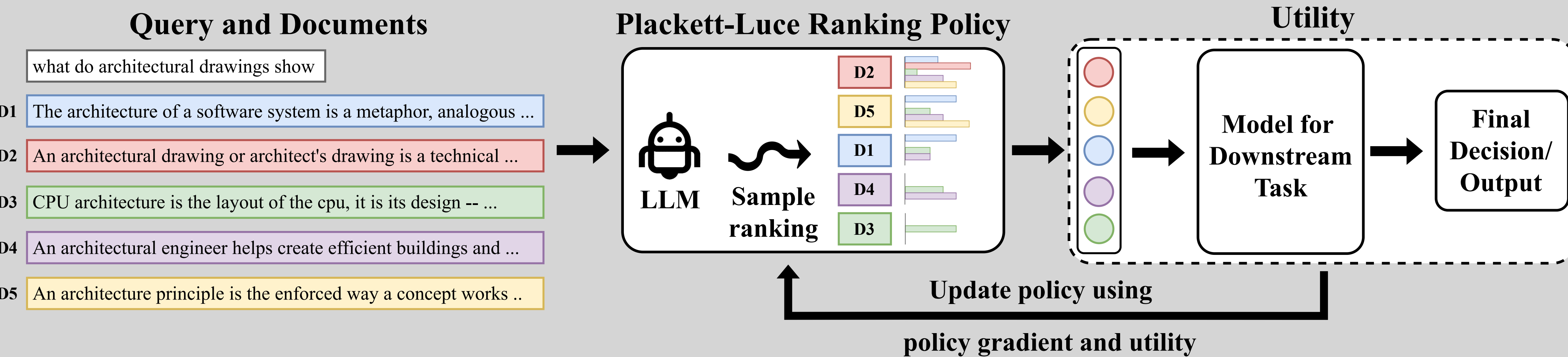# Policy-Gradient Training of Language Models for Ranking

Ge Gao,  Jonathan D. Chang, Claire Cardie, Kianté Brantley, Thorsten Joachims

## Overview

How to train LLM-based retrieval models that directly optimize downstream decision-making quality?
We introduce Neural PG-RANK:

- Learns to rank by instantiating a LLM as a Plackett-Luce ranking policy
- End-to-end training of retrieval models as part of larger pipelines via policy gradient
- Can optimize the ranker for any cardinal loss function evaluating the downstream decisions



## Method

Optimize a Plackett-Luce policy $\pi_\theta(r|q) = \prod_{i=1}^{n} \dfrac{\exp s_\theta(q, d_{r(i)})}{\sum_{j \in \{r(i),\dots,r(n)\}} \exp s_\theta(q, d_j)}$

+ REINFORCE update
+ Monte Carlo sampling with N samples
+ Variance reduction with leave-one-out baseline
+ nDCG@10 as utility function

$$\hat{\nabla}_\theta U(\pi_\theta|q) = \frac{1}{N} \sum_i \Big[ \nabla_\theta \log \pi_\theta(r_i|q)\Big(\Delta(r_i|q) - \frac{1}{N-1}\sum_{j \neq i} \Delta(r_j|q)\Big)\Big]$$

$$= \frac{1}{N}\sum_i \Big[ \sum_k \nabla_\theta \log \pi_\theta(r_{i,k}|q, r_{i,1:k-1})$$

$$\Big( \text{nDCG}(r_{i,k:}|q, r_{i,1:k-1}) - \frac{1}{N-1}\sum_{j\neq i} \text{nDCG}(r_{j,k:}|q, r_{i,1:k-1})\Big)\Big]$$

## Experimental Setup

<u>Data</u>: MS MARCO for training; BEIR for evaluation

<u>Evaluation metric</u>: nDCG@10

<u>Our ranking policy</u>: either SBERT or TAS-B as warmstart, with Neural PG-RANK method as fine-tuning

<u>Comparison systems</u>: supervised learning SOTA bi-encoder models with distilbert-base-uncased

| Method | Source of Negative Docs | | | Additional Supervision | Loss |
|---|---|---|---|---|---|
| | In-Batch | BM25 | Dense Model | | |
| SBERT (Reimers & Gurevych, 2019) | | ✓ | ✓✓✓ | ✓ | MarginMSE + NLL |
| TAS-B (Hofstätter et al., 2021) | ✓ | ✓ | ✓ | ✓✓ | MarginMSE + Distillation |
| SPLADEv2 (Formal et al., 2021) | | ✓ | ✓✓ | ✓ | MarginMSE + Sparsity |
| Neural PG-RANK (Ours) | | | ✓ | | Utility Maximization |

## Second-Stage Reranking

<u>Setup</u>: search over a candidate set of 1k documents per query

<u>In-domain Results</u>:
- Performance gains with both warmstart models

<u>Out-of-domain Results</u>:
- Comparable generalization
- Notable improvements on widely-studied QA datasets
- Weaker in the domain of science and finance

| Dataset | Domain | Comparison Systems | | | Ours: Neural PG-RANK | |
|---|---|---|---|---|---|---|
| | | SBERT* | TAS-B* | SPLADEv2* | with SBERT | with TAS-B |
| MS MARCO dev | misc. | 0.892 | 0.893 | 0.900 | **0.987** | <u>0.982</u> |
| TREC-DL 2019 | misc. | <u>0.743</u> | **0.749** | **0.749** | 0.742 | 0.741 |
| TREC-COVID | bio-medical | **0.764** | 0.711 | <u>0.731</u> | 0.690 | 0.630 |
| NFCorpus | bio-medical | <u>0.308</u> | 0.320 | **0.341** | 0.249 | 0.303 |
| NQ | Wikipedia | 0.836 | 0.836 | 0.854 | <u>0.869</u> | **0.878** |
| HotpotQA | Wikipedia | 0.747 | 0.785 | 0.834 | **0.902** | <u>0.900</u> |
| FiQA-2018 | finance | <u>0.291</u> | 0.279 | **0.342** | 0.131 | 0.139 |
| ArguAna | misc. | 0.351 | 0.479 | **0.480** | <u>0.354</u> | 0.443 |
| Touché-2020 | misc. | **0.480** | 0.423 | <u>0.460</u> | 0.363 | 0.361 |
| Quora | Quora | 0.962 | **0.982** | <u>0.967</u> | 0.963 | **0.982** |
| DBPedia | Wikipedia | 0.513 | 0.513 | **0.533** | 0.521 | <u>0.525</u> |
| SCIDOCS | scientific | 0.144 | <u>0.151</u> | **0.163** | 0.108 | 0.136 |
| FEVER | Wikipedia | **0.931** | 0.911 | <u>0.929</u> | 0.907 | 0.913 |
| Climate-FEVER | Wikipedia | <u>0.442</u> | 0.433 | **0.444** | 0.438 | 0.383 |
| SciFact | scientific | <u>0.597</u> | 0.579 | **0.696** | 0.316 | 0.410 |

## First-Stage Retrieval

<u>Setup</u>: search over all documents

<u>In-domain Results</u>: suboptimal

| Dataset | | Comparison Systems | | | Ours: Neural PG-RANK | |
|---|---|---|---|---|---|---|
| | BM25 | SBERT* | TAS-B* | SPLADEv2* | with SBERT | with TAS-B |
| MS MARCO dev | 0.228 | **0.434** | 0.407 | <u>0.433</u> | 0.416 | 0.401 |