AI vs Bots on Xiaohongshu: Modeling the Arms Race in Ecommerce Comment Ecosystems

Other possible catchy titles:

- 1. Simulating Social Bots: The Role of AI in Xiaohongshu's Comment Ecosystem
- 2. Manufacturing Trust: Coordinated Bot Interactions on Social Media
- 3. Bot-Driven Credibility: Simulating Fake Endorsements on Xiaohongshu
- 4. Bot vs Guardian: Simulating Content Moderation Arms Race on Social Platforms

Team members: Xintong (Sylvia) LING, Huanrui (Saikoro) CAO, Jiayi (Jaye) CHEN GitHub repo URL

Section 1: Phenomena of interest.

Instructions

Our project focuses on the phenomenon of bot-generated comments on Xiaohongshu (RedNote), a Chinese social media and e-commerce platform. Automated bots on Xiaohongshu (RedNote) frequently engage in comment sections, particularly under posts that mention keywords related to shipping, logistics, and specific product categories. These bots operate in coordinated ways to promote advertisements, counterfeit products, and services such as fortune-telling.

Key Dynamics

- 1. Keyword-Based Information Extraction Bots scan user-generated posts and comments for specific keywords (e.g., "shipping," "logistics") to identify relevant targets for promotional activity.
- 2. Automated Advertising Bots generate promotional comments to advertise products or services, often embedding links or contact details.
- 3. Coordinated Interaction for Credibility Some bots pose as genuine users to endorse scam products, engaging in scripted conversations. For example, one bot might claim to have purchased a counterfeit luxury item and praise its authenticity, while another bot follows up with a question asking where to buy it, creating an illusion of organic discussion.
- 4. Impact on User Experience and Trust These automated interactions manipulate user perception and can degrade trust in the platform, influencing purchasing decisions and the overall reliability of user-generated content.

The cyclical pattern where each moderation system upgrade triggers new evasion strategies, which in turn drive detection algorithm improvements. This phenomenon manifests as an escalating arms race where:

- 1. Adaptive Content Generation Bots employ adversarial neural networks to create promotional content that bypasses detection thresholds, using techniques like:
- Context-aware keyword substitution (e.g., "sh!pping" → "shipping")
- Semantic-preserving sentence restructuring
- AI-generated counterfeit user reviews
- 2. Dynamic Detection Evasion Bot networks implement reinforcement learning to:
- Analyze historical takedown patterns

- Predict moderation system update cycles
- Optimize posting timing and content distribution
- 3. Countermeasure Adaptation Moderation AIs develop:
- Graph-based bot cluster detection
- Behavioral fingerprint analysis
- Cross-modal consistency checks (text-image matching)
- 4. User-AI Symbiosis Human users unconsciously adopt bot-like behaviors:
- Mimicking verified account patterns
- Developing community-specific circumvention jargon
- Participating in crowdsourced detection testing

Why Xiaohongshu?

Xiaohongshu (RedNote) represents a unique hybrid of social media and e-commerce, where user-generated content directly drives purchasing decisions. Unlike Western platforms like Instagram, Xiaohongshu's "grass-planting" () culture emphasizes peer recommendations over brand advertisements, creating fertile ground for bot-driven manipulation. Three key factors motivate our focus:

- 1. **Platform Affordances**: Integrated "post-to-purchase" flow enables bots to directly influence conversion rates.
- 2. Cultural Specificity: High trust in UGC (User-Generated Content) amplifies bot impact compared to review-centric platforms like Amazon.
- 3. **Regulatory Context**: China's evolving internet governance policies create rapid shifts in bot evasion tactics.

Section 2: Relevant Works

References

Himelein-Wachowiak, M., Giorgi, S., Devoto, A., Rahman, M., Ungar, L., Schwartz, H. A., Epstein, D. H., Leggio, L., & Curtis, B. (2021). Bots and misinformation spread on social media: Implications for COVID-19. *Journal of Medical Internet Research*, 23(5). https://doi.org/10.2196/26933

Demonstrates pattern recognition in bot-driven misinformation campaigns

Mena, P., Barbe, D., & Chan-Olmsted, S. (2020). Misinformation on Instagram: The impact of trusted endorsements on message credibility. Social Media + Society, 6(2). https://doi.org/10.1177/2056305120935102

Presents GAN-based approaches for detection-evasion simulation

Zhang, Y., Song, W., Koura, Y. H., & Su, Y. (2023). Social Bots and information propagation in social networks: Simulating cooperative and competitive interaction dynamics. *Systems*, 11(4), 210. https://doi.org/10.3390/systems11040210

Models multi-agent adversarial dynamics in social networks