

Projekt Information Retrieval

Bericht zur Implementierung von COIL in LightningIR

Lukas Günther

Samuel Müller

Tom Schmidt

Erstellt am 02 Februar, 2025

Friedrich Schiller Universität Jena
Ferdinand Schlatt

1 Einleitung

1.1 Zielsetzung

Ziel dieser Arbeit war es, das Retrieval Model COIL aus dem Paper "COIL: Revisit Exact Lexical Match in Information Retrieval with Contextualized Inverted List" von Luyu Gao, Zhuyun Dai, Jamie Callan (<https://github.com/luyug/COIL>) in "Lightning IR: Straightforward Fine-tuning and Inference of Transformer-based Language Models for Information Retrieval" von Ferdinand Schlatt, Maik Fröbe, Matthias Hagen (<https://github.com/webis-de/lightning-ir>) zu implementieren.

1.2 COIL

Traditionelle Informationsretrieval-Systeme wie BM25 basieren auf exakten lexikalischen Übereinstimmungen und nutzen invertierte Listen für eine effiziente Suche. Neuere neuronale IR-Modelle tendieren zu weichen semantischen Übereinstimmungen zwischen allen Begriffen in Anfrage und Dokument, verlieren dabei jedoch die rechnerische Effizienz der exakten Übereinstimmungssysteme.

Der im Paper "COIL: Revisit Exact Lexical Match in Information Retrieval with Contextualized Inverted List" von Luyu Gao, Zhuyun Dai und Jamie Callan vorgestellte Ansatz führt hingegen eine kontextualisierte exakte Übereinstimmungsarchitektur ein, bei der die Bewertung auf den kontextualisierten Repräsentationen von sich überschneidenden Token in Anfrage und Dokument basiert. Diese Architektur speichert die kontextualisierten Token-Repräsentationen in invertierten Listen und vereint so die Effizienz der exakten Übereinstimmung mit der Repräsentationskraft tiefer Sprachmodelle.

Damit liegt der zentrale Unterschied zu traditionellen und neuronalen IR-Methoden liegt in der Art, wie die exakte Wortübereinstimmung verarbeitet wird:

1. Statt nur einfache Termfrequenzen oder BM25-Werte zu speichern, speichert COIL für jedes Token im Dokument eine deep-learning-basierte, kontextualisierte Vektor-Repräsentation. Diese Repräsentationen werden aus einem Transformer-Modell wie zum Beispiel BERT extrahiert, sodass sie die semantische Bedeutung des Tokens im Kontext des Dokuments einfangen.
2. Klassische IR-Methoden speichern für jeden Token nur eine Posting-Liste mit Dokument-IDs und Termfrequenzen. COIL erweitert diese Struktur, indem es kontextualisierte Vektoren für jeden vorkommenden Token speichert. Dadurch bleibt die Effizienz der invertierten Listen erhalten, während gleichzeitig semantische Informationen genutzt werden.

3. COIL führt nur Berechnungen für genau übereinstimmende Tokens (Exact Lexical Match) durch, anstatt eine vollständige semantische Ähnlichkeitsberechnung über alle Token im Dokument durchzuführen. Dies bedeutet, dass es die Effizienz klassischer lexikalischer Methoden mit der Ausdruckskraft tiefer Sprachmodelle kombiniert (Contextualised Exact Lexical Match).
4. Bei der Anfrageverarbeitung sucht COIL nur nach übereinstimmenden Token zwischen Query und Dokument. Die Bewertung erfolgt dann durch die Interaktion der gespeicherten Token-Repräsentationen aus der invertierten Liste mit den Query-Token-Repräsentationen. Dadurch vermeidet COIL eine kostspielige Berechnung über alle Token eines Dokuments, wie es bei dichten neuronalen Repräsentationen der Fall wäre.

Experimentelle Ergebnisse zeigen, dass COIL sowohl klassische lexikalische Retrieval-Methoden als auch moderne neuronale Retrieval-Modelle übertrifft und dabei eine ähnliche oder geringere Latenz aufweist. Damit kombiniert die vorgestellte Architektur die Effizienz klassischer exakter Wortübereinstimmungssysteme mit der Ausdruckskraft tiefer Sprachmodelle.

1.3 Lightning-IR

Das Paper "Lightning IR: Straightforward Fine-tuning and Inference of Transformer-based Language Models for Information Retrieval" von Ferdinand Schlatt, Maik Fröbe und Matthias Hagen stellt Lightning IR vor, ein Framework zur Anwendung von Transformer-basierten Sprachmodellen im Bereich Information Retrieval. Dieses auf PyTorch Lightning basierende Framework bietet ein modular und erweiterbares Gerüst, welches alle Phasen einer Retrieval-Pipeline unterstützt: vom Fine-Tuning und Indexieren bis hin zum Suchen und Re-Ranking. Das Framework ist darauf ausgelegt, benutzerfreundlich, skalierbar und reproduzierbar zu sein und steht als Open-Source-Projekt zur Verfügung.

2 Code

2.1 Overview

Das Python-Dokument `coil.py` implementiert das Embedding und Scoring des Coil Retrieval Model in Lightning IR. Es enthält die folgenden Klassen:

- `CoilEmbedding`
- `CoilScoringFunction`
- `CoilConfig`
- `CoilModel`

2.2 CoilEmbedding

Coil Embedding erweitert `BiEncoderEmbedding` um das CLS Embedding.

```
1 @dataclass
2 class CoilEmbedding(BiEncoderEmbedding):
3     cls_embeddings: Tensor
```

2.3 CoilScoringFunction

CoilScoringFunction erweitert `ScoringFunction`.

CoilScoringFunction.forward berechnet die Scores, sowie die CLS Similarity. Die Maske, mit welcher sich die Scores errechnen, wird, wie im Coil Paper, durch Exact Lexical Matching bestimmt. Die Summe der Werte der Arrays Scores und CLS Similarity wird dann ausgegeben.

```

1 class CoilScoringFunction(ScoringFunction):
2
3     def forward(
4         self,
5         query_embeddings: CoilEmbedding,
6         doc_embeddings: CoilEmbedding,
7         num_docs: Sequence[int] | int | None = None,
8     ) -> Tensor:
9
10        num_docs_t = self._parse_num_docs(query_embeddings, doc_embeddings,
11        num_docs)
12        query_cls_embeddings = (query_embeddings.cls_embeddings).
13        repeat_interleave(num_docs_t, 0).unsqueeze(2)
14        doc_cls_embeddings = doc_embeddings.cls_embeddings.unsqueeze(1)
15
16        query_embeddings = self._expand_query_embeddings(query_embeddings,
17        num_docs_t)
18        doc_embeddings = self._expand_doc_embeddings(doc_embeddings,
19        num_docs_t)
20
21        similarity = self._compute_similarity(query_embeddings,
22        doc_embeddings)
23
24        query = query_embeddings.encoding.input_ids.repeat_interleave(
25        num_docs_t, 0)[: , 1:]
26        docs = doc_embeddings.encoding.input_ids[: , 1:]
27
28        mask = (query[: , :, None] == docs[: , None, :]).float()
29
30        cls_similarity = self.similarity_function(query_cls_embeddings,
31        doc_cls_embeddings)
32
33        scores = self._aggregate(mask * similarity, doc_embeddings.
34        scoring_mask[: , :, 1:], "max", -1)
35        scores = self._aggregate(scores, query_embeddings.scoring_mask[: ,
36        1:, :], self.query_aggregation_function, -2)
37        return scores[... , 0, 0] + cls_similarity[... , 0, 0]

```

2.4 CoilConfig

CoilConfig erweitert **BiEncoderConfig** um die CLS-Embedding-Dimension.

```
1
2 class CoilConfig(BiEncoderConfig):
3     model_type = "coil"
4
5     def __init__(
6         self,
7         query_expansion: bool = False,
8         embedding_dim: int = 32,
9         cls_embedding_dim: int = 768,
10        projection: Literal["linear", "linear_no_bias"] | None = "
linear_no_bias",
11        **kwargs,
12    ) -> None:
13        kwargs["query_pooling_strategy"] = None
14        kwargs["doc_expansion"] = False
15        kwargs["attend_to_doc_expanded_tokens"] = False
16        kwargs["doc_pooling_strategy"] = None
17        super().__init__(
18            query_expansion=query_expansion,
19            embedding_dim=embedding_dim,
20            projection=projection,
21            **kwargs,
22        )
23        self.cls_embedding_dim = cls_embedding_dim
```

2.5 CoilModel

CoilModel erweitert **BiEncoderModel**.

CoilModel.encode encodes die Token-Sequenzen für die Queries und Documents. Zurückgegeben wird ein **CoilEmbedding**, welches die CLS-Embeddings und die Tokenembeddings von Query und Document separat enthält.

```

1 class CoilModel(BiEncoderModel):
2     config_class = CoilConfig
3
4     def __init__(self, config: BiEncoderConfig, *args, **kwargs) -> None:
5         super().__init__(config, *args, **kwargs)
6         self.scoring_function = CoilScoringFunction(self.config)
7
8         if self.config.projection is not None:
9             if "linear" in self.config.projection:
10                 self.cls_projection = Linear(
11                     self.config.hidden_size,
12                     self.config.cls_embedding_dim,
13                     bias="no_bias" not in self.config.projection,
14                 )
15             else:
16                 raise ValueError(f"Projection for COIL is not supported")
17         else:
18             if self.config.embedding_dim != self.config.hidden_size:
19                 warnings.warn(
20                     "No projection is used (embedding_dim != hidden_size)"
21                     "Output embeddings will not have embedding_size dim."
22                 )
23     def encode(
24         self,
25         encoding: BatchEncoding,
26         expansion: bool = False,
27         pooling_strategy: Literal["first", "mean", "max", "sum"] | None =
None,
28         mask_scoring_input_ids: Tensor | None = None,
29     ) -> BiEncoderEmbedding:
30         embeddings = self._backbone_forward(**encoding).last_hidden_state
31
32         if self.projection is not None:
33             cls_embeddings = self.cls_projection(embeddings[:, [0]])
34             embeddings = self.projection(embeddings[:, 1:])
35
36         embeddings = self._sparsification(embeddings, self.config.
sparsification)
37         embeddings = self._pooling(embeddings, encoding["attention_mask"],
pooling_strategy)
38
39         if self.config.normalize:
40             embeddings = normalize(embeddings, dim=-1)
41         scoring_mask = self.scoring_mask(encoding, expansion,
pooling_strategy, mask_scoring_input_ids)
42         return CoilEmbedding(embeddings, scoring_mask, encoding,
cls_embeddings)

```

3 Modell Checkpoint auf Hugging Face

Die Entwickler des COIL Models stellen vortrainierte Modelle bereit und unterscheiden dabei zwischen zwei Systemen: eines verwendet harte Negative (HN) und das andere nicht. COIL ohne HN wird mit BM25-Negativen trainiert, während COIL mit HN zusätzlich mit harten Negativen trainiert wird, die mit einem anderen trainierten COIL-Modell gewonnen wurden.

Für die Implementierung in COIL wurde der, sich in der von COIL bereitgestellten GitHub Repository befindende, Checkpoint *hn-checkpoint.tar.gz* genutzt. An diesem Checkpoint wurde für das im Paper beschriebene CLS Token noch ein neuer Projectionslayer hinzugefügt und dann mittels *module.save_pretrained()* ein neues Modell erstellt. Das Modell wurde dann auf Huggingface unter <https://huggingface.co/inforettom/InfoRet-Coil> bereitgestellt.

4 Evaluation durch ReRanking

Die Entwickler des COIL-Models haben für ihr Model schon eine ausreichende Auswertung gemacht. Trotzdem wäre es interessant zu testen, ob unsere Implementierung funktioniert.

Für die Auswertung wird ein Re-Ranking zwischen COIL und ColBERT durchgeführt. Das ursprüngliche Ranking stammt vom Modell Anserini (BM25) und basiert auf den Datensätzen *msmarco-passage-trec-dl-2019-judged* und *msmarco-passage-trec-dl-2020-judged*. Wir nutzen den $nDCG@10$ (normalized Discounted Cumulative Gain at rank 10) als Metrik für unsere Auswertung. $nDCG@10$ basiert auf dem DCG (Discounted Cumulative Gain), der die Relevanz von Dokumenten unter Berücksichtigung ihrer Position in der Ergebnisliste gewichtet, indem höher gerankte Dokumente stärker gewichtet werden. Der DCG wird anschließend normalisiert ($nDCG$), indem er durch den idealen DCG (IDCG) geteilt wird, um Werte zwischen 0 und 1 zu erhalten. Ein $nDCG@10$ von 1 bedeutet, dass die zehn relevantesten Dokumente in perfekter Reihenfolge präsentiert wurden.

4.1 Anserini (BM25)

Die Shell-Befehle für die Auswertung des gegebenen Rankings:

```
1 $ ir_measures "msmarco-passage/trec-dl-2019" "msmarco-passage-trec-dl-2019-  
   judged.run" nDCG@10  
2  
3 $ ir_measures "msmarco-passage/trec-dl-2020" "msmarco-passage-trec-dl-2020-  
   judged.run" nDCG@10
```

Das Ausgangs-Ranking erzielt folgende $nDCG@10$ -Werte:

| Test Metric | $nDCG@10$ |
|-------------------------------------|-----------|
| msmarco-passage/trec-dl-2019/judged | 0.5058 |
| msmarco-passage/trec-dl-2020/judged | 0.4796 |

Tabelle 1: $nDCG@10$ -Werte - Anserini (BM25)

4.2 COIL

Das Skript für das Re-Ranking mittels COIL:

```
1 trainer
2   logger false
3 model
4   class_path lightning_ir.BiEncoderModule
5   init_args
6     model_name_or_path ./hn-checkpoint
7     config coil.CoilConfig
8     evaluation_metrics
9       - nDCG@10
10 data
11   class_path lightning_ir.LightningIRDataModule
12   init_args
13     inference_datasets
14       - class_path: lightning_ir.RunDataset
15         init_args
16           run_path_or_id ./msmarco-passage-trec-dl-2019-judged.
17       run
18         depth 100
19       - class_path: lightning_ir.RunDataset
20         init_args
21           run_path_or_id ./msmarco-passage-trec-dl-2020-judged.
22       run
23         depth 100
24     inference_batch_size 4
```

Das Re-Ranking mit COIL erzielt folgende nDCG@10-Werte:

| Test Metric | nDCG@10 |
|-------------------------------------|---------|
| msmarco-passage/trec-dl-2019/judged | 0.6837 |
| msmarco-passage/trec-dl-2020/judged | 0.6512 |

Tabelle 2: nDCG@10-Werte - COIL

4.3 ColBERT

Das Skript für das Re-Ranking mittels ColBERT:

```
1 trainer
2     logger false
3 model
4     class_path lightning_ir.BiEncoderModule
5     init_args
6         model_name_or_path webis/colbert
7         evaluation_metrics
8             - nDCG@10
9 data
10    class_path lightning_ir.LightningIRDataModule
11    init_args
12        inference_datasets
13            - class_path: lightning_ir.RunDataset
14            init_args
15                run_path_or_id ./msmarco-passage-trec-dl-2019-judged.
16            run
17                depth 100
18            - class_path: lightning_ir.RunDataset
19            init_args
20                run_path_or_id ./msmarco-passage-trec-dl-2020-judged.
21            run
22                depth 100
23            inference_batch_size 4
```

Das Re-Ranking mit ColBERT ergibt folgende nDCG@10-Werte:

| Test Metric | nDCG@10 |
|-------------------------------------|---------|
| msmarco-passage/trec-dl-2019/judged | 0.3980 |
| msmarco-passage/trec-dl-2020/judged | 0.3161 |

Tabelle 3: nDCG@10-Werte - ColBERT

4.4 Auswertung

Das COIL-Re-Ranking verbessert den nDCG@10-Wert des ursprünglichen Rankings deutlich. ColBERT hingegen schneidet sogar schlechter als das Originalranking ab. Dies könnte dran liegen, dass wir für das Re-Ranking eine Tiefe von 100 nehmen und

nicht das volle Ranking. **Das war jedoch Aufgrund von Hardwarerestriktionen und dem Umfang des Moduls zum jetzigen Zeitpunkt nicht möglich.**