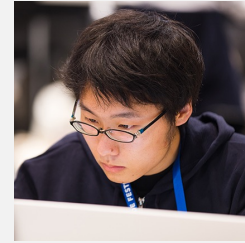


Overview of the NTCIR-16 Data Search 2 Task

Makoto P. Kato (University of Tsukuba), Hiroaki Ohshima (University of Hyogo),
Ying-Hsang Liu (Oslo Metropolitan University), Hsin-Liang Chen (Philadelphia College of
Osteopathic Medicine), Yu Nakano (University of Tsukuba)



- **The open data movement is now being accelerated by the expectations for open science and citizen science**
 - Each country strongly encourages the open data movement:
 - Data.gov (United States)
 - Data.gov.uk (United Kingdom)
 - Data.gov.au (Australia)
 - e-Stat (Japan)
- **Besides the governmental portals, there are also thousands of data repositories on the Web**

Demand for a better data search engine

(e.g. Google Dataset Search)

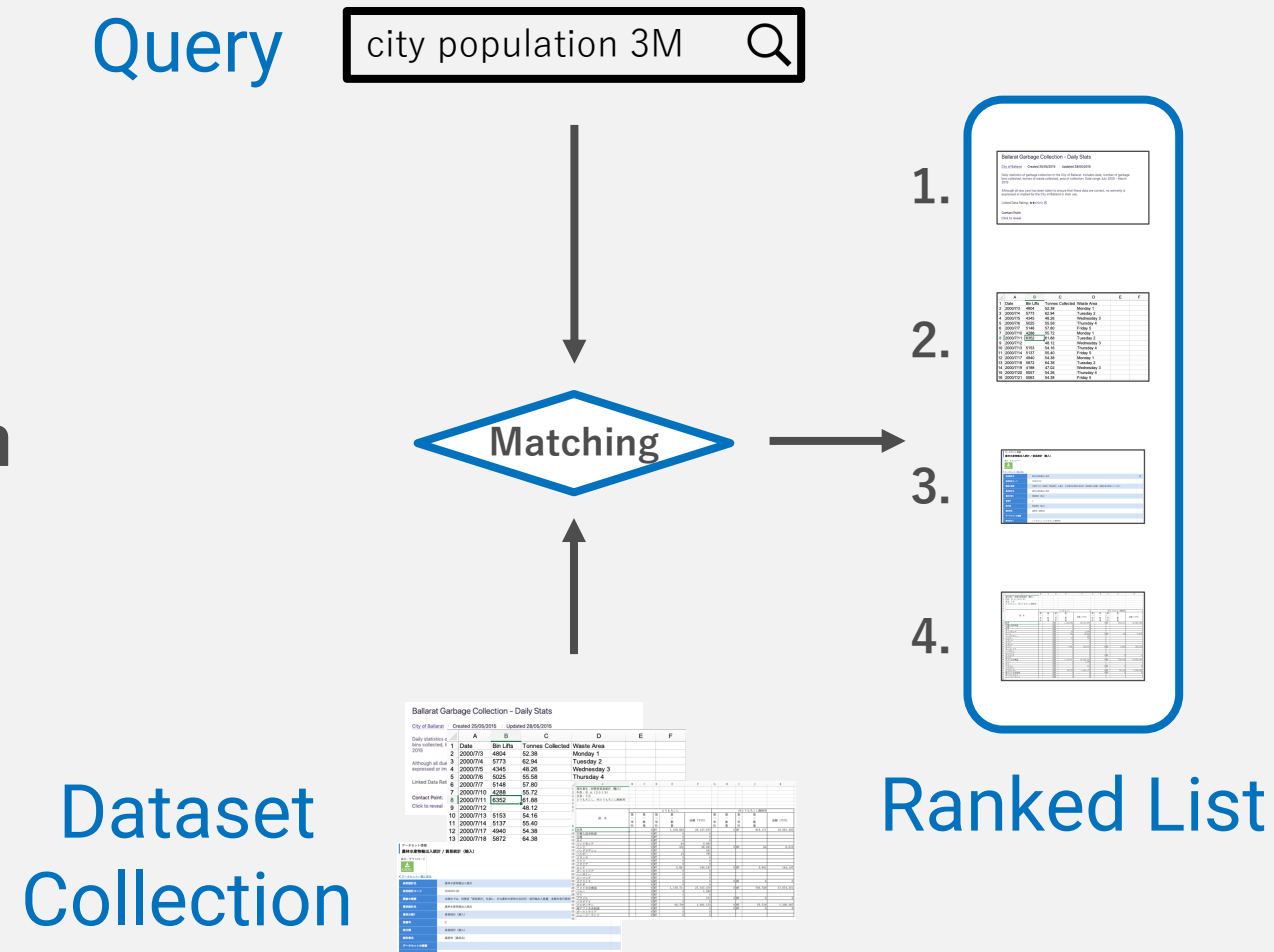
The very first IR evaluation campaign for data search

English	Documents (or <i>datasets</i>)	46,615
	Training queries	96
	Test queries	96
	Relevance judgments for training queries	2,008
	Relevance judgments for test queries	6,240

Japanese	Documents (or <i>datasets</i>)	1,338,402
	Training queries	96
	Test queries	96
	Relevance judgments for training queries	2,035
	Relevance judgments for test queries	5,719

Ad-hoc retrieval for statistical data

- **Subtasks**
 - English and Japanese
- **Input**
 - 96 queries for each of the subtasks
- **Document (or Dataset) collection**
 - Data.gov for English
 - e-Stats for Japanese
- **Output**
 - Ranked list of datasets for each query



Q1. What techniques were potentially effective?

- Not very conclusive, but possibly neural models and table understanding

Q2. What queries were difficult in dataset retrieval?

- Time-related queries are especially difficult

Q3. Is the topic variability large?

- Yes, it is much larger than the system variability

The second round of the "Data Search" task with more queries

English	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
Japanese	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

- **IR Subtask**

- Given a query and a dataset collection, a system is expected to generate a ranked list of datasets.

- ~~**QA Subtask**~~

- ~~- Given a question and a dataset, a system is expected to generate an answer to the question, mainly by extracting a part of the dataset.~~

- ~~**UI Subtask**~~

- ~~- Participants are expected to develop a search system with an effective search interface for dataset search tasks.~~

- **NTCIR-15 Data Search**

- Information needs were derived from questions in cQA

- **NTCIR-16 Data Search 2**

- Information needs were derived from webpages referring to a dataset

- Parsed Common Crawl webpages (~25B)
 - Identified 47,242 URLs including “data.gov” and 137,388 URLs including “stat.go.jp”.

- Manually extracted a potential need from the webpages

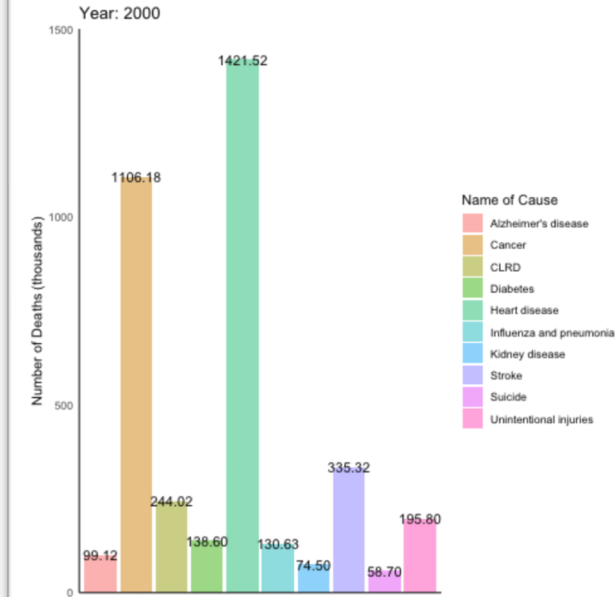
Largest Causes of Death in the United States: 1999 – 2016

December 16, 2018 10:22 am by IWB

Sharing is Caring!

What is the largest causes of death in United State in 1999-2016?

Generated need



Data was provided by the Center for Disease Control and Prevention, found [here](#). I used the [graphviz](#) library for the visualization. The code is available on [GitHub](#).

- Used crowd-sourcing services to convert information needs to queries
 - Showed a need and asked workers to input a query without looking at the need
 - Tried to simulate a more realistic situation
 - Selected the most "probable" query from 10 workers' queries
 - Built a unigram language model from those queries, and selected the one with the highest perplexity with respect to the language model

プレビュー

すべての作業を確認することはできません。プレビューは一部の作業のみ確認することができます。 [次の作業 ▶](#)

検索用キーワードの作成作業

ある情報を知りたい人、もしくは、ある疑問を解決したい人の要望が以下に提示されます。その人のために、Googleやヤフーなどで検索するための検索キーワードを考えて入力してください。

1. まず「質問・要望」をよく読んで理解してください。
2. 画面を下の方にスクロールすると「検索キーワード入力欄」がありますので、Googleやヤフーなどで検索するための検索キーワードを自分で考えて入力してください。**決して「質問・要望」を見ながら検索キーワードを入力しないでください。コピー&ペーストも厳禁です。**

例

- 表示される「質問・要望」の例：「千葉県では日本一落花生がとれますか？」
- 適切な「検索キーワード」の例：「千葉県 落花生 生産量」
- 表示される「質問・要望」の例：「最近、東京都内ではどれくらいマンションが増えているのでしょうか？」
- 適切な「検索キーワード」の例：「2019年 都内 マンション 建設数」

注意点

- **決して「質問・要望」を見ながら検索キーワードを入力しないでください。コピー&ペーストも厳禁です。**
- 検索単語を入力する場合には空白（スペース）で区切ってください。
- 1回分の検索用のキーワードを入力してください。2回に渡って検索することを前提としないでください。

質問・要望

都道府県別の有効求人倍率をおしえてください。

検索キーワード

[次の作業 ▶](#)

← Need

← Query

ID	Need	Query
DS2-E-0001	What is the largest causes of death in United State in 1999-2016?	causes of death us 1999-2016
DS2-E-0004	How much is the tuition fee at a private elementary school?	tuition fee private elementary school
DS2-E-0005	Are there hospital differences across the US states?	us states hospital differences

- **See the tutorial given by Prof. Sakai[†]**

- Computed the residual variance of $n\text{DCG}@10$ based on the NTCIR-15 Data Search results

- The minimum topic size* is

- $n = 36$ for English runs → **58**

- $n = 68$ for Japanese runs → **72**

Topic set size design based on the statistical power of the (paired) t-test (1)

2019

Data from Sakai's SIGIR 2019 Paper (Which Diversity Evaluation Measures Are "Good?")
Data from Sakai's CLEF20 book chapter: How to Run an Evaluation Task

2018

Data from Sakai's SIGIR 2018 Short Paper (Comparing Two Binned Probability Distributions for Information Access Evaluation)
Tutorial kit for Sakai SIGIR 2018: Conducting Laboratory Experiments Properly with Statistical Tools: An Easy Hands-on Tutorial (includes the five excel files listed below besides other materials)

[samplesizeANOVA2.xlsx](#) for computing topic set sizes to achieve high statistical power for one-way ANOVA (recommended: See Sakai's book)
[samplesize2SAMPLET1.xlsx](#) for computing topic set sizes to achieve a tight confidence interval for paired data (recommended: See Sakai's book)
[samplesize2SAMPLET2.xlsx](#) for computing topic set sizes to achieve high statistical power for the two-sample t-test
[samplesizeTTEST2.xlsx](#) for computing topic set sizes to achieve high statistical power for the paired t-test
[samplesizeCI2.xlsx](#) for computing topic set sizes to achieve a tight confidence interval for unpaired data

Use this tool even if you're interested in the t-test

queries used in NTCIR-16

[†] Sakai. Laboratory Experiments in Information Retrieval: Sample Sizes, Effect Sizes, and Statistical Power. 2018.

* $\text{minD} = 0.05$, $m = 10$, $\alpha = 0.05$, and $\beta = 0.20$

• English

–Data.gov

- <https://www.data.gov/>
- 46,615 (~445GB)

The screenshot shows the Data.gov interface for a dataset from the City of New York. The dataset is titled 'Demographic and Housing Profiles by Borough' and was last updated on March 29, 2019. It contains selected demographic and housing estimates for the city and its boroughs, based on five-year estimates from the Census Bureau's American Community Survey. The page includes an 'Access & Use Information' section with public access status, a non-federal disclaimer, and license information. A 'Downloads & Resources' section offers an MS Excel file with 18 views and a download button.

• Japanese

–e-Stat

- <https://www.e-stat.go.jp/>
- 1,338,402 (~100GB)

The screenshot shows the e-Stat interface for a dataset titled '作物統計調査 / 作況調査(水陸稲、麦類、豆類、かんしょ、飼料作物、工芸農作物) 速報 令和元年産一番茶の摘採面積、生葉収穫量及び荒茶生産量 (主産県)'. It includes a 'データセット情報' (Dataset Information) section with a download button for an EXCEL file. Below this is a table with details about the dataset.

データセット情報	
作物統計調査 / 作況調査(水陸稲、麦類、豆類、かんしょ、飼料作物、工芸農作物) 速報 令和元年産一番茶の摘採面積、生葉収穫量及び荒茶生産量 (主産県)	
表示・ダウンロード	
データセット一覧に戻る	
政府統計名	作物統計調査 i
政府統計コード	00500215
調査の概要	本調査は、毎年、耕地の状況、収穫量等を調査し、耕地面積、農作物の作付面積、収穫量、被害面積・被害量等を、全国、都道府県（主産県）別等に提供しています。
提供統計名	作物統計調査
提供分類1	作況調査(水陸稲、麦類、豆類、かんしょ、飼料作物、工芸農作物)
提供分類2	速報
提供分類3	令和元年産一番茶の摘採面積、生葉収穫量及び荒茶生産量（主産

- The relevance of each dataset for a given query is judged by crowd-sourcing workers

- 0: Not-relevant
- 1: Partially relevant
- 2: Highly relevant

- Inter-rater agreement

(measured by Krippendorff's α)

- English: **0.444**
- Japanese: **0.474**

(Fairly consistent with those of NTCIR-15)

Instructions

Please judge how useful a **DATASET** of a webpage is for answering a given **REQUEST**. Please carefully read a given **REQUEST**, visit a webpage describing a **DATASET**, and give a usefulness score (0, 1, or 2) to each of the datasets.

Rules

1. Carefully read a **REQUEST** (Note: this page contains a few types of requests.)
2. Make sure that you visit a webpage that describes a **DATASET**, and judge how useful the **DATASET** is for answering the **REQUEST**.
3. Usefulness score is defined as:
 - 0: (Useless) The **DATASET** is not useful to answer the **REQUEST** at all, or was not accessible for some reasons.
 - 1: (Partially useful) The **DATASET** is useful to partially answer the **REQUEST**, but cannot fully answer the **REQUEST**.
 - 2: (Highly useful) The **DATASET** is useful to fully answer the **REQUEST**.

Cautions

- You will be rejected if the website is not accessed.
- You will be rejected if the work time is too short.
- There are some **REQUEST** and **DATASET** for which a true usefulness score is known. You will be rejected if your answer is very different from the true answer.
- You will be rejected if your work result has been rejected before.

1.

REQUEST: Do people in the East Coast dislike oysters?

DATASET: [LINK](#)

0: Useless 1: Partially useful 2: Highly useful

• Applied standard retrieval models to only the metadata

Query

domestic self salt rate

Retrieve



Metadata

データセット情報
食料需給表 / 確報 平成16年度食料需給表

表示・ダウンロード
EXCEL

← データセット一覧に戻る

政府統計名	食料需給表
政府統計コード	00500300
調査の概要	本統計では、穀類等の品目別の国内生産量、輸出入量、国内消費仕向量、1人当たり供給純食料、食料自給率等を毎年提供しています。
提供統計名	食料需給表
提供分類1	確報
提供分類2	平成16年度食料需給表

Data file

(1) 自給率の推移
① 品目別自給率の推移

品目	昭和35年度	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	5
米	102	95	98	95	94	95	101	116	118	117	106	100	100	101	102	110	100	1
うち主食用																		
小麦	39	43	38	17	28	28	21	20	20	14	9	8	5	4	4	4	4	4
大麦	104	89	85	60	58	57	56	50	46	38	28	23	14	8	9	8	8	8
裸麦	112	89	90	28	119	123	90	92	115	98	73	73	64	87	111	98	87	
大・裸麦計	107	89	85	51	20	20	15	59	60	48	34	29	18	10	11	10	9	
雑穀	21										1	1	1	1	1	1	0	
いも類	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	99	98
かんしょ	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	1
ばれいしょ	101	101	101	100	100	100	100	100	100	100	100	100	100	100	100	100	99	98
でんぶん	76	76	71	85	72	67	56	55	56	44	41	36	36	30	23	24	29	
豆類	44	42	35	31	23	25	19	20	17	14	13	11	13	11	11	9	8	
大豆	28	25	21	17	13	11	9	8	7	5	4	4	4	3	4	4	3	
その他の豆類	90	90	76	74	58	70	57	72	61	56	65	49	67	60	59	45	42	
野菜	100	100	100	100	100	100	100	100	100	100	99	99	99	98	98	99	98	
果実	100	98	97	92	90	90	89	89	88	85	84	81	81	83	83	84	81	
みかん	111	113	111	109	109	109	107	109	105	106	105	106	103	104	104	102	103	1
りんご	102	102	102	101	102	102	102	102	103	102	102	101	101	101	100	100	100	1
米類	93	95	95	90	92	93	91	86	85	85	89	84	82	80	84	76	75	

Not Used

• Baseline retrieval models

- BM25, BM25 + RM3, BM25 + SDM, BM25 + BM25PRF
- Query Likelihood, Query Likelihood + RM3, Query Likelihood + SDM

- **NTCIR-16 Data Search attracted seven research groups and received 48 systems' results in total**
 - Including the organizers' team for providing baseline runs
 - 25 for English and 23 for Japanese
- **Participants**
 - UHGSIS: University of Hyogo
 - STIS: Politeknik Statistika STIS
 - WUT21: Wuhan University of Technology
 - KSU: Kyoto Sangyo University
 - NYUCIN: Universidade Federal de Pernambuco
 - OUHCIR: The University of Oklahoma

Run ID	Description	Data	Neural	Entity	Num.
KSU-E-1	Category+Table Clipping+Table Header+BERT+MLP	Y	Y	Y	N
KSU-E-3	Category+Table Header+BERT+MLP	Y	Y	Y	N
KSU-E-5	File Clipping+Table Header+BERT+MLP	Y	Y	Y	N
KSU-E-7	Table Header+BERT+MLP	Y	Y	Y	N
KSU-E-9	Category+Table Header+BM25	Y	N	Y	N
NYUCIN-E-1	BM25 and BERT	Y	Y	N	N
ORGE-E-1	bm25prf+bm25	N	N	N	N
ORGE-E-2	bm25	N	N	N	N
ORGE-E-3	bm25.accurate	N	N	N	N
ORGE-E-4	edm+qld	N	N	N	N
ORGE-E-5	rm3+bm25	N	N	N	N
ORGE-E-6	qld	N	N	N	N
ORGE-E-7	slu+rm3	N	N	N	N
ORGE-E-8	rm3+qld	N	N	N	N
OUHCIR-E-1	BM25 and TFIDF WEIGHT ADJUSTED and Sentence Transformer	N	Y	N	N
OUHCIR-E-2	BM25 and TFIDF WEIGHT ADJUSTED and Sentence Transformer	N	Y	N	N
OUHCIR-E-3	BM25 and TFIDF WEIGHT ADJUSTED	N	N	N	N
OUHCIR-E-4	BM25 and TFIDF WEIGHT ADJUSTED	N	N	N	N
OUHCIR-E-5	BM25 and TFIDF WEIGHT ADJUSTED	N	N	N	N
OUHCIR-E-6	BM25 and TFIDF WEIGHT ADJUSTED	N	N	N	N
OUHCIR-E-7	DOC2VEC	N	N	N	N
OUHCIR-E-8	DOC2VEC	N	N	N	N
STIS-E-1	prop+bert_score+bm25	Y	Y	N	N
STIS-E-2	prop+bert_score	Y	Y	N	N
wut21-E-1	LM Jelinek Mercer	Y	N	N	N

- **Data (14/48 runs):**

- Whether the data files are used

- **Neural (22/48 runs):**

- Whether neural language models (e.g., BERT) are used

- **Entity (10/48 runs):**

- Whether entities are treated differently from the other tokens

- **Num. (0/48 runs):**

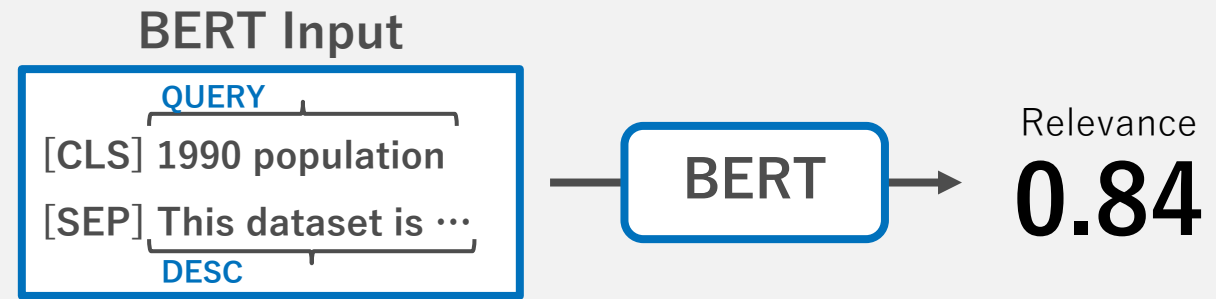
- Whether numbers are treated differently from the other tokens

Run ID	Description	Data	Neural	Entity	Num.
KSU-E-1	Category+Table Clipping+Table Header+BERT+MLP	Y	Y	Y	N
KSU-E-3	Category+Table Header+BERT+MLP	Y	Y	Y	N
KSU-E-5	Table Clipping+Table Header+BERT+MLP	Y	Y	Y	N
KSU-E-7	Table Header+BERT+MLP	Y	Y	Y	N
KSU-E-9	Category+Table Header+BM25	Y	N	Y	N
NYUCIN-E-1	BM25 and BERT	Y	Y	N	N
ORGE-E-1	bm25prf+bm25	N	N	N	N
ORGE-E-2	bm25	N	N	N	N
ORGE-E-3	bm25.accurate	N	N	N	N
ORGE-E-4	sdm+qld	N	N	N	N
ORGE-E-5	rm3+bm25	N	N	N	N
ORGE-E-6	qld	N	N	N	N
ORGE-E-7	sdm+bm25	N	N	N	N
ORGE-E-8	rm3+qld	N	N	N	N
OUHCIR-E-1	BM25 and TFIDF WEIGHT ADJUSTED and Sentence Transformer	N	Y	N	N
OUHCIR-E-2	BM25 and TFIDF WEIGHT ADJUSTED and Sentence Transformer	N	Y	N	N
OUHCIR-E-3	BM25 and TFIDF WEIGHT ADJUSTED	N	N	N	N
OUHCIR-E-4	BM25 and TFIDF WEIGHT ADJUSTED	N	N	N	N
OUHCIR-E-5	BM25 and TFIDF WEIGHT ADJUSTED	N	N	N	N
OUHCIR-E-6	BM25 and TFIDF WEIGHT ADJUSTED	N	N	N	N
OUHCIR-E-7	DOC2VEC	N	N	N	N
OUHCIR-E-8	DOC2VEC	N	N	N	N
STIS-E-1	prop+bert_score+bm25	Y	Y	N	N
STIS-E-2	prop+bert_score	Y	Y	N	N
wut21-E-1	LM Jelinek Mercer	Y	N	N	N

Run ID	Description	Data	Neural	Entity	Num.
KSU-J-10	Category+Table Header+BM25	Y	N	Y	N
KSU-J-2	Category+Table Header+BERT+MLP	Y	Y	Y	N
KSU-J-4	Category+Table Header+BERT+MLP	Y	Y	Y	N
KSU-J-6	Table Clipping+Table Header+BERT+MLP	Y	Y	Y	N
KSU-J-8	Table Header+BERT+MLP	Y	Y	Y	N
ORGJ-J-6	qld	N	N	N	N
ORGJ-J-1	bm25prf+bm25	N	N	N	N
ORGJ-J-8	rm3+qld	N	N	N	N
ORGJ-J-7	sdm+bm25	N	N	N	N
ORGJ-J-2	bm25	N	N	N	N
ORGJ-J-5	rm3+bm25	N	N	N	N
ORGJ-J-4	sdm+qld	N	N	N	N
ORGJ-J-3	bm25.accurate	N	N	N	N
UHGSIS-J-9	BM25, BERT, query modification, target 1000	N	Y	N	N
UHGSIS-J-7	BM25, BERT, query modification, target 2000	N	Y	N	N
UHGSIS-J-6	BM25, BERT, query original, target 3000	N	Y	N	N
UHGSIS-J-1	BM25, BERT, query modification, target all	N	Y	N	N
UHGSIS-J-8	BM25, BERT, query original, target 2000	N	Y	N	N
UHGSIS-J-4	BM25, query original, target all	N	N	N	N
UHGSIS-J-3	BM25, query modification, target all	N	Y	N	N
UHGSIS-J-2	BM25, BERT, query original, target all	N	Y	N	N
UHGSIS-J-5	BM25, BERT, query modification, target 3000	N	Y	N	N
UHGSIS-J-10	BM25, BERT, query original, target 1000	N	Y	N	N

- **Neural : BERT-based relevance estimation**

- Input a query and a description of a dataset into BERT for estimating the query-dataset relevance



- **Data: Table header extraction**

- Index terms in the table header together with the description of a dataset

Data file

(1) 自給率の推移

① 品目別自給率の推移

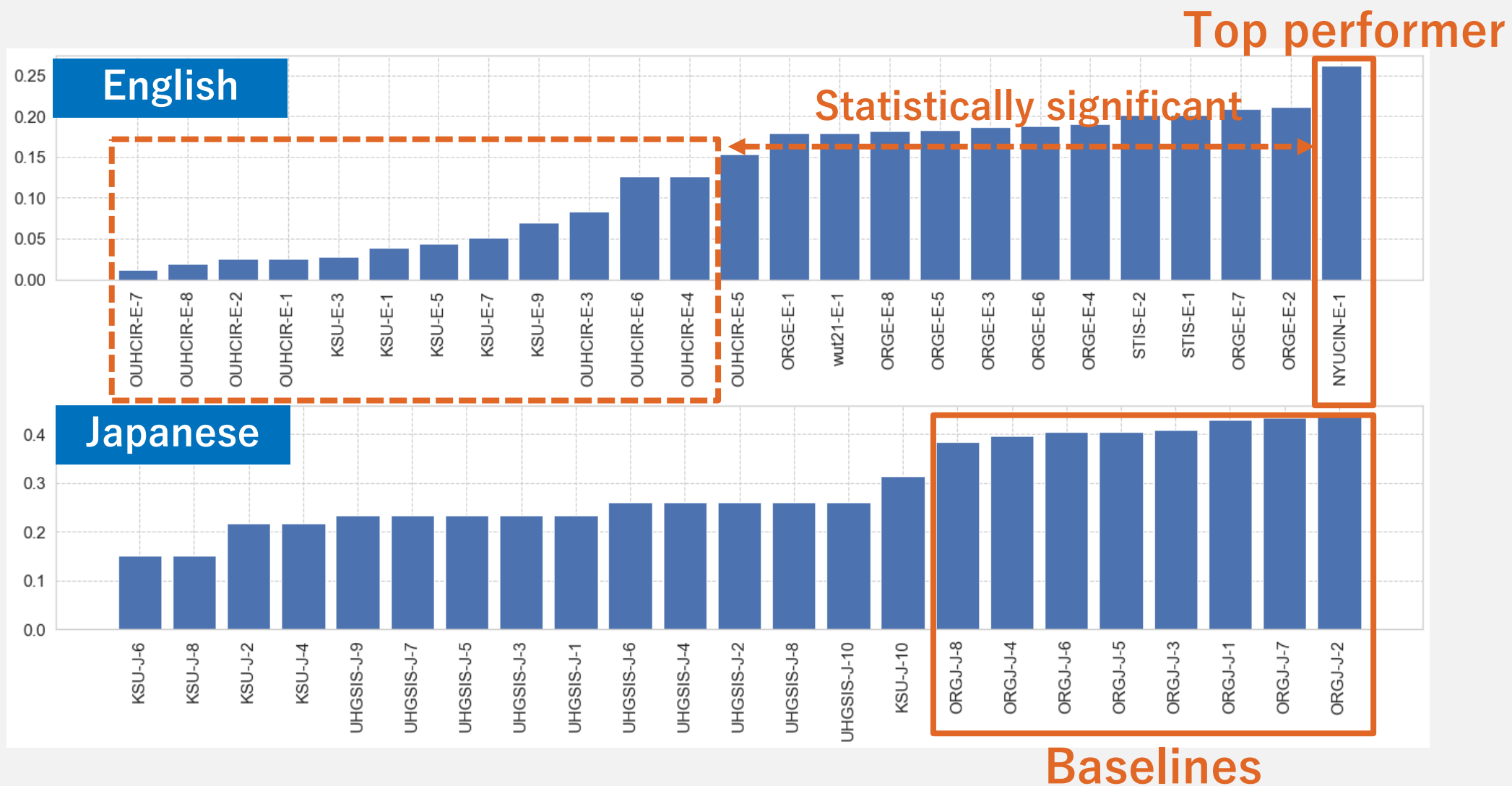
品 目	昭 和 35年度	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52
米	102	95	98	95	94	95	101	116	118	117	106	100	100	101	102	110	100	100
うち主食用																		
小麦	39	43	38	17	28	28	21	20	20	14	9	8	5	4	4	4	4	4
大麦	104	89	85	60	58	57	56	50	46	38	28	23	14	8	9	8	8	8
雑穀	112	89	90	28	119	123	90	92	115	98	73	73	64	87	111	98	87	
大・裸麦計	107	89	87	51	70	73	65	59	60	48	34	29	18	10	11	10	9	
雑穀	21	15	11	9	6	5	3	3	2	2	1	1	1	1	1	1	0	
いも類	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	99	98
かんしょ	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	1
ばれいしょ	101	101	101	100	100	100	100	100	100	100	100	100	100	100	100	99	98	
でんぷん	76	76	71	85	72	67	56	55	56	44	41	36	36	30	23	24	29	
豆類	44	42	35	31	23	25	19	20	17	14	13	11	13	11	11	9	8	
大豆	28	25	21	17	13	11	9	8	7	5	4	4	4	3	4	4	3	
その他の豆類	90	90	76	74	58	70	57	72	61	56	65	49	67	60	59	45	42	
野菜	100	100	100	100	100	100	100	100	100	100	99	99	99	98	98	99	98	
果実	100	98	97	92	90	90	89	89	88	85	84	81	81	83	83	84	81	
みかん	111	113	111	109	109	109	107	109	105	106	105	106	103	104	104	102	103	1
りんご	102	102	102	101	102	102	102	102	103	102	102	101	101	101	100	100	100	1
雑穀	93	95	95	90	91	93	91	86	85	85	89	84	89	80	84	76	75	

Index



nDCG@10

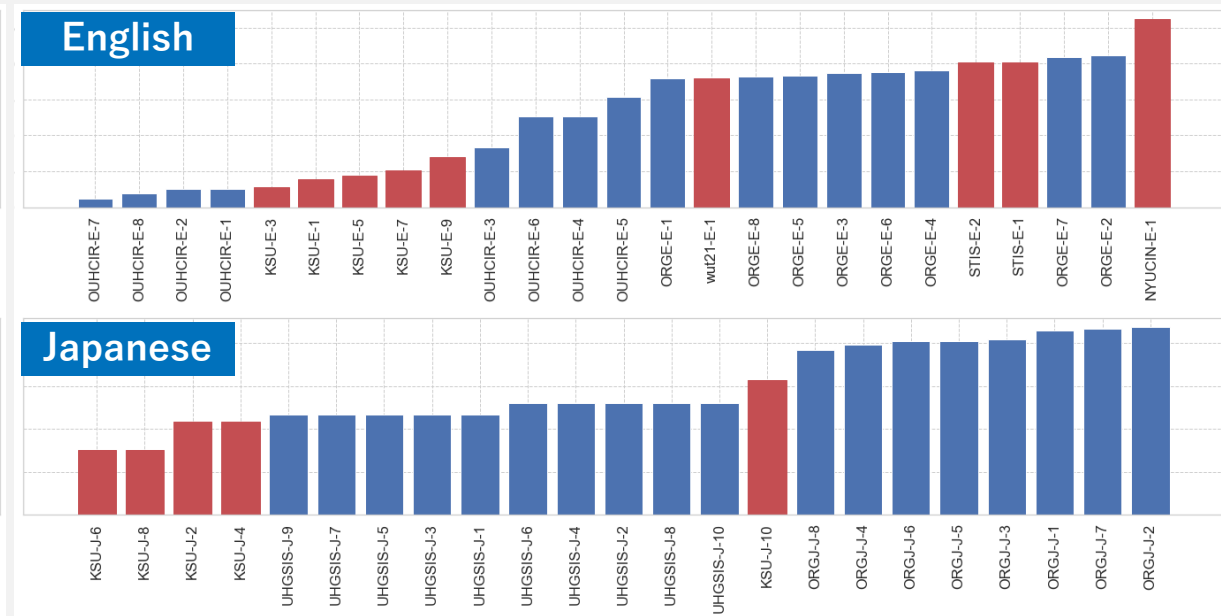
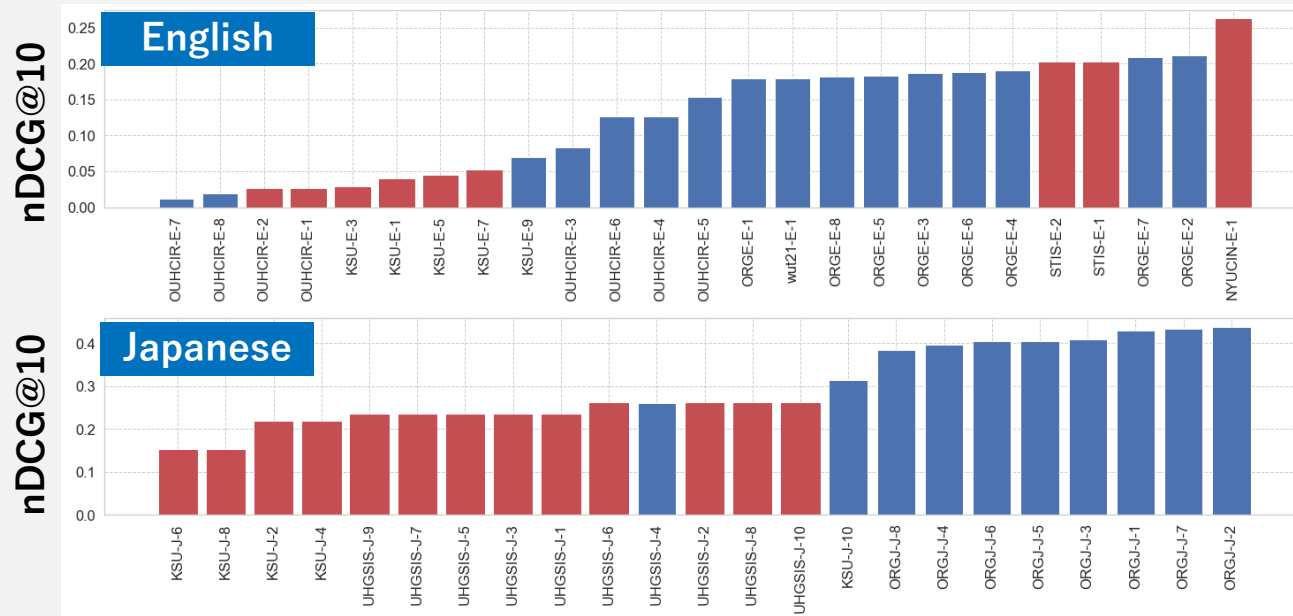
nDCG@10



NYUCIN-E-1 is the top English run, though it is not significantly different from the other top runs including baselines

"Neural" runs

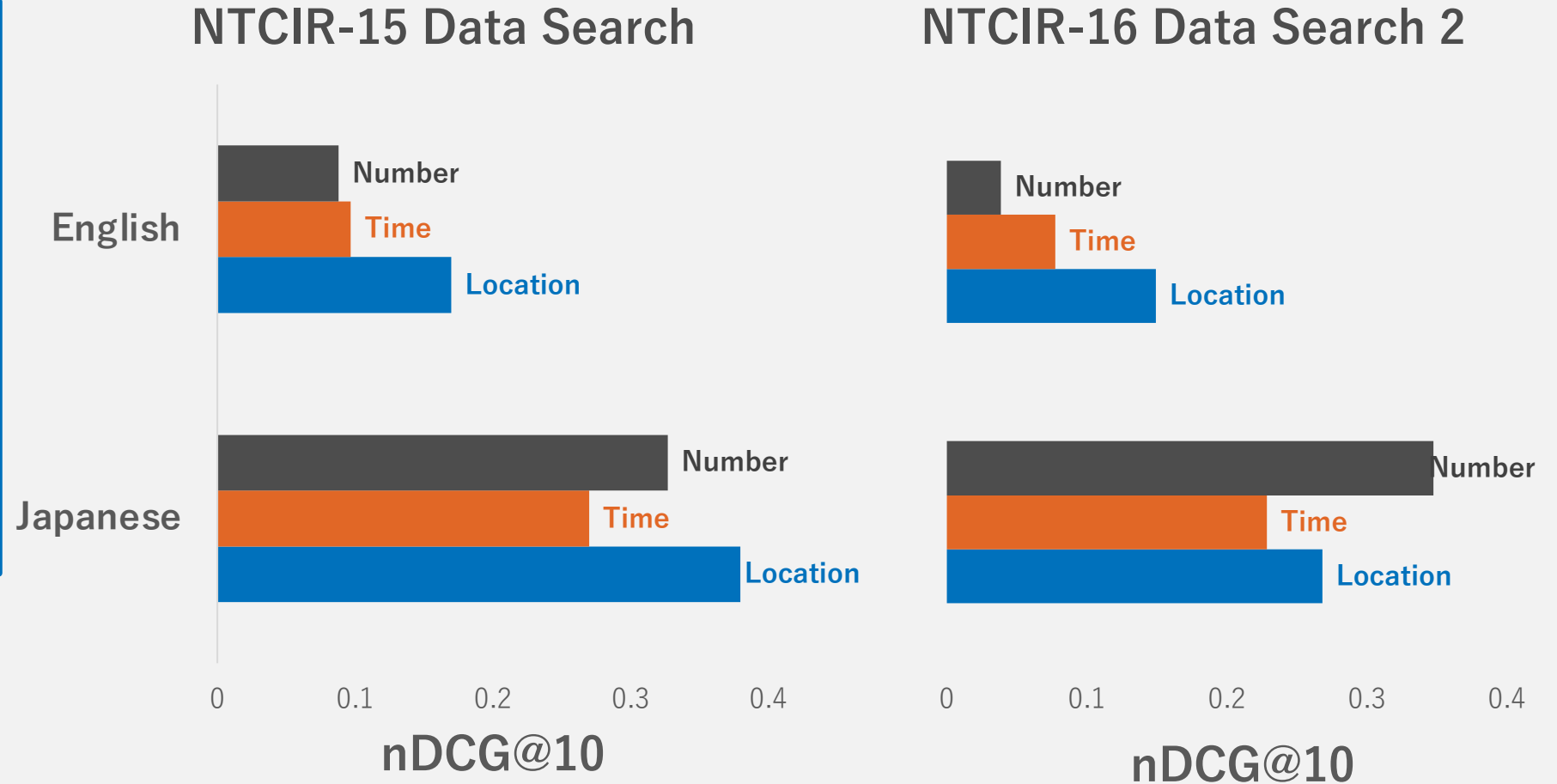
"Data" runs



Again, not conclusive, but the top English runs incorporated "Neural" and "Data"

Queries were classified into three types based on what entities are contained

- **Location** (e.g. tokyo population)
- **Time** (e.g. 1990 population)
- **Number** (e.g. 2m population city)



Processing "Time" in queries seems the most difficult
 (Fairly consistent with the finding in the previous round)

- **Data Search 2 addressed an ad-hoc retrieval task for datasets**
- **Details of the runs will be discussed at Task Session**

[JST] 10:30 ~ 11:30 on July 17; [UTC] 01:30 ~ 02:30 on July 17; [EDT] 21:30 ~ 22:30 on July 16

- 1. NYUCIN at the NTCIR-16 Dataset Search 2 Task
 - 2. KSU Systems at the NTCIR-16 Data Search2 IR Subtask
 - 3. STIS at the NTCIR-16 Data Search Task: Ad-hoc Data Retrieval Ranking with Pretrained Representative Words Prediction
- **Additionally, an invited talk will be given by Prof. Gong Cheng (Nanjing University):**

**Towards Content-Based Dataset Search:
Test Collections and Beyond**