

# Digital Threats to Democracy

# Digital Threats to Democracy



# Social media's relationship to elections and democratic deepening in emerging democracies



# **Analysis and action research during elections in Africa, Asia, and Latin America**



Increasing evidence suggests identity politics, misinformation, and hate speech is **dominant globally**



# Aggie & Case Study: Nigeria, Ghana, and Kenya



Human-in-the-loop semi-real-time tracking platform with direct facilities for escalation and response

Sources



Reports



Incidents



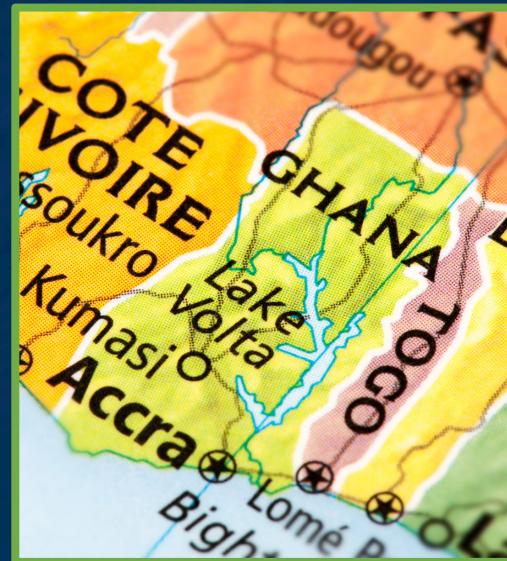
Aggie: Social Media Tracking Platform

## Nigeria



2011

## Ghana



2012

## Kenya



2013

**Country Comparisons**

Country	Election Date	Data Acquisition Time Span	Twitter Reports
Nigeria	16 April 2011	15 April 7:00 AM – 19 April 7:00 AM	192,142
Ghana	7 Dec 2012	6 Dec 7:00 AM – 10 Dec 7:00 AM	320,789
Kenya	4 March 2013	3 March 7:00 AM – 7 March 7:00 AM	254,216
Total			767,147



## Twitter: Policy vs. Identity Politics



**Policy queries:**  
most mentioned bi-grams  
from the top two contending  
political party's electoral  
campaign manifestos



**Identity queries:**  
terms associated with the  
relevant tribal or ethnic  
groups, religion, and  
geopolitical region specific

**Twitter: Policy vs. Identity Politics**

Country	Projected Number of Tweets Relevant to Policy Terms	Policy Percentage of Dataset	Projected Number of Tweets Relevant to Identity Terms	Identity Percentage of Dataset
Nigeria	4,014	2.09%	8,323	4.33%
Ghana	41,466	12.93%	6,104	1.90%
Kenya	5,929	2.33%	4,372	1.72%



## Twitter: Policy vs. Identity Politics

"5 years on, how will we measure gov effectiveness? Improved economic integration with our neighbors... #choice2013  
#kenyadecides"



**2,822 tweets** about economic growth in Kenya

**4,327 tweets** about identity in Kenya



"I'll never ever trust kikuyu hata... [REDACTED] you kikuyuz [REDACTED] you"

Example Tweets from Kenya



**Virtual space of Twitter  
is a reflection of and  
an actor in the  
democratic process**



**Increasing evidence that  
identity and false  
information flows with  
little friction**



In 2011-2012,  
**Nigeria had 2  
identity tweets for  
every 1 policy tweet;**  
**Ghana had 7 policy  
tweets for every 1  
identity tweet.**

**Where are We Headed?**

# **ML / NLP & Case Study: Myanmar**

# Challenges in Emerging & At-Risk Democracies:



Limited text tools  
and data



Diverse language  
and culture



Range of preferred  
platforms



ML with small data &  
under-resourced languages

# Machine Learning / Natural Language Processing for Hate, Disinformation, and Electoral Irregularities

A screenshot of the Aggie software interface. The top navigation bar includes 'AGGIE Reports', 'Incidents', 'Sources', and 'Analysis'. Below the navigation is a search bar with fields for 'Enter Keywords', 'Enter Tags', 'Enter Author', 'Status', 'Media', 'Source', 'Linked Incident', 'DateTime', and a 'Search' button. A sidebar on the left lists 'Tags' such as 'Anti-Buddhist', 'Anti-Muslim', 'Covid-19', 'Hate Speech', and 'Violence'. The main area displays a grid of reports. Each report card includes columns for 'Time', 'Sources', 'Thumbnail', 'Author', 'Content', 'Tags', and 'Incident/Flagged'. The content section shows snippets of text in various languages, often with flagged words like 'Hate Speech' or 'Violence'. The bottom right corner of the interface shows summary statistics: 'Total Reports 1,297', 'Reported last min 0', 'Unread 1,293', 'Flagged 10', 'Incidents 5', and 'Escalated Incidents 0'.



## Case Study: Myanmar Election November 2020

- ◆ Collaborated with **civil society partners**, The Carter Center and New Myanmar Foundation
- ◆ Deployed **distributed social media tracking** and response center
- ◆ **Military coup** in February 2021





Experienced **significant harm** from hate speech and disinformation over **social media**



Ethnic, linguistic, and religious **minorities** have been **persecuted**, including **genocide**

## Social Media in Myanmar



Model	Precision		Recall	
	Hate	Not Hate	Hate	Not Hate
Balanced Random Forest	0.08	0.99	0.82	0.58
Balanced Random Forest (Oversampled minority class)	0.84	0.97	0.98	0.82
FastText Multilingual Embeddings	0.71	0.97	0.27	1.00

Daniel Nkemelu and Harshil Shah

## Hate Speech Classifier: Most Recent Results

## Edit Report Tags



## Add and Remove Tags

Hate Speech ×

Anti-Muslim ×

Violence x Type

Type to Enter Tags

## Clear Tags

Cancel

Submit

# Trackers Tag Hate Speech, Disinformation, & Irregularities



**“Low-resourced” languages, limited training datasets, limited text tools, limited language models**



**Working across languages**



**Coordination and communication challenges**