

# Exploring Thematic Coherence in Fake News

Martins Samuel Dogo

Deepak P

Anna Jurek-Loughrey



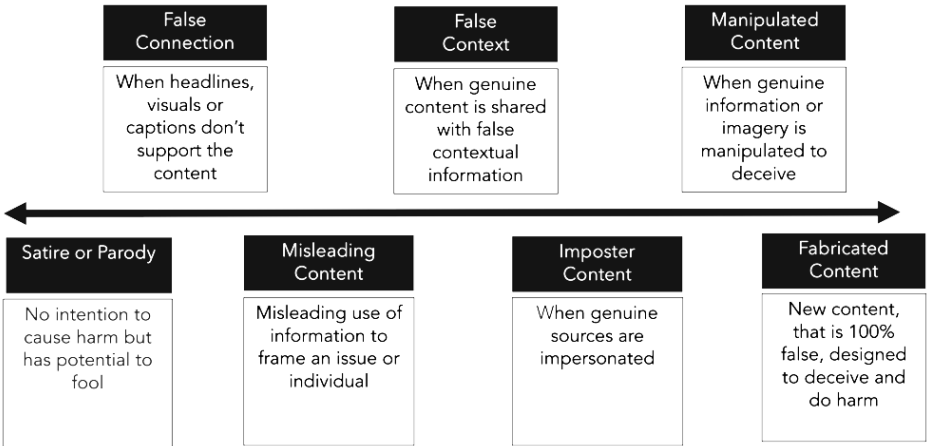
**QUEEN'S  
UNIVERSITY  
BELFAST**

# Fake news: definition

Definition: a news item that contains **deliberately** and **verifiably** falsified information.

WHAT

are the types of mis- and dis-information?



Adapted from Claire Wardle, First Draft News

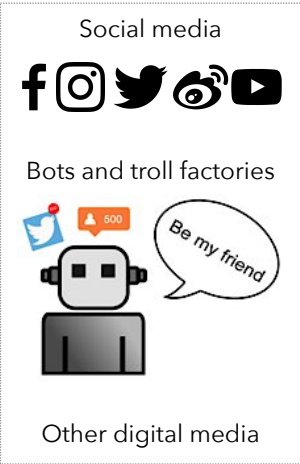
WHY

&

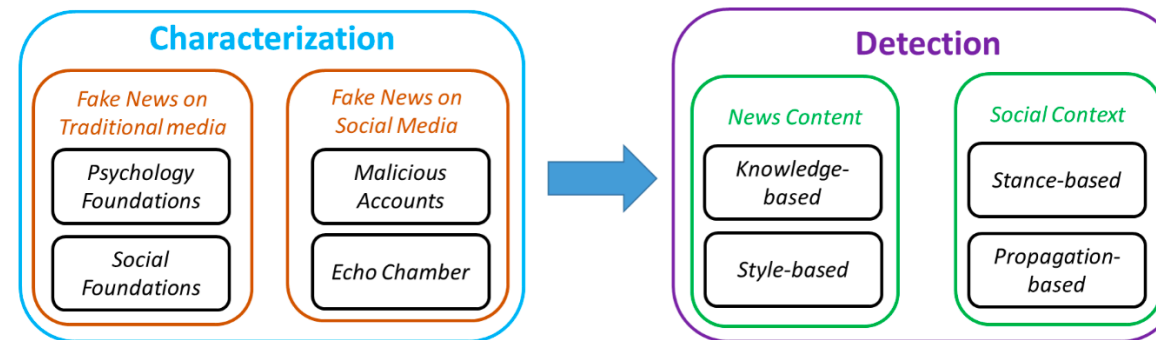
HOW

are they shared?

- Poor Journalism
- To Parody
- To provoke or 'punk'
- Passion
- Partisanship
- Profit
- Political Influence/power
- Propaganda



# Fake news: characterisation and detection



Source: Shu et al. (2017)

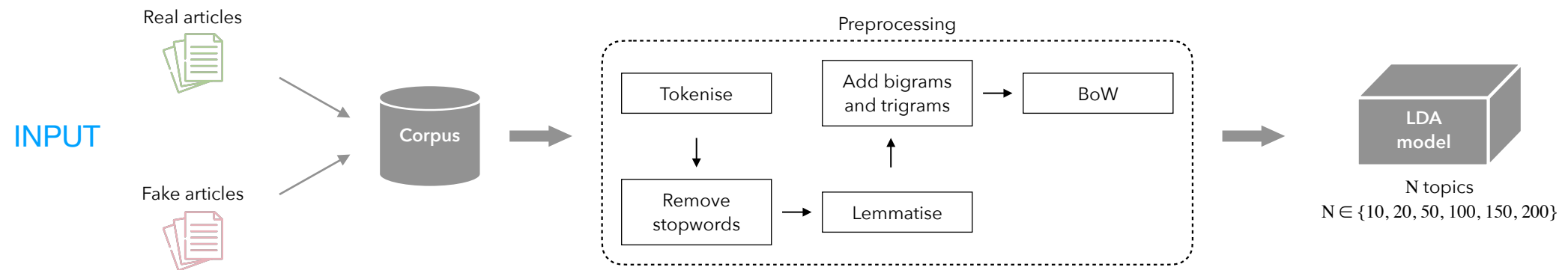
- Textual content is the most abundant type of information in fake news datasets.
- Style-based features include lexical, syntactic, and latent features such as embeddings.
- Characterisation is key for developing unsupervised fake news detection techniques.
- We focus on characterising fake news in a new way using topic modelling.

# Objectives of this study

- Assess the importance of internal consistency within articles as a high-level feature to distinguish fake from real news.
- Use latent themes/topics to analyse the coherence of articles.
- Experiment with real-world datasets to demonstrate the efficacy of our proposed approach.

# Method:

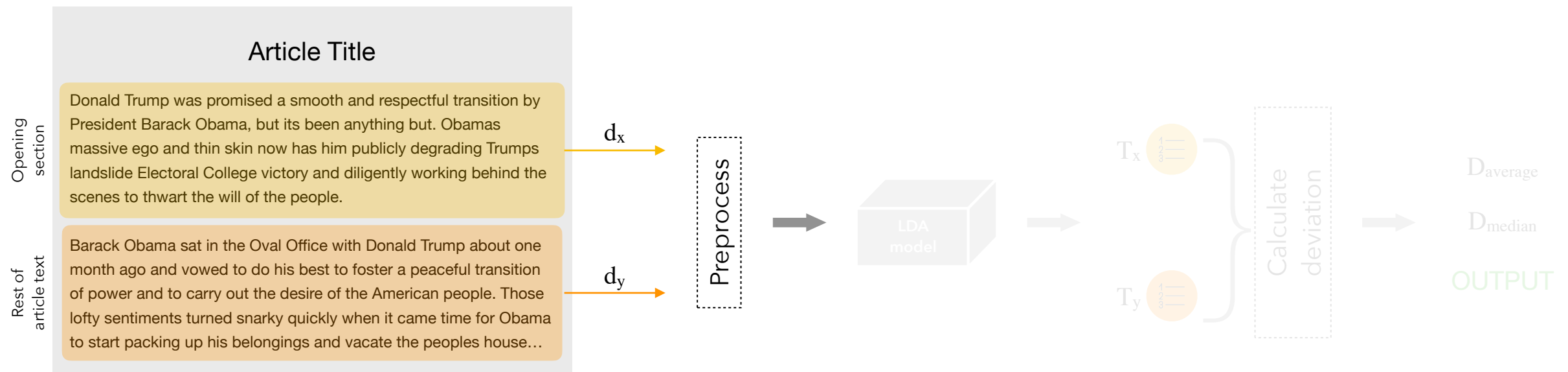
## Step 1 – build topic models



- Create a corpus consisting of entire dataset.
- Preprocess corpus and build LDA model with N topics.

# Method:

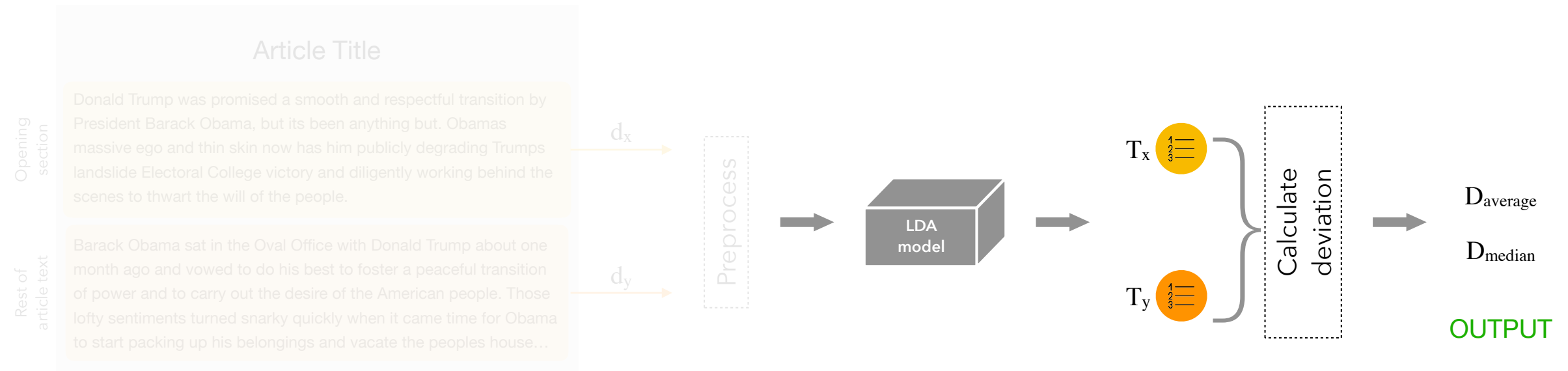
## Step 2 – split articles



- Split each article into two documents:  
its opening section (first five sentences excl. the title;  $d_x$ ), and its remainder ( $d_y$ ).
- Preprocess each document.

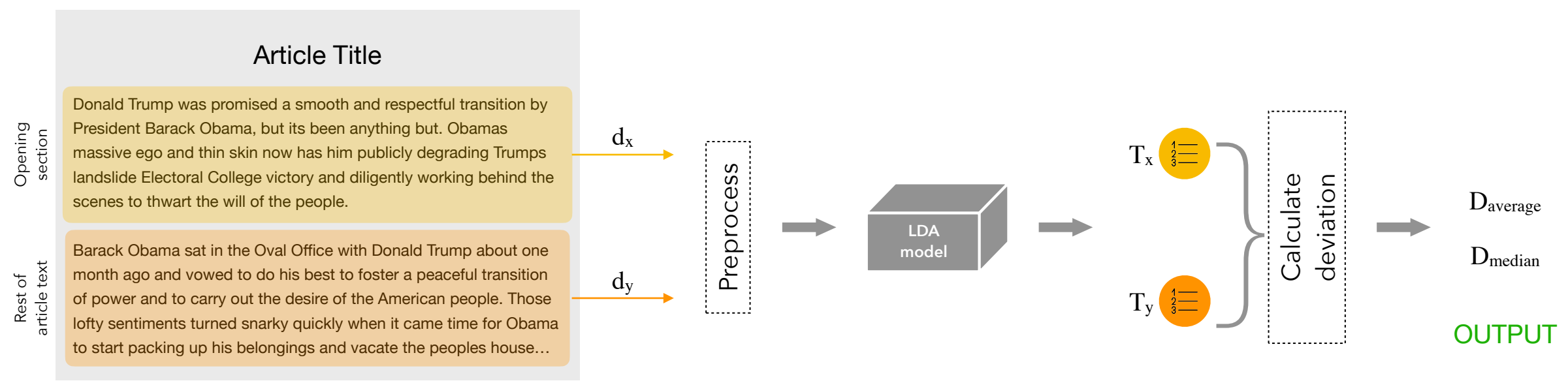
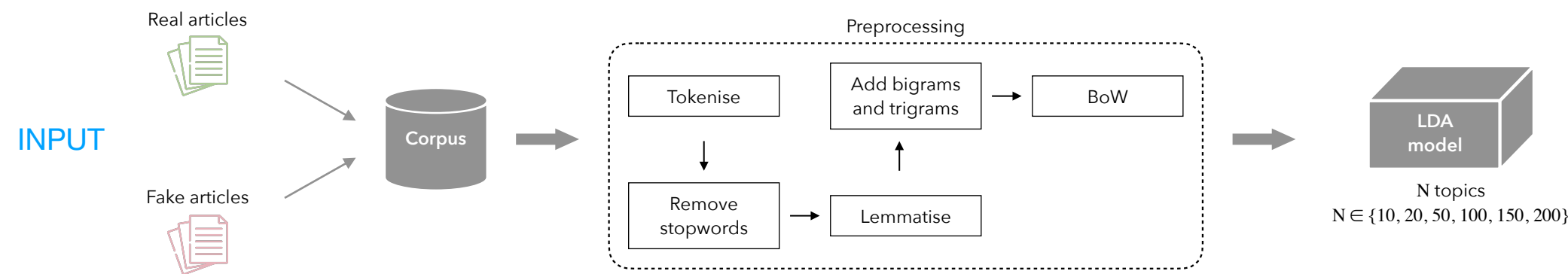
# Method:

## Step 3 – extract topics and calculate deviation



- Use LDA model to extract  $N$  topics from both documents.
- Calculate the deviation between their topic distributions ( $T_x$  and  $T_y$ ) using distance metrics. Find mean and median deviation across all values of  $N$ .

# Method: all steps





# Experimental evaluation

## ■ Datasets\*:

- FakeNewsAMT & Celebrity: 641 articles
- BuzzFeed Political: 243 articles
- BuzzFeed Webis: 1,545 articles
- George McIntyre: 5,547 articles
- ISOT Lab: 36,147 articles
- POLIT False-n-Legit: 256 articles
- Syrian Violations Documentation Center: 664 articles

## ■ Distance metrics:

- Chebyshev (Chessboard) distance
- Euclidean distance
- Squared Euclidean distance

## ■ Parameters:

- Length of article opening: 5 sentences
  - Number of topics: 10, 20, 50, 100, 150, 200
- 
- We evaluate the differences in coherence/deviation between fake and real articles using the T-test at 5% significant level.

\* Articles remaining after preprocessing dataset.

# Results and observations:

## coherence

	Fake	Real	Fake	Real
Dataset	Mean $D_{Ch}(f)$	Mean $D_{Ch}(r)$	Median $D_{Ch}(f)$	Median $D_{Ch}(r)$
AMT+C	0.2568	0.2379	0.2438	0.2285
BuzzFeed-Political	0.2373	0.2149	0.2345	0.2068
BuzzFeed-Web	0.2966	0.2812	0.2863	0.2637
GMI	0.4580	0.4241	0.4579	0.4222
ISOT	0.3372	0.2971	0.3369	0.2989
POLIT	0.2439	0.1939	0.2416	0.1894
SVDC	0.2975	0.2517	0.2934	0.2435

- Real news articles are more coherent (i.e., show less deviation) than fake ones.
- In other words, there is a greater shift in the topics the body of a fake article discusses, from its opening section.

# Results and observations:

## statistical test

Dataset	p-value ( $D_{Ch}$ )	p-value ( $D_E$ )	p-value ( $D_{SE}$ )
AMT+C	0.144	0.126	0.116
BuzzFeed-Political	0.0450	0.0147	0.0287
BuzzFeed-Web	0.209	0.209	0.207
GMI	0.0480	0.00535	0.0106
ISOT	0.00319	0.000490	0.000727
POLIT	0.000660	0.0000792	0.0000664
SVDC	0.000684	0.0000112	0.0000789

- Our results show statistical significance in most datasets.
- This suggests that it may be possible to stylistically characterise fake news using the transition of topics within its text.

# Conclusion

In this work, we:

- investigated the coherence of topics discussed in fake and real articles in seven cross-domain datasets.
- proposed a new method for characterising fake news using only textual data.
- empirically demonstrated that fake articles are thematically less coherent than real ones.

**Future work:** investigate how these characteristics can be utilised to develop unsupervised models for fake news detection.

# Thank you

If you have any questions or would like to give some feedback, please write to me:

Martins Samuel Dogo

[mdogo01@qub.ac.uk](mailto:mdogo01@qub.ac.uk) | [martinssamuel.com](http://martinssamuel.com)



**QUEEN'S  
UNIVERSITY  
BELFAST**