

# Dinamik Web Scraping ve Sayfa Analiz Aracı

HAZIRLAYAN: GAZALİ KEPENÇ

TARİH: 20.12.2025

## **İÇİNDEKİLER**

<b>1.GİRİŞ.....</b>	<b>3</b>
<b>2.TEKNİK METODOLOJİ.....</b>	<b>3</b>
<b>3. ANALİZ TABLOSU.....</b>	<b>3</b>
<b>4. TEKNİK ÇIKTI ÖRNEKLERİ.....</b>	<b>4</b>
<b>4.1 Terminal Üzerinden Çalıştırma ve Log Kayıtları (HTTP 200 – Başarılı Senaryo).....</b>	<b>4</b>
<b>4.2 Terminal Üzerinden Çalıştırma ve Log Kayıtları (HTTP 403 – Erişim Kısıtlı Senaryo).....</b>	<b>5</b>
<b>4.3 Dosya Sistemi ve Çıktı Yönetimi.....</b>	<b>5</b>
<b>4.4 Görsel Veri Analizi.....</b>	<b>6</b>
<b>4.5 403 Erişim Kısıtı Altındaki Görsel Çıktı Değerlendirmesi.....</b>	<b>7</b>
<b>5.SONUÇ VE DEĞERLENDİRME.....</b>	<b>7</b>

## 1.GİRİŞ

Bu rapor, Siber Tehdit İstihbaratı (CTI) faaliyetlerinde temel bir yetenek olan **web tabanlı veri toplama ve sayfa analizinin** otomatikleştirilmesini amaçlayan *Dinamik Web Scraping ve Sayfa Analiz Aracı* çalışmasının sonuçlarını sunmaktadır.

Modern tehdit aktörleri; forumlar, haber siteleri, açık kaynak platformlar ve kamuya açık web sayfaları üzerinden iz bırakmaktadır. Bu nedenle bir CTI analistinin, hedef web kaynaklarını **güvenilir, kanıt üretir ve tekrar edilebilir** şekilde analiz edebilmesi kritik önemdedir.

Geliştirilen araç;

- Hedef web sayfasının ham HTML içeriğini,
- Sayfanın tarama anındaki görsel durumunu (ekran görüntüsü),
- Sayfa üzerindeki harici bağlantı yapısını (link haritası)

Otomatik olarak elde ederek CTI analiz süreçleri için temel veri setini oluşturmaktadır.

## 2. TEKNİK METODOLOJİ

Araç, modern web uygulamalarında yaygın olarak kullanılan JavaScript tabanlı dinamik içerik üretimini doğru şekilde analiz edebilmek amacıyla Headless Browser yaklaşımı ile tasarlanmıştır.

Kullanılan Teknolojiler

- **Programlama Dili:** Go (Golang)
- **Tarayıcı Motoru:** chromedp (Chrome DevTools Protocol)

Chromedp tercih edilmesinin temel nedeni; yalnızca statik HTML çekimi değil, aynı zamanda **tarayıcı seviyesinde render edilmiş DOM yapısına** erişim sağlamasıdır. Bu yaklaşım, klasik HTTP tabanlı scraper'ların kaçırdığı dinamik içeriklerin analiz edilmesine olanak tanır.

**Temel Yetenekler**

- Hedef URL'ye otomatik navigasyon
- Render edilmiş DOM ağacının tamamının çıkarılması
- HTTP durum kodunun tespiti
- Mükerrer bağlantıların elenerek sayfa link haritasının oluşturulması
- Tam sayfa ekran görüntüsü alınarak görsel kanıt üretilmesi

Bu metodoloji, CTI süreçlerinde kanıta dayalı analiz yapılmasını desteklemektedir.

## 3. ANALİZ TABLOSU

Bu bölümde, geliştirilen aracın 15 farklı karakteristikteki web sitesi üzerindeki performansı ve elde edilen istihbari veriler tablolandırılmıştır.

No	Web Sitesi	Platform Türü	İçerik Kapsamı	Erişim Durumu	Teknik Not
1	btkt.gov.tr	Resmî Kurum	Mevzuat, duyuru	200	Statik + yarı dinamik
2	archive.org	Dijital Arşiv	Açık veri, arşiv	200	Yüksek veri erişilebilirliği
3	youtube.com	Sosyal Medya	Video, medya	200	JS ağırlıklı
4	google.com	Arama Motoru	Genel bilgi	200	Bot algılama mevcut
5	figma.com	Tasarım Platformu	UI/UX içerik	200	Dinamik SPA
6	telegram.org	İletişim Platformu	Mesajlaşma	200	Hafif bot koruması
7	stackoverflow.com	Teknik Forum	Yazılım, Q&A	403	İçerik yine de alındı
8	yildizcti.com	CTI Platformu	Siber tehdit raporları	200	CTI odaklı
9	imdb.com	Medya Veritabanı	Film, dizi	403	Anti-bot aktif
10	shodan.io	Siber Güvenlik	Açık servis tarama	200	Güçlü CTI kaynağı
11	github.com	Kod Deposu	Açık kaynak	200	Dinamik içerik
12	sans.org	Güvenlik Kurumu	Eğitim, rapor	200	Kurumsal yapı
13	linkedin.com	Sosyal Ağ	Profesyonel veri	200	Oturum bağımlı
14	ktun.edu.tr	Akademik Kurum	Duyuru, akademik	200	Statik yapı
15	sciencedirect.com	Akademik Yayın	Bilimsel makale	403	Yayıncı koruması

#### 4. TEKNİK ÇIKTI ÖRNEKLERİ

Sistemin çalışma prensibi ve ürettiği veriler aşağıda örnek çıktılar üzerinden gösterilmiştir. Bu bölümde, farklı HTTP durum kodlarına sahip hedefler üzerinde aracın davranışı ve ürettiği teknik çıktılar karşılaştırmalı olarak sunulmaktadır.

##### 4.1 Terminal Üzerinden Çalıştırma ve Log Kayıtları (HTTP 200 – Başarılı Senaryo)

Uygulama, komut satırı üzerinden hedef URL argümanı olarak çalıştırılmıştır. Navigasyon, HTTP durum kodu tespiti ve dosya üretim süreci terminal üzerinden anlık olarak izlenmiştir.

Çalıştırılan Komut: **go run main.go https://www.sans.org**

Terminal Çıktısı Analizi:

- Hedef URL başarıyla yüklenmiştir.
- HTTP durum kodu 200 OK olarak tespit edilmiştir.
- HTML içeriği, tam sayfa ekran görüntüsü ve bağlantı listesi hatasız şekilde oluşturulmuştur.

```
PS C:\Users\ASUS\Documents\Webscraper> go run main.go https://www.sans.org
Hedef URL işleniyor: https://www.sans.org
Durum Kodu: 200
Kaydedilen Dosyalar:
- www_sans_org.html
- www_sans_org.png
- www_sans_org_links.txt
```

**Görsel 1:** go run main.go https://www.sans.org komutunun çalıştırılması, HTTP 200 durum kodunun alınması ve çıktı dosyalarının başarıyla oluşturulduğunu gösteren terminal çıktısı.

#### 4.2 Terminal Üzerinden Çalıştırma ve Log Kayıtları (HTTP 403 – Erişim Kısıtlı Senaryo)

Bu senaryoda, bot algılama ve erişim kısıtlaması uygulayan bir platform üzerinde aracın davranışı gözlemlenmiştir.

Çalıştırılan Komut: **go run main.go https://www.imdb.com**

Terminal Çıktısı Analizi:

- Hedef URL'ye bağlantı kurulmuştur.
- HTTP durum kodu 403 Forbidden olarak alınmıştır.
- Erişim kısıtına rağmen sayfa içeriği render edilmiş ve teknik çıktılar üretilmiştir.

```
PS C:\Users\ASUS\Documents\Webscraper> go run main.go https://www.imdb.com
Hedef URL işleniyor: https://www.imdb.com
Durum Kodu: 403
Kaydedilen Dosyalar:
- www_imdb_com.html
- www_imdb_com.png
- www_imdb_com_links.txt
```






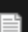
**Görsel 2:** go run main.go https://www.imdb.com komutunun çalıştırılması sonucunda HTTP 403 erişim kısıtı alınmasına rağmen dosya üretiminin devam ettiğini gösteren terminal çıktısı.

#### 4.3 Dosya Sistemi ve Çıktı Yönetimi

Her bir tarama işlemi sonucunda oluşturulan dosyalar, hedef siteye özgü şekilde isimlendirilerek proje dizinine kaydedilmektedir. Bu yaklaşım, veri bütünlüğünü korumakta ve CTI analiz süreçlerinde geriye dönük incelemeyi kolaylaştırmaktadır.

Üretilen Dosya Türleri:

- .html → Ham sayfa içeriği
- .png → Tam sayfa ekran görüntüsü
- \_links.txt → Ayırıştırılmış bağlantı listesi

	www_sans_org.html	20.12.2025 13:02	Chrome HTML Do...	439 KB
	www_sans_org.png	20.12.2025 13:02	PNG Dosyası	141 KB
	www_sans_org_links.txt	20.12.2025 13:02	Metin Belgesi	1 KB
	www_imdb_com.html	20.12.2025 13:00	Chrome HTML Do...	1 KB
	www_imdb_com.png	20.12.2025 13:00	PNG Dosyası	7 KB
	www_imdb_com_links.txt	20.12.2025 13:00	Metin Belgesi	0 KB

**Görsel 3:** Farklı hedef siteler için oluşturulmuş HTML, PNG ve bağlantı listesi dosyalarının site bazlı isimlendirme ve dizin yapısını gösteren dosya sistemi görünümü.

#### 4.4 Görsel Veri Analizi

Aşağıdaki görsel, chromedp motoru kullanılarak alınan tam sayfa ekran görüntüsünü göstermektedir. Bu çıktı, web sitesinin tarama anındaki gerçek kullanıcı görünümünü ve DOM yapısının başarıyla render edildiğini doğrulamaktadır.



**Görsel 4:** HTTP 200 durum kodu ile erişilen bir hedef sitenin, araç tarafından otomatik olarak oluşturulmuş tam sayfa ekran görüntüsü.

#### 4.5 403 Eriřim Kısıtı Altındaki Grsel ıktı Deęerlendirmesi

Bot algılama mekanizması bulunan platformlarda HTTP 403 yanıtı alınmasına rağmen, sayfa içerięinin belirli bir seviyede render edilebildięi ve grsel ıktının oluřturulabildięi gözlemlenmiřtir.

## 403 Forbidden

**Grsel 5:** HTTP 403 eriřim kısıtına sahip bir platformdan elde edilen tam sayfa ekran grnts. Bu ıktı, pasif istihbarat toplama aısından grsel kanıt nitelięi tařımaktadır.

#### 5.SONU VE DEęERLENDİRME

Geliřtirilen ara, test edilen 15 farklı web kaynaęının byk oęunluęunda HTML içerięi, grsel kanıt ve link haritası retme grevlerini bařarıyla yerine getirmiřtir. Elde edilen sonular, aracın CTI n keřif (initial reconnaissance) ařamasında etkin bir řekilde kullanılabileceęini gstermektedir.

Analiz srecinde zellikle arama motorları ve byk platformlarda **403 eriřim engelleri** ile karřılařılmıřtır. Bu durum, modern web altyapılarında bot tespit mekanizmalarının yaygın olarak kullanıldığını ortaya koymaktadır.

Bir sonraki geliřtirme ařamasında ařaęıdaki iyileřtirmeler kritik nemdedir:

1. **User-Agent rotasyonu** ile tarayıcı parmak izinin eřitlendirilmesi
2. **Headless tarayıcı tespitini azaltmaya ynelik nlemler**
3. Elde edilen linklerin ikinci seviye analiz (deep crawl) iin kullanılması

Bu geliřtirmeler ile ara, operasyonel CTI alıřmalarında daha etkin bir keřif platformuna dnřtrlebilir.