

RNN 条件生成与Attention

七月在线 张雨石

2018年9月13日

<http://blog.csdn.net/stdcoutzyx>

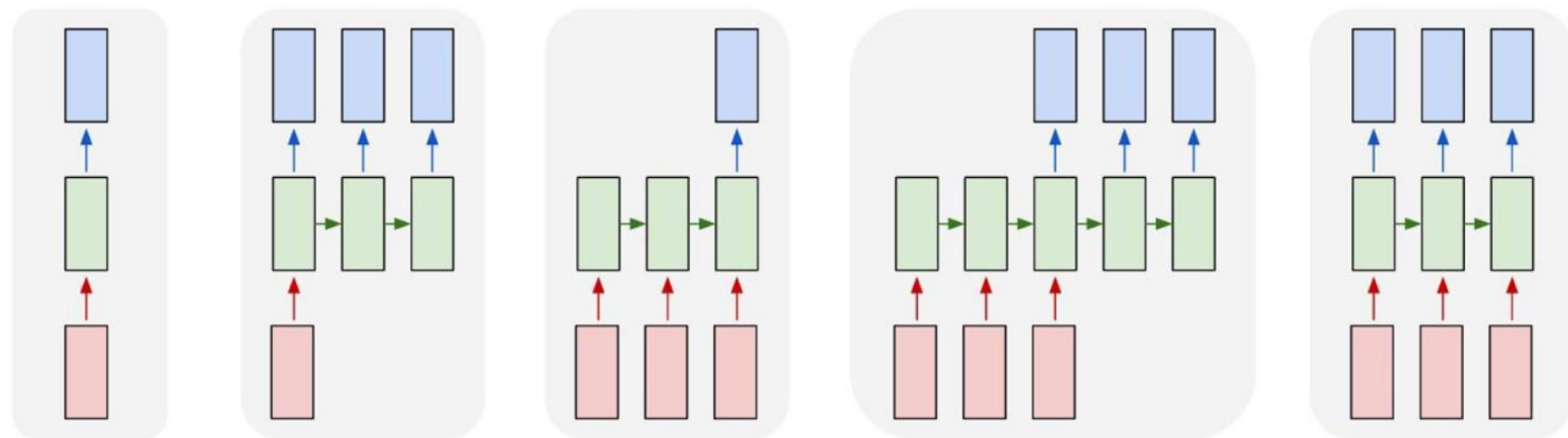
RNN条件生成与Attention

- ☐ RNN条件生成
- ☐ 机器翻译
- ☐ Attention
- ☐ 图像生成文本



RNN条件生成

□ RNN解决的问题



RNN条件生成

□ RNN解决的问题

- 图片生成描述
- 文本分类
- 机器翻译
- 视频解说

□ 条件生成问题 $P(y|x)$

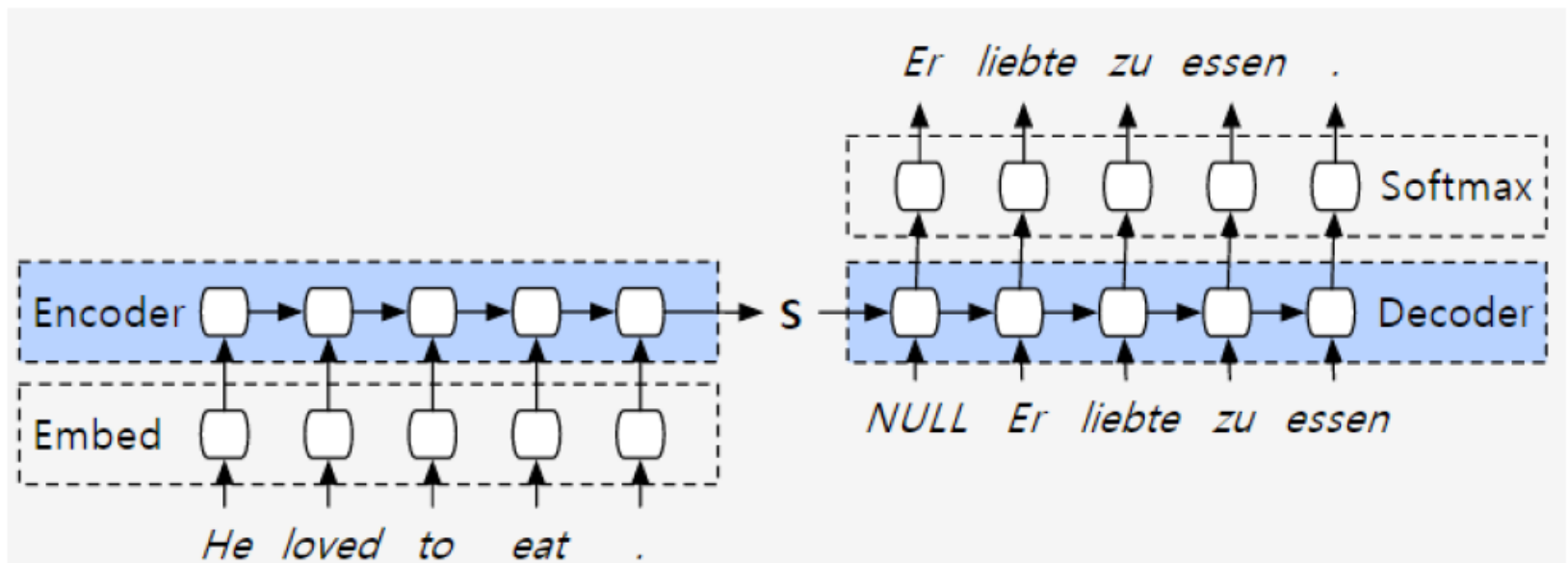


机器翻译

- ☐ V1: Encoder-Decoder
- ☐ V2: Attention-based Encoder-Decoder
- ☐ V3: Bi-directional encode layer
- ☐ V4: Residual-based Encoder-Decoder
- ☐ 总结



V1: Encoder-Decoder



V1: Encoder-Decoder

- 编码过程

- 解码过程

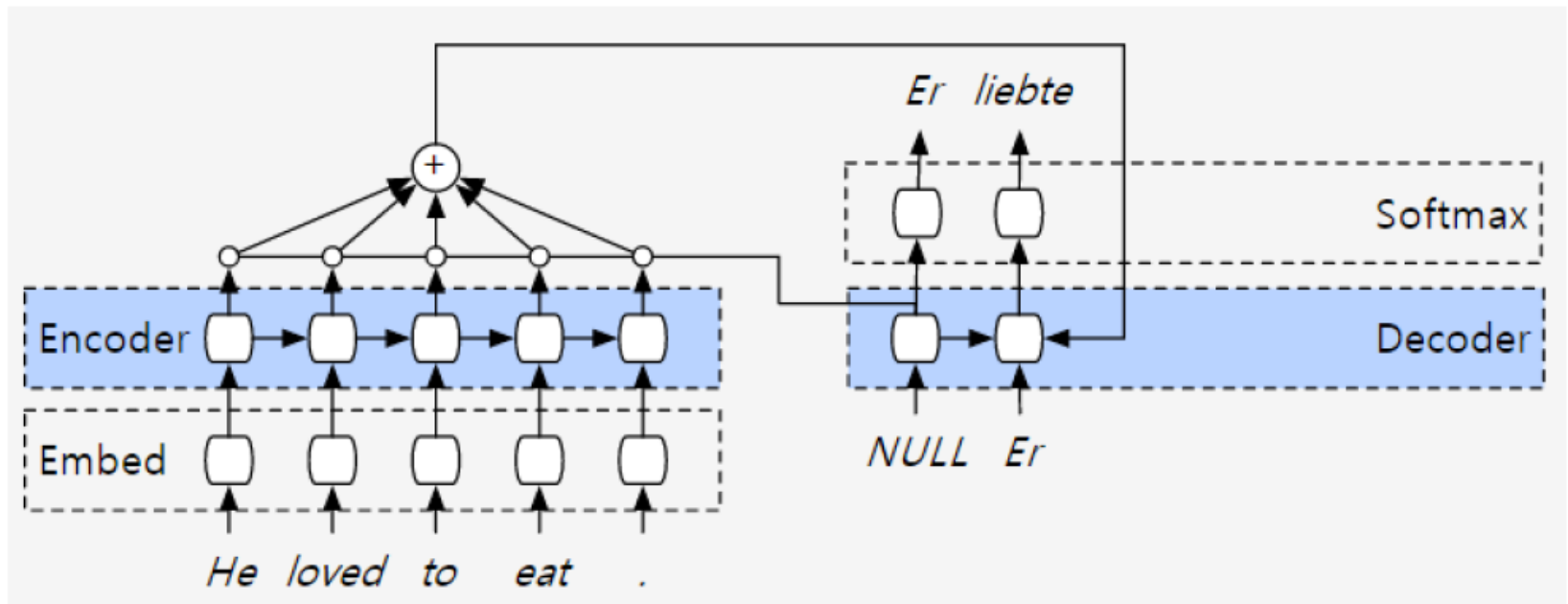
- 缺点

 - 定长编码是信息瓶颈

 - 长度越长，前面输入进RNN的信息就越被稀释



V2: Attention-based Encoder-Decoder



V2: Attention-based Encoder-Decoder

□ Attention层

■ Attention计算

□ LSTM各步的输出向量 $a=[a_1, a_2, \dots, a_n]$

□ Decoder每步的中间状态 h_i

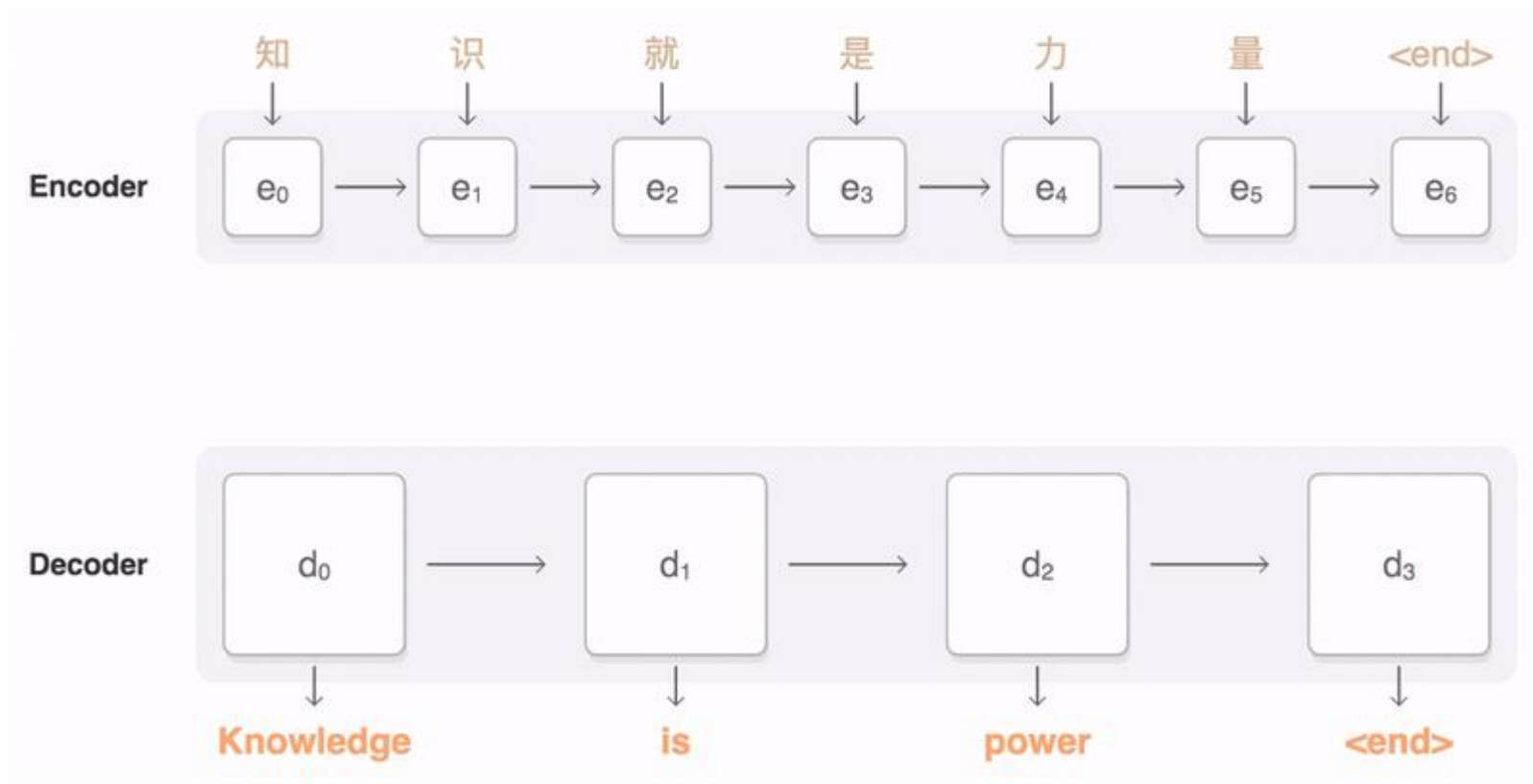
□ $\text{Alpha} = [\tanh(w_1 * a_j + w_2 * h_i) \text{ for } j \text{ in range}(n)]$

■ Attention加权

□ $\text{Alpha} * a$



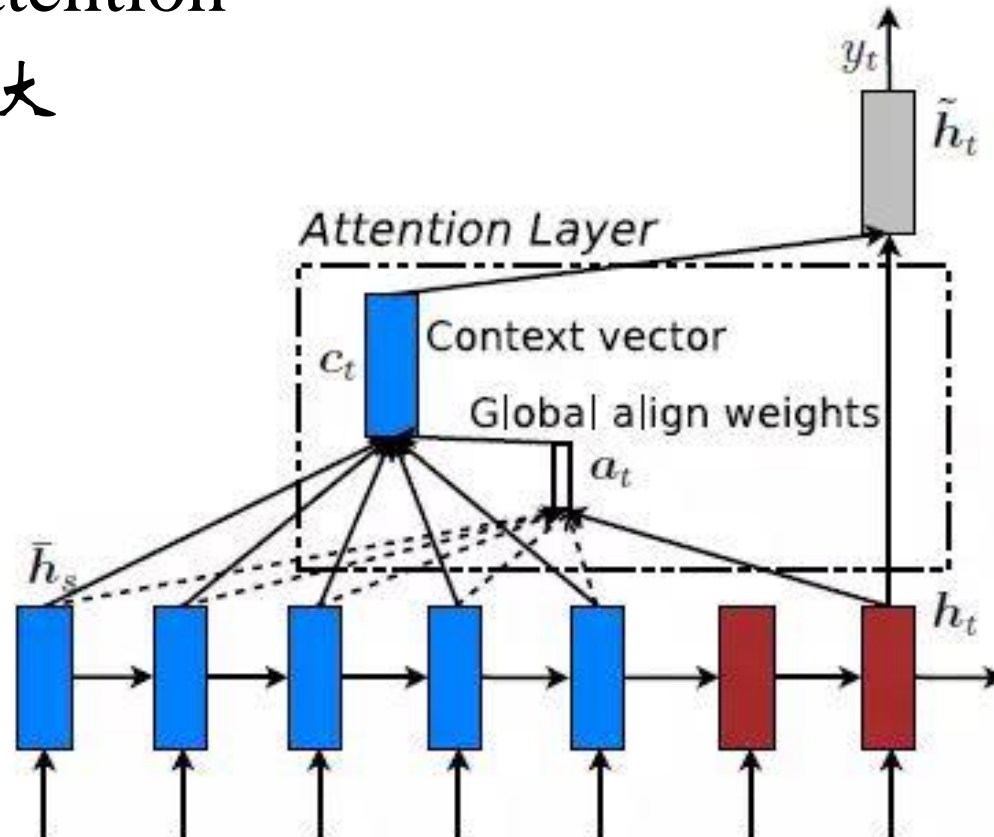
V2: Attention-based Encoder-Decoder



V2: Attention-based Encoder-Decoder

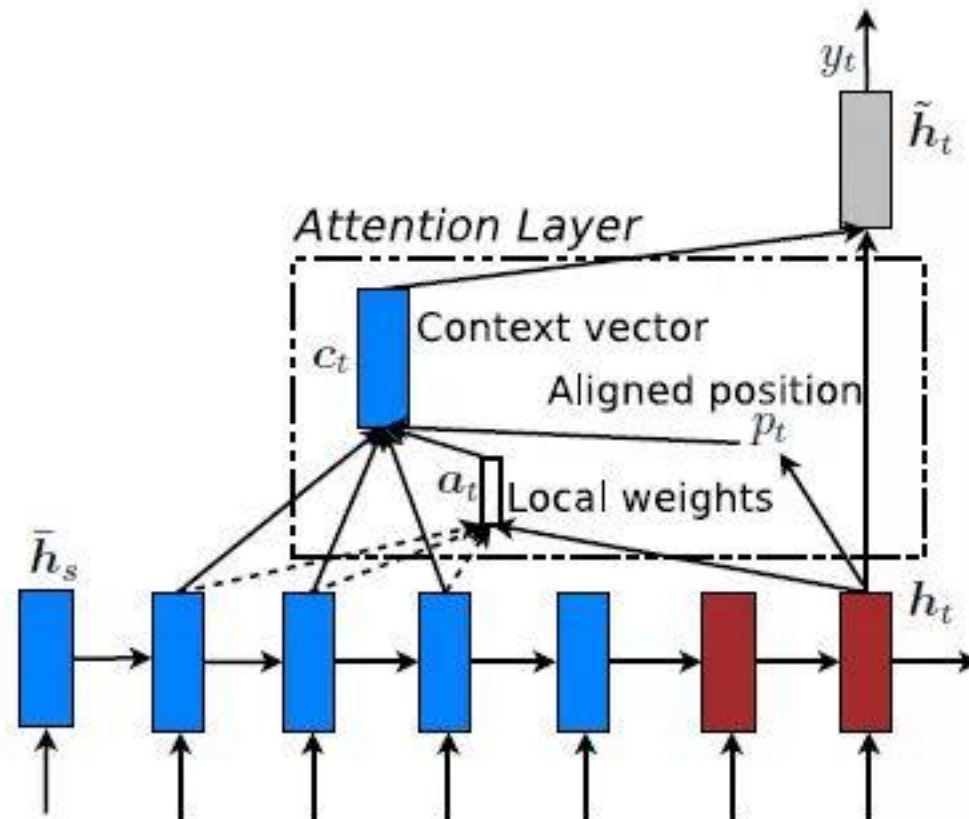
□ Global Attention

■ 计算量大

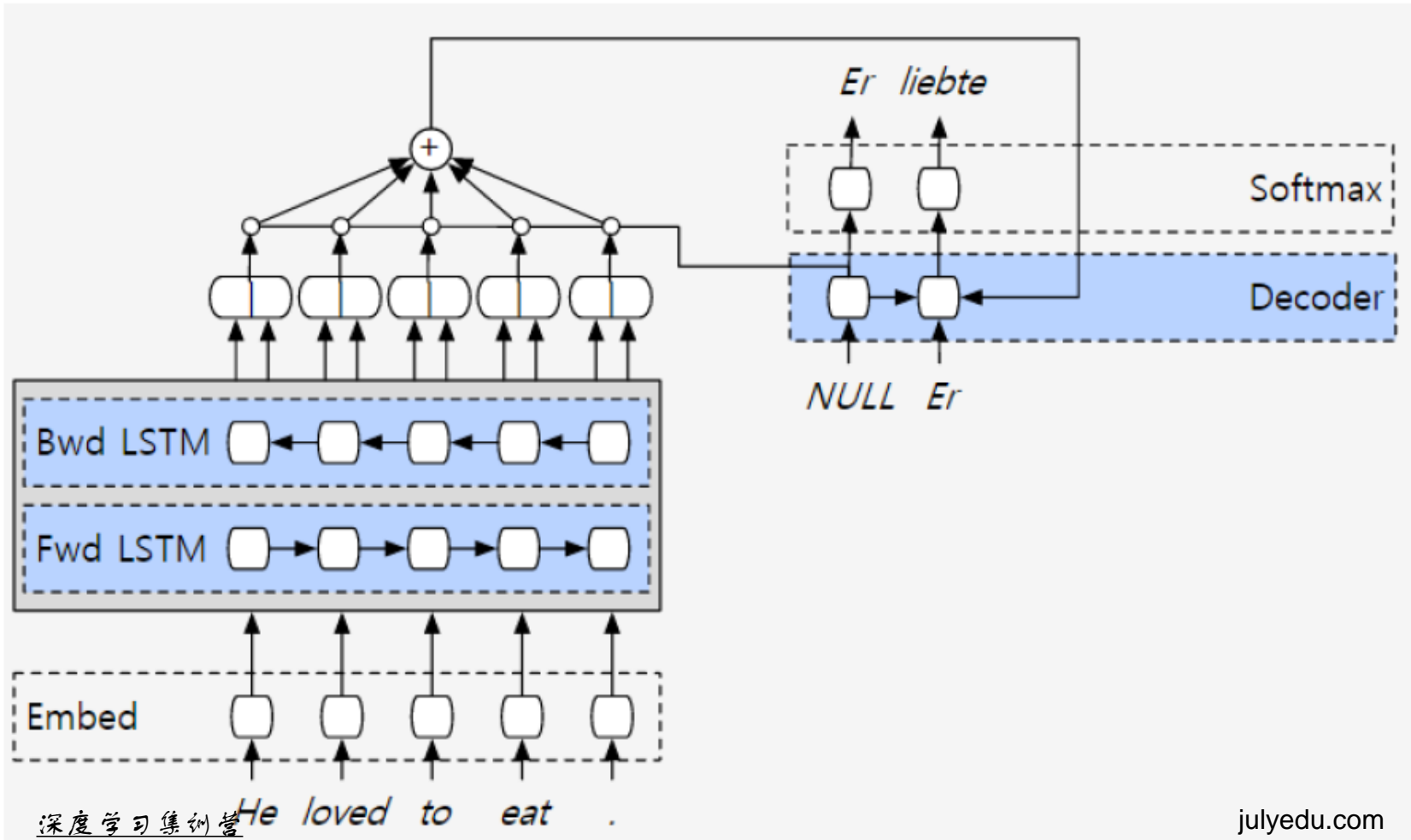


V2: Attention-based Encoder-Decoder

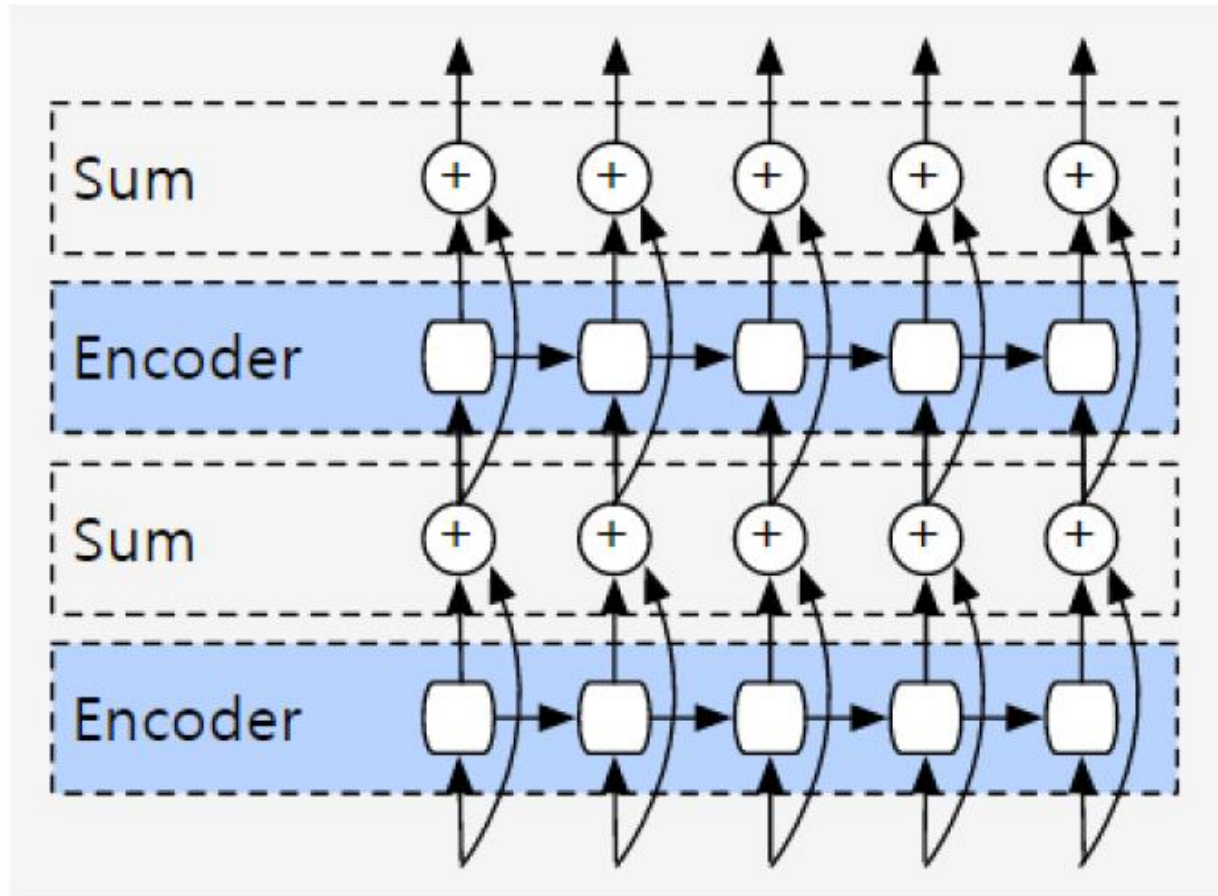
□ Local Attention



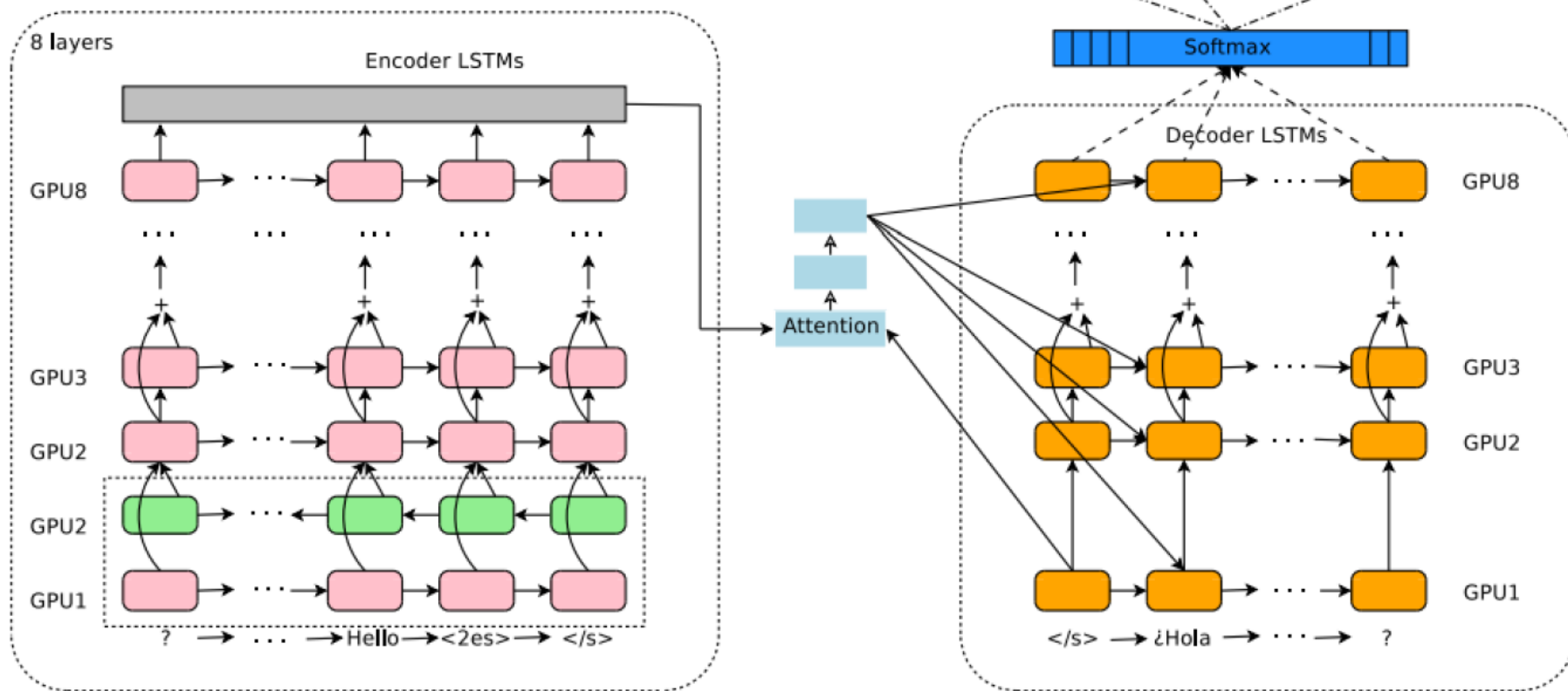
V3: Bi-Directional Encode layer



V4: Residual Encode layer



机器翻译-总结



Attention

- ☐ Global Attention
- ☐ Local Attention
- ☐ Self Attention
- ☐ Hierarchical Attention

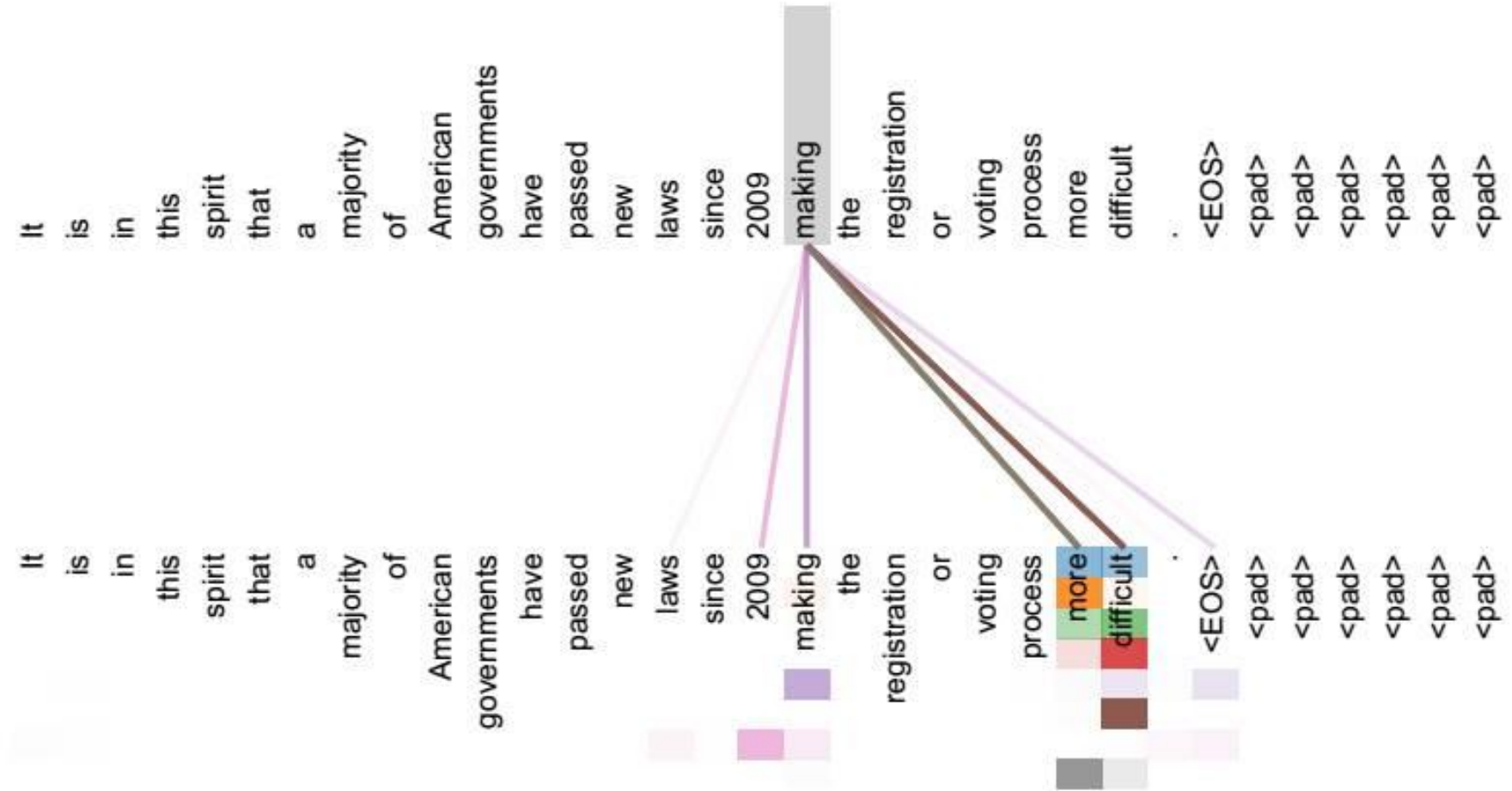


Self Attention

- ❑ 传统的attention只有source和target之间的关联关系
- ❑ 忽略了source和target端分别的关联关系



Self Attention



Self Attention

- 可以理解为source=target的encoder-decoder
- 捕捉到词与词之间的关系
- 增加并行性



Hierarchical Attention

☐ 文本分类

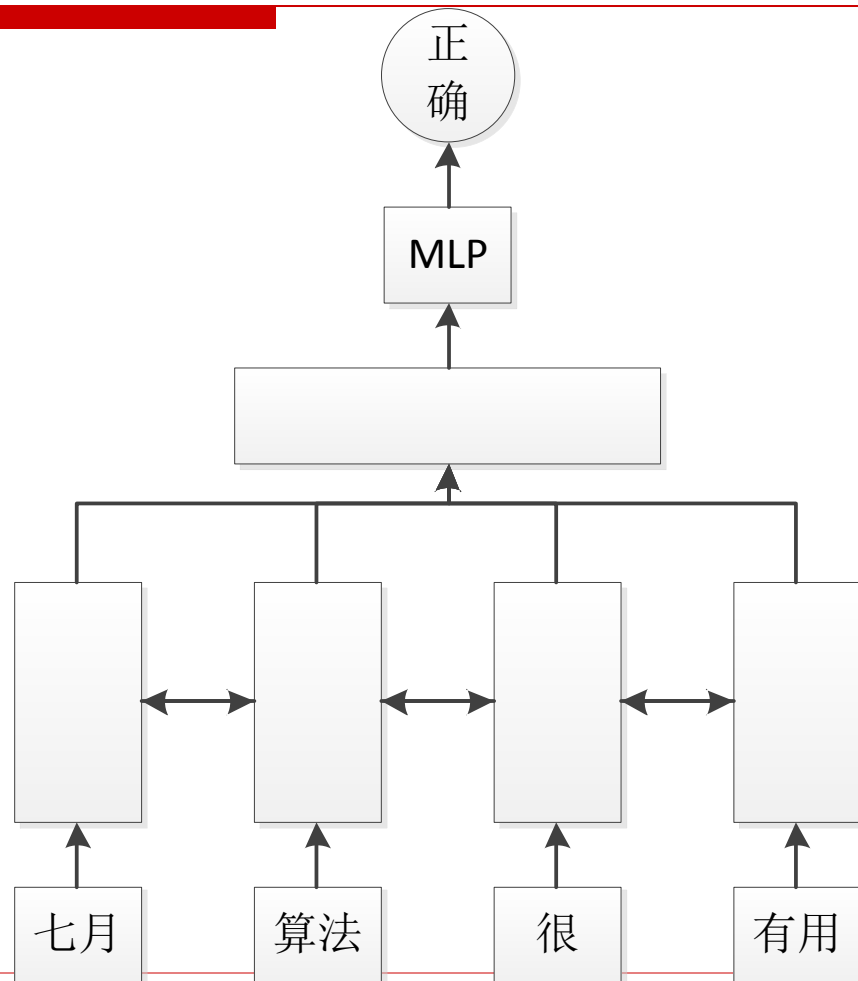
■ 双向LSTM

■ 输出合并

☐ 拼接

☐ Average

☐ max



Hierarchical

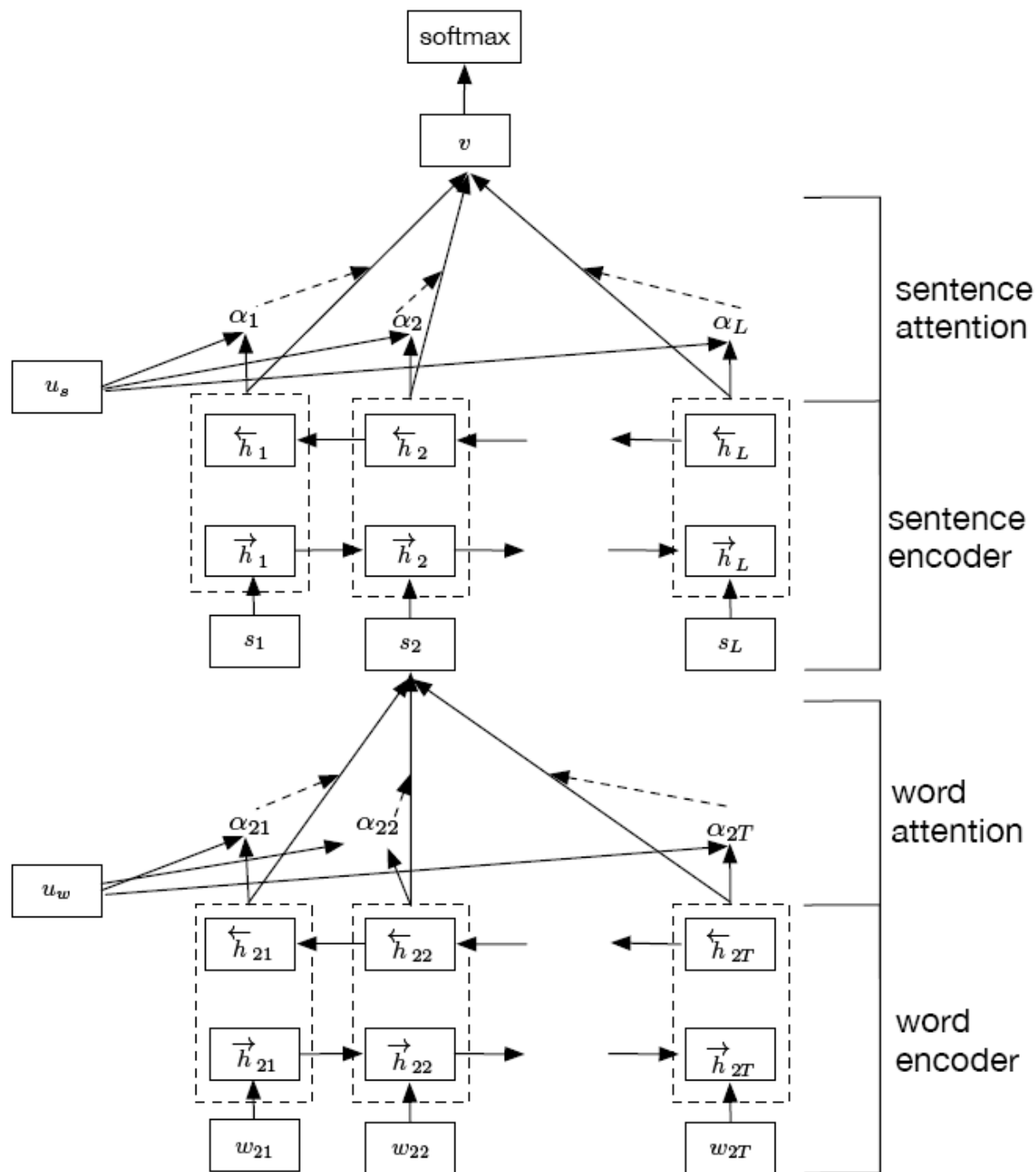
□ 文本分类

■ 分层

□ 词到句

□ 句到段落

■ Attention



图像生成文本

□ 问题引入

□ 模型变迁

- M-RNN

- Neural Image Caption

- Attention-based

问题引入

□ Deep Learning 出现之前

- 检索问题
- 不能泛化到新的数据

问题引入

☐ 图像搜索

- 找到图像对应的描述，丰富搜索元数据

☐ 盲人导航

☐ 少儿教育

- 看图说话

问题引入

□ 数据集

- IAPR TC-12: 20000张图像，平均每张图像有1.7句描述
- Flickr8k: 8000张图像，语法比IAPR简单
- Flickr30k: 31783张图像，158915句描述，平均每个五句。
- MS COCO: 16W张图像
- AI Challenger: 图像中文描述数据

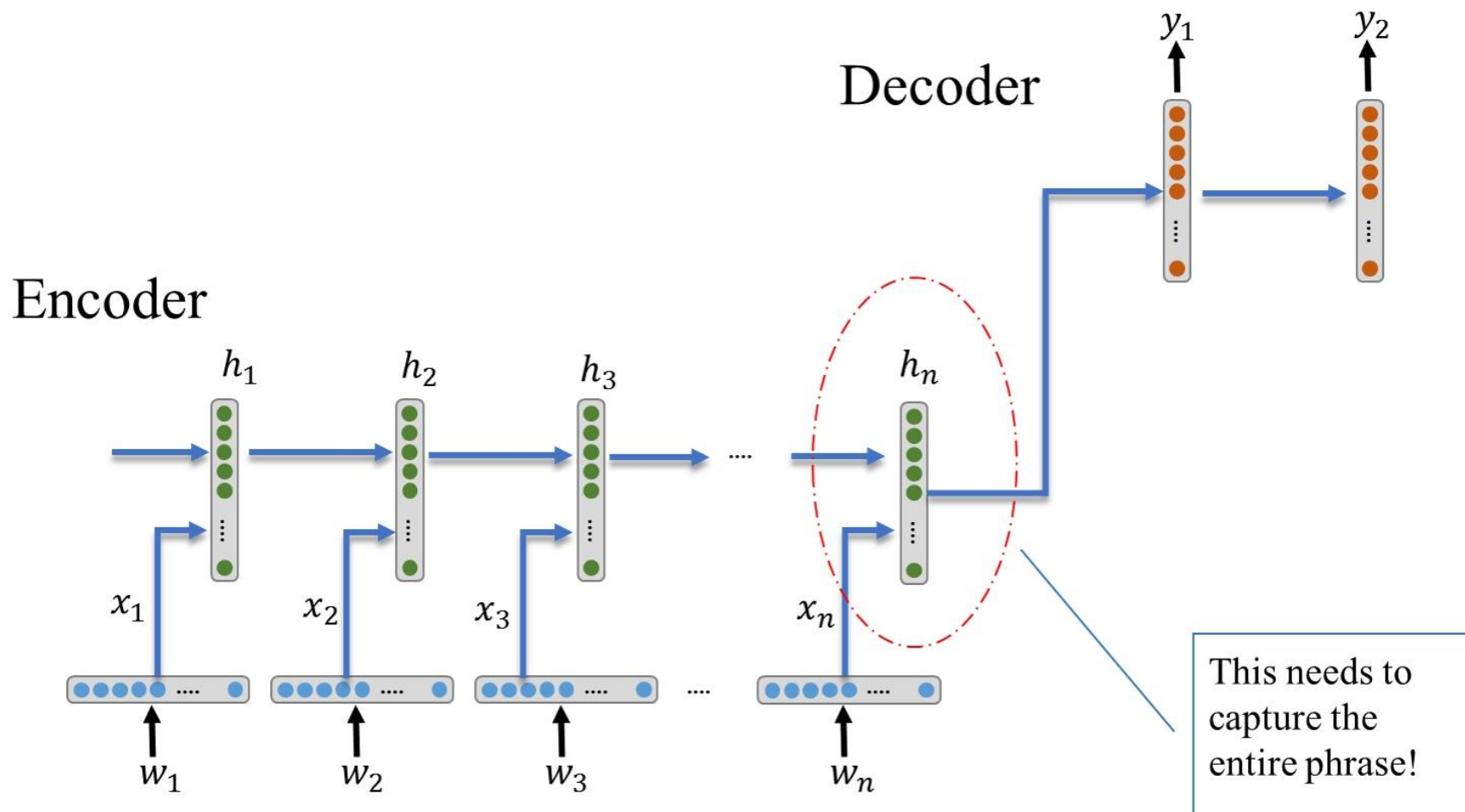
问题引入

□ 评测

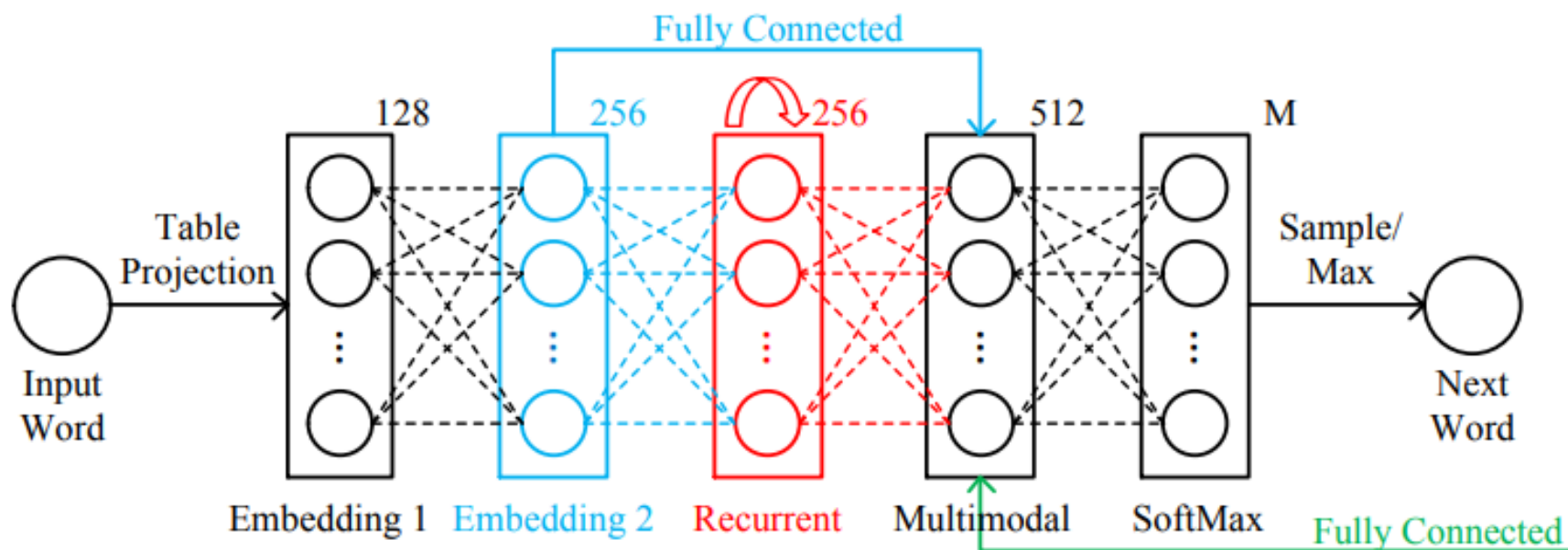
- BLEU score

- 句子检索和图像检索

模型变迁—Encoder-Decoder



模型变迁——Multi-model RNN



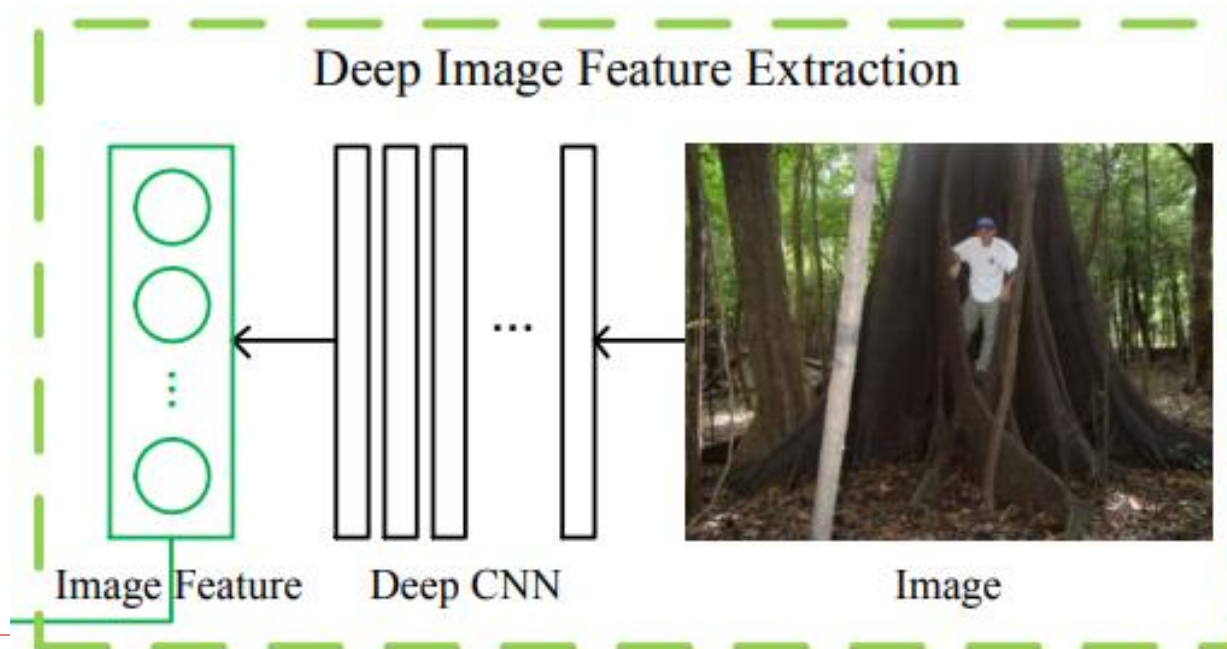
模型变迁——Multi-Model RNN

- 输入词生成embedding
- Embedding输入到RNN和Multi-Model
- RNN生成更加抽象的embedding
- 原始Embedding、RNN的embedding和图像Feature同时输入给分类层

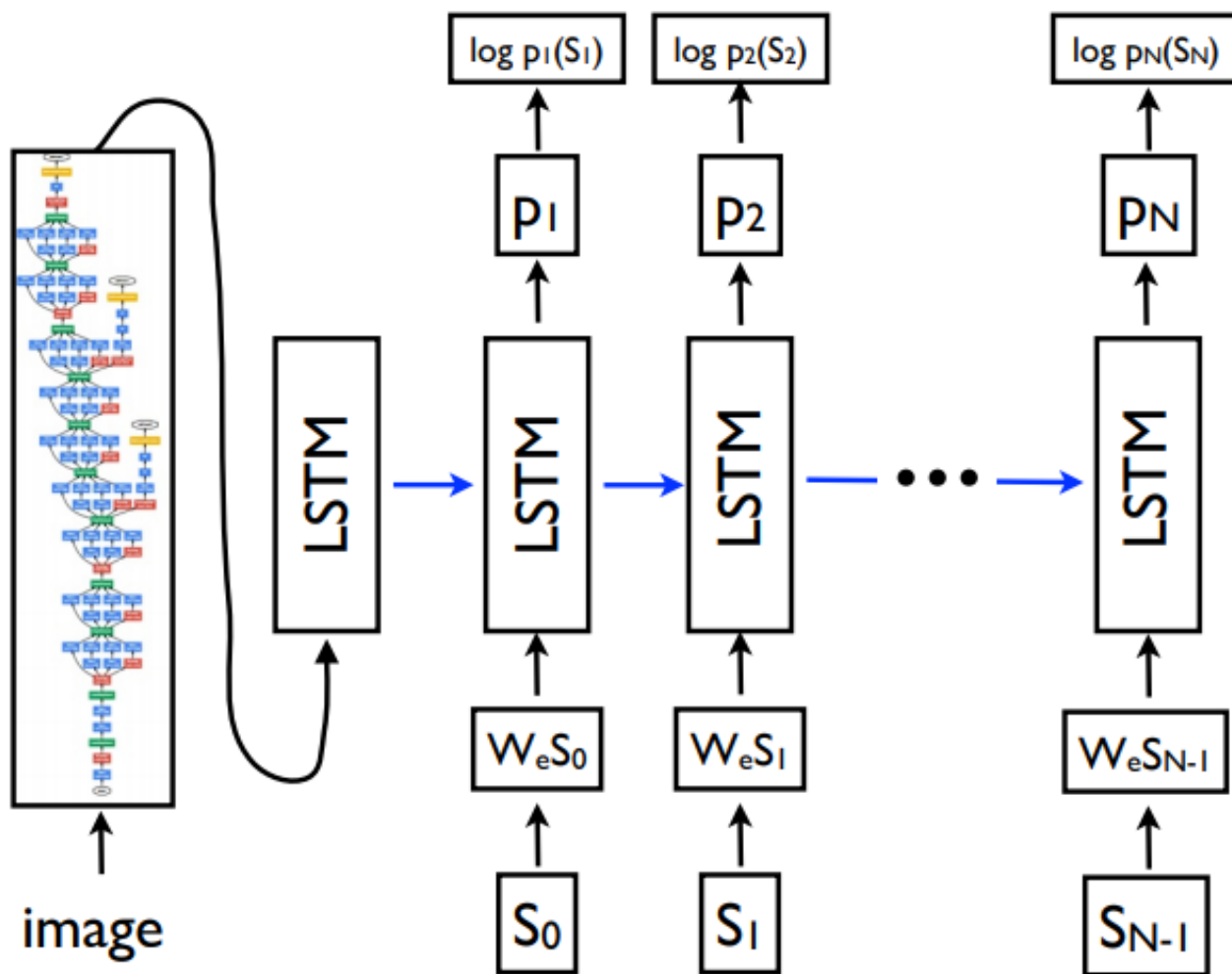


模型变迁——Multi-model RNN

- AlexNet 7th feature
- Object detection feature



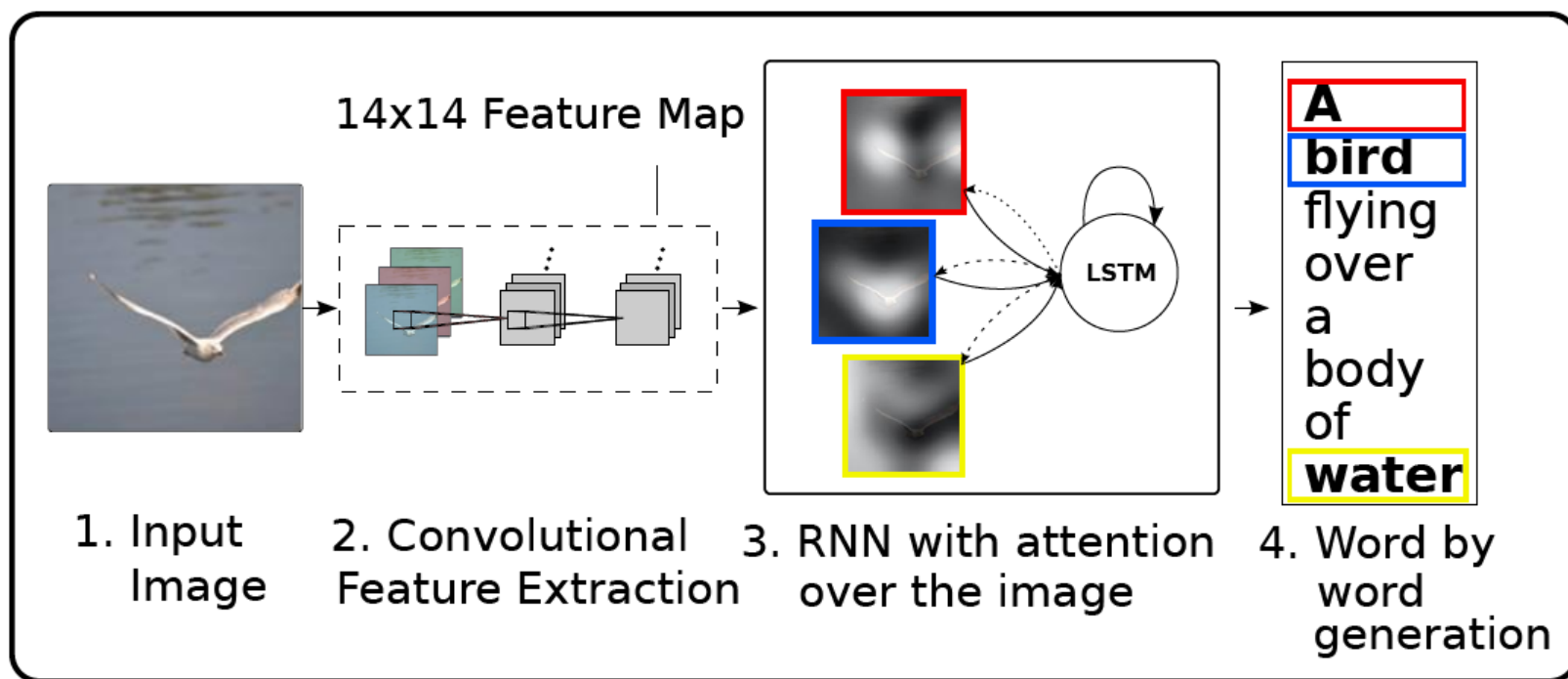
模型变迁——Show And Tell



模型变迁——Show And Tell

- 图像特征使用更强大的CNN提取
 - GoogNet、Residual等
- 图像特征只提取一次
- LSTM生成文本

模型变迁——Show Attend and Tell



模型变迁——Show Attend and Tell

□ 不使用全连接层而是某个卷积层的feature

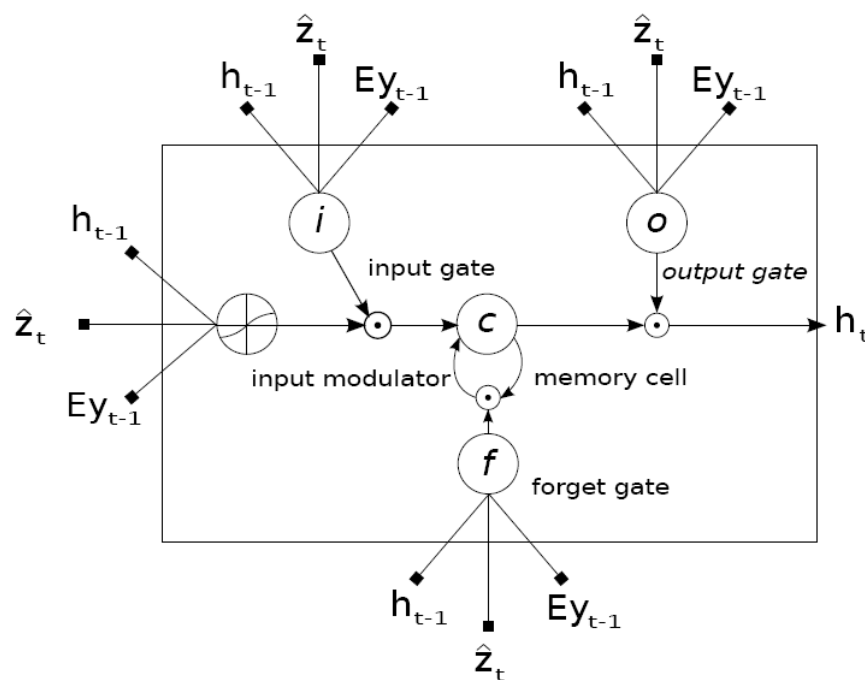
■ 不使用全连接层

■ 使用卷积层的输出

□ 有位置信息

■ LSTM输入是加权和

□ attention



模型变迁——Show Attend and Tell

□ 不使用全连接层而是某个卷积层的feature

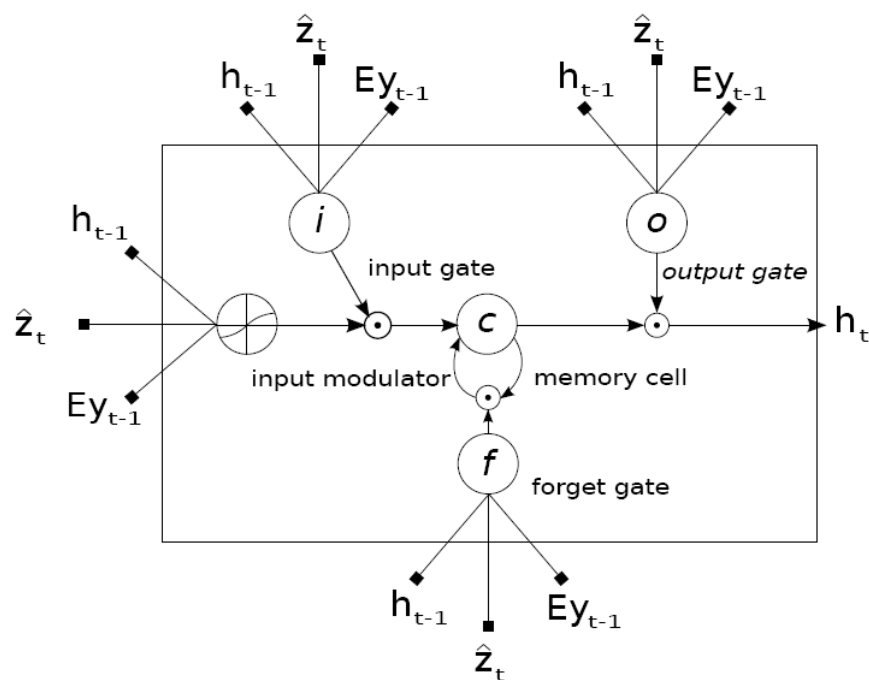
■ $14*14*256$: 有 $14*14$ 个位置可以选择

$$a = \{\mathbf{a}_1, \dots, \mathbf{a}_L\}, \mathbf{a}_i \in \mathbb{R}^D$$

$$e_{ti} = f_{\text{att}}(\mathbf{a}_i, \mathbf{h}_{t-1})$$

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{k=1}^L \exp(e_{tk})}$$

$$\hat{\mathbf{z}}_t = \phi(\{\mathbf{a}_i\}, \{\alpha_i\})$$



模型变迁——Show Attend and Tell



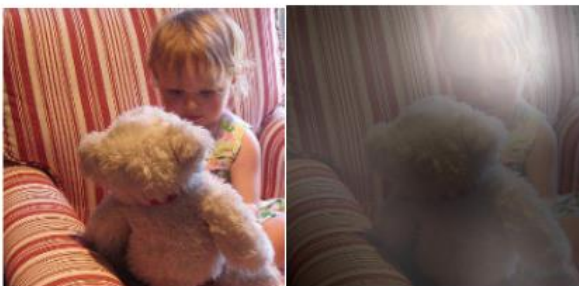
A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



A little girl sitting on a bed with a teddy bear.

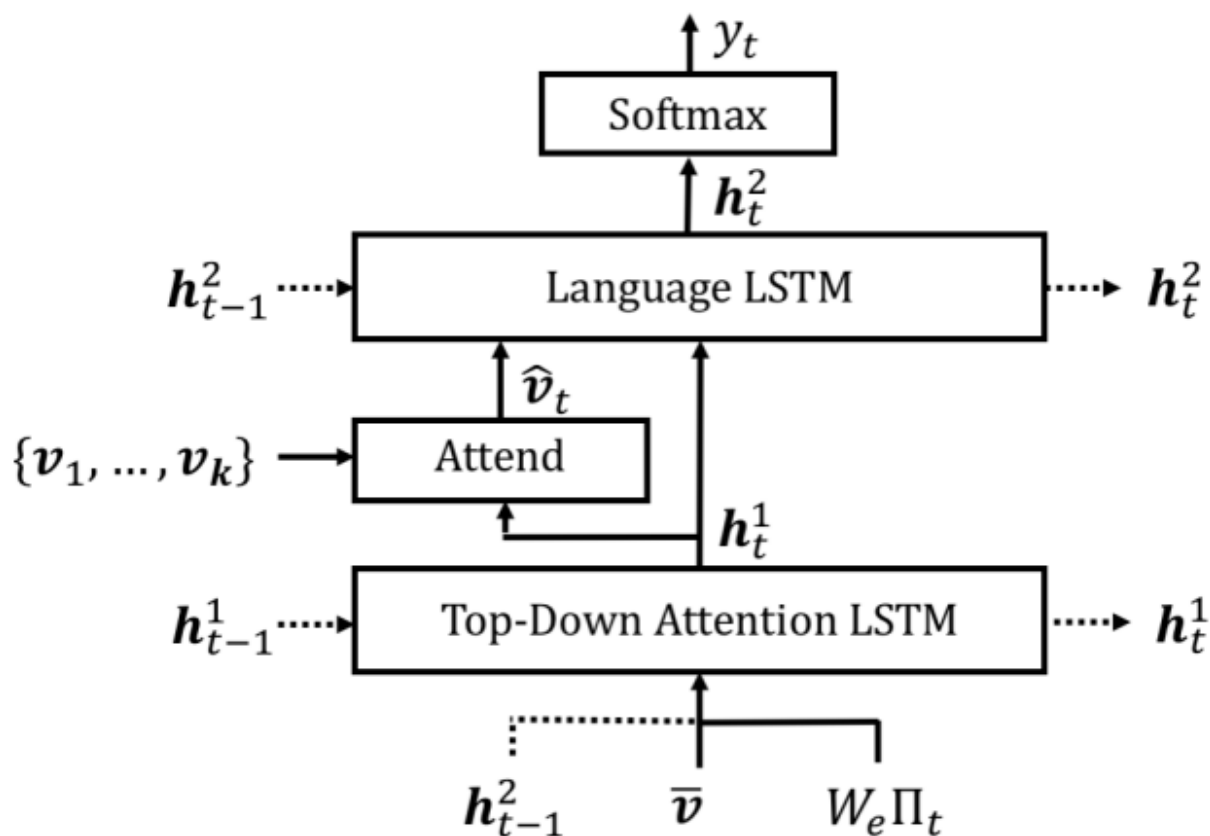


A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.

模型变迁——Top Down Attention

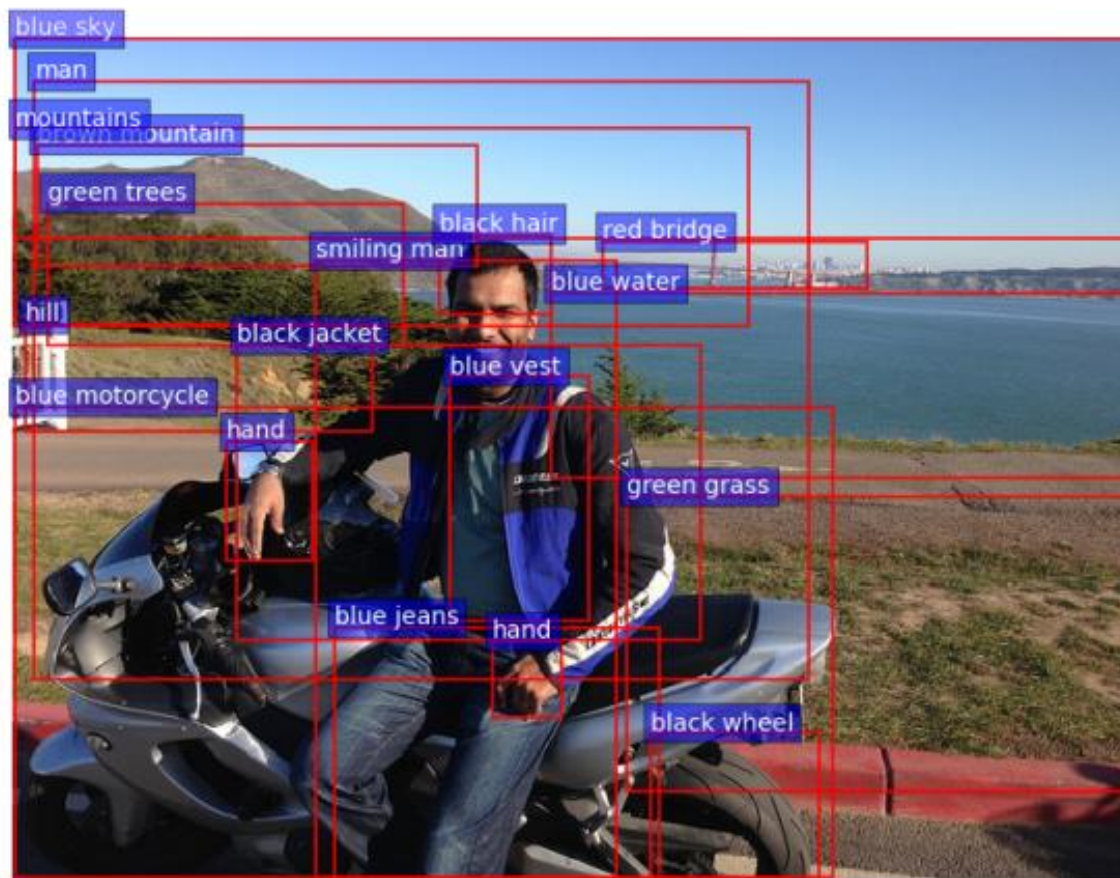


模型变迁——Top Down Attention

- 第一层LSTM负责学习 attention
 - 前一时刻Language Model的hidden state
 - 图像的全局信息
 - 所有feature map的平均
 - 输入词语的embedding

$$x_t^1 = [h_{t-1}^2, \bar{v}, W_e \Pi_t]$$

模型变迁——Top Down Attention



模型变迁——Top Down Attention

□ 第一层LSTM输出和图像feature计算attention

$$a_{i,t} = \mathbf{w}_a^T \tanh(W_{va}\mathbf{v}_i + W_{ha}\mathbf{h}_t^1)$$

$$\alpha_t = \text{softmax}(\mathbf{a}_t)$$

$$\hat{\mathbf{v}}_t = \sum_{i=1}^K \alpha_{i,t} \mathbf{v}_i$$

模型变迁——Top Down Attention

□ 第二层LSTM

- 第一层LSTM的隐含层
- 加权平均的图像feature

$$\mathbf{x}_t^2 = [\hat{\mathbf{v}}_t, \mathbf{h}_t^1]$$

$$p(y_t \mid y_{1:t-1}) = \text{softmax}(W_p \mathbf{h}_t^2 + \mathbf{b}_p)$$

总结

- RNN 条件生成问题

- 机器翻译

 - Encoder-Decoder、attention、双向、残差

- Attention

 - Self-attention、Hierarchical attention

- 图像生成文本

 - Multi-Modal、Show and Tell

 - Show attend and tell、Top-down bottom-up attention

Thanks!

Q&A