**iPS** Institute of
Policy Studies

ENGAGING MINDS, EXCHANGING IDEAS

# Digital Frontiers Seminar:
# Man and Machine (2): Online and Offline Sensing of Public Sentiments on National Day Rally 2016

## Thursday, 23 February 2017
## Seminar Room 3-1, Level 3, Manasseh Meyer Building

**LKY** Lee Kuan Yew
School *of* Public Policy

National University of Singapore

# Part I
# Big Data and Online Sensing

## Associate Professor Jiang Jing
Living Analytics Research Centre
School of Information Systems
Singapore Management University

Lee Kuan Yew
School of Public Policy
National University of Singapore

IPS Institute of
Policy Studies

**Engaging Minds, Exchanging Ideas**

# Part I
# Big Data and Online Sensing

## Dr Palakorn Achananuparp

Research Scientist
Living Analytics Research Centre
School of Information Systems
Singapore Management University

Lee Kuan Yew School of Public Policy
National University of Singapore

IPS Institute of Policy Studies

**Engaging Minds, Exchanging Ideas**

# Online Sensing for NDR 2016

# Background

- How did Singaporeans respond to the topics covered in NDR 2016?
  - Traditional survey – scale is limited
  - Analyzing social media – data is public, number of online users is big
- LARC's capabilities in social sensing
  - Continuous crawling of social media content (Twitter, blogs, forums, Facebook, etc.) since 2011
  - Development of advanced data analytics techniques

# Objectives of the Study

- By collecting and analyzing online content related to NDR 2016, we aim to understand the general public's interests and sentiments on NDR topics.

- By comparing the findings of online sensing with those from the survey done by IPS, we aim to gauge the representativeness of public interests and sentiments harnessed from social media.

# Research Questions

- Which online platforms were the most popular for discussing NDR 2016?
- Which NDR topics gained the most attention online?
- What were the online sentiments towards NDR topics?
- Which NDR topics gained the most divided sentiments?
- How do topics differ across different online platforms?
- How do sentiments toward various topics differ across different online platforms?
- How close are online sensing results to survey results?

# Data Sources

- Online content between Aug 7, 2016 and Sep 10, 2016. (Two weeks before and two weeks after NDR 2016)
- Platforms:
  - Twitter: ~150K Singapore-based users
  - Facebook: public pages identified by IPS
  - Blogs: 243 blogging sites (Mothership, Online Citizen, etc.)
  - Forums: Hardwarezone, AsiaOne, Sammy's Alfresco Coffee Forum, Reddit Singapore
  - Websites of mainstream media: Straits Times, Today, etc.
- Two types of content: posts and comments
  - Posts: content published independent of any previous content
  - Comments: content published in response to a post

# Methodology

collected online posts and comments

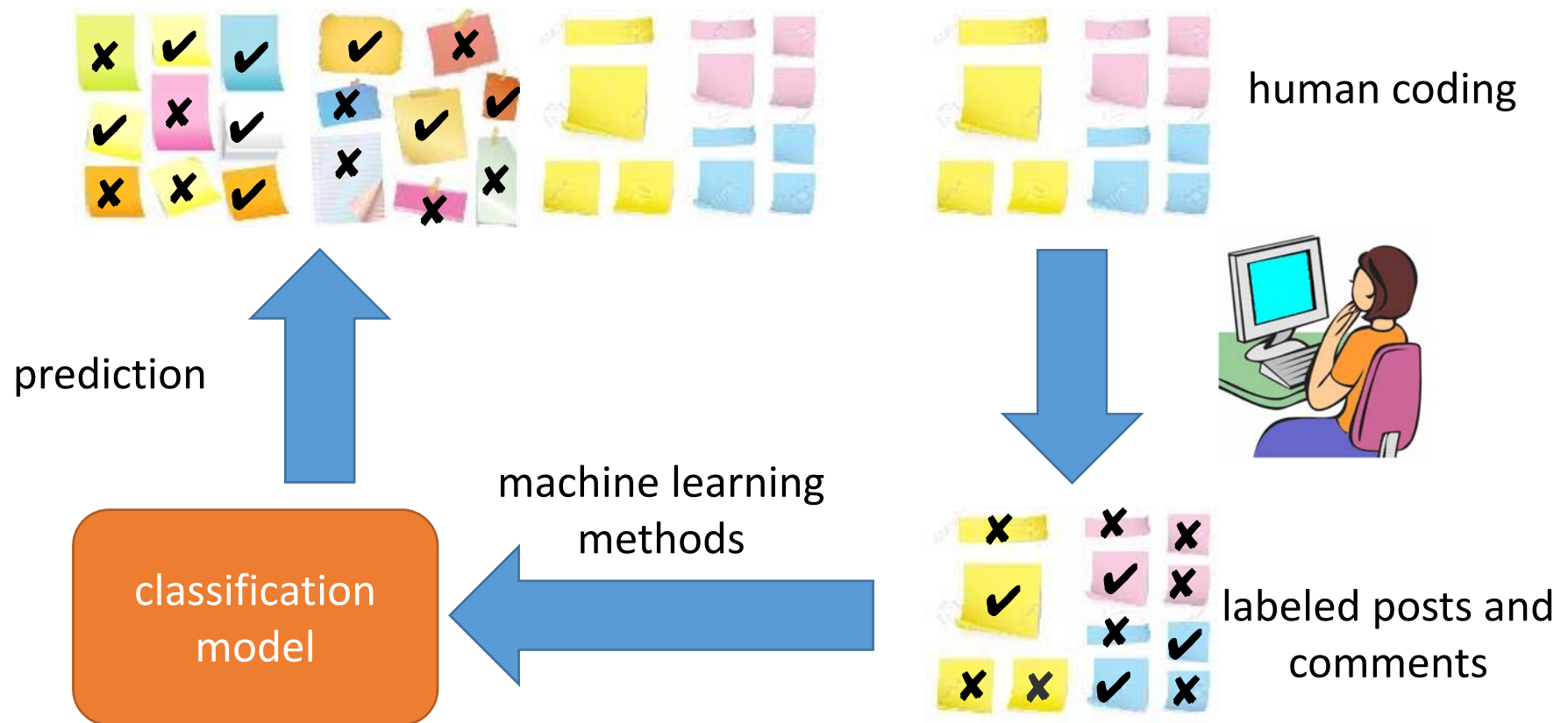Step 1: Identify NDR-related content

Step 2: Categorize posts and comments by topics

Step 3: Classify sentiments

We use a combination of human coding and machine learning to perform the tasks above.

# Methodology

# Methodology

- Step 1: Identify NDR-related content
  - Labels: relevant and irrelevant
- Step 2: Categorize posts and comments into topics
  - Labels: 26 topics identified by IPS
- Step 3: Classify sentiments
  - Labels: positive and negative

# Topics

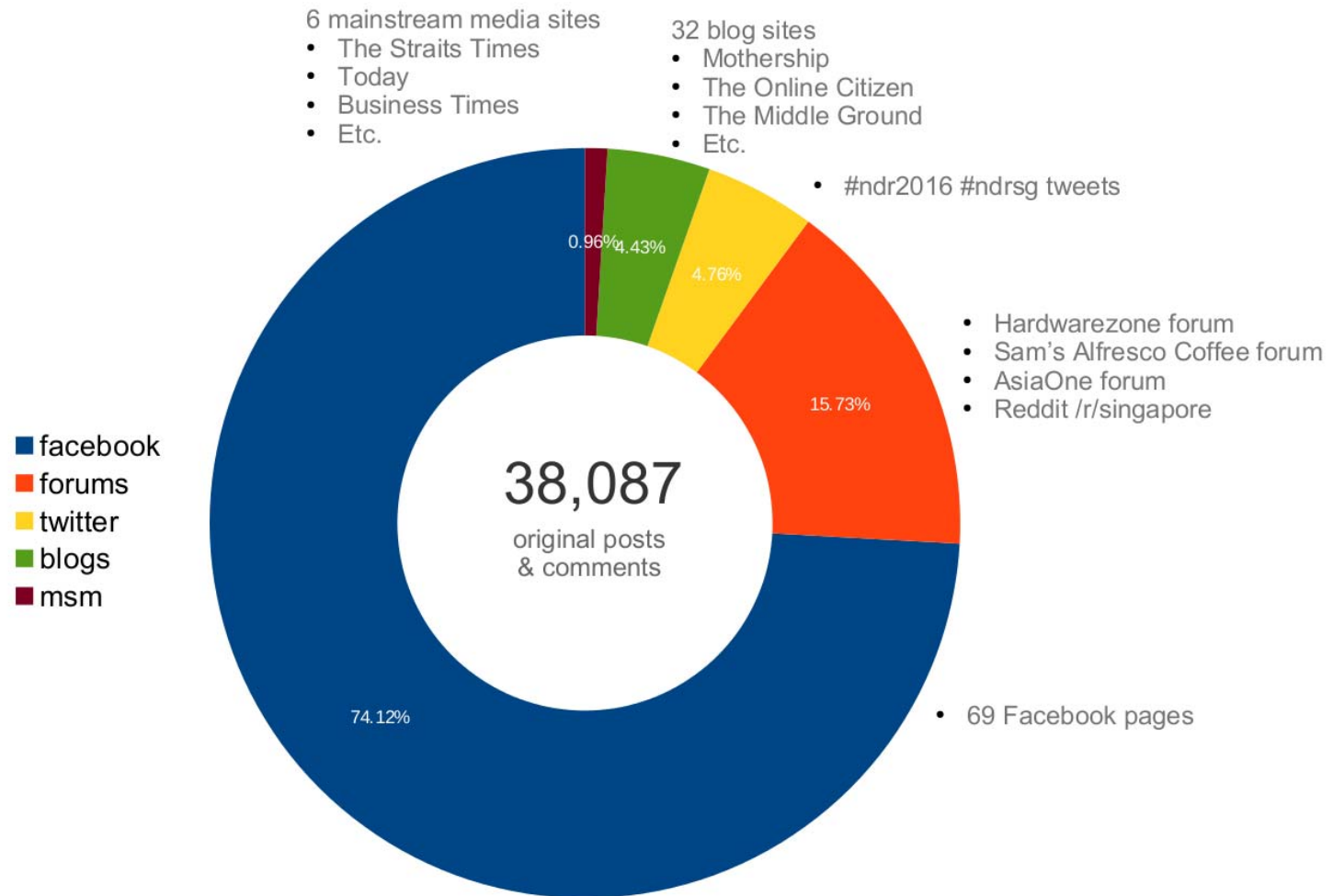| Elected Presidency | Eligibility Criteria | Minority Representation | Social Cohesion |
|---|---|---|---|
| Foreign Policy | South China Sea Dispute | Bilateral Relations | Terrorism |
| Leadership Succession | ARS | Good Politics | Progress Together |
| Entrepreneurship and Innovation | Disruption to Economy | Economy Performance | SkillsFuture |
| Heng Swee Keat | Lawrence Wong | Minister for Finance | Eldershield |
| PM Lee | Unwell | Schooling | Pokemon |
| Delivery of Speech | Others | | |

# Classification Accuracy

- Step 1 (relevant content): 0.8469
- Step 2 (topic categorization): 0.7351 (on blogs), 0.6573 (on tweets)
- Step 3 (sentiment classification): 0.84


- Typically a classification accuracy of 0.75 – 0.8 is considered good.
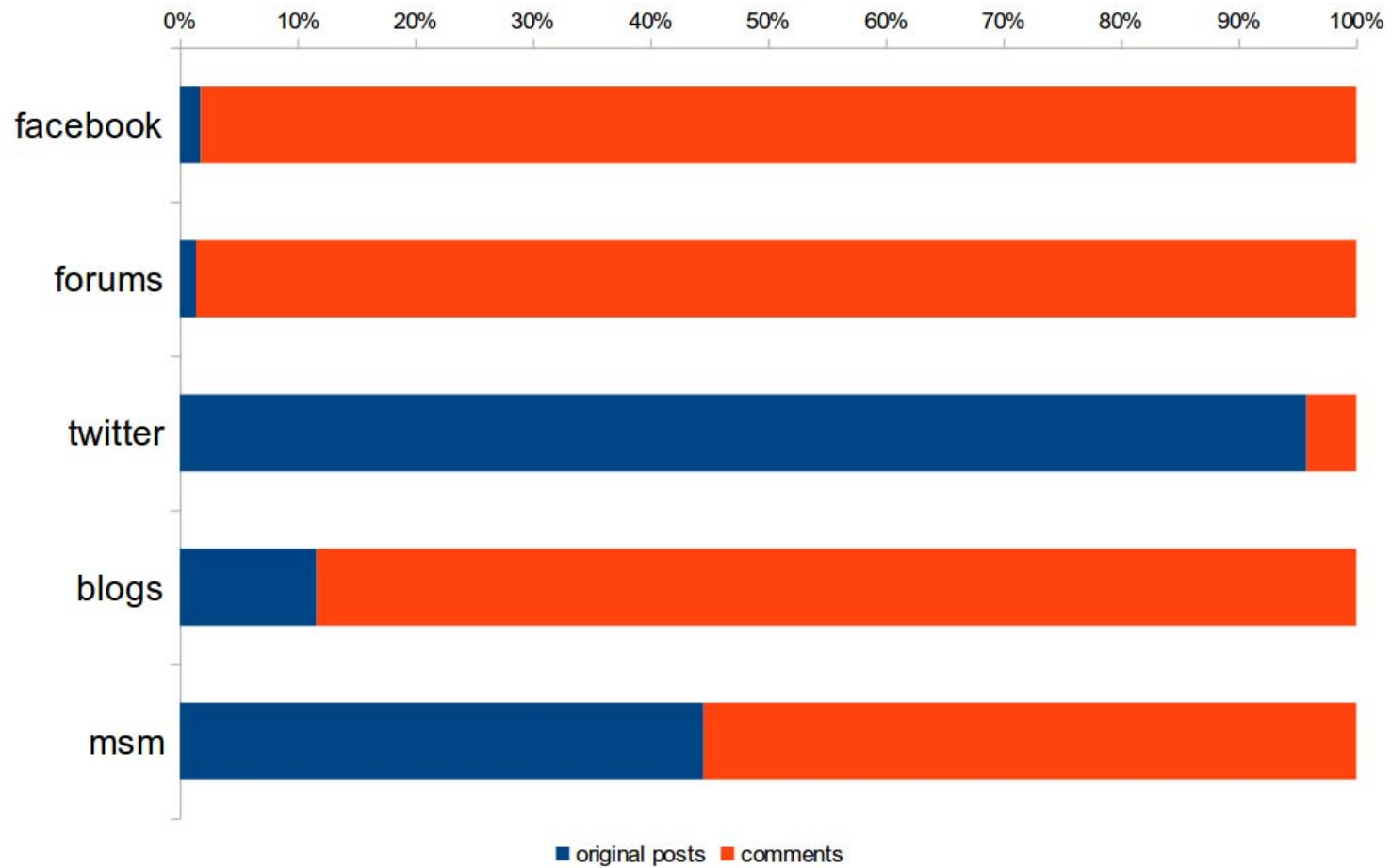- Accuracy goes down when there are many categories and/or there is little training data.

Which online platforms were the most **popular**?

**Content Volume** = Original Posts + Comments
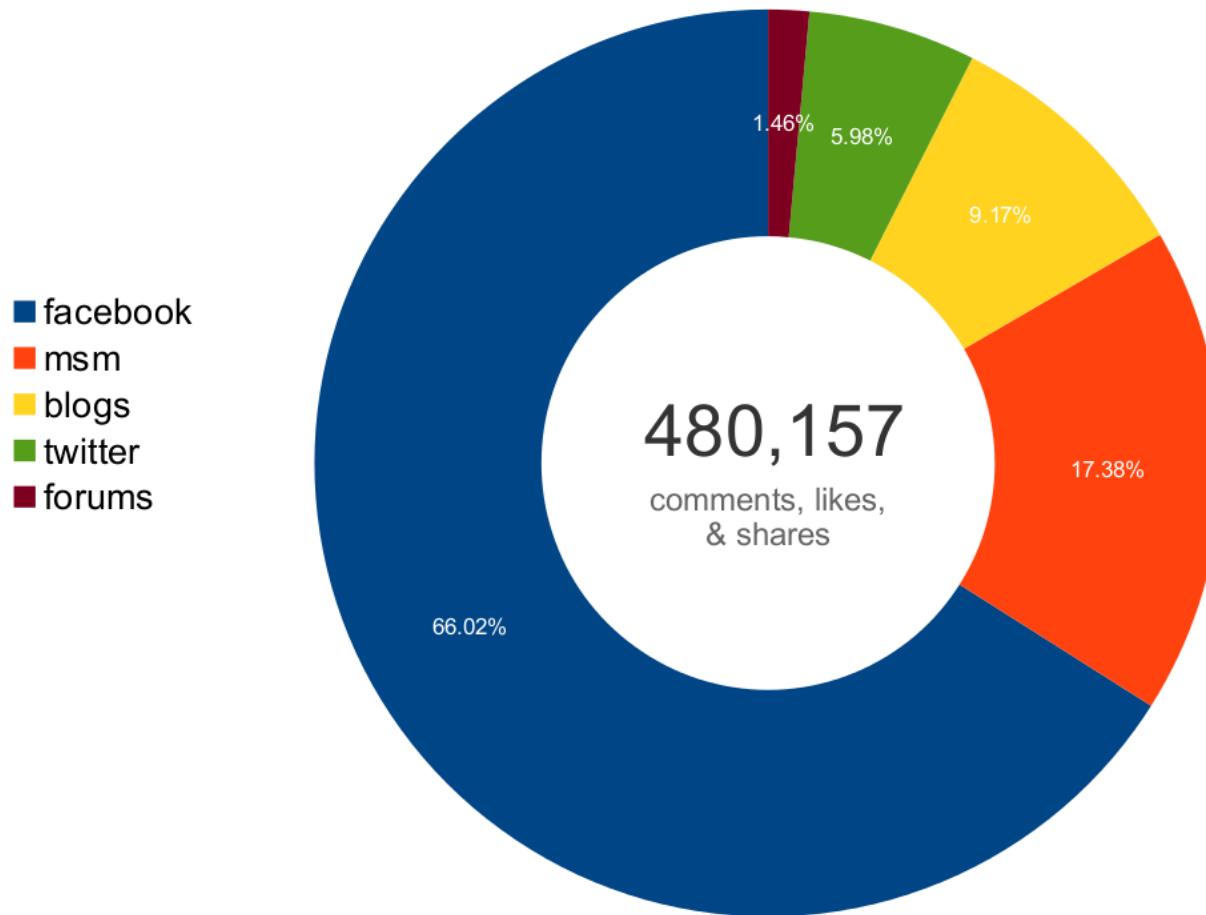
# Content Volume By Platforms

6 mainstream media sites
- The Straits Times
- Today
- Business Times
- Etc.

32 blog sites
- Mothership
- The Online Citizen
- The Middle Ground
- Etc.

- #ndr2016 #ndrsg tweets

- Hardwarezone forum
- Sam's Alfresco Coffee forum
- AsiaOne forum
- Reddit /r/singapore

- 69 Facebook pages

**38,087**
original posts
& comments

- facebook
- forums
- twitter
- blogs
- msm

0.96% 4.43%
4.76%
15.73%
74.12%

# Proportions of Content Types By Platforms

**Engagement Volume** = Comments + Likes + Shares

# Engagement Volume By Platforms



- facebook
- msm
- blogs
- twitter
- forums

480,157
comments, likes,
& shares

1.46%  5.98%
9.17%
17.38%
66.02%

# Proportions of Engagement Types By Platforms

# Numbers of Users By Platforms



facebook
twitter
forums
blogs
msm

27,674

posters and commenters

- 561 Posters
- 5334 Retweeters
- 55 Commenters

72.63%
21.50%
4.66%
0.83%
0.38%

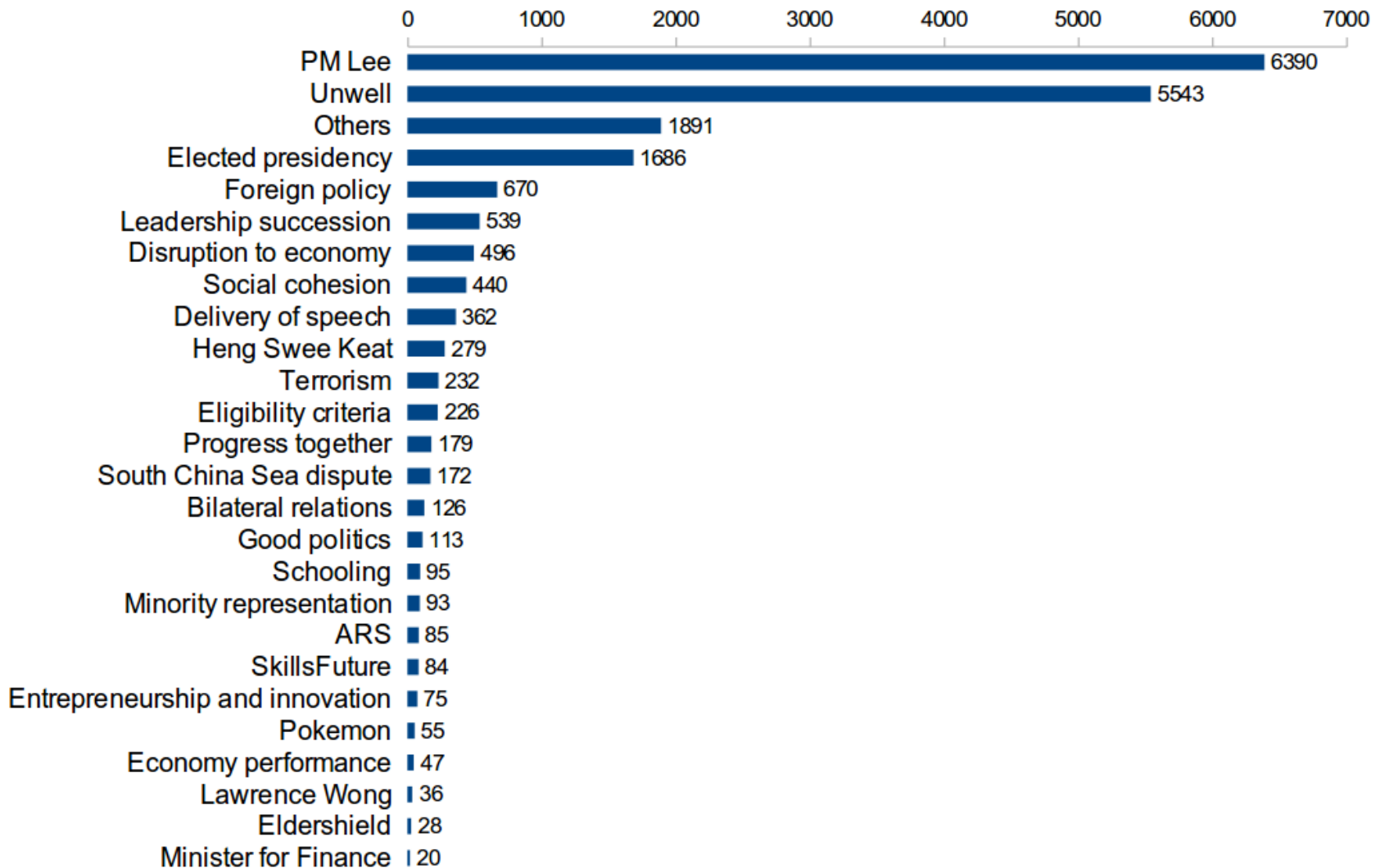# Proportions of Users Types By Platforms

# Which NDR topics gained the most attention online?

# Content Volume By Topics



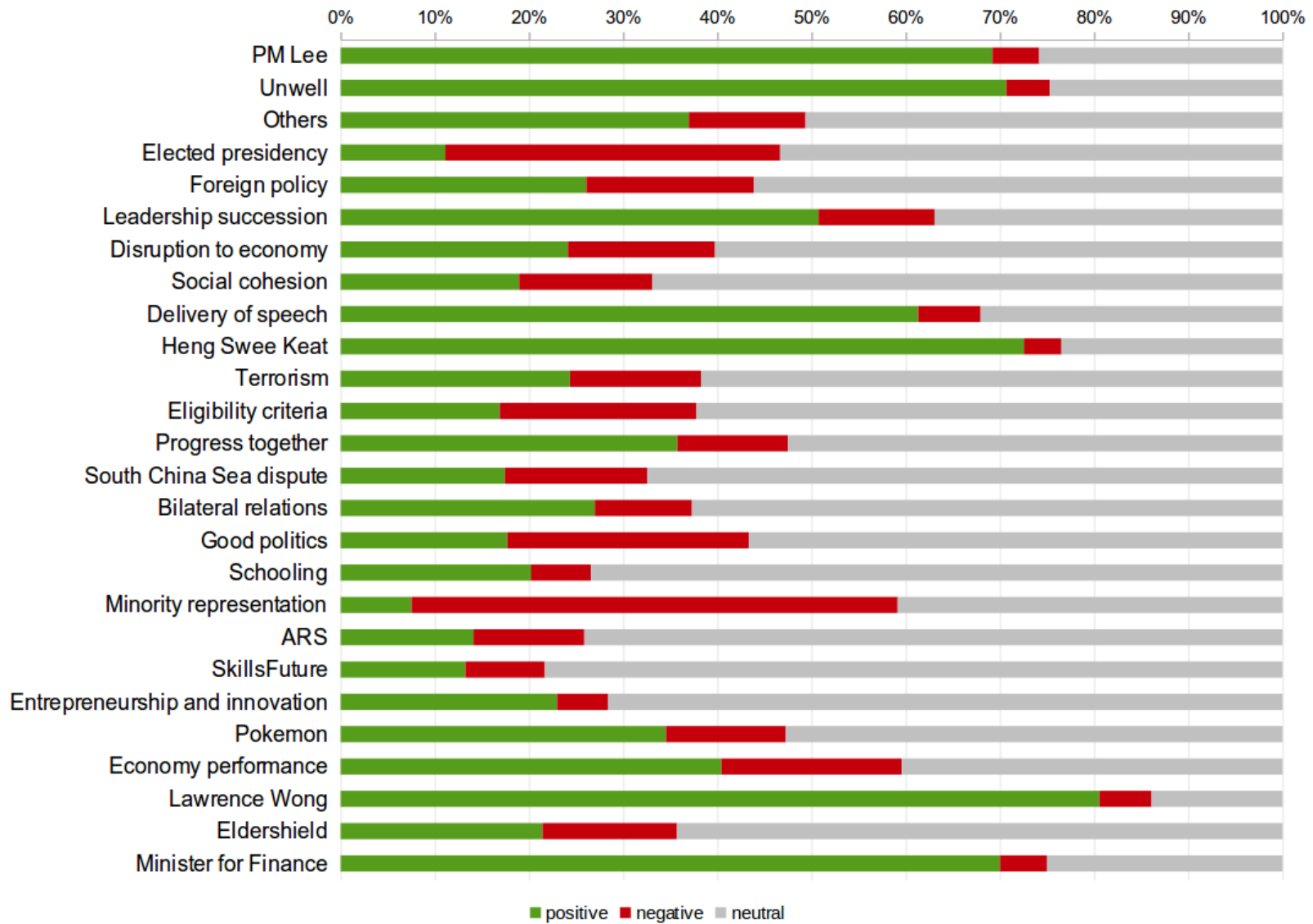| Topic | Value |
|---|---|
| PM Lee | 6390 |
| Unwell | 5543 |
| Others | 1891 |
| Elected presidency | 1686 |
| Foreign policy | 670 |
| Leadership succession | 539 |
| Disruption to economy | 496 |
| Social cohesion | 440 |
| Delivery of speech | 362 |
| Heng Swee Keat | 279 |
| Terrorism | 232 |
| Eligibility criteria | 226 |
| Progress together | 179 |
| South China Sea dispute | 172 |
| Bilateral relations | 126 |
| Good politics | 113 |
| Schooling | 95 |
| Minority representation | 93 |
| ARS | 85 |
| SkillsFuture | 84 |
| Entrepreneurship and innovation | 75 |
| Pokemon | 55 |
| Economy performance | 47 |
| Lawrence Wong | 36 |
| Eldershield | 28 |
| Minister for Finance | 20 |

# Content Volume By Survey Topics

What were the online **sentiments** toward NDR topics?

# Proportions of Sentiments By Topics



Legend: positive (green), negative (red), neutral (grey)

Topics (top to bottom): PM Lee, Unwell, Others, Elected presidency, Foreign policy, Leadership succession, Disruption to economy, Social cohesion, Delivery of speech, Heng Swee Keat, Terrorism, Eligibility criteria, Progress together, South China Sea dispute, Bilateral relations, Good politics, Schooling, Minority representation, ARS, SkillsFuture, Entrepreneurship and innovation, Pokemon, Economy performance, Lawrence Wong, Eldershield, Minister for Finance
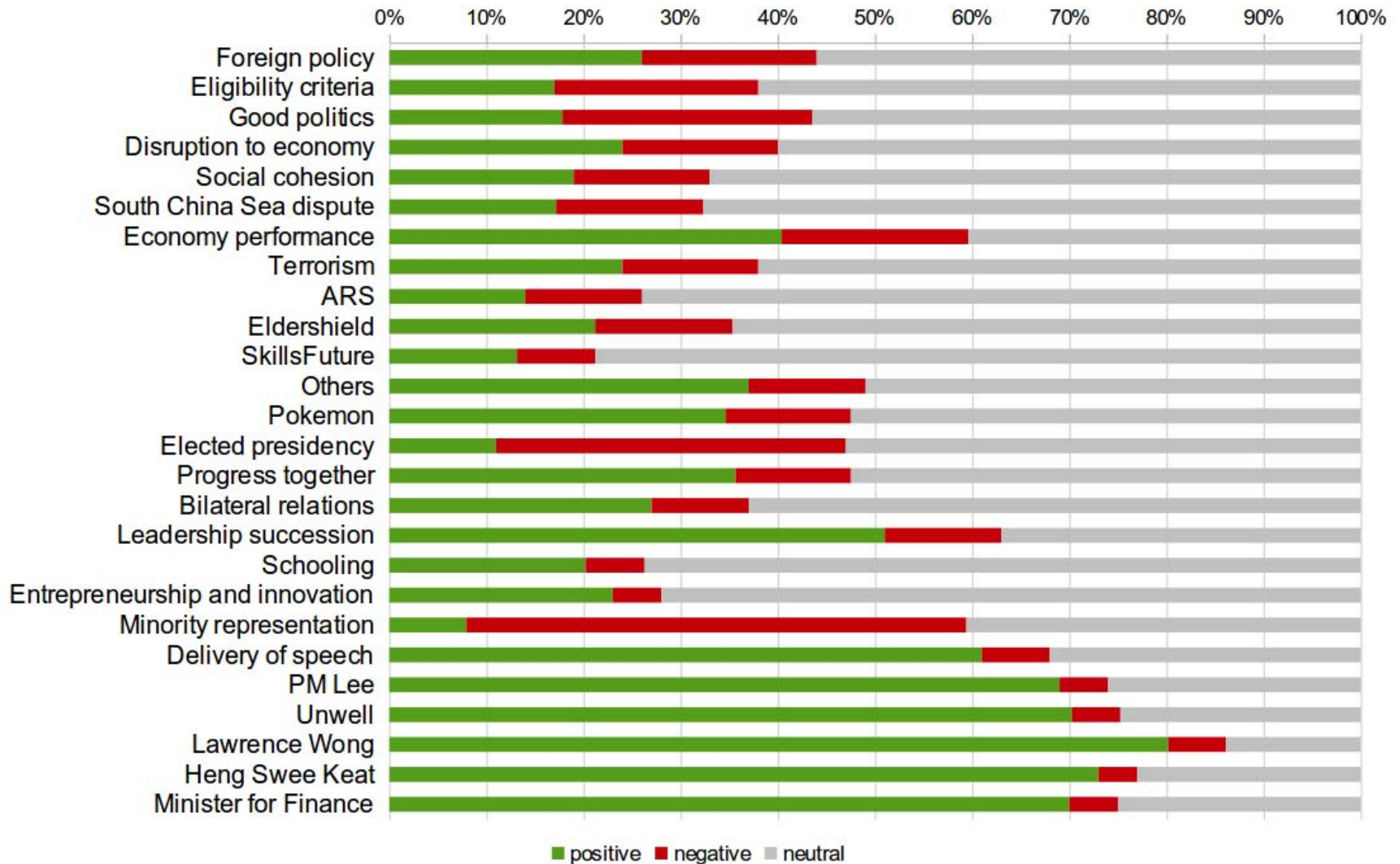
# Proportions of Sentiments By Survey Topics

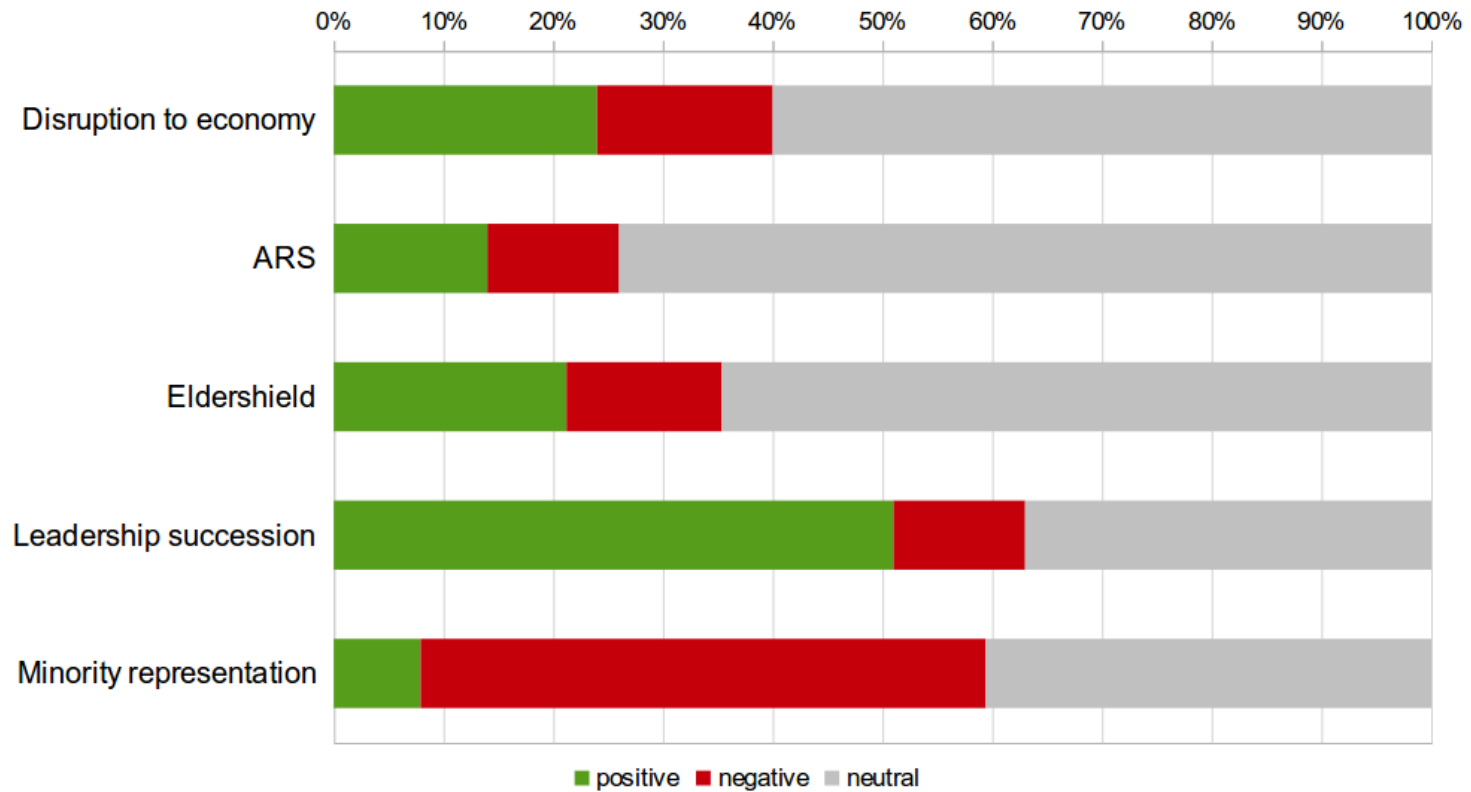Which NDR topics gained the most **divided sentiments**?

**Divisiveness assumptions:**

- High volume of subjective content (positive + negative)
- Equal portions of positive and negative sentiments

# Most Divided Topics



Horizontal stacked bar chart showing positive, negative, and neutral proportions (0%–100%) for topics: Foreign policy, Eligibility criteria, Good politics, Disruption to economy, Social cohesion, South China Sea dispute, Economy performance, Terrorism, ARS, Eldershield, SkillsFuture, Others, Pokemon, Elected presidency, Progress together, Bilateral relations, Leadership succession, Schooling, Entrepreneurship and innovation, Minority representation, Delivery of speech, PM Lee, Unwell, Lawrence Wong, Heng Swee Keat, Minister for Finance.

Legend: ■ positive ■ negative ■ neutral

# Most Divided Survey Topics

How do **topics** differ across online platforms?

# Distance Between Platforms

|  | Topic 1 | Topic 2 | Topic 3 | Topic 4 | Topic 5 |
|---|---|---|---|---|---|
| Platform A | 40% | 30% | 10% | 10% | 10% |
| Platform B | 40% | 30% | 10% | 10% | 10% |

distance(A, B) = 0

|  | Topic 1 | Topic 2 | Topic 3 | Topic 4 | Topic 5 |
|---|---|---|---|---|---|
| Platform A | 40% | 30% | 10% | 10% | 10% |
| Platform C | 20% | 20% | 20% | 20% | 20% |

distance(A, C) = 0.05

|  | Topic 1 | Topic 2 | Topic 3 | Topic 4 | Topic 5 |
|---|---|---|---|---|---|
| Platform A | 40% | 30% | 10% | 10% | 10% |
| Platform D | 0% | 0% | 0% | 0% | 100% |

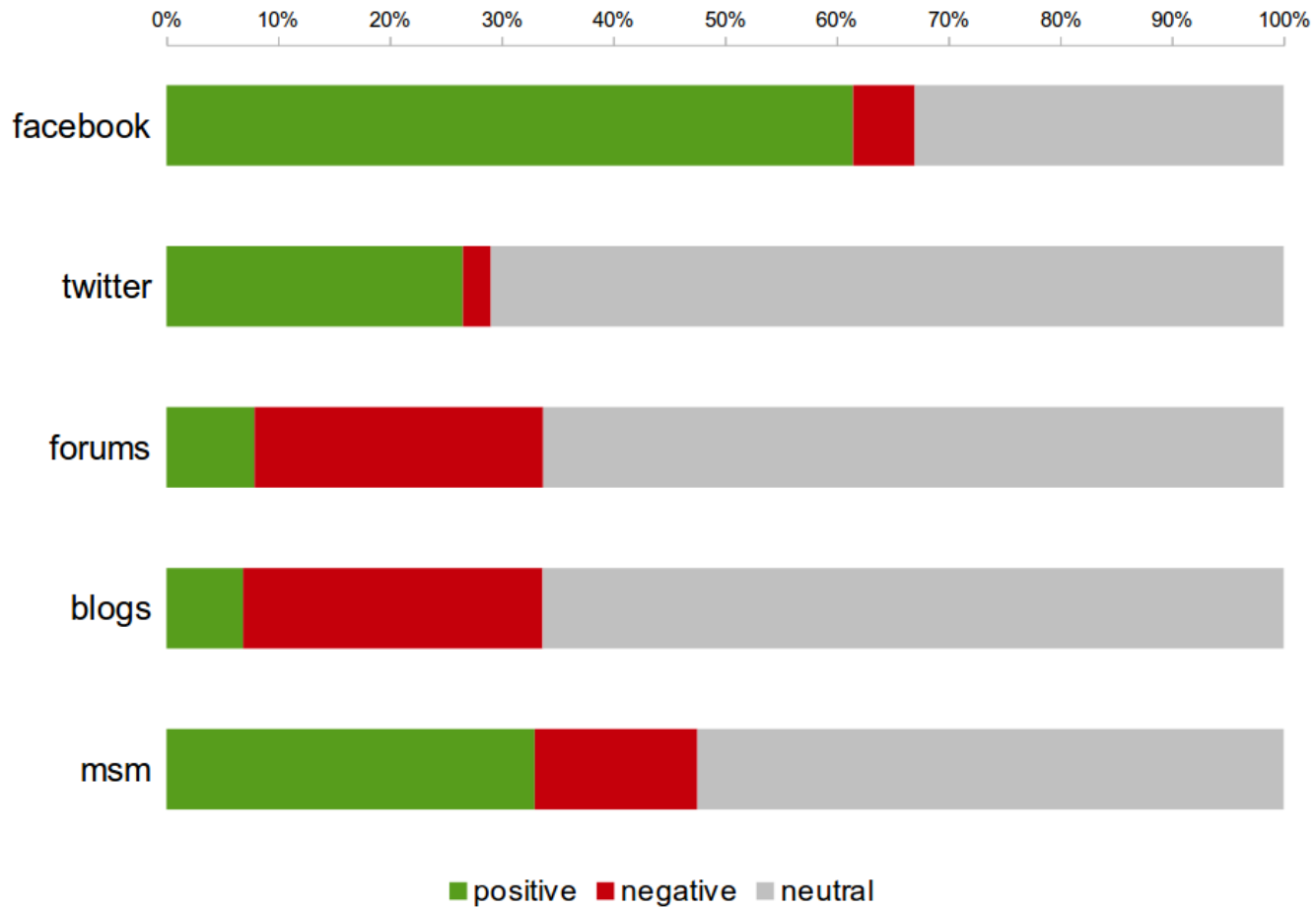distance(A, D) = 0.53

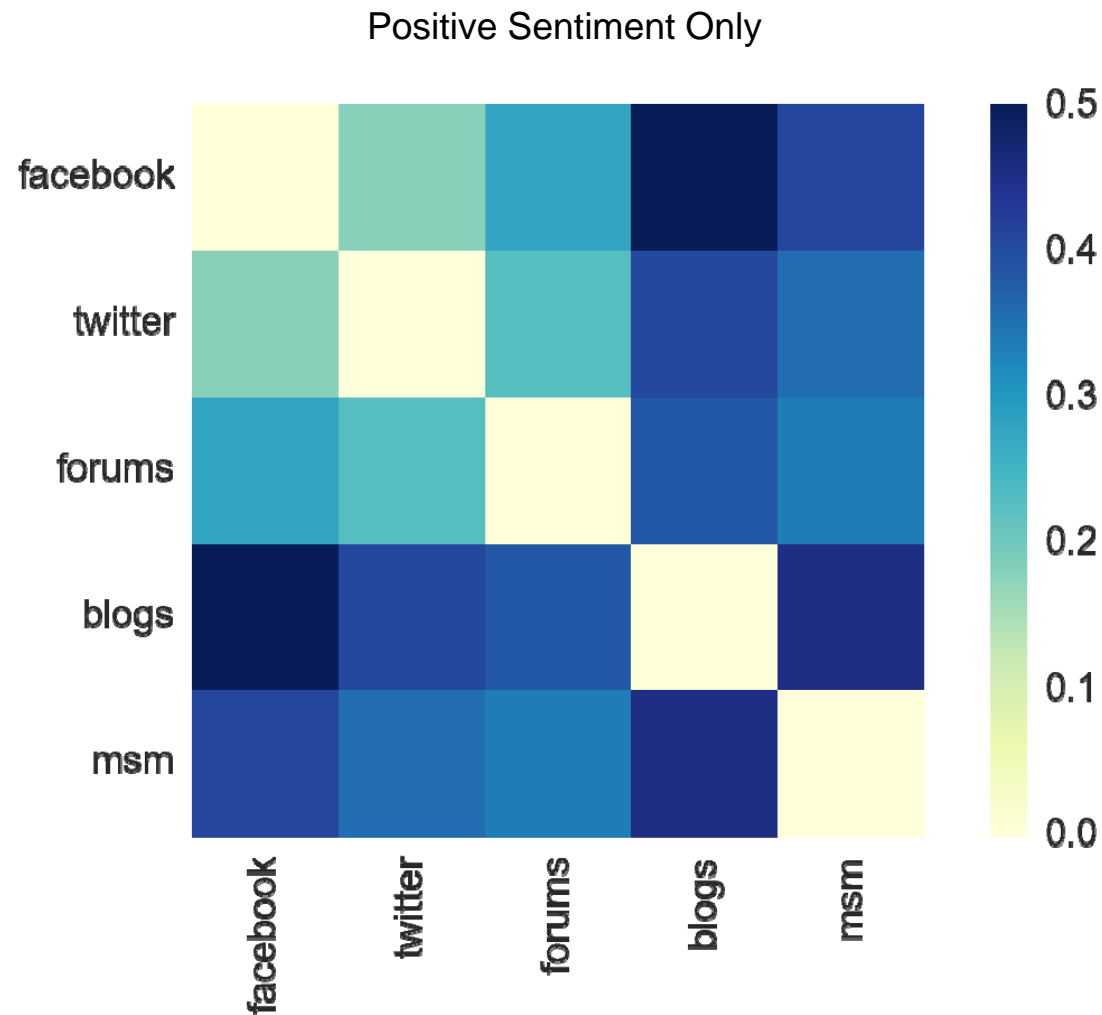# Topic Distance Between Platforms

# Top-5 Topics By Platforms

How do **sentiments** differ across platforms?
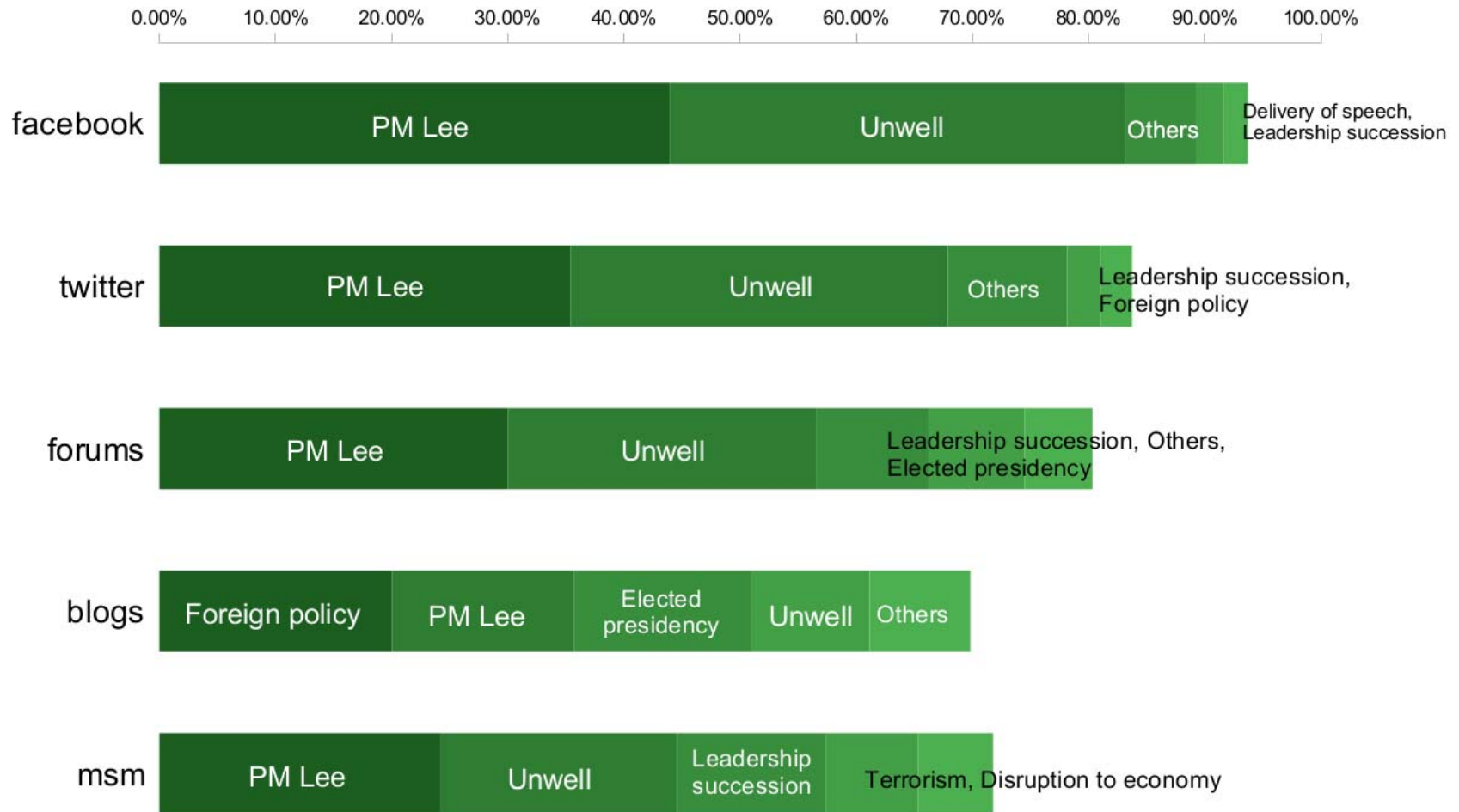
# Proportions of Sentiments By Platforms

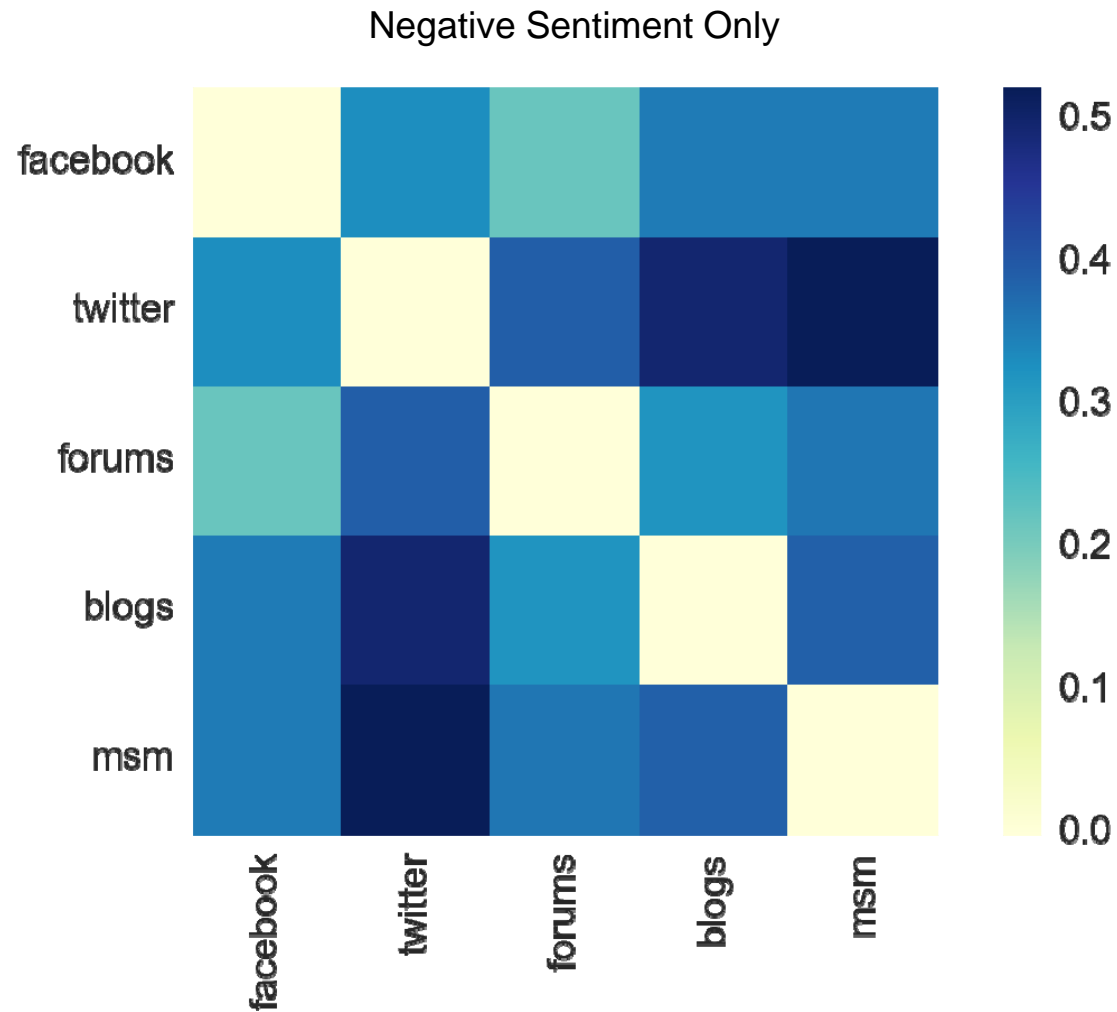How do **sentiments** toward various **topics** differ
across online platforms?

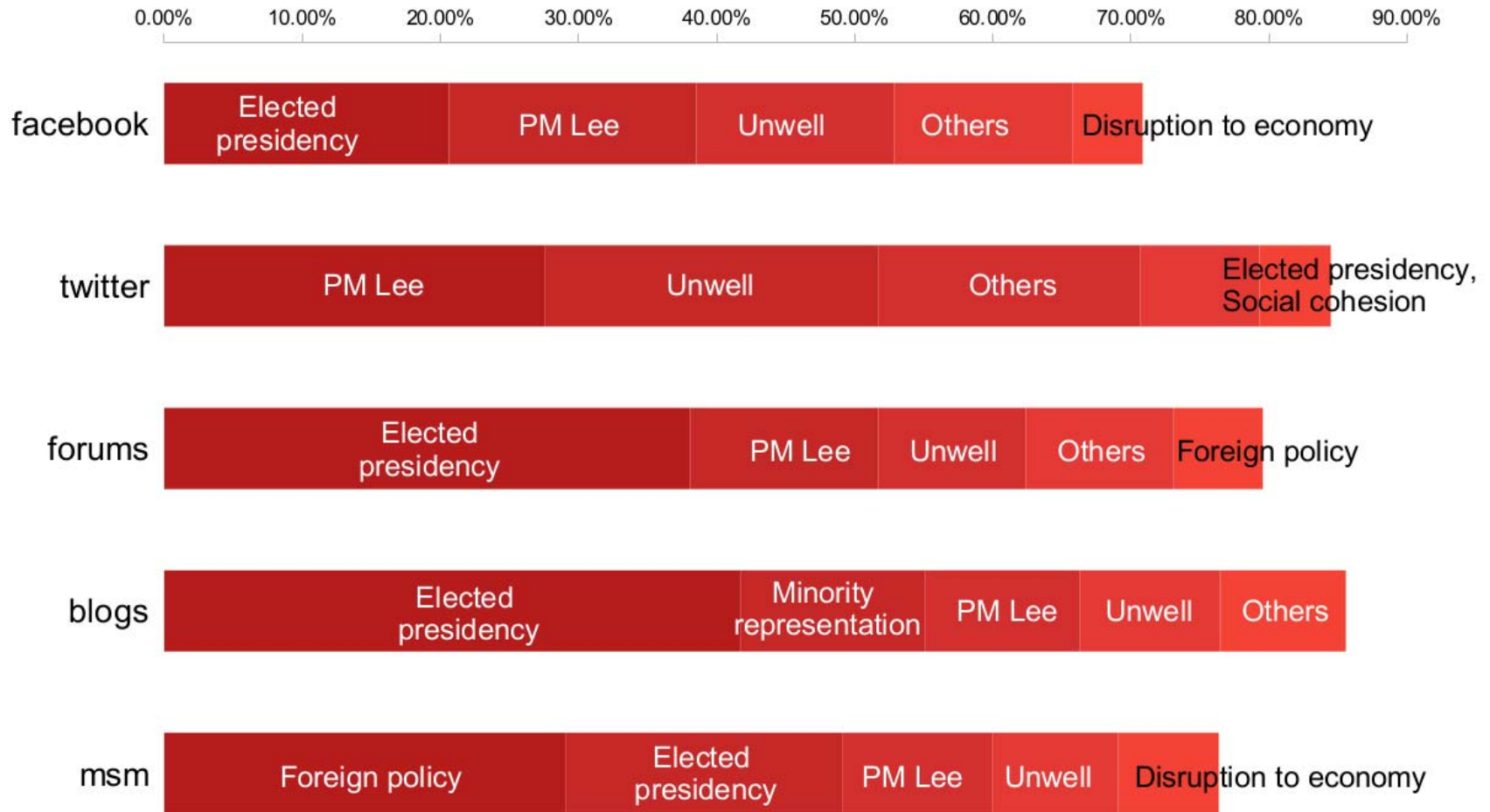# Topic-Sentiment Distance Between Platforms



Positive Sentiment Only

# Top-5 Positive Topics By Platforms



| | 0.00% | 10.00% | 20.00% | 30.00% | 40.00% | 50.00% | 60.00% | 70.00% | 80.00% | 90.00% | 100.00% |
|---|---|---|---|---|---|---|---|---|---|---|---|

**facebook** — PM Lee | Unwell | Others | Delivery of speech, Leadership succession

**twitter** — PM Lee | Unwell | Others | Leadership succession, Foreign policy

**forums** — PM Lee | Unwell | Leadership succession, Others, Elected presidency

**blogs** — Foreign policy | PM Lee | Elected presidency | Unwell | Others

**msm** — PM Lee | Unwell | Leadership succession | Terrorism, Disruption to economy

# Topic-Sentiment Distance Between Platforms



Negative Sentiment Only

# Top-5 Negative Topics By Platforms

# Key Findings

- Facebook has more online users and more content generated than other platforms.
- The topic distributions on Facebook and Twitter are clearly different from those on blogs, forums.
- Sentiments are more positive on Facebook, Twitter and mainstream media than on blogs and forums.
- Online discussions (especially on Facebook and Twitter) are dominated by the topic on PM's health.
  o The event may have shifted the public's attention from the key NDR topics.
  o The event may have stimulated more interest and discussions on NDR.

# THANK YOU!