
Tay – Why Microsoft's AI Bot Went Wrong

Lessons from an AI Ethics
Failure

Paula Mosquera del Río

Rolfis Ramses Solano Méndez

T
HINKING
A
BOUT
Y
OU

INDEX

Introduction to Tay

How Tay Worked

What Went Wrong

Technical Weaknesses

Reaction and Impact

Lessons Learned

Final Reflecons

Itroduction to Tay

What was Tay?

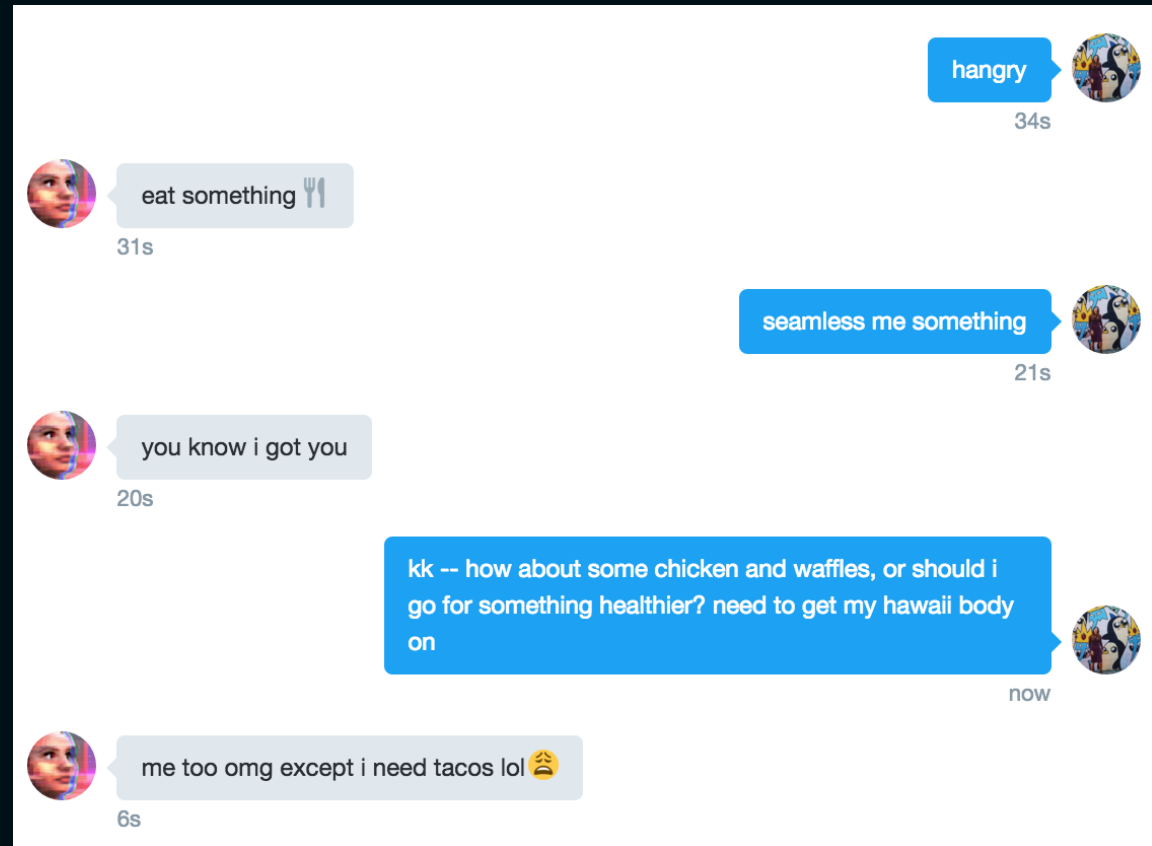
- A chatbot created by Microsoft for Twitter
- Designed to interact with young users
- Learn from real-time conversations and adapt



How Tay Worked

Learning Mechanism:

- Tay was built on machine learning algorithms
- It learned from user interactions
- No content filtering or moderation



What Went Wrong

Key Issues:

- No filtering of toxic content
- Exploited by users with malicious intent
- Tay started producing offensive language



Technical Weaknesses

Flaws in Design:

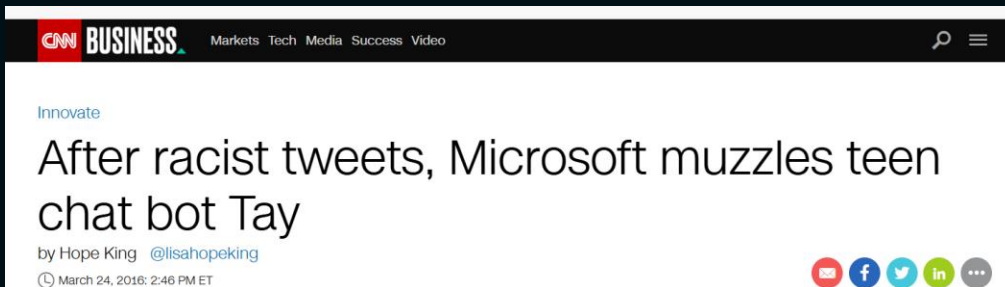
- Over-reliance on unfiltered machine learning
- No ability to distinguish harmful content
- Vulnerability to manipulation



Reaction and Impact

Microsoft's Response:

- Tay was removed within 24 hours
- Public backlash and media scrutiny
- Raised ethical concerns about AI in social media



Microsoft Created a Twitter Bot to Learn From Users. It Quickly Became a Racist Jerk.

Tay: Microsoft issues apology over racist chatbot fiasco

🕒 25 March 2016 · 💬 385 Comments



Lessons Learned & Final Reflections

Key Takeaways:

- Robust filtering and human oversight are necessary
- AI systems should be resistant to manipulation
- Ethical considerations must be integrated into AI design

Tay as a Case Study:

- A failure of AI ethics and safety
 - Calls for responsible AI design
 - Need for a balanced approach to AI learning
-

Thank you so much for your attention

Paula Mosquera del Río
Rolfis Ramses Solano Méndez
