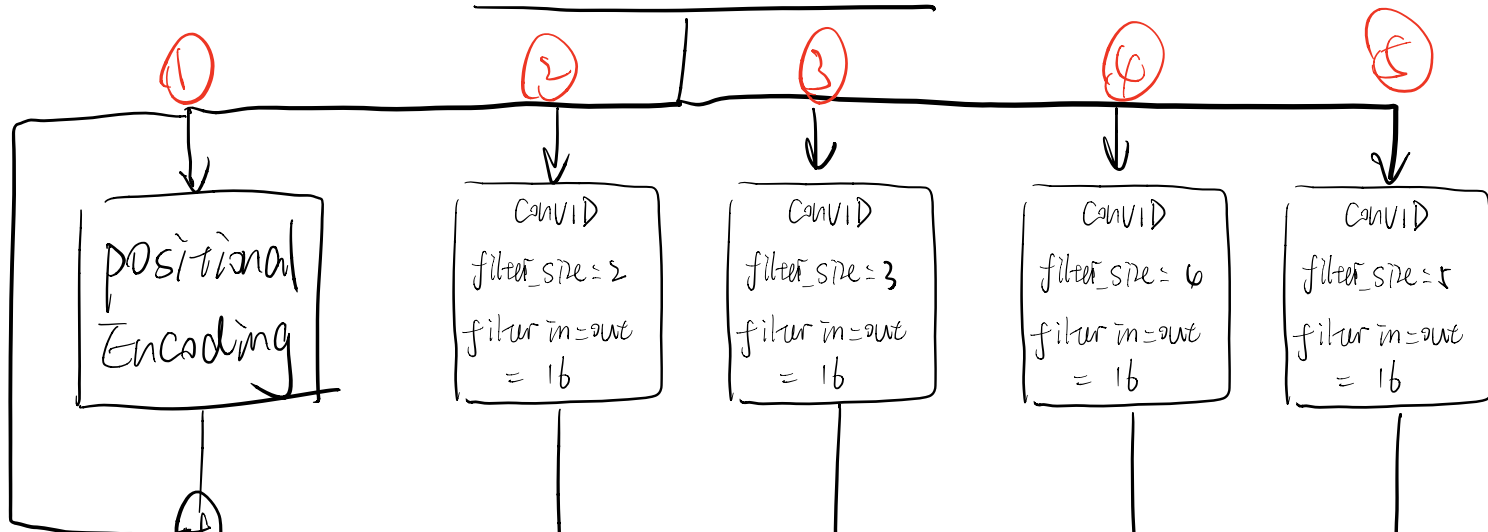
Discriminator

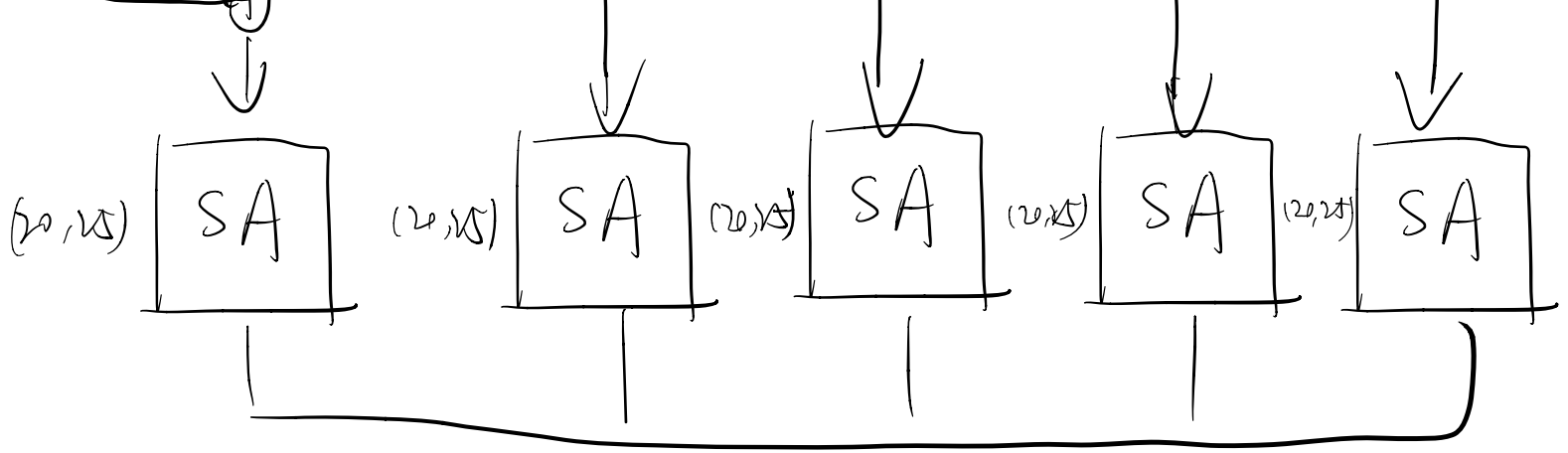


(20, 25)



将概率重新转换到语义特征矩阵上来





① 跟传统 SA 一样  
捕捉同一序列内词之  
词级比需要关注的  
权重 (positional Encoding 块词序)

②③④⑤ 用 Conv1D 的 filter  
size 为 5 长 (类似 skip gram  
思想) 捕捉相邻词的关系。  
但是这种思想需要多层 CNN  
叠加才能捕捉远距离特征。  
在后面加 SA 层用来弥补这个  
缺陷, 可以捕捉类似消息  
帧中位置  $i$  到位置  $(i + \text{filter\_size})$   
部分的数据与位置  $j$  到位置  
 $(j + \text{filter\_size})$  的关系。

这里 filter 取 2, 3, 4, 5 涵盖大部消息  
头的长度。Conv1D 是根据 stride 大小遍历

concat (100, 25)

linear

scalar

的, 故我觉得内含了顺  
序信息, 所以不需要  
positional Encoding。

① 其实以理解为 filter  
size 为 1 的 Conv1D, 但  
我觉得应该比 Conv1D 效  
果好。最后 concat  
起来, 保证 filter size  
1-5 每个占的比例一致  
linear 之后得到 一个  
scalar 表示权重。

Q1. SA 是否需要 multi-head 机制?  
head 数似乎会影响远距离特征  
捕捉能力, 但是我们的数据长度  
其实不长。