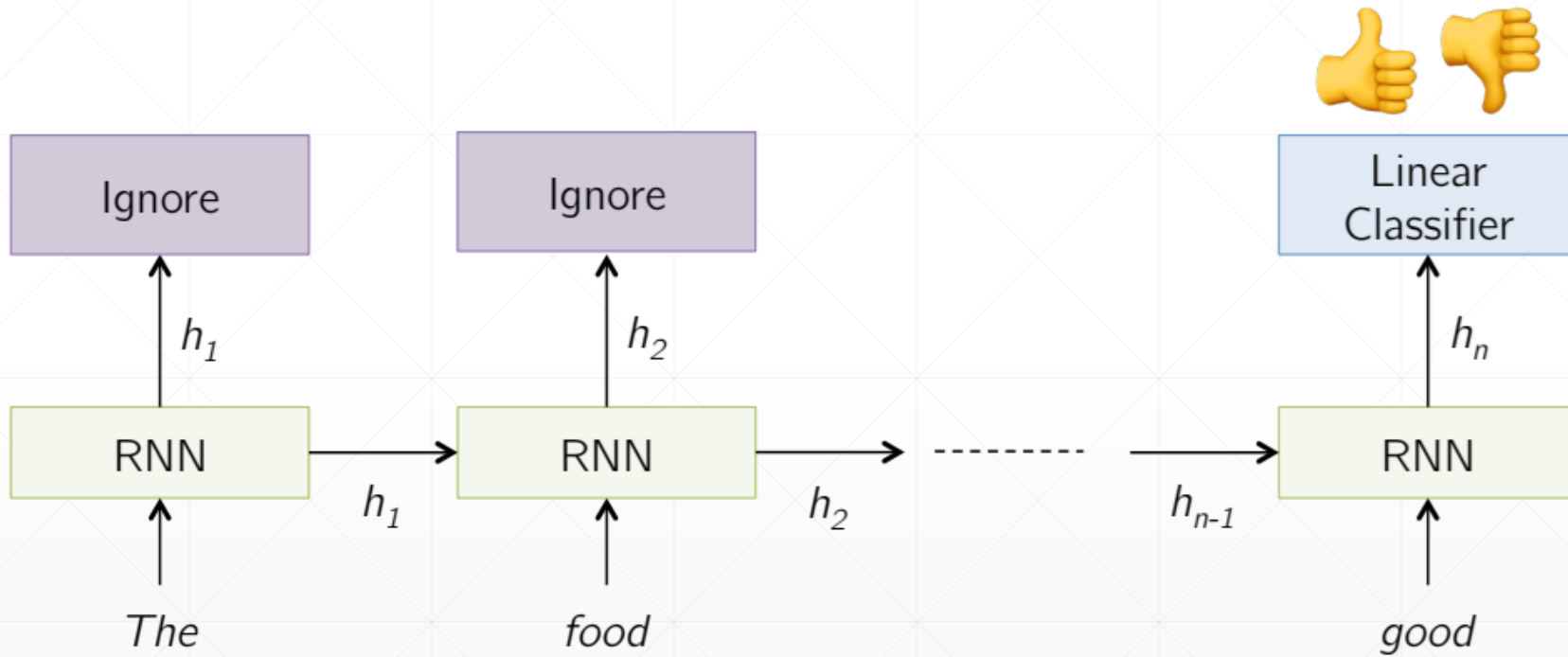


# 情感分类实战

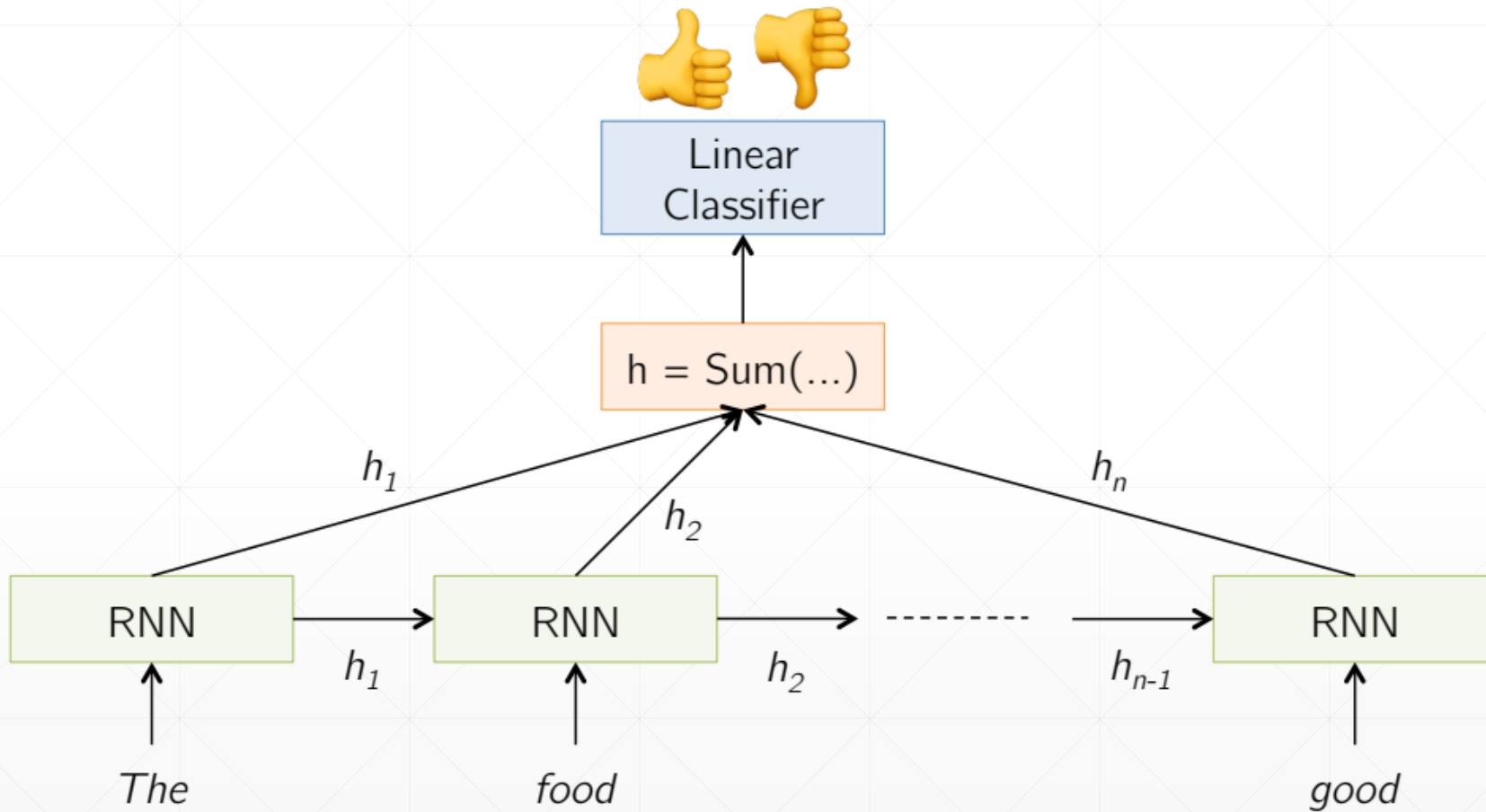
---

主讲人：龙良曲

# Sentiment Classification



# Sentiment Classification



# Google CoLab

- Continuous 12 hours
  - free K80 for GPU
  - no need to cross GFW
-

# Load Dataset

```
1 TEXT = data.Field(tokenize='spacy')
2 LABEL = data.LabelField(dtype=torch.float)
3 train_data, test_data = datasets.IMDB.splits(TEXT, LABEL)
4
5 print('len of train data:', len(train_data))
6 print('len of test data:', len(test_data))
7 len of train data: 25000
8 len of test data: 25000
9
10 print(train_data.examples[15].text)
11 print(train_data.examples[15].label)
12 ['I', 'loved', 'this', 'film', '.', 'I', 'thought', 'it', 'would', 'be', 'easy', 'to', 'watch',
    ',', 'and', 'easy', 'to', 'forget', '.', 'I', 'ran', 'out', 'after', 'watching', 'this', 'to',
    'buy', 'the', 'DVD', ',', 'obv', 'not', 'easily', 'forgotten!<br', '/><br', '/>The', 'script',
    'is', 'brilliant', ',', 'and', 'the', 'casting', 'could', "n't", 'be', 'more', 'perfect', '.',
    'Each', 'character', 'has', 'their', 'moment', ',', 'and', 'I', 'laughed', 'hard', 'throughout',
    'this', 'film', ',', 'comedic', 'timing', 'was', 'spot', '-', 'on.<br', '/><br', '/', '>']
13 pos
```

# Network

```
1 class RNN(nn.Module):
2     def __init__(self, vocab_size, embedding_dim, hidden_dim):
3         super(RNN, self).__init__()
4         # [0-10001] => [100]
5         self.embedding = nn.Embedding(vocab_size, embedding_dim)
6         # [100] => [256]
7         self.rnn = nn.LSTM(embedding_dim, hidden_dim, num_layers=2,
8                             bidirectional=True, dropout=0.5)
9         # [256*2] => [1]
10        self.fc = nn.Linear(hidden_dim*2, 1)
11        self.dropout = nn.Dropout(0.5)
12    def forward(self, x):
13        # [seq, b, 1] => [seq, b, 100]
14        embedding = self.dropout(self.embedding(x))
15        # output: [seq, b, hid_dim*2]
16        # hidden/h: [num_layers*2, b, hid_dim]
17        # cell/c: [num_layers*2, b, hid_dim]
18        output, (hidden, cell) = self.rnn(embedding)
19        # [num_layers*2, b, hid_dim] => 2 of [b, hid_dim] => [b, hid_dim*2]
20        hidden = torch.cat([hidden[-2], hidden[-1]], dim=1)
21        # [b, hid_dim*2] => [b, 1]
22        hidden = self.dropout(hidden)
23        out = self.fc(hidden)
24        return out
```

# Load word embedding

```
1 rnn = RNN(len(TEXT.vocab), 100, 256)
2
3 pretrained_embedding = TEXT.vocab.vectors
4 print('pretrained_embedding:', pretrained_embedding.shape)
5 rnn.embedding.weight.data.copy_(pretrained_embedding)
6 print('embedding layer inited.')
```

*Word Vector  
Lookup Table!*

*300 features*

*10,000 words*



# Train

```
1 def train(rnn, iterator, optimizer, criteon):
2     avg_acc = []
3     rnn.train()
4
5     for i, batch in enumerate(iterator):
6         # [seq, b] => [b, 1] => [b]
7         pred = rnn(batch.text).squeeze(1)
8         loss = criteon(pred, batch.label)
9         acc = binary_acc(pred, batch.label).item()
10        avg_acc.append(acc)
11
12        optimizer.zero_grad()
13        loss.backward()
14        optimizer.step()
```



# Test

```
1 def binary_acc(preds, y):
2     preds = torch.round(torch.sigmoid(preds))
3     correct = torch.eq(preds, y).float()
4     acc = correct.sum() / len(correct)
5     return acc
6
7 def eval(rnn, iterator, criteon):
8     avg_acc = []
9     rnn.eval()
10    with torch.no_grad():
11        for batch in iterator:
12            # [b, 1] => [b]
13            pred = rnn(batch.text).squeeze(1)
14            loss = criteon(pred, batch.label)
15            acc = binary_acc(pred, batch.label).item()
16            avg_acc.append(acc)
17    avg_acc = np.array(avg_acc).mean()
18    print('>>test:', avg_acc)
```

# 下一课时

---

GAN

**Thank You.**

---