

个性化语音合成项目

技术概览

2020.12.06

本周主要内容：

- 个性化语音合成技术整体架构。
- GE2E 的设计与实现。
- LSTM 结构设计及实现。
- BiLSTM的实现逻辑。

~~salep~~
~~movement~~
~~word flow~~

① 接收



② 输入 语音 模型
文本 → 语音

个性化 - 个方面

duration

emotion Embedding

gST

global

style

1 Tacotron 2

fast speed

DeepVoc

attention



24000

上采样

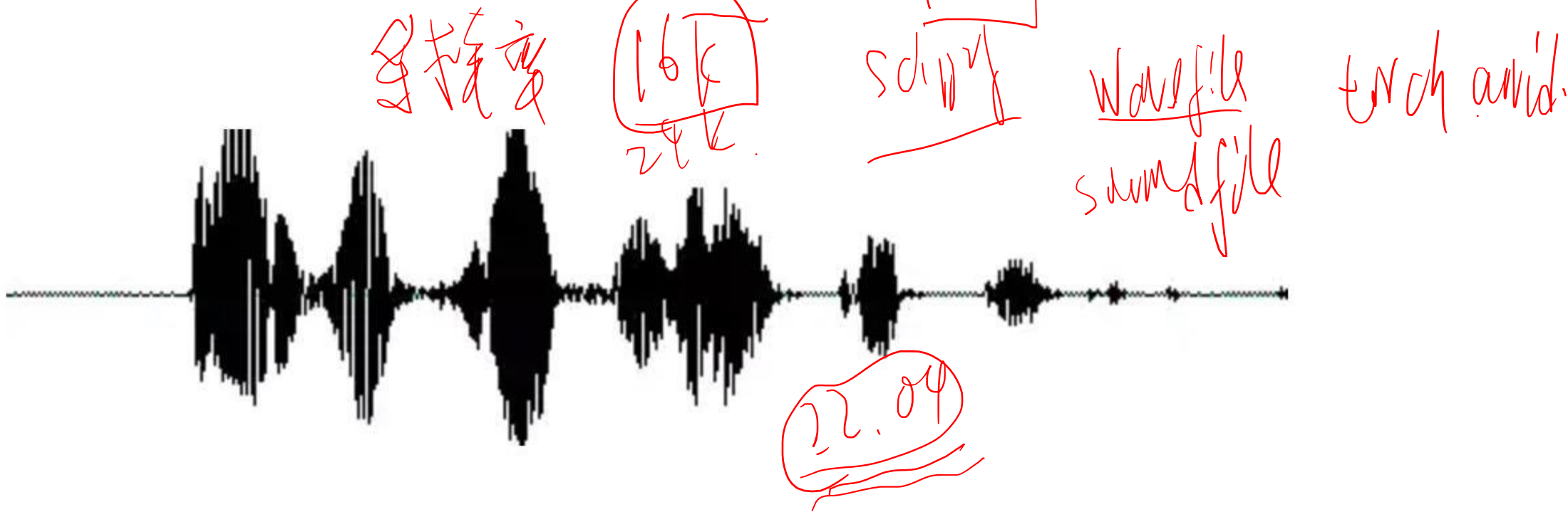
conv

3

± 1.8 20

基本概念-声音的表示

声音是一种波。常见的mp3、wmv等格式都是压缩格式，必须转成非压缩的纯波形文件来处理，比如Windows PCM文件，也就是俗称的wav文件。wav文件由存储的文件头和声音波形的采样点构成。



Web site - Valid)
 decoder ↑ 2d ELMs

之间有 $25-10=15$ 毫秒的交叠。我们称为以帧长25ms、帧移10ms分帧。

sel



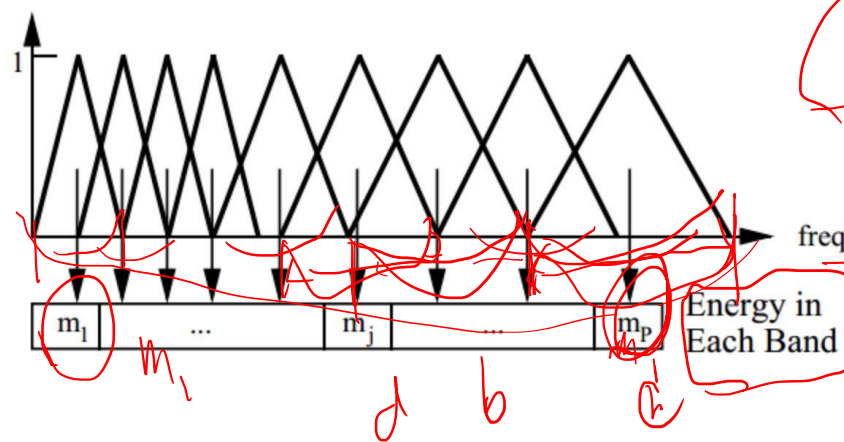
2/15 21:15

↓
每程中

周期性

基本概念

- 人耳对声音频谱的响应是非线性的。前端处理算法，以类似于人耳的方式对音频进行处理，可以提高语音识别的性能。FilterBank分析就是这样的一种算法。FBank特征提取要在预处理之后进行，这时语音已经分帧，我们需要逐帧提取FBank特征。



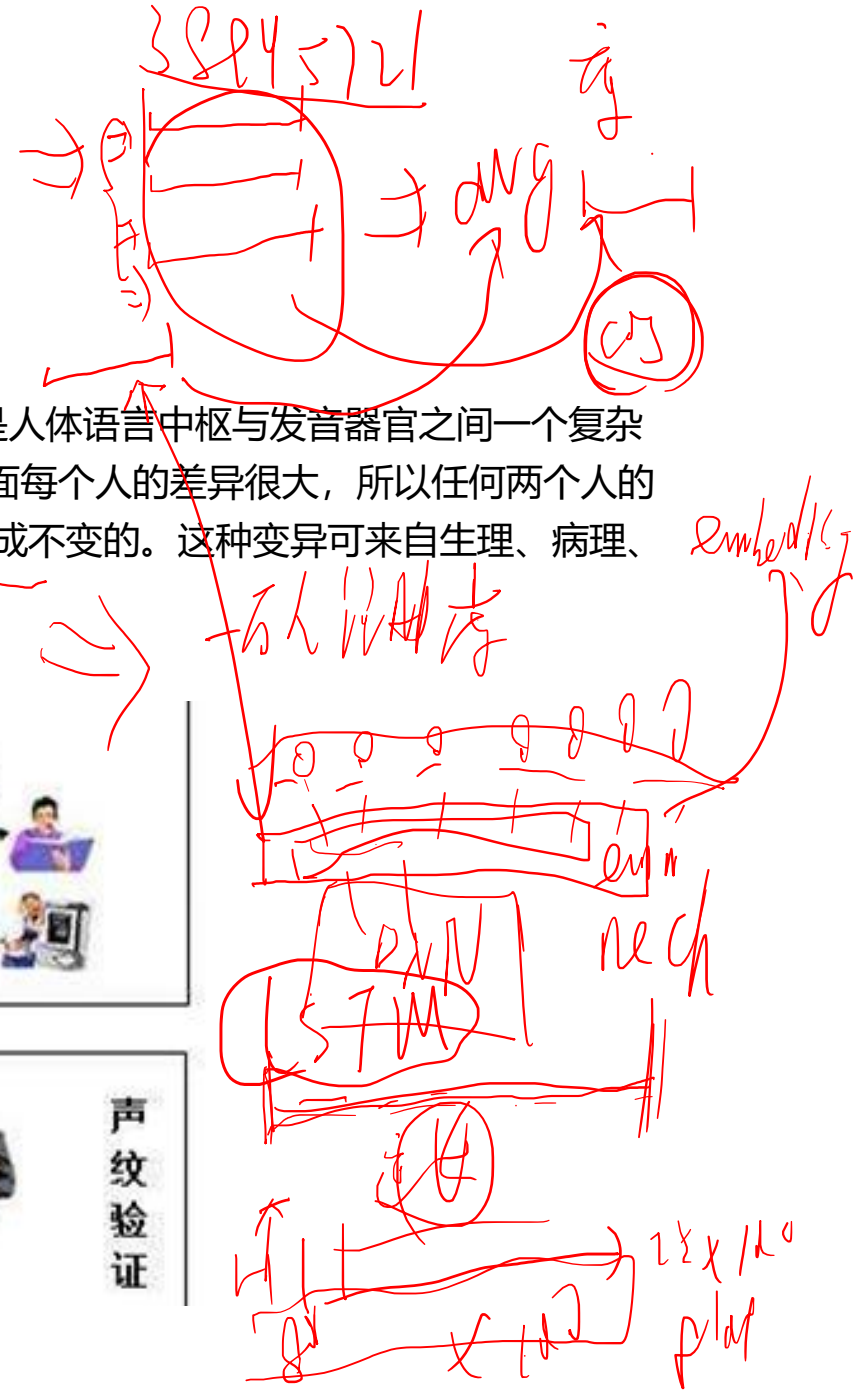
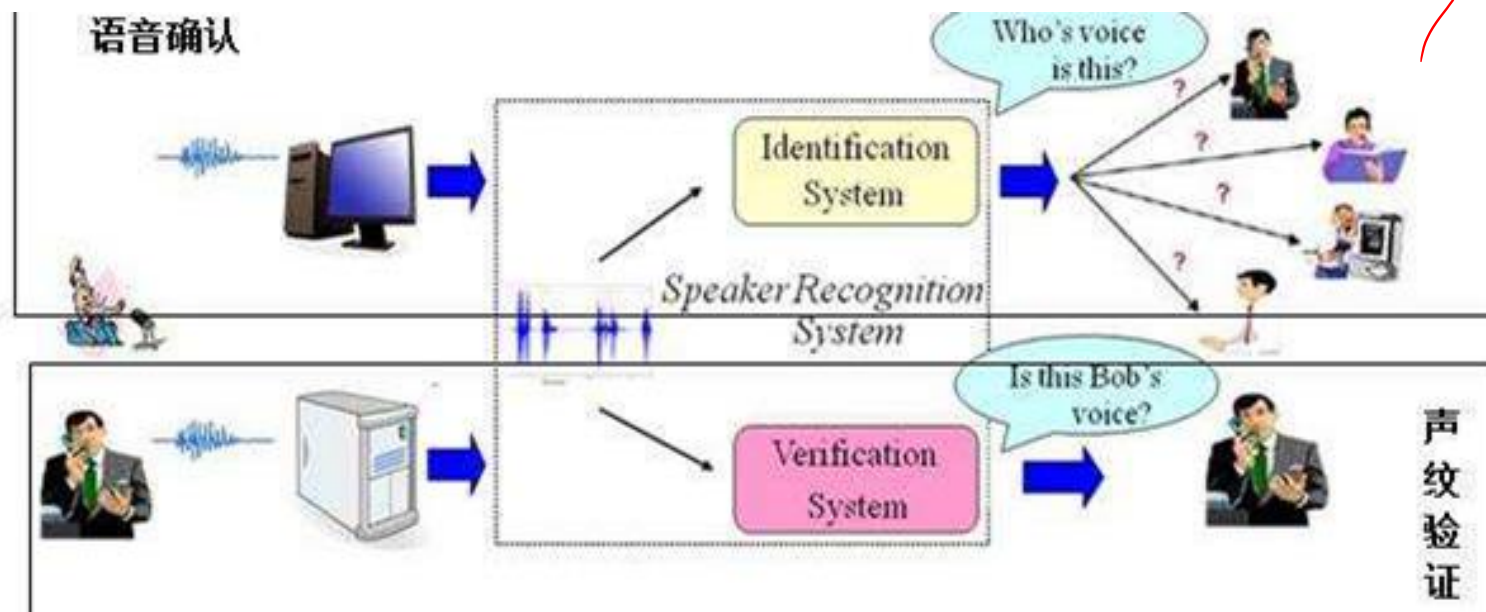
基本概念-声码器

- 声码器在发送端对语音信号进行分析，提取出语音信号的特征参量加以编码和加密，以取得和信道的匹配，经信息通道传递到接受端，再根据收到的特征参量恢复原始语音波形。分析可在频域中进行，对语音信号作频谱分析，鉴别清浊音，测定浊音基频，进而选取清-浊判断、浊音基频和频谱包络作为特征参量加以传送。分析也可在时域中进行，利用其周期性提取一些参数进行线性预测，或对语音信号作相关分析。
- 声码器可以分成：通道式声码器、共振峰声码器、图案声码器、线性预测声码器、相关声码器、正交函数声码器。

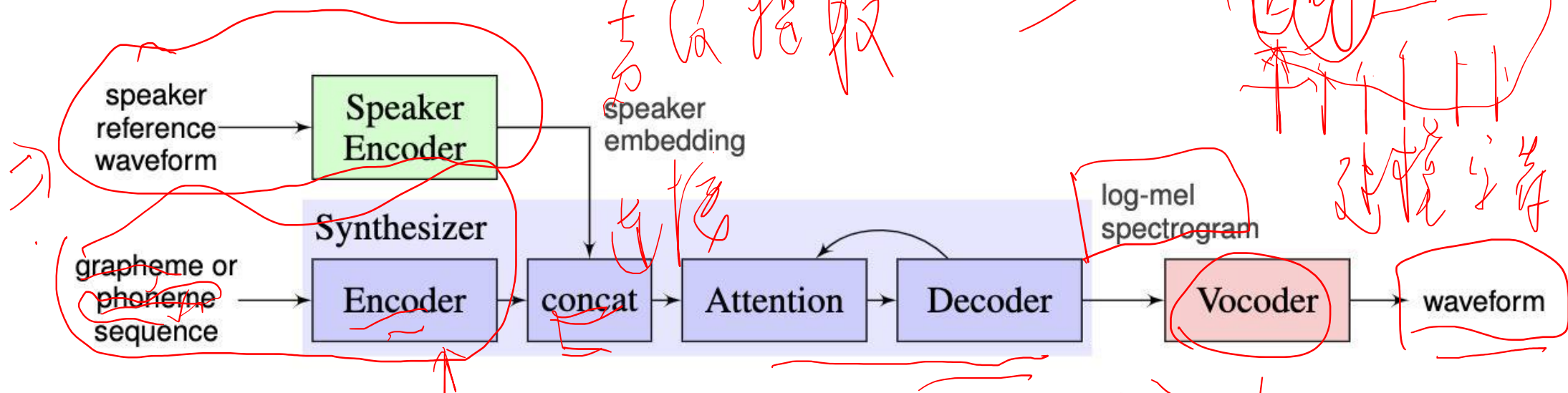


基本概念-声纹 (VoicePrint)

所谓声纹(Voiceprint), 是用电声学仪器显示的携带言语信息的声波频谱。人类语言的产生是人体语言中枢与发音器官之间一个复杂的生理物理过程, 人在讲话时使用的发声器官—舌、牙齿、喉头、肺、鼻腔在尺寸和形态方面每个人的差异很大, 所以任何两个人的声纹图谱都有差异。每个人的语音声学特征既有相对稳定性, 又有变异性, 不是绝对的、一成不变的。这种变异可来自生理、病理、心理、模拟、伪装, 也与环境干扰有关。



个性化语音合成整体架构



该系统由3个可独立训练模块构成：

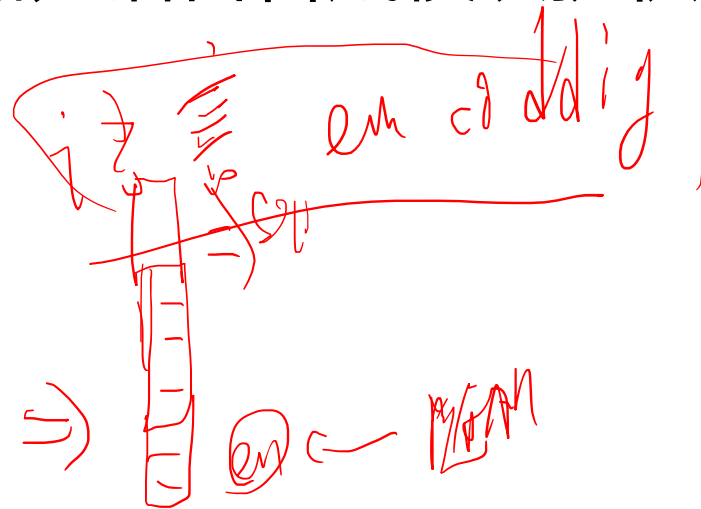
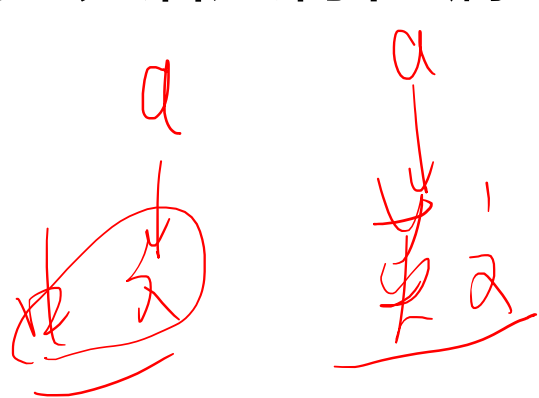
说话人特征编码器-speaker encoder：实现对目标说话人发音特征提取。

合成器-synthesizer：实现文本到梅尔谱(mel spectrogram) Seq2Seq变换。

声码器-vcoder: 实现梅尔谱 (mel spectrogram) 到波形(Waveform)变化。

Speaker Encoder

- 说话人特征编码器具有鉴别特定说话人的功能。输入为短时自适应语音信号，输出为说话人在发音特征向量空间中的特征向量。该特征编码与语音对应的文本及背景噪音独立，是对说话人发音特质的描述。说话人特征编码器为语音合成模块提供个性化定制参考。



Speaker Encoder

156 4/11

网络由3个LSTM层构成。每层LSTM包含768个cell，后接仿射层(Projection，通常目的是减少计算量和特征降维)变换到256维。最终的Embedding由该说话人最后一帧在神经网络运算结束后通过对最后一层在每个时间步的输出进行平均和L2 Norm生成。

```
def __init__(self, device, loss_device):
    super().__init__()
    self.loss_device = loss_device

    # Network definition
    self.lstm = nn.LSTM(input_size=mel_n_channels,
                        hidden_size=model_hidden_size,
                        num_layers=model_num_layers,
                        batch_first=True).to(device)

    self.linear = nn.Linear(in_features=model_hidden_size,
                            out_features=model_embedding_size).to(device)

    self.relu = torch.nn.ReLU().to(device)

    # Cosine similarity scaling (with fixed initial parameter values)
    self.similarity_weight = nn.Parameter(torch.tensor([10.])).to(loss_device)
    self.similarity_bias = nn.Parameter(torch.tensor([-5.])).to(loss_device)

    # Loss
    self.loss_fn = nn.CrossEntropyLoss().to(loss_device)
```

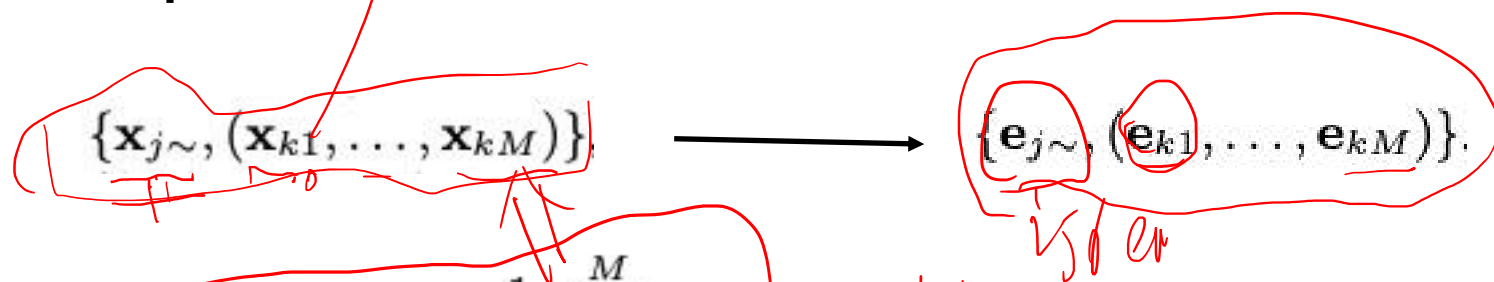
TE2E

$k \sim j, \dots, M$ 不同的 i, j

25/6
mini

~~FC~~
LSIM
LSIM
LSIM

- Tuple-based end-to-end model



$$\mathbf{c}_k = \mathbb{E}_m[\mathbf{e}_{km}] = \frac{1}{M} \sum_{m=1}^M \mathbf{e}_{km}$$

ipw

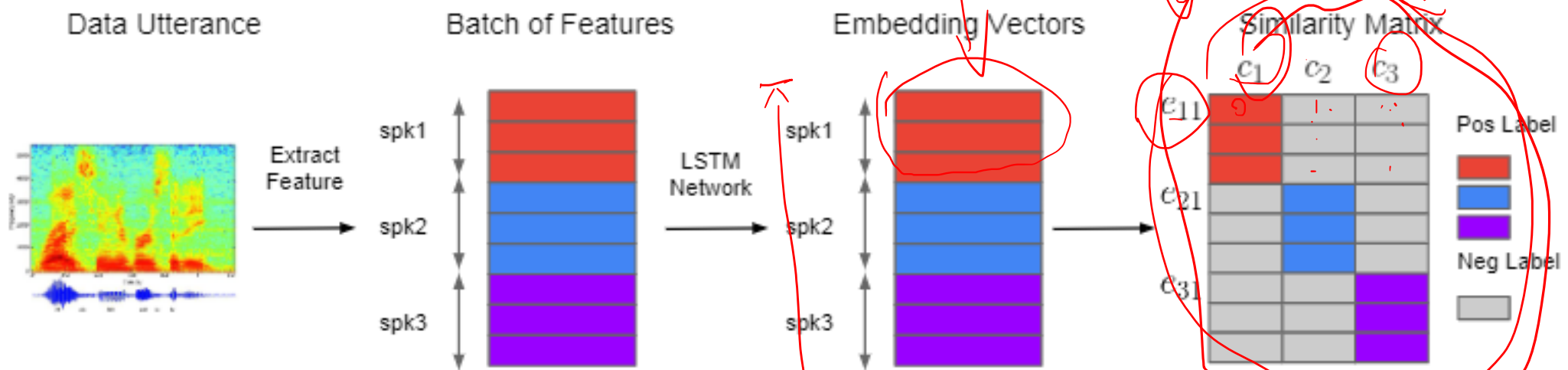
$$s = \underbrace{w}_{\text{weight}} \cdot \cos(\mathbf{e}_{j\sim}, \mathbf{c}_k) + b$$

$$L_T(\mathbf{e}_{j\sim}, \mathbf{c}_k) = \delta(j, k) (1 - \sigma(s)) + (1 - \delta(j, k)) \sigma(s).$$

$$\sigma(x) = 1/(1 + e^{-x})$$

GE2E

GE2E系统结构:



GE2E

emb'ding

$$\mathbf{e}_{ji} = \frac{f(\mathbf{x}_{ji}; \mathbf{w})}{\|f(\mathbf{x}_{ji}; \mathbf{w})\|_2}$$

CL-VI norm

$$S_{ji,k} = w \cdot \cos(\mathbf{e}_{ji}, \mathbf{c}_k) + b,$$

log max

$$L(\mathbf{e}_{ji}) = 1 - \sigma(S_{ji,j}) + \max_{\substack{1 \leq k \leq N \\ k \neq j}} \sigma(S_{ji,k}),$$

$$L(\mathbf{e}_{ji}) = -S_{ji,j} + \log \sum_{k=1}^N \exp(S_{ji,k}).$$

$$\mathbf{c}_j^{(-i)} = \frac{1}{M-1} \sum_{\substack{m=1 \\ m \neq i}}^M \mathbf{e}_{jm},$$

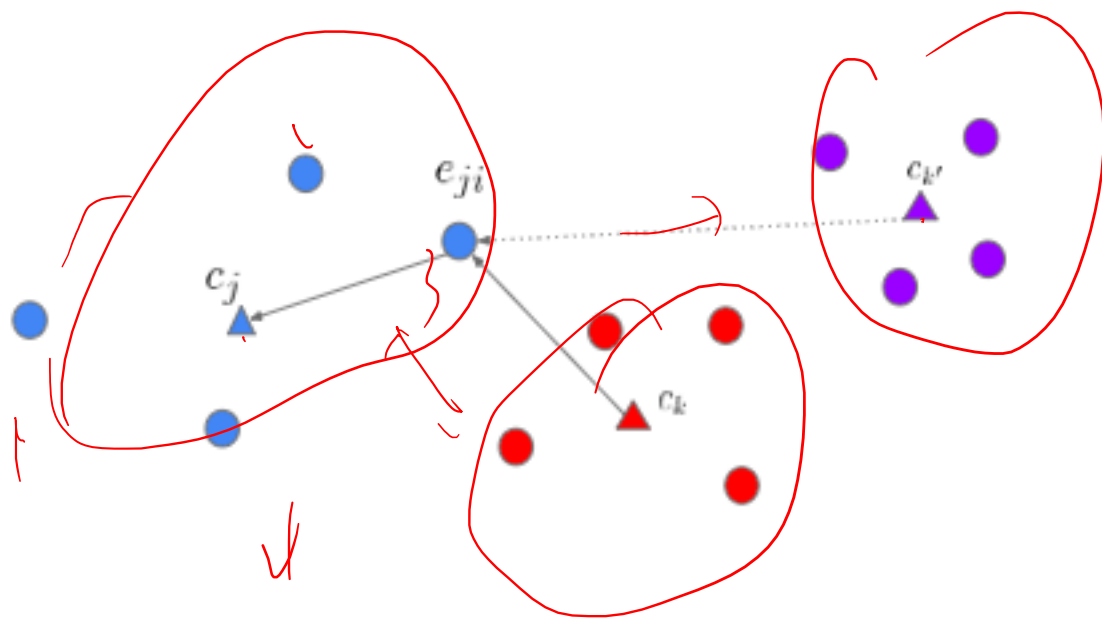
$$S_{ji,k} = \begin{cases} w \cdot \cos(\mathbf{e}_{ji}, \mathbf{c}_j^{(-i)}) + b & \text{if } k = j; \\ w \cdot \cos(\mathbf{e}_{ji}, \mathbf{c}_k) + b & \text{otherwise.} \end{cases}$$

相似

$$L_G(\mathbf{x}; \mathbf{w}) = L_G(\mathbf{S}) = \sum_{j,i} L(\mathbf{e}_{ji}).$$

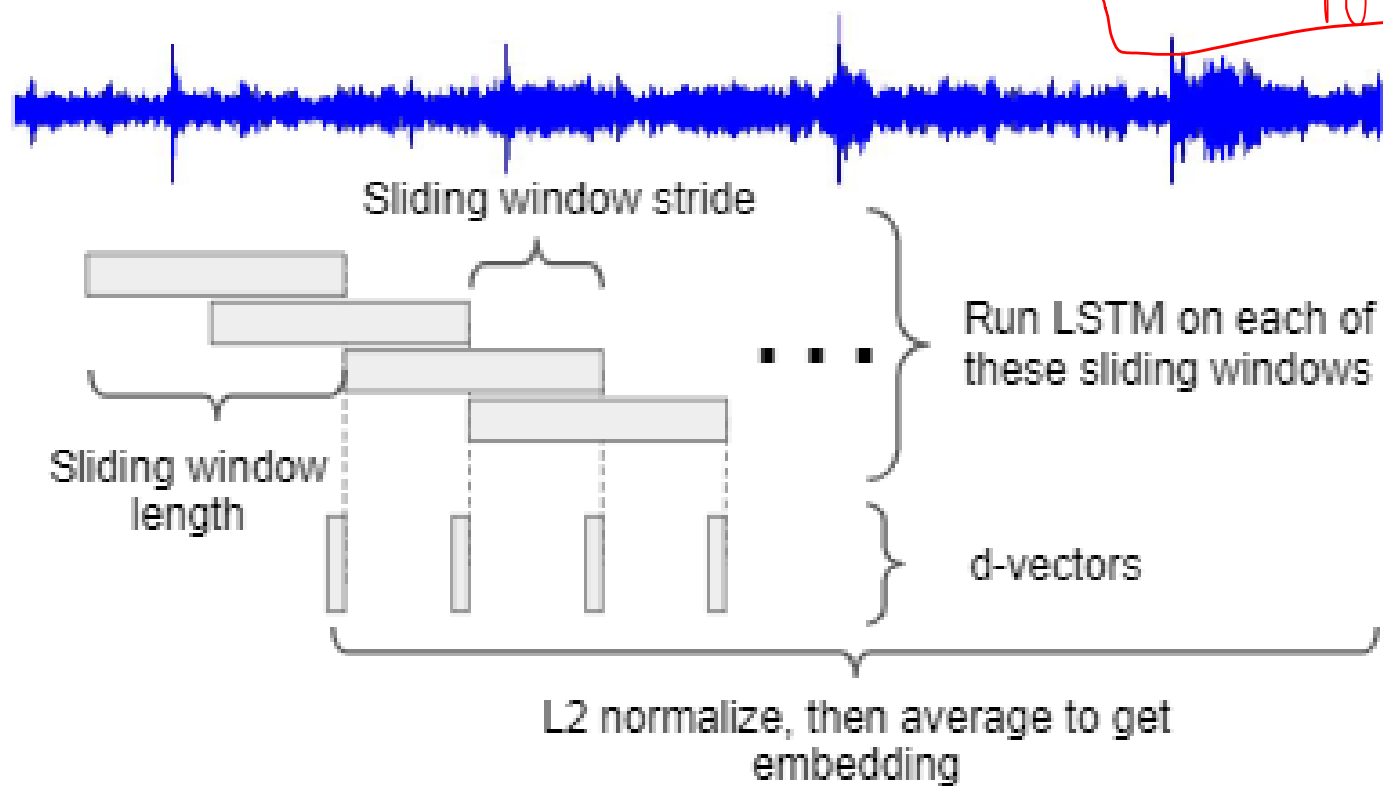
GE2E

GE2E优化过程:



GE2E

GE2E 滑动窗口机制



句
数据打乱
增强

合成器(Synthesizer)

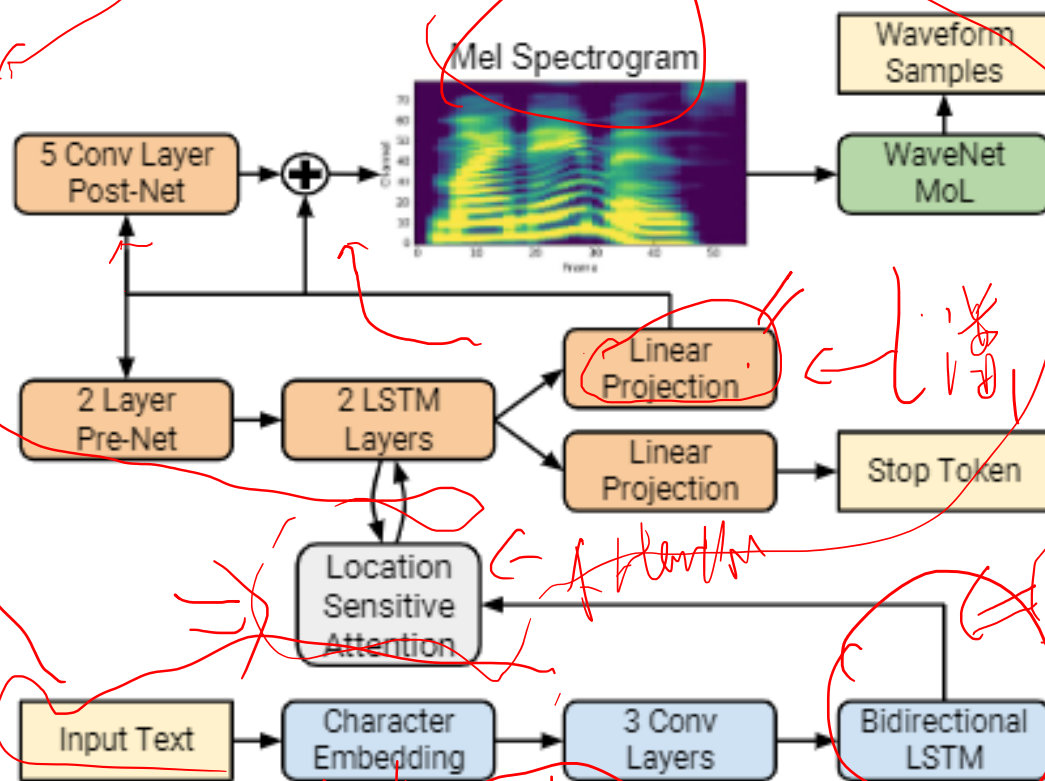
nvidia / tensorflow

合成器

decoder

1/2/3

= MFC / F_hall



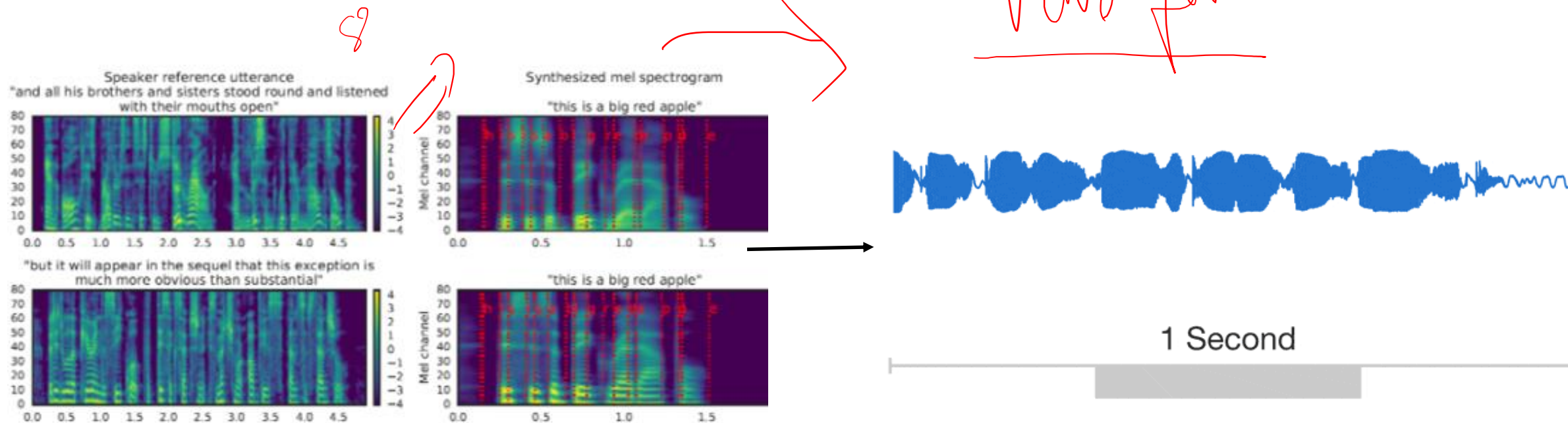
1/2/3 - MFC
1/2/3 - F_hall

en -> ca

encoder

声码器

频谱到波形的变换过程。



声码器-传统算法

- Griffin-lim算法

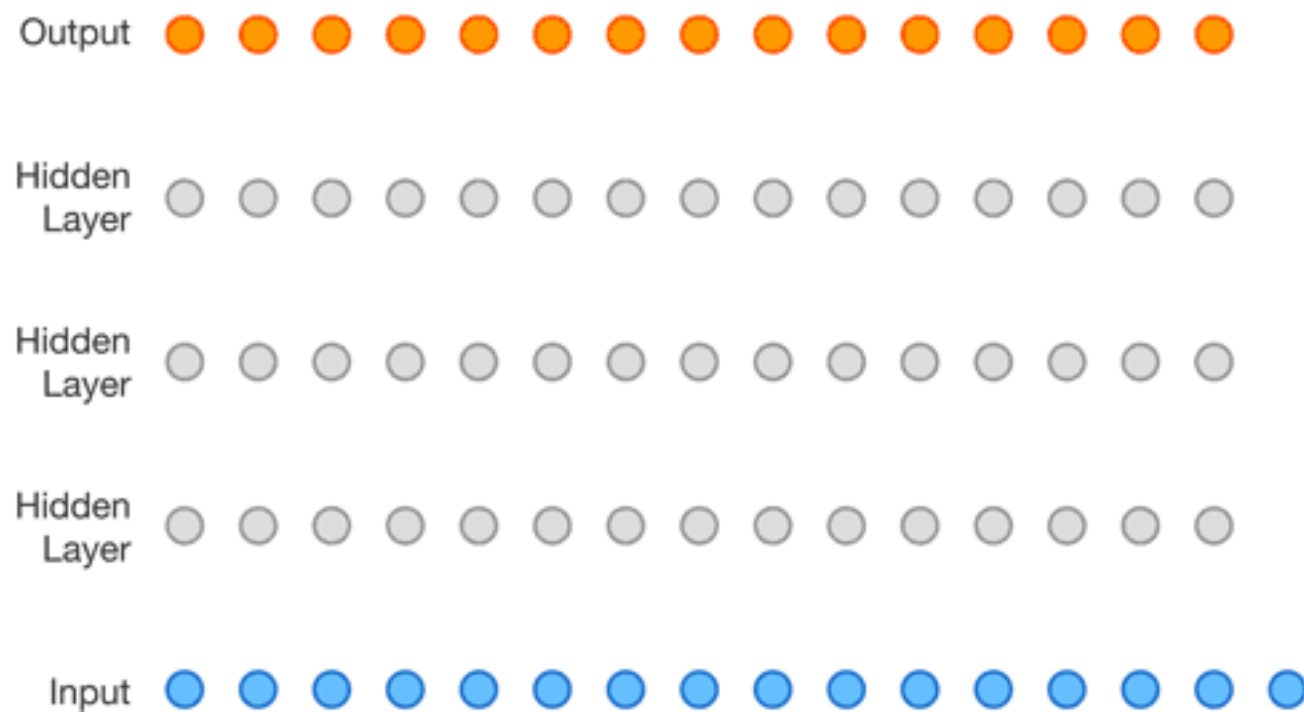
griffin-lim重建语音信号需要使用到幅度谱和相位谱。而MEL谱当中是不含相位信息的，因此griffin-lim在重建语音博形的时候只有MEL谱可以利用，利用帧与帧之间的关系估计出相位信息，重建语音波形。

这里的MEL谱可以看做是实部，而相位信息可以看做是虚部，通过对实部和虚部的运算，得到最终的结果。

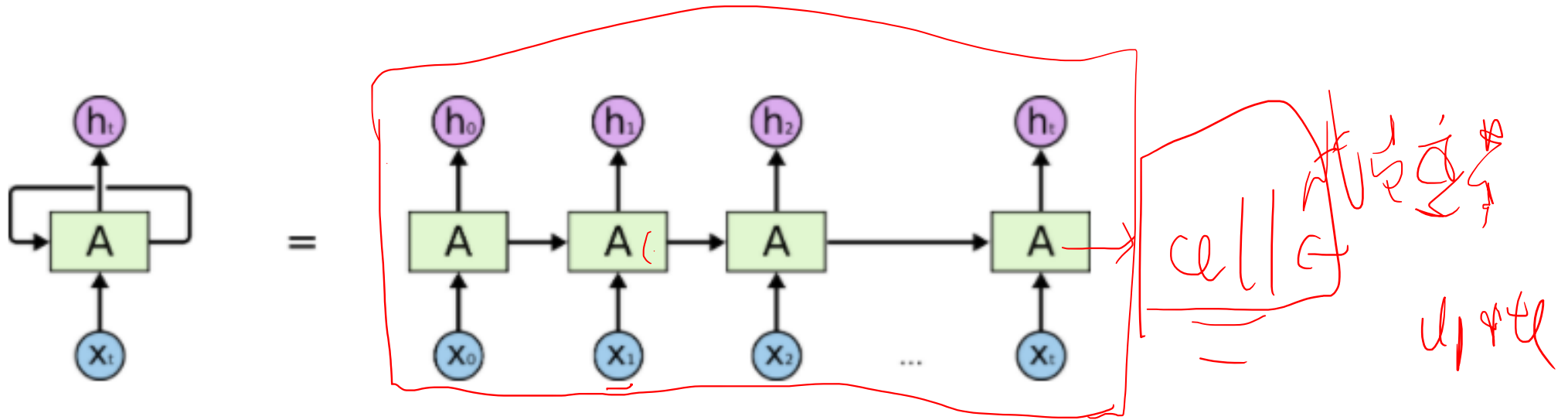
声码器-griffin-lim算法流程

- 随机初始化一个相位谱
- 用这个相位谱与已知的幅度谱（来自MEL谱）经过ISTFT（逆傅里叶变换）合成新的语音波形
- 用合成语音做STFT，得到新的幅度谱和新的相位谱
- 丢弃新的幅度谱，用已知幅度谱与新的相位谱合成新的语音
- 重复2,3,4多次，直至合成的语音达到满意的效果或者迭代次数达到设定的上限

声码器-wavenet

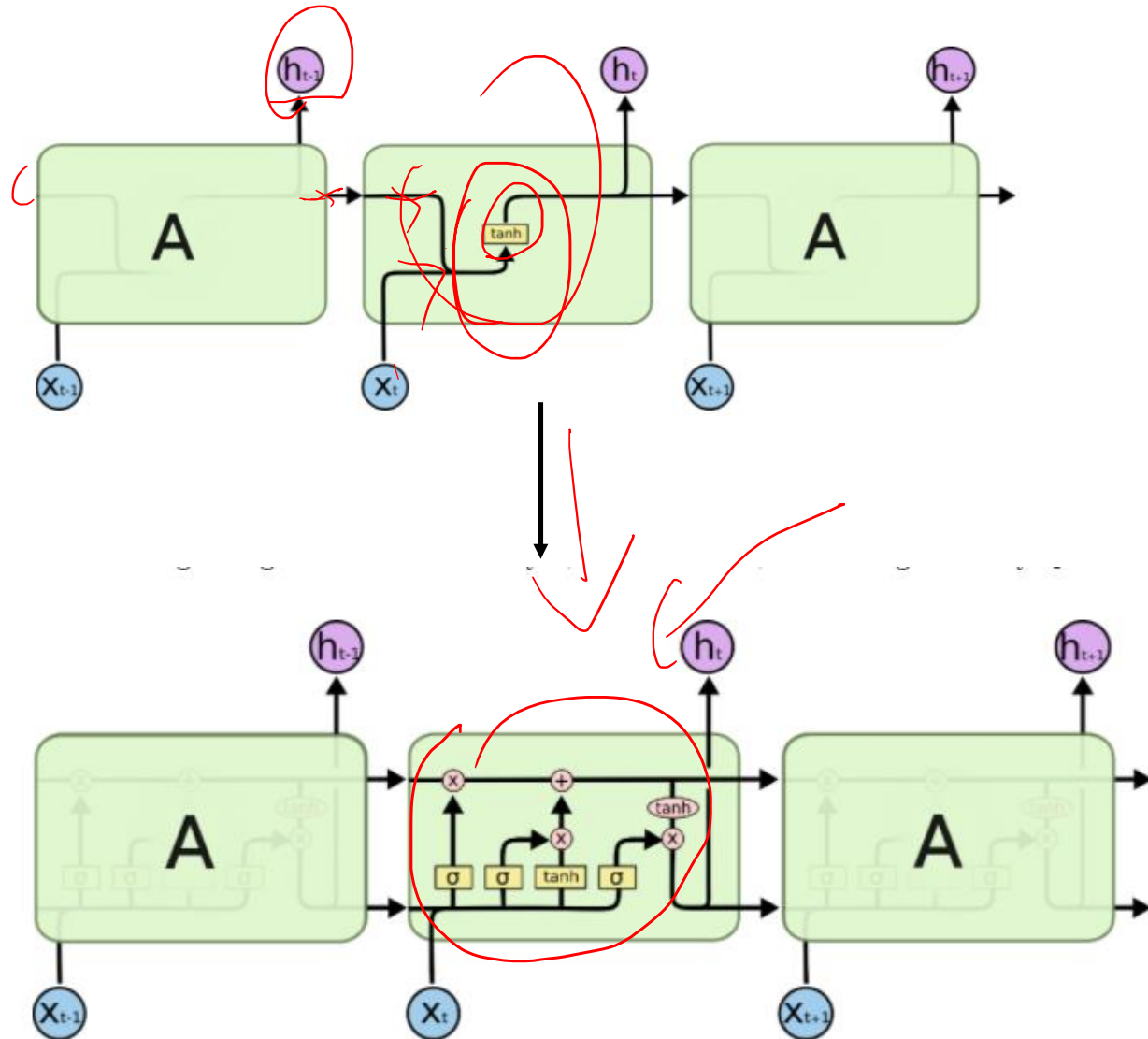


Recurrent Neural Networks

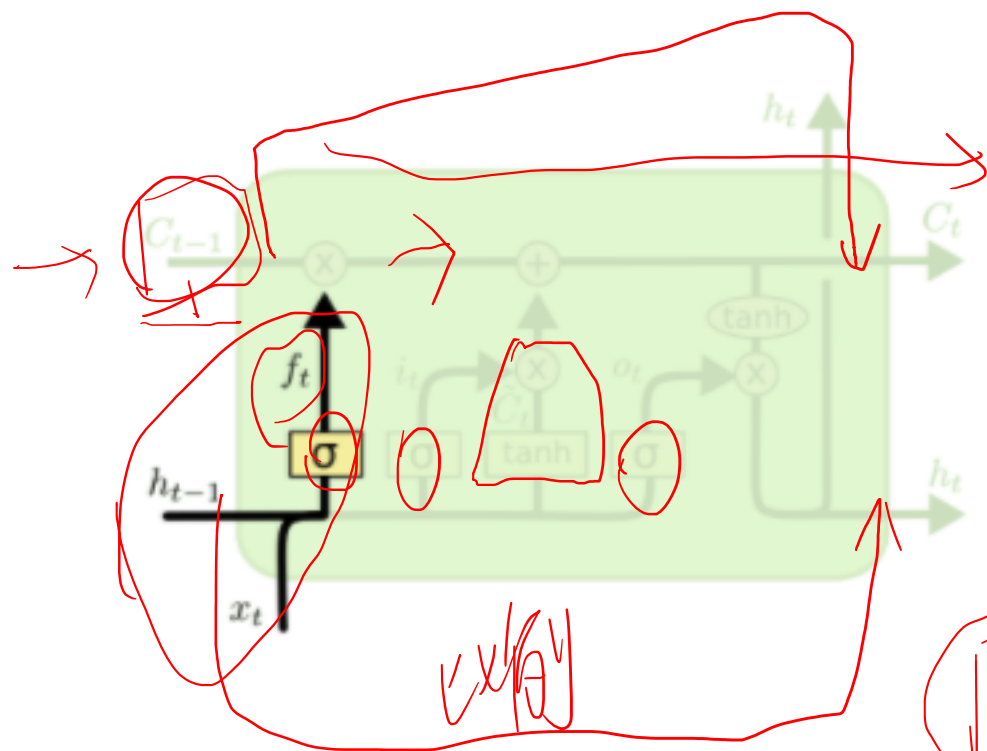


An unrolled recurrent neural network.

RNN->LSTM



LSTM遗忘门



USM GRV

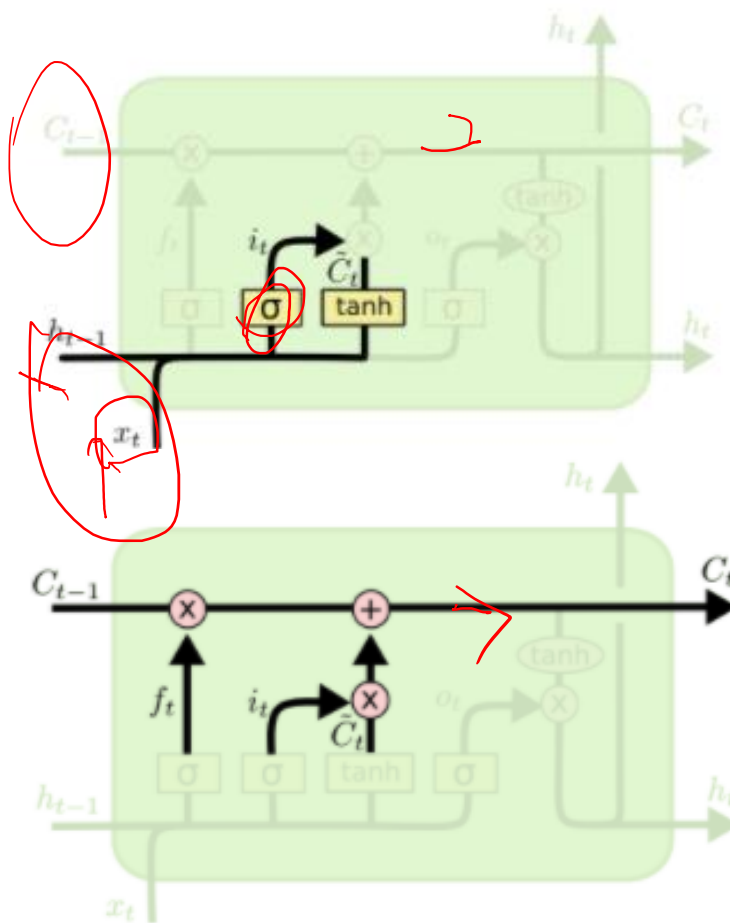
ReLU Net

$$H(y) = (F(x) + x)$$

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

projection layer
bottleneck
h
s/l

LSTM记忆门



$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

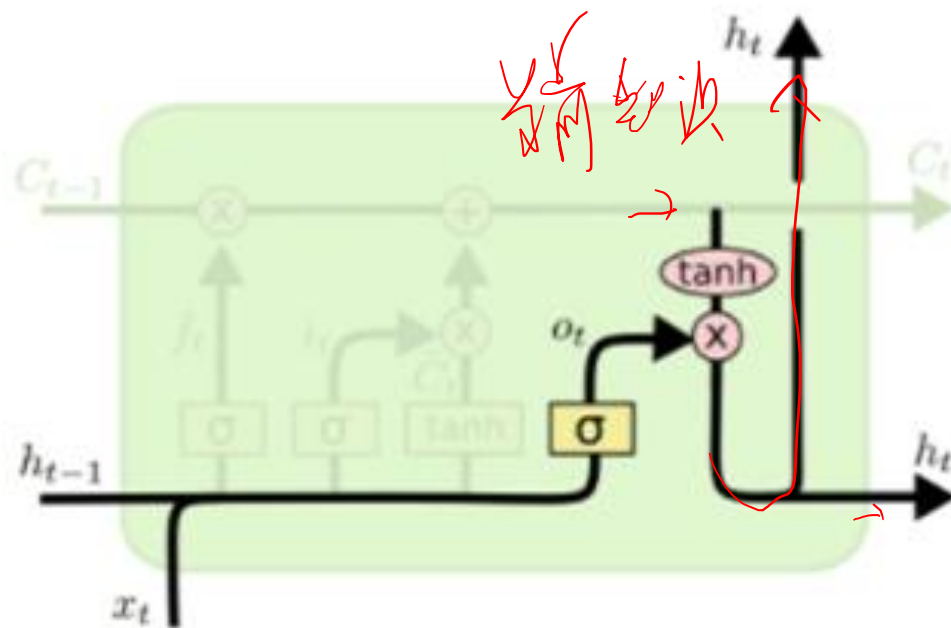
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

记忆门

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

遗忘
前
新
增
知
识

LSTM输出门



$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

LSTM实现

```
244 class Gate {
245 public:
246     int _n_sub_seq_size, _n_batch_size;
247     int _n_tbptt;
248     size_t _n_feat_dim, _n_cell_dim, _n_rec_dim;
249
250     FMatrix *_bias, *_wc;
251     Weight *_wx, *_wr;
252
253     StateMatrix _stat_oe;
254
255     Gate(int tbptt, size_t feat_dim, size_t cell_dim, size_t outDim);
256     ~Gate();
257     void set_batch_size(int sub_seq_size, int batch_size);
258     void input_forward(InOutput &in_out, IN_OUT_TYPE_T i_type, IN_OUT_TYPE_T o_type);
259     void time_forward(int t, FMatrix &oc, FMatrix &or_);
260 };
```

```
262 class OutGate : public Gate {
263 public:
264     OutGate(int tbptt, size_t feat_dim, size_t cell_dim, size_t rec_dim)
265         : Gate(tbptt, feat_dim, cell_dim, rec_dim) {}
266     void time_forward(int t, FMatrix &Oc, FMatrix &Or);
267 };
```

Gate的定义。不同的Gate实现具体的time forward算法。

LSTM实现

```
25 LstmLayer::LstmLayer(LstmConfig &conf) : Layer(conf) {
26     init();
27     _weights = static_cast<LstmWeights *>(conf._weights);
28     size_t feat_dim = conf.in_dim();
29     size_t cell_dim = conf.cell_dim();
30     size_t out_dim = conf.out_dim();
31     size_t prj_dim = conf.prj_dim();
32     size_t rec_dim = conf.rec_dim();
33
34     _n_tbptt = 1; // Tbptt;
35     _n_feat_dim = feat_dim;
36     _n_cell_dim = cell_dim;
37     _n_rec_dim = rec_dim;
38     _n_prj_dim = prj_dim;
39     _n_out_dim = out_dim;
40
41     // gate
42     int rec = (int)_n_rec_dim;
43     if (rec == 0) {
44         rec = (int)_n_out_dim;
45     }
46     _output_g = new OutGate(1, feat_dim, cell_dim, rec);
47     _input_g = new Gate(1, feat_dim, cell_dim, rec);
48     _forget_g = new Gate(1, feat_dim, cell_dim, rec);
49     _cells = new Cells(1, feat_dim, cell_dim, rec);
50
51     set_weights(static_cast<LstmWeights *>(conf._weights));
52     _cells->_min = _weights->_cec_out_limit_lo;
53     _cells->_max = _weights->_cec_out_limit_hi;
54
55     _rec_act = Activation::create(conf.rec_act());
56
57     set_batch_size(conf.sub_seq_size(), conf.batch_size());
58 }
```

prj feat

```
for (int t = 0; t < (int)sub_seq_size; t++) {
    size_t start_row = _n_tbptt + t - 1;
    size_t end_row = _n_tbptt + t;

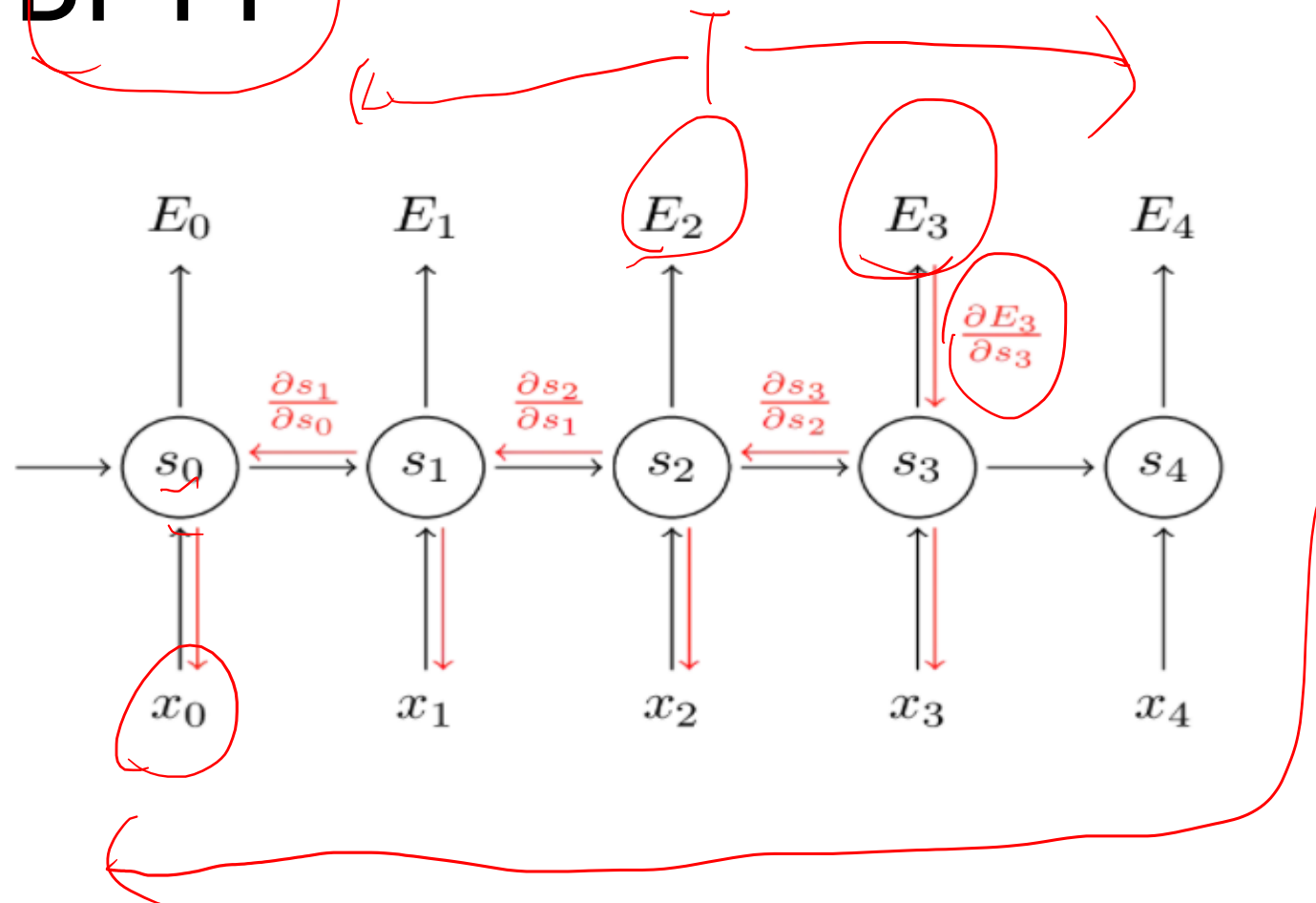
    FMatrix &r_pre_o =
        _statoe_r._o->range_row(start_row, end_row, _n_batch_size);
    _in_out.set_input(&r_pre_o);
    _in_out.set_output(&wx_out.range_row(t, t + 1, _n_batch_size));
    _weights->_wr_iofc.mul(_in_out, INOUT_TYPE, INOUT_TYPE, NULL, 1.0f, 1.0f);
    _in_out.clear_input(INOUT_TYPE);
    _in_out.clear_output(INOUT_TYPE);

    // input gate Wic * Ct-1
    FMatrix &wxwr_out = wx_out.range_row(t, t + 1, _n_batch_size);
    _i_out.copy_from(wxwr_out.range_col(0, 1 * _n_cell_dim, 1));
    _i_out.mul_diag_mat_sigmoid(
        _statoe_c._o->range_row(start_row, end_row, _n_batch_size),
        *_weights->_i_wc._f_w, _i_out, 1.0f, 1.0f);

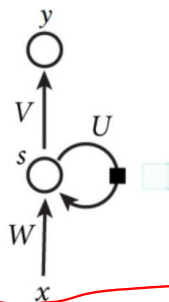
    // forget gate Wfc * Ct-1
    _f_out.copy_from(wxwr_out.range_col(1 * _n_cell_dim, 2 * _n_cell_dim, 1));
    _f_out.mul_diag_mat_sigmoid(
        _statoe_c._o->range_row(start_row, end_row, _n_batch_size),
        *_weights->_f_wc._f_w, _f_out, 1.0f, 1.0f);

    // cell
    FMatrix &c_cur_out = _statoe_c._cur_o->range_row(t, t + 1, _n_batch_size);
    FMatrix &c_pre_out = _statoe_c._o->range_row(t, t + 1, _n_batch_size);
    _g_out.copy_from(wxwr_out.range_col(3 * _n_cell_dim, 4 * _n_cell_dim, 1));
    // g()
    _g_out.tanh();
    // Ft . Ct-1 + It . g()
    c_cur_out.elem_mul_add(_g_out, _i_out, _f_out, c_pre_out);
}
```

BPTT



BPTT



$$\begin{cases} s_t = U h_{t-1} + W x_t \\ h_t = \tanh(s_t) \\ z_t = V h_t \\ \hat{y}_t = \text{softmax}(z_t) \\ E_t = -y_t^T \log(\hat{y}_t) \\ E = \sum_t E_t \end{cases}$$

$$\begin{aligned} \frac{\partial E_t}{\partial V_{ij}} &= \text{tr} \left(\left(\frac{\partial E_t}{\partial z_t} \right)^T \cdot \frac{\partial z_t}{\partial V_{ij}} \right) \\ &= \text{tr} \left((\hat{y}_t - y_t)^T \cdot \begin{bmatrix} 0 \\ \vdots \\ \frac{\partial z_t^{(i)}}{\partial V_{ij}} \\ \vdots \\ 0 \end{bmatrix} \right) \\ &= r_t^{(i)} h_t^{(j)} \end{aligned}$$

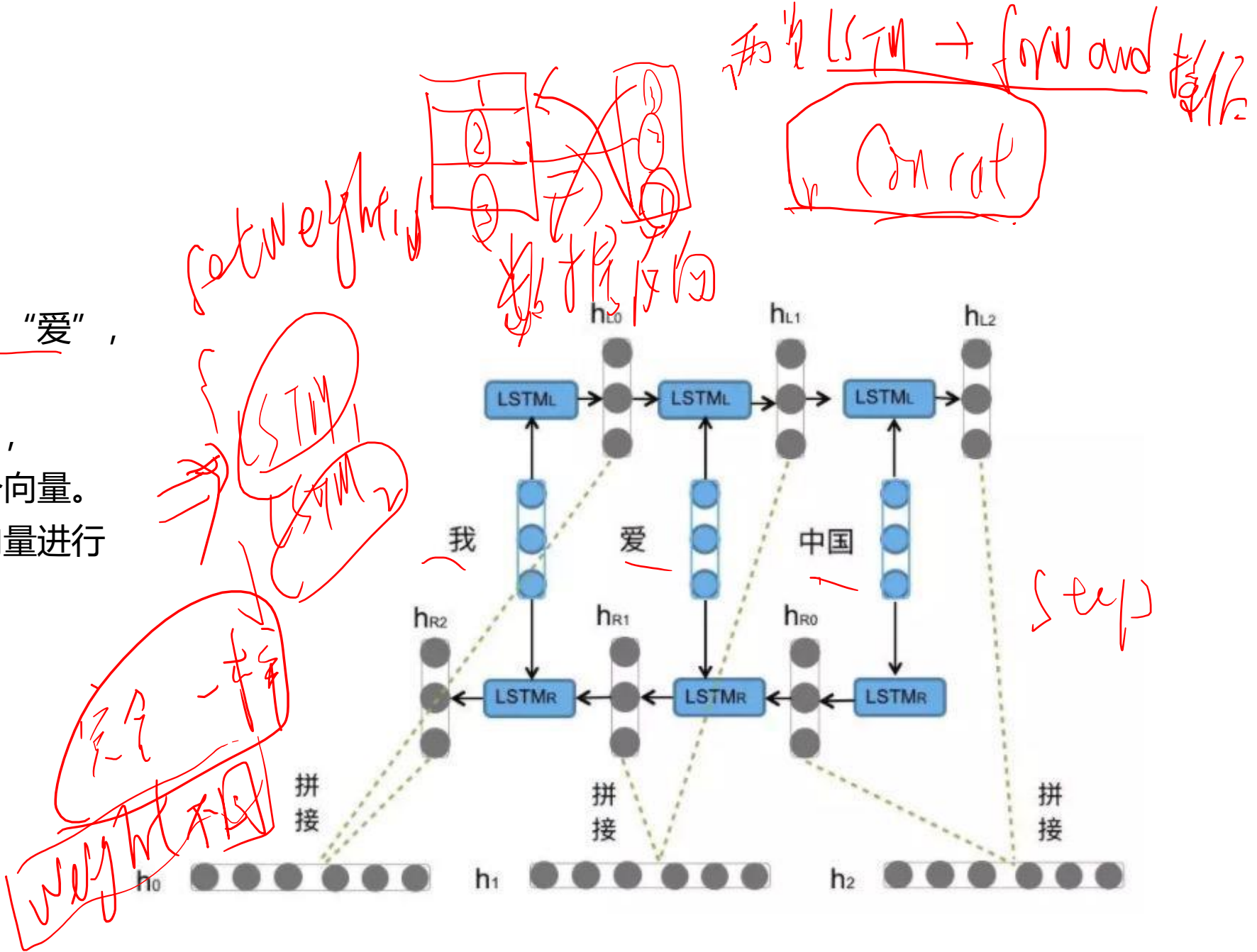
$$\begin{aligned} \frac{\partial E_t}{\partial U} &= \sum_{k=0}^t \delta_k \otimes h_{k-1} \\ \frac{\partial E}{\partial U} &= \sum_{t=0}^T \sum_{k=0}^t \delta_k \otimes h_{k-1} \end{aligned}$$

$$\frac{\partial E_t}{\partial W} = \sum_{k=0}^t \delta_k \otimes x_k$$

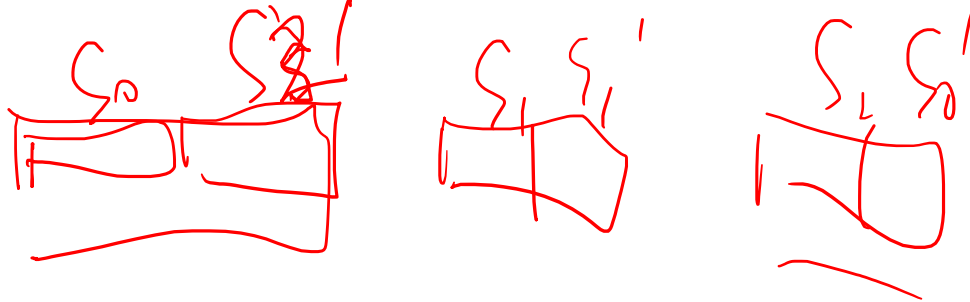
$$\begin{aligned} V &:= V - \lambda \sum_{t=0}^T (\hat{y}_t - y_t) \otimes h_t \\ U &:= U - \lambda \sum_{t=0}^T \sum_{k=0}^t \delta_k \otimes h_{k-1} \\ W &:= W - \lambda \sum_{t=0}^T \sum_{k=0}^t \delta_k \otimes x_k \end{aligned}$$

BiLSTM

- 前向的依次输入“我”，“爱”，“中国”得到三个向量。
- 后向的依次输入“中国”，“爱”，“我”得到三个向量。
- 最后将前向和后向的隐向量进行拼接得到。



Bilstm的实现



```
class FastBilstmLayer : public Layer {
public:
    FastBilstmLayer(FastBilstmConfig& config);
    ~FastBilstmLayer();
    void forward(InOutput &in_out, IN_OUT_TYPE_T i_type);
    void set_batch_size(int frame_size);
    void reset(int sent_idx);
    void store_current_out();
    inline int is_append() {
        return _is_append;
    }

protected:
    InOutput _invert_in_out;
    InOutput _in_out; // 正向输入
    FMatrix _invert_in_mat, _invert_out_mat;
    // 记录forward和backward LSTM的输入数据的映射顺序
    CpuIVector _map_id_vec1, _map_id_vec2; // used for encode and decode mapId
    int _map_length;
    int _is_append;

    // =====
    FastLstmLayer *_lstm_layer;
    FastLstmWeights *_fwd_weight;
    Container<FMatrix*> _fwd_his_o;

    FastLstmWeights *_bwd_weight;
    Container<FMatrix*> _bwd_his_o;
};
```

```
BilstmLayer::BilstmLayer(BilstmConfig &conf) : Layer(conf) {
    // 结构一致
    _lstm_layer = new LstmLayer(*conf.normal_config());
    // 前向weight和反向weight
    _fwd_weight = static_cast<LstmWeights*>(conf.normal_config()->_weights);
    _bwd_weight = static_cast<LstmWeights*>(conf.inver_config()->_weights);
    // 是否采用append模式
    _is_append = conf.is_append();
    _n_batch_size = 0;
    _fwd_his_o.clear();
    _bwd_his_o.clear();
}
```



Bilstm的forward

```
1416 /*
1417  *双向lstm
1418  */
1419 void BiLstmLayer::forward(InOutput &in_out, IN_OUT_TYPE_T i_type) {
1420     int height = (int)in_out.get_height();
1421     int width = (int)in_out.get_in_width();
1422     in_out.trans_in(i_type, INOUT_TYPE);
1423     FMatrix *in = in_out.f_in();
1424
1425     // normal-forward
1426     _in_out.set_input(in);
1427     _lstm_layer->set_weights(_fwd_weight);
1428     _lstm_layer->set_history(_fwd_his_o); //+++++
1429     _lstm_layer->forward(_in_out, i_type);
1430     _lstm_layer->store_current_out();
1431     _lstm_layer->get_history(_fwd_his_o); //+++++
1432
1433     _in_out.clear_input(INOUT_TYPE);
1434 }
```

正向执行

Bilstm的forward

```
1435 // inverse-forward
1436 _invert_in_out.resize_in(height, width, INOUT_TYPE);
1437 FMatrix *inverIn = _invert_in_out.f_in();
1438
1439 //把输入顺序进行反转
1440 for (int ii = 0; ii < height; ii++) {
1441     int pos = _map_id_vec1.get_value(ii);
1442     inverIn->range_row(pos, pos + 1).copy_from(in->range_row(ii, ii + 1, 1));
1443 }
1444 _lstm_layer->set_weights(_bwd_weight);
1445 _lstm_layer->set_history(_bwd_his_o);
1446 _lstm_layer->forward(_invert_in_out, INOUT_TYPE);
1447 _lstm_layer->store_current_out();
1448 _lstm_layer->get_history(_bwd_his_o); //+++++++
1449
1450 // translate the tow outputs
1451 _in_out.trans_out(_lstm_layer->o_type(), INOUT_TYPE);
1452
1453 _invert_in_out.trans_out(_lstm_layer->o_type(), INOUT_TYPE);
1454 FMatrix *inter_out = _in_out.f_out();
1455 FMatrix *inver_out = _invert_in_out.f_out();
1456 in_out.trans_out(_lstm_layer->o_type(), INOUT_TYPE);
1457 FMatrix *out = in_out.f_out();
1458
```

输入反转，反向lstm执行

Bilstm执行

```
1459 // according to the require to combine the output
1460 if (!_is_append) {
1461     out->resize(inter_out->get_height(), inter_out->get_width());
1462     out->copy_from(*inter_out);
1463     for (int ii = 0; ii < height; ii++) {
1464         int pos = _map_id_vec2.get_value(ii);
1465         out->range_row(pos, pos + 1).add(inver_out->range_row(ii, ii + 1));
1466     }
1467 } else {
1468     CHECK(inter_out->get_width() == inver_out->get_width(), "Not Matched");
1469     width = (int)inter_out->get_width();
1470
1471     out->resize(height, 2 * width);
1472     for (int ii = 0; ii < height; ii++) {
1473         int pos = _map_id_vec2.get_value(ii);
1474         out->range_row(ii, ii + 1)
1475             .range_col(0, width)
1476             .copy_from(inter_out->range_row(ii, ii + 1));
1477         out->range_row(pos, pos + 1)
1478             .range_col(width, 2 * width)
1479             .copy_from(inver_out->range_row(ii, ii + 1));
1480     }
1481 }
1482 // layer activation
1483 _act->forward(*in_out.f_out(), *in_out.f_out());
1484 in_out.trans_out(INOUT_TYPE, o_type);
1485 }
```

结果拼合

✓✓✓

本周作业

- 个性化合成声纹模型训练和ONNX导出。

参考：

项目地址：<https://github.com/CorentinJ/Real-Time-Voice-Cloning.git>

encoder_preprocess.py: 执行wav文件预处理成系统指定格式。

encoder_train.py: 执行声纹提取的训练

模型导出：<https://pytorch.org/docs/stable/onnx.html>

签到 & 反馈

