

# 语音合成：从入门到精通

## 第四讲：基于序列到序列的声学模型

主讲人 阳珊

西北工业大学语音与语言处理实验室  
syang@nwpu-aslp.org





## 1. 概述



## 2. 序列到序列模型及注意力机制



## 3. 序列到序列声学模型Tacotron



## 4. 序列到序列声学模型变体（一）



## 5. 实战



## 1. 概述



## 2. 序列到序列模型及注意力机制



## 3. 序列到序列声学模型Tacotron



## 4. 序列到序列声学模型变体（一）



## 5. 实战



## 概述-文本分析回顾

### 文本分析过程

- 分词、注音、词性标注、韵律预测等
- 将标准化的抄本表示为音素级别的文本特征
- 文本特征的长度 $M$ 仅与抄本本身有关，与语音时长无关



# 概述-文本分析回顾

## 文本分析过程

示例：欢迎大家参加深蓝学院语音合成课程

```
1 3500000 XX^X-;IL+ka=X X/A:X X+X/B:X-X-X-X-X-X-X/X/C:4+X+2 X/D:X-X/E:X X@X+X6X+X#X+X/F:2=1/G:X X/H:X-X^X=X|X/I:4=3/J:16+10-4$
2 3500000 4431312 X^SIL-k+ao=z@1_2/A:X X_X+X/B:4-X=2-v@1-161-4#6-0|ao/C:3+uo+2_f/D:X-X/E:2_1@1+36X+X#X+X/F:2=1/G:X X/H:4=3^1=4|X/I:6=4/J:16+10-4$
3 4431312 5680483 IL^k-ao+z=uo@2_1/A:X X_X+X/B:4-X=2-v@1-161-4#0-2|ao/C:3+uo+2_f/D:X-X/E:2_1@1+36X+X#X+X/F:2=1/G:X X/H:4=3^1=4|X/I:6=4/J:16+10-4$
4 5680483 6121631 ^k-ao+z=uo@1_2/A:ao_2+v/B:3-X=2-f@1-162-3#2-0|uo/C:4+ia+2_v/D:2-1/E:2_1@2+26X+X#X+X/F:4=2/G:X X/H:4=3^1=4|X/I:6=4/J:16+10-4$
5 6121631 7386675 ao^z-uo+x=ia@2_1/A:ao_2+v/B:3-X=2-f@1-162-3#0-2|uo/C:4+ia+2_v/D:2-1/E:2_1@2+26X+X#X+X/F:4=2/G:X X/H:4=3^1=4|X/I:6=4/J:16+10-4$
6 7386675 8153613 ^uo-x+ia=p@1_2/A:3_uo_2+f/B:4-X=2-v@1-263-2#2-0|ia/C:1+o+2_v/D:2-1/E:4_2@3+16X+X#X+X/F:4=2/G:X X/H:4=3^1=4|X/I:6=4/J:16+10-4$
7 8153613 9491280 wo^x+ia=p@2_1/A:3_uo_2+f/B:4-X=2-v@1-263-2#0-1|ia/C:1+o+2_v/D:2-1/E:4_2@3+16X+X#X+X/F:4=2/G:X X/H:4=3^1=4|X/I:6=4/J:16+10-4$
8 9491280 10457192 x^ia-p+o=lp@1_2/A:ia_2+v/B:1-X=2-v@2-164-1#1-0|o/C:X+X^X X/D:2-1/E:4_2@3+16X+X#X+X/F:4=2/G:X X/H:4=3^1=4|X/I:6=4/J:16+10-4$
9 10457192 12363202 ia^p-o+lp=q@2_1/A:ia_2+v/B:1-X=2-v@2-164-1#0-5|o/C:X+X^X X/D:2-1/E:4_2@3+16X+X#X+X/F:4=2/G:X X/H:4=3^1=4|X/I:6=4/J:16+10-4$
10 12363202 14314374 p^o-lp=q=ian@X X/A:1_o_2+v/B:X-X-X-X-X-X-X/X/C:2+ian+2_f/D:4-2/E:X X@X+X6X+X#X+X/F:4=2/G:4_3/H:X-X^X=X|X/I:6=4/J:16+10-4$
11 14314374 15470580 o^lp-q+ian=f@1_2/A:X X_X+X/B:2-X=2-f@1-261-6#5-0|ian/C:1+ang+2_f/D:4-2/E:4_2@1+46X+X#X+X/F:2=1/G:4_3/H:6=4^2=3|X/I:2=1/J:16+10-4$
```

```
1 X^X-sil+h_uan'0#0/A:0(0;0/B:0+0;0+0/C:2)3;11)29/E:0-0-0;0/F:0]0]0;0]0]0=0~X|0/G:2#6#17;0/H:X<v;0<0|0<0=0<0/I:0-0/J:1
2 X^sil-h+uan_ing'1#2/A:0(0;0/B:2+3;11+29/C:1)4;7)0/E:0-0-0;0/F:2]6]17;1]2]1]6=1|17~uan|0/G:2#4#0;1/H:X<v<r;0<2|0<5=0<0/I:0-0/J:1
3 sil^h-uan+ing_d'2#1/A:0(0;0/B:2+3;11+29/C:1)4;7)0/E:0-0-0;0/F:2]6]17;1]2]1]6=1|17~uan|0/G:2#4#0;1/H:X<v<r;0<2|0<5=0<0/I:0-0/J:1
4 h_uan-ing+d_a'1#1/A:2(0;0/B:1+3;11+29/C:2)4;7)0/E:0-0-0;0/F:2]6]17;2]1]2]5=2|16~ing|1/G:2#4#0;0/H:X<v<r;0<2|0<5=0<0/I:0-0/J:1
5 uan-ing+d+a_j'1#2/A:1(3;0/B:2+4;11+29/C:2)4;7)0/E:2-0-0;1/F:2]6]17;1]2]3]4=3|15~a|0/G:2#4#0;1/H:v<r<v;1<1|1<4=0<0/I:0-0/J:1
6 ing^d-a+j_ia'2#1/A:1(3;0/B:2+4;11+29/C:2)4;7)0/E:2-0-0;1/F:2]6]17;1]2]3]4=3|15~a|0/G:2#4#0;1/H:v<r<v;1<1|1<4=0<0/I:0-0/J:1
7 d^a-j+ia_c'1#2/A:2(3;0/B:2+4;11+29/C:2)4;7)0/E:2-0-0;0/F:2]6]17;2]1]4]3=4|14~ia|1/G:2#4#0;0/H:v<r<v;1<1|1<4=0<0/I:0-0/J:1
8 a_j-ia+c_an'2#1/A:2(3;0/B:2+4;11+29/C:2)4;7)0/E:2-0-0;0/F:2]6]17;2]1]4]3=4|14~ia|1/G:2#4#0;0/H:v<r<v;1<1|1<4=0<0/I:0-0/J:1
9 j^ia-c+an_j'1#2/A:2(4;0/B:2+4;11+29/C:2)7;7)0/E:2-0-0;1/F:2]6]17;1]2]5]2=5|13~an|0/G:4#4#0;2/H:r<v<ni;2<0|2<3=0<0/I:0-0/J:1
10 ia^c-an+j_ia'2#1/A:2(4;0/B:2+4;11+29/C:2)7;7)0/E:2-0-0;1/F:2]6]17;1]2]5]2=5|13~an|0/G:4#4#0;2/H:r<v<ni;2<0|2<3=0<0/I:0-0/J:1
11 c^an-j+ia_sh'1#2/A:2(4;0/B:2+4;11+29/C:2)7;7)0/E:2-0-0;0/F:2]6]17;2]1]6]1=6|12~ia|2/G:4#4#0;0/H:r<v<ni;2<0|2<3=0<0/I:0-0/J:1
12 an-j+ia+sh_en'2#1/A:2(4;0/B:2+4;11+29/C:2)7;7)0/E:2-0-0;0/F:2]6]17;2]1]6]1=6|12~ia|2/G:4#4#0;0/H:r<v<ni;2<0|2<3=0<0/I:0-0/J:1
13 j^ia-sh+en_l'1#2/A:2(4;11)0/B:2+7;7+29/C:2)2;11)0/E:2-6-0;2/F:4]4]17;1]4]1]4=7|11~en|0/G:2#7#0;0/H:v<ni<n;0<0|3<2=0<0/I:3-0/J:1
14 ia^sh-en+l_an'2#1/A:2(4;11)0/B:2+7;7+29/C:2)2;11)0/E:2-6-0;2/F:4]4]17;1]4]1]4=7|11~en|0/G:2#7#0;0/H:v<ni<n;0<0|3<2=0<0/I:3-0/J:1
15 sh^en-l+an_x'1#2/A:2(4;11)0/B:2+7;7+29/C:2)2;11)0/E:2-6-0;0/F:4]4]17;2]3]2]3=8|10~an|0/G:2#7#0;0/H:v<ni<n;0<0|3<2=0<0/I:3-0/J:1
16 en^l-an+x_ve'2#1/A:2(4;11)0/B:2+7;7+29/C:2)2;11)0/E:2-6-0;0/F:4]4]17;2]3]2]3=8|10~an|0/G:2#7#0;0/H:v<ni<n;0<0|3<2=0<0/I:3-0/J:1
17 l^an-x+ve_van'1#2/A:2(4;11)0/B:2+7;7+29/C:1)2;11)0/E:2-6-0;0/F:4]4]17;3]2]3]2=9|9~ve|0/G:2#7#0;2/H:v<ni<n;0<0|3<2=0<0/I:3-0/J:1
18 an^x-ve+van_v'2#1/A:2(4;11)0/B:2+7;7+29/C:1)2;11)0/E:2-6-0;0/F:4]4]17;3]2]3]2=9|9~ve|0/G:2#7#0;2/H:v<ni<n;0<0|3<2=0<0/I:3-0/J:1
19 x^ve-van+v_in'1#1/A:2(4;11)0/B:1+7;7+29/C:1)2;11)0/E:2-6-0;0/F:4]4]17;4]1]4]1=10|8~van|2/G:2#7#0;0/H:v<ni<n;0<0|3<2=0<0/I:1-0/J:1
20 ve^van-v+in_h'1#1/A:1(7;7)0/B:1+2;11+29/C:1)4;0)0/E:4-4-0;2/F:2]7]17;1]2]1]7=11|7~v|0/G:2#0#0;1/H:ni<n<v;0<1|4<1=0<0/I:1-0/J:1
21 van^v-in+h_e'1#1/A:1(7;7)0/B:1+2;11+29/C:2)4;0)0/E:4-4-0;0/F:2]7]17;2]1]2]6=12|6~in|1/G:2#0#0;0/H:ni<n<v;0<1|4<1=0<0/I:1-0/J:1
22 v^in-h+e_ch'1#2/A:1(2;7)0/B:2+4;11+29/C:2)5;0)0/E:2-4-0;1/F:2]7]17;1]2]3]5=13|5~e|0/G:3#0#0;1/H:n<v<n;1<0|5<0=0<2/I:1-0/J:1
23 in^h-e+ch_eng'2#1/A:1(2;7)0/B:2+4;11+29/C:2)5;0)0/E:2-4-0;1/F:2]7]17;1]2]3]5=13|5~e|0/G:3#0#0;1/H:n<v<n;1<0|5<0=0<2/I:1-0/J:1
24 h^e-ch+eng_k'1#2/A:2(2;7)0/B:2+4;11+29/C:2)5;0)0/E:2-4-0;0/F:2]7]17;2]1]4]4=14|4~eng|1/G:3#0#0;0/H:n<v<n;1<0|5<0=0<2/I:1-0/J:1
25 e^ch-eng+k_e'2#1/A:2(2;7)0/B:2+4;11+29/C:2)5;0)0/E:2-4-0;0/F:2]7]17;2]1]4]4=14|4~eng|1/G:3#0#0;0/H:n<v<n;1<0|5<0=0<2/I:1-0/J:1
26 ch^eng-k+e_ch'1#2/A:2(4;7)0/B:2+5;11+29/C:2)0;0)0/E:2-4-0;1/F:3]7]17;1]3]5]3=15|3~e|0/G:0#0#0;0/H:v<n<w;0<0|0<0=0<1/I:1-0/J:1
27 eng^k-e+ch_eng'2#1/A:2(4;7)0/B:2+5;11+29/C:2)0;0)0/E:2-4-0;1/F:3]7]17;1]3]5]3=15|3~e|0/G:0#0#0;0/H:v<n<w;0<0|0<0=0<1/I:1-0/J:1
28 k^e-ch+eng_pau'1#2/A:2(4;7)0/B:2+5;11+29/C:1)0;0)0/E:2-4-0;0/F:3]7]17;2]2]6]2=16|2~eng|0/G:0#0#0;3/H:v<n<w;0<0|0<0=0<1/I:1-0/J:1
29 e^ch-eng+pau_sil'2#1/A:2(4;7)0/B:2+5;11+29/C:1)0;0)0/E:2-4-0;0/F:3]7]17;2]2]6]2=16|2~eng|0/G:0#0#0;3/H:v<n<w;0<0|0<0=0<1/I:1-0/J:1
30 ch^eng-+sil_X'1#1/A:2(4;7)0/B:1+5;11+29/C:0)0;0)0/E:2-4-0;0/F:3]7]17;3]1]7]1=17|1~X|3/G:0#0#0;0/H:n<w<X;0<0|0<0=0<1/I:1-0/J:1
31 eng^pau-sil+X_X'0#0/A:1(5;11)29/B:0+0;0+0/C:0)0;0)0/E:3-7-17;3/F:0]0]0;0]0]0=0~X|0/G:0#0#0;0/H:w<X<X;0<0|0<0=0<0/I:3-0/J:1
```



## 概述-声学模型回顾

### 声学模型建模目标

建模文本序列 $X$  和语音特征 $Y$ 之间的映射关系

文本 $X = \{x_1, \dots, x_M\}$ ,  $M$ 帧音素

语音 $Y = \{y_1, \dots, y_N\}$ ,  $N$ 帧声学特征

建模目标:

$$\hat{\Lambda} = \arg \max_{\Lambda} p(Y|X, \Lambda)$$

文本特征

语音特征

声学模型参数

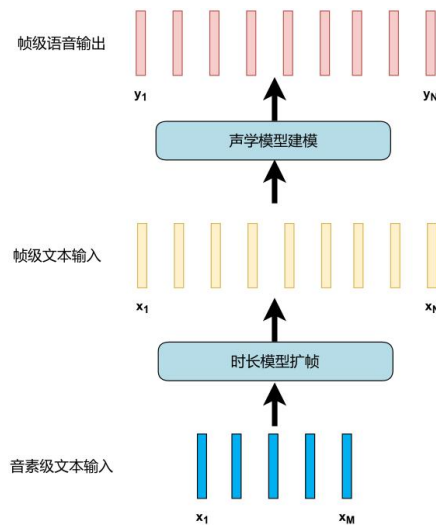


# 概述-声学模型回顾

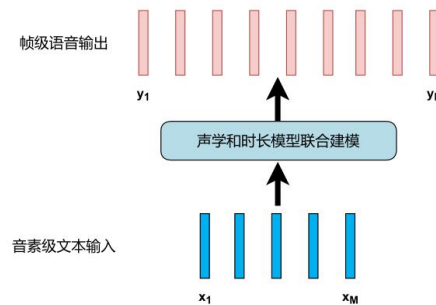
## 声学模型建模过程

- 帧级声学模型
  - 时长模型 (M帧->N帧)
  - N帧->N帧
- 序列到序列声学模型
  - M帧->N帧

$$\hat{\Lambda} = \arg \max_{\Lambda} p(Y|X, \Lambda)$$



(a) 帧级声学模型



(b) 序列到序列声学模型



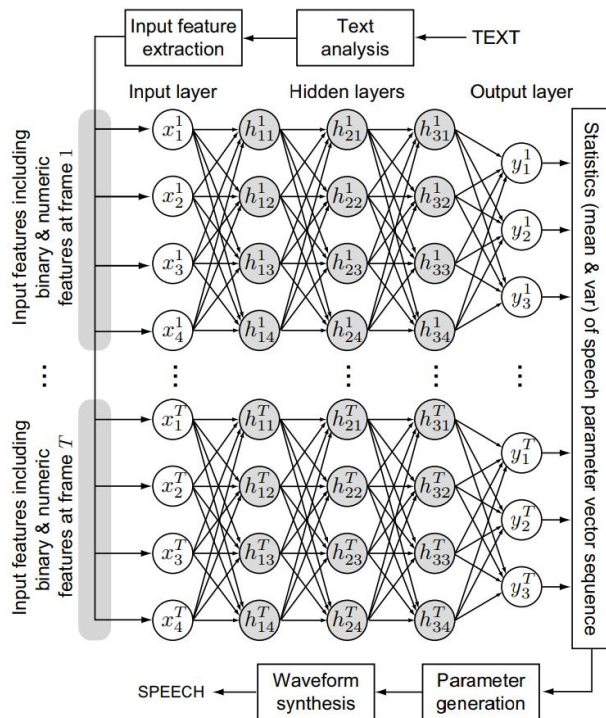
## 帧级声学模型建模过程

### 缺点

- 声学模型文本输入长度与语音序列**长度必须一致**。
- 由于文本本身不含有时序信息，因此帧级声学模型通常需要**额外的时长模型**给文本引入时间信息。
- 时长模型的误差影响后续声学模型的效果：**误差累积问题**。

### 优点

- 由于帧级声学模型是逐帧对应的，其**稳定性**很好。







## 1. 概述



## 2. 序列到序列模型及注意力机制



## 3. 序列到序列声学模型Tacotron



## 4. 序列到序列声学模型变体（一）



## 5. 实战



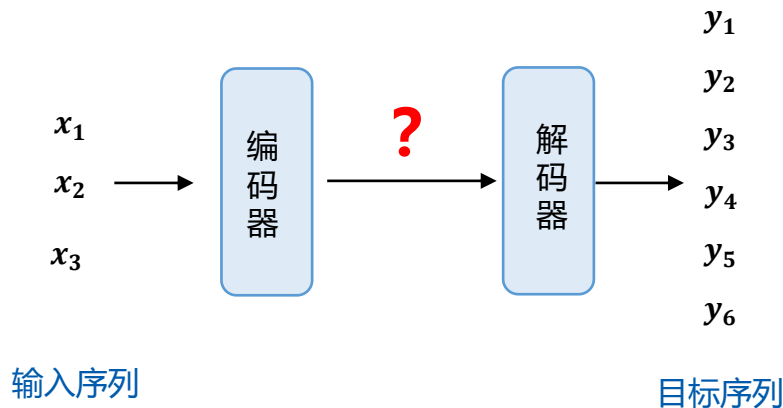
## 概述

序列到序列模型的核心思想：在学习不等长序列映射关系的同时，建模序列间的时序关系，从而达到序列到序列建模的目的。

典型的序列到序列模型包含一个**编码器网络**和一个**自回归解码器网络**，

- 编码器：建模输入序列的内在联系
- 解码器：生成目标序列

核心问题：如何对齐两个不等长序列





# 序列到序列模型结构

## 编码器Encoder

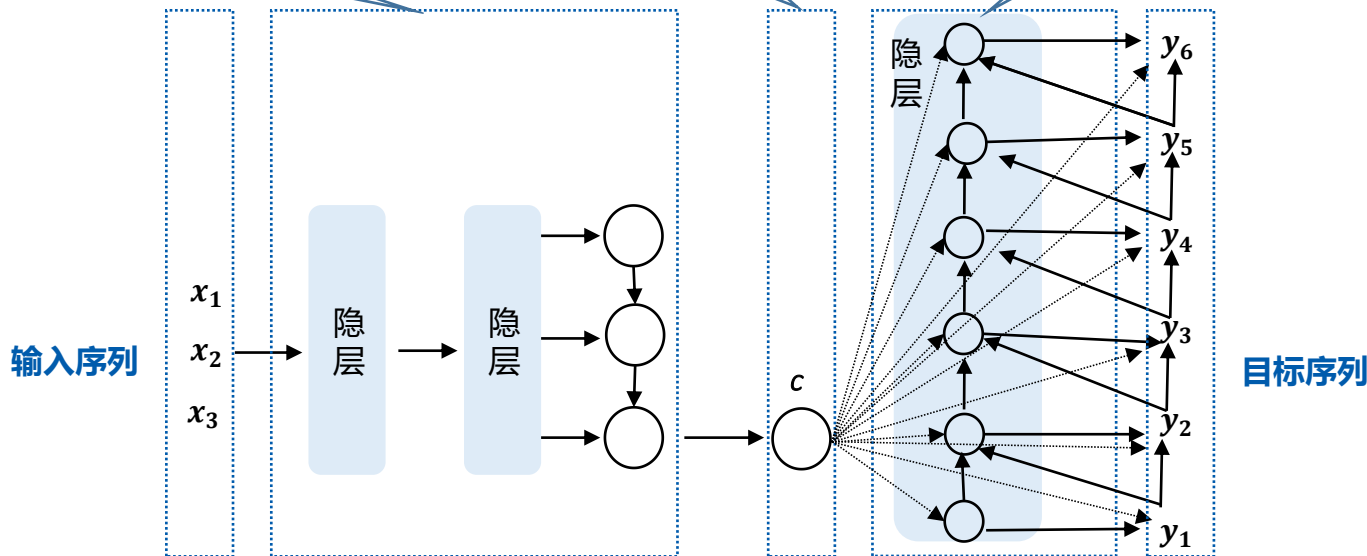
作用：学习输入  
将其编码成内容向量

## 内容向量 C

作用：连接编码器和解码  
器的中间状态向量

## 解码器Decoder

作用：通过对状态向量的学习来进行输出  
对应的序列

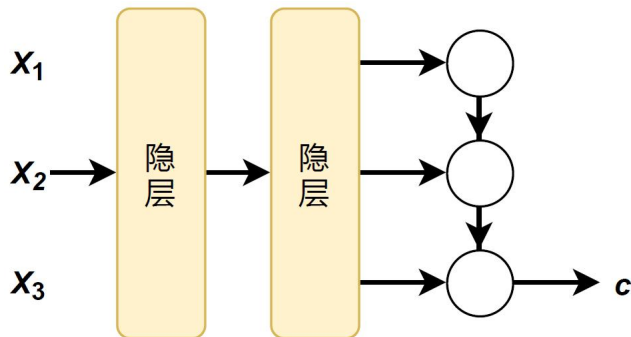




## 编码器

学习输入序列 $X$ 的固定长度隐层表示向量  $c$        $c = \text{Encoder}(X), c \in \mathbb{R}^D$

- 通常采用RNN结构，其最后一个时刻的输出作为输入序列的固定长度表示
- 固定长度的隐层向量可以看做是整个输入序列 $X$ 中所有信息的压缩表示
- 该表示（内容向量）直接用来控制解码器的生成过程





### 自回归模型

- 经典的生成式模型，建模目标数据 $x$ 的分布 $P(x)$
- 通过相邻点或前一个点来直接建模下一个点的分布，从而获得整个数据 $x$ 的分布
- 给定初始数据点，可以无限制（长度不确定）、无约束（内容不确定）的生成无穷个符合分布 $P(x)$ 的数据点 → 无条件自回归过程
- 针对语音来讲，符合语音的单调时序生成特点（语音的每一个采样点）

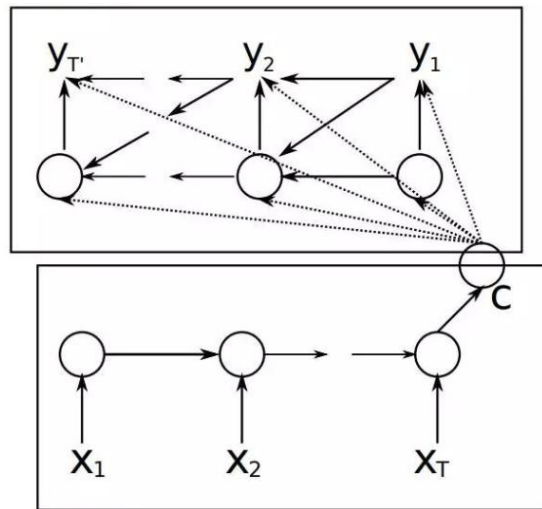
$$p_{data}(x) = \prod_{i=1}^n p(x_i | x_1, x_2, \dots, x_{i-1})$$



## 解码器

- 自回归地生成目标序列每一帧
  - 条件自回归模型  $p(Y|X) = \prod_{t=1}^N p(y_t | c, y_1, \dots, y_{t-1})$
  - 利用解码器的初始状态 $s_0$ 和编码器的内容向量 $c$ 生成目标序列第一帧
  - 通过自回归过程循环迭代生成目标序列 $Y$ 的每一帧，直到合成到EOS为止
  - 每一个时刻都需要使用编码器学习到的内容向量

$$\hat{y}_t = \text{Decoder}(c, y_{t-1}, s_t)$$





## 经典的序列到序列模型

- 给定输入序列 $X$ 和目标序列 $Y$ , 最大化给定输入 $X$ 下生成 $Y$ 的概率

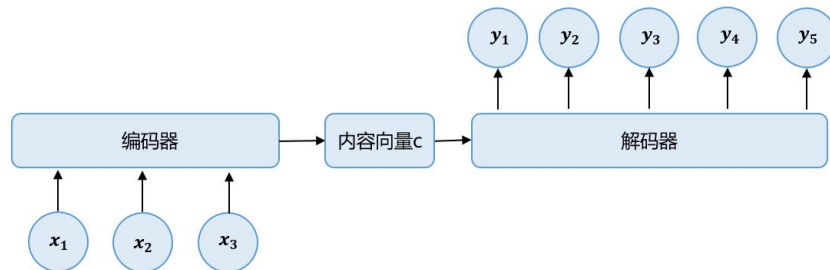
- 条件自回归 
$$p(\mathbf{Y}|\mathbf{X}) = \prod_{t=1}^N p(\mathbf{y}_t | \mathbf{c}, \mathbf{y}_1, \dots, \mathbf{y}_{t-1})$$

- 编码器: 学习输入序列 $X$ 的 $D$ 维的固定长度隐层表示向量 $\mathbf{c}$

$$\mathbf{c} = \text{Encoder}(X), \mathbf{c} \in \mathbb{R}^D$$

- 解码器: 根据内容向量 $\mathbf{c}$ 和 $t$ 时刻的隐状态 $\mathbf{s}_t$ , 自回归地生成目标序列每一帧

$$\hat{\mathbf{y}}_t = \text{Decoder}(\mathbf{c}, \mathbf{y}_{t-1}, \mathbf{s}_t)$$





## 经典的序列到序列模型

- 优点

- 不再依赖额外的时长模型，直接通过解码器的自回归生成过程来获取变长的目标序列；
- 避免了时长模型和声学模型级联时的累积误差；
- 简化了语音合成的流程。

- 缺点

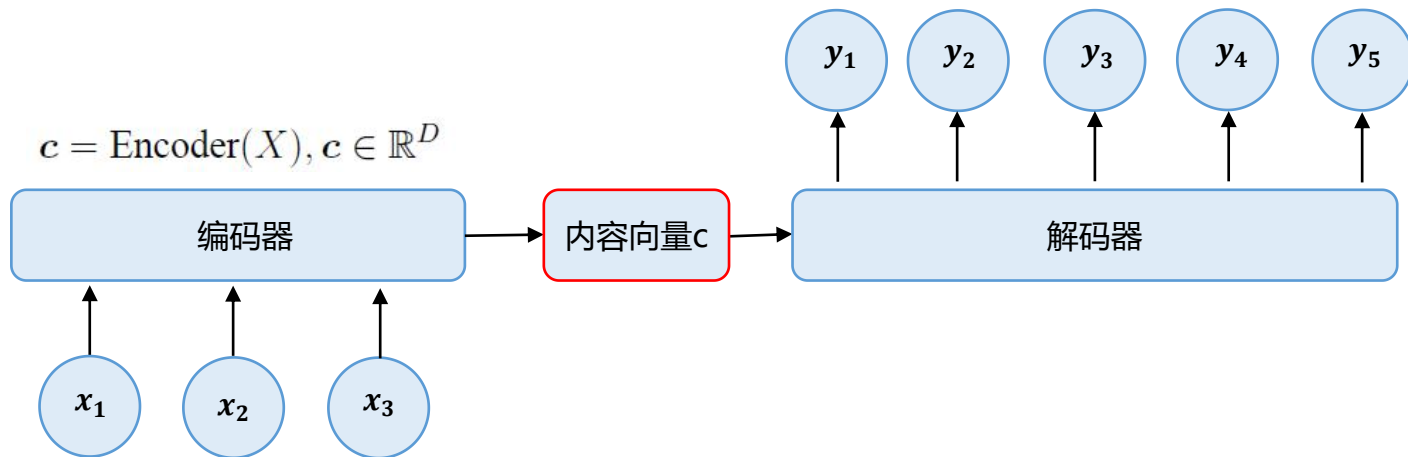
- 编码器对M帧输入序列X压缩成一帧向量，过度压缩存在信息损失；
- 解码器的每一个时刻使用的是同一个内容向量，无法考虑输入序列中不同部分对当前时刻贡献，例如语音合成中一段语音的内容主要是由某一个字所决定的。





## 经典的序列到序列模型

- 解决方法
  - 编码器不再压缩输入序列，并且解码器生成过程中的不同时刻使用**特定的内容向量**
- 固定的内容向量
  - 在编码器学习输入序列 $X$ 的隐层表示时，直接获得内容向量





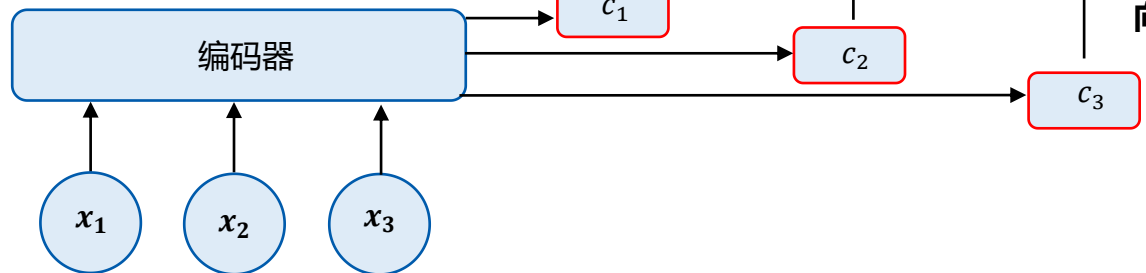
## 基于注意力机制的序列到序列模型

- 变化的内容向量
- 核心问题是求解对齐矩阵

(1) 编码器只学习输入序列 $X$ 的隐层表示, 不做时间维度的压缩

- $M$ 帧输入序列 $X$ 的隐层表示

$$L = \text{Encoder}(X), L \in \mathbb{R}^{M \times D}$$



(2) 在解码器的 $t$ 时刻提供不同的内容向量 $c_t$ , 其中 $c_t$ 为 $L$ 中每一帧的加权和

$$c_t = \sum_{s=1}^M \alpha_{ts} l_s$$

对齐矩阵



## 注意力机制

- 目标
  - 求解对齐矩阵 $\alpha$ ，以获得每一个解码器时刻所对应的内容向量

$$\alpha_{ts} = \frac{\exp(\text{score}(\mathbf{h}_t, \mathbf{l}_s))}{\sum_{s'=1}^M \exp(\text{score}(\mathbf{h}_t, \mathbf{l}_{s'}))}$$

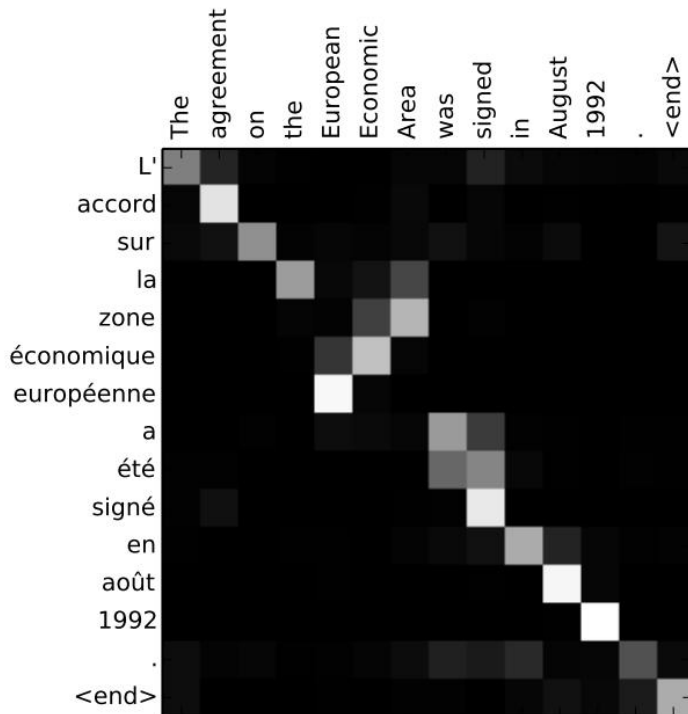
$$\mathbf{c}_t = \sum_{s=1}^M \alpha_{ts} \mathbf{l}_s$$

- 给定目标序列 $Y$ 及解码器的隐层输出 $\mathbf{h}$ ，输入序列 $X$ 的隐层表示 $L$ 
  - 计算每一个时刻 $\mathbf{h}_t$ 在每一帧 $\mathbf{l}_s$ 上的得分（打分函数）



## 注意力机制

- 对齐矩阵示例

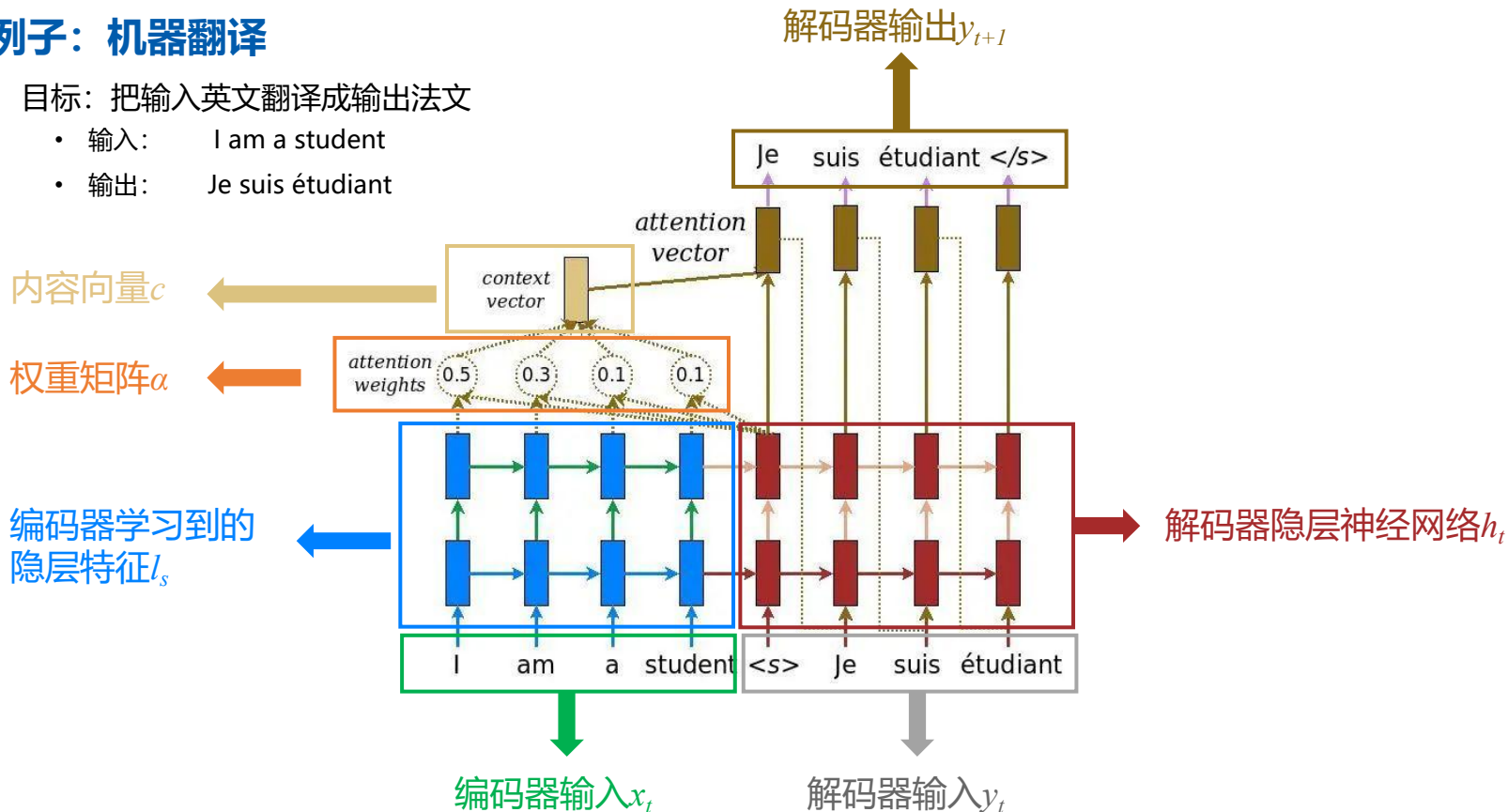




# 基于注意力机制的序列到序列模型

## 例子：机器翻译

- 目标：把输入英文翻译成输出法文
  - 输入： I am a student
  - 输出： Je suis étudiant





# 基于注意力机制的序列到序列模型

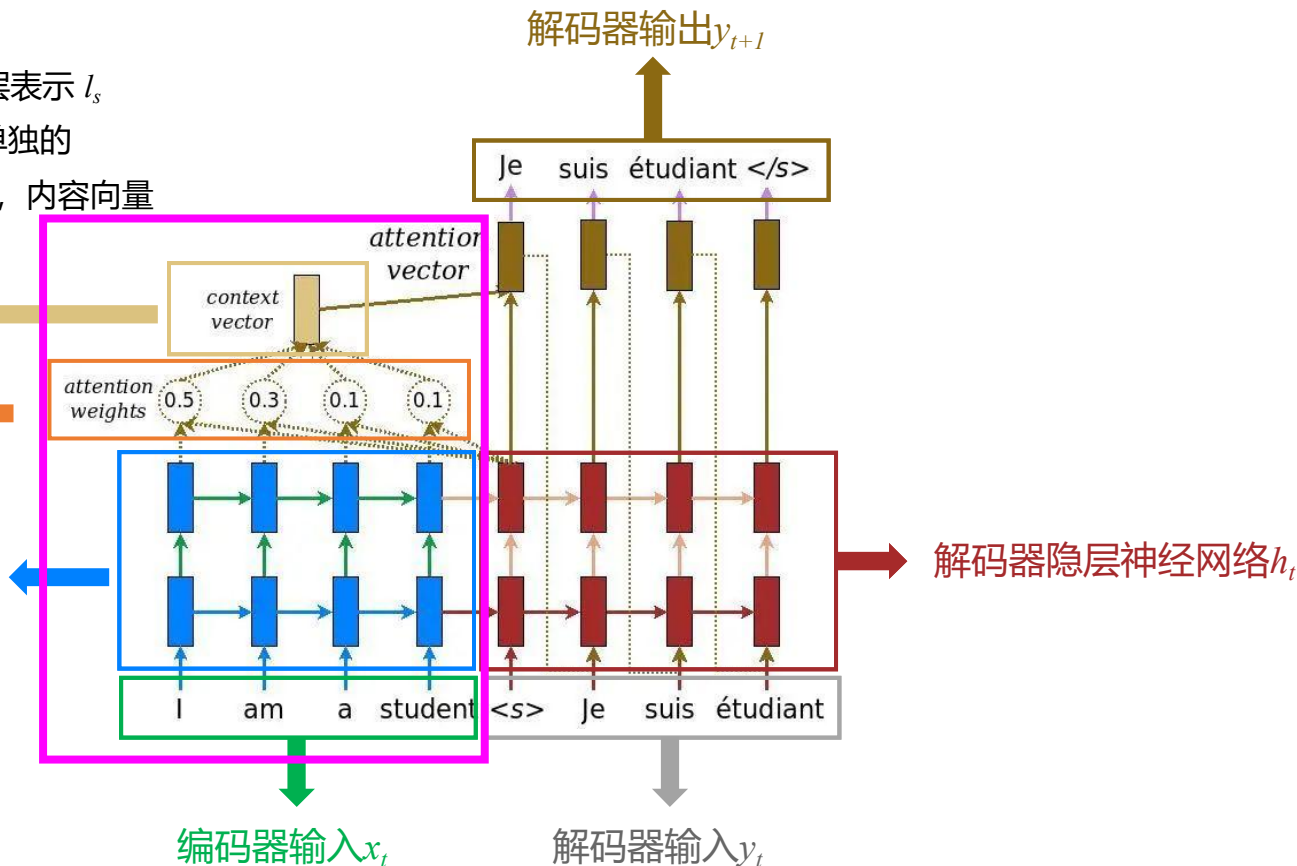
## 编码器

- 把输入  $x_t$  映射到隐层表示  $l_s$
- 目前的编码器都是单独的
- 打分函数+softmax, 内容向量

内容向量  $c$

权重矩阵  $\alpha$

编码器学习到的  
隐层特征  $l_s$

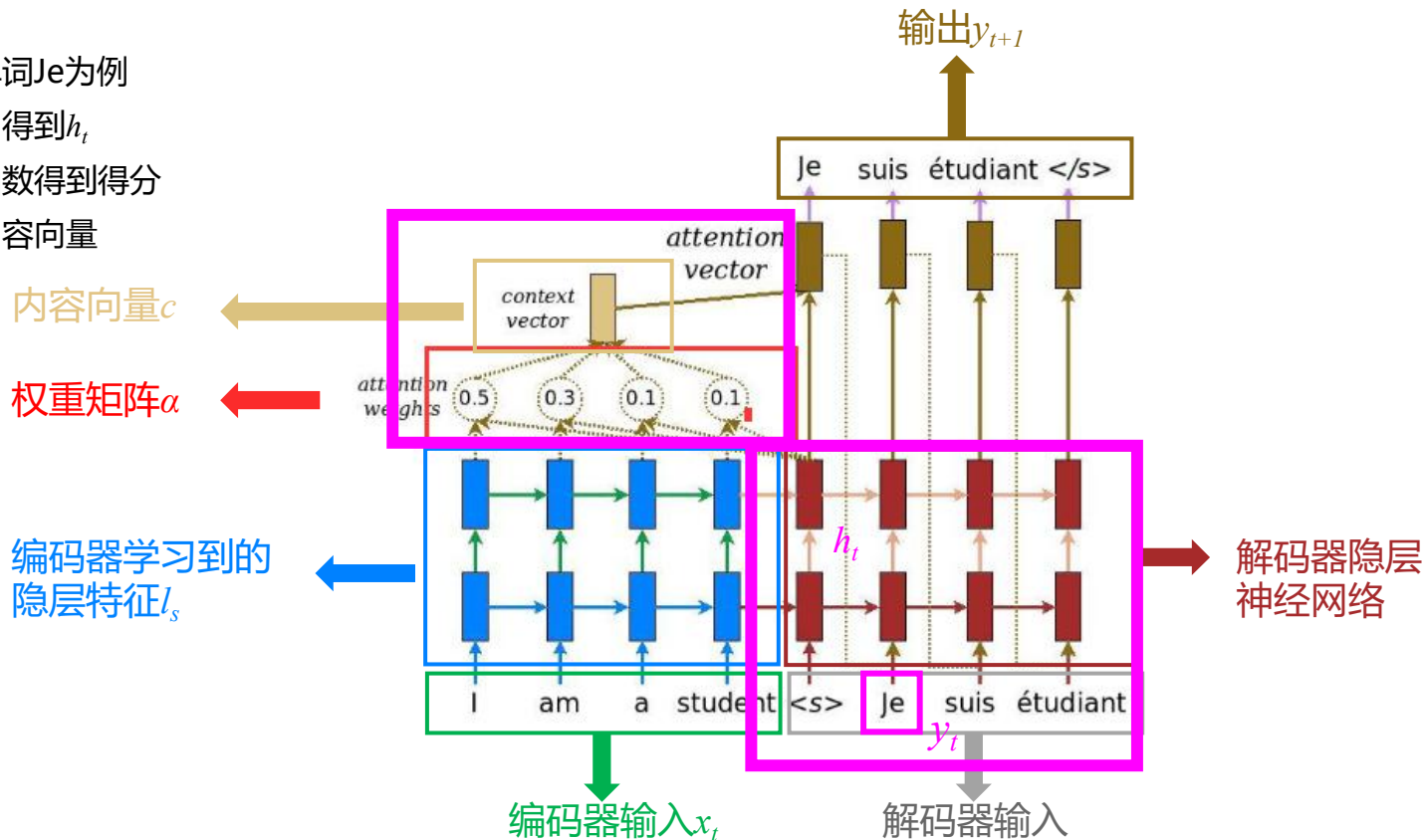




# 基于注意力机制的序列到序列模型

## 解码器

- 以第一个单词Je为例
- 隐层变换, 得到 $h_t$
- 通过打分函数得到得分
- 相乘得到内容向量

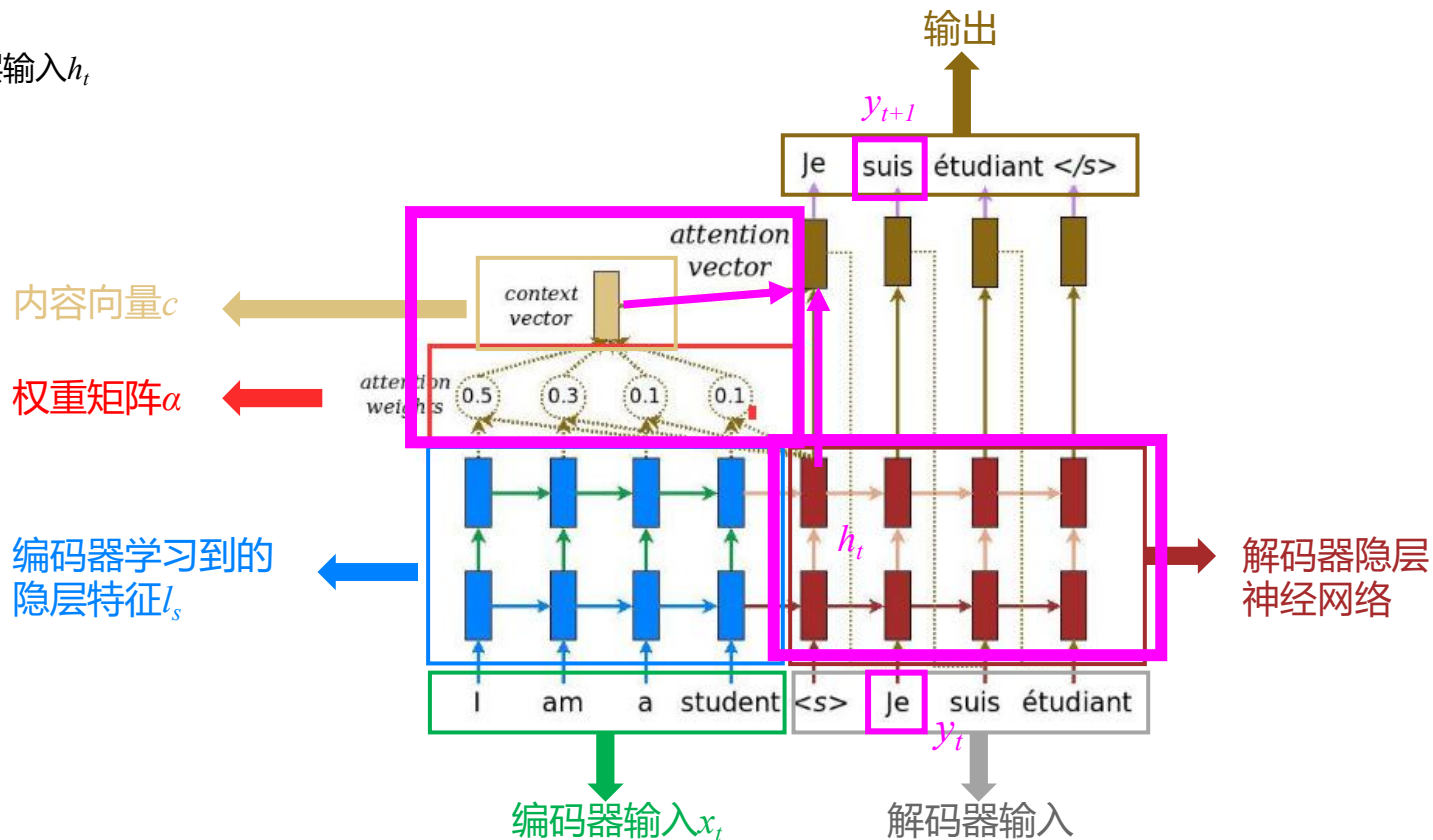




# 基于注意力机制的序列到序列模型

## 解码器

- 解码器隐层输入  $h_t$
- 内容向量  $c$
- $\rightarrow$  输出  $y_{t+1}$



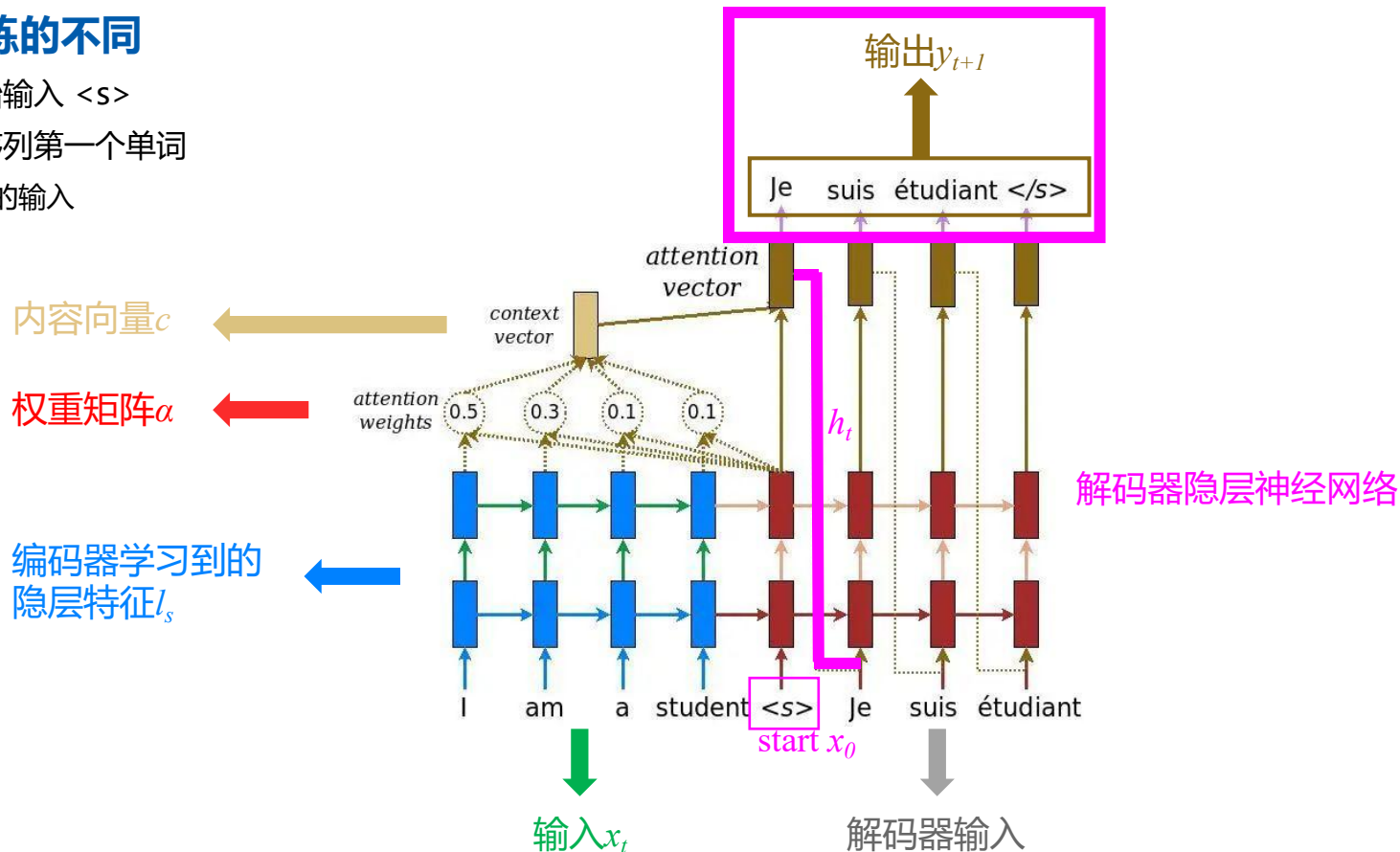




# 基于注意力机制的序列到序列模型

## 合成与训练的不同

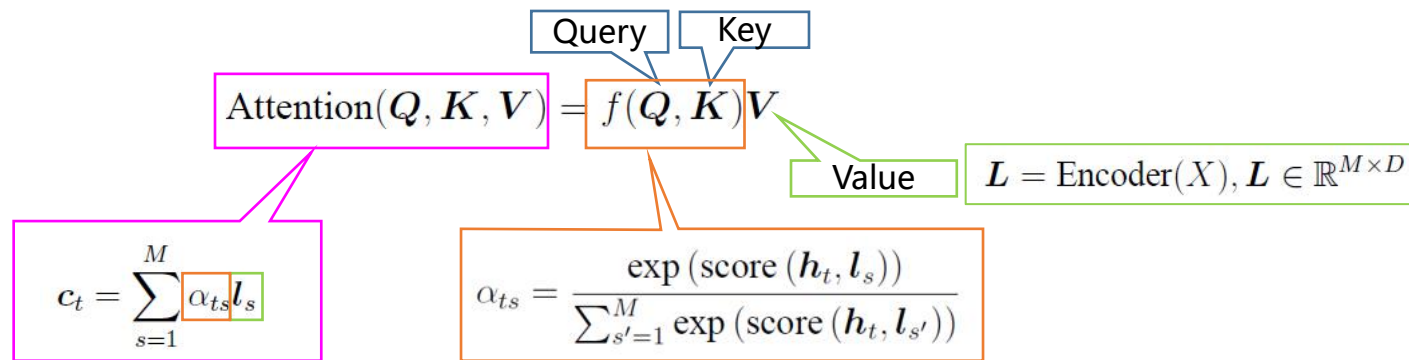
- 解码器初始输入  $\langle s \rangle$
- 得到目标序列第一个单词
- 作为下一个的输入





## 注意力机制概述

- 给定解码器 $t$ 时刻的向量，求解权重矩阵 $\alpha_t$ ，从而获得其对应的内容向量 $c_t$ 
  - Value: 值向量，由编码器编码获得，矩阵大小 $M \times D_1$ ， $M$ 为输入序列长度， $D_1$ 为维度
  - Query: 查询向量，一般为解码器隐层输出，其矩阵大小为 $N \times D_2$ ， $N$ 为目标序列长度， $D_2$ 为维度
  - Key: 键向量，作为计算权重矩阵的索引，一般由Value通过前馈网络获得，矩阵大小 $D_2 \times M$ ，
  - 最后计算得到的内容向量矩阵 $C$ 的大小为 $N \times D_1$ ，其长度和目标序列一致





## 序列到序列模型

### 注意力机制代码讲解

[https://github.com/tensorflow/tensorflow/blob/r1.12/tensorflow/contrib/seq2seq/python/ops/attention\\_wrapper.py](https://github.com/tensorflow/tensorflow/blob/r1.12/tensorflow/contrib/seq2seq/python/ops/attention_wrapper.py)



## 注意力机制分类

- Hard (stochastic) attention
  - 依靠概率采样编码器的一部分来计算内容向量
  - 蒙特卡洛采样方法来计算梯度
- Soft (deterministic) attention (主流)
  - 计算编码器表示中每一帧的确定性得分来获取内容向量
  - 直接利用反向传播来计算梯度

$$\text{Attention}(Q, K, V) = f(Q, K)V$$

Diagram illustrating the components of the Attention function:  $Q$  is labeled Query,  $K$  is labeled Key, and  $V$  is labeled Value.

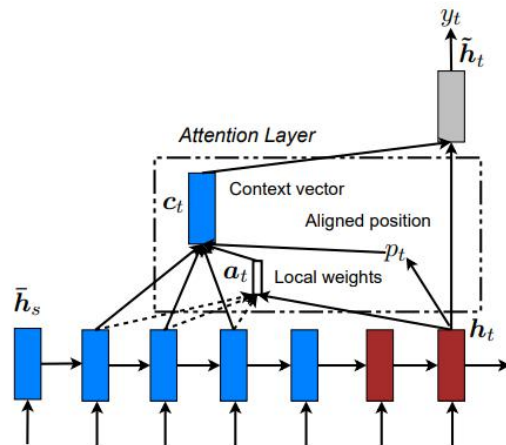


## 注意力机制分类

- Global attention

- 解码器每一个时刻的内容向量计算过程均考虑编码器的所有输入信息
- 能够考虑输入序列X的所有信息
- 序列X过长时计算量开销较大

$$\begin{aligned} a_t(s) &= \text{align}(h_t, \bar{h}_s) \\ &= \frac{\exp(\text{score}(h_t, \bar{h}_s))}{\sum_{s'} \exp(\text{score}(h_t, \bar{h}_{s'}))} \end{aligned}$$





## 注意力机制分类

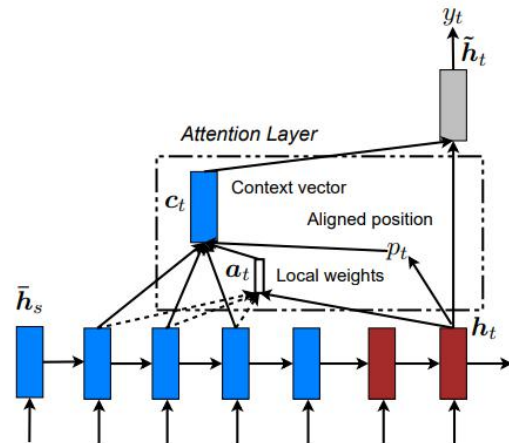
- Local attention

- 解码器 $t$ 时刻只考虑编码器输出的部分信息  $[p_t - D, p_t + D]$
- 可以看做介于soft 和hard attention之间, 更易训练

$$p_t = S \cdot \text{sigmoid}(\mathbf{v}_p^\top \tanh(\mathbf{W}_p \mathbf{h}_t)),$$

$$\mathbf{a}_t(s) = \text{align}(\mathbf{h}_t, \bar{\mathbf{h}}_s) \exp\left(-\frac{(s - p_t)^2}{2\sigma^2}\right)$$

- Etc. (GMM based)





## 1. 概述



## 2. 序列到序列模型及注意力机制



## 3. 序列到序列声学模型Tacotron



## 4. 序列到序列声学模型变体（一）

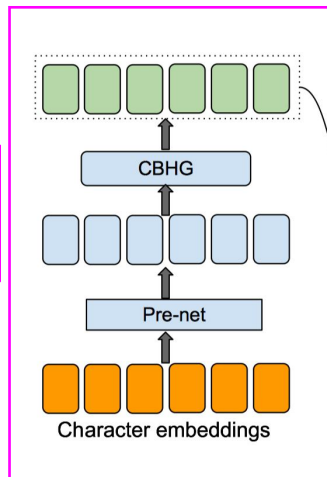


## 5. 实战

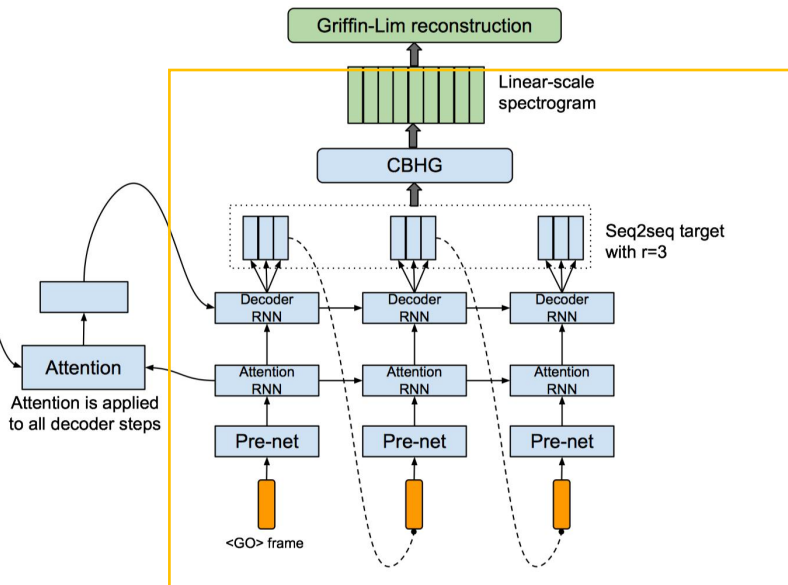


## Tacotron模型框架结构

**编码器：**  
字符编码作为输入，  
学习文本表示



**后处理：**  
Griffin-Lim声码器，把梅尔谱还原成语音信号



**解码器：**  
自回归地生成语音的梅尔谱

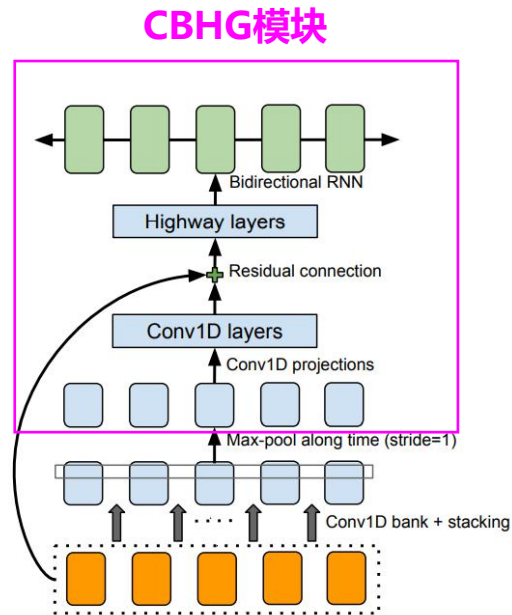




# 序列到序列声学模型Tacotron

## 编码器

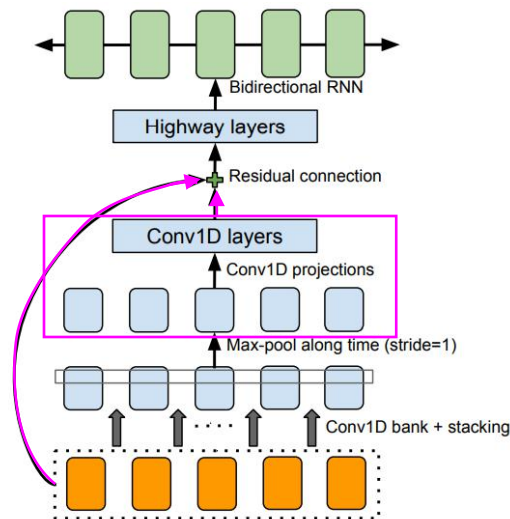
- Conv bank
  - 目的：学习不同尺度context的文本表示
  - 组成：多个卷积核大小不一的卷积神经网络
- 输出拼接stacking
- Max-pool along time
- CBHG模块
  - 一维卷积
  - Highway network
  - 双向GRU





## CBHG模块

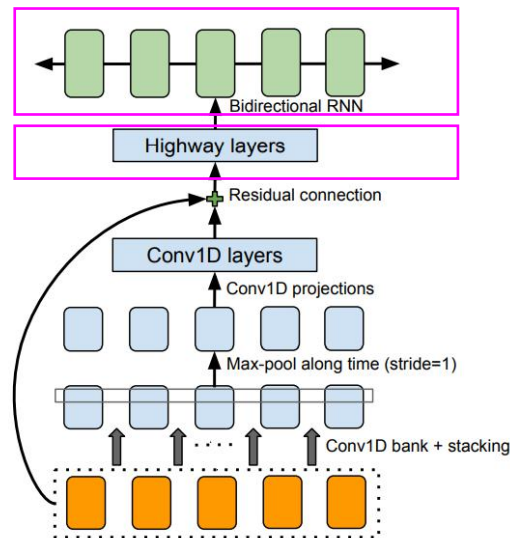
- 一维卷积
  - 原因：多个卷积神经网络的输出信息维度过高
  - 目的：压缩融合多个CNN的高维输出信息
- 残差连接
  - 输入信息直接送到隐层网络
  - 目的：减少随机梯度下降过程中的优化难度  
(深层网络梯度消失)





## CBHG模块

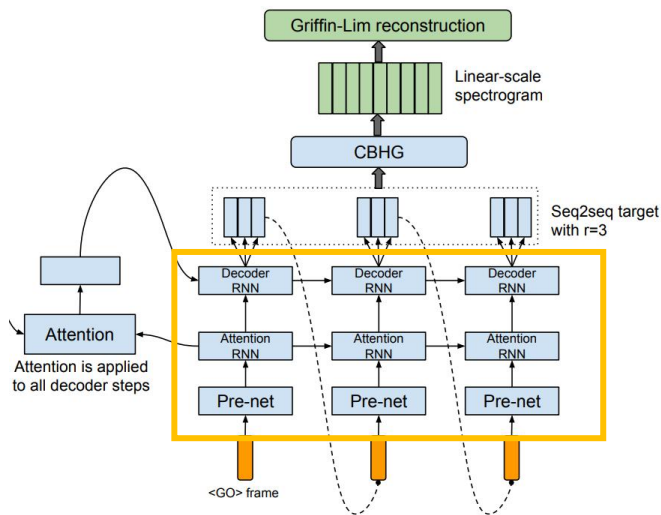
- Highway network  $y_i = H_i(\mathbf{x}) * T_i(\mathbf{x}) + x_i * (1 - T_i(\mathbf{x}))$ 
  - 隐层输入 $x$ 的一部分经过非线性变换，一部分直接传递给  $y$
  - $T(x)$ 作为门控单元， $T(x) = 0$ 时输入信息直接传给  $y$
  - 反向传播时，可以让更多的信息直接回流到输入，不需要经过非线性变换
  - 深度网络的收敛性能好
- Bi-GRU:
  - 双向RNN网络，学习序列的上下文依赖





## 自回归解码器网络

- Pre-net
  - 自回归训练阶段,  $t$ 时刻使用 $t - 1$ 时刻的真实梅尔谱
  - 生成阶段,  $t$ 时刻使用 $t - 1$ 时刻模型输出的梅尔谱
  - 利用dropout提升模型泛化能力
- Attention RNN:
  - 隐层中使用RNN建模序列的时序依赖
- Decoder RNN:
  - 建模内容向量和隐层特征的时序依赖





## 自回归解码器网络

- 注意力机制

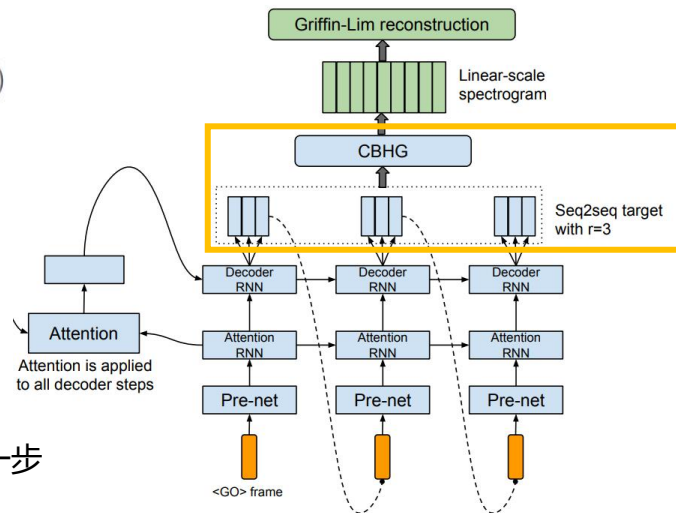
- Global soft attention  $\text{score}(h_t, l_s) = v^\top \tanh(W_1 h_t + W_2 l_s)$

- 多帧预测

- 文本序列长度和语音特征长度差异过大
- 减少输入和输出序列长度差异，易收敛

- CBHG后处理

- 自回归只能考虑之前时刻的上下文信息，后处理网络能够进一步考虑整个序列信息，最后生成线性谱



$$[\hat{y}_{t+1}, \hat{y}_{t+2}, \dots, \hat{y}_{t+K}]^\top = \text{DecRNN}(c_t, h_t)$$



## 实验

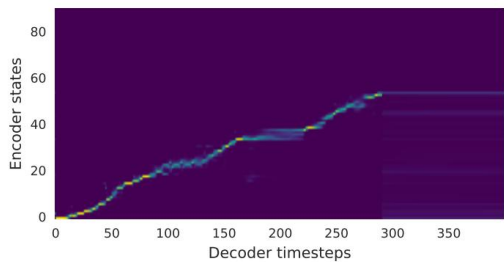
- 文本：英文字符输入，中文音素、韵律边界等
- 语音表示：梅尔谱，50ms帧长，12.5ms帧移
- 多帧预测：一次预测两帧
- 语音还原：Griffin-Lim

	mean opinion score
Tacotron	$3.82 \pm 0.085$
Parametric	$3.69 \pm 0.109$
Concatenative	$4.09 \pm 0.119$

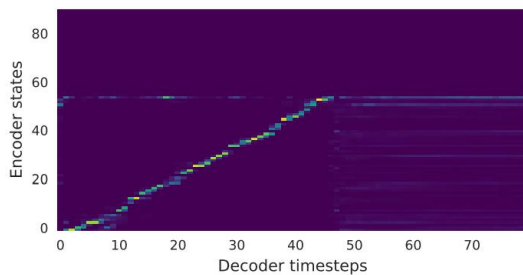
Spectral analysis	<i>pre-emphasis</i> : 0.97; <i>frame length</i> : 50 ms; <i>frame shift</i> : 12.5 ms; <i>window type</i> : Hann
Character embedding	256-D
Encoder CBHG	<i>Conv1D bank</i> : $K=16$ , conv- $k$ -128-ReLU <i>Max pooling</i> : stride=1, width=2 <i>Conv1D projections</i> : conv-3-128-ReLU → conv-3-128-Linear <i>Highway net</i> : 4 layers of FC-128-ReLU <i>Bidirectional GRU</i> : 128 cells
Encoder pre-net	FC-256-ReLU → Dropout(0.5) → FC-128-ReLU → Dropout(0.5)
Decoder pre-net	FC-256-ReLU → Dropout(0.5) → FC-128-ReLU → Dropout(0.5)
Decoder RNN	2-layer residual GRU (256 cells)
Attention RNN	1-layer GRU (256 cells)
Post-processing net CBHG	<i>Conv1D bank</i> : $K=8$ , conv- $k$ -128-ReLU <i>Max pooling</i> : stride=1, width=2 <i>Conv1D projections</i> : conv-3-256-ReLU → conv-3-80-Linear <i>Highway net</i> : 4 layers of FC-128-ReLU <i>Bidirectional GRU</i> : 128 cells
Reduction factor ( $r$ )	2



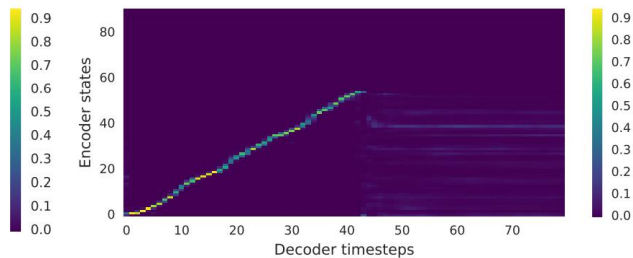
## 注意力机制分析（一次预测5帧）



(a) Vanilla seq2seq + scheduled sampling



(b) GRU encoder

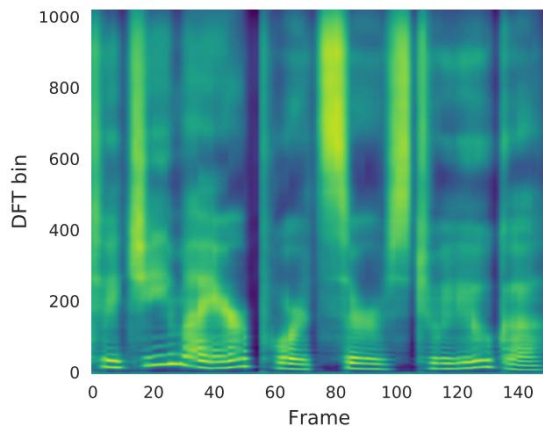


(c) Tacotron (proposed)

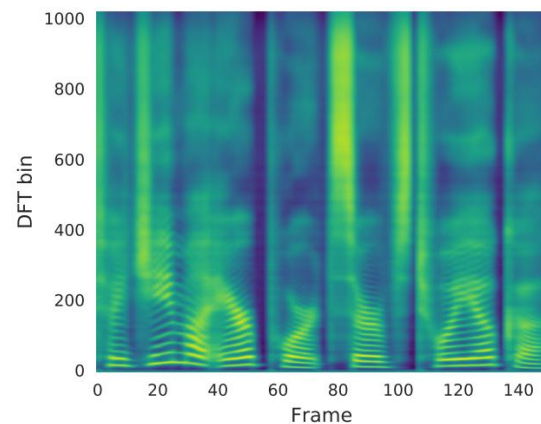
Samples: <https://google.github.io/tacotron/publications/tacotron/index.html>



## Post-net效果分析



(a) Without post-processing net



(b) With post-processing net





# 序列到序列声学模型Tacotron

## Tacotron优势

- 无需复杂的文本前端分析模块
- 无需额外的时长模型
- 极大地简化了声学模型构建流程，减少了语音合成任务对领域知识的依赖



## Tacotron存在的问题

- (1) CBHG网络比较复杂, **参数量**相对较大
- (2) **Griffin-Lim**重建过程中相位问题导致的**语音失真**
- (3) 自回归解码器生成过程中**无法停止**
- (4) Global soft **attention** 的问题
  - 针对语音合成任务稳定性较差
  - 模型训练过程中, 每个时刻预测的语音帧数越少越难训练



## 1. 概述



## 2. 序列到序列模型及注意力机制



## 3. 序列到序列声学模型Tacotron



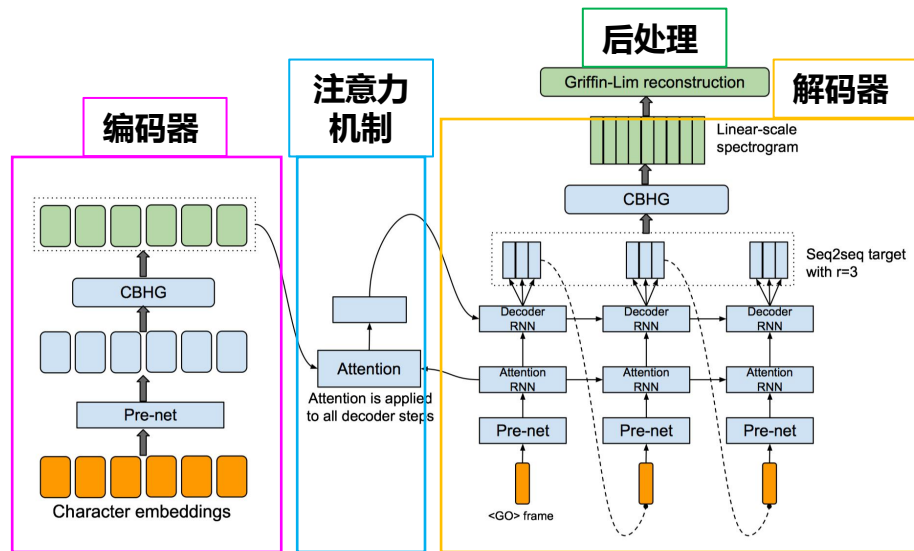
## 4. 序列到序列声学模型变体（一）



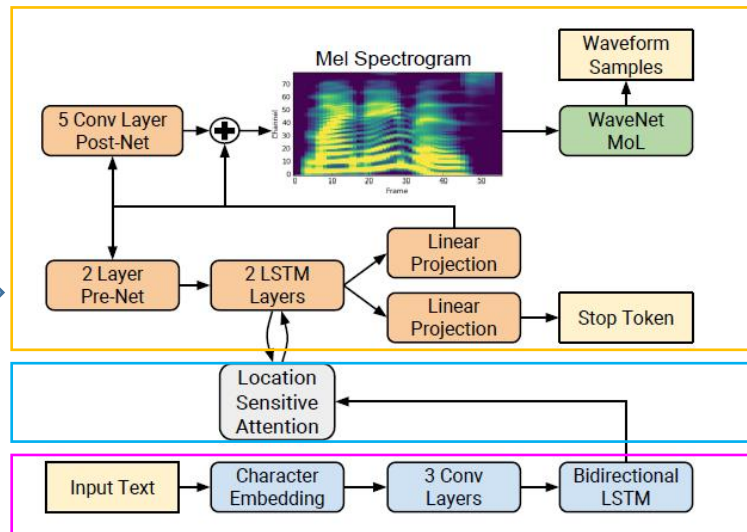
## 5. 实战



## Tacotron模型框架结构



## Tacotron2 模型框架结构



- (1) CBHG导致的参数量问题
- (2) Griffin-Lim语音失真问题

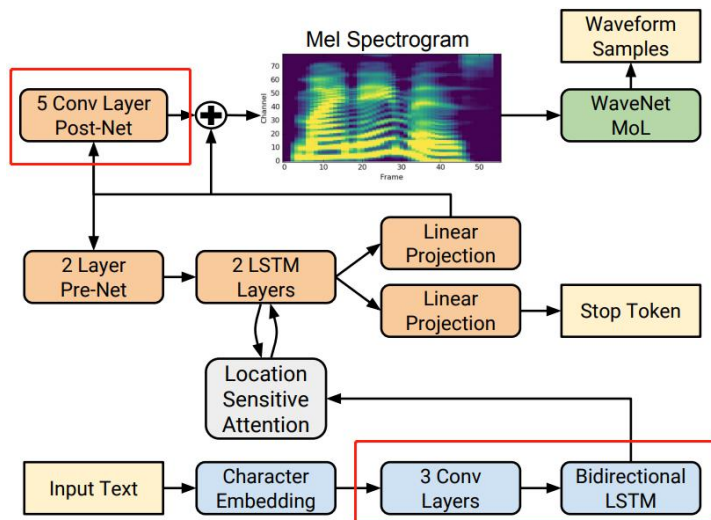
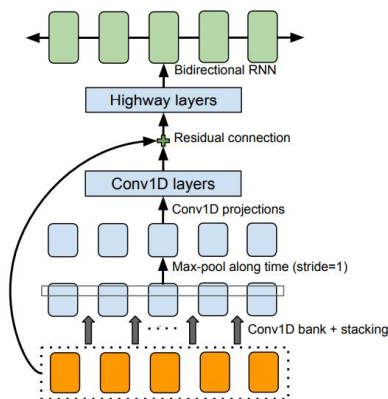
- (3) 解码器生成过程中无法停止的问题
- (4) Global soft attention稳定性问题



# 序列到序列声学模型变体（一）

## Tacotron 2

- 参数量的减少
  - CBHG -> 3 Conv layers + BLSTM 或 5 Conv layers
  - 去除Attention RNN

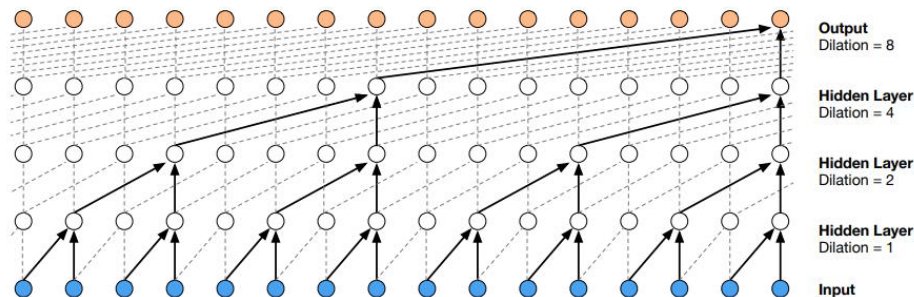




# 序列到序列声学模型变体（一）

## Tacotron 2

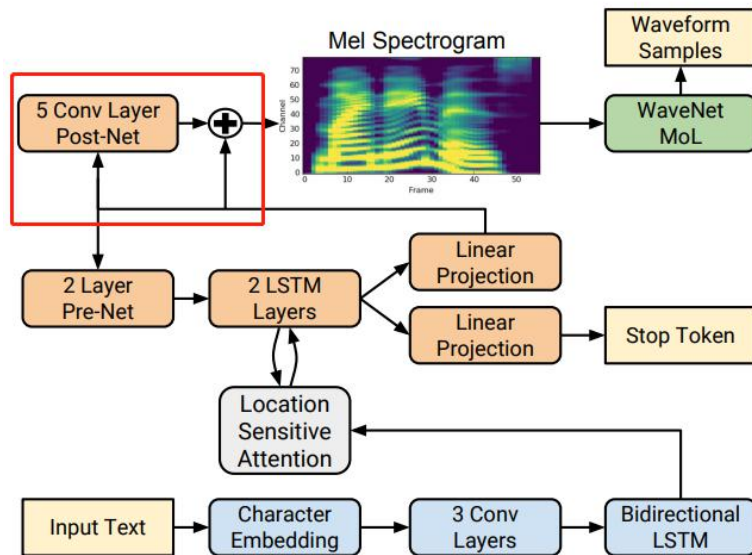
- Griffin-Lim导致的语音失真问题
- WaveNet
  - 条件自回归模型 (梅尔谱作为条件)
  - 直接建模语音采样点
  - 因果卷积：语音信号的单调特性
  - 膨胀卷积：考虑更多的上文信息





## Tacotron 2

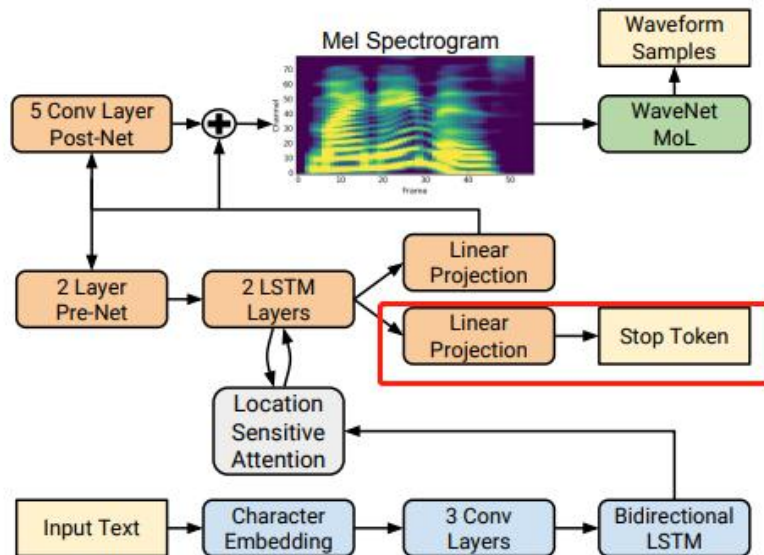
- PostNet的改进
  - CBHG -> 5层卷积
  - 后处理网络的输入和输出均和真实梅尔谱计算L2误差
  - 残差连接，帮助模型收敛





## Tacotron 2

- 自回归生成过程中无法停止的问题
  - 解码的每个时刻预测是否应该停止 (stop token)
  - 通过decoder RNN的输出和注意力机制的内容向量, 使用一层Sigmoid激活的前馈神经网络预测[0,1]之间的停止概率
  - 解码时, 设定阈值判断是否停止生成





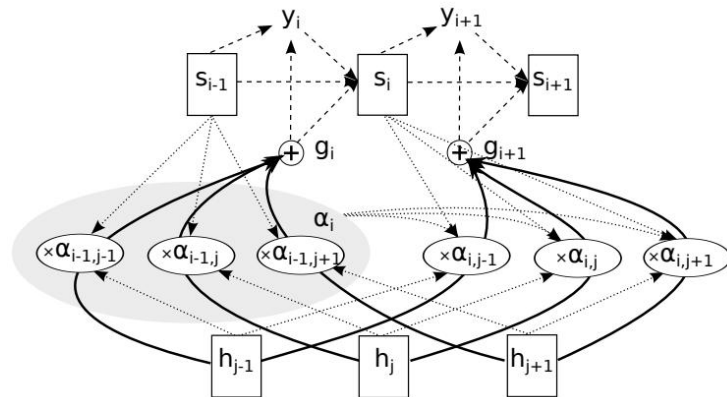


## Tacotron 2

- 注意力机制稳定性问题
  - 传统注意力机制中，每个时刻使用编码器哪些部分是无约束的
  - 语音的产生过程是单调的
- Location sensitive attention (LSA)
  - 在解码器 $t$ 时刻考虑前面时刻的对齐矩阵

$$e_{i,j} = \text{Score}(s_{i-1}, h_j)$$

$$\alpha_{i,j} = \exp(e_{i,j}) / \sum_{j=1}^L \exp(e_{i,j})$$



$$f_i = F * \alpha_{i-1}$$

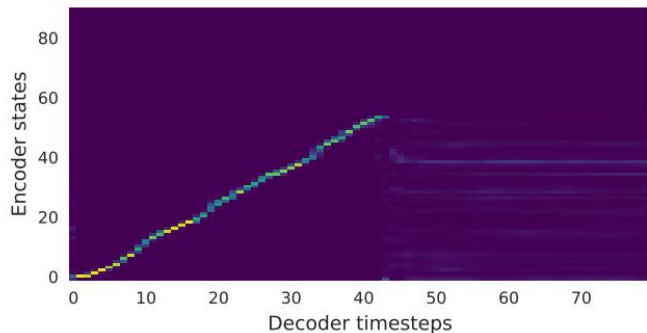
$$e_{i,j} = w^\top \tanh(W s_{i-1} + V h_j + b) \longrightarrow e_{i,j} = w^\top \tanh(W s_{i-1} + V h_j + U f_{i,j} + b)$$



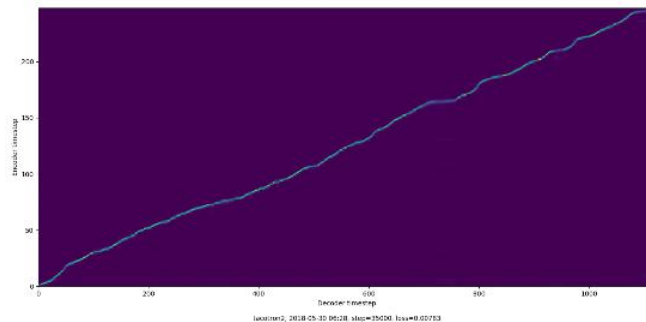
# 序列到序列声学模型变体（一）

## Tacotron 2

- 不同注意力机制对比：



(a) 传统注意力机制



(b) LSA

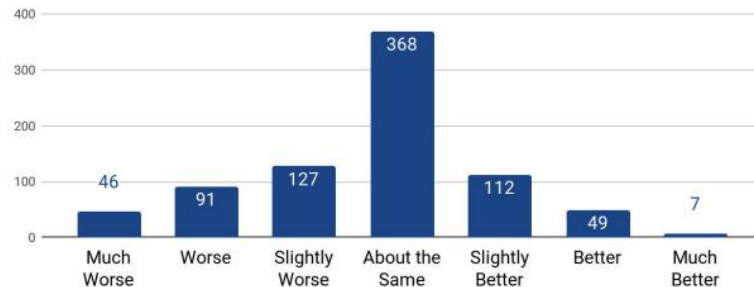


## Tacotron 2

- 实验

System	MOS
Parametric	$3.492 \pm 0.096$
Tacotron (Griffin-Lim)	$4.001 \pm 0.087$
Concatenative	$4.166 \pm 0.091$
WaveNet (Linguistic)	$4.341 \pm 0.051$
Ground truth	$4.582 \pm 0.053$
Tacotron 2 (this paper)	<b><math>4.526 \pm 0.066</math></b>

**Table 1.** Mean Opinion Score (MOS) evaluations with 95% confidence intervals computed from the t-distribution for various systems.



**Fig. 2.** Synthesized vs. ground truth: 800 ratings on 100 items.



### Tacotron 2

- 声学模型与声码器对接时的不匹配问题

- 声学模型训练过程中，每个时刻解码器输入使用真实的梅尔谱
- 声学模型解码过程中，每个时刻使用解码器预测的梅尔谱
- 声码器WaveNet训练时，使用的是真实的梅尔谱
- 声码器合成时，使用的是声学模型预测出的梅尔谱

$$\hat{y}_t = \text{Decoder}(c, y_{t-1}, s_t)$$



## Tacotron 2

- GTA生成 (Ground-truth alignment)
  - 声学模型训练完成后, 使用该模型对训练数据做一次解码
  - 解码时每个时刻使用上一时刻真实的Mel谱, 保证生成谱和真实谱的长度一致
  - 缓解声学模型与声码器对接中的特征不匹配

Training	Synthesis	
	Predicted	Ground truth
Predicted	$4.526 \pm 0.066$	$4.449 \pm 0.060$
Ground truth	$4.362 \pm 0.066$	$4.522 \pm 0.055$



## 总结

- 序列到序列声学模型在发音自然度方面优于传统帧级声学模型
- 序列到序列声学模型能够大大简化语音合成系统的构建流程，但是其训练效率相对帧级模型较低，参数量也较大。
- 基于序列到序列声学模型能够便捷地应用到风格合成、说话人自适应等方面
- 就目前来看，序列到序列声学模型的一个关键问题是生成稳定性问题，其根源为注意力机制本身的稳定性问题



## 1. 概述



## 2. 序列到序列模型及注意力机制



## 3. 序列到序列声学模型Tacotron



## 4. 序列到序列声学模型变体（一）



## 5. 实战



## 实战：基于中文Tacotron系统

尝试按照README.md中的步骤，实现关键代码，从而构建中文Tacotron合成系统  
针对开源标贝中文数据集，已经给出准备的文本特征。首先，针对文本对应的语音数据提取语音特征，实现Tacotron系统中的CBHG模块，完成模型训练和测试。

Repo: [https://github.com/nwpuaslp/TTS\\_Course](https://github.com/nwpuaslp/TTS_Course)



感谢聆听 !  
Thanks for Listening

