



CS221 Final Report

Xử lý ngôn ngữ tự nhiên (Trường Đại học Công nghệ thông tin, Đại học Quốc gia Thành phố Hồ Chí Minh)



Scan to open on Studocu

ĐẠI HỌC QUỐC GIA TP.HCM
TRƯỜNG ĐẠI HỌC CÔNG NGHỆ THÔNG TIN



KHOA KHOA HỌC MÁY TÍNH

MÔN HỌC: XỬ LÝ NGÔN NGỮ TỰ NHIÊN

**Tiêu đề: Efficient Document
Retrieval by End-to-End
Refining and Quantizing BERT
Embedding with Contrastive
Product Quantization**

Nhóm 7

Giảng viên:
Ts. Nguyễn Thị Quý

1 Giới thiệu

1.1 Bài toán truy vấn tài liệu (Document Retrieval)

Truy xuất tài liệu được định nghĩa là việc khớp một số truy vấn của người dùng đã nêu với một tập hợp các bản ghi văn bản tự do. Những bản ghi này có thể là bất kỳ loại văn bản phi cấu trúc nào, chẳng hạn như các bài báo, hồ sơ bất động sản hoặc các đoạn văn trong sổ tay. Truy vấn của người dùng có thể bao gồm từ mô tả đầy đủ nhiều câu về nhu cầu thông tin cho đến một vài từ.

1.2 Dataset

NYT là dataset chứa những bài báo được xuất bản bởi The New York Times. Dataset chứa 3 tập train, val, test trong đó tập train chứa 9221 bài báo, tập val chứa 1152 bài báo và tập test chứa 1154 bài báo.

Dataset là file **csv** với 2 cột: cột thứ 1 là label (category của bài báo) và cột thứ 2 là nội dung bài báo.

Các bài báo được chia vào 26 categories: dance, immigration, gun_control, abortion, economy, tennis, football, basketball, movies, federal_budget, international_business, military, stocks_and_bonds, law_enforcement, golf, baseball, environment, energy_companies, hockey, television, cosmos, soccer, surveillance, the_affordable_care_act, music, gay_rights.

index	label	content
0	22	manchester — manchester united striker wayne Rooney is in good shape mentally and physically as the premier league champions prepare to face bitter rivals Liverpool this weekend. he was with me at everton as a young boy and his career has gone exactly the way that people expected it to go. Chelsea's public pursuit of Rooney has cooled off since Monday night, but that's average, and i don't look around this locker room and see average inconsistency has been our issue, and here it is again showing its ugly head today. the dolphins and we didn't come up with the plays we needed, Wallace said. Tammehill had passer ratings of 6.6 and 4.1 in the final two games.
1	7	miami gardens, Fla. — rex Ryan is coming back next year. the Miami Dolphins are sailing out another postseason. the New York Jets eliminated the Dolphins from playoff contention.
2	22	berlin — Germany's storied Bayen Munich soccer club is sticking by its president, despite a state court's announcement on Monday that he will stand trial on charges of tax evasion. the team to build its financial resources to a level rarely seen in European club soccer. his decision to go public with his story triggered a national discussion over social justice and
3	8	waco, Texas — odyssey Sims scored 29 points and had 13 assists to become the Baylor career leader in that category as the No. 9 Lady Bears beat San Jose State 113-73 on Tues night after the Baylor and Kentucky women finish their game at the site of the men's NCAA final four next April, the men's teams from the
4	7	east Rutherford, N.J. — the talk all week by both teams was of desperation, a must-win showdown between the Miami Dolphins and New York Jets one of the AFC East rivals brought round pick to sit? or, should the jets let him try to work himself out of his funk? with the playoffs an unlikely scenario, New York might opt to find out if Smith could be the quarterback the jets couldn't get anything going all game. leaky defense: the jets' defense couldn't get off the field in the first half, and things weren't much better in the last two quarters. they
5	7	the NFL loves to trumpet its calendar, from draft day to kickoff weekend to Super Bowl Sunday, then there is Black Monday, the day after the regular season ends, when a dozen competitive environment, there's a race to get things done because there is the same group of candidates out there, said Mike Lannenbaum, the former general manager of the Jets. he ran an article that said that the day after the regular season ends traditionally is called 'Black Monday' in the coaching profession, though when this so-called tradition started was
6	14	Centennial, Colo. — the ghostly self-portraits, taken with a cellphone hours before he opened fire at an Aurora movie theater, show James E. Holmes in the final sprint of his death. there were never going to be enough ambulances, enough police officers to get everyone out of there fast enough, she added. he didn't care who he killed and how many he killed.
7	5	washington — Americans spent more in May as their income rose, encouraging signs after a slow start to the year. the Commerce Department said on Thursday that consumer spending rose 0.4 percent in May, up from a revised 0.2 percent increase in April.
8	16	Corey Kluber pitched eight shutout innings, and the Cleveland Indians spoiled Stephen Strasburg's return with a 2-0 win over the visiting Washington Nationals on Sunday. Kluber es visiting Milwaukee. marlins, Cardinals 2. Ricky Nolasco allowed one run and three hits in seven innings, and host Miami won to take two of three games from St. Louis. the Cardinals
9	1	the poetically named Shantala Shivalingappa was born in India, but raised in Paris, where she trained in Kuchipudi, a South Indian classical dance style, with her mother, the well-known dancer and choreographer Meena Bhagat. She began her training at the age of five under the tutelage of her mother, who was a student of the famous dancer Uday Shankar. She completed her training in Paris and then moved to the United States, where she performed in various theaters and dance studios. She has also performed in India and abroad, including in London and New York City. She has won several awards for her performances, including the National Award for Best Dancer in 1985. She has also written a book on her life and experiences as a dancer, titled "Kuchipudi: The Art of Indian Classical Dance".

Figure 1: Minh họa dataset

1.3 Độ đo đánh giá

Với mỗi query từ tập test, tìm trong tập train top 100 document phù hợp nhất với query. Sau đó, độ chính xác truy xuất được tính bằng tỷ lệ của các tài liệu liên quan. Các tài liệu được tính là liên quan nếu có cùng label với label của query.

2 Quy trình thực hiện

2.1 Install and import dependencies

Model yêu cầu 2 thư viện chính là **torch** và **transformers**.

```
1 pip install -r requirements.txt
```

Ngoài ra ta cần import một số thư viện khác:

```
1 import os
2 import argparse
3 import re
4 import pandas as pd
5 import numpy as np
```

2.2 Tạo Data Loader và xử lý dữ liệu

Trước tiên, ta cần đọc dữ liệu từ file .csv, tiền xử lý dữ liệu, sau đó tạo Data loader để load dữ liệu theo batch size.

```
1 class LabeledDocuments(Data):
2     def __init__(self, hparams):
3         super().__init__(hparams=hparams)
4
5     def load_datasets(self):
6         self.X_train, self.Y_train = self.read_data(
7             f"{self.file_path}/train.csv", max_length=self.
max_length)
8         self.X_val, self.Y_val = self.read_data(
9             f"{self.file_path}/val.csv", max_length=self.
max_length)
10        self.X_test, self.Y_test = self.read_data(
11            f"{self.file_path}/test.csv", max_length=self.
max_length)
12
13        print(
```

```

14         f"train size: {len(self.X_train)}, validation
size: {len(self.X_val)}, test size: {len(self.X_test)}")
15
16     def read_data(self, path, max_length):
17         def label_fn(x):
18             return x - 1
19
20         rows = pd.read_csv(
21             path,
22             sep=",",
23             on_bad_lines='skip',
24             header=None,
25             skiprows=None,
26             quoting=0,
27             keep_default_na=False,
28             encoding="utf-8",
29         )
30
31         label_fn = label_fn if label_fn is not None else (
32             lambda x: x)
33         labels = rows[0].apply(lambda x: label_fn(x))
34         sentences = rows[1]
35         sentences = sentences.apply(
36             lambda x: clean_tokenize_truncate(x, max_length))
37         return sentences.tolist(), labels.tolist()
38
39     def get_loaders(self, batch_size, num_workers,
40                     shuffle_train=False,
41                     get_test=True):
42         train_dataset = My_Dataset(
43             self.X_train,
44             self.Y_train,
45             self.max_length)
46         database_dataset = My_Dataset(
47             self.X_train, self.Y_train, self.max_length)
48         val_dataset = My_Dataset(self.X_val, self.Y_val, self
49 .max_length)
50         test_dataset = My_Dataset(self.X_test, self.Y_test,
51 self.max_length)
52
53         # DataLoader
54         train_loader = DataLoader(
55             dataset=train_dataset,
56             batch_size=batch_size,
57             shuffle=shuffle_train,

```

```

55                                         num_workers=num_workers)

56
57     database_loader = DataLoader(
58         dataset=database_dataset,
59         batch_size=512,
60         shuffle=False,
61         num_workers=num_workers)
62
63     val_loader = DataLoader(
64         dataset=val_dataset,
65         batch_size=512,
66         shuffle=False,
67         num_workers=num_workers)
68
69     test_loader = DataLoader(
70         dataset=test_dataset,
71         batch_size=512,
72         shuffle=False,
73         num_workers=num_workers) if
74     get_test else None
75     return train_loader, database_loader, val_loader,
76     test_loader,

```

Một số hàm hỗ trợ tiền xử lý dữ liệu.

```

1 def clean_string(string):
2 """
3     Tokenization/string cleaning for yelp data set
4     Based on https://github.com/yoonkim/CNN_sentence/blob/
5     master/process_data.py
6 """
7     string = re.sub(r"[^A-Za-z0-9(),.?\\\'\\\"\\]", " ", string)
8     string = re.sub(r"\'s", " 's", string)
9     string = re.sub(r"\'ve", " 've", string)
10    string = re.sub(r"\n\\'t", " n't", string)
11    string = re.sub(r"\'re", " 're", string)
12    string = re.sub(r"\\'d", " 'd", string)
13    string = re.sub(r"\\'ll", " 'll", string)
14    string = re.sub(r"\,", " , ", string)
15    string = re.sub(r"\!", " ! ", string)
16    string = re.sub(r"\(\", " \(", string)
17    string = re.sub(r"\\"", ' \" ', string)
18    string = re.sub(r"\)\", " \) ", string)
19    string = re.sub(r"\?", " \? ", string)
20    string = re.sub(r"\s{2,}", " ", string)
21    return string.lower().strip()

```

```

1 def clean_tokenize_truncate(x, max_length):
2     x = clean_string(x)
3     x = x.split(" ")
4     x = x[:max_length]
5     return x

```

2.3 Xây dựng model

- Xây dựng Product Quantization:

```

1 class PQ_head(nn.Module):
2     """
3         The pq implementaion refers to SPQ(https://github.com/youngkyunJang/SPQ).
4     """
5     def __init__(self, N_words, N_books, L_word, gumbel_temp,
6                  dist_metric, sample_method = "gumbel_softmax"):
7         # N_words = K; number of codewords
8         # N_books = M; number of codebooks
9         # L_word = Dimension of the codewords
10
11         super(PQ_head, self).__init__()
12
13         self.C = nn.Parameter(torch.randn(N_words, N_books * L_word))
14
15         self.N_books = N_books
16         self.L_word = L_word
17         self.gumbel_temp = gumbel_temp
18         self.dist_metric = dist_metric
19
20         self.sample_method = sample_method
21
22     def sample_gumbel(self, shape, eps = 1e-20):
23         U = torch.rand(shape)
24         U = U.cuda()
25         return -torch.log(-torch.log(U + eps) + eps)
26
27     def gumbel_softmax_sample(self,logits, temperature = 1):
28         y = logits + self.sample_gumbel(logits.size())
29         return F.softmax(y / temperature, dim = -1)
30
31     def gumbel_softmax(self, logits, temperature = 1, hard = True):
32         y = self.gumbel_softmax_sample(logits, temperature)

```

```

32
33     if not hard:
34         return y
35     raise NotImplementedError()
36
37
38     def defined_sim(self, x, c):
39         if self.dist_metric == "euclidean":
40             diff = x.unsqueeze(1) - c.unsqueeze(0)
41             return - torch.sum(diff * diff, -1)
42         else:
43             raise NotImplementedError()
44
45
46
47     def forward(self, X):
48         # input shape: (bsz, N_books * L_word)
49         x = torch.split(X, self.L_word, dim=1) # tuple; ele
        of the tuple has the shape like (batch_size, L_word)
N_books element
50         c = torch.split(self.C, self.L_word, dim=1) # tuple
; ele of the tuple has the shape like (N_words, L_word)
N_books elements.
51         prob_list = []
52         for i in range(self.N_books):
53             logits = self.defined_sim(x[i], c[i])
54
55             if self.sample_method == "softmax":
56                 soft_prob = F.softmax(logits / self.
gumbel_temp, dim=1)
57                 prob_list.append(soft_prob)
58             elif self.sample_method == "gumbel_softmax":
59                 logits = F.softmax(logits, dim=1)
60                 logits = torch.log(logits + 1e-9)
61                 soft_prob = self.gumbel_softmax(logits, self.
gumbel_temp, False)
62                 prob_list.append(soft_prob)
63
64             if i == 0:
65                 Q = torch.mm(soft_prob, c[i])
66             else:
67                 Q = torch.cat((Q, torch.mm(soft_prob, c[i])), dim=1)
68         return Q, prob_list

```

- Xây dựng Mutual Information Improved Contrastive Product Quantization

(MICPQ) model:

+ Khởi tạo tham số:

```
1 def define_parameters(self):
2     bert_config = AutoConfig.from_pretrained('bert-base-
3         uncased')
4     bert_config.attention_probs_dropout_prob = self.hparams.
5         dropout_rate
6     bert_config.hidden_dropout_prob = self.hparams.
7         dropout_rate
8     self.bert_model = AutoModel.from_pretrained('bert-base-
9         uncased', config = bert_config)
10    for param in self.bert_model.parameters():
11        param.requires_grad = False
12
13    assert self.hparams.encode_length == self.hparams.N_books
14        * int(np.log2(self.hparams.N_words))
15
16    self.pq_head = PQ_head(
17        N_words = self.hparams.N_words,
18        N_books = self.hparams.N_books,
19        L_word = self.hparams.L_word,
20        gumbel_temp = self.hparams.gumbel_temperature,
21        dist_metric = self.hparams.dist_metric,
22        sample_method = self.hparams.sample_method
23        )
24
25
26    self.criterion = NtXentLoss(self.hparams.temperature)
```

+ Tinh chỉnh BERT embedding:

```
1 def get_embeddings(self, inputs, pooling="mean"):
2     bert_output = self.bert_model(**inputs)[0] # (batch_size,
3         sequence_length, hidden_size) Sequence of hidden-states
4         at the output of the last layer of the model.
5
6     if pooling == 'mean':
7         attention_mask = inputs['attention_mask'].unsqueeze
8             (-1) # (batch_size, sequence_length, 1)
```

```

6         output = torch.sum(bert_output*attention_mask, dim =
7             1) / torch.sum(attention_mask, dim = 1)
8     elif pooling == 'cls':
9         output = bert_output[:,0]
10    else:
11        raise NotImplementedError()
12    return output

```

+ Hàm Mutual Information Loss:

```

1 def compute_prob_loss(self, prob0, prob1):
2     def get_entropy(probs):
3         q_ent = -(probs.mean(0) * torch.log(probs.mean(0)
4 + 1e-12)).sum()
5         return q_ent
6
7     def get_cond_entropy(probs):
8         q_cond_ent = - (probs * torch.log(probs + 1e-12)).sum(
9 (1).mean()
10    return q_cond_ent
11
12    q_ent0 = get_entropy(prob0)
13    q_cond_ent0 = get_cond_entropy(prob0)
14
15    q_ent1 = get_entropy(prob1)
16    q_cond_ent1 = get_cond_entropy(prob1)
17
18    wim_0 = q_ent0 - self.hparams.cond_ent_weight *
19    q_cond_ent0
20    wim_1 = q_ent1 - self.hparams.cond_ent_weight *
21    q_cond_ent1
22
23    wim = (wim_0 + wim_1) / 2.
24    prob_loss = - wim
25    return prob_loss

```

+ Hàm forward để tính tất cả các hàm loss:

```

1 def forward(self, inputs):
2     embd_0 = self.pro_layer(self.get_embeddings(inputs,
3         pooling=self.hparams.pooler_type))
4     embd_1 = self.pro_layer(self.get_embeddings(inputs,
5         pooling=self.hparams.pooler_type))
6
7     Q_0, prob0_list = self.pq_head(embd_0)
8     Q_1, prob1_list = self.pq_head(embd_1)

```

```

8     codebooks_losses = []
9     for i in range(len(prob0_list)):
10        prob0 = prob0_list[i]
11        prob1 = prob1_list[i]
12        codebooks_losses.append(self.compute_prob_loss(prob0,
13                               prob1))
14
15     prob_loss = sum(codebooks_losses)
16     contra_loss = self.criterion(Q_0, Q_1)
17
18     loss = self.hparams.code_weight * contra_loss + self.
19     hparams.prob_weight * prob_loss
20
21     return {'loss': loss, 'contra_loss': contra_loss, 'prob_loss': prob_loss}

```

+ Hàm run_training_session để train model. Tham số **bad_epochs** để model có thể dừng train nếu số lượng **bad_epochs** lớn hơn số lượng bad epochs được quy định trước (**num_bad_epochs**).

```

1 def run_training_session(self, run_num, logger):
2     self.train()
3
4     # Scramble hyperparameters if number of runs is greater
5     # than 1.
6     if self.hparams.num_runs > 1:
7         logger.log('RANDOM RUN: %d/%d' % (run_num, self.
8         hparams.num_runs))
9         for hparam, values in self.get_hparams_grid().items():
10            assert hasattr(self.hparams, hparam)
11            self.hparams.__dict__[hparam] = random.choice(
12                values)
13
14            random.seed(self.hparams.seed)
15            torch.manual_seed(self.hparams.seed)
16
17            self.define_parameters()
18            # logger.log(str(self))
19            logger.log('%d params' % sum([p.numel() for p in self.
20                parameters()])))
21            logger.log('hparams: %s' % self.flag_hparams())
22
23            device = torch.device('cuda' if self.hparams.cuda else 'cpu')
24            self.to(device)

```

```

21
22     optimizer = self.configure_optimizers()
23     gradient_clippers = self.configure_gradient_clippers()
24     train_loader, database_loader, val_loader, _ = self.data.
25     get_loaders(
26         self.hparams.batch_size, self.hparams.num_workers,
27         shuffle_train=True, get_test=False
28     )
29     best_val_perf = float('-inf')
30     best_state_dict = None
31     bad_epochs = 0
32
33     # logger.log('Time: %s' % str(timedelta(seconds=round(
34         timer() - start))))
35     try:
36
37         for epoch in range(1, self.hparams.epochs + 1):
38             epoch_start = timer()
39
40             forward_sum = {}
41             num_steps = 0
42             for batch_num, batch in enumerate(train_loader):
43                 optimizer.zero_grad()
44
45                 inputs = squeeze_dim(
46                     move_to_device(batch[0], device), dim=1)
47                 forward = self.forward(inputs)
48
49                 for key in forward:
50                     if key in forward_sum:
51                         forward_sum[key] += forward[key]
52                     else:
53                         forward_sum[key] = forward[key]
54                     num_steps += 1
55
56                     if math.isnan(forward_sum['loss']):
57                         logger.log('Stopping epoch because loss
58                         is NaN')
59                         break
60
61                     forward['loss'].backward()
62                     for params, clip in gradient_clippers:
63                         nn.utils.clip_grad_norm_(params, clip)
64                     optimizer.step()

```

```

63         if math.isnan(forward_sum['loss']):
64             logger.log('Stopping training session because
65             loss is NaN')
66             break
67
68         val_perf = self.evaluate(
69             database_loader, val_loader, device, is_val=
70             True)
71             logger.log('End of epoch {:3d}'.format(epoch)
72             , False)
73             logger.log(' | '.join([' | {:.8f}' * num_steps]
74             .format(*[key, forward_sum[key] / num_steps]
75             for key in forward_sum]), False)
76             logger.log(' | val perf {:.8f}'.format(val_perf)
77             , False)
78             logger.log(' | time : {}'.format(
79                 str(timedelta(seconds=round(timer() -
80                 epoch_start)))))

81             if val_perf > best_val_perf:
82                 best_val_perf = val_perf
83                 bad_epochs = 0
84                 logger.log('\t *Best model so far, deep
85                 copying*')
86                 best_state_dict = deepcopy(self.
87                 state_dict())
88             else:
89                 bad_epochs += 1
90                 logger.log('\t Bad epoch %d' % bad_epochs)

91             if bad_epochs > self.hparams.num_bad_epochs:
92                 break
93         except KeyboardInterrupt:
94             logger.log('-' * 89)
95             logger.log('Exiting from training early')
96
97     return best_state_dict, best_val_perf

```

+ Để load các tham số từ pretrained model, ta sử dụng hàm `load()`:

```

1 def load(self):
2     device = torch.device('cuda' if self.hparams.cuda else 'cpu')
3     checkpoint = torch.load(self.hparams.model_path) if self.
4     hparams.cuda \
5     else torch.load(self.hparams.model_path,

```

```

5                         map_location=torch.device('cpu'))
6     if checkpoint['hparams'].cuda and not self.hparams.cuda:
7         checkpoint['hparams'].cuda = False
8     self.hparams = checkpoint['hparams']
9     self.define_parameters()
10    self.load_state_dict(checkpoint['state_dict'])
11    self.to(device)

```

+ Hàm `run_test` để chạy model, trả về kết quả top 100 documents và độ chính xác

```

1 def run_test(self):
2     device = torch.device('cuda' if self.hparams.cuda else 'cpu')
3     _, database_loader, val_loader, test_loader = self.data.
4     get_loaders(
5         self.hparams.batch_size, self.hparams.num_workers,
6         shuffle_train=True, get_test=True)
7     val_perf, val_res = self.evaluate(
8         database_loader, val_loader, device, is_val=True)
9     test_perf, test_res = self.evaluate(
10        database_loader, test_loader, device, is_val=False)
11    return val_perf, val_res, test_perf, test_res

```

2.4 Hàm đánh giá

- Hàm evaluate để đánh giá hiệu quả của model:

```

1 def evaluate(self, database_loader, eval_loader, device,
2             is_val=True):
3     self.eval()
4     with torch.no_grad():
5         perf, top_k = compute_retrieval_precision(
6             database_loader, eval_loader,
7             device, self.encode_continuous,
8             self.pq_head.C,
9             self.hparams.N_books,
10            self.hparams.num_retrieve)
11    self.train()
12    return perf, top_k

```

Trong đó, precision được tính bằng tỷ lệ số document liên quan đến query trên top 100 documents tìm kiếm được.

```

1 def compute_retrieval_precision(train_loader, eval_loader,
2                                 device, encode_continuous,

```

```

3                                     Codebooks , N_books ,
4                                     num_retrieve=100):
5
6     def extract_target(loader):
7         encoding_chunks = []
8         label_chunks = []
9         for (docs, labels) in loader:
10             docs = squeeze_dim(move_to_device(docs, device),
11                                 dim=1)
12             encoding_chunks.append(docs if encode_continuous
13                                     is None else
14                                     encode_continuous(docs))
15             label_chunks.append(labels)
16
17             encoding_mat = torch.cat(encoding_chunks, 0)
18             label_mat = torch.cat(label_chunks, 0)
19             label_lists = [[label_mat[i].item()] for i in range(
20                           label_mat.size(0))]
21             return encoding_mat, label_lists
22
23
24     def extract_database_code(loader, Codebooks):
25         code_chunks = []
26         label_chunks = []
27         for (docs, labels) in loader:
28             docs = squeeze_dim(move_to_device(docs, device),
29                                 dim=1)
30             encodings = encode_continuous(docs)
31             codes = Indexing(Codebooks, N_books, encodings)
32             code_chunks.append(codes)
33
34             label_chunks.append(labels)
35
36             code_mat = torch.cat(code_chunks, 0).type(torch.int)
37             label_mat = torch.cat(label_chunks, 0)
38             label_lists = [[label_mat[i].item()] for i in range(
39                           label_mat.size(0))]
40             return code_mat, label_lists
41
42
43     src_codes, src_label_lists = extract_database_code(
44     train_loader, Codebooks)
45
46
47     import time
48     start = time.time()
49     tgt_encodings, tgt_label_lists = extract_target(
50     eval_loader)
51     prec, top_k = compute_topK_average_precision(

```

```

41             tgt_encodings, tgt_label_lists,
42             src_codes, src_label_lists,
43             Codebooks, N_books,
44             num_retrieve)
45     return prec, top_k

```

3 Thực nghiệm

3.1 Thông tin các hyperparameter

- Số chiều của codeword $\frac{D}{M}$ trong codebook nhỏ là C^m cố định ở 24.
- Số lượng codeword K trong mỗi codebook C^m cố định ở 16.
- Số lượng codebook nhỏ M có thể là {4, 8, 16, 32}. Trong phần thực nghiệm này, số lượng được chọn là 4.
- Các codeword trong các codebook C có thể biểu diễn dưới dạng {16, 32, 64, 128} bit theo công thức $B = M \log_2 K$. Do đã chọn $M = 4$ và $K = 16$ nên số bit là 16.
- Learning rate được chọn là 0.001.
- Drop rate p_{drop} để tạo các cặp positives cho contrastive learning được chọn là 0.3.
- Gumbel-softmax temperature τ được chọn là 10 cho 16 bit và 5 cho code dài hơn.
- Trade-off coefficient α được chọn là 0.1.
- Coeficient được chọn trong {0.1, 0.2, 0.3} tùy theo hiệu năng trên tập validation.

3.2 Kết quả

Sau khi train với 100 epochs (tuy nhiên bị dừng sau epoch thứ 24 do số lượng bad epoch vượt quá số lượng bad epoch cho phép) và đánh giá trên tập test ta được kết quả :

	Validation	Test
Precision	81.29	80.39

So với kết quả do tác giả cung cấp trong paper, precision của thực nghiệm thấp hơn khoảng 2%. Nguyên nhân có thể là do quá trình train bị dừng sớm ở epoch 24. Tuy nhiên, kết quả thực nghiệm vẫn cao hơn đa số các phương pháp cũ.