



# AI Bot Crawling and Broadcasters: Strategy, Tradeoffs and Common Approaches

## Deciding Whether or Not to Block AI Bots

Artificial Intelligence (AI) systems increasingly act as intermediaries between audiences and information sources. In this role, AI systems use automated crawlers (“bots”) to systematically browse and collect online content, which is then used by AI models, including large language models, to generate responses to user queries.

In deciding whether to allow or restrict automated crawling, broadcasters must balance several strategic considerations, including **content control, audience reach, infrastructure costs, risks related to accuracy and brand perception, potential monetization opportunities and broader long-term technology and AI strategy.**

Broadcasters take many approaches to managing automated AI crawlers. Some restrict crawler access to control how their content is used, while others allow access to support discoverability in AI-driven information tools. Decisions typically reflect an organization’s individual assessment of its strategy, infrastructure capacity and views about emerging AI technologies.

## Strategic Considerations

### Control of Content Use

AI systems may collect text, audio and video to generate summaries, answers or other outputs. Restricting AI crawler access can provide broadcasters with some control over how their content is accessed and incorporated into AI systems; however, such measures are not always fully effective. Many have characterized bot blocking as a “whack-a-mole” problem, as crawlers can adapt to blocking methods or find alternative ways to access content. To the extent these measures are effective, broadcasters may be better positioned to support appropriate attribution, manage downstream use and maintain editorial control over how content appears in automated outputs.

### Visibility and Discoverability

Allowing access to crawlers may increase the likelihood that content appears in AI-generated responses or summaries, while restricting access may limit visibility in AI-driven discovery environments.

## Infrastructure and Technical Costs

High volume automated scraping can generate large numbers of server requests and data transfers, increasing bandwidth usage, server load and compute costs. Significant bot traffic may also trigger autoscaling in cloud infrastructure environments, resulting in additional operational expenses and system monitoring requirements.

## Accuracy, Attribution and Reputational Risk

AI-generated outputs may contain inaccuracies, incomplete context or incorrect attribution, which may affect how reporting is represented and create potential reputational risks for broadcasters. AI systems also may not reflect real-time updates, particularly for breaking news, resulting in outdated or incomplete information. At the same time, restricting access to high quality, reputable sources could contribute to inaccuracies by leading AI systems to rely more heavily on lower quality or less reliable material.

## Monetization and Licensing Considerations

The automated collection of content raises questions about whether and how organizations may capture value from the use of their material in AI systems. Restricting crawler access may be evaluated as one mechanism for preserving potential licensing or partnership opportunities. On the other hand, enabling crawler access may provide an opportunity to demonstrate that value through its use and visibility in AI-generated outputs.

## Long Term Technology Strategy

AI technologies are rapidly evolving and increasingly integrated into search tools, newsroom workflows and audience engagement platforms. Decisions regarding crawler access may therefore be considered as part of broader planning around how organizations interact with emerging AI ecosystems.

### Common Reasons That May Prompt Blocking

- **Large volumes of automated traffic** that increase server load or infrastructure cost (ex. One broadcaster experienced up to 10x more bot traffic than legitimate traffic).
- **Unusual traffic patterns** indicating automated scraping activity.
- **Use of content without attribution, referral traffic or compensation.**
- Concerns about **AI-generated inaccuracies or misattribution of content.**
- **Risk of enabling voice cloning or synthetic media** due to the collection of audio and video content.
- **Loss of control of content** and future use.

### Common Reasons That May Discourage Blocking

- **Reduced visibility** if AI systems become an important content discovery channel.
- **Loss of referral traffic** from search engines or AI-driven search features.
- **Technical complexity** of implementing and maintaining advanced blocking systems.
- **Difficulty distinguishing** between **legitimate and malicious crawlers.**
- **Limited enforcement** for certain blocking mechanisms.
- **Inaccurate AI-generated information** erodes trust and misleads audiences who expect accurate news (ex. One broadcaster faced public backlash and station protests after inaccurate information was attributed to their reporting).

# Common AI Bot Blocking Approaches and Solutions

---

Blocking AI bots is not as simple as pressing a “block” button; it requires a strategic, layered approach. Most organizations managing automated crawler access utilize a combination of website rules, server level protections and specialized bot management platforms. These approaches vary in complexity, effectiveness, cost and level of enforcement.

## 1. Website Level Controls (Robots Exclusion Protocol)

Website-level controls help manage crawler access by publishing directives that communicate a site’s access preferences to automated systems. The primary mechanism for doing so is the Robots Exclusion Protocol (REP), a standardized framework for communicating crawler permissions through machine-readable rules. Because REP does not enforce access restrictions, its effectiveness depends on whether AI bots accurately identify themselves and choose to follow the rules.

### Solution Example: robots.txt

REP is implemented through a robots.txt file located at the root of a website, which specifies allow and disallow rules for different crawler user agents. It allows site operators to define access at a granular level (e.g., by directory or page.)

Example robots.txt configuration:

```
User-agent: GPTBot
Disallow: /

User-agent: ClaudeBot
Disallow: /
```

### Key Capabilities:

- Controls crawler access at the page or site level.
- Signals permission to compliant bots.

### Advantages

- Simple to implement.
- Widely recognized by legitimate crawlers.
- No infrastructure changes required.

### Limitations

- Compliance is voluntary.
- Does not technically prevent scraping.
- Malicious bots may ignore it.

## 2. Server or Firewall Controls

Server or firewall level controls manage crawler access by enforcing rules at the network or infrastructure layer, evaluating and filtering requests before they reach the application. These controls evaluate traffic based on factors such as IP reputation, request patterns, and known bot characteristics.

### Solution Examples: AWS WAF Bot Control, Google Cloud Armor

These solutions enable organizations to block, filter or rate-limit automated traffic based on configurable rules and traffic analysis, helping mitigate unwanted bot activity and reduce infrastructure load.

#### Key Capabilities:

- Prevents unwanted requests at the infrastructure level.
- Blocks high volume automated traffic.
- Restricts suspicious IP ranges.
- Identifies scraping behavior patterns.
- Enables customizable, rules-based enforcement.

#### Advantages

- Strong enforcement compared to website level controls.
- Reduces server load and infrastructure costs.
- Integrates with cloud hosting environments.

#### Limitations

- Requires engineering configuration and monitoring.
- Effectiveness depends on rule tuning.

## 3. Bot Detection and Management Platforms

Bot detection and management platforms manage crawler access by identifying and classifying automated traffic using a combination of behavioral analysis, device characteristics, and network-level signals.

### Solution Examples: Cloudflare Bot Management, DataDome Bot Protect

These solutions enable organizations to distinguish between human and automated activity and apply mitigation measures based on that classification.

#### Key Capabilities:

- Identifies automated traffic using IP reputation, device fingerprinting and behavioral analysis.
- Analyzes request patterns and network-level traffic signals.
- Classifies traffic as human or automated.
- Automatically blocks, challenges or allows requests.

#### Advantages

- Advanced bot detection accuracy due to large traffic dataset.
- Capable of integrating with CDN services.
- Automated mitigation tools.

#### Limitations

- Additional cost.
- Time and effort required for infrastructure integration.
- Configuration can be complex.

## 4. Traffic Monitoring and Analytics

Traffic monitoring and analytics tools help organizations make informed decisions about whether and how to implement blocking or mitigation measures by providing visibility into automated activity and traffic patterns.

### Solution Examples: Server logs, CDN analytics tools

These tools provide detailed data on request patterns, traffic volume, and sources, enabling analysis of crawler behavior and detection of anomalies. While they can help identify potential scraping activity, they do not directly block or restrict access.

#### Key Capabilities:

- Analyzes traffic patterns and bot behavior.
- Identifies spikes or anomalies in automated activity.

#### Advantages

- Improves understanding of bot traffic.
- Enables more targeted mitigation strategies.

#### Limitations

- Does not prevent scraping on its own.
- Requires technical expertise to interpret data.

### Additional Resources

Robots Exclusion Protocol: <https://www.rfc-editor.org/rfc/rfc9309.html>

Google Robots.txt Documentation: <https://developers.google.com/search/docs/crawling-indexing/robots/intro>

OpenAI Crawler Information (GPTBot): <https://platform.openai.com/docs/bots>

Anthropic ClaudeBot Crawling Policy: <https://support.claude.com/en/articles/8896518>

AWS WAF Bot Control Documentation: <https://docs.aws.amazon.com/waf/latest/developerguide/waf-bot-control.html>

DataDome Bot Protection Documentation: <https://docs.datadome.co>

Common Crawl CCBot Information: <https://commoncrawl.org/ccbot>