



Department of Information Engineering

Master thesis in Artificial Intelligence and Data Engineering

# **Learning What to Crawl: Metadata-Enriched LLMs for Crawling Frontier Prioritisation**

Supervisors:

**Prof: Nicola Tonellotto**

**Dott.ssa: Francesca Pezzuti**

Candidate :

**Niccolò Settimelli**

Academic Year 2024/2025

## **Abstract**

The exponential growth of the Web requires crawlers to efficiently discover valuable content. Classical strategies such as PageRank rely on global link analysis, which, while effective, is computationally expensive and less practical for real-time frontier management.

This thesis proposes a shift in perspective: from global graph-based scoring to semantic-aware frontier prioritization. Leveraging transformer-based neural models, we estimate the potential of a page to lead toward relevant content through its outlinks, based solely on its textual and metadata features. This reframing from page-centric quality estimation to outlink utility prediction aligns crawling with the broader transition from keyword-based retrieval to semantic search.

The proposed approach integrates semantic modelling and metadata fusion to generate real-time priority scores that guide the crawler toward semantically rich content. By embedding neural estimation directly into the crawling process, our methodology improves efficiency, surfaces higher-quality documents earlier, and avoids the prohibitive costs of repeated graph-wide computations.

This work highlights how modern neural architectures can redefine the principles of web crawling, paving the way for next-generation crawlers that are lighter, adaptive, and better aligned with semantic search paradigms.

# Chapter 1

## Code and Implementation

This section documents the practical setup used to reproduce our finetuning experiments, the repository layout, and the exact launch commands. The goal is to make every step auditable, from dataset construction to model training and score extraction for crawler integration.

### 1.1 Repository and Environment

The codebase is publicly available on GitHub as [nikisetti01/crawler-quality-scoring-model](#). All experiments were run in a Python  $\geq 3.10$  environment to ensure compatibility with the `pyterrier-quality` stack used for QualScore conversion.

- **Virtual environment.**

```
python3 -m venv .venv
source .venv/bin/activate
pip install --upgrade pip
pip install -r requirements.txt
```

- **Accelerate configuration (one-time).**

```
accelerate config
# select: multi-GPU (if available), mixed precision (fp16 or bf16),
# and gradient accumulation as required by GPU memory.
```

### Repository Structure

The repository is organized by experimental task; each top-level folder may contain dedicated subfolders/files for *text-only* (prefix `only_text`) and *metadata* (prefix `metadata_`).

- **Dataset\_generation:** scripts to build the two main training corpora (sampling, cleaning, joins).

- **Finetuning:** training scripts for all model families:
  - `finetune_only_text.py`: BERT text-only baseline (binary classification).
  - `finetune_metadata.py`: BERT with metadata fusion (additive Uni-Attention).
- **Test:** evaluation scripts on held-out sets to verify that fine-tuning improves intrinsic metrics.
- **Results:** notebooks for plots/tables reporting learning curves, ROC/AUC, calibration, and ablations.
- **Inference:** batch inference utilities to score large collections and convert scores into *QualScore* via `pyterrier-quality` (local QualCache hub).

## Pipeline Orchestration (Dataset Creation)

Dataset construction is automated by a lightweight pipeline driver that sequentially executes preprocessing stages and logs outputs to disk. The following example mirrors the style of our internal scripts:<sup>1</sup>

```
#!/usr/bin/env bash
set -euo pipefail

# Select which models/datasets to prepare (arrays can be expanded)
models=("bert-base-uncased")
datasets=("train-split-a")

for MODEL_PATH in "${models[@]}"; do
  for DATASET_NAME in "${datasets[@]}"; do

    source ./scripts/models_config.sh    # load per-model hyperparameters

    python process_links.py \
      --dataset ${DATASET_NAME} \
      --out tmp/links_${DATASET_NAME}.csv

    python division_fast.py \
      --in tmp/links_${DATASET_NAME}.csv \
      --test_out data/test_set.csv \
      --links_out data/link_labels2_cleaned.csv

    python add_text.py \
      --links data/link_labels2_cleaned.csv \
      --text_corpus /data/text \
      --train_out data/dataset_train.csv \
```

---

<sup>1</sup>The actual filenames may differ depending on your local dataset sources; the orchestration pattern remains identical.

```

--val_out    data/dataset_val.csv

python balancer.py \
  --in  data/dataset_train.csv \
  --out data/dataset_balanced.csv \
  --strategy label_ratio --pos_ratio 0.5

done
done

```

**Scripts and roles.** *(i)* `process_links.py`: extracts/normalizes link-level labels; *(ii)* `division_fast.py`: builds a balanced test split and cleans link labels; *(iii)* `add_text.py`: joins labels with cleaned page text to produce `dataset_train.csv` and `dataset_val.csv`; *(iv)* `balancer.py`: optionally re-balances the training set according to a target positive ratio.

## Finetuning Launch Commands

We use `accelerate launch` as a drop-in replacement for `python`, enabling multi-GPU and mixed precision where available.

### Text-only BERT (binary classification).

```

accelerate launch --mixed_precision=fp16 finetune_only_text.py \
  --model_name bert-base-uncased \
  --train_csv data/dataset_train.csv \
  --val_csv   data/dataset_val.csv \
  --out       ckpts/best_bert.pt \
  --batch_size 32 \
  --lr 2e-5 \
  --epochs 3 \
  --max_len 512 \
  --sep ", "

```

### BERT + Metadata (additive Uni-Attention).

```

accelerate launch --mixed_precision=fp16 finetune_metadata.py \
  --model_name bert-base-uncased \
  --train_csv data/dataset_train_metadata.csv \
  --val_csv   data/dataset_val_metadata.csv \
  --out       ckpts/best_meta.pt \
  --batch_size 24 \
  --lr 2e-5 \
  --epochs 3 \
  --max_len_text 512 \
  --max_len_meta 64 \
  --H 128 --A 64 --V_dom 50000 \
  --num_cols num_inlinks,num_outlinks,outlinks_domains_count,inlinks_domains_count \
  --sep ", "

```

## User-Selectable Parameters

The following arguments are surfaced to the user at launch time; they control data access, model choice, and training budget:

- `-model_name`: HF encoder backbone (e.g., `bert-base-uncased`).
- `-train_csv`, `-val_csv`: input files with columns `text`, `label` and, for metadata runs, `anchor/domain/numeric` fields as specified in the repository readme.
- `-out`: destination path for the best `state_dict`.
- `-batch_size`, `-lr`, `-epochs`: standard optimization knobs.
- `-max_len` / `-max_len_text`, `-max_len_meta`: token limits for text and metadata encoders.
- `-H`, `-A`, `-V_dom`: metadata latent dim, additive-attention hidden size, and hashing space for domains.
- `-num_cols`: comma-separated list of numeric metadata columns to include (automatically projected and fused).
- `-sep`: CSV separator (`,` or `\t`).
- `accelerate` options: device visibility (`CUDA_VISIBLE_DEVICES`), mixed precision (`fp16/bf16`), gradient accumulation.

## 1.2 Finetuning Experiment Setup and Optimization

The finetuning stage constitutes the core of our experimental framework and is organized into two main subdirectories: `only_text_finetuning` and `metadata_finetuning`. In the text-only case (`binary_classification` and `score_estimation`), the training pipeline differs structurally from the metadata setting. For text-only runs, we leverage existing modules from Hugging Face `transformers` and `datasets` to manage textual inputs, tokenize them, and feed them into either encoder-based models with a classification head (e.g., BERT, DeBERTa) or encoder-decoder models (QualT5, derived from T5).

The only customization added in this setting was a `WeightedTrainer` class, extending Hugging Face’s standard `Trainer` to apply a class-weighted cross-entropy loss, thereby accounting for the imbalance in the label distribution:

```
1 class WeightedTrainer(Trainer):
2     def compute_loss(self, model, inputs, return_outputs=False):
3         labels = inputs.pop("labels")
4         outputs = model(**inputs)
5         logits = outputs.logits
6         loss_fct =
7         torch.nn.CrossEntropyLoss(weight=class_weights_tensor)
8         loss = loss_fct(
9             logits.view(-1, model.config.num_labels),
```

```

9         labels.view(-1)
10     )
11     return (loss, outputs) if return_outputs else loss

```

Listing 1.1: Weighted cross-entropy trainer for binary classification.

In addition, an auxiliary function `compute_metrics` was employed during validation and testing. This function evaluates not only the loss trajectory but also key classification metrics such as accuracy and F1-score. Once defined, the parameters for Hugging Face’s `TrainingArguments` and the custom `WeightedTrainer` were set as summarized in Table 1.1, using models `bert-uncased` and `deberta-v3`.

As an additional optimization, we parallelized training across two GPUs and enabled `**mixed precision training (FP16) **` [micikevicius2017mixed]. By performing most operations in 16-bit floating point, this approach significantly reduces GPU memory footprint and computational overhead, while retaining numerical stability through automatic casting. All metrics during both training and validation were continuously monitored through direct integration with the W&B platform, providing reliable tracking for large-scale finetuning tasks.

Table 1.1: Training setup for text-only binary classification experiments.

<b>Tokenizer</b>	AutoTokenizer (RoBERTa / DeBERTa), max_length=512, truncation + padding
<b>Data split</b>	Train/Validation 90/10 (stratified, seed=42); Test from CSV
<b>Preprocessing</b>	Remove invalid rows (empty / non-string text)
<b>Batch size</b>	Train = 64 /device; Eval = 256 /device
<b>Learning rate</b>	$2 \times 10^{-5}$
<b>Epochs</b>	3
<b>Weight decay</b>	0.01
<b>Precision</b>	FP16 (mixed precision enabled)
<b>Loss function</b>	Weighted Cross-Entropy (balanced class weights)
<b>Freezing policy</b>	All frozen except <code>encoder.layer.10</code> , <code>encoder.layer.11</code> , and <code>classifier</code>
<b>Evaluation</b>	Every 20,000 steps
<b>Checkpointing</b>	Save every 20,000 steps (keep last 3); best model loaded at end
<b>Best model metric</b>	Accuracy
<b>Data collator</b>	<code>DataCollatorWithPadding(tokenizer)</code>
<b>Optimizer / Scheduler</b>	Hugging Face defaults (AdamW + scheduler)
<b>Logging</b>	Every 500 steps, reported to W&B
<b>Training status</b>	~200k steps completed

### 1.2.1 Score Estimation (Seq2Seq, QualT5)

For score estimation we adopt a generative reformulation: given the document text, the model produces the token "true" or "false". The main deviation from the binary classification setup lies in the input preprocessing, where we apply minimal prompt engineering to create natural language instructions and very short targets.

```
1 def format_prompt(example):
2     # Prompt: provide the document, ask for relevance (true/false)
3     prompt = f"Document: {example['text']} Relevant:"
4     # Target: "true" if label = 1, else "false"
5     target = "true" if int(example["label"]) == 1 else "false"
6     return {"input_text": prompt, "target_text": target}
```

Listing 1.2: Prompt formatting for QualT5 score estimation.

At evaluation time, predictions are decoded and mapped to {0,1} to compute accuracy (and optional F1). Training is managed via `Seq2SeqTrainer` with `predict_with_generate=True`, ensuring that validation mirrors the inference pathway.



Table 1.2: Training setup for text-only score estimation (QualT5).

<b>Tokenizer</b>	AutoTokenizer (QualT5)
<b>Task format</b>	Input: Document: <text> Relevant: → Target: "true" / "false"
<b>Max lengths</b>	Input max_length = 512 (truncation + padding); Target max_length = 5
<b>Data split</b>	Train/Validation 90/10 (stratified by label, seed = 42); Test from CSV
<b>Preprocessing</b>	Drop rows with empty / non-string text
<b>Batch size</b>	Train = 64 /device; Eval = 256 /device
<b>Learning rate</b>	$2 \times 10^{-5}$
<b>Epochs</b>	3
<b>Weight decay</b>	0.01
<b>Precision</b>	FP32 (no mixed precision set)
<b>Loss / Objective</b>	Seq2Seq Cross-Entropy via decoder (generate "true"/"false")
<b>Predict with generate</b>	True (evaluation uses decoding)
<b>Freezing policy</b>	Encoder frozen except last 2 layers; decoder trainable
<b>Evaluation</b>	Every 20,000 steps
<b>Checkpointing</b>	Save every 20,000 steps (keep last 3); load best at end
<b>Best model metric</b>	Accuracy
<b>Metrics function</b>	Decode outputs → map to 0/1 → compute Accuracy
<b>Data collator</b>	DataCollatorForSeq2Seq(tokenizer, model)
<b>Optimizer / Scheduler</b>	Hugging Face defaults (AdamW + scheduler)
<b>Training status</b>	~180k steps completed

### 1.2.2 Multimodal Experimentation Setup

The multimodal experiments extend the text-only pipeline by introducing a dedicated metadata fusion layer on top of a pre-trained encoder (BERT). The architecture is implemented in three PyTorch modules: (i) a *MetadataEncoder* that projects heterogeneous metadata (anchor texts, domain lists, and numeric link statistics) into a shared hidden space; (ii) an *Additive Unified Attention* that aligns metadata tokens over the contextual textual sequence; and (iii) a *MultiModal-WebClassifier* that fuses the enriched metadata representation with the textual [CLS] vector for binary prediction. The implementation follows the codebase in `model_bert.py`.

## Model Components

**MetadataEncoder.** Each metadata source is mapped through a small MLP into a common dimensionality and stacked as a short “meta sequence”. This yields a compact representation amenable to cross-modal attention.

```
1 class MetadataEncoder(nn.Module):
2     """
3     Encode heterogeneous metadata sources (anchors, domains, numerics)
4     into a unified hidden representation. Each feature type is projected
5     through a small feed-forward layer into the same hidden space.
6     """
7     def __init__(self, hidden_dim):
8         super().__init__()
9         self.anchor_encoder = nn.Sequential(nn.Linear(32, hidden_dim),
10 nn.ReLU())
11         self.domain_proj = nn.Sequential(nn.Linear(64, hidden_dim),
12 nn.ReLU())
13         self.numeric_proj = nn.Sequential(nn.Linear(8, hidden_dim),
14 nn.ReLU())
15
16     def forward(self, anchor_out, anchor_in, domain_out, domain_in,
17 numerics):
18         # Stack different metadata encodings along a pseudo-sequence
19         axis
20         return torch.stack([
21             self.anchor_encoder(anchor_out.float()),
22             self.anchor_encoder(anchor_in.float()),
23             self.domain_proj(domain_out.float()),
24             self.domain_proj(domain_in.float()),
25             self.numeric_proj(numerics)
26         ], dim=1) # shape [B, 5, H]
```

Listing 1.3: MetadataEncoder: project and pack heterogeneous metadata.

**Additive Unified Attention.** A Bahdanau-style additive attention lets each metadata token attend over the full textual sequence. Residual connections and LayerNorm stabilize optimization; dropout regularizes the cross-modal match.

```
1 class AdditiveUniAttention(nn.Module):
2     """
3     Unified additive attention (Bahdanau-style) that aligns metadata
4     tokens
5     with the textual representation. Each metadata token queries the
6     entire
7     text sequence to extract relevant contextual information.
8     """
9     def __init__(self, meta_h=HIDDEN_DIM, txt_h=BERT_H, attn_h=64):
10         super().__init__()
11         self.Wq = nn.Linear(meta_h, attn_h, bias=True) # project
12         metadata tokens
```

```

10         self.Wk = nn.Linear(txt_h, attn_h, bias=True)      # project
    textual tokens
11         self.v = nn.Linear(attn_h, 1, bias=True)            # scoring
    vector
12         self.Vv = nn.Linear(txt_h, meta_h, bias=True)      # project
    values into meta space
13         self.norm = nn.LayerNorm(meta_h)
14         self.dropout = nn.Dropout(0.1)
15
16     def forward(self, meta_tokens, bert_hidden, bert_attn_mask):
17         B, M, H = meta_tokens.size()
18         L = bert_hidden.size(1)
19
20         # Linear projections
21         Qe = self.Wq(meta_tokens)      # queries from metadata
22         Ke = self.Wk(bert_hidden)      # keys from text
23
24         # Compute additive attention scores:  $v^T \tanh(Q + K)$ 
25         Qe_exp = Qe.unsqueeze(2).expand(-1, -1, L, -1)    # [B, M, L, A]
26         Ke_exp = Ke.unsqueeze(1).expand(-1, M, -1, -1)    # [B, M, L, A]
27         e_ij = torch.tanh(Qe_exp + Ke_exp)
28         scores = self.v(e_ij).squeeze(-1)                  # [B, M, L]
29
30         # Apply mask to ignore padding tokens
31         if bert_attn_mask is not None:
32             mask = (bert_attn_mask == 1).unsqueeze(1).expand(-1, M, -1)
33             scores = scores.masked_fill(~mask, float('-inf'))
34
35         # Normalize into attention weights
36         attn = torch.softmax(scores, dim=-1)                # [B, M, L]
37
38         # Weighted sum of textual values projected into metadata space
39         Vproj = self.Vv(bert_hidden)                        # [B, L, H]
40         context = torch.bmm(attn, Vproj)                    # [B, M, H]
41
42         # Residual connection + layer normalization
43         out = self.norm(meta_tokens + self.dropout(context))
44         return out

```

Listing 1.4: AdditiveUniAttention: Bahdanau-style cross-modal alignment.

**MultiModalWebClassifier.** The final head concatenates the textual [CLS] vector and the metadata vector (mean-pooled over meta tokens) and predicts a binary logit. A masked-mean over domain embeddings avoids bias from padding.

```

1 class MultiModalWebClassifier(nn.Module):
2     """
3     Integrates:
4         1) BERT text encoding,
5         2) metadata tokens (anchors, domains, numerics),
6         3) additive cross-modal attention.

```

```

7   Outputs a binary relevance logit.
8   """
9   def __init__(self, encoder: BertModel, hidden_dim=HIDDEN_DIM,
10  use_lora=False):
11       super().__init__()
12       self.encoder = encoder
13       self.domains_embedding = nn.Embedding(30522, 64)
14       self.meta_encoder = MetadataEncoder(hidden_dim)
15       self.uni_attn = AdditiveUniAttention(meta_h=hidden_dim,
16  txt_h=BERT_H, attn_h=64)
17       self.head = nn.Sequential(
18           nn.Linear(hidden_dim + BERT_H, 128),
19           nn.ReLU(),
20           nn.Dropout(0.2),
21           nn.Linear(128, 1)  # binary logit
22       )
23
24   def forward(self, input_ids, attention_mask,
25               anchor_out_ids, anchor_in_ids,
26               anchor_out_mask, anchor_in_mask,
27               domains_out_ids, domains_in_ids,
28               numerics):
29
30       enc = self.encoder(input_ids=input_ids,
31  attention_mask=attention_mask)
32       last_hidden = enc.last_hidden_state          # [B, L, 768]
33       cls = last_hidden[:, 0]                     # [CLS] [B, 768]
34
35       # Masked mean pooling over domain IDs (skip padding)
36       def masked_mean(emb, ids):
37           mask = (ids != 0).unsqueeze(-1)
38           emb = emb * mask
39           return emb.sum(1) / mask.sum(1).clamp(min=1e-6)
40
41       dom_out = masked_mean(self.domains_embedding(domains_out_ids),
42  domains_out_ids)
43       dom_in = masked_mean(self.domains_embedding(domains_in_ids),
44  domains_in_ids)
45
46       meta_tokens = self.meta_encoder(anchor_out_ids, anchor_in_ids,
47  dom_out, dom_in, numerics) # [B, 5, H]
48       attn_out = self.uni_attn(meta_tokens, last_hidden,
49  attention_mask)           # [B, 5, H]
50       meta_vec = attn_out.mean(dim=1)
51
52       logits = self.head(torch.cat([cls, meta_vec],
53  dim=-1)).squeeze(-1)
54       return logits

```

Listing 1.5: with metadata and classify.]MultiModalWebClassifier: fuse [CLS] with metadata and classify.

*Design notes.* The model includes residual connections and LayerNorm inside the cross-modal attention block, a masked-mean for domain embeddings, and a modest MLP head with dropout. These choices proved effective to stabilize gradients and mitigate overfitting under heterogeneous inputs. :contentReference[oaicite:1]index=1

## Dataset Creation (Tokenization and Normalization)

The dataset creation script tokenizes clean text and textual metadata while normalizing numeric features. The core tokenization routine is reported below.

```

1 def tokenize_fn(batch):
2     def to_str_list(x):
3         return [str(e) if e is not None else "" for e in x]
4
5     text_enc      = tokenizer(to_str_list(batch['text']),
6 truncation=True, padding="max_length", max_length=512)
7     anchor_out_en =
8 tokenizer(to_str_list(batch['outlink_outlink_anchors']),
9 truncation=True, padding="max_length", max_length=32)
10    anchor_in_enc =
11 tokenizer(to_str_list(batch['inlink_inlink_anchors']),
12 truncation=True, padding="max_length", max_length=32)
13    domain_out_en = tokenizer(to_str_list(batch['domains_out']),
14 truncation=True, padding="max_length", max_length=64)
15    domain_in_enc = tokenizer(to_str_list(batch['domains_in']),
16 truncation=True, padding="max_length", max_length=64)
17
18    return {
19        'input_ids':      text_enc['input_ids'],
20        'attention_mask': text_enc['attention_mask'],
21        'anchor_out_ids': anchor_out_en['input_ids'],
22        'anchor_out_mask': anchor_out_en['attention_mask'],
23        'anchor_in_ids':  anchor_in_enc['input_ids'],
24        'anchor_in_mask': anchor_in_enc['attention_mask'],
25        'domains_out_ids': domain_out_en['input_ids'],
26        'domains_out_mask': domain_out_en['attention_mask'],
27        'domains_in_ids': domain_in_enc['input_ids'],
28        'domains_in_mask': domain_in_enc['attention_mask'],
29    }

```

Listing 1.6: Tokenization function for multimodal inputs.

Training proceeds in two phases. Phase-1 freezes the encoder and trains the fusion head at large batch size; Phase-2 unfreezes the encoder and injects LoRA adapters into self-attention projections for parameter-efficient adaptation. The salient training code is shown next, following `train.py`. :contentReference[oaicite:2]index=2

```

1 def run_eval_bert(model, val_loader, criterion):
2     model.eval()
3     val_loss, correct, total = 0.0, 0, 0
4     with torch.no_grad():

```

```

5         for batch in val_loader:
6             outputs = model(**{k: v for k, v in batch.items() if k !=
"label"})
7             loss = criterion(outputs, batch["label"])
8             val_loss += loss.item()
9             preds = (outputs > 0.0).long()
10            correct += (preds == batch["label"].long()).sum().item()
11            total    += batch["label"].size(0)
12    return (val_loss / max(len(val_loader), 1)), (correct / max(total,
1))

```

Listing 1.7: BERT evaluation loop (loss + accuracy with 0.5 threshold).

```

1 def train_phase_bert(model, train_loader, val_loader, epochs, lr,
max_steps,
2                     freeze_encoder=True, use_lora=False):
3     # Freeze/unfreeze backbone
4     for p in model.encoder.parameters():
5         p.requires_grad = not freeze_encoder
6
7     # Optional: add LoRA adapters to the encoder (after freeing)
8     if use_lora:
9         accelerator.wait_for_everyone()
10        model.encoder = apply_lora_to_bert_encoder(model.encoder)
11        model = accelerator.prepare(model)
12
13    optimizer = make_optimizer(model, lr=lr)
14    optimizer = accelerator.prepare(optimizer)
15
16    best_val = float("inf")
17    for epoch in range(epochs):
18        model.train()
19        running, steps_this_epoch = 0.0, 0
20
21        with tqdm(total=min(max_steps, len(train_loader)),
22                  desc=f"Epoch {epoch+1}",
23                  disable=not accelerator.is_local_main_process) as
pbar:
24            for batch in train_loader:
25                if steps_this_epoch >= max_steps:
26                    break
27                optimizer.zero_grad(set_to_none=True)
28                outputs = model(**{k: v for k, v in batch.items() if k
!= "label"})
29                loss = criterion(outputs, batch["label"])
30                accelerator.backward(loss)
31                torch.nn.utils.clip_grad_norm_(model.parameters(), 1.0)
32                optimizer.step()
33
34                running += loss.item()
35                steps_this_epoch += 1
36                pbar.update(1)

```

```

37     val_loss, val_acc = run_eval_bert(model, val_loader, criterion)
38     train_loss = running / max(1, min(max_steps, len(train_loader)))
39
40
41     if accelerator.is_local_main_process:
42         wandb.log({"epoch": epoch + 1, "train_loss": train_loss,
43                  "val_loss": val_loss, "val_accuracy": val_acc})
44         if val_loss < best_val:
45             best_val = val_loss
46             accelerator.save(model.state_dict(),
"best_model_bert.pt")

```

Listing 1.8: Two-phase BERT training (freeze → unfreeze + LoRA).

Table 1.3: Phase-1: Base Multimodal (frozen encoder).

Tokenizer	BertTokenizerFast
<b>Inputs</b>	Text (512), Out-anchors (32), In-anchors (32), Out-domains (64), In-domains (64), Numerics (8 standardized)
<b>Metadata prep</b>	Train-only StandardScaler (numerics); anchors/domains cleaned & normalized; domain lists joined with space
<b>Split</b>	Train/Validation = 90/10 (stratified by label, seed=42)
<b>Batch size</b>	256
<b>Optimizer</b>	AdamW (lr = $10^{-3}$ , weight_decay=0.01)
<b>Loss</b>	BCEWithLogitsLoss(pos_weight) (from label distribution)
<b>Freezing policy</b>	Encoder frozen, train metadata fusion head only
<b>Evaluation metric</b>	Validation accuracy (threshold 0.5 on sigmoid)
<b>Max steps</b>	25,000
<b>Infrastructure</b>	accelerate (device placement + AMP), W&B logging

Table 1.4: Phase-2: LoRA Fine-Tuning (unfrozen encoder).

<b>Batch size</b>	32
<b>Optimizer</b>	AdamW (lr = $10^{-4}$ , weight_decay=0.01)
<b>Freezing policy</b>	Encoder unfrozen
<b>LoRA injection</b>	Self-attention projections: Q, K, V in <code>encoder.layer.*.attention.self</code>
<b>LoRA config</b>	$r = 8$ , $\alpha = 16$ , dropout = 0.05, bias="none"
<b>Loss / Metric</b>	Same as Phase-1 (BCEWithLogitsLoss + validation accuracy)
<b>Max steps</b>	10,000
<b>Checkpointing</b>	Save <code>best_model_bert.pt</code> on best validation loss/accuracy

### 1.2.3 Testing and Inference Implementation

The testing and inference phase is implemented to ensure scalability and reproducibility across both the text-only and the multimodal branches. For binary classification with BERT/DeBERTa, evaluation is embedded in the training script itself, while for the multimodal BERT and text-only score estimation with QualT5, dedicated testing scripts are provided under the `test/` directory. In all cases, pre-tokenized shards are loaded directly from disk to avoid redundant preprocessing, checkpoints are restored with non-strict weight loading, and batch-level inference is executed with pinned-memory data loaders and large batch sizes to maximize throughput.

Model restoration follows the training setup. In the BERT multimodal branch, the encoder and metadata fusion head are re-instantiated; if LoRA adapters were used during fine-tuning, they are automatically attached to the self-attention projections before loading the state dictionary. In the QualT5 branch, the encoder-decoder is loaded together with its tokenizer to ensure consistency in decoding the target tokens "true" and "false". All models are placed on the selected device (cuda or cpu) and run fully in `no_grad` mode.

Scoring logic differs across branches. In the BERT multimodal case, the classifier outputs a single logit per sample; applying the sigmoid yields a probability  $\hat{p} = \sigma(\text{logit})$ , with hard predictions obtained at a fixed threshold of 0.5. For QualT5 score estimation, only the first decoding step is considered, restricting logits to the tokens corresponding to "true" and "false". A local softmax over these two logits yields the probability of "true", ensuring consistency with the generative training objective. In both cases, hard predictions are obtained via a threshold, while soft probabilities are retained for computing ranking metrics and proper scoring rules.

Outputs are exported as CSV files. For text-only score estimation and multimodal inference, minimal tuples such as `clueweb_id`, `url`, `score` are produced, while binary classification additionally logs the predicted labels. These outputs can optionally be converted into a compact on-disk *QualCache*, a binary in-



dex mapping document identifiers to float quality scores, consisting of two files: `quality.f4` (dense float array) and `docno.npids` (identifier mapping). This conversion is performed through the `pyterrier_quality` library and is particularly useful for static pruning or crawler frontier prioritization.

The following commands illustrate reproducible runs for the testing phase and optional cache conversion:

```
1 python3 only_text_testing.py \  
2   --model pyterrier-quality/qt5-small \  
3   --test_csv dataset_generation/dataset/dataset_test.csv \  
4   --max_length 512 \  
5   --ckpt checkpoints/QualT5_finetuned/final_model \  
6   --batch_size 256 --device cuda
```

Listing 1.9: Only-text (QualT5) — testing on CSV.

```
1 python3 metadata_testing.py \  
2   --model bert-base-uncased \  
3   --test_dir /mnt/ssd_data/tokenized_bert_test \  
4   --chunks 126 \  
5   --ckpt /path/to/best_model_bert.pt \  
6   --batch_size 128 --num_workers 4 \  
7   --device cuda
```

Listing 1.10: Multimodal (BERT + metadata) — testing on pre-tokenized shards.

```
1 python3 qual_cache_build.py \  
2   --csv /path/to/scores.csv \  
3   --outdir /path/to/qual_cache_out \  
4   --chunksize 1000000  
5 # Produces: qual_cache_out/quality.f4 and qual_cache_out/docno.npids
```

Listing 1.11: Optional: convert scores CSV to an on-disk QualCache.