

Why Whitelist?

How Common Approaches to Beating Bot
Traffic Fall Short

Introduction

Within the digital advertising industry, waste caused by non-human traffic is a well-documented problem. One 2014 study found that over half of all internet traffic was generated by bots.¹ By another estimate, in 2015 marketers will waste \$6.2 billion on false clicks generated by bots – to say nothing of the dollars wasted on CPM impressions served to non-humans.²

There are many technology providers attempting to solve this multi-billion dollar problem, and most of them take the same approach: blacklisting.

Like a brand might blacklist low-quality or family-unfriendly sites on an exchange, fraud prevention firms use blacklists in an attempt to weed out non-human impressions and counter waste. Yet although this is the most common methodology (and certainly an improvement over not fighting fraud), blacklisting is not the best approach available today.

To return to our website example, when a brand blacklists sites on an exchange, there are often far too many sites cropping up every day for the blacklist to remain comprehensive. Offending sites often manage to slip through the cracks. The same is true with blacklisting to fight bot fraud.

Now imagine the brand has a whitelist rather than a blacklist, and has pre-approved a list of high-quality, brand-safe sites. In this instance, there's no way a low-quality or offensive site will make an appearance. Whitelisting to fight fraud, though less common, is a superior approach because it ensures that all impressions are served to a human.

This guide will review four distinct approaches to fight fraud and explain the benefits and pitfalls of each.

First, let's take a deeper look at the two most common approaches to blacklisting: IP blacklisting and user-level blacklisting.

Blacklisting for Fraud Prevention

What is IP Blacklisting?

One of the more common approaches is IP-based blacklisting, which relies on finding patterns that represent fraud and, when found, blocking traffic from that IP.

Why IP Blacklisting Falls Short

While blacklisting by IP address is certainly a much more robust method than say, site-level blacklisting, it remains an imperfect methodology due to two major shortcomings: it is imprecise and more importantly, it only detects robust, malicious attacks with lag time.

IP-based targeting is often used in digital advertising, most commonly as an approximation of targeting an individual user. Much the same way IP targeting attempts to isolate a user, IP-level blacklisting is an approximation of a user-level blacklist, and is sometimes even referred to as such, though this is a misnomer.

As it stands, multiple people and multiple computers often use a single IP address. Public IP addresses, for example, can support more than 50,000 connections simultaneously. Think of a coffee shop or an office where many people are using the same network. If one computer in the network is infected, the entire IP is still blocked, and all of those users (many of whom are almost certainly humans) will also be blocked as well.

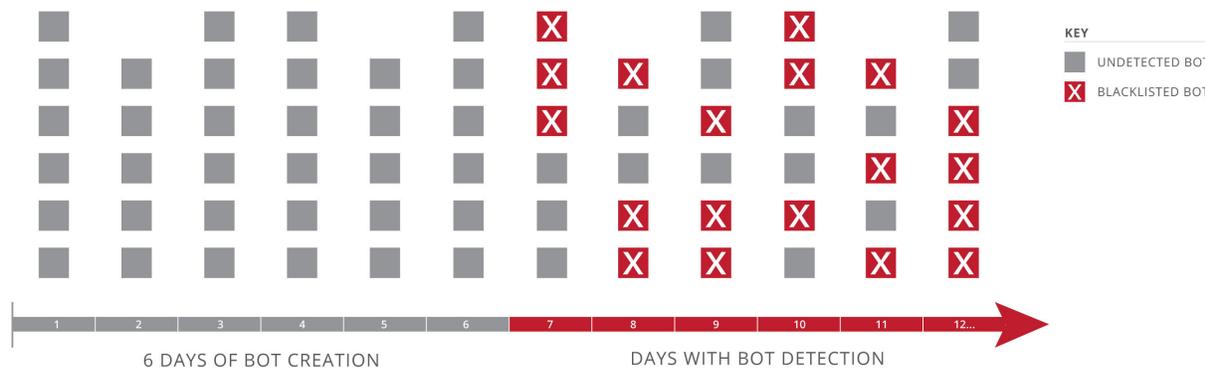
Secondly, and more importantly, this methodology is only able to identify and recognize malicious attacks with some lag time. Ad fraudsters only employ bots for a very short period of time -- on average six days. Six days isn't an arbitrary number, that's approximately how long it takes for the leading providers' algorithms to identify a new bot. As a rule, most blacklists are not

1. <http://www.incapsula.com/blog/bot-traffic-report-2014.html>

2. <http://www.nytimes.com/2014/12/10/business/media/study-puts-a-price-tag-on-digital-ad-click-fraud-.html>

recognizing bot traffic until the bot has already ceased to exist. Even in the instance where a bot is recognized in a very, very short amount of time--let's say two days--you will still be serving ads to bots over the course of those two days.

Every day, bad actors pioneer new ways to approximate user behavior, and firms who rely on blacklisting play a never-ending game of catch up.



IP blacklisting also does not protect brands from impressions generated from the “good bots,” that make up 27% of web traffic. “Good bots” can include crawlers like GoogleBot that index pages for search engines, scrapers that index data like product listings or pricing or weather information. While many good bots, like Google, do identify themselves in order to block ad impressions, many serve purposes unrelated to advertising do not. Since these bots aren't creating deliberate or malicious attacks, they're often missed by blacklists designed to root out intentional advertising fraud.

What is User-Level Blacklisting?

User-level blacklisting is a more precise methodology than IP-level blacklisting, however it suffers from the same major shortcoming: it responds only to recognizable, malicious attacks and there is always a degree of lag time.

There is a second issue with user-level blacklisting. The most common form of bot is a real person's computer that's been infected with malicious code. Over time, many people will at some point have malware on their computer, but most people also will eventually clean their system. If they are recognized as a bot and blacklisted, that legitimate user is blacklisted even after they clean their system.

Whitelisting

At a conceptual level, blacklisting assumes that a user is human, and only removes users from an audience pool once they're proven to be a bot. However, given that humans are responsible for less than half of online traffic, and billions of dollars are at stake, that assumption isn't the right place to start. Whitelisting, on the other hand, takes the opposite tack, and assumes a user is a bot until proven otherwise.

There are two primary methods for whitelisting human traffic: **site-level whitelisting** and **user-level whitelisting**.

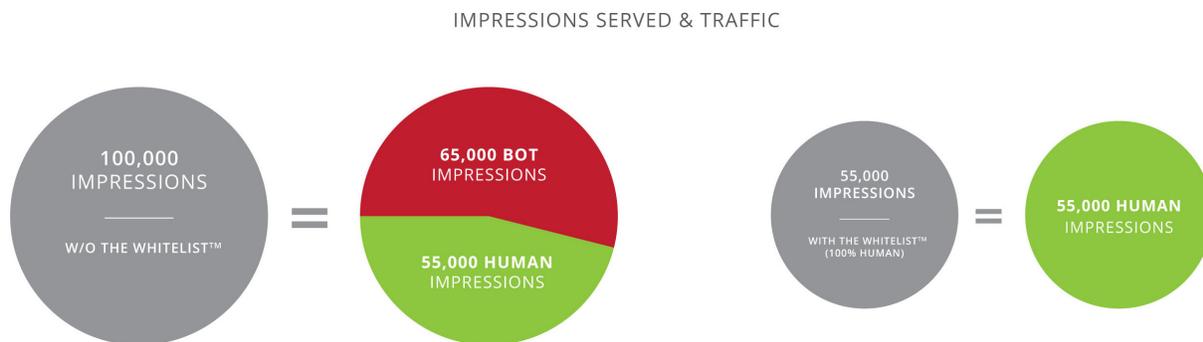
Site-Level Whitelisting

One of the more common approaches is IP-based blacklisting, which relies on finding patterns that represent fraud and, when found, blocking traffic from that IP.

Why Site-Level Whitelisting Falls Short

That said, even the highest performing sites have some level of fraud, which means marketers relying on site-level whitelisting are still wasting money on impressions served to bots. Top-ranked Alexa sites that see upwards of one million daily are still plagued by fraud: by one estimate, 50% of their traffic in 2014 was non-human.³ Even if the figures for many top sites are much lower, say 20%, you're still wasting a fifth of your budget on non-human traffic.

Secondly, by whitelisting only prominent sites, you could be missing valuable opportunities to reach your target audience on smaller sites that likely charge less and could be highly relevant and aligned with your messaging.



User-Level Whitelisting

User-level whitelisting is a more precise method of focusing advertising dollars on human traffic. Using algorithms that detect uniquely human behavior, this method creates a whitelist of verified human traffic.

This method is particularly tricky, because the more sophisticated bots can behave in ways that closely mimic human behavior—they can click on ads, mimic page scrolling or even add products to shopping carts.⁴ That said, there are many behaviors such as completing a purchase that, taken together over time, can indicate with certainty that a given user is a human. What bots don't do is act like humans over a long period of time in many, diverse scenarios. With those indications in place, users are added to a whitelist that is regularly updated. The key to an effective whitelist is gathering data from users in many situations, across many sites, many times a day.

Conclusion

User-level whitelisting is the only method that can conclusively ensure that no advertising impression is wasted on non-human traffic. Whitelists can never include every single human, which means there are real humans who you won't serve ads to. That concern is not without merit. However, the risk of missing out on some humans seems small compared to the certainty of wasting money on meaningless impressions when you rely on blacklists.

If you can't find verified humans, it doesn't make sense to backfill a campaign with impressions that are just as likely as not to be served to a bot. Even if you serve fewer overall impressions, you know with 100% confidence that those impressions are being served to a human. Given that around half of ad impressions are served to bot traffic, even if you were to dial back your impression volume by half you're still serving ads to roughly the same number of people.

Massive amounts of marketers' budgets are wasted on ad impressions that are served to bots. The most common methods to fight fraud—the different forms of blacklisting—will limit how many leave a lot of gaps. Whitelisting is the only conclusive way to focus spend on human traffic and human traffic only.

3. <http://www.incapsula.com/blog/bot-traffic-report-2014.html>

4. <http://www.mediapost.com/publications/article/246824/kill-the-bots-first-then-tackle-viewability.html>