

# HCT NLP Week 2

问答摘要与推理  
词向量实践

# Outline

1. 词向量计算优化方法
2. 初识深度学习框架
3. Gensim代码实践
4. 作业

3, 1, 2, 3  
n

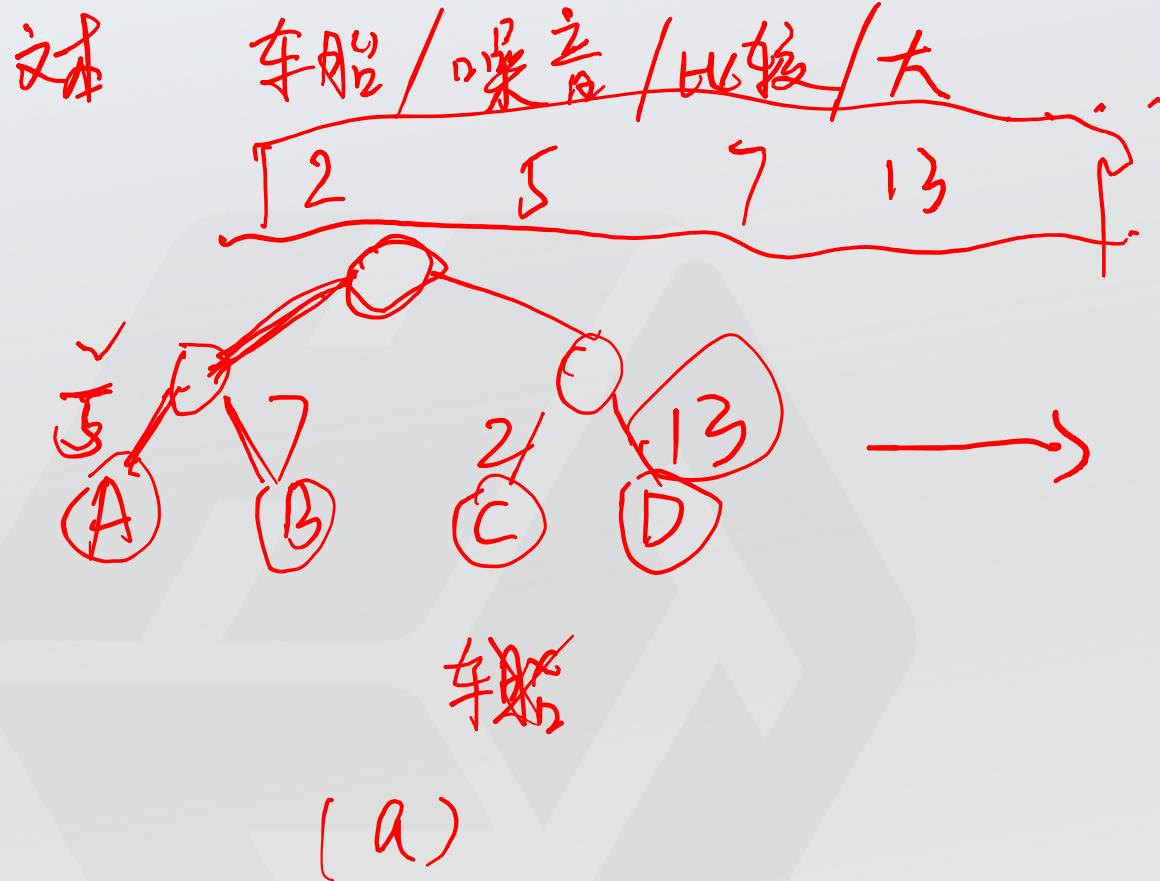
Python  
~~log~~ sorted 分治

## 1. 词向量计算优化方法

$$O(v) \rightarrow O(\log v)$$

# Hierarchical Softmax

## Huffman Tree (哈夫曼树)



$$\text{Path\_a} = 5 \times 2 + 7 \times 2 + 2 \times 2 + 13 \times 2 = 54$$

$$\text{Path\_b} = 13 \times 1 + 7 \times 2 + 2 \times 3 + 5 \times 3 = 48$$

哈夫曼树（b）：

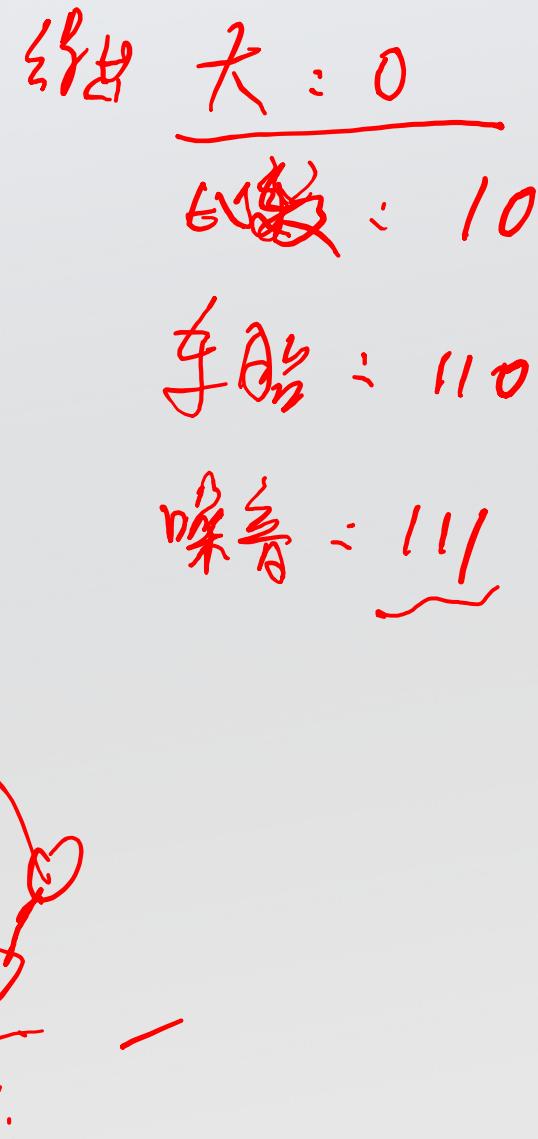
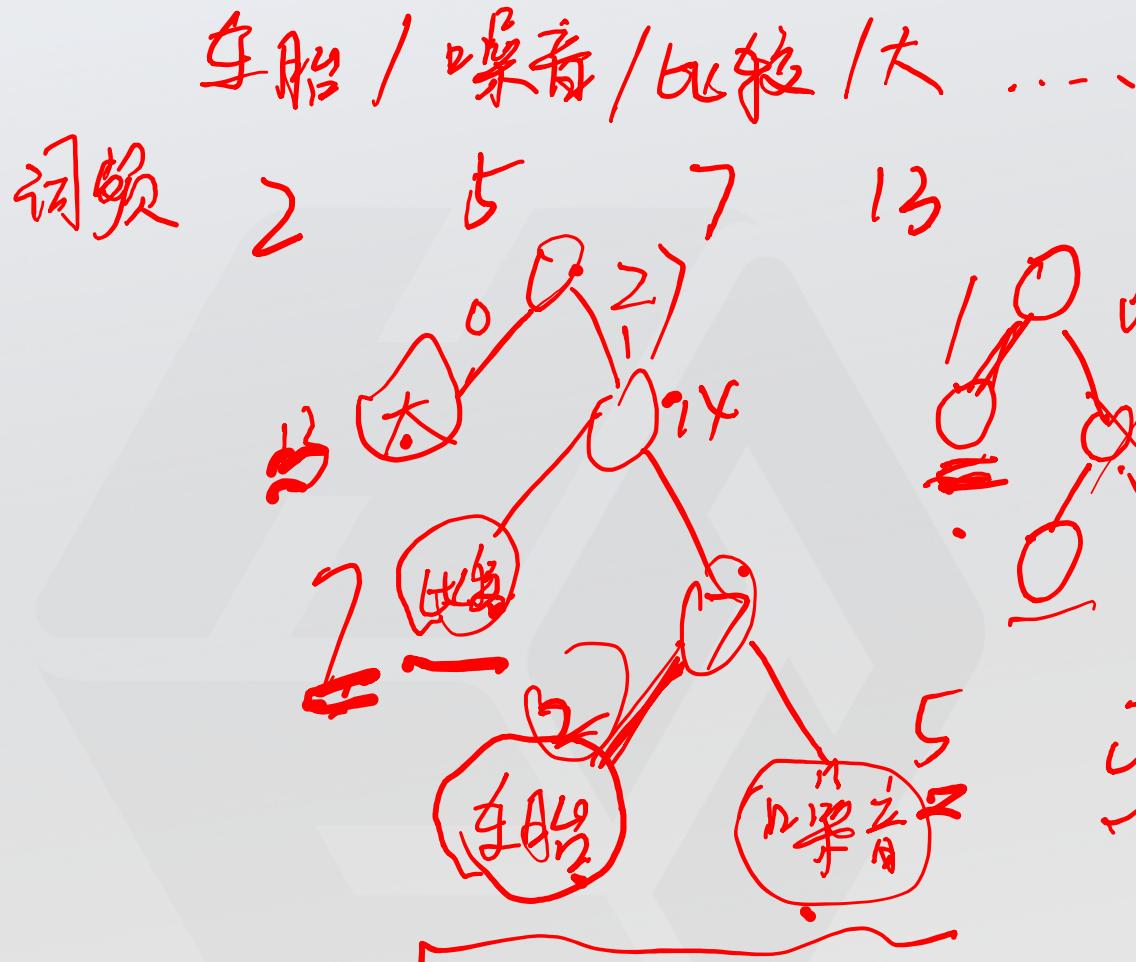
```

graph TD
    Root(( )) --- Node1(( ))
    Root --- Node2(( ))
    Node1 --- Node3(( ))
    Node1 --- Node4(( ))
    Node2 --- Node5(( ))
    Node2 --- Node6(( ))
    Node3 --- Node7((B))
    Node3 --- Node8((C))
    Node4 --- Node9((A))
    Node4 --- Node10((D))
    
```

(b)

# Hierarchical Softmax

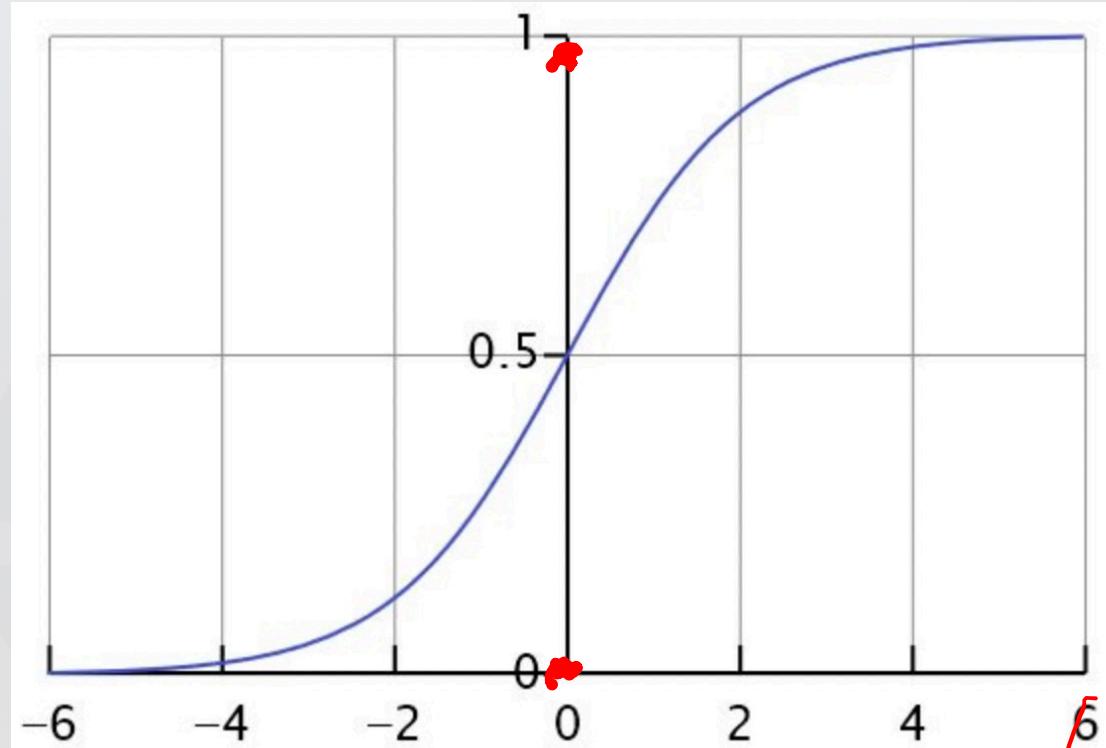
Huffman Tree (哈夫曼树)



- 1)  $V$  个叶子节点
- 2)  $V-1$  个中间节点
- 3) 一一对应

# Hierarchical Softmax

## Logistic Regression LR



$$\hat{y} = P(y=1|x) \quad 0 \leq \hat{y} \leq 1$$

$$\hat{y} = w^T x + b$$

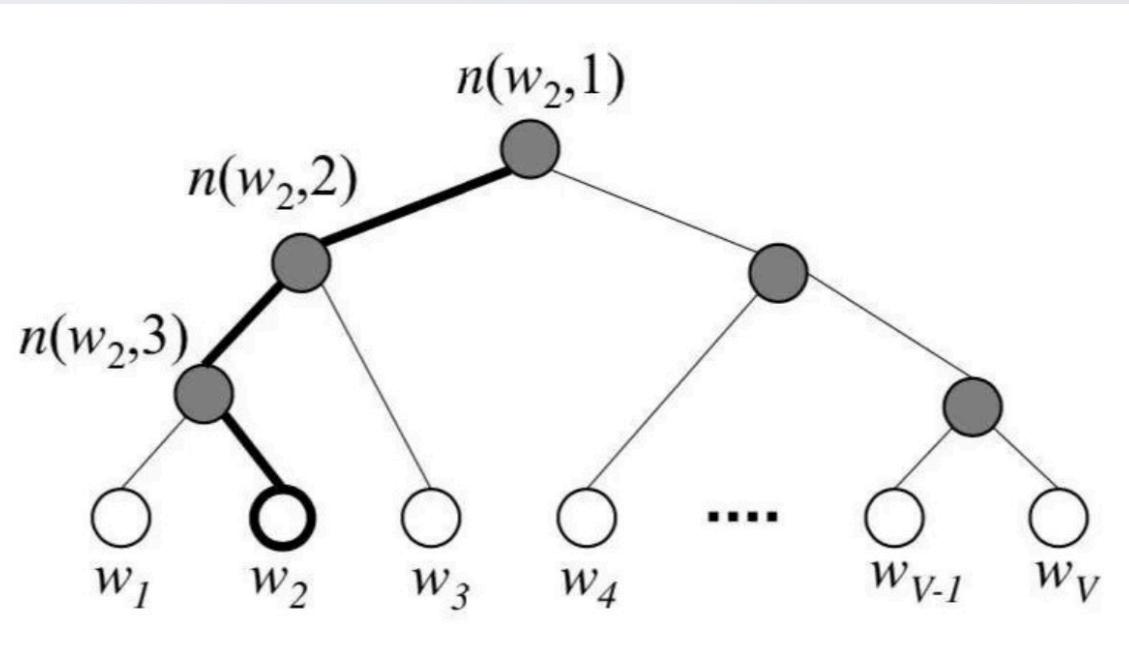
$$w^T x = \sum_{i=1}^n w_i x_i = w_0 + w_1 x_1 + \dots + w_n x_n$$

$$\sigma(z) = \frac{1}{1+e^{-z}} = \frac{1}{1+e^{-(w^T x + b)}}$$

$$\text{Loss} \quad P(y|x) = \begin{cases} \hat{y} & y=1 \\ 1-\hat{y} & y=0 \end{cases}$$

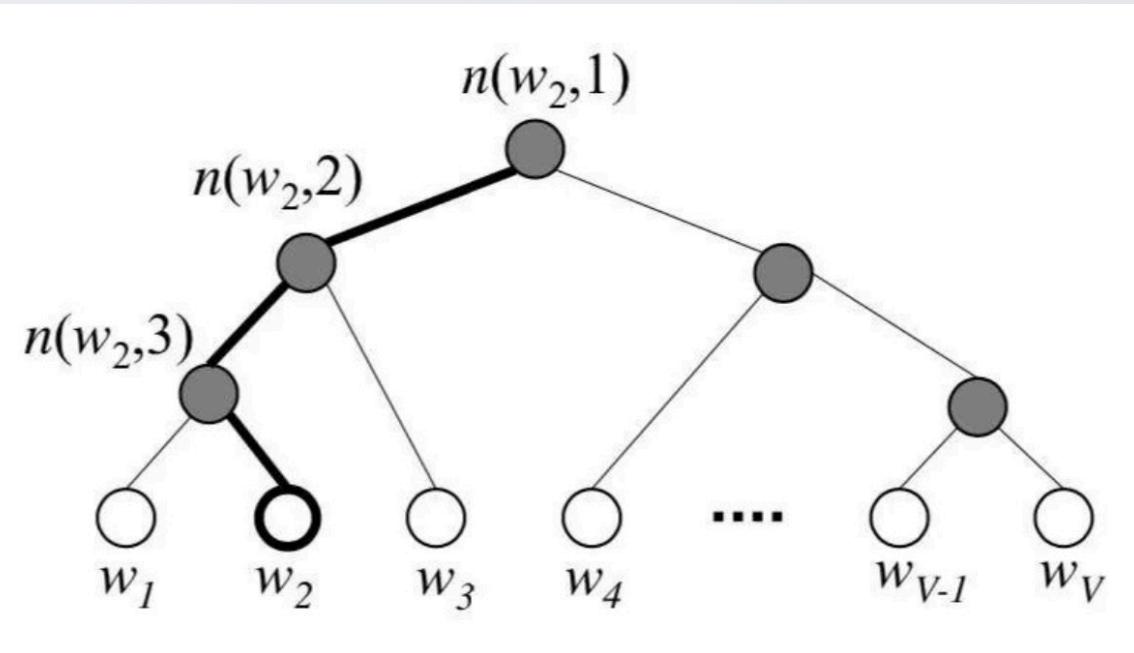
$$\log(P(y|x)) = \underbrace{\hat{y} \log \hat{y}}_{y \log y} + \underbrace{(1-\hat{y}) \log(1-\hat{y})}_{(1-y) \log(1-y)}$$

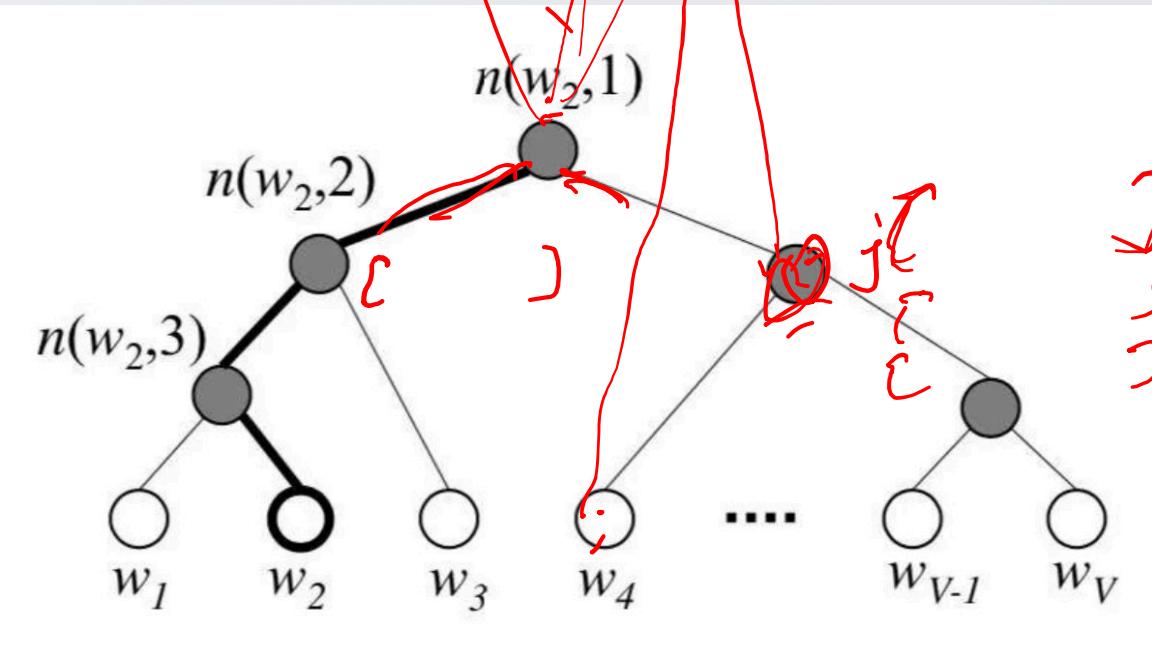
# Hierarchical Softmax



1

# Hierarchical Softmax





$P(w | \text{context}(w))$

$$= \prod_{j=2}^L P(d_j^w | x_w, \theta_{j-1}^w)$$

$d_j^w$  表示路径第  $j$  个节点对应编码

$\theta_{j-1}^w$  : 路径中非叶子节点对应的参数向量

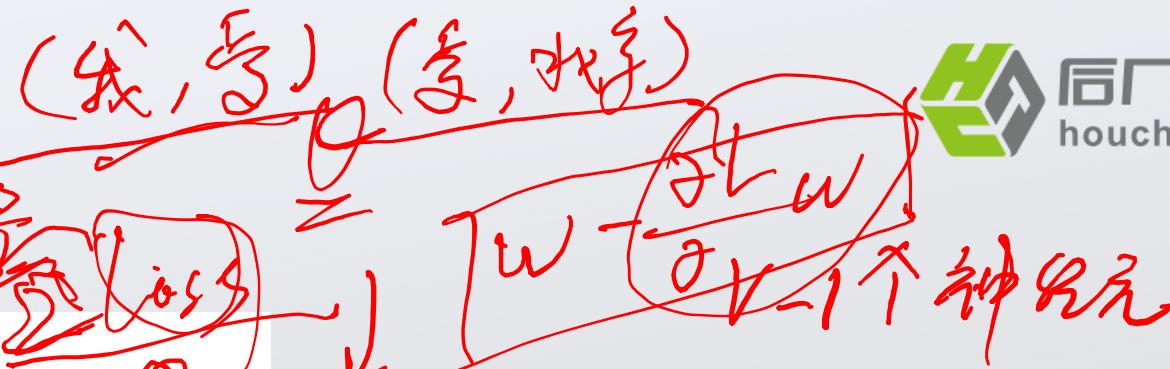
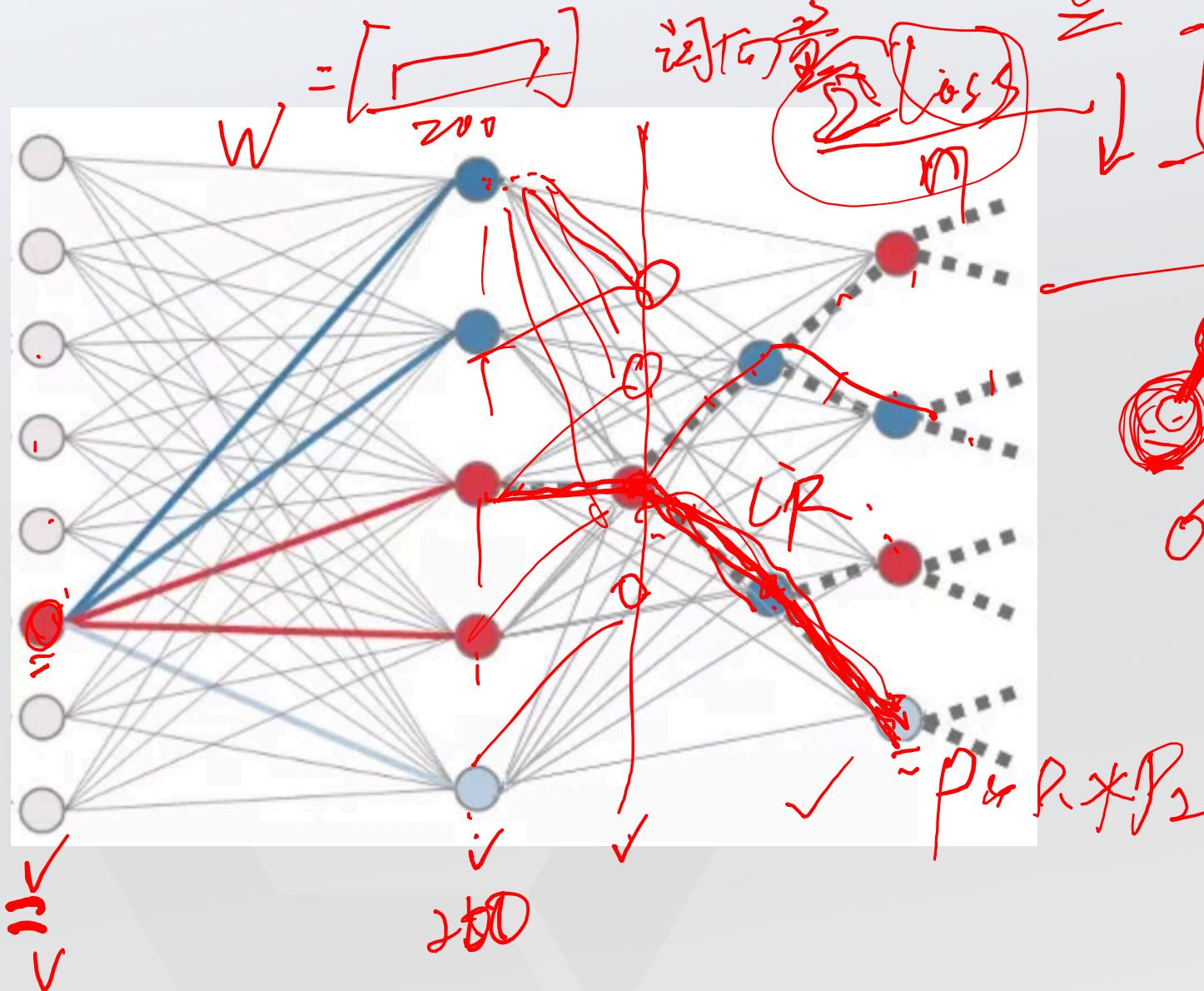
$L^w$ : 路径中包含节点的个数

左侧 增加 缓慢

$$\checkmark P(d_j^w | x_w, \theta_{j-1}^w) = \begin{cases} G(x_w^T \theta_j^w) & d_j^w = 0 \\ -G(x_w^T \theta_j^w) & d_j^w = 1 \end{cases}$$

$(1-d_j^w) \log [ ] + d_j^w \cdot \log [ ]$

# Hierarchical Softmax



优点: ①  $O(V)$   
 $\rightarrow O(\log_2 V)$



$$P_1 * P_2 * P_3$$

# Negative Sampling

uni bi

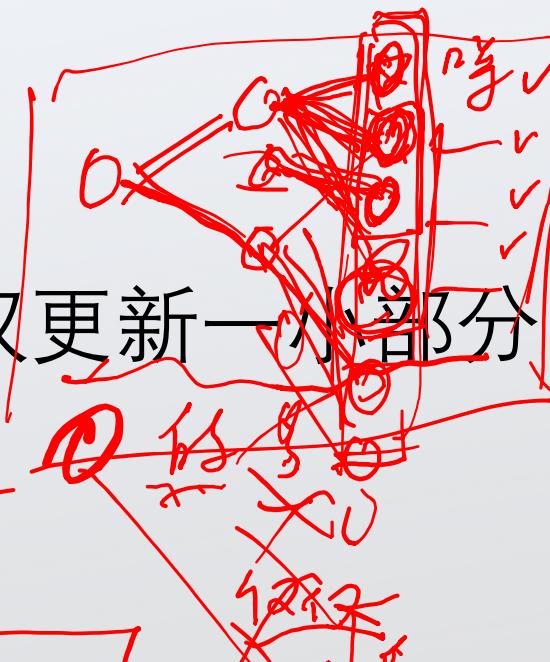
每次让一个训练样本仅仅更新一小部分的权重

车胎 / 烟 / 比较 / 大  
轮胎 / 烟 / 比较 / 大

更新小部分权重

(input: 车胎, output 烟) ✓ 5

negative words ✗ 12



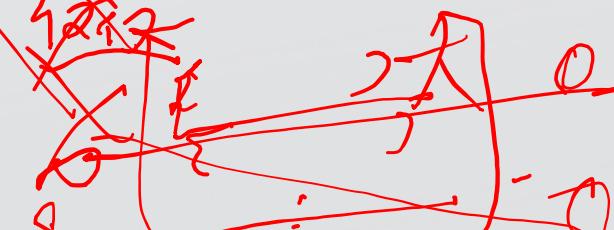
10000



后厂理工学院  
houchangtech.com

6-21

n datasets 5-20个



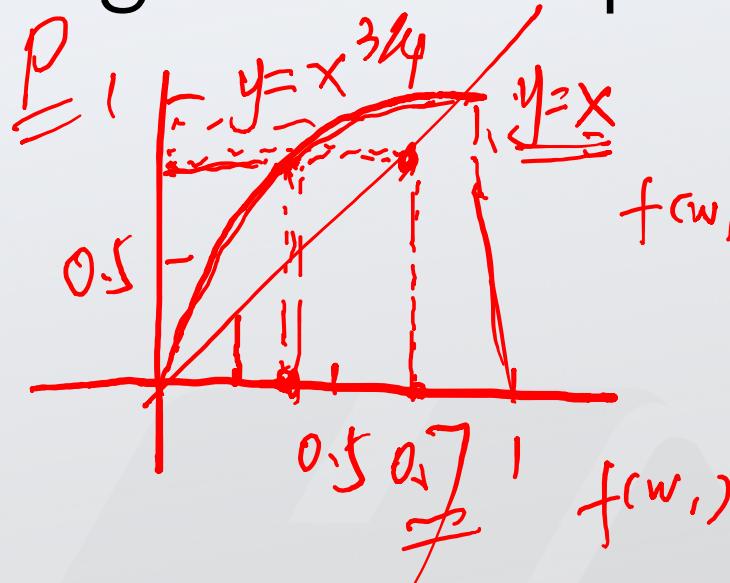
$$300 \times 10000 = 300W$$

$$300 \times 6 = 1800$$

0.06%

$$P(w_i) = \frac{(f(w_i)^{3/4})}{\sum_{j=0}^n (f(w_j)^{3/4})}$$

# Negative Sampling



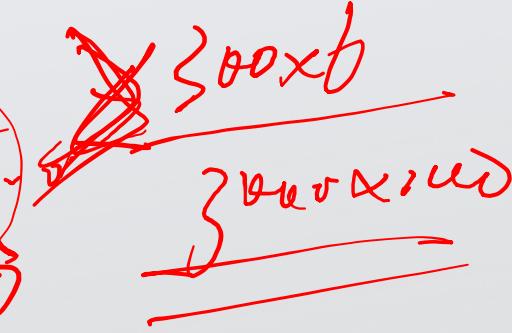
(1)

(2) 将长边1线段分为 $M$ 等份

$300 \times 10000$

~~$200 \times 10000$~~

~~$200 \times 6$~~



高频词正采样

$$M \gg |V|$$



$2 - 10$ 倍

① 对高频词效果更好

② 向量维度较低时效果更好

③ 维度高时近似误差比较大

$$\text{len}(w) = P(w_i)$$

$M \quad 10^8$

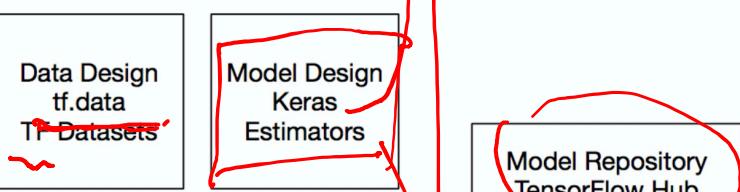
$$P(w_i) = \frac{\text{Sample}}{\sqrt{\text{freq}(w_i)}}$$

## 2. 初识深度学习框架

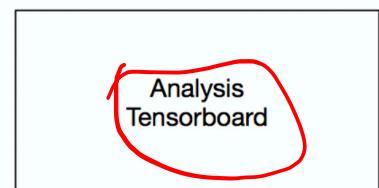
# 框架产品结构

tf, keras, layer, rnn  
TensorFlow

Training



Training  
Distribution Strategy  
CPU GPU TPU



C++ LV → pytorch → Tensorflow  
gather\_nd\_op → shape.  
(64, 16, 200) → .

知识多端  
量化 float32 → 16 int  
Paddle lab

量化 float32  
2.1-rc 版本

Deployment

Cloud  
TensorFlow Extended

Android, iOS  
TensorFlow Lite

Browser and Node  
TensorFlow.JS

Other Language Bindings  
C, Java, Go



On CUDA  
RAND

# 动态图机制

## DyGraph

更加灵活便捷的代码组织结构：使用python的执行控制流程和面向对象的模型设计  
更加便捷的调试功能

```
import paddle.fluid as fluid
import numpy as np

class MyLayer(fluid.dygraph.Layer):
    def __init__(self, name_scope):
        super(MyLayer, self).__init__(name_scope)
        self.fc = fluid.dygraph.nn.FC(self.full_name(), size=12)

    def forward(self, inputs):
        x = self.fc(inputs)
        x = fluid.layers.relu(x)
        self._x_for_debug = x
        x = fluid.layers.elementwise_mul(x, x)
        x = fluid.layers.reduce_sum(x)
        return [x]
```

pytorch

## Eager execution

tf 2.0

高级 → CNN RNN

```
class MyModel(tf.keras.Model):
    def __init__(self):
        super(MyModel, self).__init__()
        self.conv1 = Conv2D(32, 3, activation='relu')
        self.flatten = Flatten()
        self.d1 = Dense(128, activation='relu')
        self.d2 = Dense(10, activation='softmax')

    def call(self, x):
        x = self.conv1(x)
        x = self.flatten(x)
        x = self.d1(x)
        return self.d2(x)
```

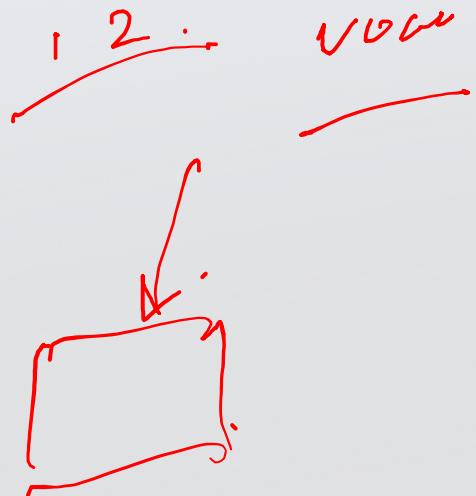
# Data

对比项	同步Feed方式 ✓	异步DataLoader接口方式 ✓
API接口	<code>executor.run(feed=...)</code>	<code>fluid.io.DataLoader</code>
数据格式	Numpy Array或LoDTensor	Numpy Array或LoDTensor
数据增强	Python端使用其他库完成	Python端使用其他库完成
速度	慢	快
推荐用途	调试模型	工业训练

```
import paddle.fluid as fluid
import numpy as np

# sample级reader
def fake_sample_reader():
    for _ in range(100):
        sample_image = np.random.random(size=(784, )).astype('float32')
        sample_label = np.random.random_integers(size=(1, ), low=0, high=9).astype('int64')
        yield sample_image, sample_label
~~~~~

# batch级reader
def fake_batch_reader():
    batch_size = 32
    for _ in range(100):
        batch_image = np.random.random(size=(batch_size, 784)).astype('float32')
        batch_label = np.random.random_integers(size=(batch_size, 1), low=0, high=9).astype('int64')
        yield batch_image, batch_label
```



# Data

```
import itertools
tf.compat.v1.enable_eager_execution()

def gen():
    for i in itertools.count(1):
        yield (i, [1] * i)

ds = tf.data.Dataset.from_generator(
    gen, (tf.int64, tf.int64), (tf.TensorShape([]), tf.TensorShape([None])))
```

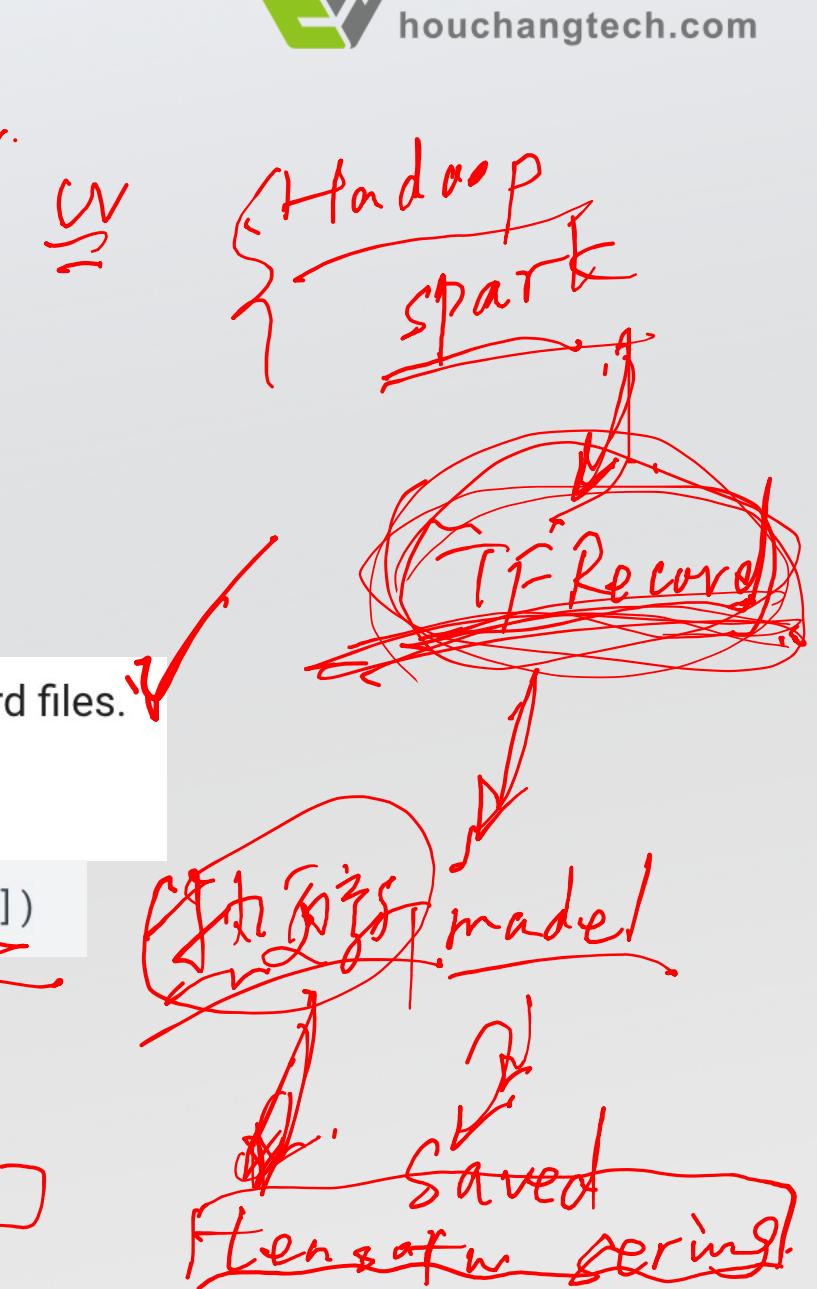
class TFRecordDataset: A Dataset comprising records from one or more TFRecord files.

class TextLineDataset: A Dataset comprising lines from one or more text files.

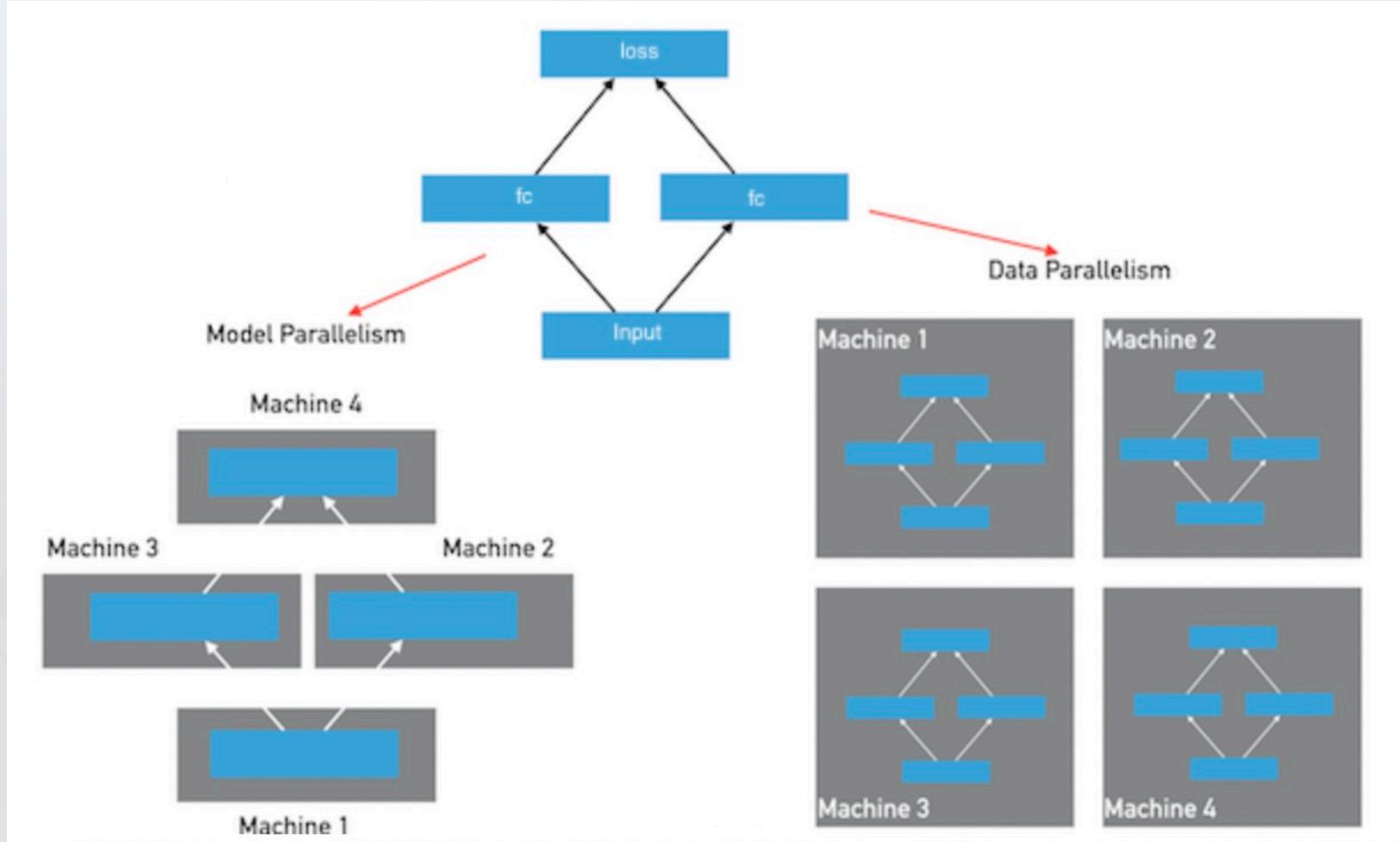
```
d = tf.data.Dataset.from_tensor_slices([1, 2, 3])
```

CSV  
txt.

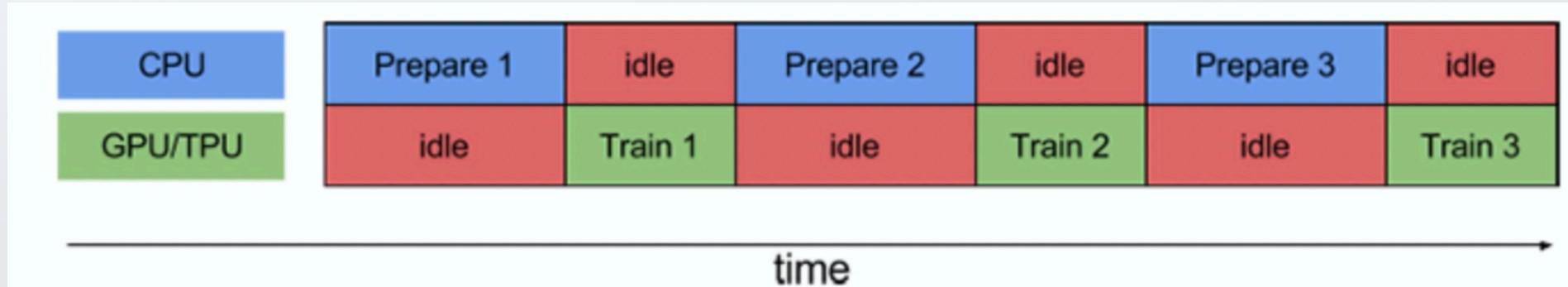
[ ] [ ] [ ] [ ]



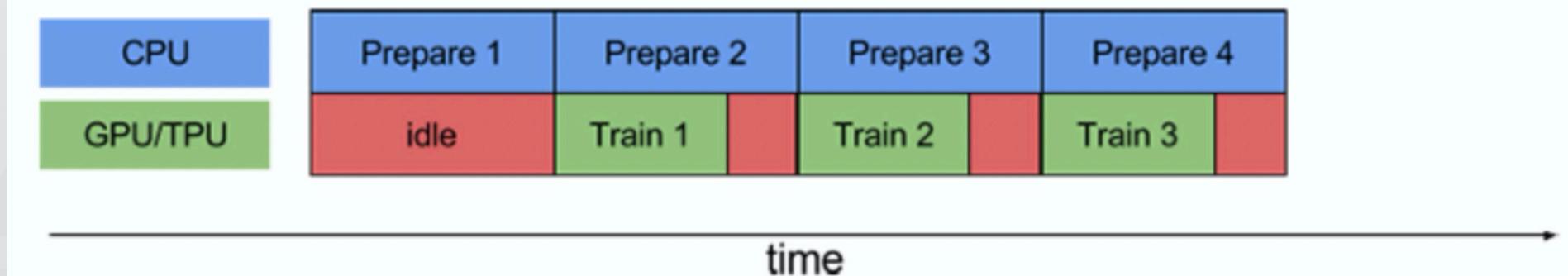
# 分布式训练



# 分布式训练



使用流水线可显著减少空闲时间：



### 3. Genism代码实践

10:20

## 4. 作业

Bye !