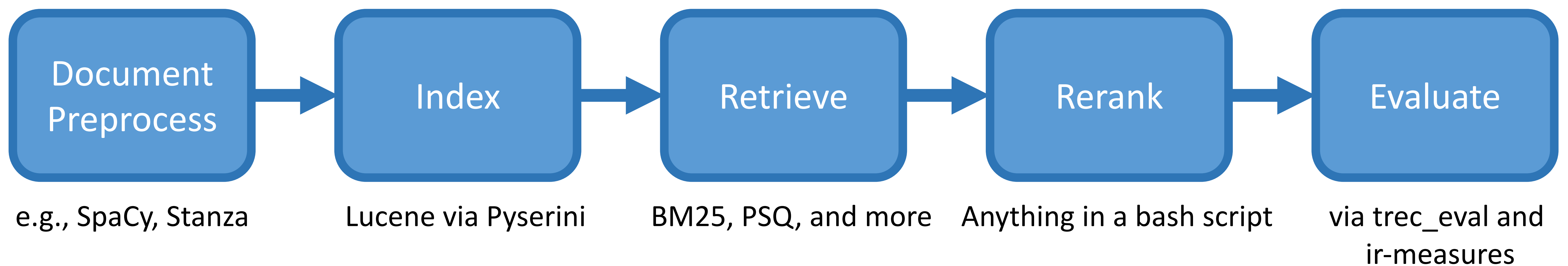


# Patapasco: A Python Framework for Cross-Language Information Retrieval Experiments

Cash Costello Eugene Yang Dawn Lawrie James Mayfield  
Human Language Technology Center of Excellence, Johns Hopkins University



```
1 run:
2   name: hc4_zho_test
3
4 documents:
5   input:
6     format: irds
7     lang: zho
8     path: hc4/zh/test
9     fields: title+text
10  process:
11    normalize:
12      lowercase: true
13      stem: false
14      stopwords: lucene
15      strict_check: true
16      tokenize: spacy
17
18 database:
19   name: sqlite
20   output: true
21
22 index:
23   name: lucene
24   output: true
25
26 topics:
27   input:
28     format: irds
29     lang: eng
30     source: original
31     path: hc4/zh/test
32     fields: title
33
34 queries:
35   output: true
36   process:
37     normalize:
38       lowercase: true
39       stem: false
40       stopwords: lucene
41       tokenize: spacy
42   psq:
43     lang: eng
44     normalize:
45       lowercase: true
46       report: false
47     path: ./eng_zho_transtable.dict
48     stem: false
49     stopwords: lucene
50     threshold: 0.97
51
52 retrieve:
53   b: 0.4
54   k1: 0.9
55   name: bm25
56   number: 1000
57   output: retrieve
58   psq: true
59
60 score:
61   input:
62     format: irds
63     path: hc4/zh/test
```

Select dataset (points to line 8)

Select query (points to line 32)

While there are high-quality software frameworks for information retrieval experimentation, they do not explicitly support cross-language information retrieval (CLIR). **To fill this gap, we have created Patapasco, a Python CLIR framework.** This framework specifically addresses the complexity that comes with running experiments in multiple languages. Patapasco is designed to be **extensible to many language pairs**, to be **scalable to large document collections**, and to support **reproducible experiments** driven by a configuration file.

	CLEF Persian			CLEF Russian			NTCIR Chinese		
	MAP	nDCG	R@1k	MAP	nDCG	R@1k	MAP	nDCG	R@1k
PSQ	0.1370	0.3651	0.4832	0.2879	0.4393	0.7441	0.1964	0.3752	0.5867
QT	0.2511	0.5340	0.6945	0.3857	0.5527	0.9268	0.2186	0.3953	0.6201
DT	0.3420	0.6808	0.8501	0.3408	0.5151	0.8881	0.3413	0.5627	0.8110
HT	0.4201	0.7476	0.9175	0.3975	0.5623	0.9401	0.4810	0.6840	0.9125

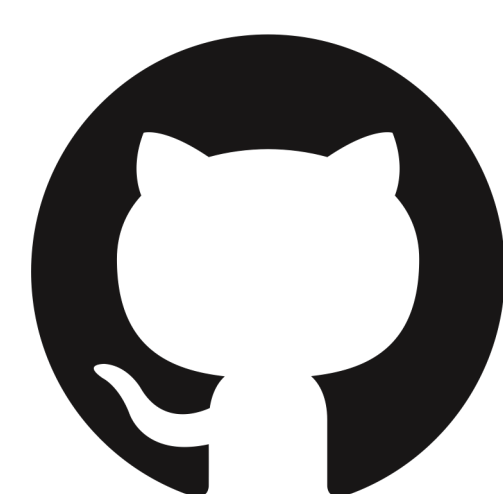
Retrieval Effectiveness Results by Using Patapasco on English Queries and Persian/Russian/Chinese Documents with PSQ, Query Machine Translation (QT), Document Machine Translation (DT), Human Query Translation (HT)

## TREC 2022 NeuCLIR Track



<https://neuclir.github.io/>

- Documents: Chinese / Persian / Russian (released)
- Queries: English (TBA in June)
- Submission due in July
- **Please consider participation! Patapasco makes it easy!**



[github.com/hltcoe/patapasco](https://github.com/hltcoe/patapasco)



human language technology  
center of excellence



JOHNS HOPKINS  
UNIVERSITY

Example Configuration file