

# THE LOGIC OF AI & DISINFORMATION ON SOCIAL MEDIA PLATFORMS

Dr. Courtney C. Radsch

SENIOR FELLOW, CENTER FOR INTERNATIONAL GOVERNANCE INNOVATION

SENIOR TECHNOLOGY POLICY ADVISOR, GLOBAL FORUM FOR MEDIA DEVELOPMENT

SENIOR US AND TECH POLICY ADVISOR, ARTICLE19

FELLOW, CENTER FOR MEDIA, DATA & SOCIETY

VISITING SCHOLAR, CENTER FOR MEDIA AT RISK, ANNENBERG, UNIV. OF PENNSYLVANIA

ADVISORY COUNCIL, RANKING DIGITAL RIGHTS

ADVISORY BOARD, DANGEROUS SPEECH PROJECT

# SNAPSHOT OF AI & DISINFO ON SOCIAL MEDIA

MANIPULATION, DIVISIVENESS,  
INAUTHENTICITY

IDENTITY GROUPS, SOCIETAL  
TENSIONS/DIVISIONS, MISOGYNY,  
RACISM & ETHNIC/RELIGIOUS HATRED

GENDERED DISINFORMATION

RISE OF ALT- AND FAR-RIGHT POPULISM

OVERREPRESENTATION OF FAR-RIGHT  
PARTIES IN EUROPE, US

INFORMATION WARFARE & FOREIGN  
MANIPULATION TO INFLUENCE DOMESTIC  
PUBLIC OPINION

FACEBOOK CONTENT REMOVALS USING  
DETECTION WITH AI TOOLS:

- 99.5% OF TERRORIST-RELATED
- 98.5 % OF FAKE ACCOUNTS
- 96% OF ADULT NUDITY & SEXUAL  
ACTIVITY
- 86% OF GRAPHIC VIOLENCE-RELATED

GLOBAL INTERNET FORUM TO COUNTER  
TERRORISM (GIFCT) HASH DATABASE

DIGITAL MILLENNIUM COPYRIGHT ACT  
(DMCA)

## COORDINATED DISINFORMATION

**70+ STATES/GOVTS** ENGAGE IN INFORMATION  
OPERATIONS/COMPUTATIONAL PROPAGANDA

2017-2020: FB IDENTIFIED 150+ COORDINATED  
INAUTHENTIC BEHAVIOR CAMPAIGNS IN MORE  
THAN **50 COUNTRIES**

2020-2021: TWITTER REMOVES TENS OF  
THOUSANDS OF **STATE-ALIGNED INFO OPS**

COMMERCIALIZATION + INDUSTRIALIZATION OF  
INFORMATION & INFLUENCE OPERATIONS

# AI-FUELED ATTACKS ON PUBLIC INTEREST: ELECTIONS, PUBLIC HEALTH, SAFETY & SECURITY

## INDIVIDUAL AND COLLECTIVE HARMS

- DISCRIMINATION, EXPLOITATION, INTIMIDATION, SELF-CENSORSHIP, SABOTAGE, RECRUITMENT, EXTREMISM, BREAKDOWN OF SOCIAL TIES
- CURTAIL DEMOCRATIC DELIBERATION, REPRESENTATION
- TRUST DEFICIT & "LIAR'S DIVIDEND"
- SELF-CENSORSHIP AND ALTERED BEHAVIOR RESULTING FROM PERVASIVE SURVEILLANCE → ↓ DIVERSITY → ↓ PLURALISM
- MODERATION & CENSORSHIP OF AUTHENTIC/LEGITIMATE INFO
- ↑ POLARIZATION → ↓ PLURALISM (LACK OF CENTER)

## TARGETS

- VULNERABLE POPULATIONS
- JOURNALISTS; FEMALE POLITICIANS; PUBLIC FIGURES
- GENDERED DISINFORMATION

## ELECTIONS, PARTISAN POLITICS, PROTEST & DEMOCRACY MOVEMENTS

- **US:** 75% OF USERS EXPOSED TO CONTENT FROM CLICKBAIT FARMS IN **MACEDONIA & KOSOVO** HAD NEVER FOLLOWED THE PAGE
- **GERMANY:** AfD HAD 5x FB ENGAGEMENT AS OTHER PARTIES; 63% OF WOMEN DON'T EXPRESS POLITICAL VIEWS ONLINE

## VIOLENT CONFLICT

- UKRAINE, ETHIOPIA, MYANMAR
- ETHNIC, RACIAL, RELIGIOUS DIVISIONS

## COVID19, ANTI-VAXXER MOVEMENT

- 12 SUPER-SPREADERS OF DISINFORMATION
- TRUST & INTEGRITY
- FUTILITY OF COMO AT SCALE

# IMPACTS OF AI-FUELED DISINFORMATION ON JOURNALISM

- TARGETING JOURNALISTS / MEDIA = DENYING THEM **AGENDA-SETTING POWER**
- **FRAMING** JOURNALISTS & PROFESSION AS "FAKE NEWS"
- **DECERTIFYING** JOURNALISTS, PARTICULARLY WOMEN, AS LEGITIMATE ACTORS IN THE PUBLIC SPHERE
- **GOAL = UNDERMINE VIABILITY & SUSTAINABILITY OF INDEPENDENT JOURNALISM**
  - REINFORCE DISTRUST IN THE MEDIA, INOCULATE THOSE IN POWER FROM OVERSIGHT + ACCOUNTABILITY
  - DELEGITIMIZE REPORTING ON CERTAIN TOPICS & BY "TYPES" OF PEOPLE
    - CONTROVERSIAL ISSUES, SOCIAL MEDIA MANIPULATION, INFORMATION OPERATIONS
    - GENDER, MINORITIES & INTERSECTIONALITY
- FORCING LOGIC OF PLATFORM AI CHOICES
  - NEWSROOM RESOURCES & PRIORITIES
    - FB VIDEO EMPHASIS → NEWS ORGS PIVOT
  - ONLINE HARASSMENT & GENDER-BASED VIOLENCE

Ukraine: Avg Troll Farm Salary = Avg FT employee (\$365/mo)

Russia: Budget of Internet Research Agency est. \$400,000/mo

RT online 24hr/da, >3 languages; reaches 700 million people in 100 countries

# THE LOGIC OF AI IN THE INFORMATION ECOSYSTEM

- AI IS SHAPED BY A CORPORATE LOGIC
  - PLATFORMS & PROFIT, WHOEVER CAN PAY + GROWTH MAXIMIZATION
  - WHERE RESOURCES ARE DEVOTED: COUNTRIES, LANGUAGES, ISSUES, ETC.
  - RESEARCH PRIORITIES: PRIVATE VS PUBLIC INTEREST
  - HARM IDENTIFICATION, PRIORITIZATION, REDUCTION
- ALGORITHMIC INTERMEDIATION
  - DESIGNED BASED ON THIS LOGIC
  - PERSONALIZATION
  - ONLY “SEE” WHAT YOU LOOK FOR OR DISCOVER
- AGENDA-SETTING POWER
- FRAMING
- FILLING INFORMATION VOIDS

# LIMITED PLURALISM & NEW GATEKEEPERS

A FEW PLATFORMS DETERMINE THE **LOGIC** OF THE PUBLIC SPHERE & FREEDOM OF EXPRESSION

- DOMINANCE OF A FEW SOCIAL MEDIA INTERMEDIARIES = LIMITED PLURALISM; THEY DECIDE HOW INFORMATION GETS TO BE SEEN AND SHARED – LACK OF ALTERNATIVE LOGICS
- SUSCEPTIBILITY TO INFLUENCE BY A HANDFUL OF INFLUENTIAL GOVERNMENTS & THEIR PRIORITIES
- PRONE TO MANIPULATION: DRIVEN BY POLITICAL, ECONOMIC, SOCIAL FACTORS

GATEKEEPING: OF EXPRESSION [CHANNELING, CENSORING, ADDING VALUE, INFRASTRUCTURE, USER INTERACTION, EDITORIAL; RELEVANCE, RECOMMENDING, TRENDING TOPICS] + SCIENTIFIC RESEARCH & INNOVATION

- CONTENT MODERATION – CHALLENGING AND PRONE TO ERROR & ABUSE AT SCALE
- SHIFTS IN WHO THE GATEKEEPERS ARE AND WHAT THEIR LOGIC OF GATEKEEPING IS
  - JOURNALISM NORMS -> PLATFORM NORMS = ADTECH

# THE POLITICAL-ECONOMIC LOGIC OF AI & DISINFO

MONETIZATION OPPORTUNITIES+ AI → RISE OF TROLL FARMS, SPAMMERS, CLICKBAIT FARMS, PLAGIARISM, MICRO-PRENEURS

- WITHOUT INTEGRITY EVALUATION, MONETIZATION ATTRACTS & REWARDS LOW-QUALITY SOURCES
- DROWN OUT ORGANIC, LEGITIMATE NEWS, OPINION, AND ENGAGEMENT
- EASE OF ACCOUNT CREATION, CONTENT POSTING & REGURGITATION ACROSS PLATFORMS
- CERTIFICATION AND EQUIVALENCY
- DEEP & SHALLOW FAKES:
  - MACHINE LEARNING TECHNIQUES → ↑ SOPHISTICATION = MORE REALISTIC, RESISTANT TO DETECTION
  - CHEAPER, EASIER TO CREATE → MORE UBIQUITOUS

## CASE STUDY: PHILIPPINES

- 2016: DUTERTE SUBJECT OF 68% OF ALL ELECTION-RELATED DISCUSSIONS (46% FOR CLOSEST RIVAL)
- MARIA RESSA @ RAPPLER REPORTS ON INFO OPS, FB MANIPULATION → HARASSMENT & LEGAL PERIL
  - RECEIVES 2021 NOBEL PRIZE
- FIRM HE HIRED MADE \$8 MILLION FROM FB & GOOGLE BEFORE BEING TAKEN DOWN IN 2019

# ZOOM IN: ECONOMIC LOGIC OF ADTECH & DISINFO

- PRE-2015: FB CLICK → PUBLISHERS (ADS SERVED BY GOOGLE) → FB INSTANT ARTICLES (FB SERVES ADS)
- **2015**: INSTANT ARTICLES IN US & EUROPE
- **2016**: INSTANT ARTICLES IN GLOBAL SOUTH
- **2018**: \$1.5 B PAID OUT
- **2019**: FB BEGINS CHECKING PUBLISHERS FOR CONTENT ORIGINALITY & DEMONETIZING
- **2021**: ADS EMBEDDED IN LIVE VIDEOS

## CASE STUDY: MYANMAR

- **2015**: 6 OF TOP 10 WEBSITES W/MOST FB ENGAGEMENT = LEGITIMATE PUBLISHERS
- **2016**: INSTANT ARTICLES ROLLED OUT
- **2017**: 2 OF 10 TOP PUBLISHERS = LEGITIMATE
- **2017 MILITARY CRACKDOWN & ANTI-MUSLIM PROPAGANDA**
- **2018**: 0 OF TOP 10 PUBLISHERS = LEGITIMATE
- **2018 UN DETERMINES VIOLENCE = GENOCIDE**
- **2021 MILITARY COUP & INFORMATION WARFARE, SOCIAL MEDIA BANS INEFFECTIVE**
- **2021**: YOUTUBE CHANNELS CONVERTED TO FB ARTICLES & REDISTRIBUTED; SOME POSTED “LIVE”

# LEVERAGING LOGIC OF PUBLISHING & PLATFORMS

## CORPORATE/PROFIT MOTIVE, NETWORK EFFECTS, DIFFUSION, VIRALITY, ENTRENCHMENT

- 1:1 + FEW:MANY + MANY:MANY
- PERVASIVE PRIVATIZED SURVEILLANCE + DATAFICATION → ITERATIVE COLLECTION AND PROCESSING OF MASSIVE AMOUNTS OF DATA ABOUT PEOPLE → BIG DATA, CLASSIFICATION, MACHINE LEARNING, NEURAL NETWORKS → PREDICTION, PROMOTION, DISCRIMINATION
- RAPID DIFFUSION WITH LIMITED CONSTRAINTS VS. SLOW, OFTEN INEFFECTIVE MITIGATION
- VIRALITY → “INFORMATION CASCADES”
- FILTER BUBBLES → CONFIRMATION BIAS
- ECHO CHAMBERS → POLARIZATION, ↓ DIVERSITY
- PREDICTIVE ENGAGEMENT ALGORITHMS BASED ON LOOK-ALIKE AUDIENCES → SPREAD OF DISINFORMATION + EXTREMISM

# AI ENABLED DISINFORMATION: IMPACT ON THE INFORMATION ECOSYSTEM

- INPUTS: DATA (SOURCES/WHAT COUNTS AS DATA/LANGUAGE), DIVERSITY & PLURALISM, FAIRNESS, BIAS
- IDENTIFICATION, TARGETING, SURVEILLANCE
- SCOPE, SCALE & SPEED + SOPHISTICATION, NUANCE
- NEURAL NETWORKS + MACHINE LEARNING + SOCIAL MEDIA NETWORK EFFECTS
  - “GENERATIVE ADVERSARIAL NETWORKS”
- ALGORITHMS IDENTIFY, MAKE VISIBLE, AMPLIFY CONTENT/ACCOUNTS + PREDICT ENGAGEMENT WITH CONTENT + TARGETED ADVERTISING & PROMOTION
- ALGORITHMIC IMPACTS: BOTS + HUMAN ATTRACTION TO NEGATIVE/NOVEL INFORMATION
  - PEOPLE & THEIR ENGAGEMENT BEHAVIOR
  - BOTS: EXPOSURE + CAN MANIPULATE PREDICTIVE CONTENT ENGAGEMENT ALGORITHMS
  - SKEWING PREVALENCE, DROWNING OUT, FILLING INFORMATION VOIDS

THANK YOU

DR. COURTNEY C. RADSCH

CRADSCH@GMAIL.COM

[WWW.MEDIATEDSPEECH.COM](http://WWW.MEDIATEDSPEECH.COM)

@COURTNEYR