



task report



BUPT

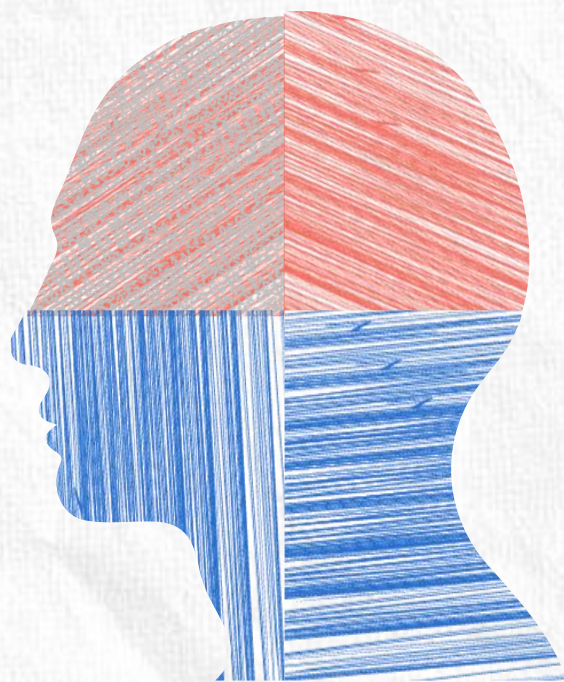
Zhang Yang

个人介绍

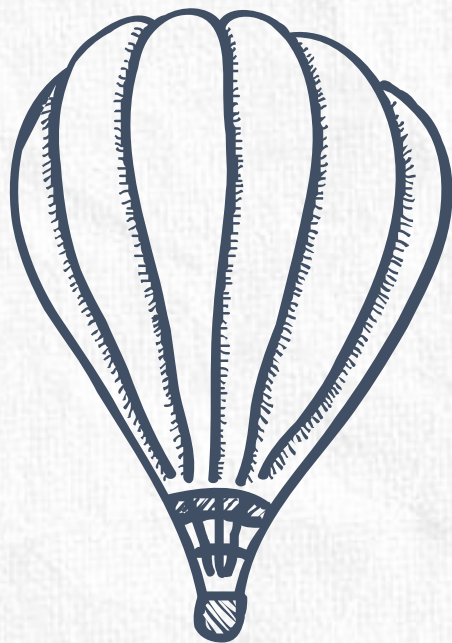
BYR

北邮信通院通信专业大三学生

张阳



CONTENT



1

ML and DL

2

Speaker Recognition and kaldi

3

test

1

Machine Learning and Deep Learning

了解了机器学习和深度学习的基本概念，学习了深度学习的基本模型和训练方法

Relationship between ML and DL

深度学习与机器学习的关系

人工智能是一个大概念，机器学习则是人工智能领域的一个小分支，如果说AI是一个合集，那么ML就是AI的子集。

任何通过数据训练的学习算法的相关研究都属于机器学习，包括很多已经发展多年的技术，比如线性回归（**Linear Regression**）、K均值（**K-means**，基于原型的目标函数聚类方法）、决策树（**Decision Trees**，运用概率分析的一种图解法）、随机森林（**Random Forest**，运用概率分析的一种图解法）、PCA（**Principal Component Analysis**，主成分分析）、SVM（**Support Vector Machine**，支持向量机）以及ANN（**Artificial Neural Networks**，人工神经网络）。

而人工神经网络则是深度学习的起源。

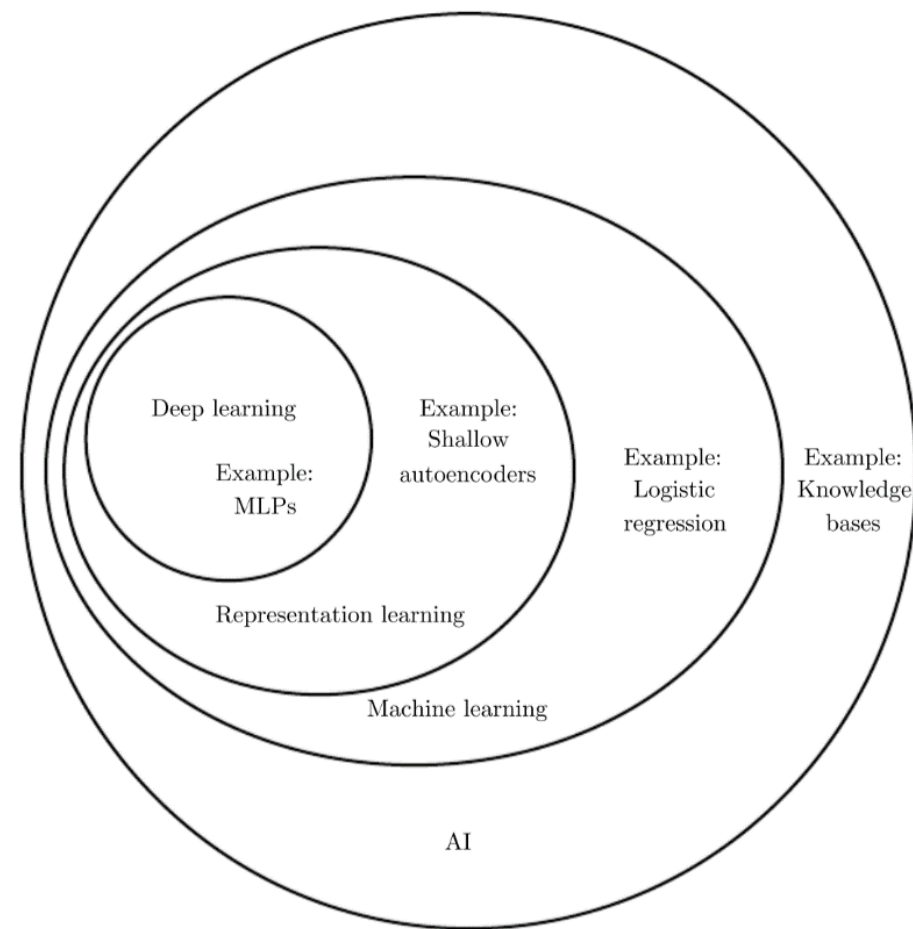
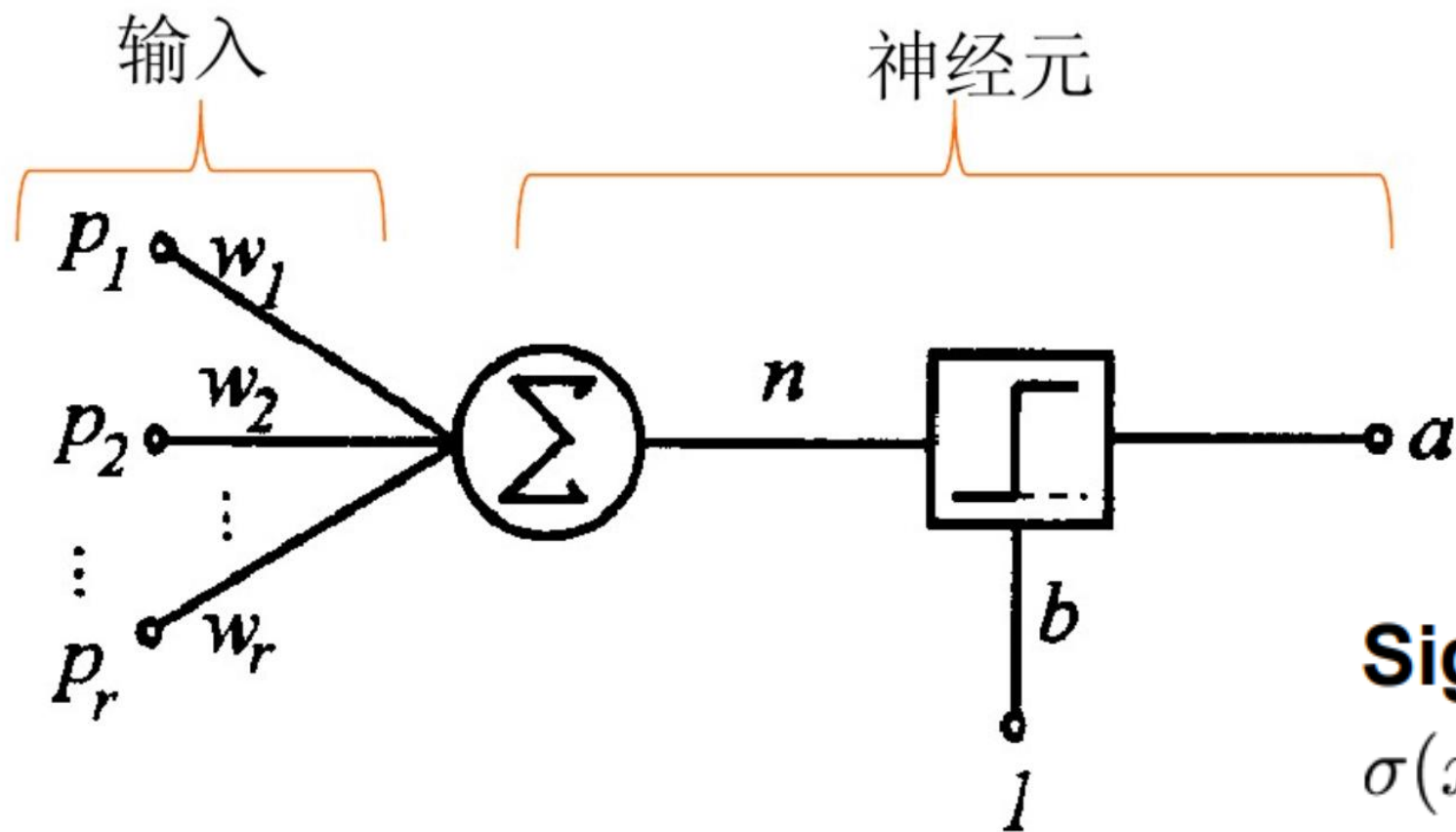


Figure 1.4: A Venn diagram showing how deep learning is a kind of representation learning, which is in turn a kind of machine learning, which is used for many but not all approaches to AI. Each section of the Venn diagram includes an example of an AI technology.

神经网络模型与激活函数

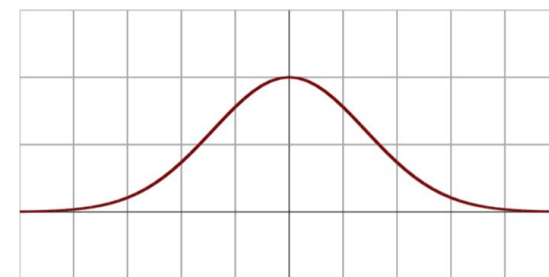
训练神经网络就是计算神经元的权值和阈值



常用的激活函数



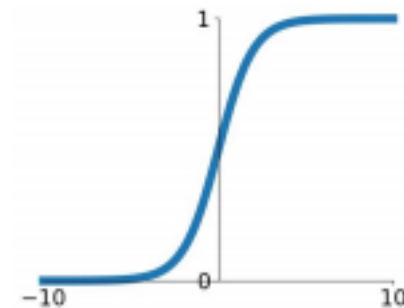
Rectifier



Gaussian

Sigmoid

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



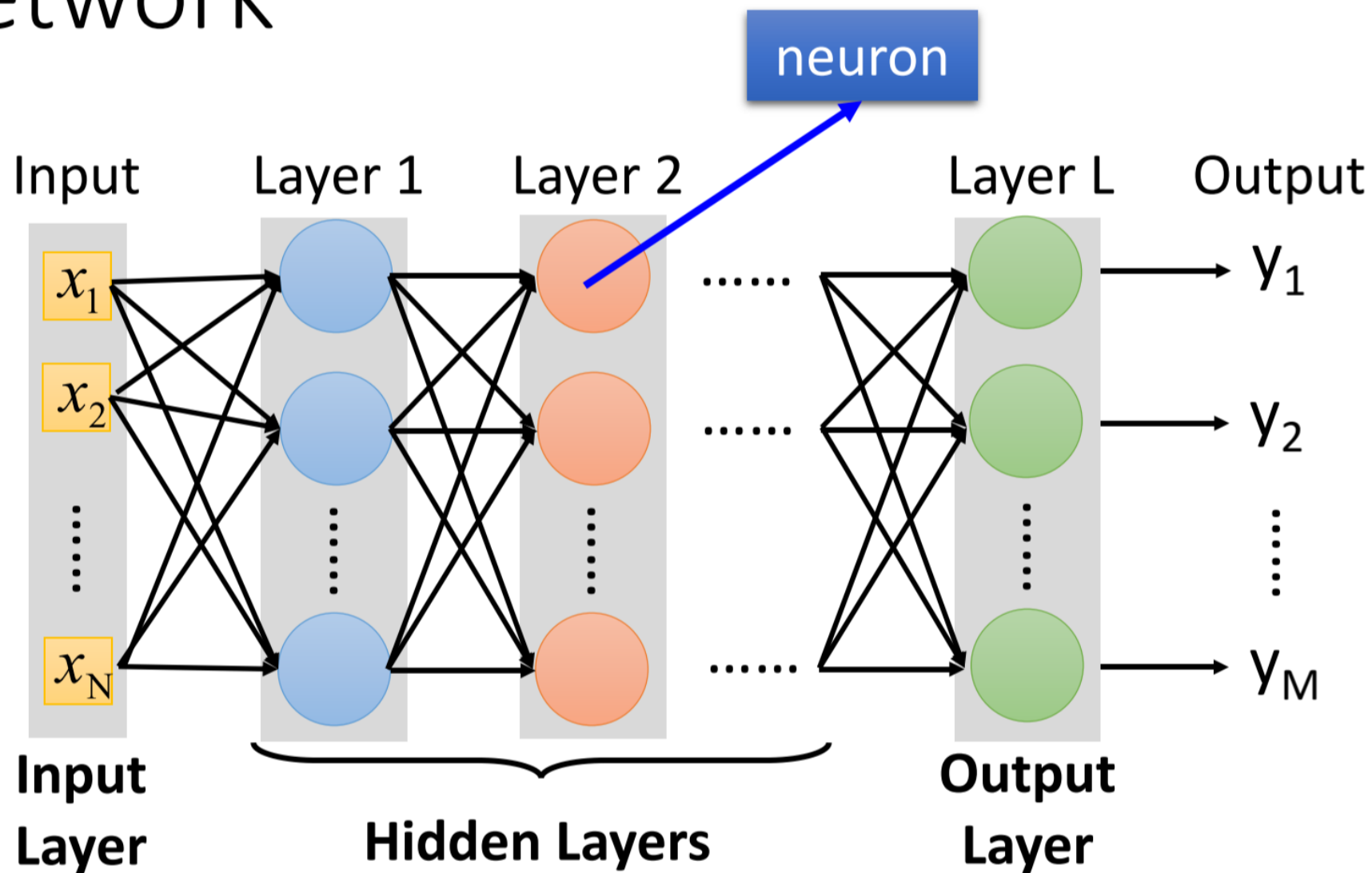
Multi-Layer Perceptron

Network

多层神经网络

神经网络的学习过程，
就是根据训练数据来调整
神经元之间的连接权
以及每个功能神经元的
阈值

换言之，神经网络学
到的东西，蕴含在连接
权和阈值之间



CNN 卷积神经网络

输入二维矩阵or图片 (input)

卷积 (Convolution)

使用滤波器 (卷积核)

池化 (pooling)

降低图像的大小

Flatten

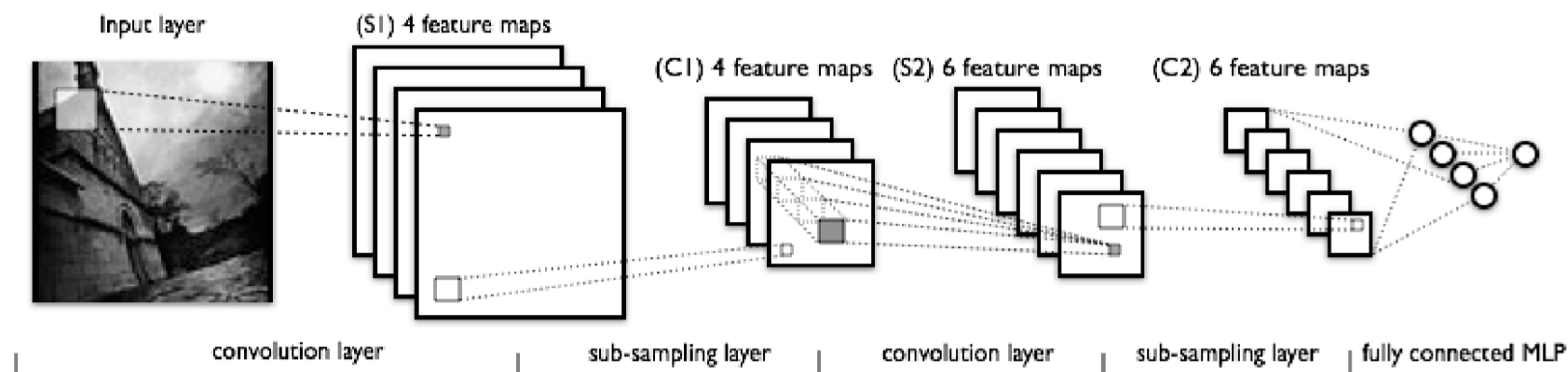
把二维特征映射图连接成一个列向量,
作为训练分类器的输入特征

卷积和池化可以
重复多次

卷积神经网络

是一个适用于图像识别的神经网络。

它模仿人类识别图像的多层过程：瞳孔
摄入像素；大脑皮层某些细胞初步处理，
发现形状边缘、方向；抽象判定形状
(如圆形、方形)；进一步抽象判定
(如判断物体是气球)。



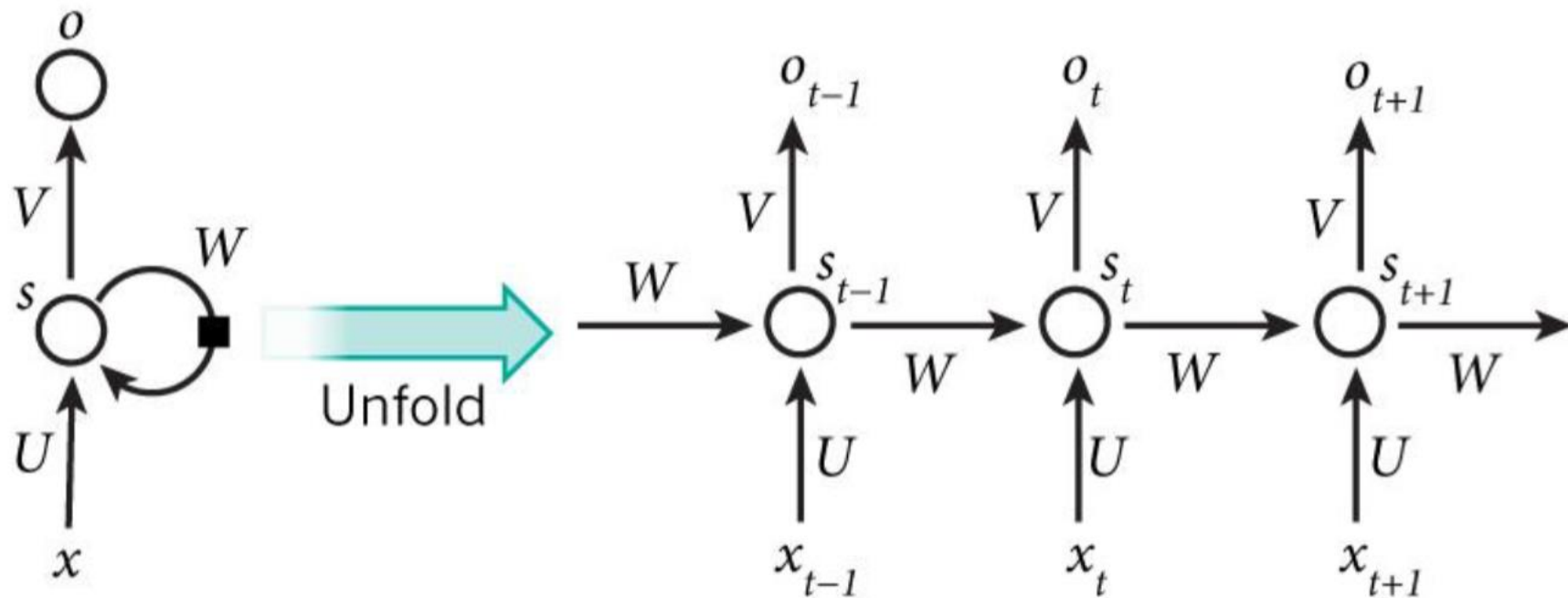
RNN 循环神经网络

序列的当前输出与之前的输出有关，

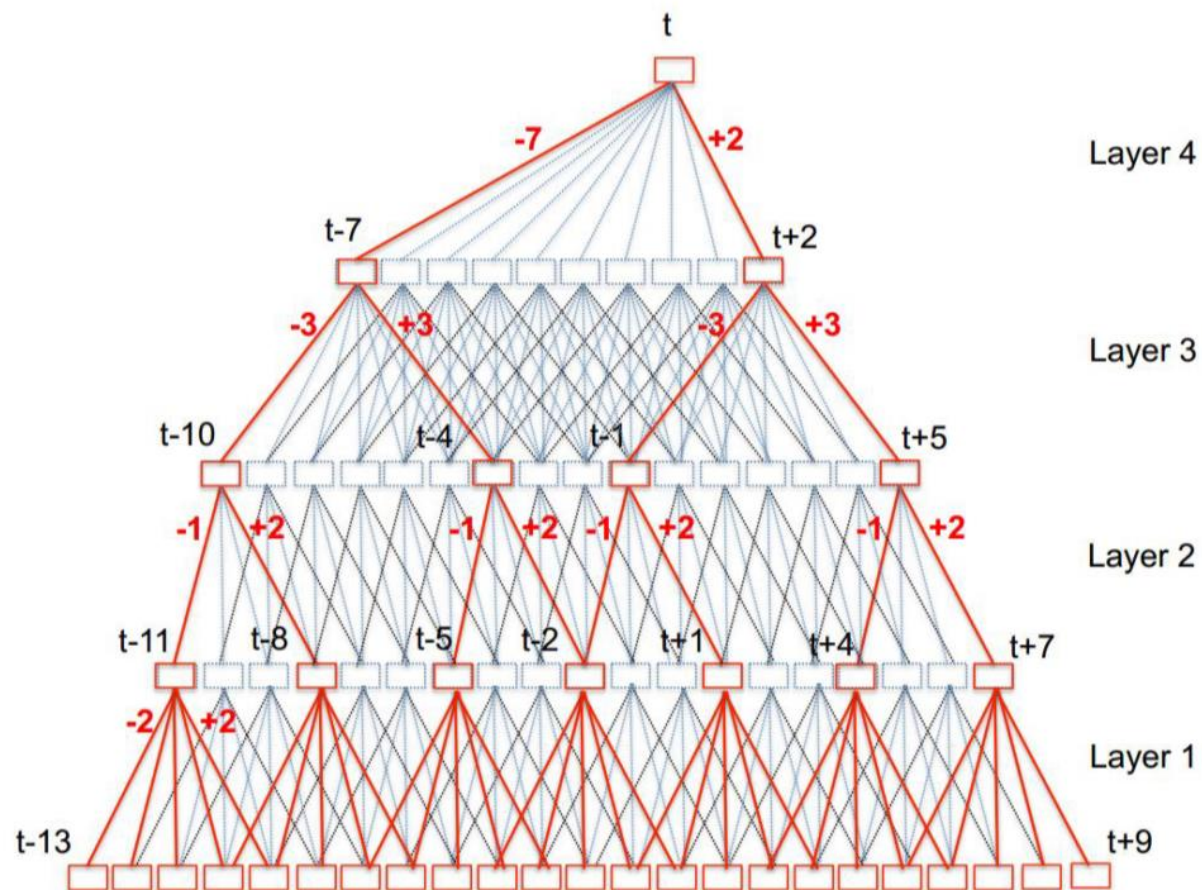
是一个有**反馈**的神经网络结构

适用于自然语言处理 (NLP)

RNN在时域上展开



TDNN 延时神经网络



输入层输入的数据包含多帧

例如图中延时为4, 同时输入5帧数据

该模型适用于语音识别

BP 算法（误差逆向传播）

BP算法是训练神经网络（计算神经网络中的权值和阈值）的一个关键算法，其基本思想分为两步：

1.正向传播时，输入样本从输入层进入网络，经隐层逐层传递至输出层，得到输出层。计算输出误差

2.反向传播时，将输出误差按原通路反传计算，通过隐层反向，直至输入层，在反传过程中将误差分摊给各层的各个单元，获得各层各单元的误差信号，并将其作为修正各单元权值的根据修正权值

2

Speaker Recognition and kaldi

了解与学习了语音识别的原理和训练模型的方法

说话人识别

分类

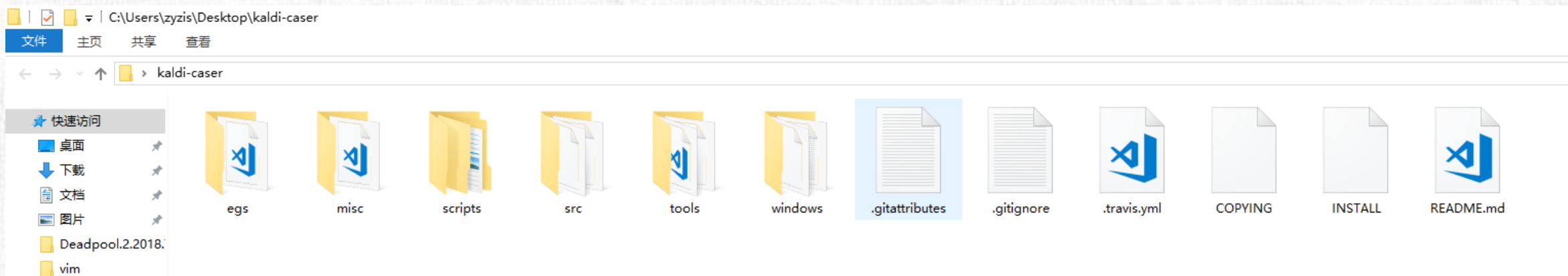
1. 有文本相关的(Text-Dependent)
2. 文本无关的(Text-Independent)

识别方法:

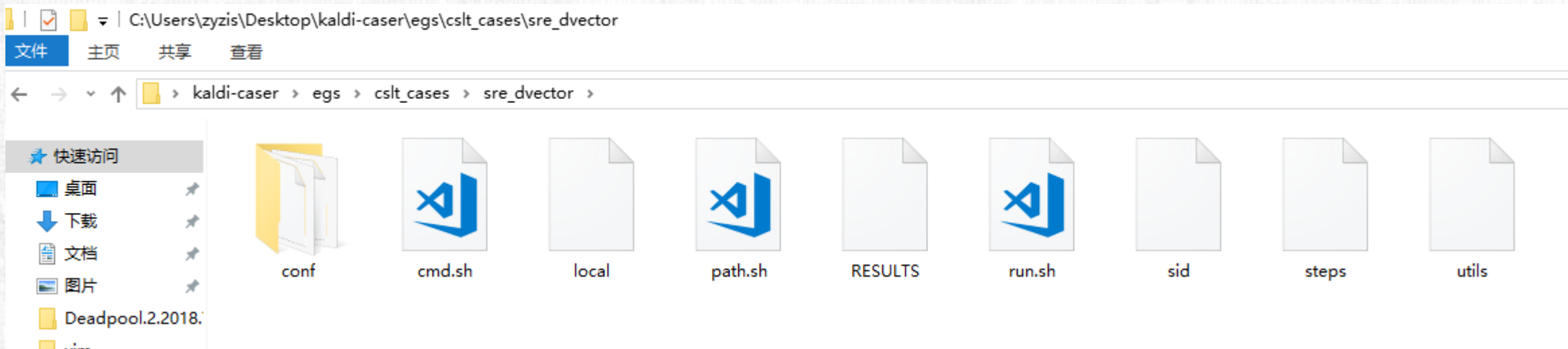
1. i-vector
2. D-vector
3. end-to-end

kaldi

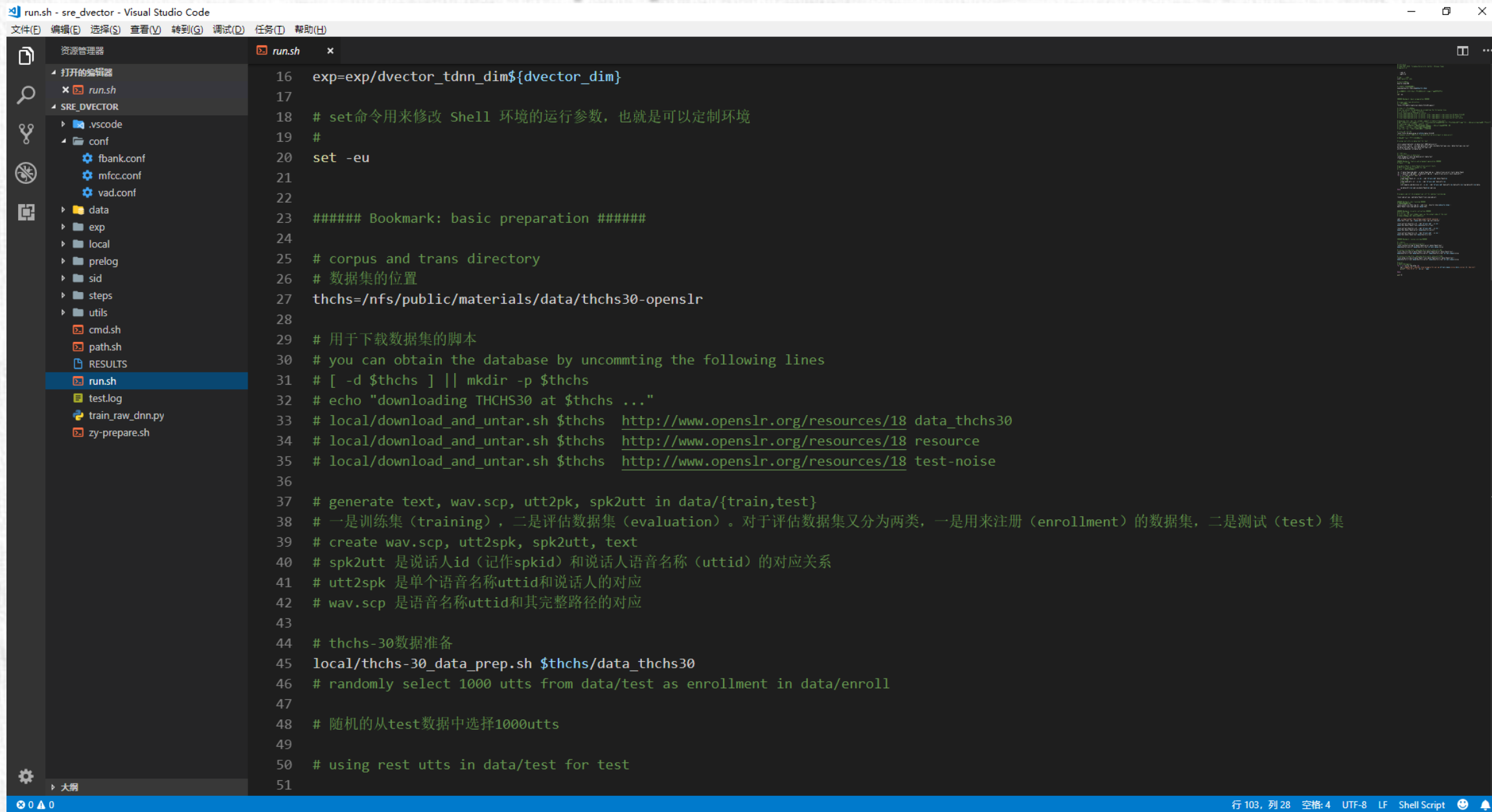
kaldi



egs文件夹下的代码



Data preparation数据准备



```
run.sh - sre_dvector - Visual Studio Code
文件(F) 编辑(E) 选择(S) 查看(V) 转到(G) 调试(D) 任务(T) 帮助(H)

资源管理器
  打开的编辑器
    x run.sh
  SRE_DVECTOR
    .vscode
    conf
      fbank.conf
      mfcc.conf
      vad.conf
    data
    exp
    local
    prelog
    sid
    steps
    utils
    cmd.sh
    path.sh
    RESULTS
    run.sh
    test.log
    train_raw_dnn.py
    zy-prepare.sh

16 exp=exp/dvector_tdnn_dim${dvector_dim}
17
18 # set命令用来修改 Shell 环境的运行参数，也就是可以定制环境
19 #
20 set -eu
21
22
23 ##### Bookmark: basic preparation #####
24
25 # corpus and trans directory
26 # 数据集的位置
27 thchs=/nfs/public/materials/data/thchs30-openslr
28
29 # 用于下载数据集的脚本
30 # you can obtain the database by uncommting the following lines
31 # [ -d $thchs ] || mkdir -p $thchs
32 # echo "downloading THCHS30 at $thchs ..."
33 # local/download_and_untar.sh $thchs http://www.openslr.org/resources/18 data_thchs30
34 # local/download_and_untar.sh $thchs http://www.openslr.org/resources/18 resource
35 # local/download_and_untar.sh $thchs http://www.openslr.org/resources/18 test-noise
36
37 # generate text, wav.scp, utt2pk, spk2utt in data/{train,test}
38 # 一是训练集 (training)，二是评估数据集 (evaluation)。对于评估数据集又分为两类，一是用来注册 (enrollment) 的数据集，二是测试 (test) 集
39 # create wav.scp, utt2spk, spk2utt, text
40 # spk2utt 是说话人id (记作spkid) 和说话人语音名称 (uttid) 的对应关系
41 # utt2spk 是单个语音名称uttid和说话人的对应
42 # wav.scp 是语音名称uttid和其完整路径的对应
43
44 # thchs-30数据准备
45 local/thchs-30_data_prep.sh $thchs/data_thchs30
46 # randomly select 1000 utts from data/test as enrollment in data/enroll
47
48 # 随机的从test数据中选择1000utts
49
50 # using rest utts in data/test for test
51
```

行 103, 列 28 空格: 4 UTF-8 LF Shell Script

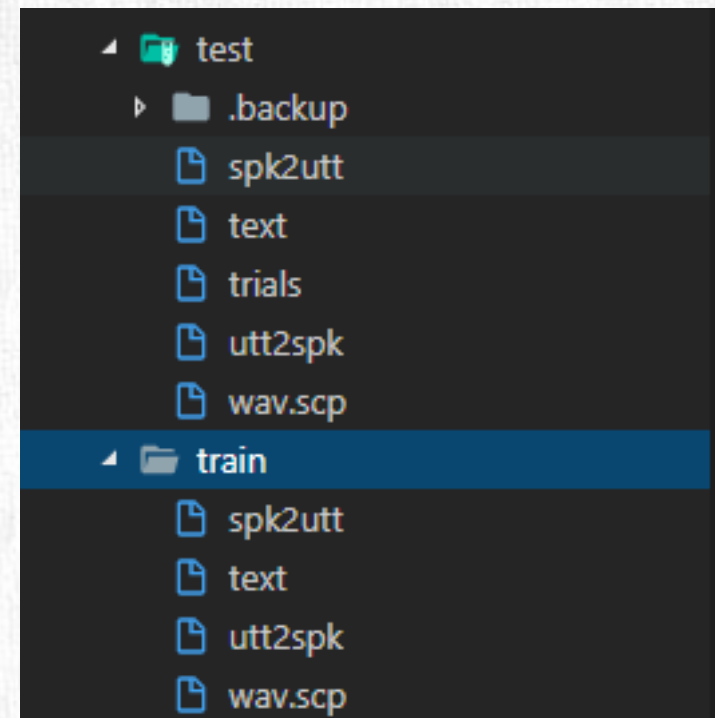
Data preparation数据准备

对于这些数据集中，一是训练集（training），二是评估数据集（evaluation）。对于评估数据集又分为两类，一是用来注册（enrollment）的数据集，二是测试（test）集

准备的过程中生成了以下类型的文件：

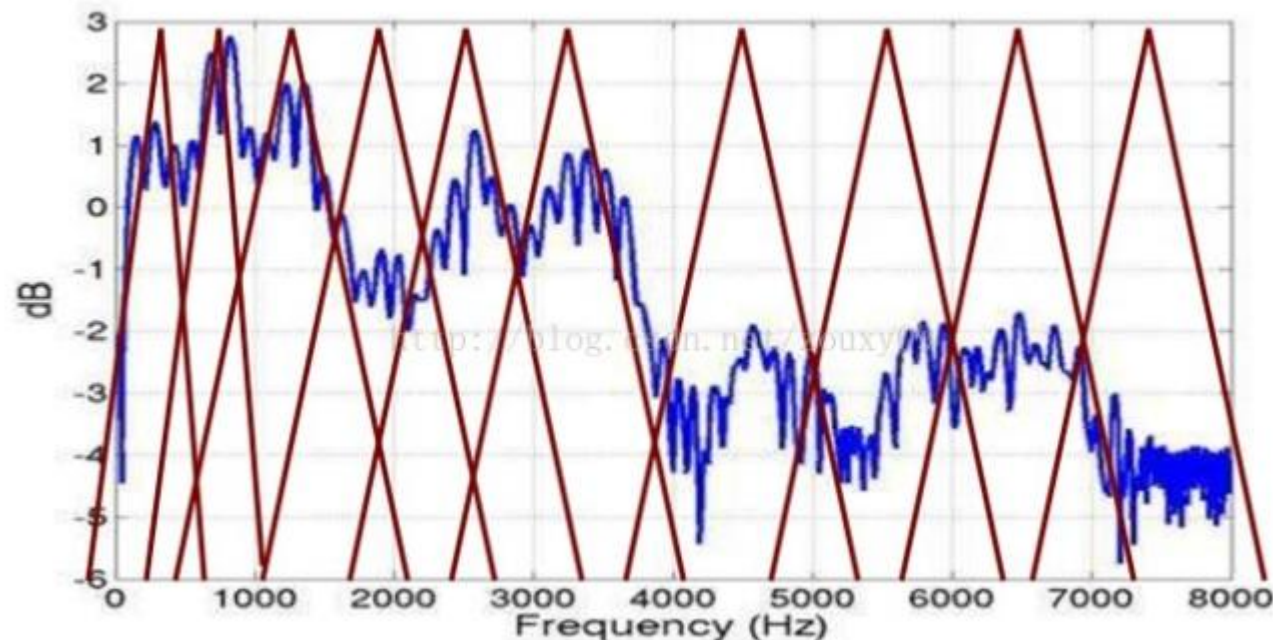
根据kaldi官方文档，我们可以得知，图中所示的四种文件格式和其作用。

1. **Text** 包含了每个语音数据的文字信息
2. **spk2utt** 是说话人id（记作spkid）和说话人语音名称（uttid）的对应关系
3. **utt2spk** 是单个语音名称uttid和说话人的对应
4. **wav.scp** 是语音名称uttid和其完整路径的对应



feature and alignment generation特征提取与对齐

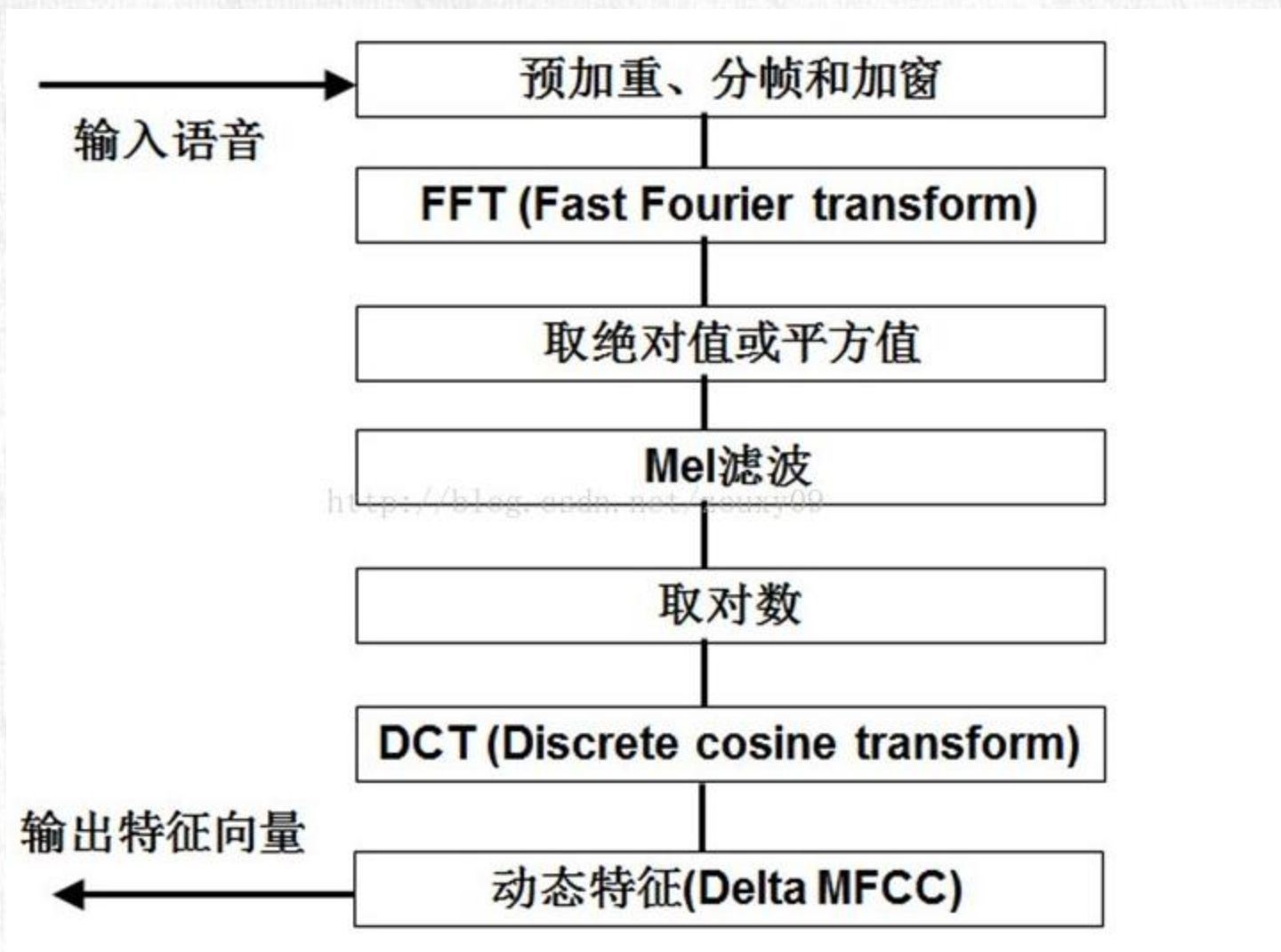
特征提取的好坏直接影响系统的效能。根据人耳的听觉特性提取具有辨识度的音频信息，使用shell脚本对已经准备好的train enroll test数据进行遍历，运行make_mfcc.sh脚本对语音进行mfcc特征提取，运行make_fbank.sh脚本对数据进行fbank滤波器带特征提取，运行compute_vad_decision.sh进行vad（语音活动检测）



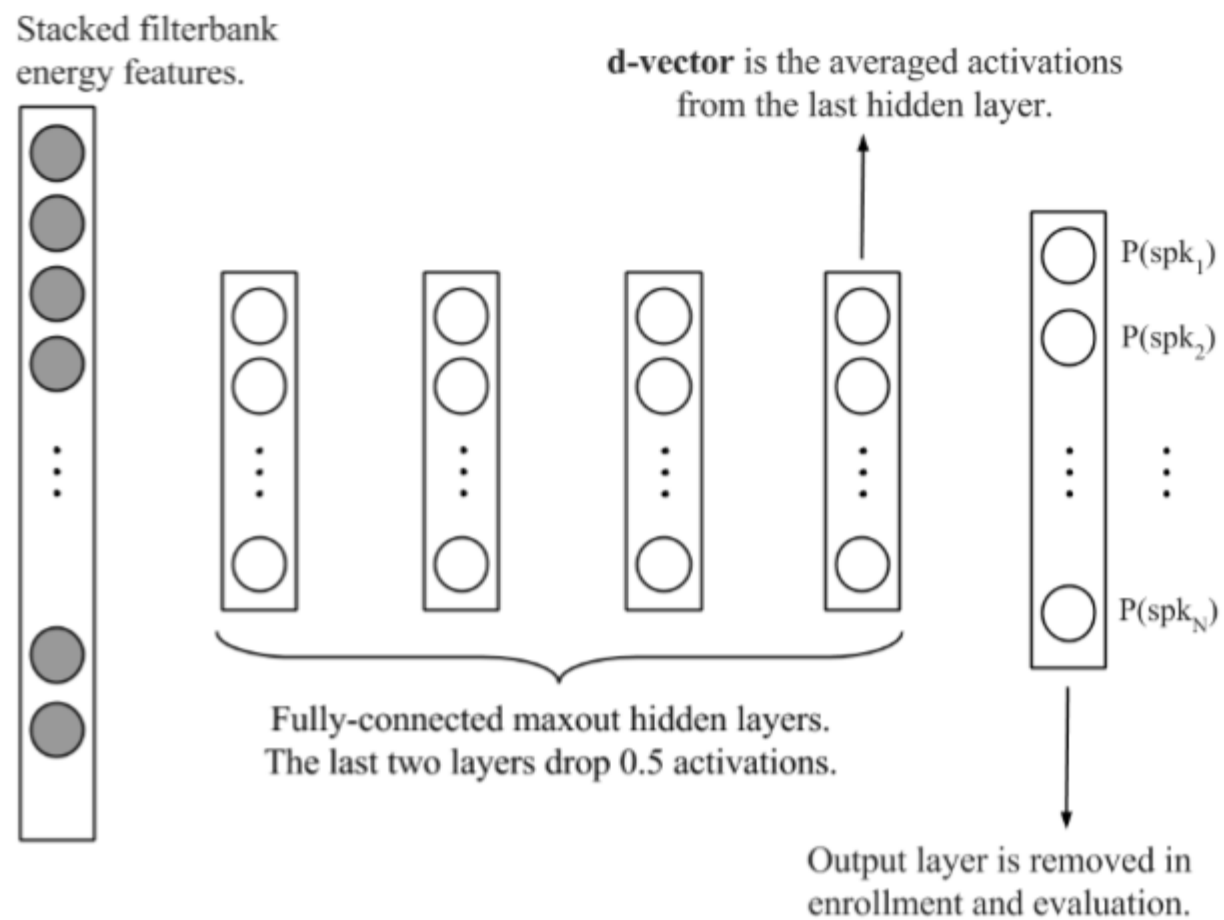
feature and alignment generation特征提取与对齐

对信号进行

1. 预加重——提高语音的高频部分
2. 分帧——由于傅里叶变换要求的是平稳的语音信号，因此要将语音信号分成小段小段的
3. 加窗处理——加上一个窗函数，例如汉明窗，消除吉布斯效应



DNN-training 训练深度神经网络

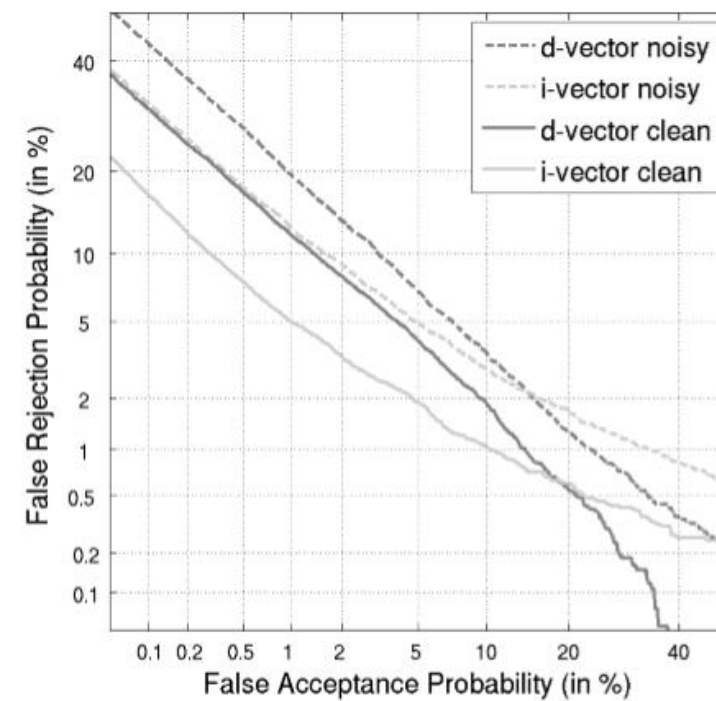
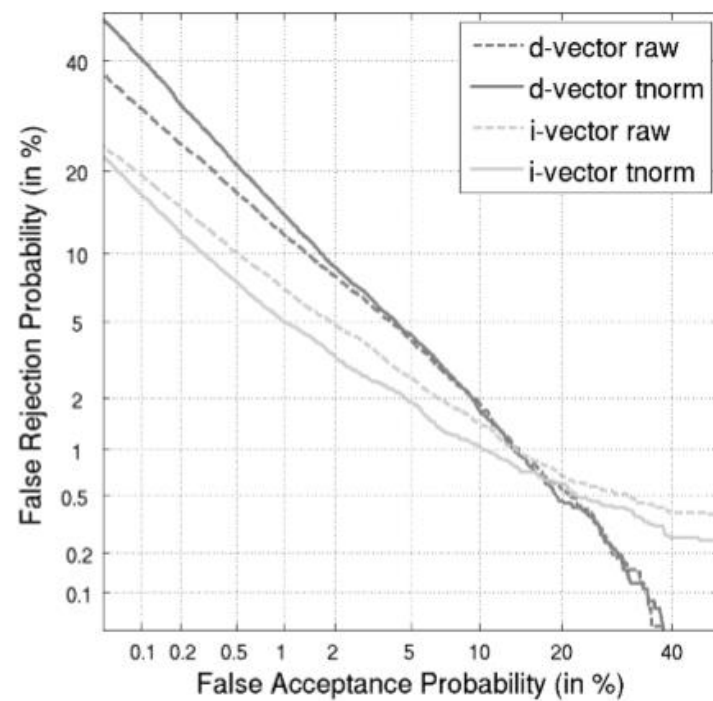


DNN-training训练深度神经网络

```
train_raw_dnn.py x
22 import libs.nnet3.train.frame_level_objf as train_lib
23 import libs.nnet3.report.log_parse as nnet3_log_parse
24
25
26 logger = logging.getLogger('libs')
27 logger.setLevel(logging.INFO)
28 handler = logging.StreamHandler()
29 handler.setLevel(logging.INFO)
30 formatter = logging.Formatter("%(asctime)s [%(pathname)s:%(lineno)s - " ...
32 handler.setFormatter(formatter)
33 logger.addHandler(handler)
34 logger.info('Starting raw DNN trainer (train_raw_dnn.py)')
35
36
37 def get_args(): ...
121
122
123 def process_args(args): ...
179
180
181 def train(args, run_opts): ...
503
504
505 def main(): ...
520
521
522 if __name__ == "__main__":
523     main()
524
```

train_raw_dnn.py

1. 参数获取函数，获取参数，控制训练的运行方式、是否使用GPU、训练的次数等信息，
2. 训练函数，控制具体的训练流程和算法，
3. Main函数，使用try和except语句进行整体函数的调用和流程的控制。



```
cosine eer: 2.74  
| lda eer: 1.34  
| plda eer: 2.34
```


End-to-End

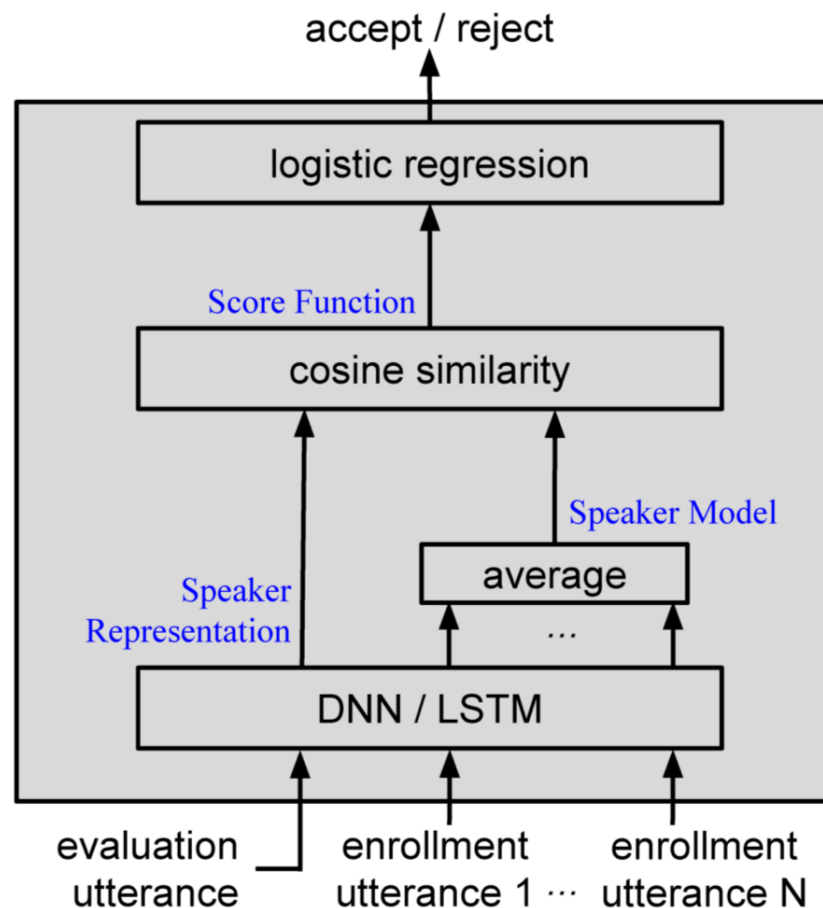


Figure 2: *End-to-end architecture, combining training of speaker representations (DNN/LSTM box), estimation of a speaker model based on up to N "enrollment" utterances (in blue), and verification (cosine similarity/logistic regression boxes).*

传统机器学习的流程往往由多个独立的模块组成

由于如今计算机硬件和深度神经网络的成熟，我们可以使用“端到端”只有一个模块的实现识别
训练时使用BP算法对模型进行修正

如图所示，神经网络采用深度神经网络和长短期记忆网络。
训练时通过一个评分函数，对于结果只有reject和accept

3

test

在服务器上运行和修改了部分代码，重现了老师的研究结果和尝试训练自己的模型

中文语音识别

输入 `bash run.sh --test-mode simulated`

Termius - csIt

Logged out

Port Forwarding

Hosts

csIt

History

tecent-root (3)

tecent-zy (2)

tecent-root (2)

tecent-zy (2)

tecent-root

45.32.128.89

tecent-root (4)

raspberry (2)

tecent-root

csIt (2)

Show more...

```
-bash-4.2$ pwd
/work4/preintern/zhangy/kaldi/egs/thchs30/online_demo
-bash-4.2$ ./run.sh --test-mode simulated

SIMULATED ONLINE DECODING - pre-recorded audio is used

The (bigram) language model used to build the decoding graph was
estimated on an audio book's text. The text in question is
"King Solomon's Mines" (http://www.gutenberg.org/ebooks/2166).
The audio chunks to be decoded were taken from the audio book read
by John Nicholson(http://librivox.org/king-solomons-mines-by-haggard/)

NOTE: Using utterances from the book, on which the LM was estimated
is considered to be "cheating" and we are doing this only for
the purposes of the demo.

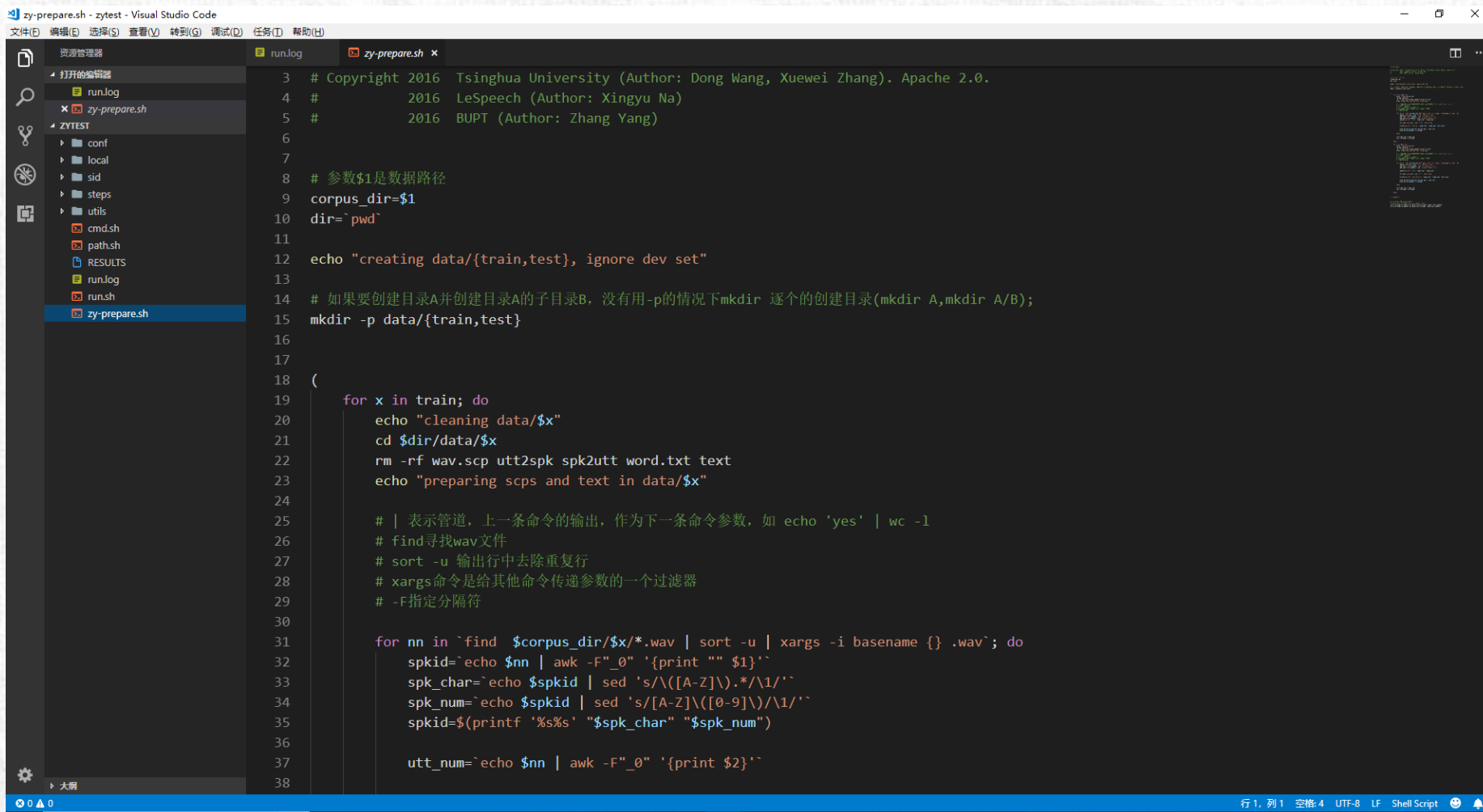
You can type "./run.sh --test-mode live" to try it using your
own voice!

online-wav-gmm-decode-faster --verbose=1 --rt-min=0.8 --rt-max=0.85 --max-active=4000 --beam=1
ls/tril/final.mdl online-data/models/tril/HCLG.fst online-data/models/tril/words.txt 1:2:3:4:5
File: A11_176
金融 水运 口岸 八十 一个 空运 口岸 二十五 个 铁路 口岸 时 个 公路 口岸 三十二 的

File: D11_750
中美 间 的 一些 爱国 将 是 巴 镇 山 梨 一 度 曾 居 五 艘 彼 岸 单 体 没 等 也 奋 起 抗 战

compute-wer --mode=present ark,t:./work/ref.txt ark,t:./work/hyp.txt
%WER -nan [ 0 / 0, 0 ins, 0 del, 0 sub ]
%SER -nan [ 0 / 0 ]
Scored 0 sentences, 0 not present in hyp.
-bash-4.2$
```

维吾尔族语言声纹识别



```
zy-prepare.sh - zytest - Visual Studio Code
文件(F) 编辑(E) 选择(S) 查看(V) 转到(G) 调试(D) 任务(T) 帮助(H)

资源管理器
  打开的编辑器
    run.log
    zy-prepare.sh
  ZYTEST
    conf
    local
    sid
    steps
    utils
    cmd.sh
    path.sh
    RESULTS
    run.log
    run.sh
    zy-prepare.sh

3 # Copyright 2016 Tsinghua University (Author: Dong Wang, Xuewei Zhang). Apache 2.0.
4 #       2016 LeSpeech (Author: Xingyu Na)
5 #       2016 BUPT (Author: Zhang Yang)
6
7
8 # 参数$1是数据路径
9 corpus_dir=$1
10 dir=`pwd`
11
12 echo "creating data/{train,test}, ignore dev set"
13
14 # 如果要创建目录A并创建目录A的子目录B, 没有用-p的情况下mkdir 逐个的创建目录(mkdir A,mkdir A/B);
15 mkdir -p data/{train,test}
16
17
18 (
19   for x in train; do
20     echo "cleaning data/$x"
21     cd $dir/data/$x
22     rm -rf wav.scp utt2spk spk2utt word.txt text
23     echo "preparing scps and text in data/$x"
24
25     # | 表示管道, 上一条命令的输出, 作为下一条命令参数, 如 echo 'yes' | wc -l
26     # find寻找wav文件
27     # sort -u 输出行中去重重复行
28     # xargs命令是给其他命令传递参数的一个过滤器
29     # -F指定分隔符
30
31     for nn in `find $corpus_dir/$x/*.wav | sort -u | xargs -i basename {} .wav`; do
32       spkid=`echo $nn | awk -F"_" '{print "" $1}`
33       spk_char=`echo $spkid | sed 's/\([A-Z]\)/./\1/'`
34       spk_num=`echo $spkid | sed 's/[A-Z]\([0-9]\)/\1/'`
35       spkid=$(printf '%s%s' "$spk_char" "$spk_num")
36
37       utt_num=`echo $nn | awk -F"_" '{print $2}``
38     done
2025/05/15 10:16
```

```
cosine eer: 3.61
lda eer: 0.94
plda eer: 2.21
```


我的收获

最后一个部分是我的收获和体会。

在这“不算很短的”预实习阶段中，这是我第一次接触到Linux远程服务器，并且还是多个高性能的服务器集群。在好奇心的驱动下，我也在腾讯云上购买了属于自己的服务器和树莓派等尝试进行云计算。在暑假闲暇的时候尝试搭建了自己的个人博客和git服务器。尝试做在线语音识别，但由于自己的服务器计算性能不足，失败了。

同时，因为自己是通信专业的学生，很多知识比如声学模型和傅里叶变换、马尔科夫链等，在我先前的专业课程中有所了解和学习，在运行kaldi中进行说话人识别的过程中，我更加深入的了解了这些通信专业课所学知识在语音识别和说话人识别领域的应用。

tecent-root (3)

tecent-zy (2)

tecent-root (2)

tecent-zy (2)

tecent-root

45.32.128.89

tecent-root (4)

raspberry (2)

tecent-root

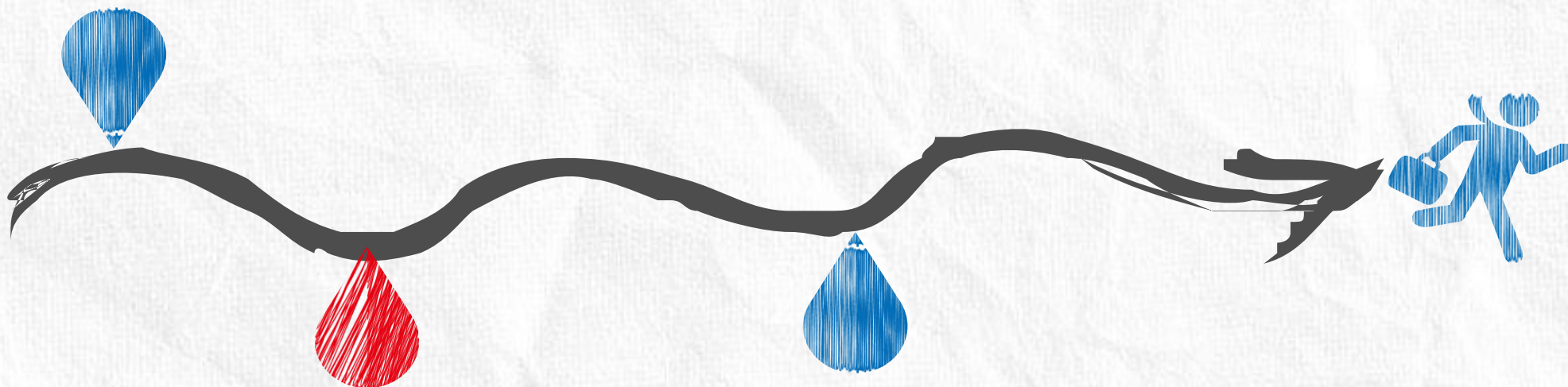
cslt (2)

Show more...

最后

感谢老师给我不止有2周的实习机会和服务器的使用机会

我的能力与水平有限，还有很多知识不会，但未来我将会认真学习



THANK S

BUPT



张阳