# Evaluating Temporal Persistence of Information Retrieval Models at CLEF 2024 LongEval-Retrieval Track
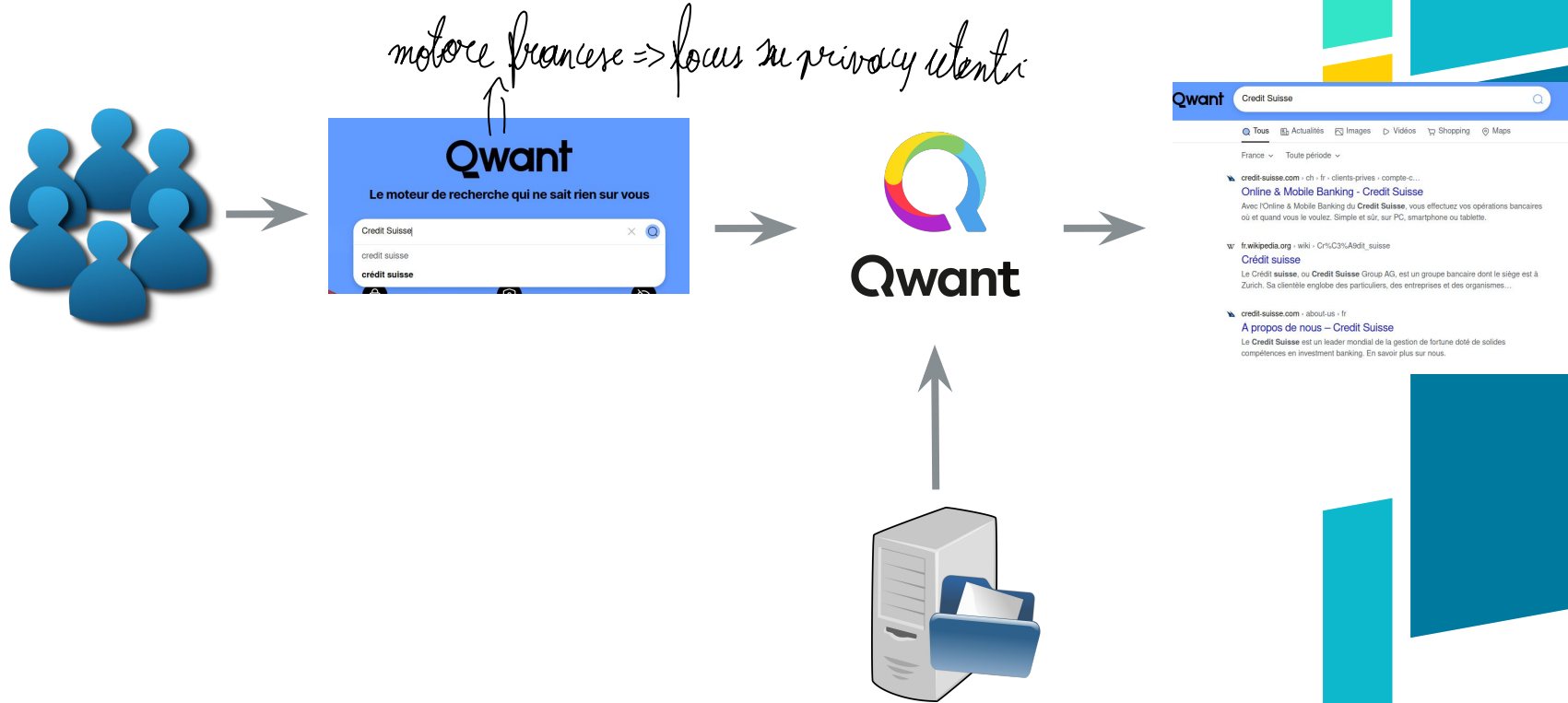
Petra Galuščáková
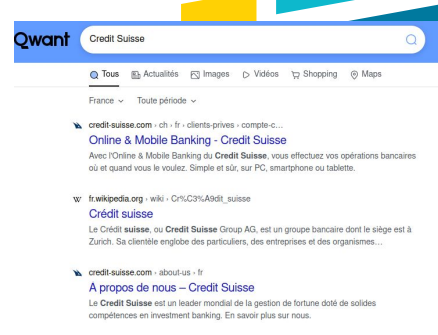petra.galuscakova@uis.no

University of Stavanger

5/4/2024

# Outline

1) **Evaluating Temporal Persistence**
2) LongEval-Retrieval CLEF Campaign
   a) Task
   b) Collection
   c) Participation
3) What (did not) work in 2023?
4) How to evaluate relevance?

# Web Retrieval



*motore francese => focus su privacy utenti*

# Evolving Collection (Queries)
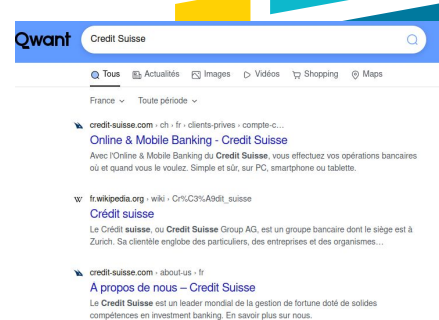
*ricerche e collezioni cambiano con tempo*



**March 20, 2023:**
1) Motion de censure
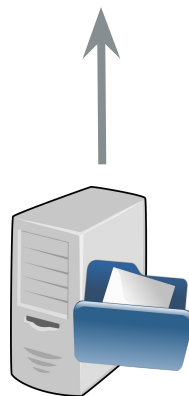2) Assemblée nationale
3) Printemps
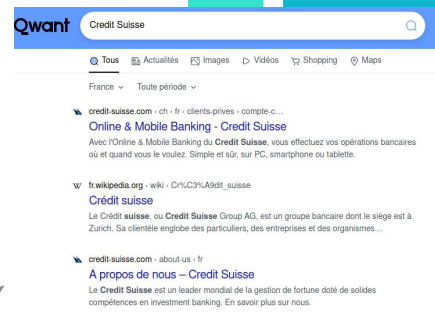4) Grand Prix f1
5) Rennes

# Evolving Collection (Queries)



**March 21, 2023:**
1) Paul Grant
2) Muriel Robin Pierre Palmade
3) Bruno Martini
4) Actualite
5) Robbie Williams

# Evolving Collection (Documents)



**March 20, 2023:**
1) Motion de censure
2) Assemblée nationale
3) Printemps
4) Grand Prix f1
5) Rennes

# Changes in the IR Framework

- Topics
  - Trends in the topics
- Documents
  - Documents are added to the index
  - Content of the documents changes
  - Documents are removed from the index
- Relevance assessments
  - Trends is general interests (what is considered to be relevant)
  - Subjective relevance for users
- …

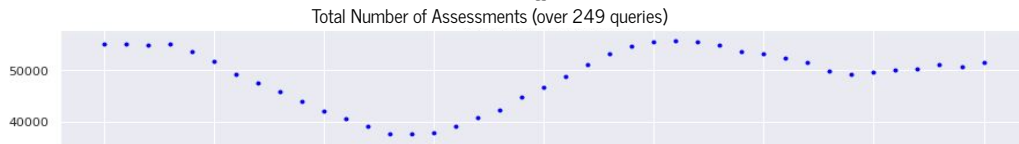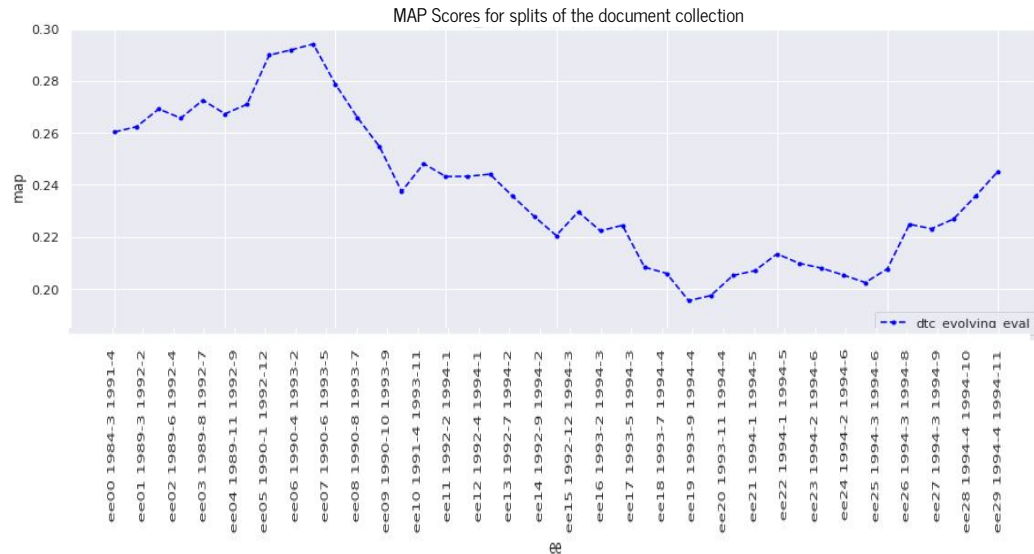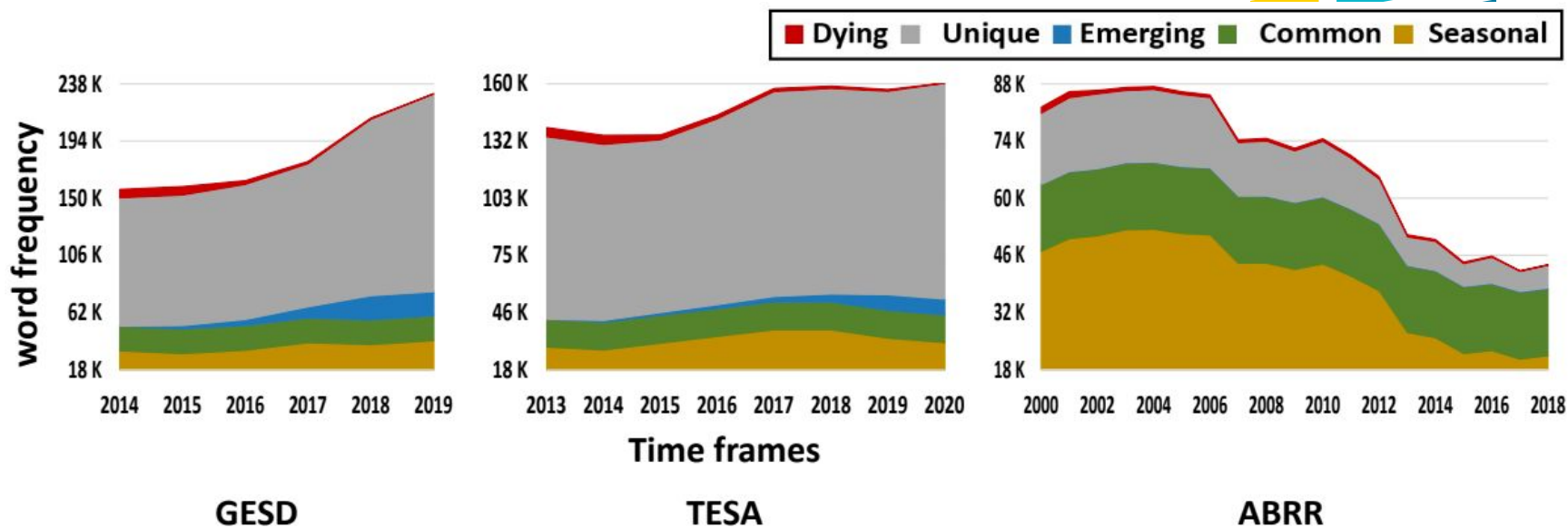# Changes in the IR Framework

- Topics
  - Trends in the topics
- Documents
  - **Documents are added to the index**
  - Content of the documents changes
  - **Documents are removed from the index**
- Relevance assessments
  - Trends is general interests (what is considered to be relevant)
  - Subjective relevance for users
- …

# Evolution and IR Performance



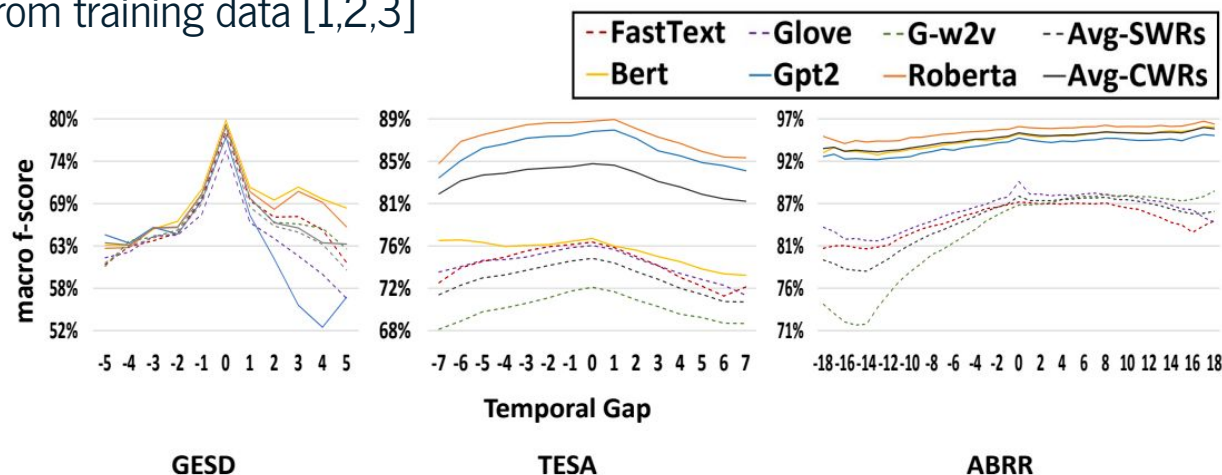MAP Scores for splits of the document collection

Total Number of Assessments (over 249 queries)

Total Number of Relevant Documents (over 249 queries)

TREC Robust 04 collection (splits of ⅛ of the collection, with 1/10 of the split shift), BM25 run

# Evolving Vocabulary



Temporal usage of different word types [1]

[1] Rabab Alkhalifa et al: Building for tomorrow: Assessing the temporal persistence of text classifiers, 2022

# Evolution and Neural Models

▸ The performance of **neural models** in **Text Classification** systems **drops over time** as patterns observed in data change, due to linguistic and societal changes.

▸ This drop is more pronounced when the testing data is **further away** in time from training data [1,2,3]



Temporal performance of different language representations across the three datasets. [1]

[1] Rabab Alkhalifa et al: Building for tomorrow: Assessing the temporal persistence of text classifiers, 2022
[2] Komal Florio et al.: Time of Your Hate: The Challenge of Time in Hate Speech Detection on Social Media, 2020
[3] Jan Lukes and Anders Søgaard: Sentiment analysis under temporal shift, 2018
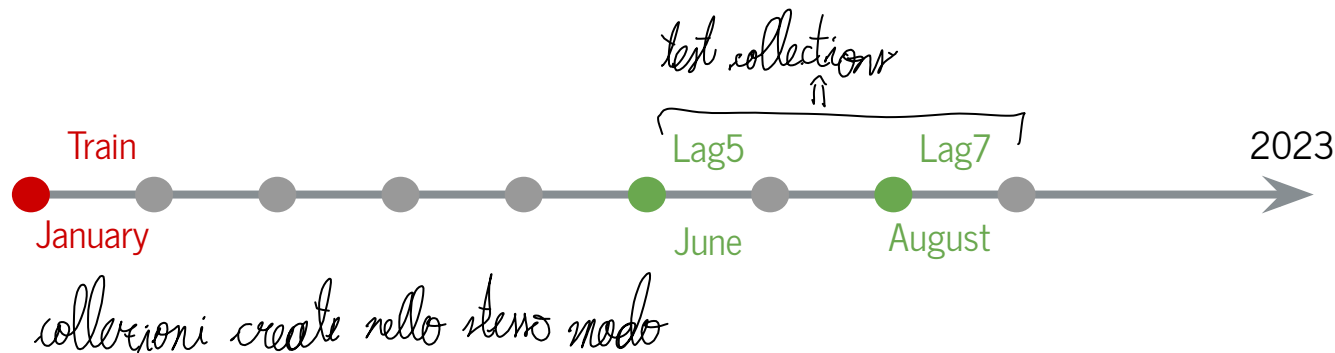
# Outline

# Questions

1) How does **search engine** behave as the **collection evolves**?

2) What IR systems are the **most robust** against the **changes** in the collection (and still perform well)?

3) When do we need to **update an IR** system as the **collection evolves?**

# LongEval-Retrieval CLEF Campaign

▸ Build a **succession of train/test collections**

▸ **Each** of them composed of a set of **documents** from Qwant's
index **queries** from Qwant's users

▸ Designed to reflect the **changes** of the Web across time



test collections

Train                                                                2023
January                June          August

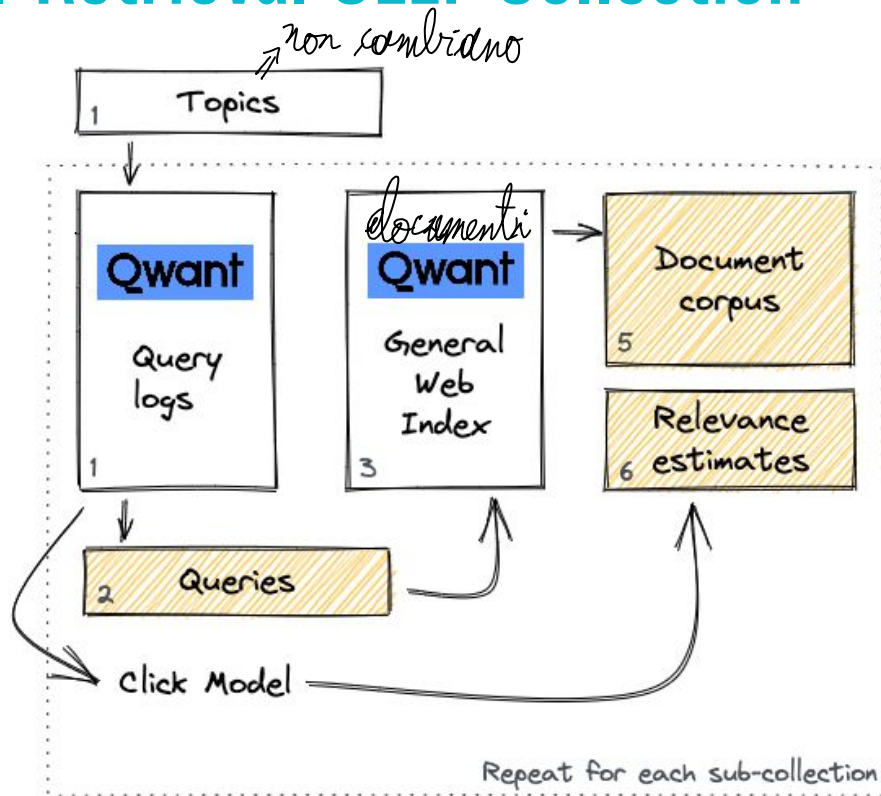collezioni create nello stesso modo

# LongEval-Classification CLEF Campaign

*non il nostro*

▸ Explore temporal persistence on **classification of sentiment of Tweets**

▸ TM-Senti Dataset:

  ▹ Large-scale Tweets sentiment dataset

  ▹ Tweets are binary labelled for sentiment as either "positive" or "negative"

▸ Long-term and short-term subtasks

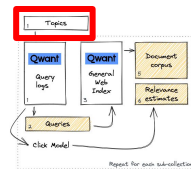▸ Evaluation: Macro-averaged F1-score and Relative Performance Drop

# Outline

1) Evaluating Temporal Persistence
2) **LongEval-Retrieval CLEF Campaign**
   a) Task
   b) **Collection**
   c) Participation
3) What (did not) work in 2023?
4) How to evaluate relevance?

16

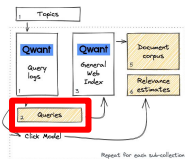# LongEval-Retrieval CLEF Collection

# Topics

▸ Selected to ensure a balance between <u>balance between</u> **popularity**, **stability**, **generality**, and **diversity**.

▸ Selected <u>from the Web and social media based on</u> **trends**.

▸ Performed <u>only</u> **once** for the entire LongEval-Retrieval collection.

| | Topic | English Description |
|---|---|---|
| 1 | eau | water |
| 2 | nourriture | food |
| 3 | espace | space |
| 4 | voiture | car |
| 5 | argent | money |
| 6 | manifestation | protest |
| 7 | virus | virus |
| 8 | terre | earth |
| 9 | énergie | energy |
| 10 | police | police |

...

# Queries



▶ In the context of the LongEval-Retrieval test collection, a query is a **multi-word chain** of characters that is **related** to one or more **topics**.

▶ Topics are used for filtering and selecting the Qwant actual user's queries.
  ▷ Select the queries which **contain topic words**
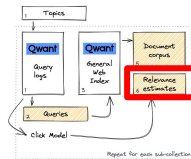
▶ Extracted for **each** sub-collection

# Translations

▸ Initially, all the queries and documents were created in French.

▸ French-English translation system CUBBITT:

   ▷ https://lindat.mff.cuni.cz/services/translation

▸ The quality of the translations of the queries is much lower than the quality of the document translation.

▸ **Multiple** possible **translations** provided for the queries

# Query Translations

| | | | | | |
|---|---|---|---|---|---|
| q01235 | **a chapeau** | a hat | a Chief | a chapeau | a chapeau |
| q01236 | **a quoi sert la prise de terre** | What is the purpose of the land grab? | What is the use of the land grab? | a what is the purpose of the land grab | a what is the purpose of the land grab? |
| q01238 | **abonnement ter bordeaux** | subscribing bordeaux | subscribing to bordeaux | subscribing to the bordeaux | bordeaux subscription |
| q012319 | **affutage couteaux** | sharpening knives | knife sharpening | knife-sharpening | Knife sharpening |
| q012331 | **aide a la reprise d'emploi pole emploi** | senior employment assistance | Senior Employment Assistance | Senior employment assistance | assistance senior employment |

possibile usare sia una che tutte traduzioni

# Relevance Assessments

- Based on user **clicks** ⇒ *creati automaticamente*
- Raw clicks cannot be used:
  - ▷ Noisy and biased
- Use click **models**
  - ▷ Goal: better model the behaviour of users
    *e conservare privacy*

- 0 = not relevant, 1 = relevant, 2 = highly relevant

# Documents

▶ 2 approaches to selection:

  ▷ Documents that were **displayed in SERPs** for the queries that we selected

  ▷ Potentially non-relevant documents are **randomly sampled** from Qwant index in order to better represent the nature of a Web test collection (~100,000 documents per topic)

▶ Web documents, **downloaded and cleaned**

*risultati ricerche*

# Documents Example

| Highly Relevant Document | Non-Relevant Document |
|---|---|
| Bordeaux Airport - Merignac - Official website Re-opening of the billi terminal: all the information you need for your next trip Our billi terminal has reopened to operate some easyJet and Ryanair flights. Shops were also reopened, in particular to offer a catering service to travellers. This decision will make it possible to deploy on the entire infrastructure the material and human resources necessary for the smooth running of the journey. Important: All other airlines: concourse A Flight boarding: check your concourse and boarding gate in real time on our website or on the terminal screens when you arrive. With large crowds, think ahead of your visit: Due to the heavy summer traffic, passengers are advised to anticipate their arrival well in advance. Wings for the World News and events Published: 10 June 2022 | Parking near Bordeaux Airport Merignac \| Beep Valet Parking contact@beep-valet-parking.com BEEP VALET Choose your seat type Parking secure + Shuttle offered + Shuttle Airport You are brought to the airport from the car park Reservation Airport Free Parking + Shuttle Parking Available Your windscreen needs replacing? We offer parking Reservation Parking Offered Secure payment by bank CIC Benefits of your car park Merignac airport RESERVE RAPIDE GAREZ EASY TALK SEREIN Close airport Mérignac Our tarmac car park is located 3 km from Bordeaux-Mérignac airport. Take 5 minutes of transfer time in our comfortable shuttle. If needed, child boosters and baby seats are offered. Airport Shuttle |

Query "**aeroport bordeaux"**, which was translated as "**airport**"

# Train and Test Sets

Train — January

Lag5 — June

Lag7 — August

2023

- **Train collection**: ~2M Web pages, 30k queries, 133k assessments (60% are non-relevant, 19% are relevant, and 21% are highly relevant)
- **Lag5 Test Collection**: ~3M Web pages, ~20k queries
- **Lag7 Test Collection**: ~3M Web pages, ~25k queries

6% were highly relevant in 2023

# Evaluation

Lag5      Lag7      2023

January      June      August

- **nDCG** (Lag5 and Lag7)
- **Relative nDCG Drop** between Lag5 and Lag7

# Outline

1) Evaluating Temporal Persistence
2) **LongEval-Retrieval CLEF Campaign**
   a) Task
   b) Collection
   c) **Participation**
3) What (did not) work in 2023?
4) How to evaluate relevance?

# Registration to CLEF 2024



CLEF Web site: https://clef2024.imag.fr

# Data Download

# Data Download



Austrian Science Fund

Project code: I4471-N

Project name: Kodicare

**Subject(s)**    information retrieval   parallel corpus   search   automatic evaluation

**Collection(s)**    LINDAT / CLARIAH-CZ Data & Tools

Show full item record

📎 Files in this item

This item is **Publicly Available** and licensed under:
Qwant LongEval Attribution-NonCommercial-ShareAlike License

| | |
|---|---|
| **Name** | longeval-train-v2.tgz |
| **Size** | 11.7 GB |
| **Format** | application/x-gzip |
| **Description** | data |
| **MD5** | e34cf8b5e9b2de98628759bbd621a4ca |

⊙ Download file

# Data Download

*loggare per scaricare*

# Submission Requirements

▸ **For each system submitted, the participants need to submit the results acquired by this system on both Lag5 and Lag7 queries**

▸ Participants also need to provide a short **description** of each of the submitted systems.

▸ Each team can submit up to **5 systems**.

# Submission Format for a System

Each system should be submitted in a **single zipped file**
consisting of a following tuple:

- team_system.L5 (Lag5 test set)
- team_system.L7 (Lag7 test set)
- team_system.meta (run description)

=> codice deve generarli, poi carichiamo anche su TIRA

system.** is a **TREC run** in TREC format (max 1000 documents)

If run is BM25, team UGA, in-time set: **UGA_BM25.L5**

# Submissions

▸ Will be done using **git**

▸ Each team will have its own private repository created on bitbucket

# LongEval at CLEF 2024

▸ **April:** Data release
▸ **May:** Participants' submissions (Still time to register and participate !)
▸ **June:** Participants' papers submissions
▸ **June:** Evaluation results release
▸ **July:** Camera ready paper submissions
▸ **September:** CLEF Conference in Grenoble

LongEval Web site: https://clef-longeval.github.io

# Outline

1) Evaluating Temporal Persistence
2) LongEval-Retrieval CLEF Campaign
   a) Task
   b) Collection
   c) Participation
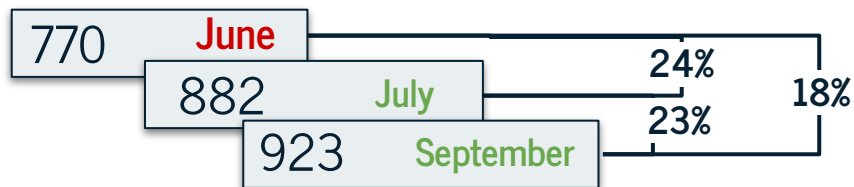3) **What (did not) work in 2023?**
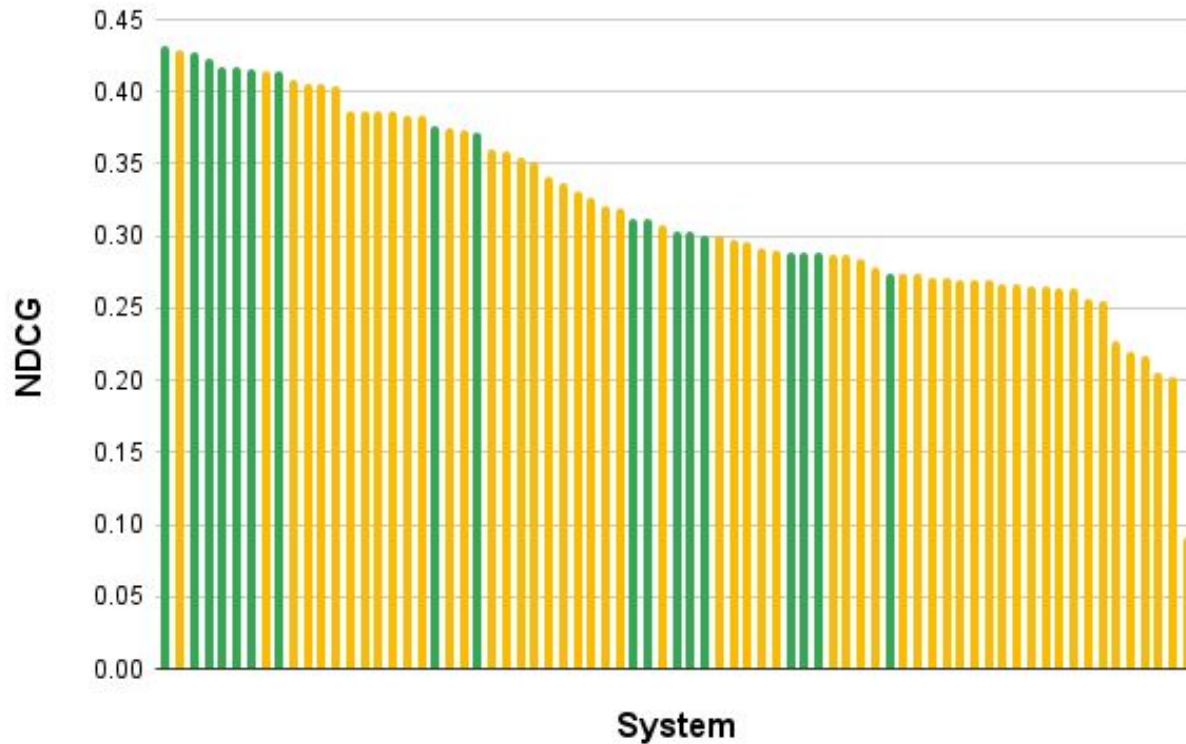4) How to evaluate relevance?

# 2023 Collection



Train / Heldout
In-time Test

Short-term
Test

Long-term
Test

2022

June    July    August    September

▶ **Documents**

| | |
|---|---|
| 1,570,734 **June** | |
| 1,593,376 **July** | 97% |
| 1,081,334 **September** | 96% |

94%

▶ **Queries**

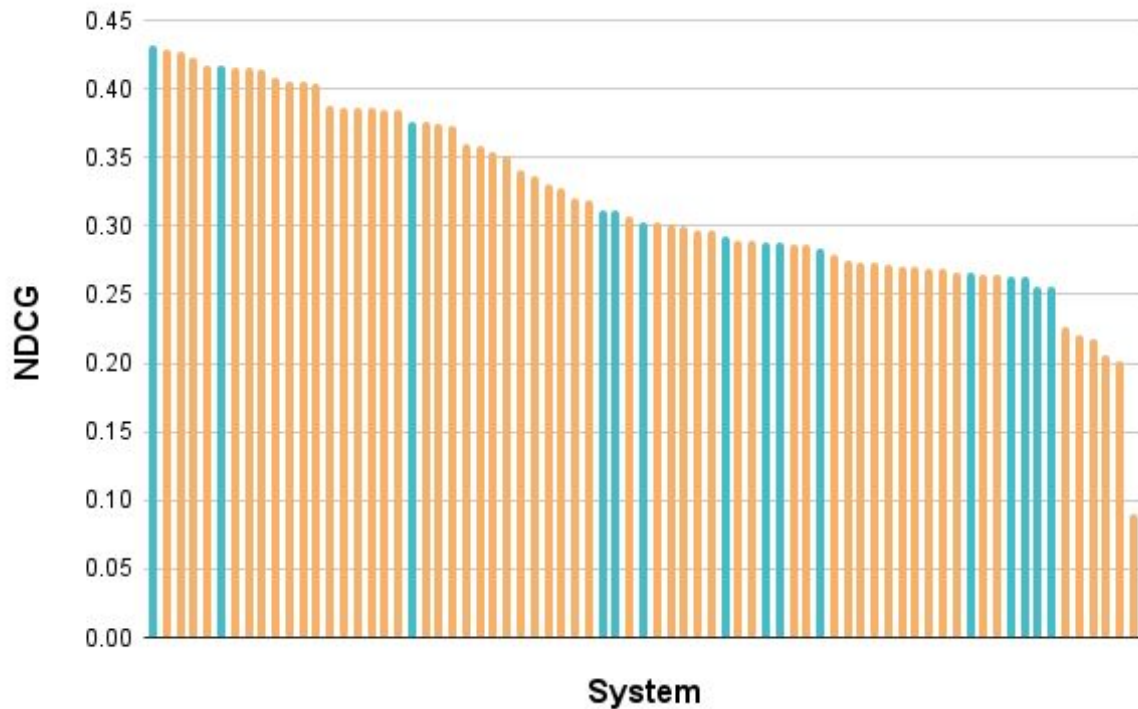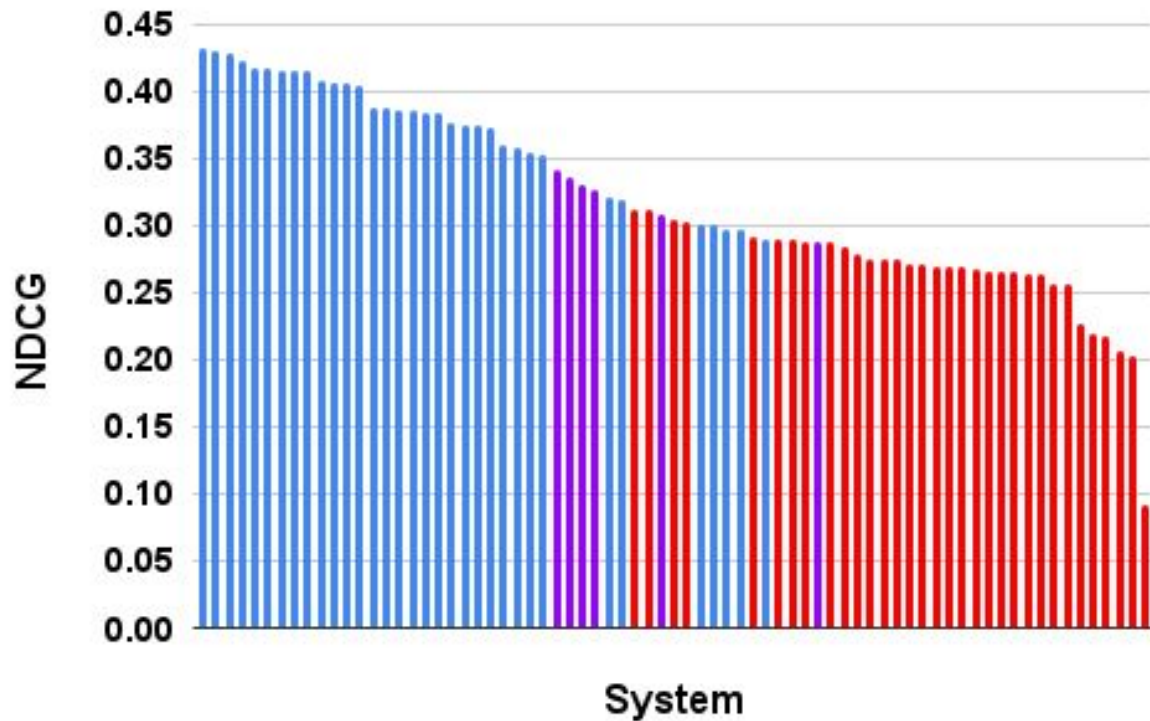| | |
|---|---|
| 770 **June** | |
| 882 **July** | 24% |
| 923 **September** | 23% |

18%

# Neural vs Statistical Models



Performance of neural (green) and non-neural (yellow) approaches on In-Time test

# Single System vs. Combination



Performance of single approach (orange) and combinations (cyan) on In-Time test

# French vs. English



Performance of systems which use French (blue), English (red) or both languages (purple) on In-Time test
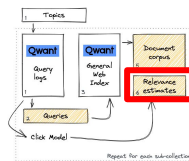
# Ideas for 2024

- Use overlaps and changes between data
- Use larger training data
  - Fine-tune models on train data
- Use n-best translations
- Make mixture of models effective
- …

# Outline

1) Evaluating Temporal Persistence
2) LongEval-Retrieval CLEF Campaign
   a) Task
   b) Collection
   c) Participation
3) What (did not) work in 2023?
4) **How to evaluate relevance?**

# Relevance Assessments



- ▶ Based on user **clicks**
- ▶ Raw clicks cannot be used:
  - ▷ Noisy and biased
- ▶ Use click **models**
  - ▷ Goal: better model the behaviour of users

- ▶ 0 = not relevant, 1 = relevant, 2 = highly relevant

**We also aim to acquire manual relevance judgements**

https://dscapp01.researchstudio.at/

# DocTAG



**Topic:** q06221312    **Annotated:** 0    **Total:** 27    ⬇ Download

admin    Logout ➔

**TOPIC**    7 - q06221312    ‹  ›

Labels    Passages    Linking    Concepts

**Title:**    groupama customer space

This is the default configuration. Choose an annotation type and start the annotation.

**Description:**

**DOCUMENT**    1 - doc062200100041    ‹  ›

**Documents' order:**    Lexicographical    Annotated docs

**last update:**

↗  🪄

space client groupama keyword analysis keyword research people searched espace client groupama also searched keyword cpc pcc volume score search results related espace client groupama search engine copyright rights reserved
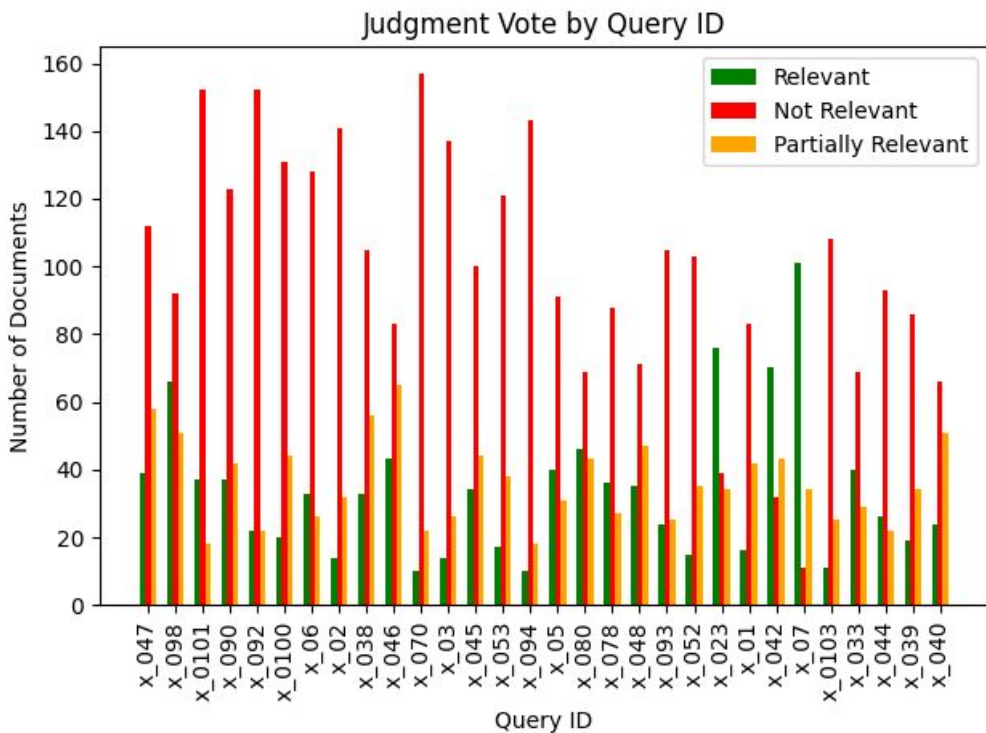
# DocTAG

# Document Relevance

- If you are a user, you search using the query terms of the topic and find the Web page, **would you stop your search** and find the search **successful**?
  - ▷ Yes -> **Relevant**.
  - ▷ Continue searching -> **Not Relevant**
- If you found some information relevant to your query, but you would still continue searching
  - ■ -> **Partially Relevant**
- The labeling option '**I don't know**' should be used only if it is not possible to do any assessment.

# Human Relevance Judgments in 2023

- 37 judges
- 150 queries
- ~130 jud/q
- 14,953 jud.



Judgment Vote by Query ID

# LongEval-Retrieval Organizers

Romain Deveaud; Qwant

Alaa El-Ebshihy; Research Studios Austria

Gabriela Nicole Gonzalez Saez; Université Grenoble-Alpes

Lorraine Goeuriot; Université Grenoble-Alpes

David Iommi; Research Studios Austria

Philippe Mulhem; Université Grenoble-Alpes

Florina Piroi; TU Wien

Martin Popel; Charles University

# Question?

https://clef-longeval.github.io/

longeval-ir-task@univ-grenoble-alpes.fr