



UTNet: 基于UNet和Transformer架构 的医学图像分割探索



展示成员：鹿逸文



目录

PART 01 项目背景与意义

PART 02 主要贡献

PART 03 数据处理

PART 04 模型构建

PART 05 训练技巧

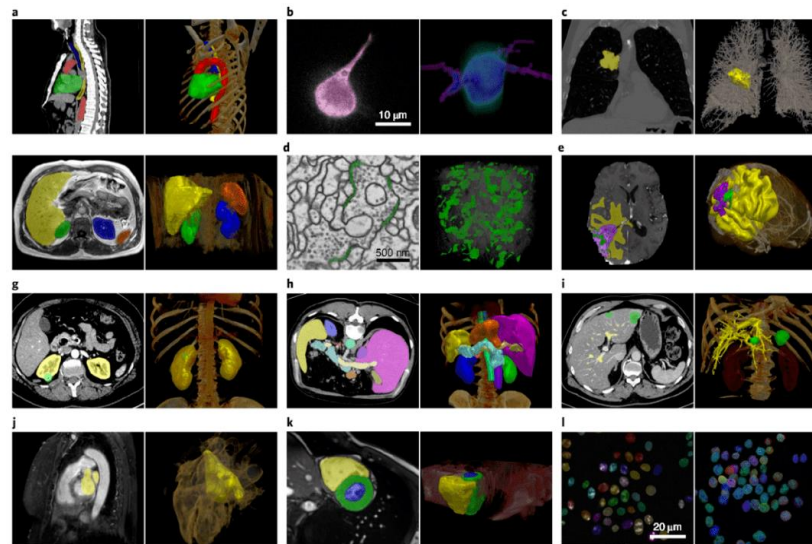
PART 06 实验



1. 项目背景与意义

项目背景

- 医学图像分割是医学领域中至关重要的任务，它在诊断、治疗规划和疾病监测等方面发挥着关键作用
- 传统的医学图像分割受限于主观性，且人力成本较高
- Transformer在自然语言处理领域取得了显著成就，但在医学视觉领域的应用尚未深入探索。



意义

UTNet有望为疾病的早期检测、治疗规划和病情预测提供有力的支持，提高医疗诊断的精确性和效率，最终改善患者的治疗结果和生存率。



2. 主要贡献



整合了卷积和注意力策略的优势



提出了一种高效的自注意机制



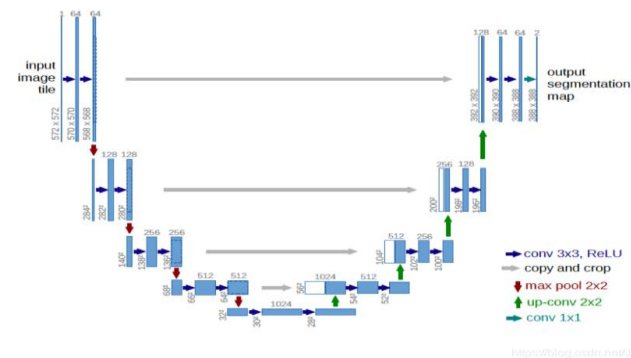
引入了全新的自注意解码器

U-Net网络

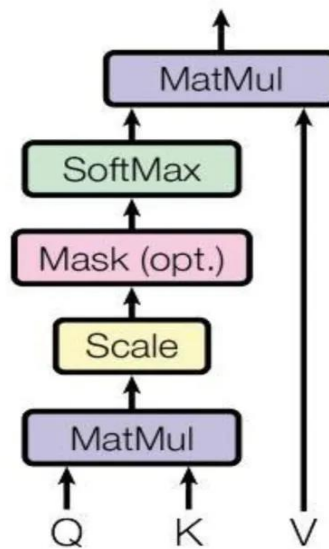
From FCN to U-Net

Architecture of U-Net

- U-Net
 - 采用编码器和解码器的U形结构
 - 输入输入大小不变
 - Skip 结合方式: Concatenation



Scaled Dot-Product Attention

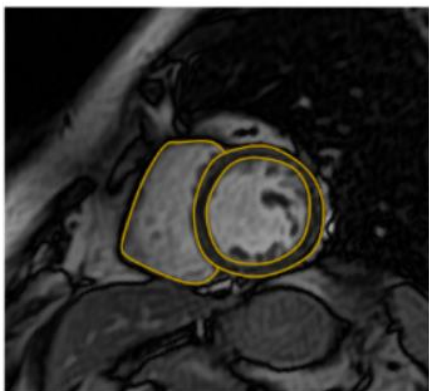




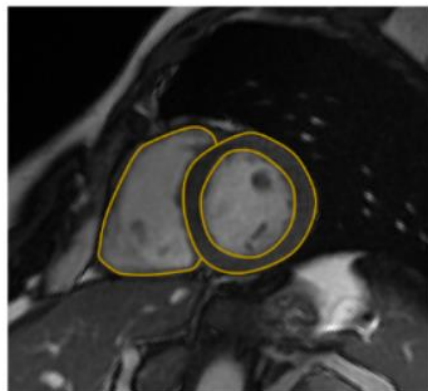
3. 数据处理

数据集选择

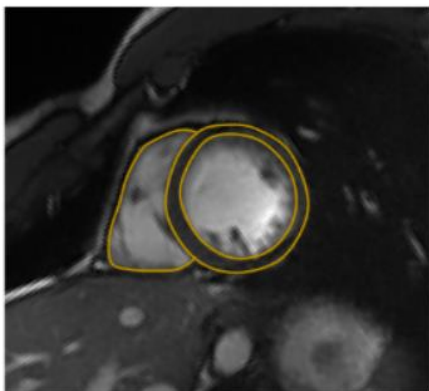
Multi-Centre, Multi-Vendor & Multi-Disease Cardiac Image Segmentation Challenge



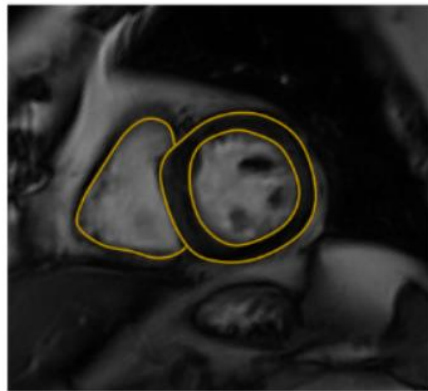
(1) Canon



(2) General Electric



(3) Philips



(4) Siemens

所用的MRI图像是心脏电影成像序列，是带时间维度的三维心脏图像。

来源

数据集来源于多中心，多供应商和多疾病心脏图像分割挑战赛。由375名肥厚性和扩张型心肌病患者以及健康受试者组成。



标签

训练集包含150个带标注的图像。由来自各个机构经验丰富的临床医生对CMR图像进行了分割，包括左（LV）和右心室（RV）血池以及左心室心肌（MYO）的轮廓。标签为：0（背景），1（LV），2（MYO）和3（RV）。





3. 数据处理

数据预处理

01

缩放

- 将所有数据重新采样为 x-y 平面中 1.2x1.2 mm 的间距

02

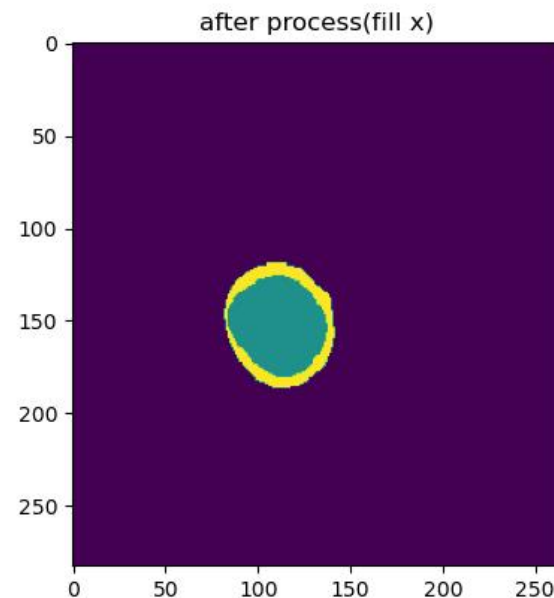
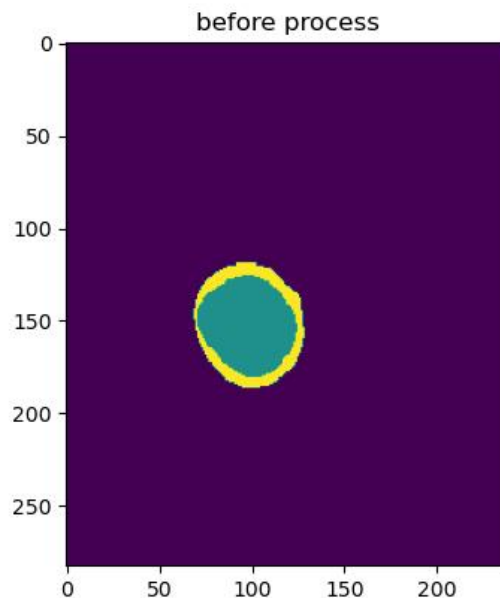
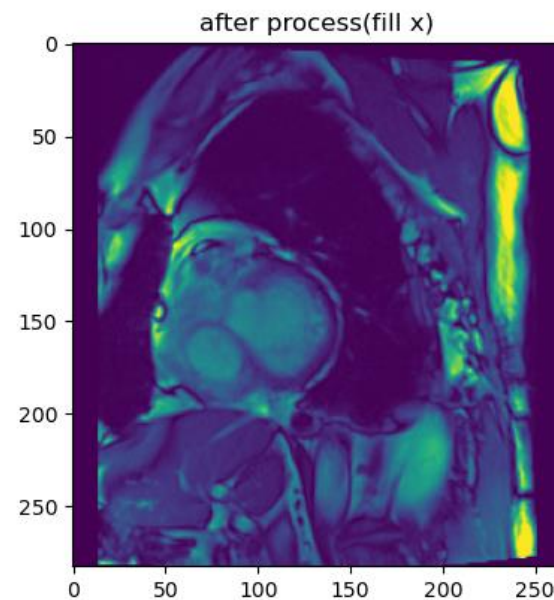
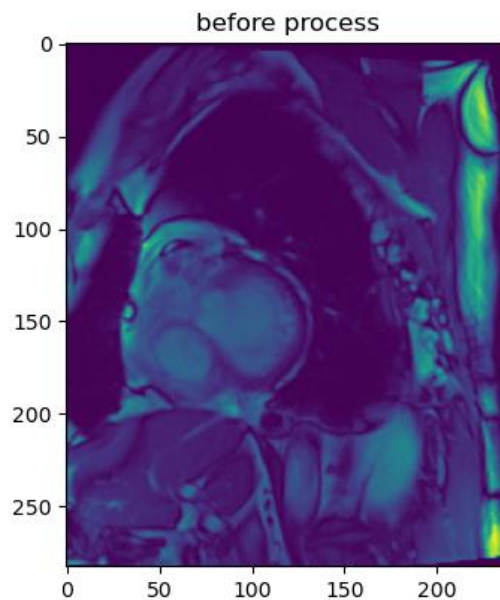
截断

- 对缩放后的原始数据进行异常值截断处理，将98%以上的灰度值和5%以下的灰度值进行截断

03

归一化

- 采用均值为0，方差为1的方式对原始图像进行归一化处理





3. 数据处理

数据增强

添加高斯噪声

1

3

Gamma校正
对比度

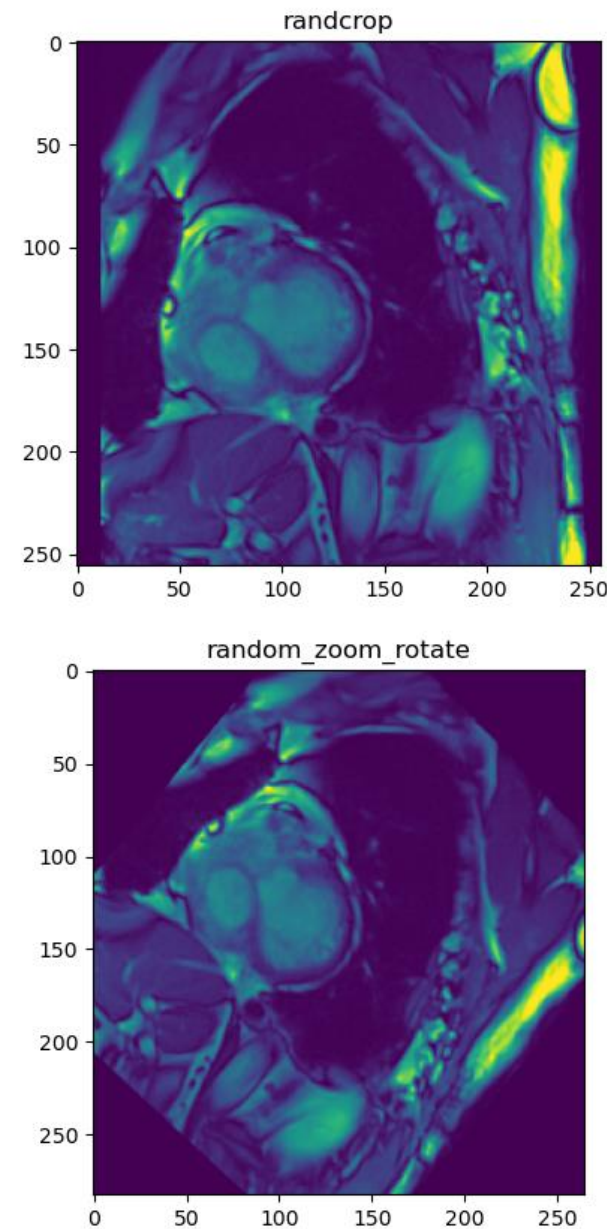
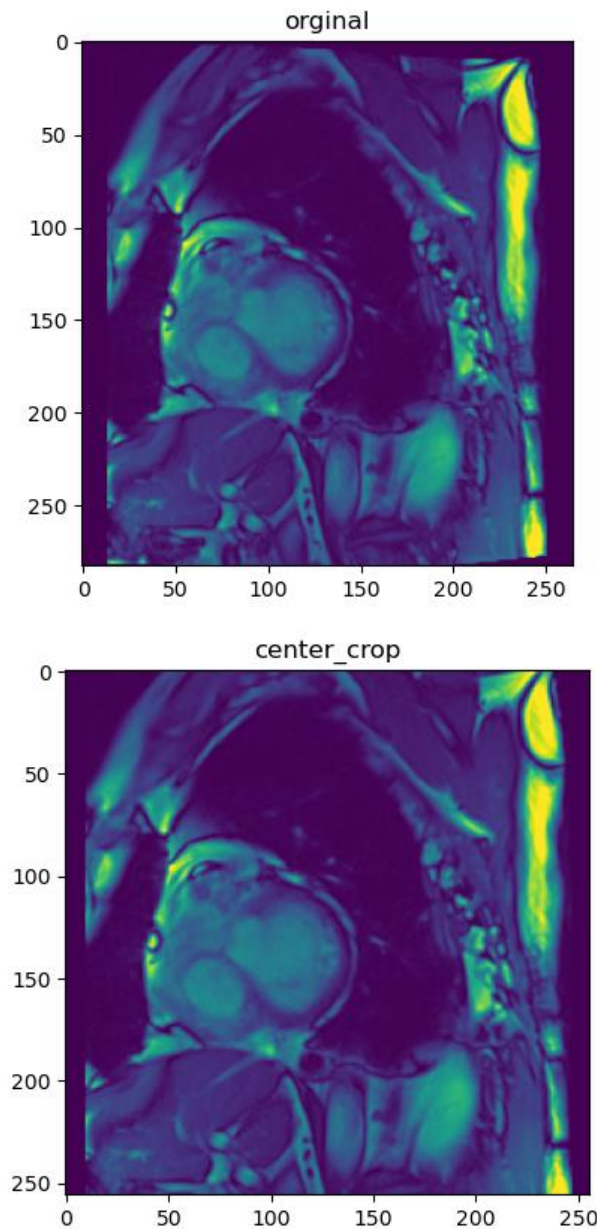
数据增强

亮度随机
调整

2

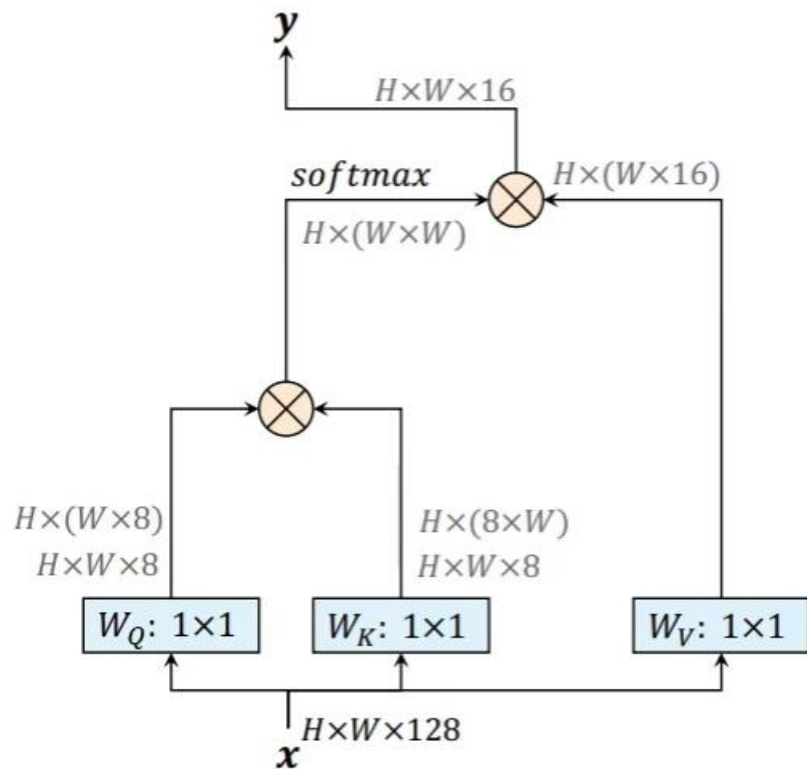
4

随机缩放、
旋转、裁剪



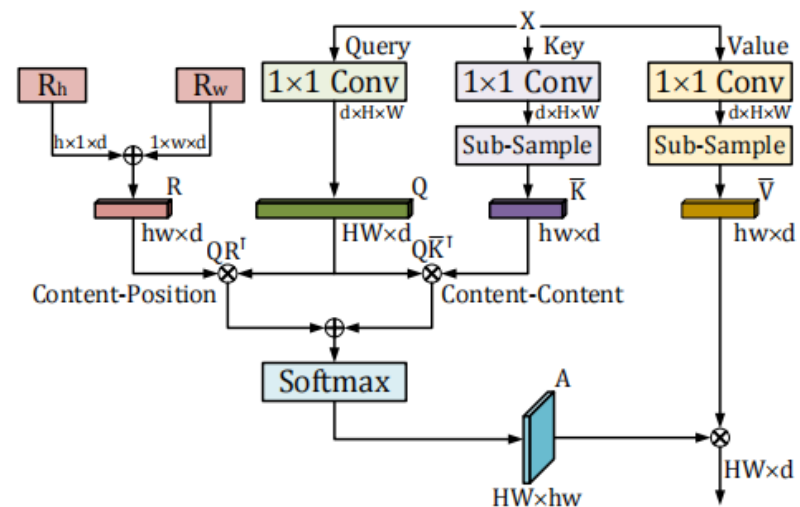


4. 模型构建：高效自注意力机制

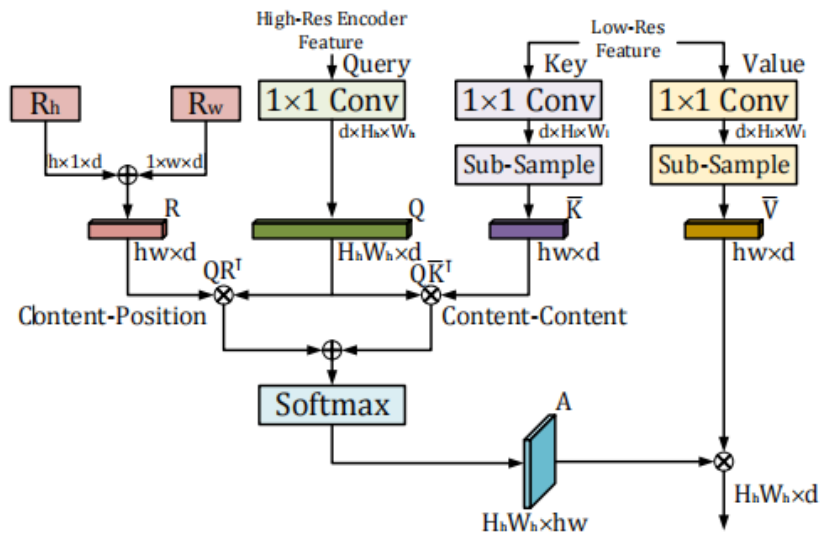


原生注意力机制

$$\text{Attention}(\mathbf{Q}, \mathbf{\bar{K}}, \mathbf{\bar{V}}) = \underbrace{\text{softmax}\left(\frac{\mathbf{Q}\mathbf{\bar{K}}^T}{\sqrt{d}}\right)}_{\mathbf{\bar{P}}: n \times k} \underbrace{\mathbf{\bar{V}}}_{k \times d}$$



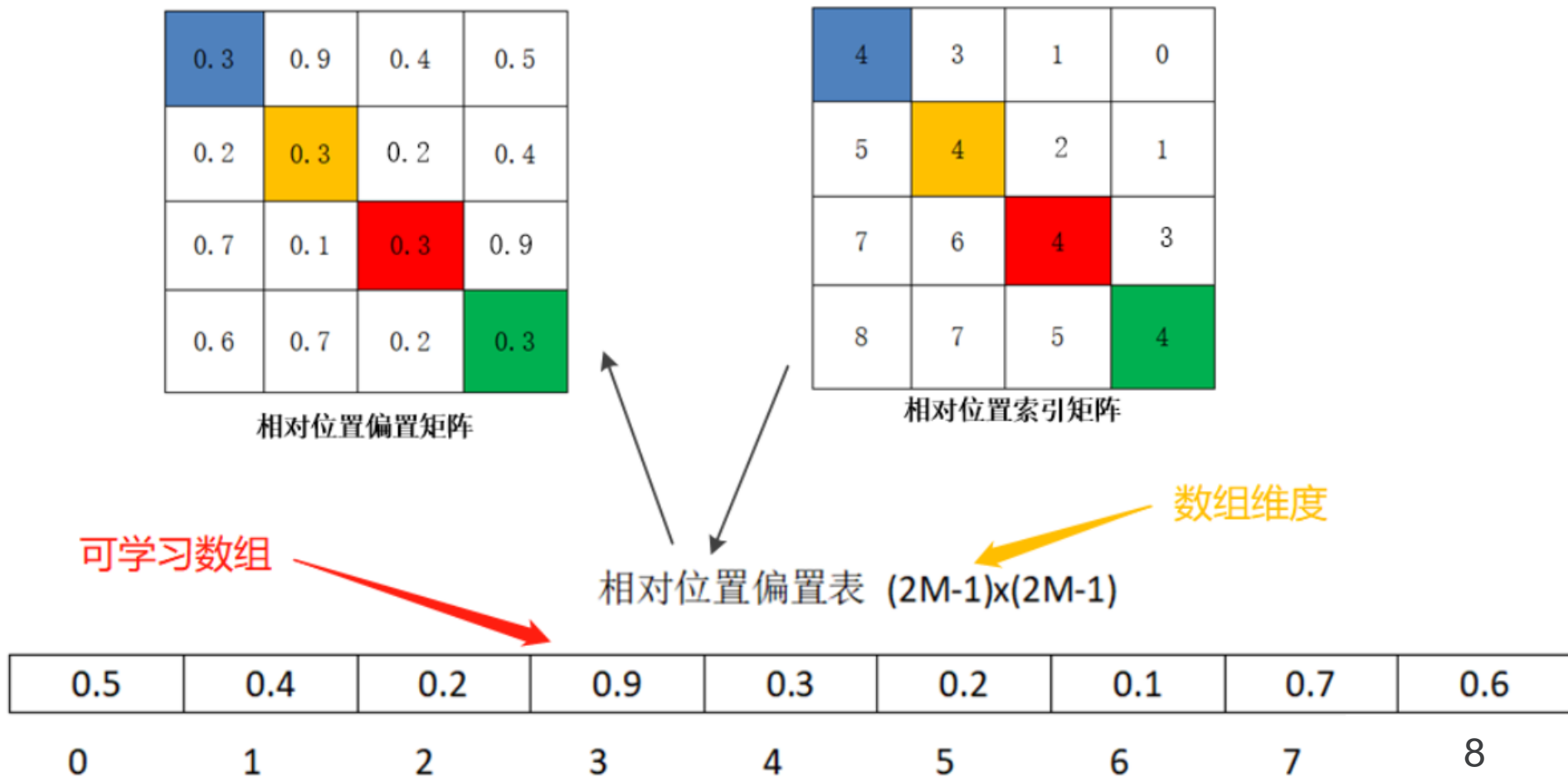
多头注意力编码器



多头注意力解码器



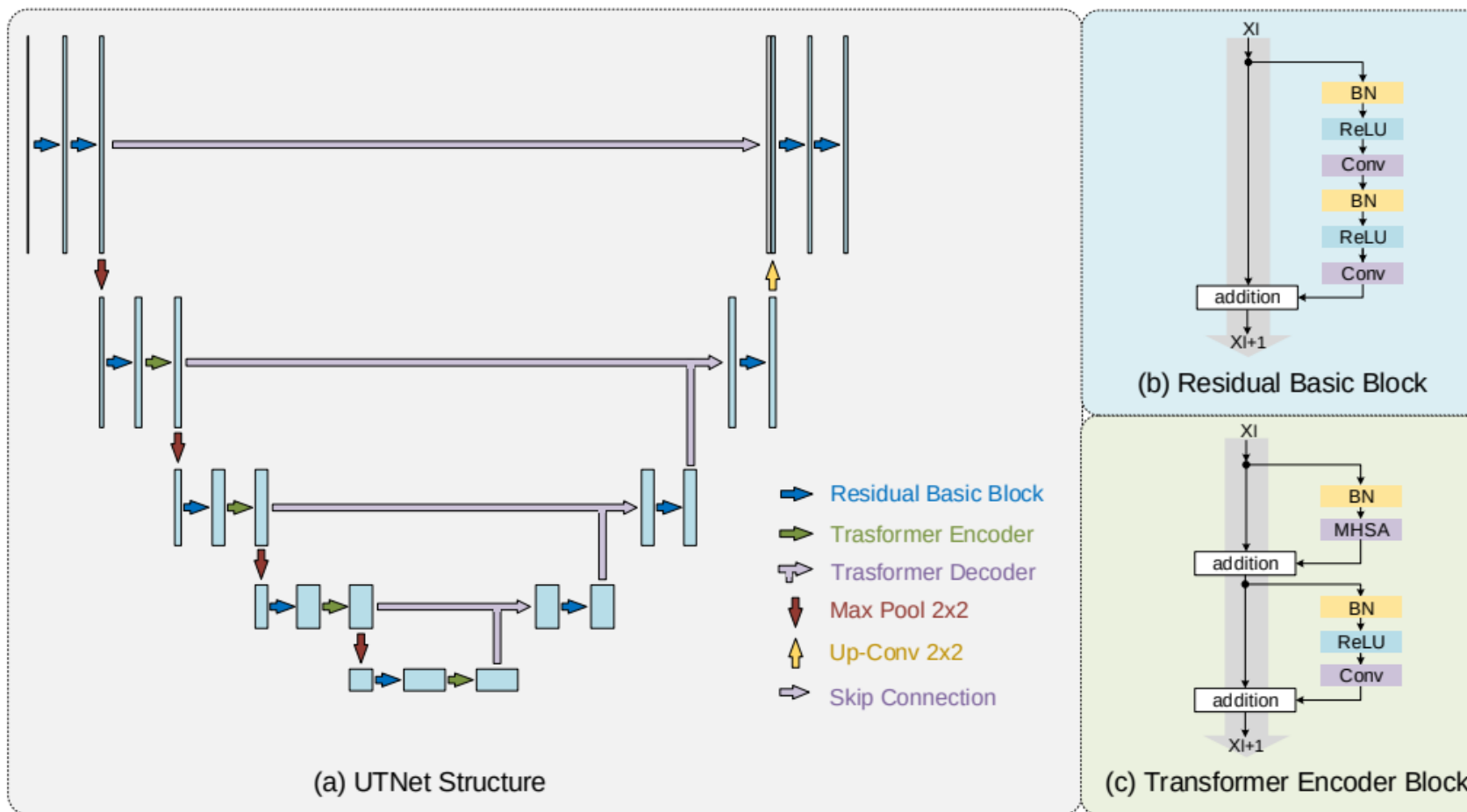
4. 模型构建：相对位置编码



$$\text{Attention}(\mathbf{Q}, \mathbf{\bar{K}}, \mathbf{\bar{V}}) = \text{softmax}\left(\underbrace{\frac{\mathbf{Q}\mathbf{\bar{K}}^T + \mathbf{S}_H^{rel} + \mathbf{S}_W^{rel}}{\sqrt{d}}}_{\mathbf{\bar{P}}: n \times k}\right) \underbrace{\mathbf{\bar{V}}}_{k \times d}$$



4. 模型构建: UNet+Attention



UNet架构

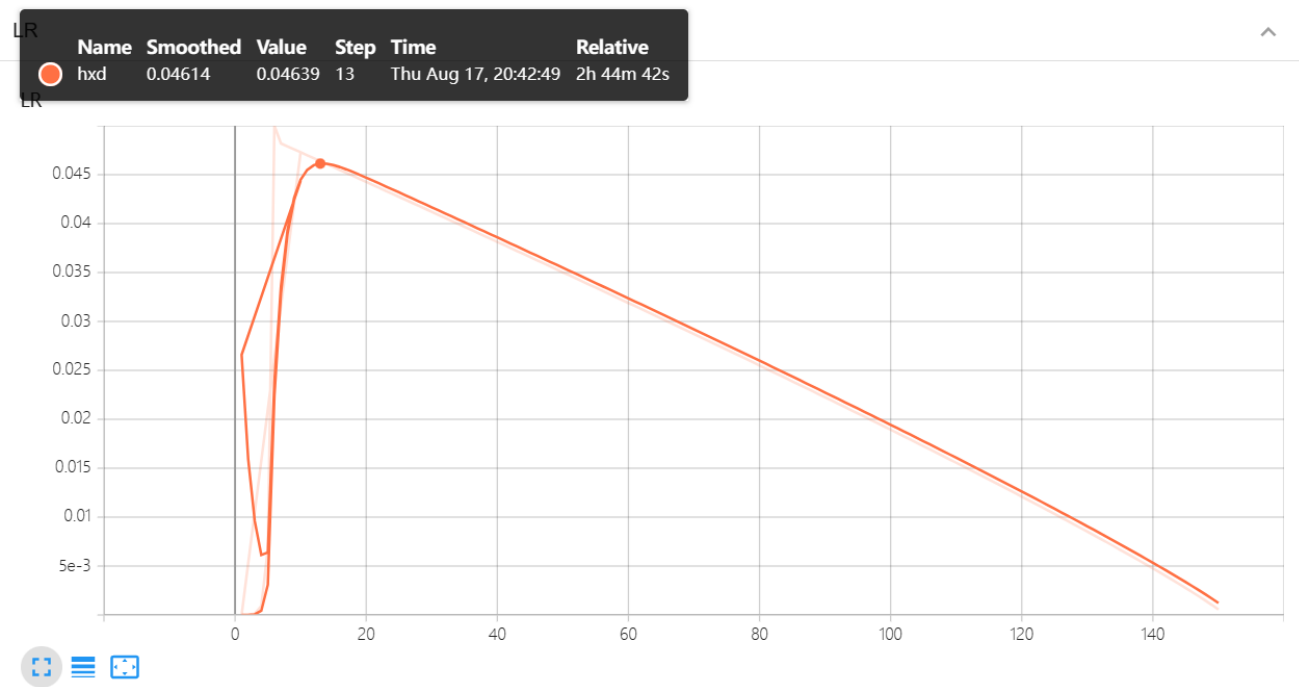


5. 训练技巧:学习率调度器

1 学习率的大小很重要

2 衰减速率同样很重要

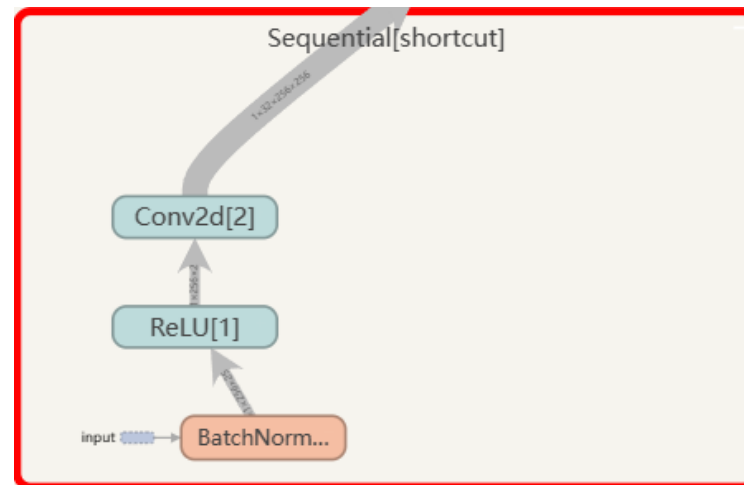
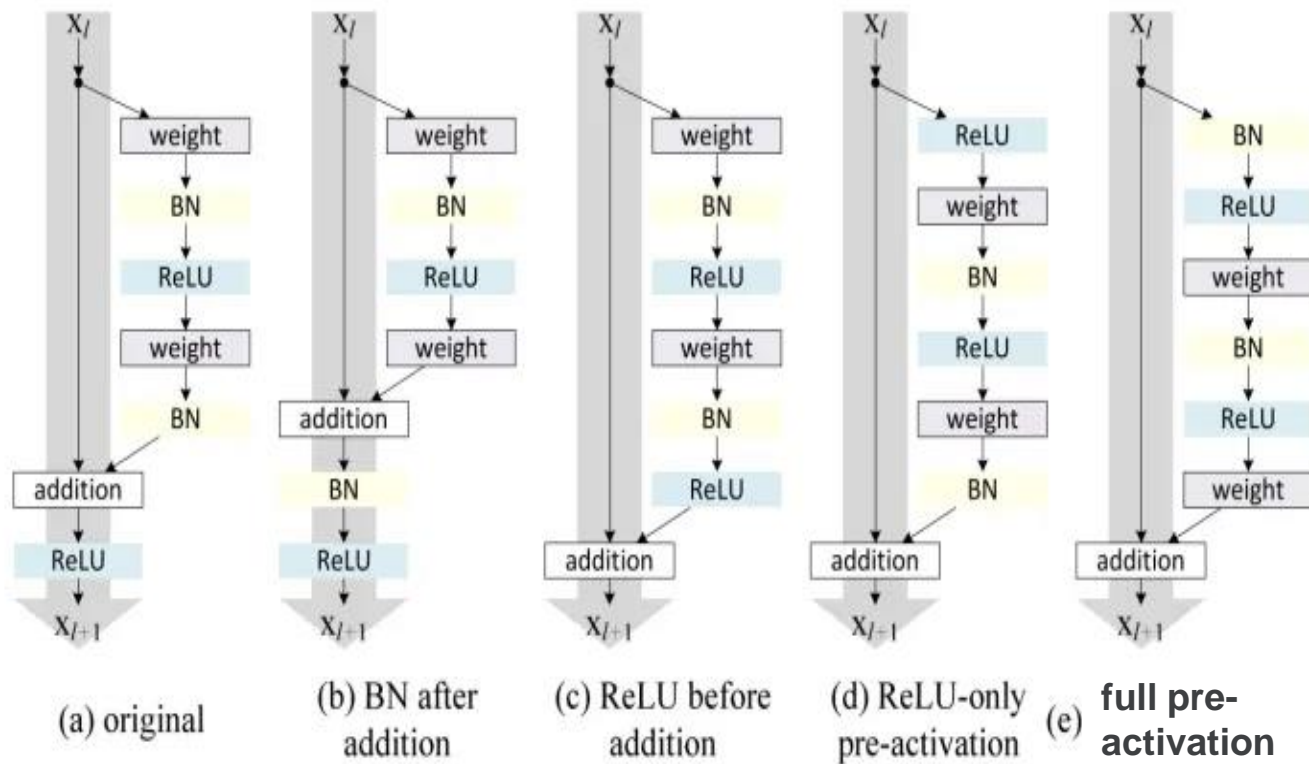
3 预热阶段



训练过程中**学习率**变化



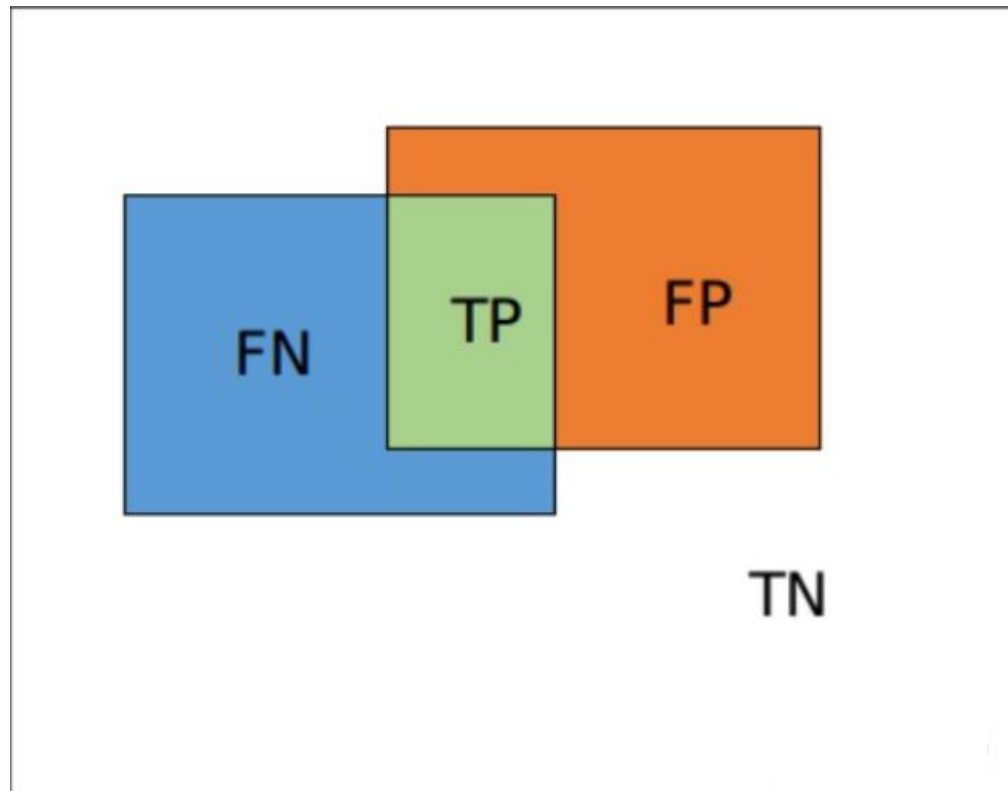
5. 训练技巧:pre-activation



UNet中的残差部分

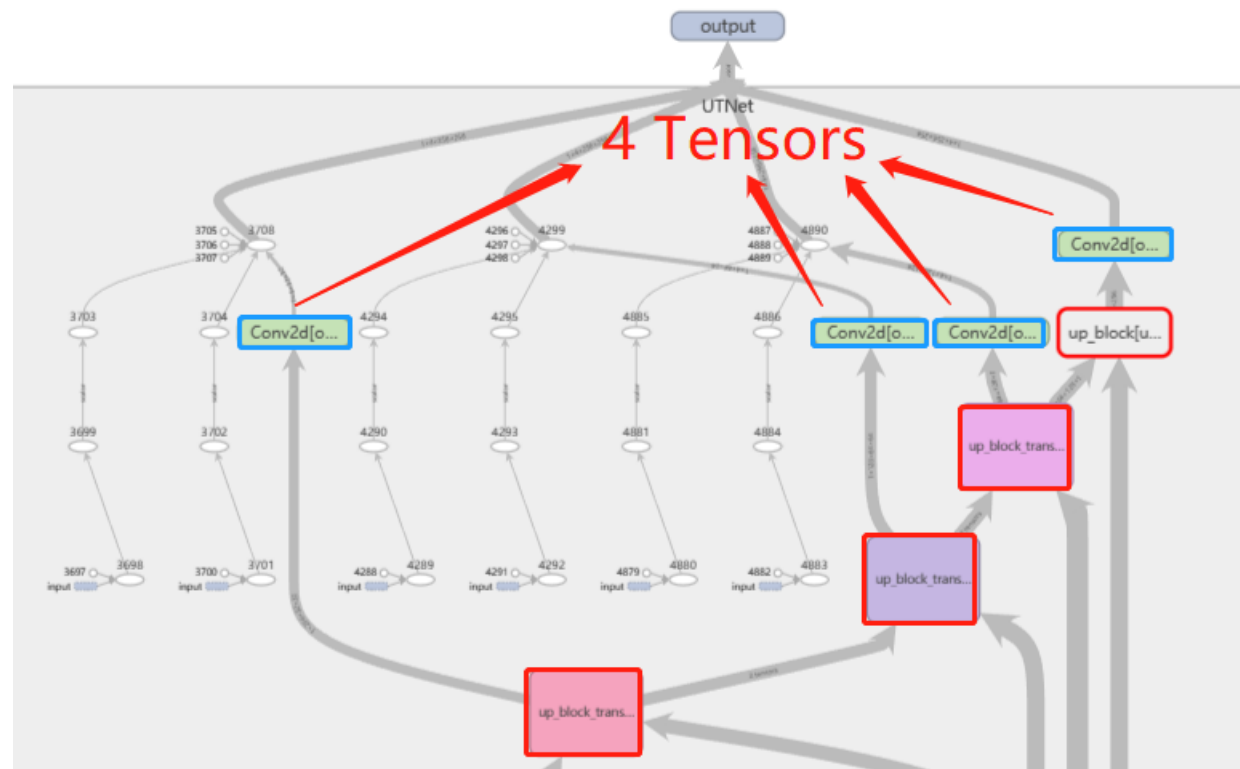


5. 训练技巧:损失函数



$$dice = \frac{2TP}{2TP+FP+FN}$$

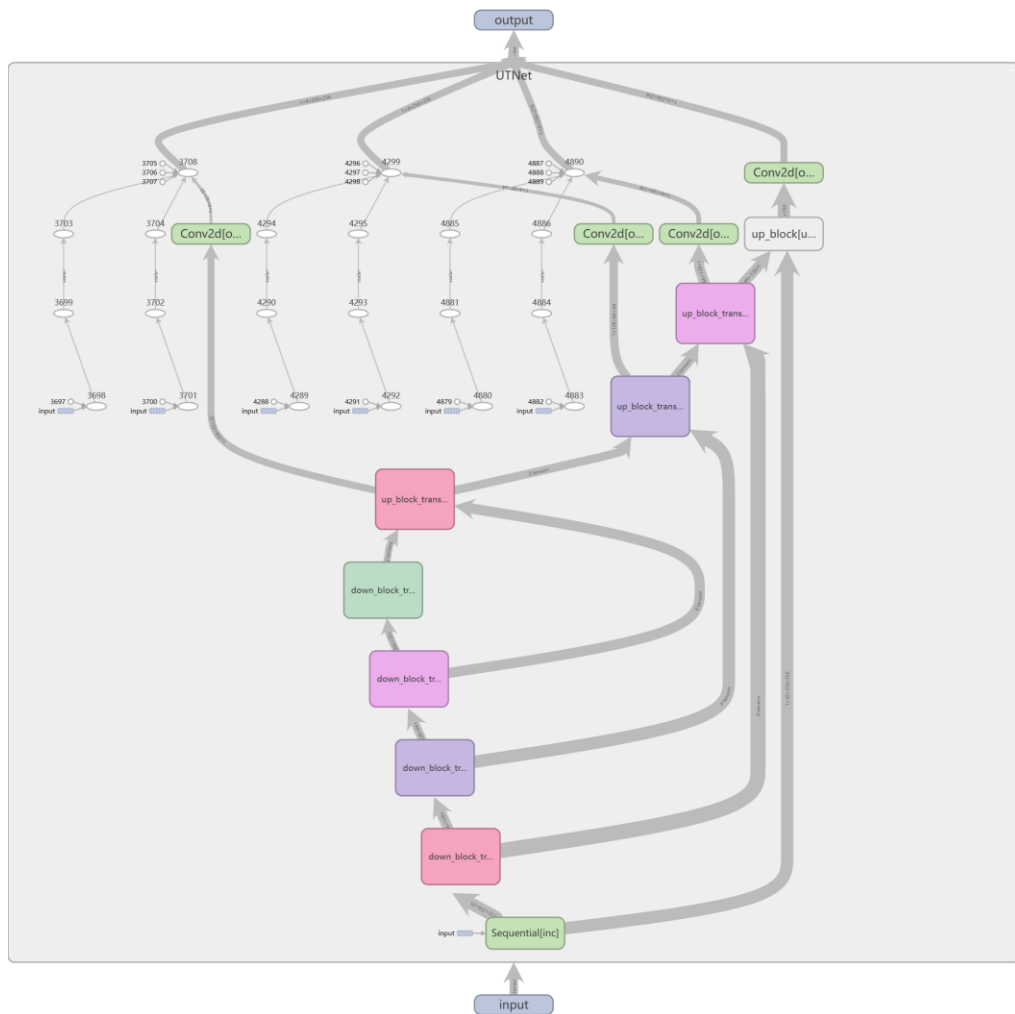
$$DiceLoss = 1 - Dice$$



Aux_loss



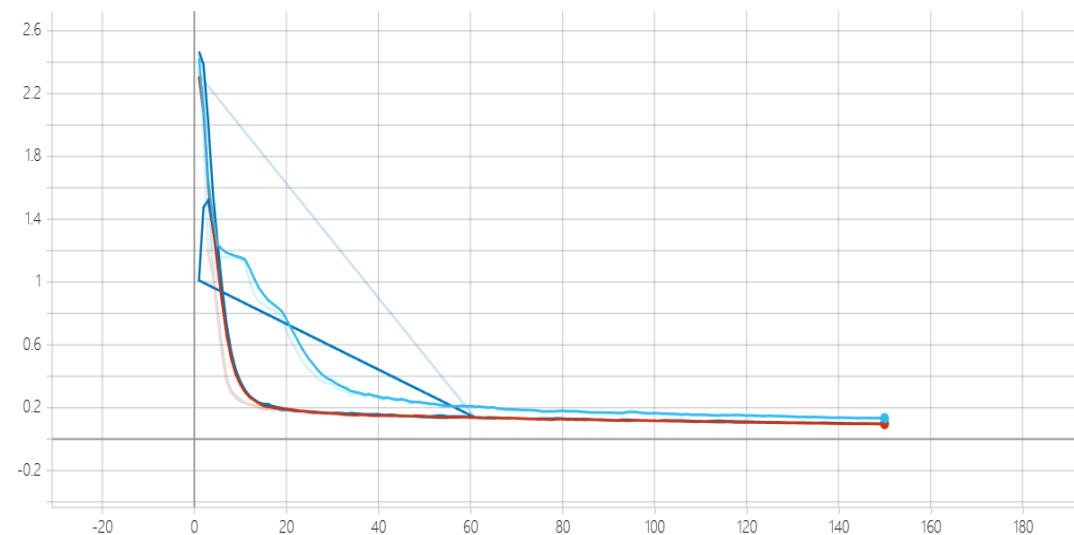
6. 实验：模型训练



UTNet的TensorBoard可视化

- ✓ ○ UTNet
- ✓ ○ TransUNet
- ✓ ○ ResNet50-UNet

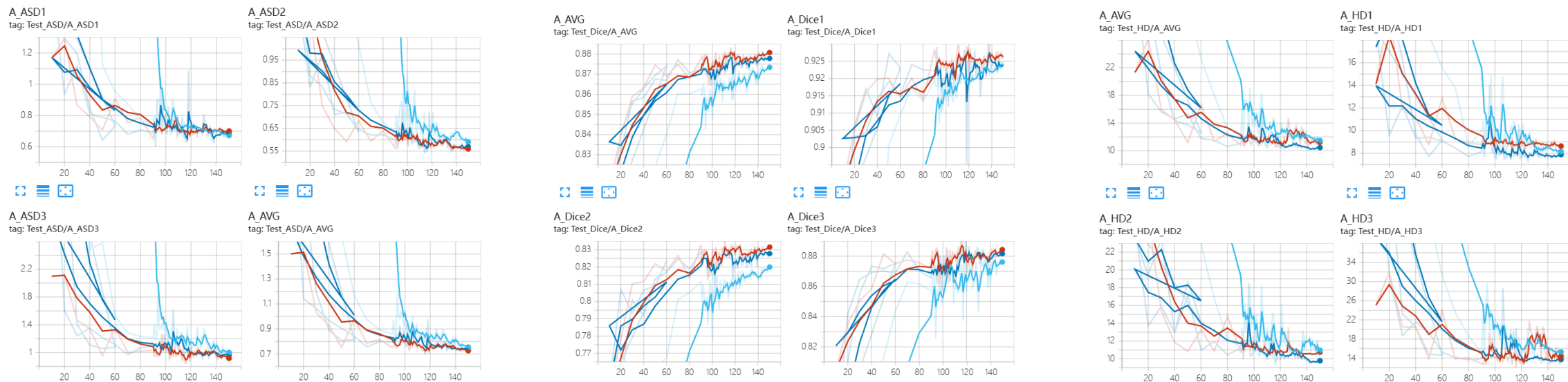
Loss
tag: Train/Loss



Train_loss



6. 实验：模型测试--指标评估



平均表面距离

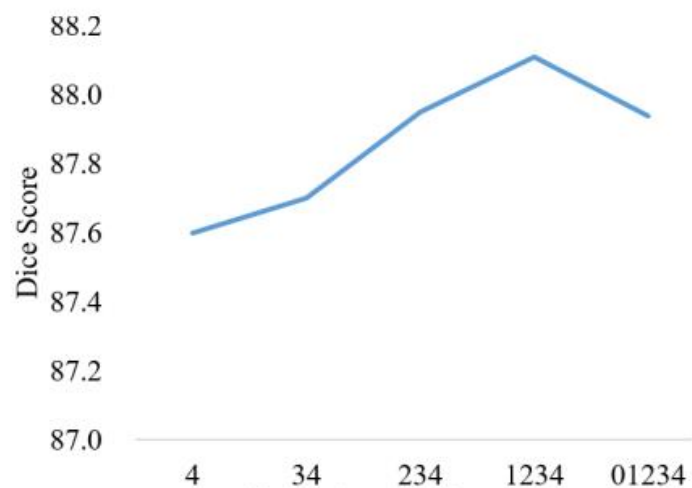
Dice系数

鲁棒豪斯多夫距离

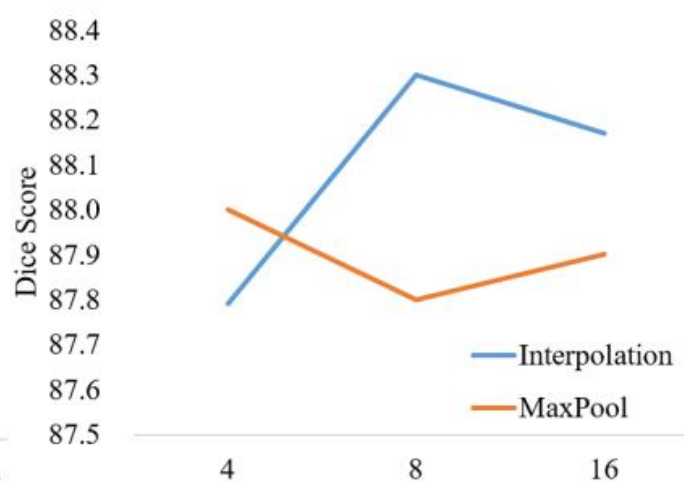


6. 实验：模型测试--消融实验

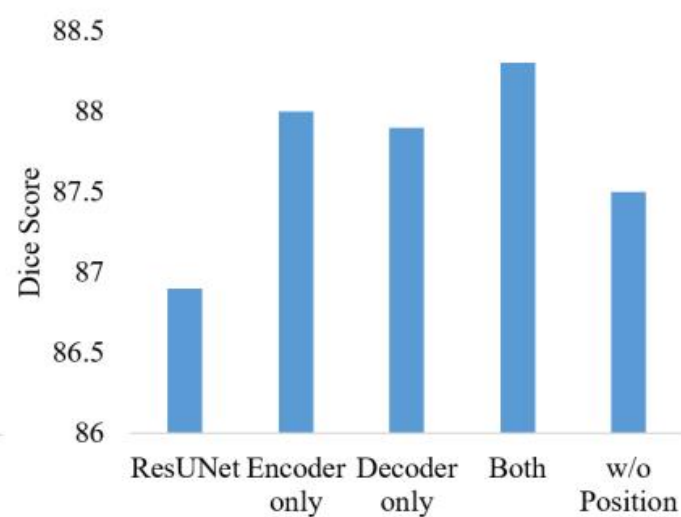
消融实验：



注意力机制放置位置



缩小尺寸



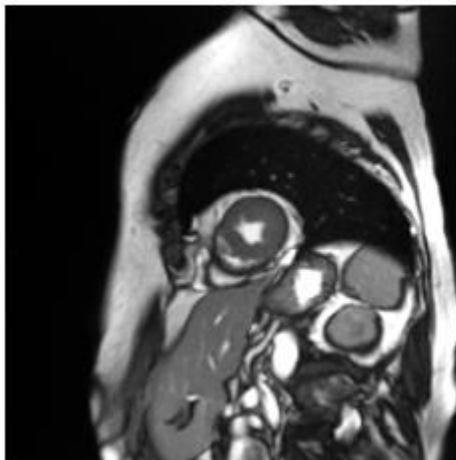
模块组合



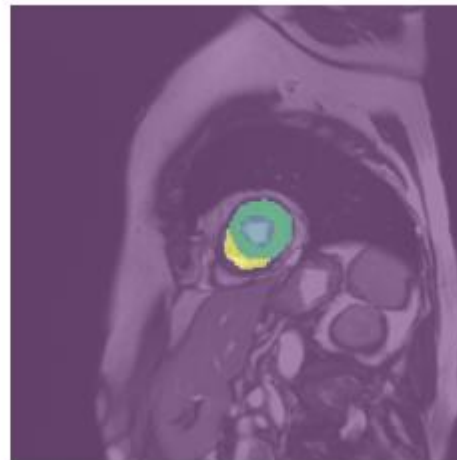
6. 实验：模型效果

供应商C

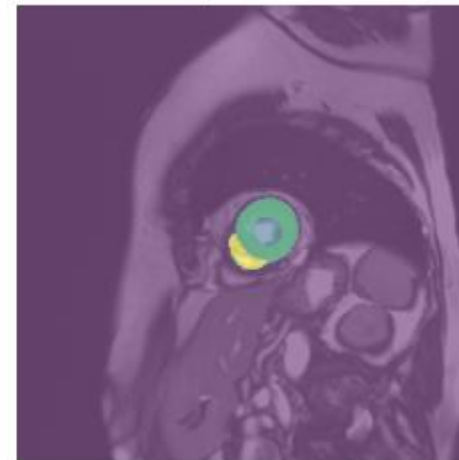
Original



Ground Truth

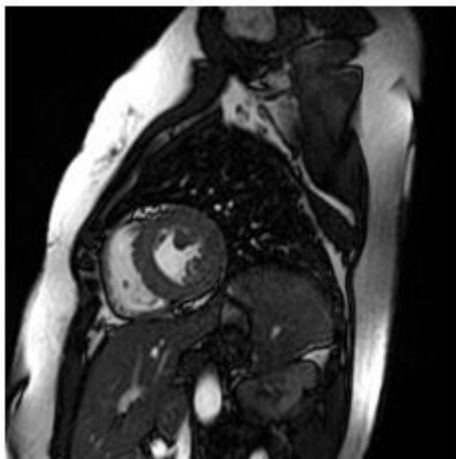


Predition(dice=0.9120)

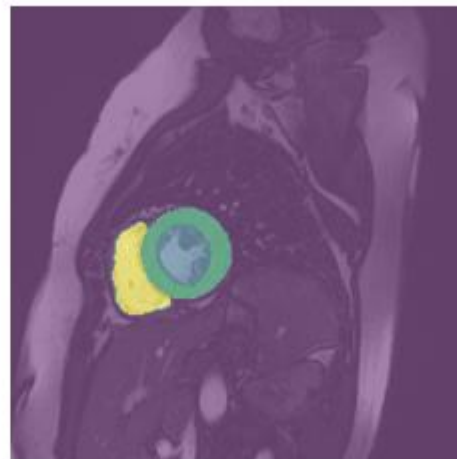


供应商D

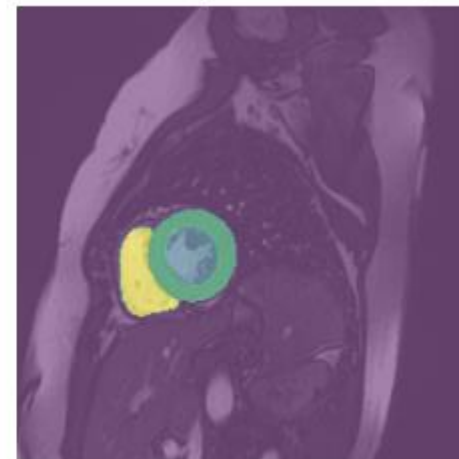
Original



Ground Truth



Predition(dice=0.9288)





总结

```

class UNet(nn.Module):
    def __init__(self, in_chan, base_chan, num_classes, reduce_size=8, block_list='14', num_blocks=1, 2, 3):
        super().__init__()
        self.in_chan = in_chan
        self.base_chan = base_chan
        self.num_classes = num_classes
        self.reduce_size = reduce_size
        self.block_list = block_list
        self.num_blocks = num_blocks

        # Encoder
        self.encoder = nn.Sequential(*[
            nn.Conv2d(in_chan, base_chan, kernel_size=3, padding=1),
            nn.BatchNorm2d(base_chan),
            nn.ReLU(inplace=True),
            nn.Conv2d(base_chan, base_chan, kernel_size=3, padding=1),
            nn.BatchNorm2d(base_chan),
            nn.ReLU(inplace=True),
            nn.MaxPool2d(kernel_size=2, stride=2),
        ])

        # Bottleneck
        self.bottleneck = nn.Sequential(*[
            nn.Conv2d(base_chan, base_chan, kernel_size=3, padding=1),
            nn.BatchNorm2d(base_chan),
            nn.ReLU(inplace=True),
            nn.Conv2d(base_chan, base_chan, kernel_size=3, padding=1),
            nn.BatchNorm2d(base_chan),
            nn.ReLU(inplace=True),
            nn.MaxPool2d(kernel_size=2, stride=2),
        ])

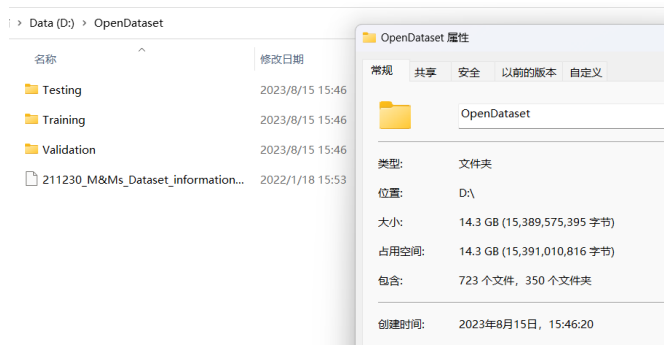
        # Decoder
        self.decoder = nn.Sequential(*[
            nn.Upsample(scale_factor=2, mode='bilinear', align_corners=True),
            nn.Conv2d(base_chan, base_chan, kernel_size=3, padding=1),
            nn.BatchNorm2d(base_chan),
            nn.ReLU(inplace=True),
            nn.Conv2d(base_chan, base_chan, kernel_size=3, padding=1),
            nn.BatchNorm2d(base_chan),
            nn.ReLU(inplace=True),
        ])

        # Output
        self.output = nn.Conv2d(base_chan, num_classes, kernel_size=1, padding=0)

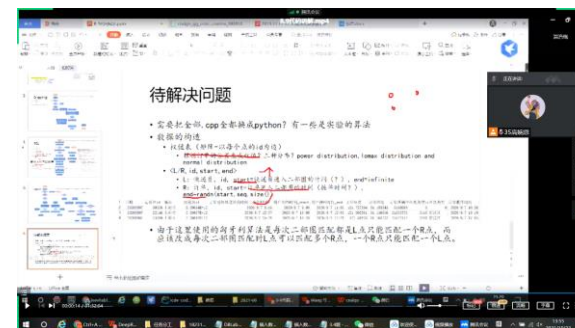
    def forward(self, x):
        x = self.encoder(x)
        x = self.bottleneck(x)
        x = self.decoder(x)
        x = self.output(x)
        return x

```

建立模型，编写
代码**1000+**



处理数据
14G+



开展会议讨
论**10**余次

阅读文献
5+篇

UTNet: A Hybrid Transformer Architecture for Medical Image Segmentation

Yunde Gao¹, Ma Zhou^{1*}, and Daisuke Matsui²

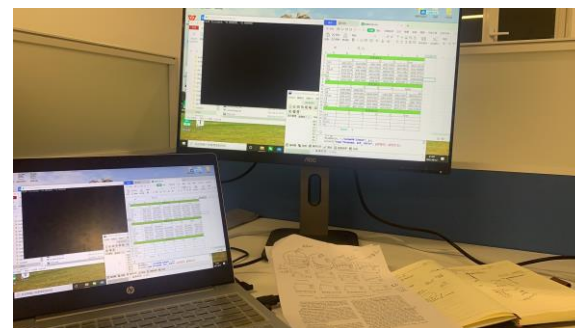
¹Department of Computer Science, Tsinghua University
²Shanghai AI Laboratory and Center for Perceptual and Interactive Intelligence

Abstract. Transformer architecture has emerged to be successful in a number of natural language processing tasks. However, its application to medical vision remains largely unexplored. In this study, we present UTNet, a single-segment hybrid Transformer architecture that integrates self-attention into a convolutional neural network for enhancing medical image segmentation. UTNet applies self-attention module in both encoder and decoder for capturing long-range dependency at different scales with minimal overhead. To this end, we propose an efficient self-attention mechanism along with relative position encoding that reduces the complexity of self-attention operation significantly from $O(n^2)$ to approximate $O(n)$. A new self-attention decoder is also proposed to restore fine-grained details from the dilated connections in the encoder. Our approach addresses the dilemma that Transformer requires large amounts of data to learn vision induction bias. Our hybrid type design allows the introduction of Transformer into convolutional networks without a need of pre-training. We have evaluated UTNet on the multi-label, multi-modal cardiac magnetic resonance imaging cohort. UTNet demonstrates superior segmentation performance and robustness against the state-of-the-art approaches, holding the promise to generalize well on other medical image segmentation tasks. Code is available.¹

1 Introduction

Convolutional networks have revolutionized the computer vision field with outstanding feature representation capability. Currently, the convolutional encoder-decoder architectures have made substantial progress in position-sensitive tasks, like semantic segmentation [14, 13, 20, 17, 6]. The used convolutional operations capture texture features by gathering local information from neighborhood pixels. To aggregate the local filter responses globally, these models stack multiple convolutional layers and connect the receptive field through dense connections. The

20h+ 实验,
修改代码





请老师批评指正!