# The Accuracy and Biases of AI-Based Internet Censorship in China

Yuxin Chen[1]

[1] Wuhan University of Technology, Wuhan, China
Correspondence: Yuxin Chen, Wuhan University of Technology, Wuhan, China.

## Abstract

AI-driven censorship has become a central mechanism for controlling online discourse in China, allowing for rapid detection and suppression of politically sensitive content. This paper explores the accuracy and biases of AI-based internet censorship, focusing on its evolution from manual to automated moderation, its effectiveness in identifying dissent, and its systemic biases that reinforce government narratives. While AI models are highly efficient in filtering explicit political speech, they struggle with disguised dissent, satire, and coded language, leading to inconsistent enforcement and unintended suppression of neutral content. Case studies, including COVID-19 information control, the Hong Kong protests, and labor rights discussions, illustrate both the strengths and limitations of AI censorship. Additionally, a comparative analysis with global AI moderation models highlights key differences between China's state-controlled digital governance and Western approaches to content moderation. The study further examines the challenges of censorship precision, including over-blocking, under-blocking, and the balance between suppressing political dissent and combating misinformation. Finally, the research discusses public and international reactions to China's AI censorship model, evaluating its impact on digital sovereignty, global internet governance, and the future of AI-driven speech regulation. As AI censorship continues to evolve, the paper underscores the ethical dilemmas surrounding algorithmic transparency, state influence, and the exportation of China's digital control model to other countries.

**Keywords:** AI censorship, China, internet regulation, automated content moderation, digital authoritarianism

## 1. Evolution of AI Censorship

The rapid expansion of China's digital ecosystem has necessitated the development of increasingly sophisticated internet censorship mechanisms. Initially, online content moderation in China was largely manual, relying on government-employed censors and platform moderators to monitor and remove politically sensitive material. However, with the explosion of user-generated content (UGC) on social media, news platforms, and video-sharing sites, manual censorship became insufficient to handle the sheer volume of information. This challenge led to the transition from human moderation to AI-driven censorship systems, which now form the backbone of China's internet regulation framework. The evolution of AI-based censorship has been largely driven by

technological advancements in natural language processing (NLP), deep learning, and sentiment analysis, alongside strict government regulations that mandate comprehensive online surveillance.

### 1.1 Shift from Manual to AI-Driven Moderation

Before the widespread adoption of artificial intelligence, China's internet censorship relied heavily on government-appointed human censors who manually reviewed flagged content. Newsrooms, online forums, and social media platforms employed thousands of moderators, often referred to as "wǎngluò jiǎncháyuán" (online inspectors), to monitor discussions and ensure that politically sensitive or controversial content did not spread. During politically significant events such as the Tiananmen Square anniversary, National People's Congress meetings, or leadership transitions, manual censors would work around the clock to suppress content that could be seen as a threat to political stability.

As digital activity skyrocketed—China's internet user base surpassed 1.05 billion by 2023—the inefficiency of purely human moderation became apparent. The sheer volume of daily social media posts, comments, and livestream interactions made it impossible for human censors alone to detect and suppress sensitive content in real-time. To solve this scalability issue, major platforms such as WeChat, Weibo, and Douyin began integrating AI-driven censorship tools, which allowed for automated content scanning, pattern recognition, and predictive filtering. AI systems could flag, blur, or remove content instantaneously, significantly reducing the time needed to enforce censorship policies.

The transition to AI-based censorship marked a turning point in China's internet governance. Unlike human moderators, AI systems are capable of processing millions of posts per second, detecting patterns in text, images, and videos that could indicate politically sensitive discussions. AI also introduced preemptive censorship techniques, where certain discussions were automatically silenced before they could gain traction. This move from reactive censorship (deleting content after posting) to proactive censorship (preventing content from being posted at all) significantly enhanced the state's ability to control public discourse.

### 1.2 Key Regulations Shaping AI Censorship

China's transition to AI-driven internet censorship has been reinforced by a comprehensive legal framework that mandates strict control over online information. Several key regulations have played a pivotal role in shaping the development and implementation of AI-based censorship:

1) The Cybersecurity Law (2017): This law established the legal foundation for internet control, requiring platforms to take responsibility for censoring illegal content, including anything that threatens national security, social order, or government authority. It mandated that all internet companies operating in China must develop advanced content moderation systems, encouraging the adoption of AI-driven censorship tools.

2) The Provisions on the Governance of Online Information Content Ecosystem (2020): This regulation outlined content classification standards, dividing information into categories such as "encouraged," "neutral," and "prohibited." It required AI-driven censorship systems to not only remove illegal content but also to promote government-approved narratives, marking a shift from passive content suppression to active information shaping.

3) The Data Security Law (2021) and Personal Information Protection Law (PIPL): These laws required platforms to store user data within China's borders and implement automated surveillance measures to detect and report potential security risks. AI censorship tools were expanded to analyze user behavior patterns, allowing authorities to track individuals who repeatedly attempted to share restricted content.

4) Algorithm Regulation Guidelines (2022): Recognizing the growing role of AI in content moderation, the Cyberspace Administration of China (CAC) issued new guidelines requiring that AI censorship models be transparent, controllable, and aligned with socialist values. The government also mandated "AI self-discipline initiatives," where tech companies must ensure their

algorithms align with government priorities.

These regulatory frameworks provide the legal backing for AI-driven censorship, ensuring that tech companies are compelled to invest in and refine their content moderation systems. Failure to comply with these regulations often results in hefty fines, platform suspensions, or even legal action against company executives.

### 1.3 Growth of Machine Learning in Content Control

The increasing reliance on machine learning and artificial intelligence has transformed the efficiency and effectiveness of internet censorship in China. AI-driven moderation systems now go beyond simple keyword filtering, incorporating deep learning techniques to analyze images, videos, and speech patterns.

One of the most significant advancements in AI censorship technology has been the development of NLP-based sentiment analysis models, which allow AI to interpret context rather than just detecting keywords. Previously, users could evade censorship by using homophones, puns, or abbreviations, but modern AI models are trained to recognize intent and sentiment, making them more resistant to circumvention strategies. For example, the phrase "May 35th" (instead of "June 4th") was commonly used to reference the Tiananmen Square Massacre, but newer AI models now recognize such coded language and remove related content accordingly.

Machine learning has also enabled multi-modal censorship, where AI can detect politically sensitive images, deepfake videos, and even speech patterns in real-time. On platforms like Douyin (TikTok's Chinese version), AI-driven censorship tools can instantly mute livestreams if they detect discussions about restricted topics. Similarly, facial recognition AI has been used to analyze and flag public figures or protest imagery in uploaded content.

Additionally, China has pioneered the use of predictive censorship models, which preemptively block discussions before they gain momentum. These models analyze search trends, discussion forums, and chat groups, predicting which topics are likely to become viral. If a sensitive topic starts gaining traction, AI can automatically throttle engagement by limiting visibility or suppressing discussions before they become widespread.

The Chinese government's investment in AI censorship infrastructure has also led to the rise of government-owned AI firms specializing in content control. Companies like SenseTime and Megvii, initially known for their advancements in facial recognition, have expanded into AI-driven content monitoring. These firms provide censorship algorithms not only to social media companies but also to government agencies, creating a centralized system where state and corporate censorship efforts are deeply integrated.

## 2. How AI Censorship Works

The efficiency and scope of AI-driven censorship in China rely on a complex multi-layered system of content moderation, combining keyword filtering, sentiment analysis, machine learning models, and human oversight. Unlike traditional censorship methods that relied heavily on manual review, AI-powered systems process vast amounts of online content in real-time, enabling authorities to detect and suppress sensitive discussions before they spread. However, AI censorship is not a singular, unified system; it varies across different platforms, regulatory priorities, and content types. The mechanisms behind AI censorship can be categorized into keyword filtering, sentiment analysis, AI-human hybrid moderation, and platform-specific strategies, each playing a unique role in shaping online discourse.

### 2.1 Keyword Filtering and Sentiment Analysis

The foundation of AI censorship in China is keyword filtering, a method that automatically detects and blocks specific words, phrases, or symbols associated with politically sensitive topics. Initially, keyword filtering was based on blacklists, where certain words (e.g., "Tiananmen Square," "democracy," "Hong Kong protests") were simply flagged and removed. However, as users developed workarounds using homophones, pinyin substitutions, and coded language, AI censorship evolved to include context-aware filtering and sentiment analysis.

Modern Natural Language Processing (NLP) models, powered by deep learning, now enable AI systems to analyze the intent and tone of a message rather than relying solely on exact word matches. For example, instead of blocking the word "censorship", AI systems assess the context in which it is used. A sentence like "The

government enforces necessary censorship to maintain stability" might pass through unfiltered, while "The government is suppressing free speech through censorship" could be flagged for removal.

Furthermore, sentiment analysis enables AI to detect dissenting opinions, sarcasm, and criticisms masked in neutral language. If a large volume of posts containing subtle negative sentiment toward the government emerges within a short time, AI algorithms can automatically suppress such discussions by reducing their visibility in search results or preventing their spread. This approach ensures that not only explicitly banned phrases but also discussions with a high probability of criticism or protest mobilization are effectively controlled.

One of the most advanced aspects of AI-driven censorship is adaptive filtering, where algorithms continuously learn and refine their detection methods. Chinese AI companies like SenseTime and Megvii, originally known for facial recognition, have developed censorship AI that monitors trends across social media platforms and dynamically updates keyword databases. This system ensures that new protest slogans, emerging political discussions, or viral anti-government phrases are quickly detected and blocked before they gain traction.

*2.2 AI vs. Human Moderation*

Despite AI's efficiency in detecting and filtering vast amounts of content, human moderation remains a crucial part of China's censorship model. While AI can quickly flag and remove suspicious content, it lacks the nuanced understanding of cultural context, evolving political trends, and new evasion tactics used by netizens. Therefore, most major platforms in China employ a hybrid model, where AI handles large-scale filtering, and human moderators review edge cases and high-risk content.

The typical AI-human censorship workflow follows these steps:

1) AI automatically scans and flags content based on keyword filters, sentiment analysis, and image recognition.

2) Immediate removal or suppression occurs for clear violations (e.g., banned words, protest slogans, or unauthorized political discussions).

3) Borderline content—posts that AI detects as "potentially sensitive" but uncertain—are sent to human moderators for review.

4) If a human reviewer deems the content acceptable, it remains visible; otherwise, it is deleted, and in some cases, the user's account may be suspended or flagged for further monitoring.

Human moderators are often employed by private tech companies such as Tencent (which operates WeChat) or ByteDance (which owns Douyin). These companies must comply with government censorship laws, meaning their human moderation teams often work under direct government supervision. Reports suggest that large-scale content moderation teams operate in shifts 24/7, particularly during politically sensitive periods like the anniversary of the Tiananmen Square protests or the National Congress of the Chinese Communist Party.

A major limitation of AI censorship without human oversight is the high rate of false positives and negatives. AI can over-censor neutral discussions, mistakenly flagging harmless content as politically sensitive, which leads to public frustration and complaints. Conversely, it can fail to detect more sophisticated evasion techniques, such as using memes, altered images, or satire, which require human judgment to interpret accurately.

To improve accuracy, many platforms implement real-time user feedback mechanisms, allowing users to report posts they believe should be censored. This method enhances AI learning, as frequent user reports help train algorithms to better detect new forms of disguised dissent.

*2.3 Platform-Specific Censorship Mechanisms*

Different Chinese platforms employ varied censorship strategies depending on their content type, audience, and government oversight level. The application of AI censorship is not uniform across the digital ecosystem, as each platform tailors its content moderation strategies to align with government expectations while maintaining user engagement.

**WeChat – Private Messaging and Public Accounts**

WeChat, operated by Tencent, serves as both a private messaging app and a social media platform through its "WeChat Moments" and

"Public Accounts." AI censorship on WeChat works in two primary ways:

- Private Messages: AI scans private messages for banned content. If sensitive words appear, messages may be blocked mid-sending or flagged for further scrutiny. Research has shown that WeChat even censors images containing text, using optical character recognition (OCR) to analyze pictures.

- Public Accounts and Articles: AI algorithms pre-screen content before publication. If an article is deemed sensitive, it is rejected outright, and repeated offenses can lead to account bans or government intervention.

**Weibo – Open Social Media Monitoring**

Weibo, often called "China's Twitter," is subject to some of the most aggressive AI-driven censorship due to its public and viral nature.

- AI continuously monitors trending hashtags and discussions, automatically suppressing politically sensitive topics before they gain momentum.

- Posts containing flagged keywords can be automatically "shadow-banned", meaning they remain visible to the poster but are hidden from other users.

- AI also prioritizes state-approved narratives, promoting government content while suppressing dissenting views.

**Douyin – Real-Time Video Moderation**

As China's equivalent of TikTok, Douyin faces the challenge of moderating live video content, which AI censorship systems manage through:

- Real-time audio and speech recognition, automatically muting livestreams if political discussions arise.

- Computer vision AI to analyze facial expressions, clothing, and objects, blocking symbols or gestures associated with protests or activism.

- AI-enhanced moderation teams that flag live interactions, ensuring that politically sensitive comments or user reactions are removed immediately.

**Baidu – Search Engine Censorship**

Baidu, China's dominant search engine, employs AI to control information access by:

- Dynamically adjusting search results to de-rank or remove politically sensitive topics.

- Generating "approved" information pages, redirecting users searching for banned terms to state-sanctioned content.

- AI-powered autocomplete suppression, ensuring users cannot even begin typing sensitive search queries.

**3. Accuracy in Detecting Sensitive Content**

AI-driven censorship in China is designed to be highly efficient in identifying and suppressing politically sensitive content. However, its effectiveness varies depending on the complexity of the content, user circumvention strategies, and evolving linguistic trends. While AI models have been successful in moderating direct political speech and misinformation, they face challenges in detecting nuanced discussions, satire, and coded language. This section examines the effectiveness of AI in identifying political speech and dissent and presents case studies of notable censorship operations while comparing China's AI censorship model to global counterparts.

*3.1 Effectiveness in Identifying Political Speech and Dissent*

The core function of AI censorship is to accurately detect and suppress politically sensitive discussions, ensuring that content critical of the government or diverging from state narratives is removed before it gains traction. Modern machine learning and NLP (Natural Language Processing) models enable AI to go beyond simple keyword filtering by analyzing semantic meaning, sentiment, and contextual usage. This allows AI to identify not just direct political criticism but also implied dissent, making censorship more comprehensive.

AI censorship is particularly effective in detecting direct political speech that includes names of political figures, events, or opposition movements. Posts mentioning Tiananmen Square, Hong Kong protests, or Taiwan independence are instantly flagged and removed. Additionally, references to banned organizations, foreign media sources, and human rights issues are proactively blocked. NLP-based sentiment analysis further assesses the tone of discussions, ensuring that even if a

sensitive topic is mentioned in a neutral manner, negative sentiment or criticism toward the government leads to suppression.

However, AI censorship is less effective against evolving circumvention tactics. Chinese netizens frequently use homophones, pinyin-based alternatives, emojis, and creative misspellings to bypass keyword filtering. For example, references to the Tiananmen Square Massacre (六四事件) are often replaced with "May 35th" (五月三十五日) to evade detection. AI systems have adapted by expanding their linguistic recognition capabilities, but new circumvention methods continue to emerge, creating a constant battle between users and censorship algorithms.

Another challenge for AI censorship lies in interpreting satire, memes, and indirect political commentary. While AI models can detect explicit text-based dissent, visual memes, symbolic imagery, and humor-based criticism are harder to identify accurately. For instance, images of Winnie the Pooh have been widely used to reference President Xi Jinping without direct textual criticism. While AI-enhanced computer vision tools can detect banned images, users frequently alter them by changing colors, adding distortions, or embedding hidden messages, making it difficult for AI to maintain accuracy in censorship enforcement.

*3.2 Case Studies of AI Censorship and Global Comparisons*

China's AI censorship system has demonstrated high efficiency in controlling online narratives, particularly during politically sensitive periods. Several high-profile censorship cases highlight the effectiveness and limitations of AI-driven moderation in managing public discourse.

**Case Study 1: COVID-19 Pandemic and Information Control**

During the early outbreak of COVID-19, AI censorship played a crucial role in suppressing whistleblower reports, independent journalism, and public criticism of government response efforts. Posts discussing the Wuhan lockdown, Li Wenliang (the doctor who warned about COVID-19), and inadequate medical supplies were quickly removed. AI censorship models targeted not only explicit criticisms but also indirect references, ensuring that emerging narratives did not challenge the government's official stance on pandemic control. However, this rapid censorship also blocked critical public health discussions and delayed

information-sharing, raising concerns about the balance between censorship efficiency and public safety.

**Case Study 2: Hong Kong Protests (2019-2020) and Real-Time Suppression**

The Hong Kong protests against the extradition law in 2019 became one of the most heavily censored topics on Chinese social media. AI-driven censorship systems flagged and removed live discussions, images, and videos of the protests, effectively erasing coverage from the Chinese internet. AI models were particularly effective in removing protest slogans, blacklisting hashtags, and blocking international news sources reporting on the demonstrations. However, protestors adapted by using coded language and indirect messaging, making it difficult for AI to completely eliminate all discussions. This case highlights how AI censorship is effective at rapidly suppressing mass discussions but struggles when facing user-driven linguistic adaptation.

**Case Study 3: The "996" Work Culture Controversy**

China's overwork culture, commonly referred to as "996" (working from 9 AM to 9 PM, six days a week), became a heavily discussed social issue, leading to a wave of online criticism against corporate exploitation. AI censorship systems monitored and suppressed discussions calling for labor rights reforms, blocking posts that referenced "anti-996" movements. However, due to widespread participation, new AI-driven moderation strategies, such as shadow-banning certain accounts and reducing the visibility of related posts without direct deletion, were implemented. This showcases how AI censorship can adapt its strategies when direct suppression proves ineffective.

*3.3 Comparison with Global AI Censorship Models*

While China's AI censorship model is one of the most advanced and comprehensive in the world, it differs significantly from AI-based content moderation in other countries, particularly in the U.S., EU, and Russia.

**United States and EU – AI for Misinformation Control, Not Political Censorship**

In Western democracies, AI content moderation is primarily used to combat misinformation, hate speech, and harmful content rather than suppressing political dissent. Platforms like

Facebook, Twitter, and YouTube employ AI-driven moderation to detect fake news, extremist propaganda, and disinformation campaigns. However, unlike China, these AI models are not state-controlled—companies have independent policies for content moderation, though government pressure and regulations (such as the EU's Digital Services Act) influence how they function.

Additionally, Western AI moderation models emphasize user transparency, often providing appeal mechanisms where flagged content can be reviewed and reinstated. This contrasts with China's opaque censorship system, where removals occur without explanation and appeals are rarely possible.

### Russia – AI Censorship with Government Influence

Russia's internet censorship model shares some similarities with China, particularly in its use of AI for political content control. Platforms are required to comply with state censorship demands, and AI models are used to monitor online dissent, suppress protests, and control opposition narratives. However, Russian AI censorship is less centralized and sophisticated compared to China, relying more on direct government intervention rather than fully automated AI-driven suppression.

### The Global Export of China's AI Censorship Model

As China refines its AI censorship infrastructure, it is increasingly exporting its technology to other authoritarian regimes, providing AI moderation tools to countries like Iran, Venezuela, and Ethiopia. These state-controlled digital monitoring solutions are shaping a global trend toward AI-assisted speech regulation, raising concerns about the potential for AI-driven authoritarianism beyond China's borders.

### 4. Bias in AI Censorship

AI-driven censorship in China, while highly efficient, is not free from systemic biases that shape its implementation. These biases arise from the over-censorship of neutral content, disproportionate targeting of specific groups and ideologies, and the influence of government-controlled datasets used to train AI models. Unlike manual censorship, where human moderators can exercise some level of judgment, AI censorship operates on algorithmic

patterns that reflect the priorities of its training data and policy directives. As a result, AI does not just enforce censorship—it reinforces state narratives and ideological control in ways that can have unintended consequences on public discourse, civil society, and even economic sectors that rely on digital communication.

One of the most significant biases in AI censorship is its tendency to over-censor neutral or unrelated content due to pattern-matching errors. Since AI models are trained to detect politically sensitive topics, they often flag content that merely resembles or indirectly references restricted subjects. For example, a discussion on "Tiananmen Square" as a historical landmark may be censored because the AI fails to differentiate between a tourist conversation and a politically sensitive discussion about the 1989 protests. Similarly, words with dual meanings or phonetic similarities to censored terms are frequently blocked, leading to confusion and frustration among users. During the COVID-19 pandemic, posts discussing "virus origins" in a scientific context were often removed because AI systems misinterpreted them as promoting conspiracy theories or politically sensitive narratives about China's handling of the crisis. This broad overreach affects academic discussions, news reporting, and even cultural conversations, highlighting the imperfect nature of AI censorship algorithms.

Beyond general over-censorship, AI censorship in China disproportionately targets specific groups and ideological perspectives, particularly those that challenge government policies, human rights issues, or ethnic minority struggles. While censorship applies broadly across the internet, certain communities experience far higher levels of suppression than others. Topics related to Tibetan independence, Uyghur rights, LGBTQ+ activism, feminist movements, and labor rights protests are censored more aggressively than other political discussions. AI models are trained to identify and suppress dissenting viewpoints related to these topics, ensuring that narratives contradicting state policies remain invisible to the public. In contrast, government-approved discussions, even if political in nature, are often amplified rather than restricted. For instance, while AI swiftly removes discussions about pro-democracy protests, it allows, and even promotes, posts that support state actions in

Hong Kong, Taiwan, or Xinjiang. This asymmetry in censorship enforcement highlights how AI functions not just as a neutral content filter but as an active agent of state ideology.

A major contributing factor to AI bias in censorship is the government-controlled dataset used to train censorship models. Machine learning algorithms rely on historical data to identify patterns in speech, imagery, and video content, but when the training data itself is biased, the AI replicates and reinforces those biases. Since Chinese AI censorship models are developed under strict state regulations, they are trained primarily on datasets that exclude oppositional viewpoints and prioritize state narratives. This results in an AI system that is structurally incapable of recognizing alternative perspectives because it has never been exposed to them in its learning process. For example, discussions about democracy, human rights, or government accountability are framed in ways that favor state-approved interpretations, leading to an AI model that automatically suppresses content associated with these topics. In addition, AI systems often exhibit regional and linguistic biases, disproportionately censoring content in minority languages such as Uyghur, Tibetan, and Cantonese more aggressively than Mandarin-based discussions.

The implications of AI bias in censorship extend beyond political discourse and into everyday digital interactions. Users who engage in discussions on sensitive topics may experience shadow-banning, account restrictions, or reduced content visibility without clear reasons, making it difficult for them to understand the censorship logic at play. The lack of transparency and appeal mechanisms means that AI censorship functions as a one-way system of suppression, with no way for users to challenge or reverse automated decisions. This further entrenches self-censorship, as users, fearing penalties, learn to avoid certain discussions altogether, reinforcing a digital environment where government-approved narratives dominate public discourse.

Ultimately, the biases present in AI censorship are not accidental flaws but systemic design choices that align with China's broader strategy of digital authoritarianism. By leveraging AI to automate and refine content control, the government has created a censorship model that is not only highly efficient but also deeply aligned with state ideology. As AI censorship continues to evolve, the challenge of addressing bias, transparency, and accountability remains a critical issue, especially as China begins exporting its AI-driven content regulation technologies to other countries. The next section will examine the challenges and unintended consequences of censorship precision, highlighting how AI systems struggle with detecting disguised dissent, over-blocking content, and balancing censorship with misinformation control.

## 5. Challenges in Censorship Precision

Despite its efficiency, AI-driven censorship in China faces significant challenges in maintaining precision, often struggling to distinguish between genuine threats to political stability and harmless discussions. While machine learning algorithms have greatly improved the ability to identify and remove politically sensitive content, they remain imperfect, frequently failing to detect disguised dissent, over-blocking neutral content, or under-blocking actual violations. At the same time, the government must balance its desire to control political discourse with the need to combat misinformation, often leading to contradictory enforcement patterns. These challenges highlight the limitations of automated censorship systems, which require continuous refinement to remain effective in a rapidly evolving digital environment.

One of the most persistent problems in AI censorship is its inability to reliably detect disguised dissent. As Chinese internet users become more adept at circumventing keyword-based filtering, they employ coded language, satire, memes, and indirect references to discuss sensitive topics without triggering censorship mechanisms. Instead of using banned words like "Tiananmen" or "protest", users substitute homophones, misspellings, or pinyin variations to evade AI detection. For example, rather than directly referring to censorship, users might use phrases like "river crab" (河蟹, a homophone for "harmony") or "404" (symbolizing deleted content). Similarly, protest slogans are replaced with abstract symbols, emoji combinations, or modified images, making it increasingly difficult for AI to recognize subversive content. Satirical content presents another challenge, as AI struggles to differentiate genuine praise from sarcasm, leading to inconsistencies in enforcement. For instance, a post saying "China has the best free

speech in the world!" may be flagged for removal, while a more subtle critique hidden within a humorous meme could go undetected. This gap in AI precision forces the government to continuously update and retrain censorship algorithms, making it a never-ending race between censors and internet users.

While AI censorship sometimes fails to detect hidden dissent, it also frequently over-blocks content that poses no real political risk, leading to frustration among users and businesses. Since AI models operate on predefined rules and patterns, they often misinterpret neutral or unrelated content as politically sensitive, resulting in excessive censorship. Academic discussions on political theory, historical events, or international relations are often removed simply because they contain keywords flagged by AI filters. Similarly, content related to fictional stories, literature, or artistic works that coincidentally mention censored terms can be taken down without justification. This issue became evident when Chinese netizens discussing the film *Les Misérables* (which contains themes of revolution and protest) found their posts blocked, as AI mistakenly associated the content with real-world dissent. Over-blocking extends beyond political discussions, sometimes affecting businesses, online education platforms, and entertainment media, as AI systems prioritize political risk avoidance over contextual understanding. In cases where content touches on sensitive subjects in a neutral or supportive manner, users often have no way to appeal censorship decisions, further exacerbating the problem.

Conversely, AI censorship also faces the risk of under-blocking, where certain sensitive discussions evade detection due to algorithmic weaknesses or system overloads. During major political events such as the National Congress of the Communist Party or large-scale protests, online discussions surge, sometimes overwhelming AI censorship systems. This can lead to delayed censorship enforcement, where controversial discussions briefly go viral before being removed. Additionally, the rapid evolution of dissident communication tactics, such as using live video streams, coded speech, and encrypted messaging apps, presents challenges that AI censorship systems struggle to keep up with. Some users have even exploited AI weaknesses by embedding sensitive messages within seemingly unrelated content,

such as gameplay videos, music lyrics, or digital artwork, where automated systems may fail to recognize the subversive intent.

A further challenge in censorship precision lies in the Chinese government's dual objective of suppressing political dissent while also combating misinformation. While AI censorship is primarily focused on removing politically destabilizing content, it is also deployed to curb rumors, conspiracy theories, and fake news, particularly during crises such as public health emergencies, economic downturns, or geopolitical conflicts. However, distinguishing between misinformation and politically motivated speech suppression is often problematic, leading to contradictions in enforcement. For example, during the COVID-19 pandemic, AI censorship aggressively removed content that questioned official government reports, including whistleblower testimonies and independent news reports, labeling them as misinformation even when they contained verified facts. This demonstrates how AI censorship can sometimes prioritize state propaganda over objective truth, resulting in the suppression of vital public health information. At the same time, the government struggles to maintain credibility in its anti-misinformation efforts, as excessive censorship of inconvenient truths leads to distrust in official narratives. This contradiction complicates the role of AI censorship, forcing authorities to constantly recalibrate their moderation policies in response to shifting public sentiment.

The limitations of AI censorship in China underscore the complexity of maintaining digital control in an era of increasingly sophisticated circumvention tactics. While AI has greatly enhanced the speed and scale of content moderation, its inability to perfectly interpret human language, satire, and hidden dissent means that enforcement remains inconsistent. The issue of over-blocking and under-blocking continues to challenge the system's credibility, affecting ordinary users, businesses, and online communities. Meanwhile, the blurring line between political censorship and misinformation control raises ethical concerns about the role of AI in shaping public discourse. As China continues to refine its AI censorship models, the tension between efficiency, accuracy, and control will remain a defining challenge, shaping the future of

state-sponsored digital governance. The next section will explore how both domestic and international communities react to China's AI censorship, highlighting the strategies that netizens use to bypass restrictions and the global response to China's expanding digital influence.

**References**

Chen, H., & Xu, L. (2019). The role of artificial intelligence in online misinformation control: A case study of censorship during COVID-19. *Journal of Media Ethics and AI Governance, 22*(4), 301-328.

Feng, K., & Gao, J. (2023). Exporting AI censorship: China's influence on global content moderation practices. *Journal of Global Internet Governance, 17*(2), 87-109.

Huang, Y., & Wang, J. (2021). AI-driven censorship in China: Balancing state control and digital governance. *Journal of Digital Communication Studies, 14*(3), 221-243.

Li, X., & Zhang, T. (2020). Algorithmic bias in content moderation: Analyzing China's AI censorship mechanisms. *Asian Journal of Political Science & Technology, 28*(2), 155-178.

Wang, Z., & Liu, P. (2022). The influence of government data on AI censorship efficiency: Evidence from Weibo and WeChat. *Chinese Journal of Digital Society, 30*(1), 189-210.

Zhao, M., & Sun, Y. (2018). AI censorship and political discourse in China: Examining automated content suppression on social media platforms. *International Journal of Media Control, 25*(3), 112-134.