# BEIR: A Heterogeneous Benchmark for Zero-shot Evaluation of Information Retrieval Models

Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, Iryna Gurevych

TECHNISCHE
UNIVERSITÄT
DARMSTADT

**Nandan Thakur**
**UKP**

**Nils Reimers**
**Hugging Face**

**Andreas Rücklé**
**Amazon**

**Abhishek**
**Srivastava, UKP**

**Iryna Gurevych**
**UKP**

**Ubiquitous Knowledge Processing Lab**
**Technische Universität Darmstadt**

**https://www.ukp.tu-darmstadt.de/**

Beir
Benchmarking IR

UBIQUITOUS
KNOWLEDGE
PROCESSING

UKP

# 🍺🍺 What is 🔍 Information Retrieval?

🔍 **Which football club does Lionel Messi play for?**

natural language query

**OR**

🔍 **Messi football club**

keyword-based query

WIKIPEDIA
The Free Encyclopedia

**5.5M Articles**

**Lionel Messi**

Lionel Andrés Messi (born 24 June 1987), also known as Leo Messi, is an Argentine professional footballer who plays as a forward for Ligue 1 club **Paris Saint-Germain** and captains the Argentina national team. Often considered the best player in the world and widely regarded as one of the greatest players of all time, Messi has won a record six Ballon d'Or awards, a record six European Golden Shoes, and in 2020 was named to the Ballon d'Or Dream Team.
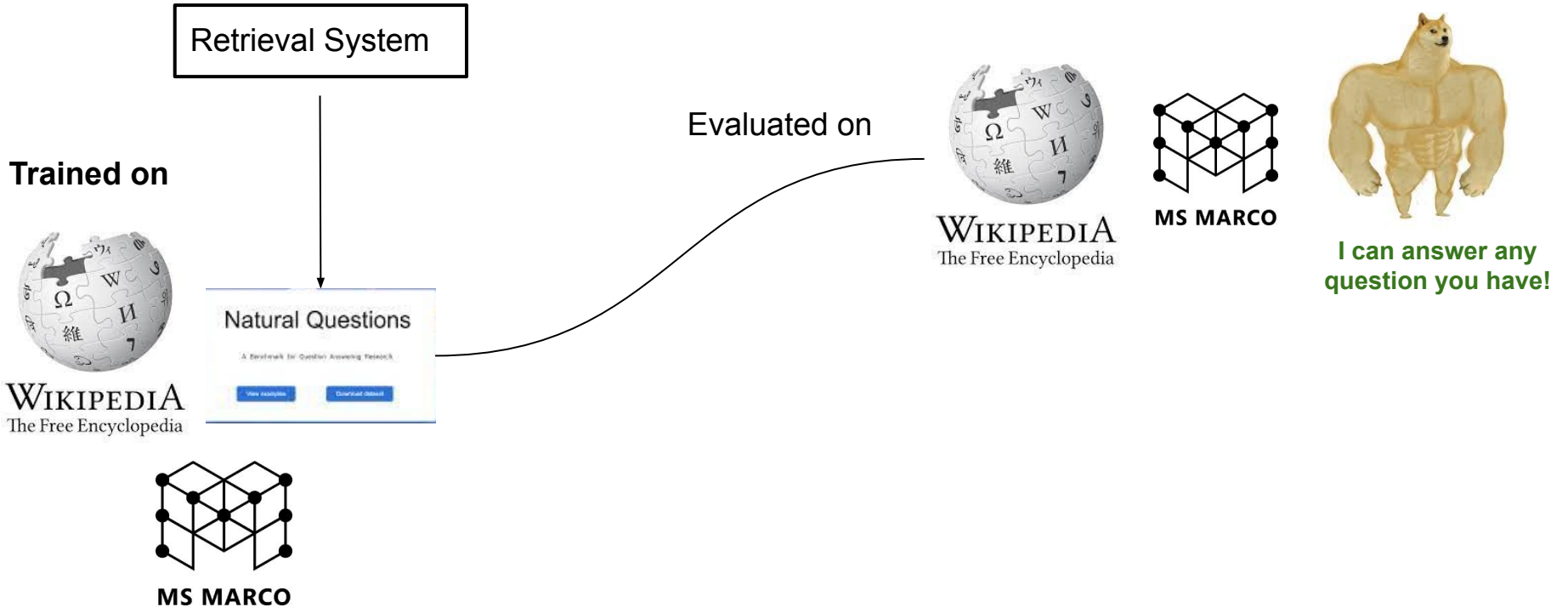
**Ubiquitous**
present, appearing, or found everywhere.

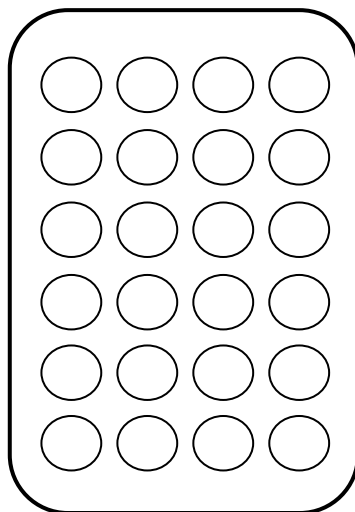# Existing Neural IR Setups

**In-domain** (Training data is available)

Retrieval System

Evaluated on

**Trained on**

WIKIPEDIA
The Free Encyclopedia

Natural Questions

A Benchmark for Question Answering Research

View examples    Download dataset

WIKIPEDIA
The Free Encyclopedia

MS MARCO

WIKIPEDIA
The Free Encyclopedia

MS MARCO

**I can answer any question you have!**

MS MARCO

# Annotating Training Data is expensive!

**No** Annotation Reqd.

**Lots** Annotation Reqd.

**Few** Annotation Reqd.

**Labeled Test Data**
Typically in **~100** pairs

**Labeled Training Data**
Typically in **~100k** pairs

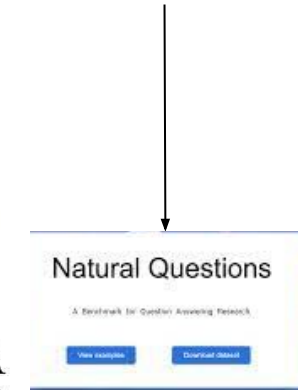**Unlabeled Data**
Typically in ~Millions

# Do these retrieval models generalize?

**In-domain** (Training data is available)

Retrieval System
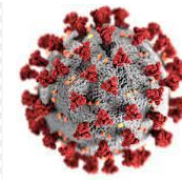
Evaluated on

**Trained on**

**I can answer any question you have!**

Natural Questions

A Benchmark for Question Answering Research

**Out-of-domain** (Training data is unavailable)
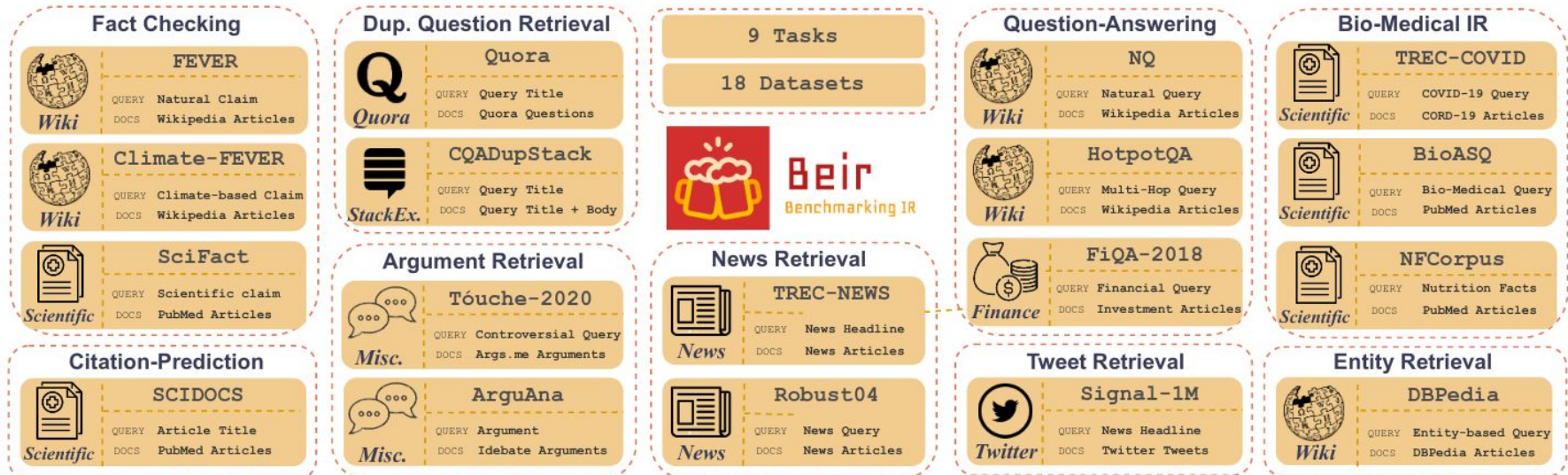
Evaluated on

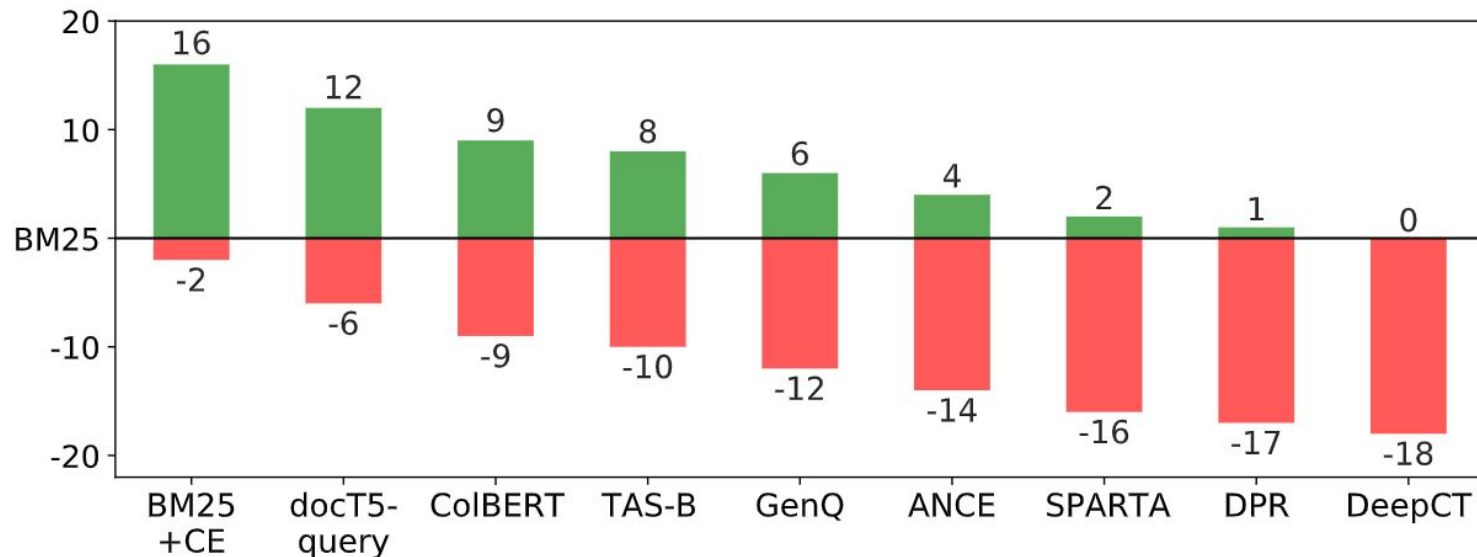**I'm sorry, I do not understand your question!**

- Robust evaluation of a retrieval system across 18 diverse datasets and 10+ domain types

- Contains datasets covering broad topics (like Wikipedia) and specific domains (COVID-19)

- Contains datasets with different corpus sizes (3k to 15Mil), query and document sizes, and different text types (Tweets vs. News articles)

# Zero-shot Results on the BEIR Benchmark

| Lexical | Sparse | Dense | Late-Interaction | Reranking |
|---------|--------|-------|------------------|-----------|
| BM25 (Anserini) | DeepCT, SPARTA, DocT5query | DPR, ANCE, TAS-B, GenQ | ColBERT | BM25+ CE (MiniLM) |

# To Summarize

**Motivation for creating the BEIR Benchmark**

- Existing neural information models have been studied in limited or narrow settings.

- To robustly evaluate model generalization, we propose a zero-shot retrieval benchmark.

- The BEIR benchmark contains over 18 publicly available datasets for evaluation, spanning across 10 different retrieval tasks and domains.

**Experimental results of diverse retrieval architectures on BEIR**

- Generalization with models is quite a difficult task and there is no free lunch!

- In-domain performances cannot be a good indicator for zero-shot performances.

- BM25 is a robust baseline, and performs competitively across several zero-shot datasets.

- Cross-Encoders or rerankers achieve the best zero-shot performances, but are slow at inference.

- Dense retrievers and sparse models suffer from out-of-distribution generalization.

# 🍺 Thank You!

- If you liked our work on BEIR benchmark, you can find more details in the GitHub repository.

- We actively maintain a leaderboard with diverse models and their zero-shot retrieval scores.

- For more interesting results, we would suggest you to read our NeurIPS publication.

### I look forward to meet you virtually and answer your questions at NeurIPS'21

📖 **ukplab/beir**

A Heterogeneous Benchmark for Information Retrieval. Easy to use, evaluate your models across 15+ diverse IR datasets.

🔵 Python   ★ 298   ⑂ 40

deepset

**BEIR: A Heterogeneous Benchmark for Zero-shot Evaluation of Information Retrieval Models**

Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, Iryna Gurevych
Ubiquitous Knowledge Processing Lab (UKP-TUDA)
Department of Computer Science, Technische Universität Darmstadt
www.ukp.tu-darmstadt.de
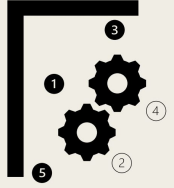
**ML Helpful Libraries**

Beir
Benchmarking IR

NLU & IR:
NEURAL IR (III)

Omar Khattab

CS224U: Natural Language Understanding
Spring 2021

Google
colab

https://colab.research.google.com/drive/1HfutiEhHMJLXiWGT8pcipxT5L2TpYEdt?usp=sharing

UKP