

Dense vs. Sparse Representations for News Stream Clustering

Todor Staykovski
Sofia University
Sofia, Bulgaria

Alberto Barrón-Cedeño* Giovanni Da San Martino Preslav Nakov

Qatar Computing Research Institute, HBKU
Doha, Qatar

* Università di Bologna, Italy



Overview

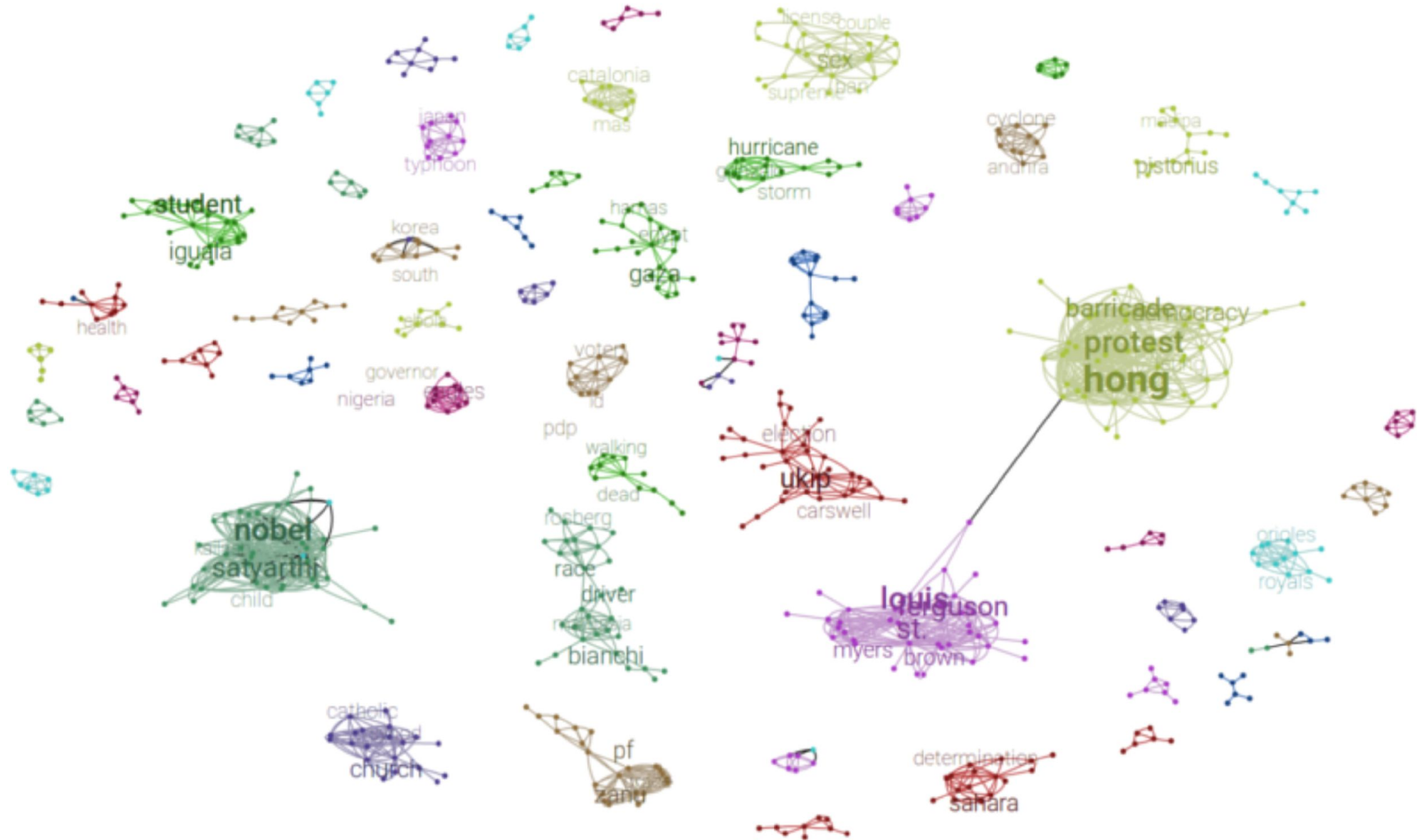
- News stream clustering
- Compare
 - dense vectors
 - sparse vectors
 - combinations
- State-of-the-art results
- B-cubed F1 for evaluation
- Integrated into a real system



Related Work

- **newsLens** [Laban & Hearst, 2017]
 - build clusters and storylines for news
 - used: keywords, graph, community detection
 - we use vectors instead
 - *We are inspired by their pipeline*
- **[Miranda et al., 2018]**
 - clustered a stream of news articles in English, Spanish and German
 - TF-IDF, timestamp
 - we use a local graph instead
 - *We use their dataset*

NewsLens Clustering Example



Our Model

1. identify local topics

- sliding window over time with overlaps
- using article vectors: TF-IDF, doc2vec, combination
- graph edge if $\text{sim}(d_i, d_j) \geq T_1$
- Louvain's community detection

2. merge long-term topics

- using the mean of all news vectors belonging to the topic
- merge topics if
 - *also, use machine* $\text{sim}(t_i, t_j) \geq T_2$ and T_2

Representation

TF-IDF

Doc2Vec

- trained on Signal Media One-million news articles corpus:
 - 265,512 blog articles
 - 734,488 news articles
 - from 93k unique sources
 - over a period of one month

Combined Representations

Combination (unsupervised)

- using T1 and T2

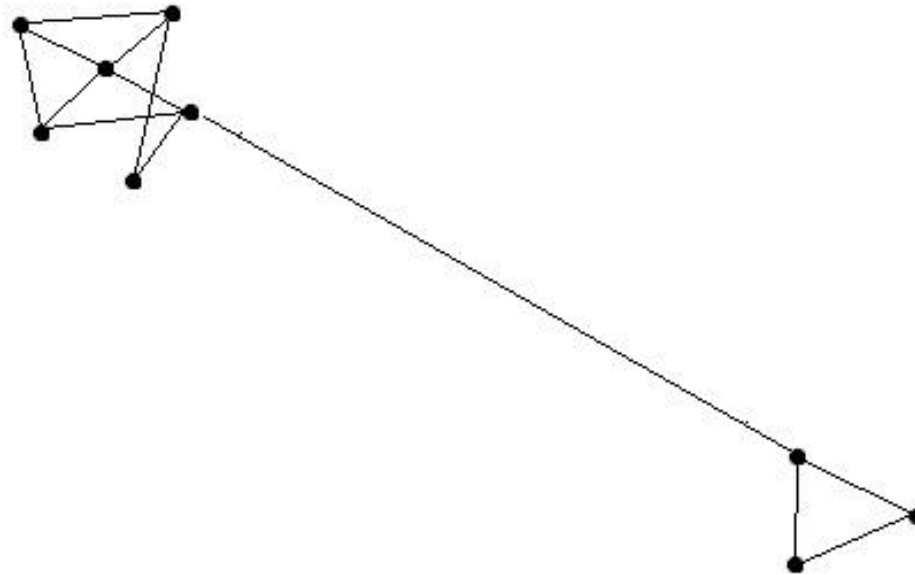
Combination (supervised)

- logistic regression with TF-IDF and doc2vec scores as input
- and whether a pair of articles are in the same cluster in the training data

Merging Long-Term Topics: Louvain's Community Detection

Topic 1

Protesters in England call for change to cricket governance



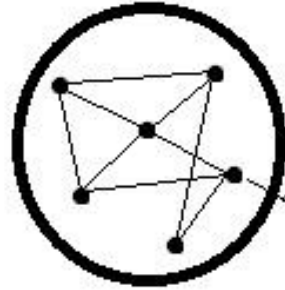
Topic 2

Gujarat quota protests turn violent

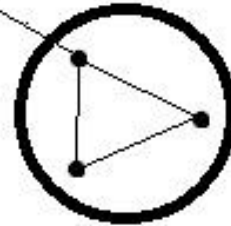
Merging Long-Term Topics: Louvain's Community Detection

Topic 1

Protesters in England call for change to cricket governance



Community detection correctly assigns the two topics to different communities



Topic 2

Gujarat quota protests turn violent

Dataset

Partition	Docs	Tokens	Clusters	Cluster Size
Train	12,233	434 \pm 364	593	21 \pm 32
Test	8,726	521 \pm 495	222	39 \pm 88

Sebastiao Miranda, Arturs Znotins, Shay B. Cohen, Guntis Barzdins. *Multilingual clustering of streaming news*. EMNLP'18.

Train

- December 18, 2013 to February 2, 2014
- No time gaps

Test:

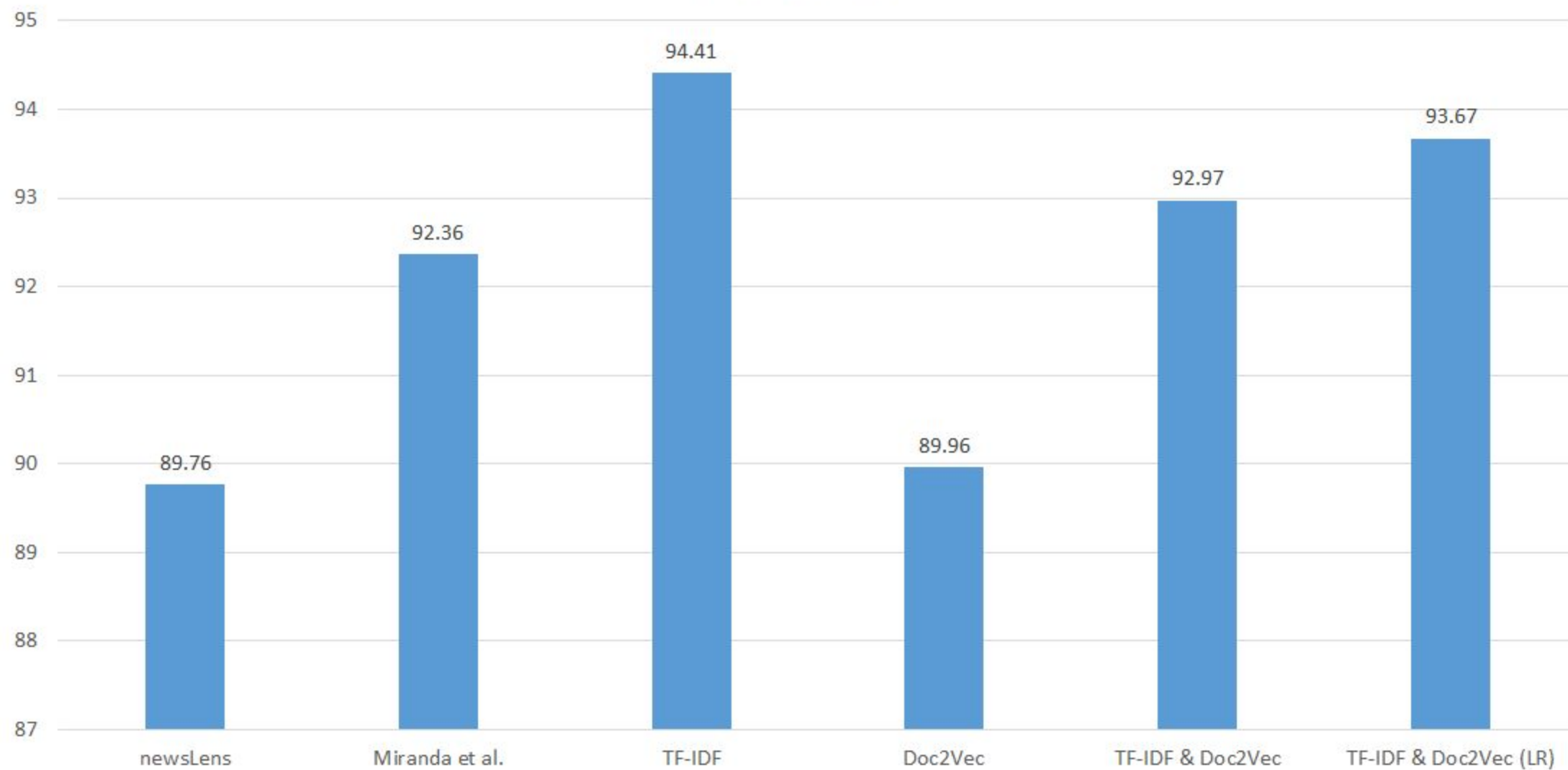
- November 2, 2014 to August 25, 2015
- Gaps as long as 3 months --> need to find long-term topics

Evaluation: BCubed measures

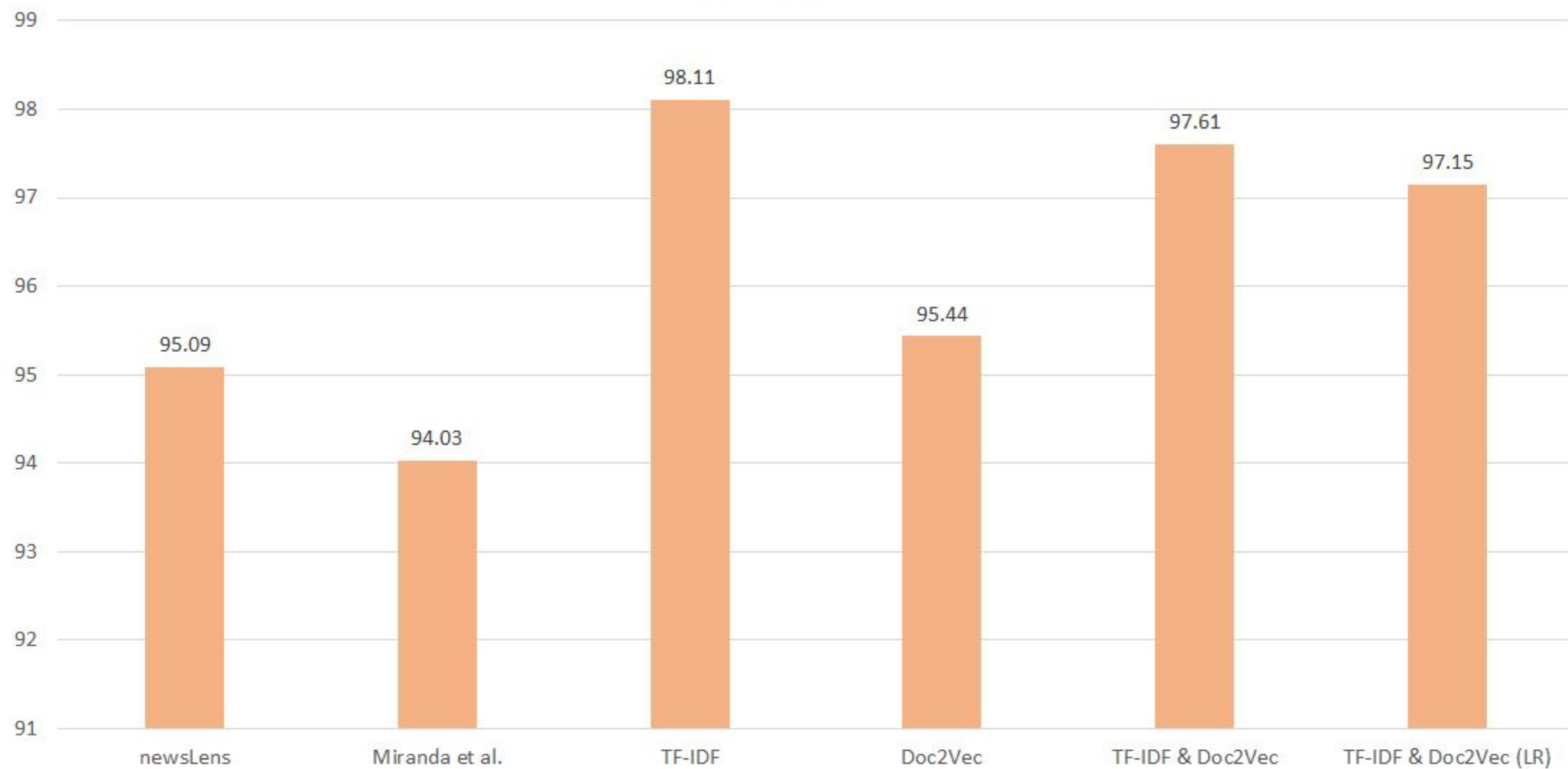
Favor clusterings that

- split a cluster that mixes two categories into two **pure clusters** (**cluster homogeneity**)
- unify two clusters that contain only items from the same category (**cluster completeness**)
- **add** an item of a different category **to an already noisy cluster** instead of a pure one
- **make small errors in a big cluster** rather than a large number of small errors in small clusters

Results BCubed F_1



Results Standard F_1



Evaluation Results

model	BCubed						clusters
	F ₁	P	R	F ₁	P	R	
baselines							
newsLens [LH17]	89.76	94.37	85.58	95.09	95.90	94.30	873
Miranda et al. [MZCB18]	92.36	94.57	90.25	94.03	98.14	90.25	326
unsupervised							
TF-IDF	94.41	95.16	93.66	98.11	97.60	98.63	484
doc2vec	89.96	93.00	87.12	95.44	95.55	95.34	785
TF-IDF & doc2vec	92.97	95.75	90.34	97.61	97.73	97.48	663
supervised (LR)							
TF-IDF	94.30	94.87	93.73	98.08	97.46	98.71	485
TF-IDF & doc2vec	93.67	92.71	94.65	97.15	95.39	98.98	431

Expected: 222 clusters

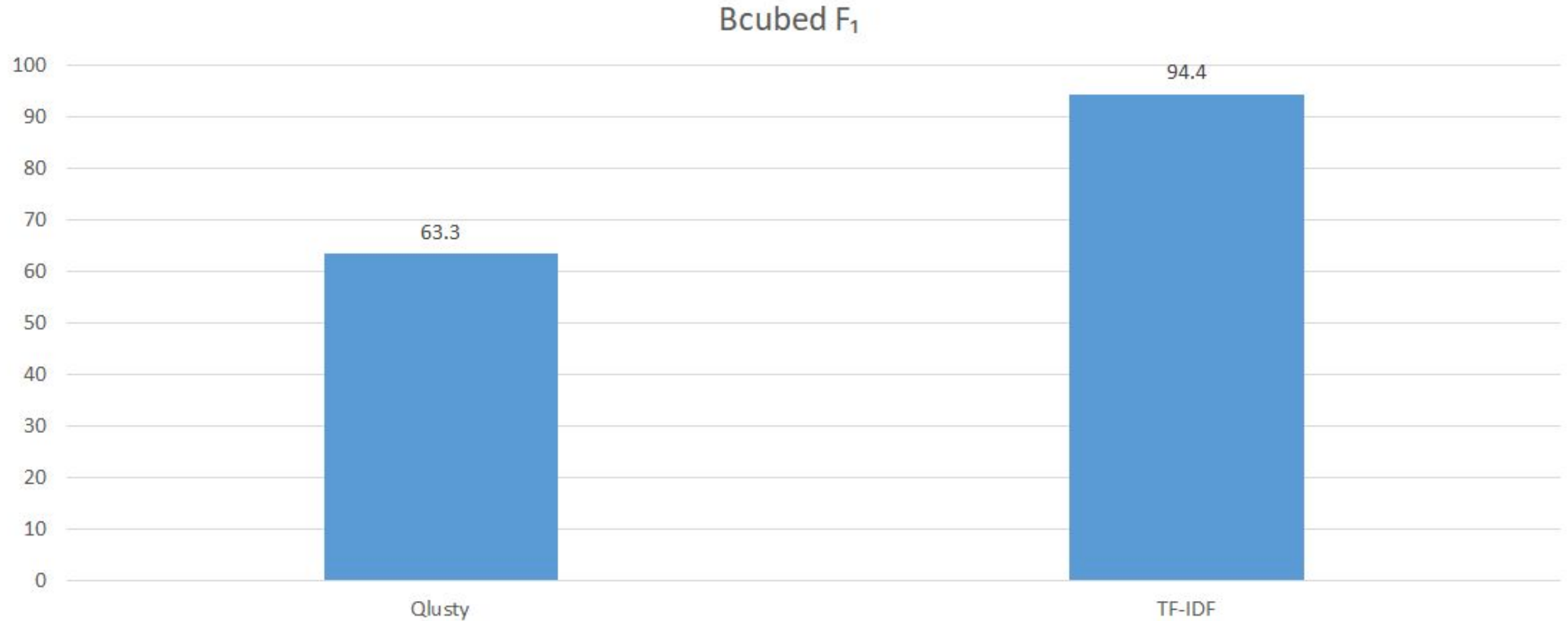
Discussion:

Impact of Louvain's Community Detection

	BCubed						
	F ₁	P	R	F ₁	P	R	clusters
before	91.55	93.76	89.44	94.36	95.55	93.20	488
after	94.41	95.16	93.66	98.11	97.60	98.63	484

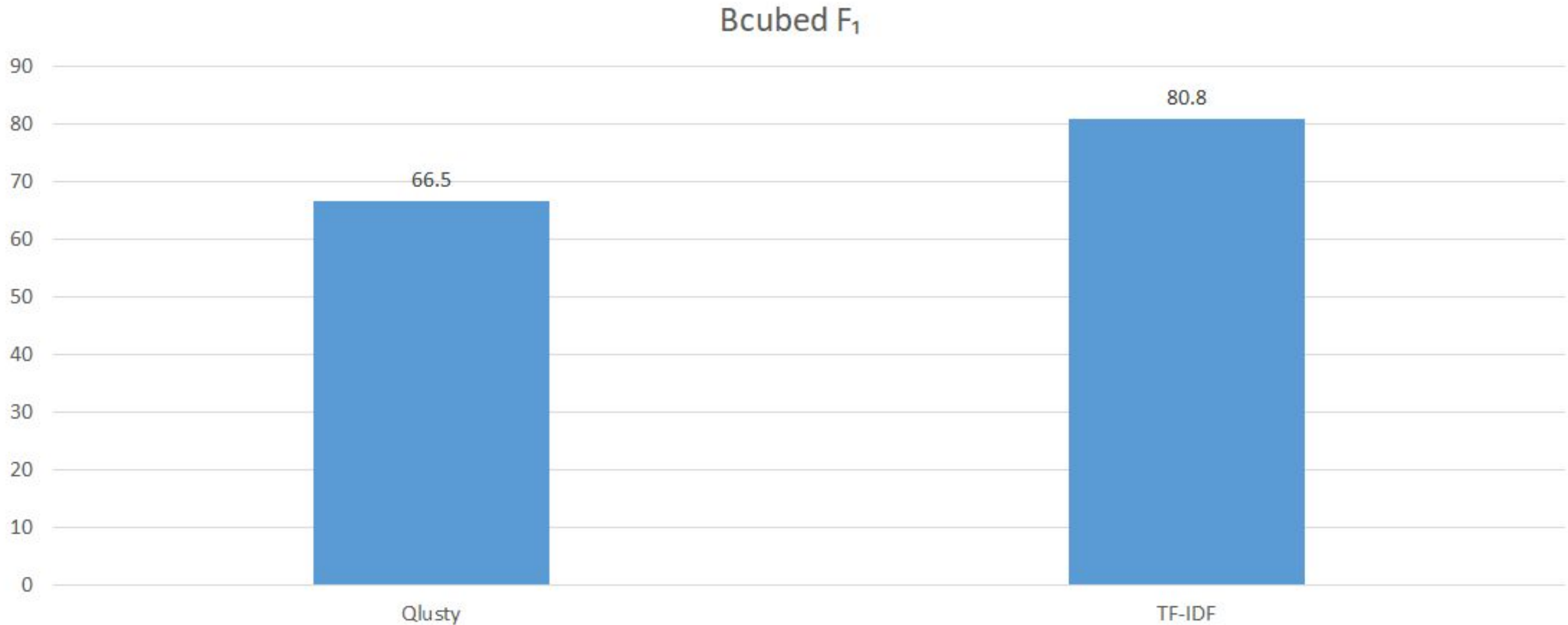
Discussion:

DBSCAN vs. Our Model: Miranda&al.-2018 corpus



Discussion:

DBSCAN vs. Our Model: Meter Corpus



Our Clustering Algorithm (w/ TF-IDF) Powers

Tanbih: <https://www.tanbih.org/>




BUSINESS

SPORTS

POLITICS

ARTS & ENT

TECH & SCI



A news card featuring a photo of a man in a dark shirt gesturing with his hands. In the top left corner are Facebook and Twitter icons. In the top right is a blue tag labeled 'CULTURAL-IDENTITY'. In the bottom left is a red circle with a white 'M'. In the bottom right is a grey button 'Propaganda?' and a green button 'unlikely'.

Joao Felix transfer: Man Utd and Chelsea 'to rival Liverpool for £60m wonderkid'

1 of 6 SPORT FOOTBALL TRANSFER NEWS

4 hours ago



A news card featuring the word 'تنبيه' (TANBIH) in large black Arabic script. Below it is the 'TANBIH' logo with 'T' in blue, 'AN' in grey, and 'BIH' in red. In the top left are Facebook and Twitter icons. In the top right is a blue tag labeled 'POLITICAL'. In the bottom left is a black circle with a yellow 'SPUTNIK' logo. In the bottom right is a grey button 'Propaganda?' and a green button 'unlikely'.

Venezuelan FM Warns Foreign Journalists Against Working Without Accreditation

1 of 3 LATAM

4 hours ago



Each story on TANBIH will have multiple articles from multiple news outlets. If you are interested in a story, you can go through all articles of a story from these previous and next buttons

Conclusion and Future Work

Summary

- Compared dense vs. sparse vectors
- State-of-the-art results
- B-cubed F1 for evaluation
- Integrated into a real system

Future work

- Improve topic matching
- Cross-language extension



<http://www.tanbih.org/>