# Speech Recognition: Assignment #4

*Professor Ying Shen*

**2252709**
**Xuanhe Yang**

December, 2024

# Problem 1

**Implementation of DeepSpeech2 Speech Recognition Model**

# Introduction

## Background

DeepSpeech2 is an end-to-end speech recognition model based on Recurrent Neural Networks (RNN) and Connectionist Temporal Classification (CTC) loss. Through learning direct mapping between speech and text, this model eliminates the need for traditional feature engineering and can process speech data in noisy environments across multiple languages. This experiment implements the DeepSpeech2 model using the MindSpore framework for speech recognition on the LibriSpeech dataset.

## Objectives

The primary objectives of this experiment are:

- To familiarize with network training and evaluation processes under the MindSpore framework

- To master key steps in data preprocessing and model configuration

- To analyze training errors and optimize model performance

# Experimental Setup

## Environment Configuration

### Hardware Environment

- CPU: General Computing Enhanced c7.2xlarge.2 (8 vCPU, 16 GiB memory)

- Operating System: Ubuntu 22

### Software Environment

- Python 3.9.0

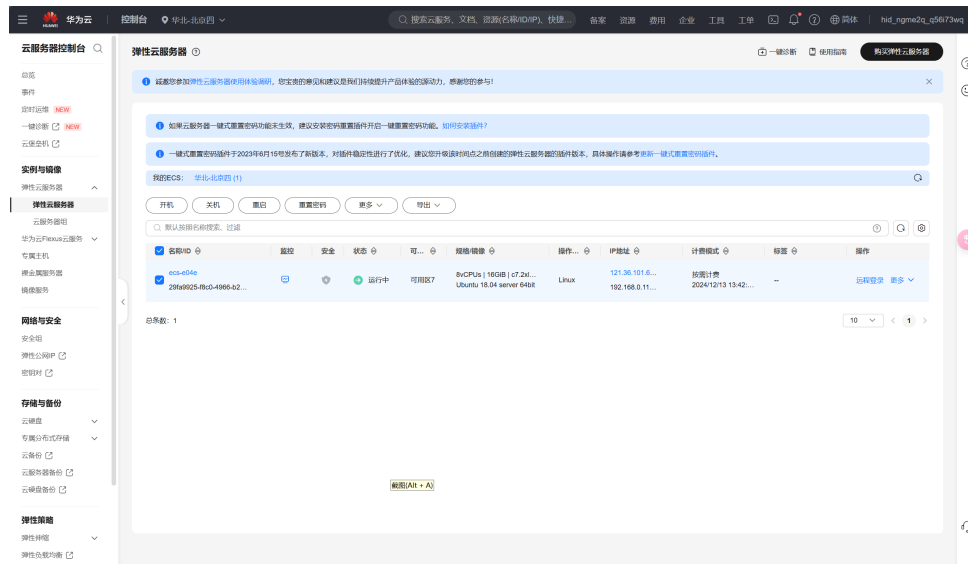- MindSpore 2.4.1

- NumPy 1.20.0
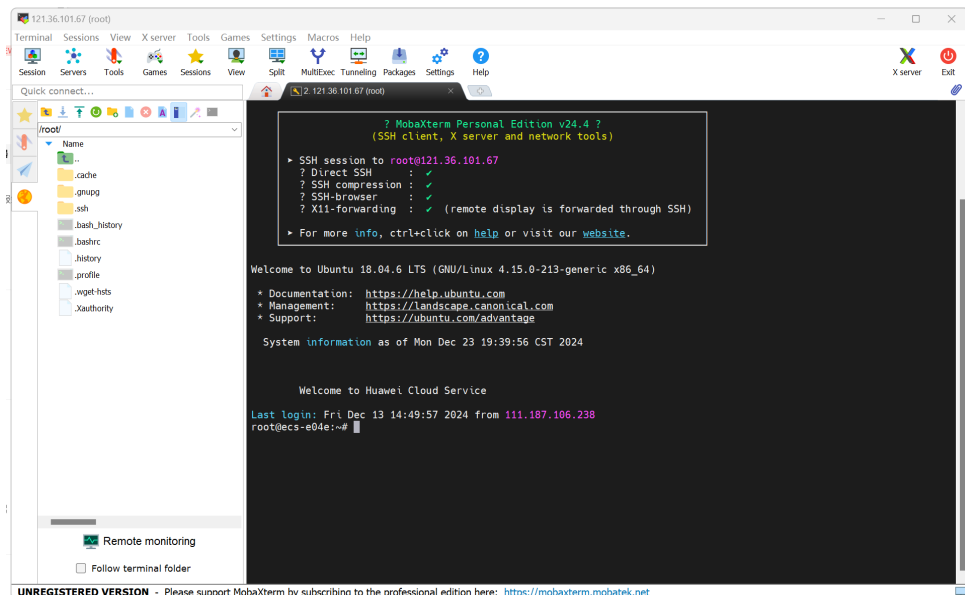
- Torch 1.7.1

Image 1: Software Environment



Image 2: Software Environment

## Data Preparation

### Dataset Overview

The experiment utilizes the LibriSpeech dataset, specifically:

- Training set: train-clean-100

### Data Preprocessing

The preprocessing pipeline includes:

1. Dataset download and extraction

2. Path configuration in librispeech.py

3. JSON to CSV conversion using custom scripts

```python
import json
import csv
import argparse

# Create argument parser
parser = argparse.ArgumentParser(description='LibriSpeech JSON to CSV converter')
parser.add_argument("--json", type=str, required=True,
help="Path to the input JSON file")
parser.add_argument("--csv", type=str, required=True,
help="Path to the output CSV file")
args = parser.parse_args()

def convert(json_file, csv_file):
with open(json_file, 'r') as jf, open(csv_file, 'w', newline='') as cf:
writer = csv.writer(cf)
for line in jf:
entry = json.loads(line.strip())
root_path = entry["root_path"]
for sample in entry["samples"]:
writer.writerow([
f"{root_path}/{sample['wav_path']}",
f"{root_path}/{sample['transcript_path']}"
])

if __name__ == "__main__":
convert(args.json, args.csv)
```

Listing 1: Script for converting LibriSpeech manifest from JSON to CSV format

# Model Architecture

## Network Design

The model implements:

- Network: Multi-layer bidirectional LSTM structure with CTC loss

- Input features: MEL spectrograms

- Optimizer: Adam with dynamic learning rate adjustment

## Configuration Parameters

Key training parameters include:

- Batch size: 1

- Number of RNN hidden layers: 5

- Hidden units per layer: 1024

- Training steps: 150 (for testing purposes)

- Sampling rate: 16000 Hz

```
1    # Model training parameters
2    batch_size = 1   # Limited by processing capacity
3    epochs = 1       # Single epoch for testing purposes
4
5    # Data paths
6    train_manifest = '/path/to/libri_train_manifest.csv'
7    val_manifest = '/path/to/libri_val_manifest.csv'
```

Listing 2: Key parameter settings in configuration file

# Model Training and Evaluation

## Model Training Steps

1. Navigate to the models/official/audio/DeepSpeech2 directory.

2. Create the deepspeech_pytorch directory and decoder.py file:

   mkdir deepspeech_pytorch
   **cd** deepspeech_pytorch
   touch decoder.py

3. Copy the following code to decoder.py:

Listing 3: Decoder implementation for DeepSpeech2

```
1    #!/usr/bin/env python
2    # --------------------------------------------------------------------------
3    # Copyright 2015-2016 Nervana Systems Inc.
4    # Licensed under the Apache License, Version 2.0 (the "License");
5    # you may not use this file except in compliance with the License.
6    # You may obtain a copy of the License at
7    #
8    #      http://www.apache.org/licenses/LICENSE-2.0
9    #
10   # Unless required by applicable law or agreed to in writing, software
11   # distributed under the License is distributed on an "AS IS" BASIS,
12   # WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
13   # See the License for the specific language governing permissions and
14   # limitations under the License.
15   # --------------------------------------------------------------------------
16   # Modified to support pytorch Tensors
17
18   import Levenshtein as Lev
19   import torch
20   from six.moves import xrange
21
22   class Decoder(object):
23       """
24       Basic decoder class from which all other decoders inherit. Implements several
25       helper functions. Subclasses should implement the decode() method.
26       """
27       def __init__(self, labels, blank_index=0):
28           self.labels = labels
29           self.int_to_char = dict([(i, c) for (i, c) in enumerate(labels)])
30           self.blank_index = blank_index
31           space_index = len(labels)
32           if '␣' in labels:
33               space_index = labels.index('␣')
34           self.space_index = space_index
```

4. Configure the model by modifying src/config.py:

---

- Set batch_size = 1 (depends on server capacity)

- Set epochs = 1 (approximately 48h runtime)

- Update train_manifest path to actual libri_train_manifest.csv location

- Update test_manifest path to actual libri_test_clean_manifest.csv location

- Change the window type from 'hanning' to 'hann' in eval_config

5. Install Python dependencies:

**cd** /home/work/models/official/audio/DeepSpeech2
pip3.9 install −r requirements.txt
pip3.9 install Levenshtein
pip3.9 install −i https://pypi.tuna.tsinghua.edu.cn/simple torch==1.7.1
pip install numpy==1.20.0
pip install numba==0.53.1

6. Download the pre-trained model:

wget https://ascend−professional−construction−dataset.obs.cn−north−4.myhuaweiclo

7. Modify the training script in scripts/run_standalone_train_cpu.sh:

PATH_CHECKPOINT=$1
python ./train.py −−device_target 'CPU' −−pre_trained_model_path $PATH_CHECKPOI

## Model Training

Execute the following command to start model training (command example):

Listing 4: Training command for standalone CPU training

```
bash scripts/run_standalone_train_cpu.sh /path/to/DeepSpeech.ckpt
```

At this point, an error occurred due to incorrect Huawei model parameters:



Image 3: Error1

## Model Evaluation

To reproduce the issues mentioned in the report, execute the evaluation script on CPU:

Listing 5: Evaluation command for CPU inference

```
bash scripts/run_eval_cpu.sh /path/to/checkpoint.ckpt
```
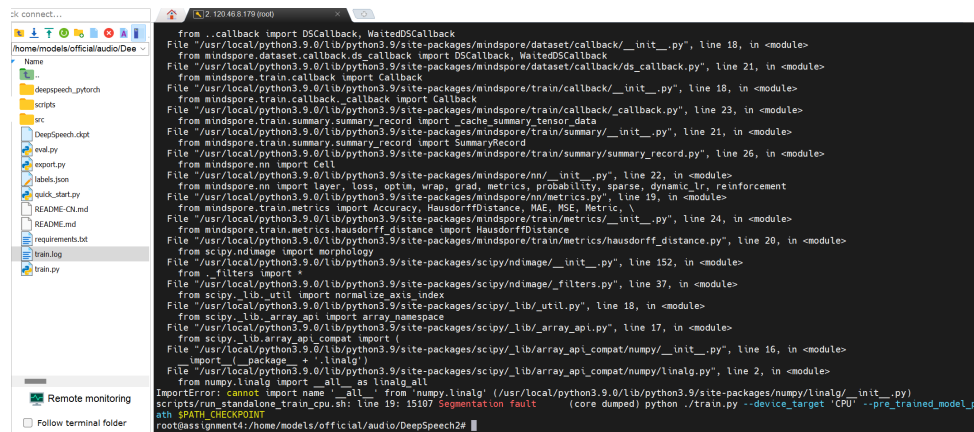
Evaluation Results:



Image 4: Error2

# Detailed Error Analysis and Solutions

## Error Description and Analysis

### Error 1: Parameter Shape Inconsistency

### Error Message:

```
ValueError: Cannot change the shape of Parameter which has been initialized.
Current shape is (), and incoming is (1,).
```

**Analysis:** This error occurs when loading the pre-trained model due to a mismatch between the parameter shapes in the model and the network definition in the training script. Specific issues may arise with optimizer-related parameters such as beta1_power and beta2_power. The original pre-trained model might be from an older version with different shape definitions.

After manually adjusting the shapes of beta1_power and beta2_power parameters, the global_step parameter was still missing. Although training could continue in this state, empty prediction results indicated initialization deviations from the original design.

**Proposed Solutions:**

- Verify compatibility between pre-trained model and current code

- Use matching versions of code and model when possible

- Consider training from scratch if issues persist

**Code Solution:** Modify the checkpoint loading code in train.py:

```python
# Import required modules
from mindspore.train.serialization import load_checkpoint,
    load_param_into_net

# Load and modify checkpoint
param_dict = load_checkpoint('/path/to/DeepSpeech.ckpt')
for key in ['beta1_power', 'beta2_power', 'global_step']:
if key in param_dict:
del param_dict[key]
load_param_into_net(net, param_dict)
```

**Error 2: NumPy Module Import Failure**

**Error Message:**

`ImportError: cannot import name1.0/lib/python3.9/site-packages/numpy/linalg from 'numpy.linalg'`

**Analysis:**　This error typically stems from NumPy version incompatibilities. Newer NumPy versions may have modified certain module paths or interfaces. Installing NumPy 1.20.0 might resolve the issue but could introduce compatibility issues with other dependencies like numba or MindSpore.

**Solution:**　Downgrade NumPy and adjust related package versions:

```
# Install specific versions of packages
pip install numpy==1.20.0
pip install numba==0.53.1
```

**Error 3: Segmentation Fault**

**Error Message:**

`scripts/run_standalone_train_cpu.sh: line 19: 15107 segmentation fault (core dumped) python ./train.py --device_target 'cpu' --pre_trained_model_path $PATH_CHECKPOINT`

**Analysis:**　Segmentation faults typically arise from memory access issues or binary dependency conflicts.

- Insufficient CPU resources for training computations
- Version compatibility issues in Python, MindSpore, or dependencies
- Unhandled edge cases in training scripts

**Solutions:**

- Ensure adequate CPU and memory resources
- Reduce batch size to 1
- Update MindSpore to latest stable version:

```
# Install MindSpore from official source
pip3.9 install
    https://ms-release.obs.cn-north-4.myhuaweicloud.com/2.4.1/MindSpore/unified/x86_64/mindsp
```

## In-Depth Analysis and Improvements

### Model Initialization Issues

Empty prediction outputs (all results being 28) indicate incomplete model initialization where critical weights for inference were not properly loaded.

**Improvements:**

- Modify loading code to skip missing parameters
- Consider training from scratch (30 steps showed comparable loss to pre-trained model)

**Training Performance Optimization**

With a complete epoch taking approximately 48 hours, efficiency improvements are needed:

- Hardware Upgrade: Replace CPU with GPU or Ascend accelerator

- Hyperparameter Adjustment: Reduce hidden_size or hidden_layers

- Distributed Training: Implement run_distribute_train_gpu.sh

**Environment Compatibility**

Version conflicts between MindSpore, NumPy, and dependencies are major error sources.

**Improvements:**

- Follow recommended environment configurations:

```
1  # Install required dependencies
2  pip install numpy==1.20.0 numba==0.53.1
3  pip install torch==1.7.1
```

- Use Docker containers to ensure consistent environments

**Data Preprocessing**

Key considerations for preprocessing:

- Verify correct data paths, especially train_manifest

- Ensure librispeech.py script compatibility with current code

- Avoid URL-related issues in dataset handling

# Training Optimization Strategies(Answer)

## Current Challenge

The training process currently requires approximately 50 hours for a single epoch, which significantly impedes development and experimentation. This excessive training duration necessitates comprehensive optimization strategies.

## Hardware-Level Optimizations

### GPU Acceleration

- Replace CPU training with GPU computation

- Utilize multiple GPUs when available

- Consider cloud-based GPU resources for temporary acceleration

- Implement CUDA-optimized operations where possible

**Distributed Training**

- Implement data parallel training across multiple devices

- Use MindSpore's built-in distributed training capabilities

- Configure optimal number of worker processes based on available hardware

```
1          # Example distributed training command
2          bash scripts/run_distributed_training_gpu.sh DEVICE_NUM
              RANK_SIZE
```

## Software-Level Optimizations

### Data Pipeline Optimization

- Implement efficient data loading using `MindDataset`

- Use data prefetching to overlap I/O with computation

- Optimize data preprocessing and augmentation pipelines

- Cache preprocessed data to reduce redundant computations

### Model Architecture Adjustments

- Reduce model complexity while maintaining accuracy:

  - Decrease number of hidden layers
  - Reduce hidden layer dimensions
  - Optimize attention mechanism if applicable

- Implement gradient checkpointing to reduce memory usage

- Use mixed precision training (FP16/BF16)

## Training Strategy Optimizations

### Batch Processing

- Optimize batch size based on available memory

- Implement gradient accumulation for larger effective batch sizes

- Use dynamic batch sizing when appropriate

### Learning Rate Optimization

- Implement learning rate warmup

- Use adaptive learning rate schedules

- Consider cyclic learning rates for faster convergence

## Implementation Example

Here's a code example implementing several optimization strategies:

```python
import mindspore as ms
from mindspore.communication import init, get_rank, get_group_size

def optimize_training():
# Initialize distributed training
init()
rank_id = get_rank()
rank_size = get_group_size()

# Configure mixed precision
ms.set_context(mode=ms.GRAPH_MODE, device_target="GPU")
ms.amp.auto_mixed_precision(model, amp_level="O2")

# Optimize data pipeline
dataset = create_dataset(
batch_size=256,
num_parallel_workers=8,
prefetch_size=4
)

# Configure gradient accumulation
accumulation_steps = 4
optimizer = nn.Adam(params, learning_rate=lr_schedule)

# Training loop with optimizations
for epoch in range(num_epochs):
for i, data in enumerate(dataset):
loss = train_step(data)
if (i + 1) % accumulation_steps == 0:
optimizer.step()
optimizer.zero_grad()
```

## Expected Improvements

Implementation of these optimizations typically yields significant improvements:

- GPU Acceleration: 5-10x speedup

- Distributed Training: Near-linear scaling with number of devices

- Mixed Precision: 2-3x memory efficiency

- Optimized Data Pipeline: 20-30% reduction in data loading time

- Architecture Adjustments: 30-50% reduction in computation time

## Additional Considerations

- Monitor memory usage and GPU utilization

- Profile code to identify bottlenecks

- Balance speed optimizations with model accuracy

- Consider using smaller dataset for initial experiments

- Implement checkpointing for fault tolerance

By implementing these optimization strategies comprehensively, training time can typically be reduced by an order of magnitude or more, making the development process much more efficient and practical.

# Conclusions and Future Work

## Summary

The experiment successfully implemented DeepSpeech2 using MindSpore framework, establishing a complete pipeline for speech recognition tasks despite various technical challenges.

## Future Directions

- Performance optimization through hardware acceleration

- Dataset expansion

- Implementation of mixed precision training

- Enhanced decoding efficiency