

# **From Shared Tasks to Benchmark Datasets: Lessons from NTCIR's 26-Year Journey**

## **Part 2: Participants, Experiments, and Key Findings**

**Makoto P. Kato** — National Institute of Informatics / University of Tsukuba

- Associate Professor at the University of Tsukuba and National Institute of Informatics (NII), Japan
- Research interests: **information retrieval**, information access, evaluation methodologies, user behavior analysis
- Co-organizer of multiple NTCIR tasks, including INTENT, 1CLICK, IMine, MobileClick, OpenLiveQ, and Data Search
- Active participant in the NTCIR community
  - PC Co-chair for NTCIR-12 to NTCIR-15
  - General Co-chair for NTCIR-16 to NTCIR-19
- Also served as:
  - SIGIR 2023 PC Co-chair
  - SIGIR-AP 2026 PC Co-chair



# What this hour covers

---

1. **Breadth of NTCIR tasks** — what we have evaluated, and where NTCIR is distinctive
2. **Spotlight on tasks I co-organized** — INTENT, 1CLICK, OpenLiveQ, DataSearch
3. **Lessons learned**
4. **NTCIR vs TREC vs CLEF** — a General Chair perspective
5. **Running an NTCIR task** — a Program Committee Co-chair perspective

Part 1 covered the *what* and *how* of NTCIR. **Part 2 is about what the community learned** — more than 25 years, from dozens of tasks, and hundreds of teams.

The most durable contribution of an evaluation campaign is rarely a single winning system. It is the **tasks, test collections, and metrics** that outlive every submission.

# **1. Breadth of NTCIR Tasks**



- **Cross-Lingual IR:** Evaluates document retrieval across multiple languages.
- **CLQA:** Finds answers to questions across different languages.
- **Patent Retrieval:** Searches prior art in patent document collections.
- **Patent Machine Translation:** Evaluates specialized machine translation for patent documents.
- **Web Search:** Evaluates Web document retrieval and ranking quality.
- **GeoTime:** Handles search involving geographic and temporal information.
- **Temporalia:** Handles search queries with temporal intent.
- **Math:** Evaluates retrieval of formulas and mathematical documents.
- **Spoken Document Retrieval:** Searches spoken documents using speech recognition results.

## Examples of NTCIR tasks (2/2)

---

- **Recipe Search:** Handles search and recommendation of cooking recipes.
- **Lifelog:** Searches personal lifelog images and records.
- **Data Search:** Evaluates discovery and retrieval of datasets.
- **Dialogue Evaluation (STC, DialEval):** Evaluates the quality of dialogue system responses.
- **MedNLP:** Handles medical text processing and diagnostic support.
- **FinArg / FinNum:** Analyzes numerical expressions and arguments in financial documents.
- **AEOLLM:** Evaluates automatically generated reports by large language models.
- **FairWeb:** Evaluates fairness in Web search results.
- **Tip-of-the-Tongue:** Searches for targets from vague or partial memories.

# What is distinctive about NTCIR's portfolio

---

- **Asian-language IR** — many tasks focusing on multiple languages including Japanese, Chinese, Korean, and English
- **First, or among the first, to evaluate:**
  - **Patent search and patent MT** — NTCIR-3 (2002)
  - **Math formula retrieval** — NTCIR-10 (2013)
  - **Lifelog retrieval** — NTCIR-12 (2016)
  - **Mobile / one-click access** — NTCIR-9 (2011)
  - **Dataset search** — NTCIR-15 (2020)

## **2. Spotlight: NTCIR Tasks I Co-Organized**

**NTCIR-9, 10 INTENT** — search intent diversification (subtopic mining + diversified ranking)

**NTCIR-9, 10 1CLICK** — direct, immediate textual answer to a query

**NTCIR-11, 12 IMine** — unified intent mining; vertical-aware diversification

**NTCIR-11, 12 MobileClick** — two-layered summary for mobile zero-click access

**NTCIR-13, 14 OpenLiveQ** — online evaluation against real Yahoo! Chiebukuro users

**NTCIR-15, 16 Data Search** — first IR campaign for dataset retrieval (e-Stat / Data.gov)

These are the parts of NTCIR I can answer in depth

**INTENT**

# INTENT — Ambiguous queries

A query *"apple"* is an **ambiguous** query that can be interpreted in many ways:



# INTENT — Underspecified queries

A query "*harry potter*" is an underspecified query:


Underspecified query

harry potter


Harry Potter films!



Harry Potter books!

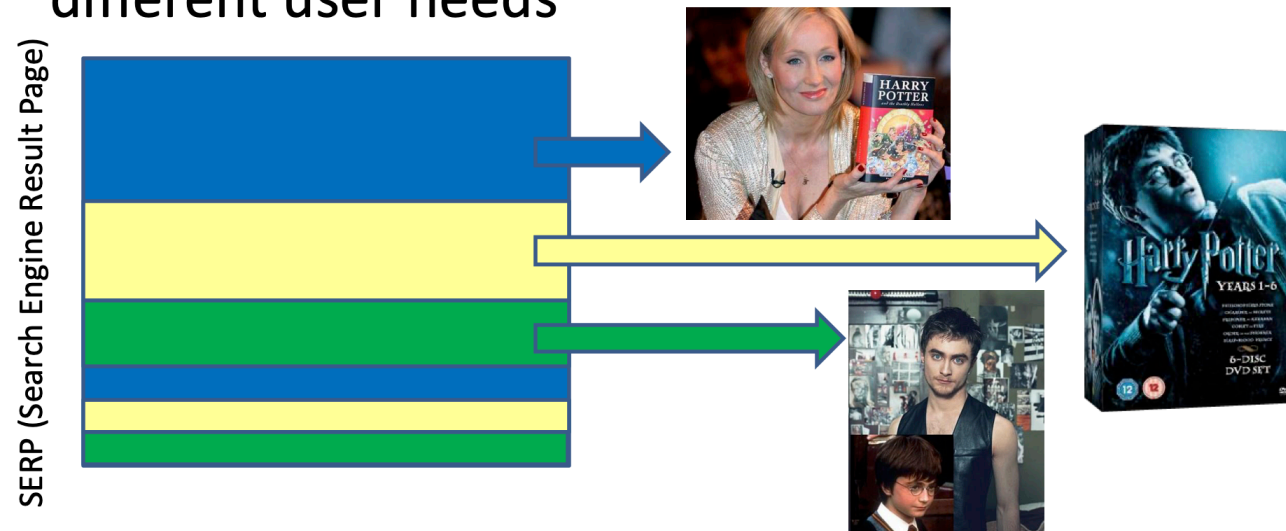


Harry Potter the character!



## Why Search Result Diversification?

- If query  $q$  has multiple **intents**  $i$  and we have no knowledge of the user, we should cover different  $i$ 's in the first SERP to accommodate different user needs



Systems must:

1. **Subtopic Mining (SM)** — return a ranked list of subtopic strings covering the query's intents
2. **Document Ranking (DR)** — return a diversified ranked list satisfying multiple intents

- **Data**

- Bilingual task (Chinese and Japanese)
- Intents were created by manual clustering of submitted subtopic strings
- Documents were judged for relevance to each intent on a 3-point scale (irrelevant, relevant, highly relevant)

# Intent voting interface

NTCIR Query Labeling Tool

> vote subtopic:イタイイタイ病 Welcome,JA\_user1! [logout](#)

Subtopics

	subtopic	subtopic description
<input type="checkbox"/>	イタイイタイ病 写真	イタイイタイ病 写真#459,イタイイタイ病 画像#460,
<input checked="" type="checkbox"/>	イタイイタイ病 原因	イタイイタイ病 原因#449,イタイイタイ病 カドミウム#451,
<input type="checkbox"/>	イタイイタイ病について	イタイイタイ病について#453,イタイイタイ病とは#462,
<input checked="" type="checkbox"/>	イタイイタイ病 症状	イタイイタイ病 症状#450,
<input checked="" type="checkbox"/>	イタイイタイ病 対策	イタイイタイ病 対策#452,
<input checked="" type="checkbox"/>	イタイイタイ病 裁判	イタイイタイ病 裁判#454,
<input type="checkbox"/>	四大公害病 イタイイタイ病	四大公害病 イタイイタイ病#455,
<input type="checkbox"/>	イタイイタイ病 歴史	イタイイタイ病 歴史#456,
<input type="checkbox"/>	イタイイタイ病 患者数	イタイイタイ病 患者数#457,
<input type="checkbox"/>	神通川 イタイイタイ病	神通川 イタイイタイ病#458,
<input type="checkbox"/>	公害 イタイイタイ病	公害 イタイイタイ病#461,

© 2011 - Privacy

10 assessor's votes used for estimating intent probabilities

# Per-intent graded relevance assessment interface

Query: QID -41 / 100 : 花样年华

Start Query: 41 Start Url: 0 GO

Meanings Operation:

花样年华【电影】(such as 电影花样年华 花样年华电影...)

Relevant  Highly Relevant

花样年华【电视剧】(such as 电视剧花样年华 免费电视剧花样...)

花样年华【音乐】(such as 花样年华音乐 花样年华插曲...)

花样年华【网络论坛、博客】(such as 花样年华论坛 花样年华时...)

花样年华【小说、期刊、文章】(such as 花样年华电子书 潘金莲...)

花样年华【在线视听、下载】(such as 花样年华在线观看 在线观...)

花样年华【导演、演员、工作人员】(such as 花样年华梁朝伟 花...)

花样年华【其他视频、日剧、韩剧、舞剧】(such as 芭蕾舞剧 花...)

花样年华【图片写真】(such as 花样年华其他图片 花样年华图...)

花样年华【简介资料】(such as 花样年华基本信息 花样年华资...)

Relevant  Highly Relevant

花样年华【解释、年龄】(such as 花样年华是什么意思 花样年...)

花样年华【其他品牌】(such as 花样年华 品牌 花样年华儿童家...)

花样年华【楼盘小区】(such as 鄞陵花样年华 鄞陵+花样年...)

花样年华【重庆同性恋】(such as 重庆花样年华 花样年华重庆...)

花样年华【游戏】(such as 阿达连连看 4.32花样年华破解版 阿...)

Others(such as 安卓美达 信息 麦田圈 外星人...)

**Intents**

Create Edit

Meaning:

load image load html UrlID - 0 / 199 : htm\abf455d587714341-3d807783d4cca6e0.htm

sina 影音娱乐世界 影音娱乐 | 新浪首页 | 导航 | 请输入关键词 搜索 Google 网

sina 影音娱乐 新浪首页 > 影音娱乐 > 电影宝库 > 正文

图片精选

《花样年华》

导演: 王家卫  
主演: 梁朝伟、张曼玉  
摄影: 杜可风、李屏宾  
美术: 张叔平

身处遥远的异国, 周慕云仍无法忘记过去与苏丽珍之间的种种。如果当天她真的答应跟他会不会还在一起? 抑或注定分离, 各分东西? ... >>点击查看详细介绍 【发表评论】

>> 最新消息

- > 《花样年华》入选百大电影最佳原声大碟(附图) (08/02/12)
- > 王家卫电影经典女性形象之《花样年华》张曼玉 (08/01/07)
- > 王家卫新作《蓝莓之夜》创意来自《花样年华》 (07/11/22)
- > 上海女人潘迪华 王家卫为其曾推翻《花样年华》 (07/11/05)
- > 王家卫御用摄影点亮新天地 曾拍《花样年华》等 (07/10/15)
- > 《色, 戒》海报曝光 被指抄袭《花样年华》(图) (07/07/12)
- > 《色, 戒》海报终曝光 疑抄袭《花样年华》(图) (07/07/08)
- > 王家卫新片《蓝莓之夜》不如《花样年华》(图) (07/05/18)

网友投票

★★★★☆

评分

评分

**Table 7: Statistics of the INTENT-1 topics and intents.**

		Subtopic Mining	Document Ranking
Chinese	topics	100	100
	intents	917	917
	subtopics	20,354	–
	unique rel docs	–	23,571
Japanese	topics	100	100
	intents	1,091	1,091
	subtopics	4,103	–
	unique rel docs	–	19,841

**Table 6: Statistics of the INTENT-2 topics and intents. Those for the revised Subtopic Mining data are shown in parentheses.**

		Subtopic Mining	Document Ranking
English	topics	50	–
	intents	392	–
	subtopic strings	4,157 (5,410)	–
Chinese	topics	98	97
	nav topics	23	22
	amb/faceted topics	23/52	23/52
	shared topics	21	21
	reused topics	19	19
	intents	616	615
	nav intents	–	125
	inf intents	–	490
	subtopic strings	6,251 (6,253)	–
	unique rel docs	–	9,295
	Japanese	topics	100
nav topics		33	28
amb/faceted topics		27/40	27/40
shared topics		21	21
reused topics		33	33
intents		587	582
nav intents		–	259
inf intents		–	323
subtopic strings		2,979 (2,989)	–
unique rel docs		–	5,085

**Table 2. NTCIR-9 INTENT Chinese and Japanese Document Collections.**

	Chinese	Japanese
Collection name	SogouT	ClueWeb-JA
#Pages	138 million	67.3 million
Size and Storage Medium	7z files on a 500GB hard disk	tarred/gzipped files on a 500GB hard disk

# INTENT — Evaluation Metric

---

- $D\#-nDCG = \gamma I-rec + (1 - \gamma) D-nDCG$ 
  - Official measure of NTCIR-9 INTENT
  - Combines *intent recall* with D-nDCG
    - Intent recall: fraction of subtopics covered by the ranked list
    - D-nDCG: nDCG with a global gain,  $\sum_i P(i|q)g_i(d)$ , where
      - $P(i|q)$ : probability of intent  $i$  for a given query  $q$
      - $g_i(d)$ : gain (relevance score) of a document  $d$  in terms of intent  $i$
- **D#-nDCG family went on** to be reused by IMine, Temporalia — i.e. the metric outlived the task

The technical legacy of INTENT is not a winning system — it is a **family of diversification metrics** still used a decade later.

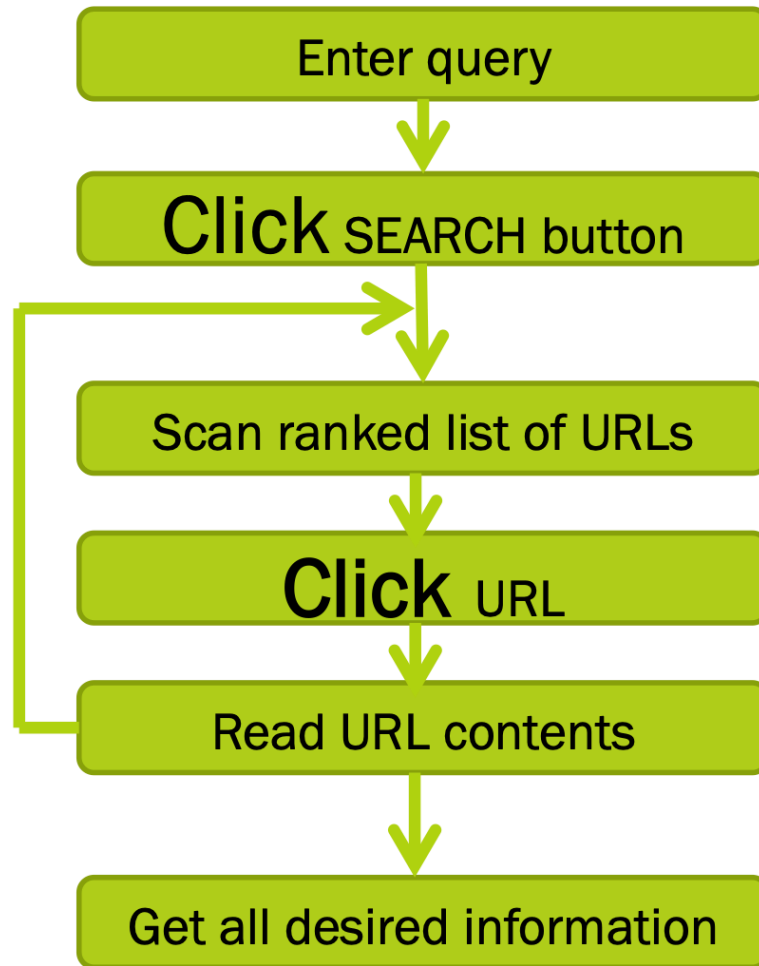
**1CLICK**

## Suppose that ...

- Finding answers for a question  
“what’s the difference between PDP and LCD?”

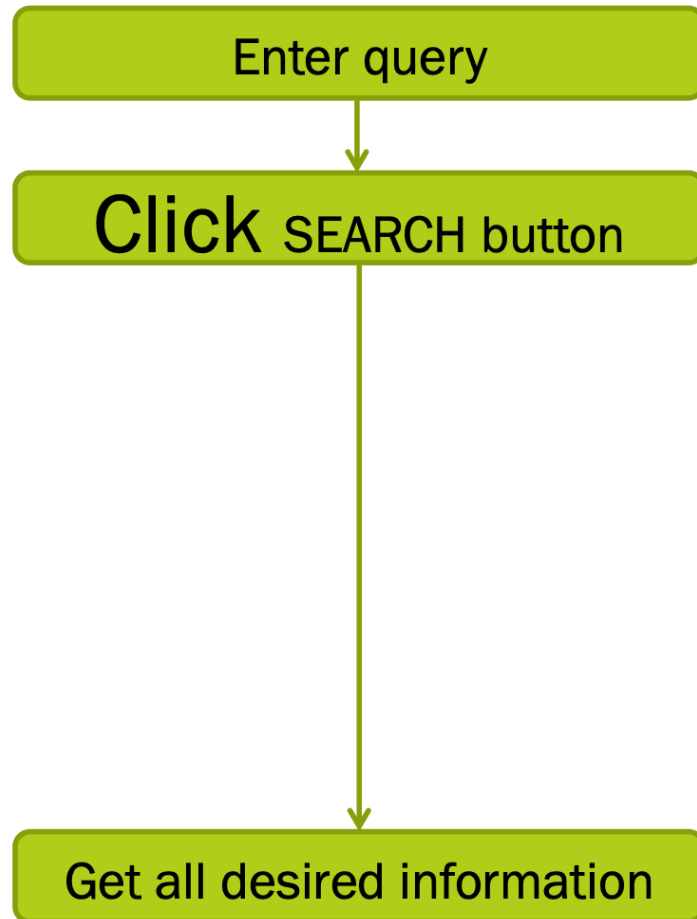


# In the "ten-blue-link" paradigm

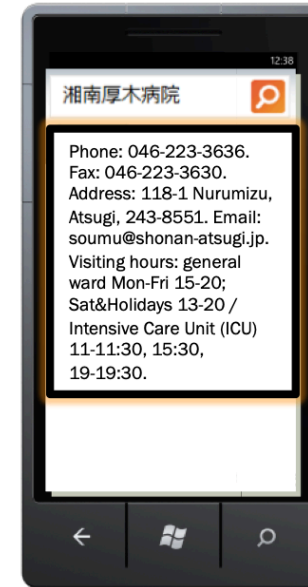


More than one clicks needed before being satisfied

# One Click Access



The system outputs *X-string*



## Task:

Given a search query, return  
a single textual output (*X-string*)

Go beyond the "ten-blue-link" paradigm, and tackle *information* retrieval rather than document retrieval

# Evaluation of 1CLICK Systems

- Manual/automatic matching between the **X-string** and **nuggets**

Phone: 046-223-3636. Fax: 046-223-3630.  
Address: 118-1 Nurumizu, Atsugi, 243-8551.  
Email: soumu@shonan-atsugi.jp

## **X-string**

- Phone number: 046-223-3636
- Fax number: 046-223-3630
- Address: 118-1 Nurumizu, Atsugi

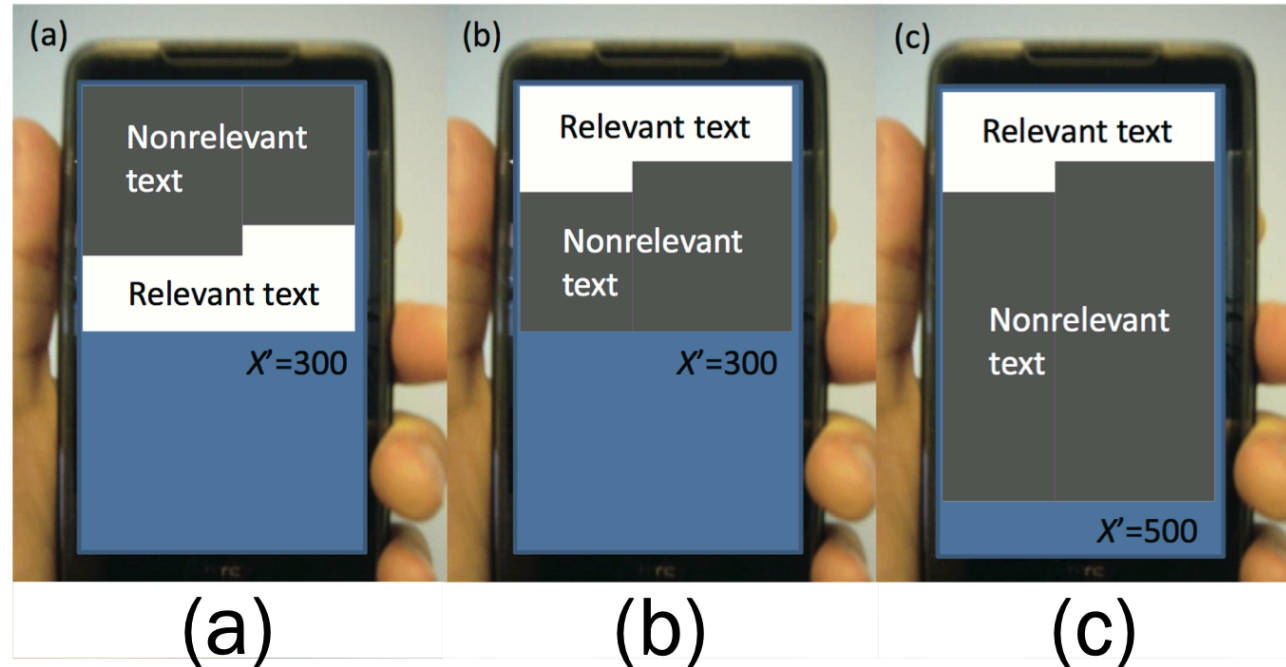
## **Nuggets**

a sentence relevant to the information need for a query

- Systems are required to present

more important information earlier

# Evaluation Metrics for 1CLICK



- Unlike nugget precision/recall, **S-measure** (position-aware weighted recall) says (a)<(b). **T-measure** (a kind of precision) says (b)>(c). **S#** (official evaluation metric) combines S and T

I believe that this task, which started in 2011, is precursor to the direct answer function of current search engines.

Instead of a ranked list of URLs, return a **single text** that answers the query directly.

*Given a query and an output window of size  $X$  characters, return an  $X$ -string that:*

- 1. presents important nuggets first, and*
- 2. minimises the amount of text the user has to read.*

- **NTCIR-9**: 60 Japanese queries, 4 query types (CELEBRITY / LOCATION / DEFINITION / QA); 2,839 nuggets
- **NTCIR-10**: 8 query types, English + Japanese; 3927 nuggets for 100 Japanese queries

**S-measure:** discount the value of a nugget by *how far the user has to read* to find it. (or a position-aware version of BLEU or ROUGE)

$$S = \frac{1}{\mathcal{N}} \sum_{u \in M} w(u) \cdot \max\left(0, 1 - \frac{\text{offset}(u)}{L}\right)$$

- $M$ : a set of nuggets for a query (useful information pieces)
- $w(u)$ : weight of a nugget  $u$
- $\text{offset}(u)$ : position of  $u$
- $L$ : a **patience parameter** (e.g. 500 chars  $\approx$  one minute of reading at 500 chars/min for Japanese)
- **Findings:** simple Wikipedia snippet baselines were surprisingly hard to beat for celebrity queries; participants were better for facility, definition, and QA queries

**OpenLiveQ**

Improve

Performance evaluated by  
REAL users

the **REAL performance** of  
question retrieval systems in a  
**production environment**

Yahoo! Chiebukuro  
(a CQA service of Yahoo! Japan)

- Given a query, return a ranked list of questions
  - Must satisfy many REAL users in Yahoo! Chiebukuro (a CQA service)



Effective for Fever **INPUT** X Q&A 検索

## Three things you should not do in fever

While you can easily handle most fevers at home, you should call 911 immediately if you also have severe dehydration with blue .... Do not blow your nose too hard, as the pressure can give you an earache on top of the cold. ....

10 Answers Posted on Jun 10, 2016

## Effective methods for fever

Apply the mixture under the sole of each foot, wrap each foot with plastic, and keep on for the night. Olive oil and garlic are both wonderful home remedies for fever. 10) For a high fever, soak 25 raisins in half a cup of water.

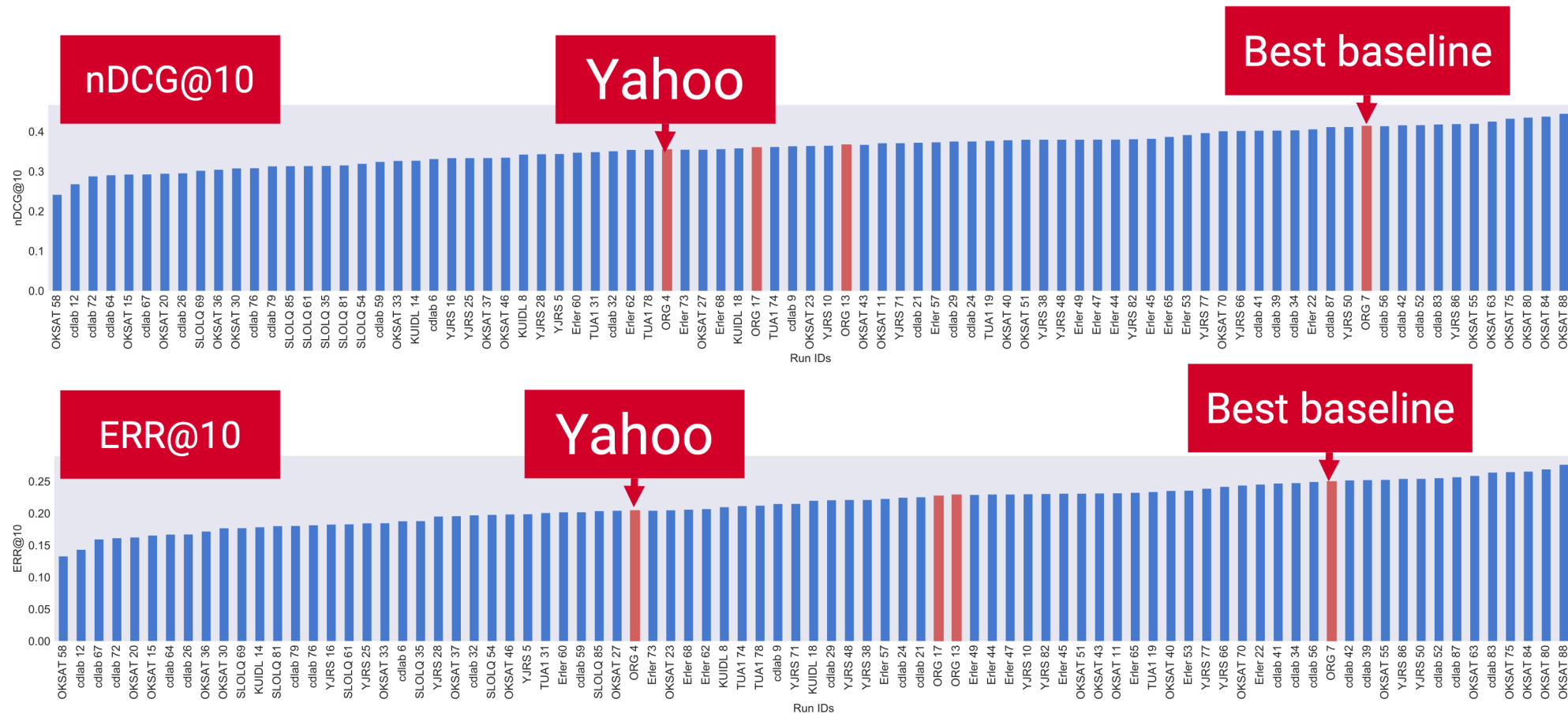
2 Answers Posted on Jan 3, 2010



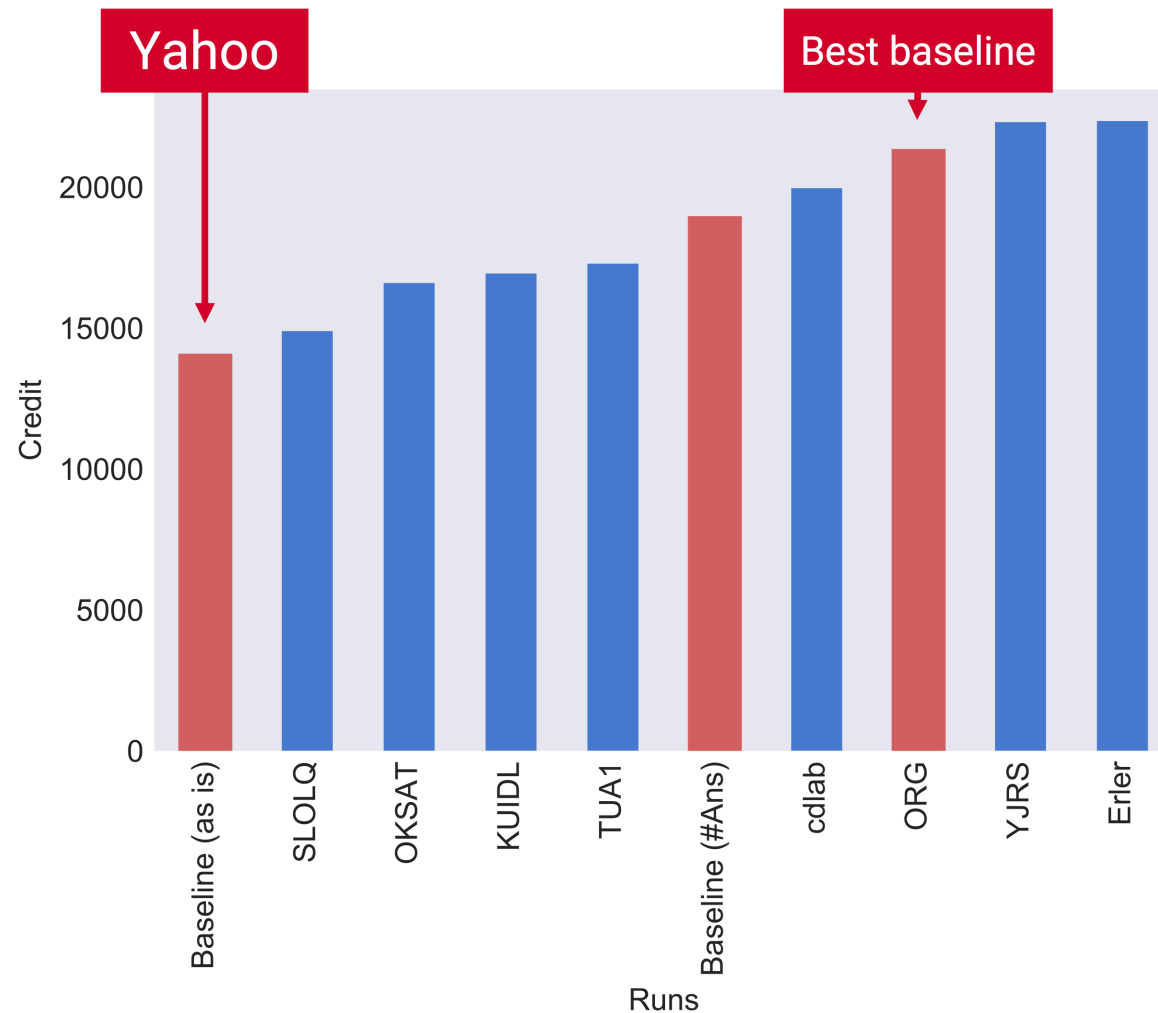


Ranked lists of questions from participants' systems are **INTERLEAVED**, presented to real users, and evaluated by their clicks

- Participants were given training data (queries + questions + clicks) and test queries (no clicks)
  - 1,000 queries sampled from Yahoo! Chiebukuro (Japanese Yahoo! Answers)
  - ~1M questions
  - 3-month click data
- Participants were asked to submit runs for the test queries, which were then evaluated both offline (using manual relevance labels) and online (against real users)
- **Online evaluation**
  - Multileaved comparison (Optimized Multileaving<sup>1</sup>) was used
    - We found it the best in our experiments<sup>2</sup>
- We intended to see a strong correlation between offline and online evaluation



Offline evaluation showed that the top 3 teams (OAKSAT, YJRS, cdlab) were significantly better than the Yahoo! production baseline.



Online evaluation showed a slightly different results: YJRS and Erler performed well but the others did not.

- **The offline and online evaluation results were not strongly correlated**
  - This suggests that offline evaluation may not always predict online performance: IR evaluation needs both offline and online methods to get a complete picture of system performance.
- **Advantages of offline evaluation:**
  - Controlled environment, reproducibility, and detailed analysis
  - Allows for rapid iteration and comparison of different approaches without the variability of real user interactions
- **Advantages of online evaluation:**
  - Real user interactions, capturing user satisfaction and engagement, and evaluating system performance in a live setting
  - Many services want to maximize user satisfaction, which is best measured through online evaluation

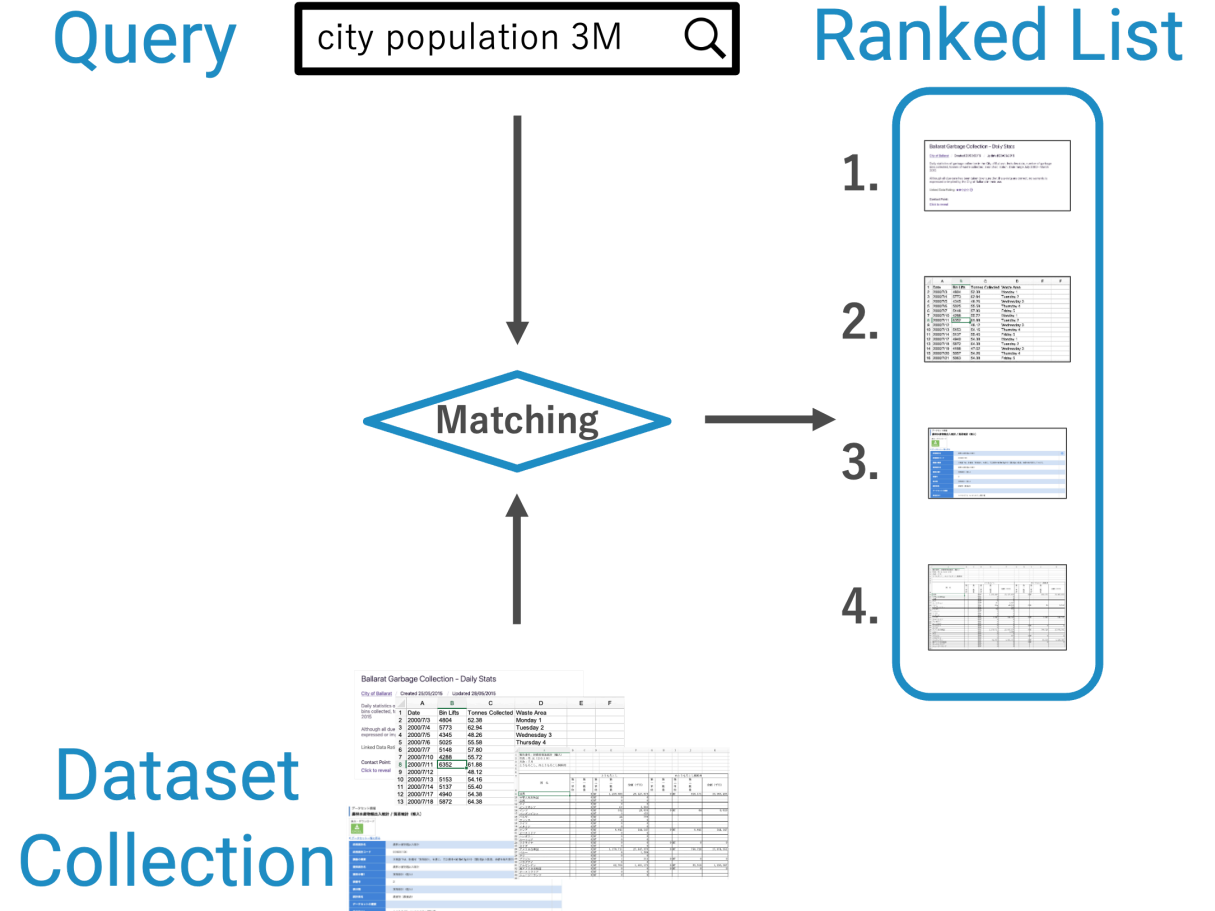
# Data Search

# Data Search — the first IR campaign for datasets

**Setting:** ad-hoc retrieval over **statistical / governmental datasets**, not Web documents.

Information need	Query
I am looking for evidences of domestic self-sufficiency rate of salt	domestic self salt rate
Are there many people who can't drive large trailers?	people can't drive large trailers
How many people have a second house?	many people second house

- **English dataset collection**
  - Data.gov <https://www.data.gov/>
- **Japanese dataset collection**
  - e-Stat <https://www.e-stat.go.jp/>



<b>English</b>	Documents (or <i>datasets</i> )	46,615
	Training queries	192
	Test queries	58
	Relevance judgments for training queries	8,248
	Relevance judgments for test queries	6,550
<b>Japanese</b>	Documents (or <i>datasets</i> )	1,338,402
	Training queries	192
	Test queries	72
	Relevance judgments for training queries	7,754
	Relevance judgments for test queries	4,035

Data.gov users! We welcome your [suggestions](#) for improving Data.gov and federal open data.

# The home of the U.S. Government's open data

Here you will find data, tools, and resources to conduct research, develop web and mobile applications, design data visualizations, and [more](#).

For information regarding the Coronavirus/COVID-19, please visit [Coronavirus.gov](#).

## GET STARTED

SEARCH OVER 335,221 DATASETS

Federal Student Loan Program Data

## HIGHLIGHTS

### Rivers of Data – Inland Electronic Navigation Charts



Nautical charts provide critical information to mariners in support of safe navigation. Historically these charts have been printed and distributed on paper, but modern communications systems allow for electronic charts that are able to be updated as new information becomes available. The National Oceanic and Atmospheric Administration (NOAA) Office of Coast Survey produces charts for coastal and Great Lakes areas, and

## 255,226 datasets found

**Lottery Powerball Winning Numbers: Beginning 2010** [1892 recent views](#)

State of New York — Go to <http://on.ny.gov/1GpWiHD> on the New York Lottery website for past Powerball results and payouts.

CSV RDF JSON XML

State

**FDIC Failed Bank List** [1449 recent views](#)

Federal Deposit Insurance Corporation — The FDIC is often appointed as receiver for failed banks. This list includes banks which have failed since October 1, 2000.

CSV HTML

Federal

**Electric Vehicle Population Data** [1253 recent views](#)

State of Washington — This dataset shows the Battery Electric Vehicles (BEVs) and Plug-in Hybrid Electric Vehicles (PHEVs) that are currently registered through Washington State Department...

CSV RDF JSON XML

State

**National Student Loan Data System** [1135 recent views](#)

Department of Education — The National Student Loan Data System (NSLDS) is the national database of information about loans and grants awarded to students under Title IV of the Higher...

XLS XLS XLS XLS XLS XLS 11 more in dataset

Federal

**Alzheimer's Disease and Healthy Aging Data** [967 recent views](#)

U.S. Department of Health & Human Services — 2015-2020. This data set contains data from BRFSS.

CSV RDF JSON XML

Federal

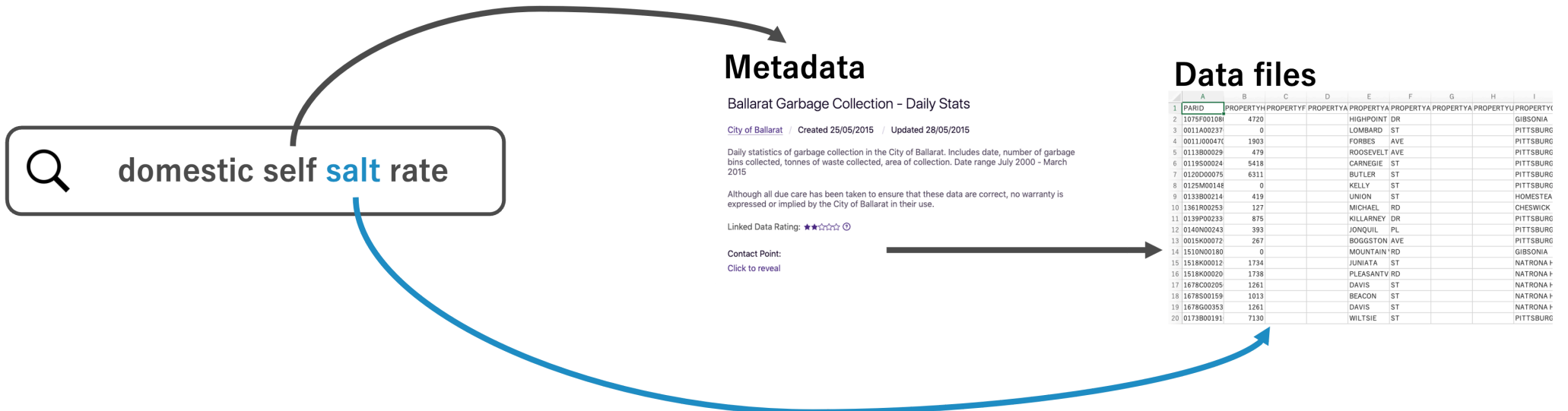


# Why is it difficult?

Required to match text with data (tabular data)

## What we are given

## What we want to retrieve



As the metadata often lack details,  
not all keywords in the query may not match the metadata

- **Followed by other data search test collections**
  - Lin et al. "ACORDAR: a test collection for ad hoc content-based (RDF) dataset retrieval." SIGIR 2022.
  - Chen et al. "ACORDAR 2.0: A test collection for ad hoc dataset retrieval with densely pooled datasets and question-style queries." SIGIR 2024.
  - Kolyada et al. "A Test Collection for Dataset Retrieval", ECIR 2025.
- **Included in a larger test collection**
  - Sun et al. "MAIR: A massive benchmark for evaluating instructed retrieval." EMNLP 2024.
- **Re-used for developing a new test collection**
  - Shi et al. "DSEBench: A Test Collection for Explainable Dataset Search with Examples"

## **4. Lessons Learned**

# Cross-task lessons from these tasks

---

- 1. Good metrics outlive any single submission.** S-measure (1CLICK) → U-measure → M-measure (MobileClick); D#-measure → still in use.
- 2. Reusable test collections are the durable contribution.** The deep value of running NTCIR-X is the data + baselines + judgments left behind for the next decade.
- 3. Strong baselines force ambition.** 1CLICK's Wikipedia baseline and OpenLiveQ's Yahoo production baseline both reframed what "winning" meant.
- 4. Online + offline triangulation reveals what either alone cannot.** OpenLiveQ's offline-vs-online divergence was a result *and* a methodological contribution.

Measure	First used	Idea
nDCG / Q-measure	NTCIR-3 Web (2003), CLIR-6 (2007)	Graded-relevance ranked retrieval
D# -nDCG, I-rec, DIN-nDCG, P+Q	INTENT (2011–13)	Diversification with intent recall
S-measure, T-measure, S#	1CLICK-1/2 (2011–13)	Position-aware text evaluation
U-measure	Sakai & Dou, SIGIR 2013	Generalised S to ranked lists / sessions
M-measure, H-measure	MobileClick / IMine	Multi-layer / hierarchical / vertical eval

# 5. NTCIR vs TREC vs CLEF

— A General Chair Perspective —

- **NTCIR** — since **1997**
- Hosted by **NII, Tokyo**
- Anchor: **Asian languages** (J, C, K)
- **Sesquiannual** (18-months cycle)
- **Tasks:** heavy on novel tasks, multilingual, domain-specific NLP tasks

- **TREC** — since **1992**
- Hosted by **NIST, USA**
- Anchor: **English**
- **Annual**
- **Tasks:** heavy on IR. Strongly connected to SIGIR topics

- **CLEF** — since **2000**
- Hosted by European universities
- Anchor: **European languages**
- **Annual**
- **Tasks:** initially focused on cross-lingual IR, later extended to multimedia and multimodal

They are not competitors. They are **geographically and linguistically complementary**.

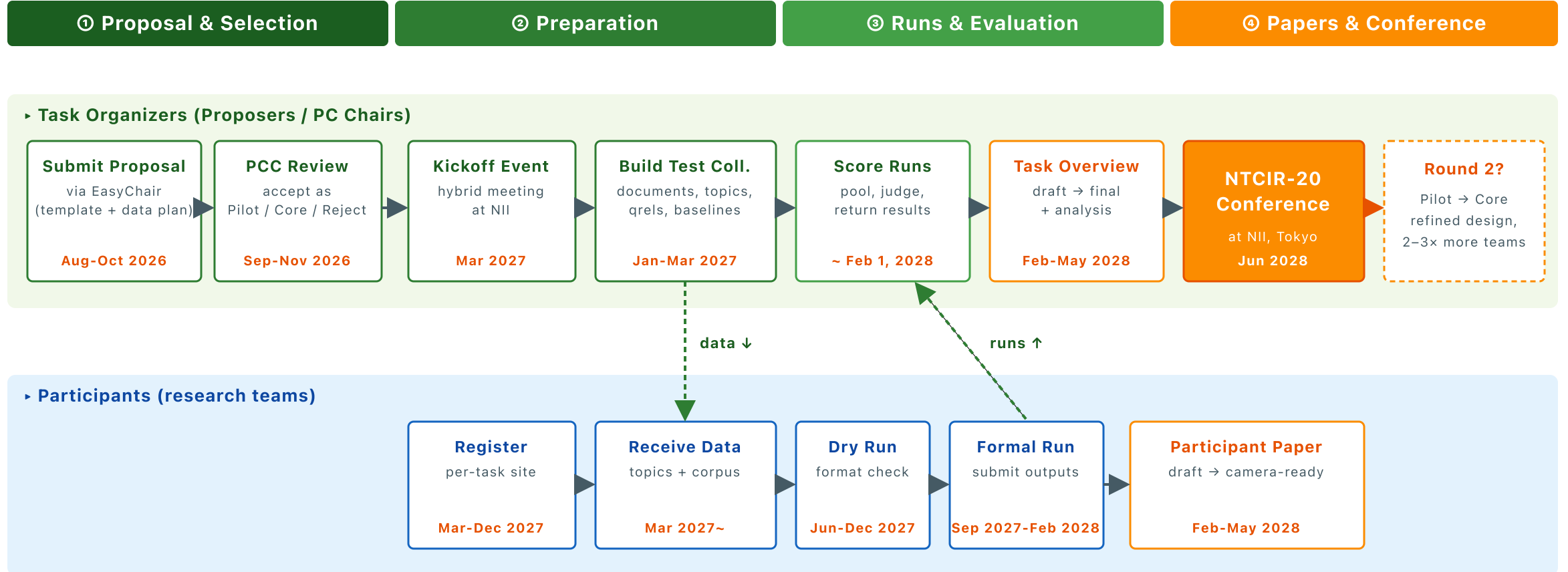
- **Language / region split** — NTCIR: East Asian (NII, Tokyo); TREC: English Web (NIST, USA); CLEF: European multilingual (CLEF Initiative + rotating European host).
- **NTCIR — pilot-core paradigm** — Patent (NTCIR-3), Math (NTCIR-10), Lifelog (NTCIR-12), 1CLICK / MobileClick (NTCIR-9–12), Data Search (NTCIR-15).
- **TREC — large-scale English benchmarks** — Deep Learning, Web, Total Recall, MS MARCO, RAG via NIST infrastructure.
- **CLEF — social & applied IR** — CheckThat!, EXIST, Touché, JOKER, LongEval, LifeCLEF, eHealth — federated lab proposals might enable societally oriented tasks.

- **Different language collections** — same task, different test beds (CLEF + NTCIR + TREC have run parallel CLIR for years)
- **Reproducibility tracks bridge the two** — **NTCIR-14 CENTRE** (CLEF-NTCIR-TREC REproducibility) explicitly required participants to reproduce NTCIR-13 WWW and TREC 2013 Web Track runs
- **Invited talks at each other's conferences** e.g. NTCIR conferences regularly featured invited talks by Ian Soboroff (TREC) and Nicola Ferro (CLEF).

# 6. Running an NTCIR Task

— A PC Chair Perspective —

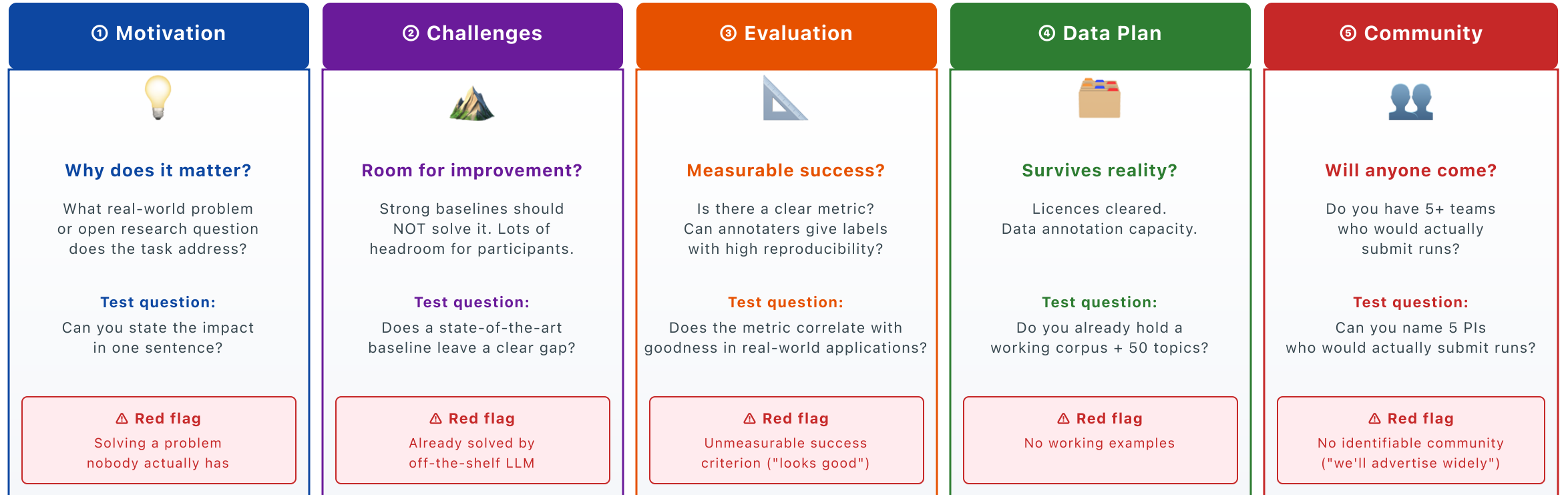
# How a task gets onto NTCIR (e.g., NTCIR-20 (TENTATIVE))



## Two rounds of proposals per cycle

- NTCIR-19 (July 2025 → Dec 2026, ~18-month cycle): 14 tasks accepted = 11 initial + 3 additional.
- Proposals are reviewed for clear motivation, clear challenges for participants, workable data plan, defensible evaluation, and likely  $\geq 5$  teams.

## Five pillars of a strong NTCIR task proposal



### PCC perspectives

- ✓ All five pillars defensible → **accept as Core**
- ~ 3-4 pillars clear, 1-2 immature → **accept as Pilot** (smaller, scoped, iterate to round 2)
- × Weak motivation · solved by baseline · unmeasurable success · unclear licence · no community → **reject**

# Advice for prospective organisers

---

- 1. Build a dataset before you propose.** Even a small prototype makes the proposal concrete and surfaces problems early.
- 2. Design with participants in mind.** Overly complex tasks scare teams away — and a task no one joins is a bad outcome for organisers too. Lower the bar with baseline code, starter kits, and a clear tutorial.
- 3. Plan for multiple rounds from day one.** NTCIR strongly encourages running a task for **at least two rounds**. Scope round 1 knowing that round 2 is where the task matures.
- 4. It does not have to be an IR task.** NTCIR has hosted plenty of NLP tasks — sentence classification, NER, argument mining, and more. Applied NLP almost always reduces to a form of information access.
- 5. Involve many people, across many places.** Strong tasks are cross-country and multilingual — in both the organising team and the participants.

**Take-home**

## Take-home messages

---

- 1. Tasks and test collections often outlast systems.** Winning runs fade, but the data, metrics, and baselines left behind tend to be reused for years.
- 2. Be open to trying new tasks.** Patent, Math, Lifelog, 1CLICK, Dataset Search — NTCIR has often hosted tasks that later grew into wider research areas.
- 3. Metrics and user models deserve careful design,** just as much as the systems they evaluate. Good ones can be reused across tasks.
- 4. Multilingual and cross-country collaboration helps,** both in the organising team and among participants.
- 5. The community matters.** Tasks come and go, but the people who gather around NTCIR are what keep the field moving forward.