

불법 · 유해 사이트 차단 프로그램

인터넷 위의 작은 영웅, 당신의 안전한 길잡이 HitAnt

중부대학교 정보보안S/W융합전공 | 정진호 | 박우경 | 최수민 | 홍준희

1

프로젝트 선정 이유
프로젝트 목표

2

프로젝트 소개
프로젝트 진행 과정

3

프로그램 검증
프로그램 배포

4

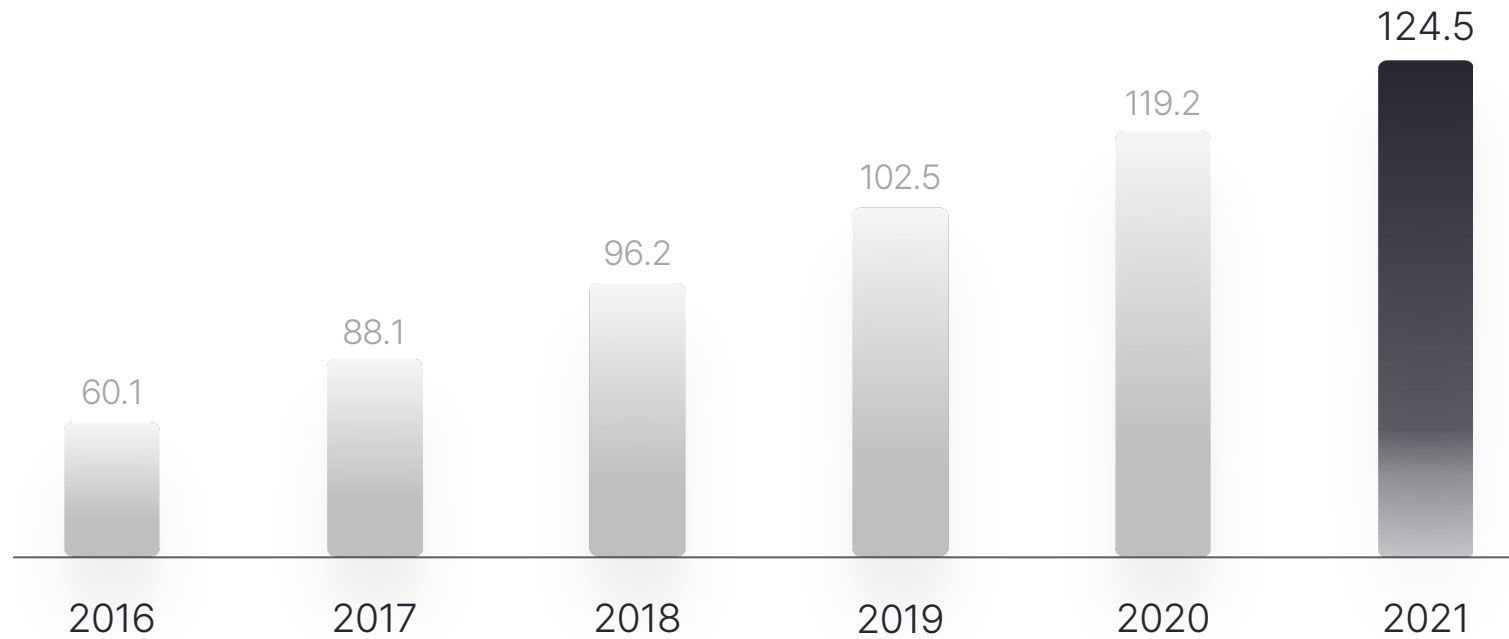
마무리

1

프로젝트 선정 이유 및 목표

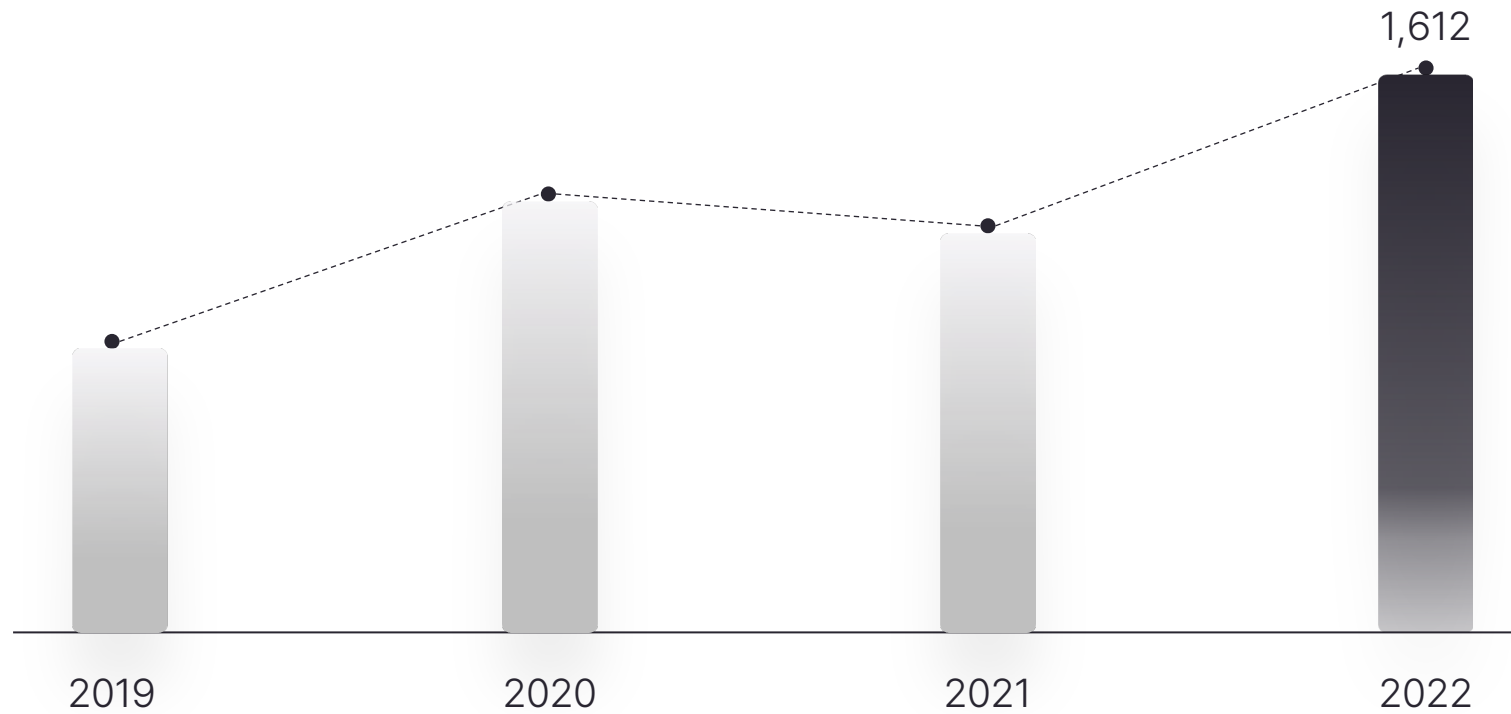
2021년 기준 K-콘텐츠 산업조사

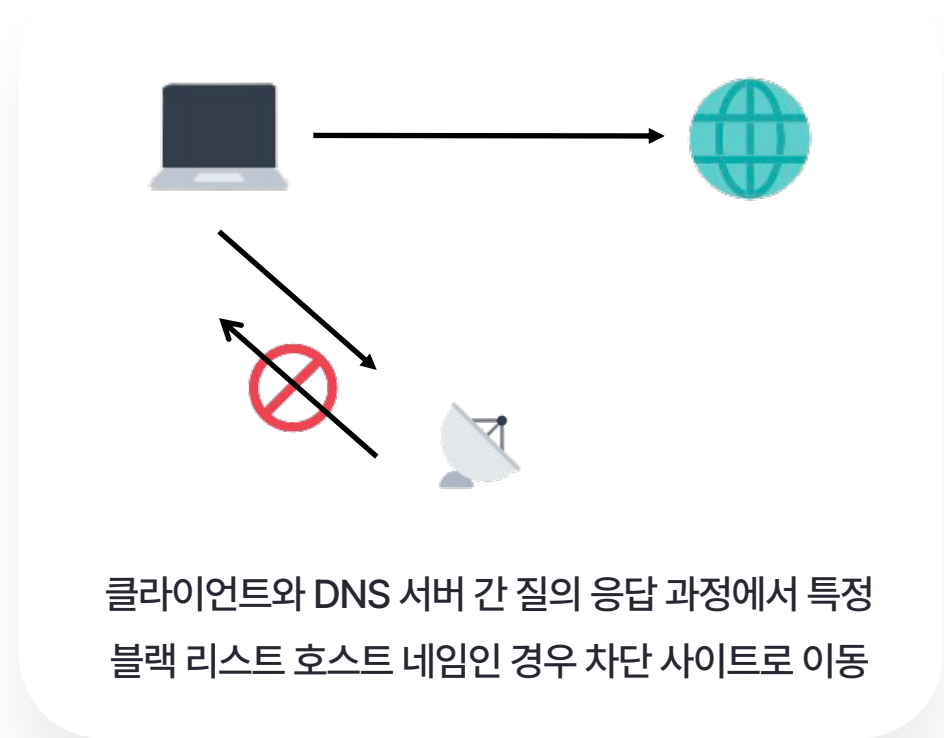
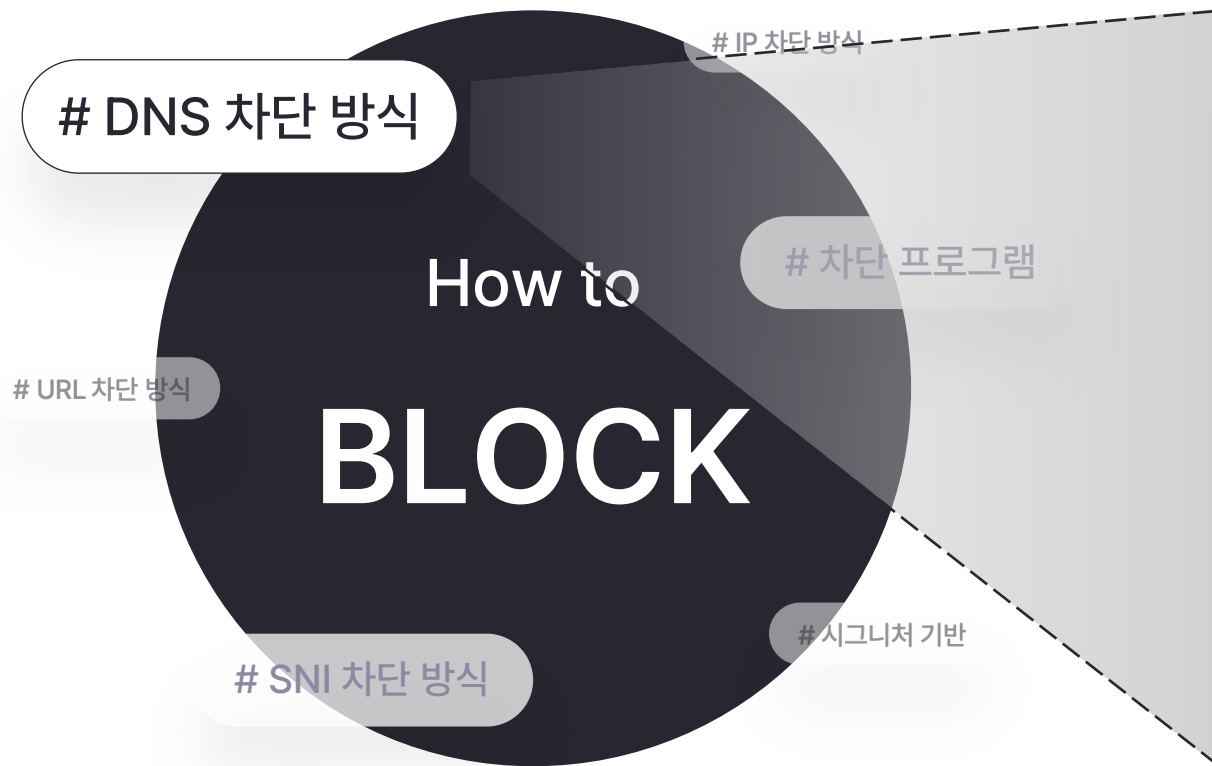
수출액 (단위: 억 달러)

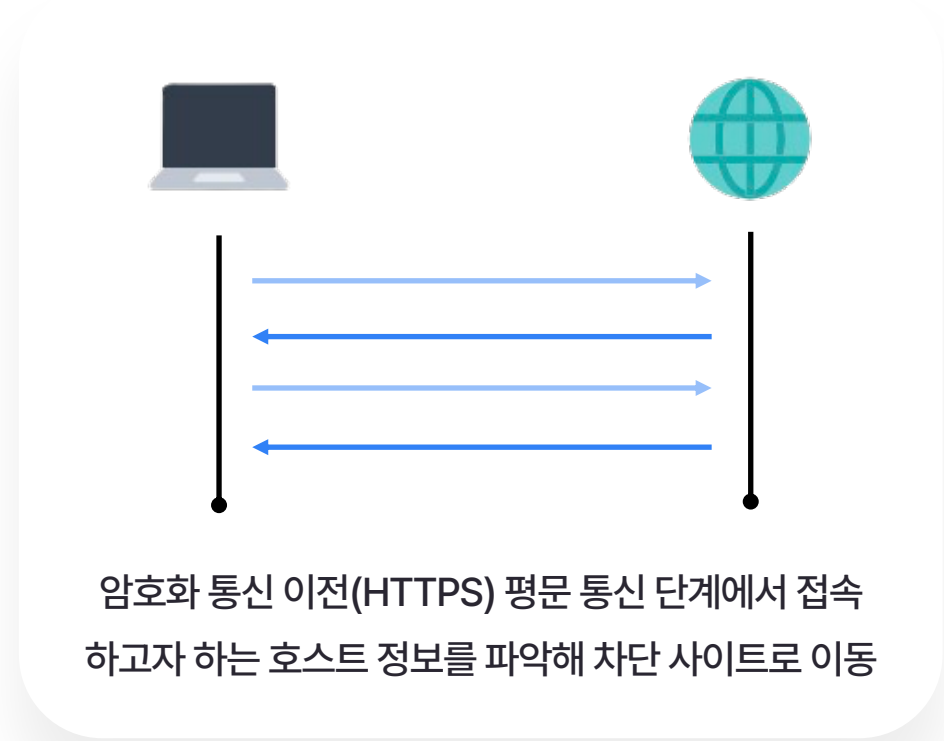
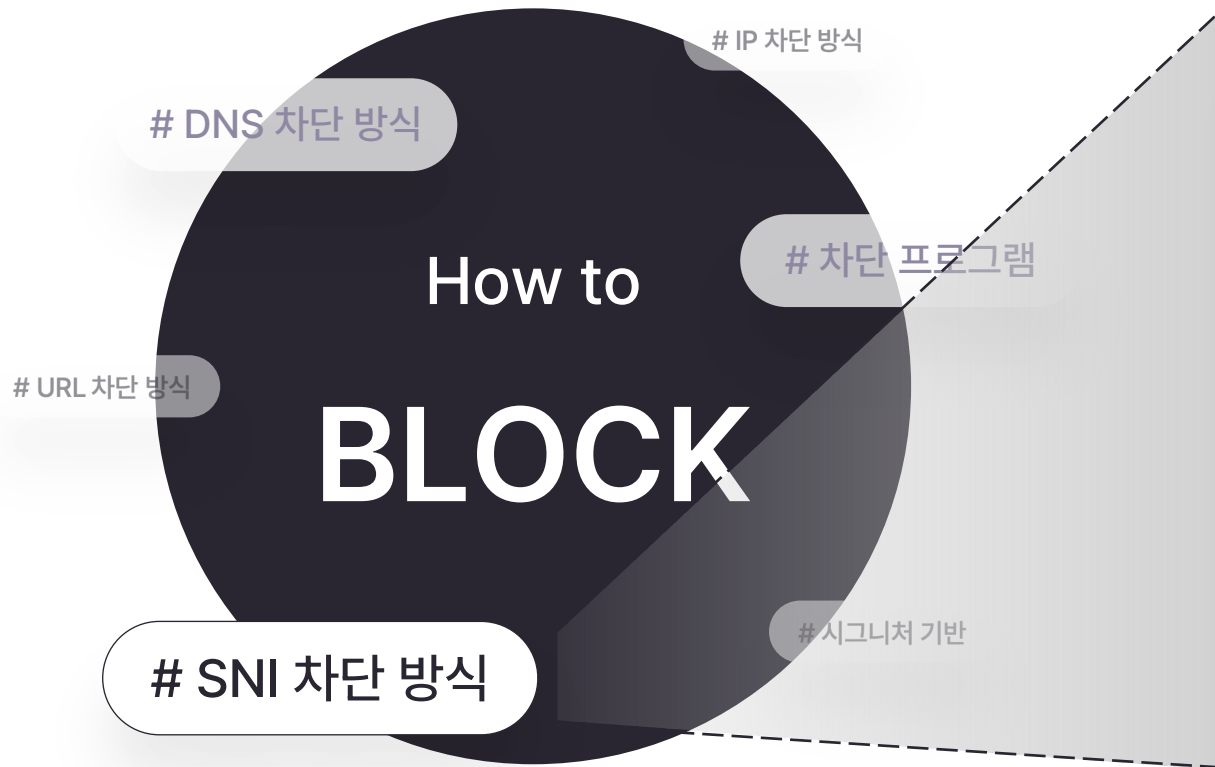


2023년도 불법 복제물 이용량

(단위: 천 개)







DNS 차단 방식

IP 차단 방식

How to

차단 프로그램

BLOCK

URL 차단 방식

SNI 차단 방식

시그니처 기반



자살·도박 정보 판치는데...
불법사이트 차단 규정 부족
2023. 10. 03.

청소년 27% 이상 불법 유해
사이트 접속

설날연휴특선영화 '서울의 봄'
누누티비로 볼 수 있나
2024. 02. 09.

K-웹툰 절반 침해... 해외 불법
사이트와 질긴 숨바꼭질
2024. 05. 01.

불법사이트 홍보하려 경북궁
낙서 테러...
2024. 05. 31.

누누티비 방지법 국회 본회의
통과... "접속차단 의무화"
2023. 12. 21.

스위트 홈 시청이 무료?
누누티비 시즌 2 버젓이 활개
2023. 12. 12.

규제의 필요성
차단에 대한 규정 부족 문제

접속차단 의무화
그러나, 현실적인 문제 해결 X

K-콘텐츠의 위협
불법·유해 사이트를 통한 유포

1

링크 모음 사이트 및 연관 링크 사이트

👤 웹문	🎬 무료드라마	19 성인사이트	19 오피/유흥	🏆 스포츠중계
🔥 문고	🔥 티베위키	🔥 AVseeTV	🔥 오피가이드	🔥 볼락TV
🔥 불역문	🔥 티베온	🔥 MissAV	🔥 부산달리기	🔥 볼TV
🔥 뉴토크	🔥 HOOHOO TV	🔥 AV19	🔥 오피스타	🔥 굿라이브TV
④ 조아문	④ TVHOT	④ AV쏘걸	④ 오피아트	④ 놀고가닷컴
⑤ 북도문	⑤ 다이소티베	⑤ 아플코리아	⑤ 아이러브밤	⑤ VIP TV
⑥ 아지문	⑥ 콕콕티비	⑥ AV탐걸	⑥ 대방	⑥ 각TV
⑦ 색문	⑦ 마징가VOD	⑦ 다크걸	⑦ 건마바다	⑦ 코난TV
⑧ 아문	⑧ 소나기티베	⑧ 아플공감	⑧ 편초이스	⑧ 올림픽TV
⑨ 편비	⑨ 무비킹	⑨ 아플란	⑨ 건마시티	⑨ 해골TV
🔥 뉴문	🔥 TV다시보자	🔥 조계로티	🔥 돌거운달리기	🔥 쉐스티비
👤 먹튀검증	🏠 토렌트	🌐 커뮤니티	19 성인용품	🇰🇷 인인교인
🔥 슈어맨	🔥 토렌트큐큐	🔥 디시인사이드	🔥 바나나물	🔥 [미국] 뉴욕코리아
🔥 먹튀검개소	🔥 토렌트씨	🔥 보배드림	🔥 조어맨조어	🔥 [오주] 코리아타운
🔥 먹튀검증소	🔥 토렌트원	🔥 일간베스트	🔥 토렌스물	🔥 [일본] 제팬인포
④ 배팅노리	④ 토렌트알지	④ 불부	④ 나이트물	④ [태국] 타이홀릭
⑤ 먹튀안내소	⑤ 토다와	⑤ 게드림	⑤ 약광	⑤ [영국] 영국사랑
⑥ 토토닥터	⑥ 토렌트힐	⑤ 인벤	⑤ 카미그라	⑥ [러시아] 코리아스
⑦ 토프세이	⑦ 토렌트탑	⑦ 루리웹	⑦ 물물삼	⑦ [캐나다] 한인광터
⑧ 토토군	⑧ 색토렌트	⑧ 해플코리아	⑧ 링크백스삼	⑧ [베트남] 배한타임즈
🔥 에방어테	🔥 토렌트킹	🔥 플리양	🔥 전바나나물	🔥 [네덜란드] 데일라

2

HTML의 a Tag 속 링크 크롤링

3

수집한 데이터 셋 정보

1

링크 모음 사이트 및 연관 링크 사이트

2

HTML의 a Tag 속 링크 크롤링

The image shows a screenshot of the Newtoki website, which is a link aggregation site for various online casinos. The website displays numerous promotional banners for different casinos, each with details about bonuses, deposit requirements, and game types. To the right of the website screenshot, the browser's developer tools are open, showing the HTML source code. The code highlights the `<a href=...` tags for each promotional link, demonstrating how the links are structured in the HTML.

3

수집한 데이터 셋 정보

1

링크 모음 사이트 및 연관 링크 사이트

2

HTML의 a Tag 속 링크 크롤링

3

수집한 데이터 셋 정보

```

URL
1 https://www.naver.com
2 https://www.daum.net
3 https://www.google.com
4 https://comic.naver.com/index
5 https://webtoon.kakao.com
6 https://www.netflix.com/kr/
7 https://page.kakao.com
8 https://www.lezhin.com/ko
9 https://www.toomics.com
10 https://www.mrbblue.com
11 https://bufftoon.plaync.com
12 https://tooptoon.com
13 https://www.myktoon.com/
14 https://www.tving.com/
15 https://www.coupangplay.com/
16 https://www.wavve.com
17 https://www.disneyplus.com/ko-kr
18 https://watcha.com
19 https://ko-kr.facebook.com
20 https://www.instagram.com
21 https://www.threads.net/
22 https://www.youtube.com/
23 https://www.tiktok.com/ko-KR/
24 https://section.blog.naver.com/
25 https://www.tistory.com
26 https://www.notion.so/
27 https://www.pinterest.co.kr
28 https://twitter.com/
29 https://www.afreecatv.com
30 https://chzk.naver.com
31 https://news.naver.com
32 https://www.coupang.com
33 https://www.ioongbu.ac.kr
34 https://everytime.kr
35 https://wrt9.ai
36 https://chat.openai.com
37 https://www.hancomdocs.com/
38 https://top.cafe.daum.net
39 https://cafe.naver.com/
40 https://finance.naver.com
41 https://search.naver.com/
42 https://papago.naver.com/
43 https://edu.ioongbu.ac.kr
44 https://aftel.net
45 https://www.mrbblue.com
46 https://page.kakao.com/
47 https://www.yna.co.kr

```

비유해 사이트

```

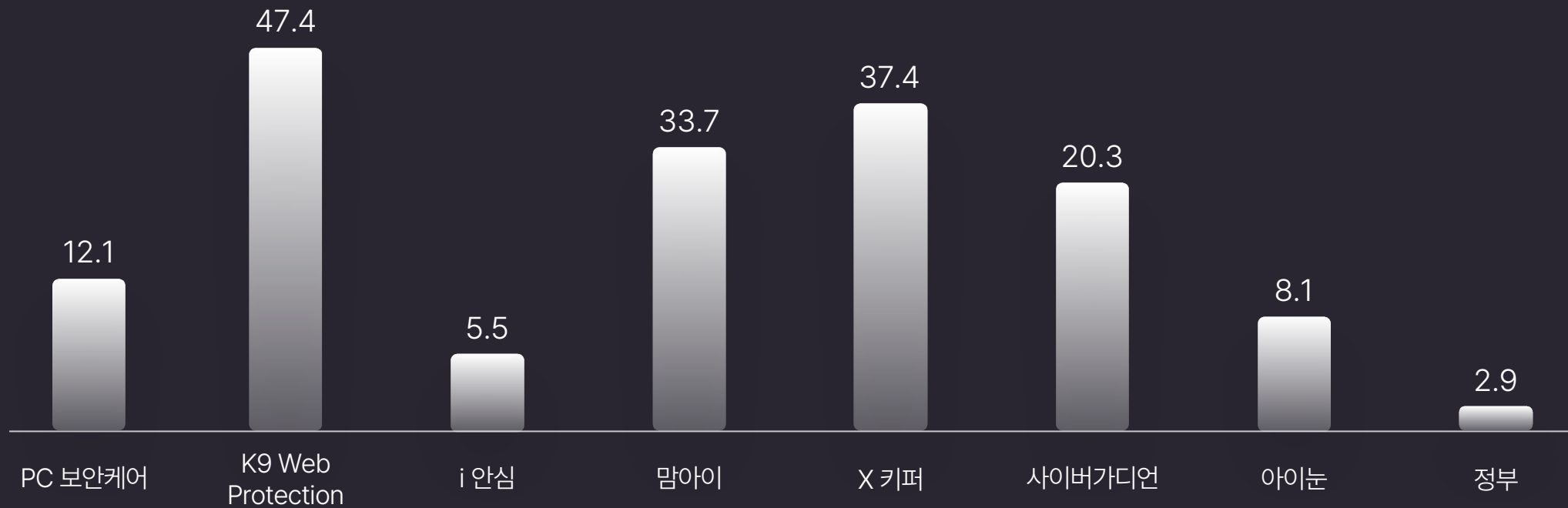
URL
1 http://1p-888.com/
2 http://블링주소.com/
3 http://이프로넷.com/
4 http://인디넷주소03.com/
5 http://b-time511.com/
6 http://bbox1212.com/
7 http://bn2222.com/#/
8 http://bp-cc.com
9 http://bwd07.com/
10 http://bz-0101.com/
11 http://cg-mvp10.com/
12 http://dada-777.com/front
13 http://fst-234.com/
14 http://fsw-333.com
15 http://fsw-555.com/
16 http://fw-0001.com/
17 http://gcity-222.com/
18 http://ggb-333.com/
19 http://gn1020.com
20 http://jws42.com/
21 http://h-two41.com/
22 http://jg-5555.com/login.asp
23 http://ka-01.com/login
24 http://jg-1010.com
25 http://jg0011.com
26 http://mmcc234.com/
27 http://mz1-one.com/
28 http://pgss1122.com
29 http://k9284.com/
30 http://site99.com/
31 http://spo-one.com/
32 http://spst-1111.com
33 http://sun-9909.com/?regcode=4557
34 http://we-324fg.com/
35 http://woori-333.com/main
36 http://www.fu2024.com/#/
37 http://www.gbqx93.com/#/
38 http://www.mxr-36.com/#/
39 http://www.pnt357.com/front
40 http://www.snx91.com/#/
41 http://www.solslotgm111.com/#/
42 http://www.xvs-49.com/#/
43 http://17b-00.com/
44 https://365kor.bet/
45 https://888-sm.com/?ref=4545

```

유해 사이트

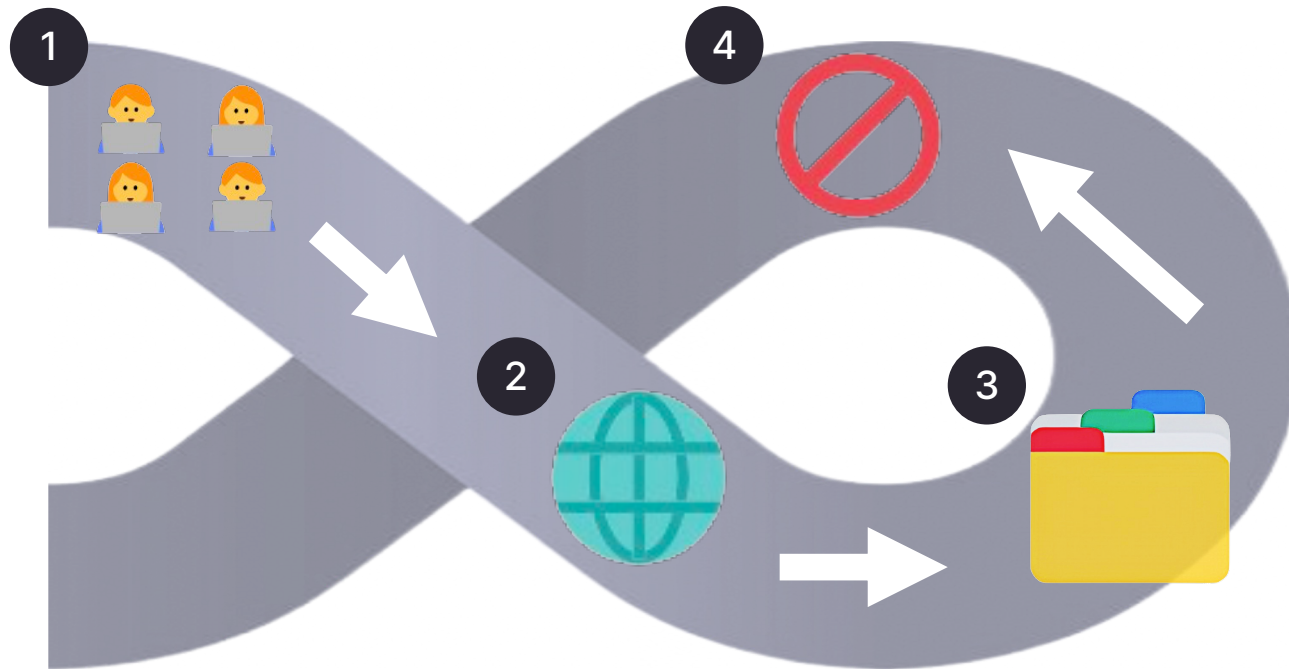
현 차단 프로그램의 차단율

(단위: %)



시그니처 방식

수동적 차단 방식



1

사용자

3

불법 · 유해 데이터베이스

2

인터넷 접속

4

차단 유무 판정

주기적인

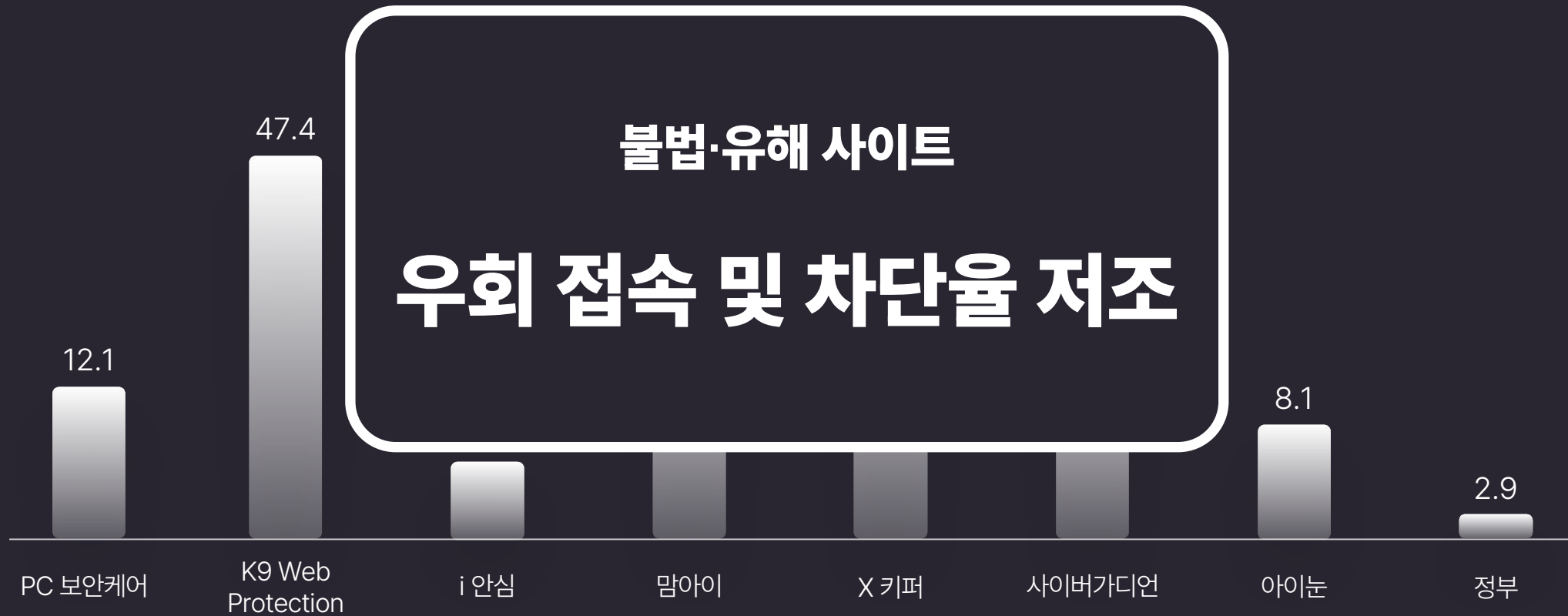
**데이터 베이스
업데이트 필수**

불법 · 유해 사이트

**도메인 주소 변경 시
차단 회피**

현 차단 프로그램의 차단율

(단위: %)



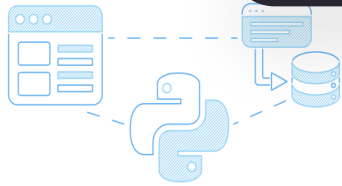
프로젝트 목표



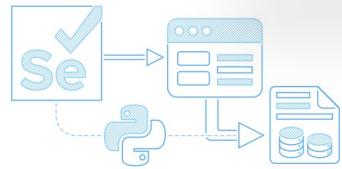
2

프로젝트 소개 및 진행 과정

requests



selenium



HTML



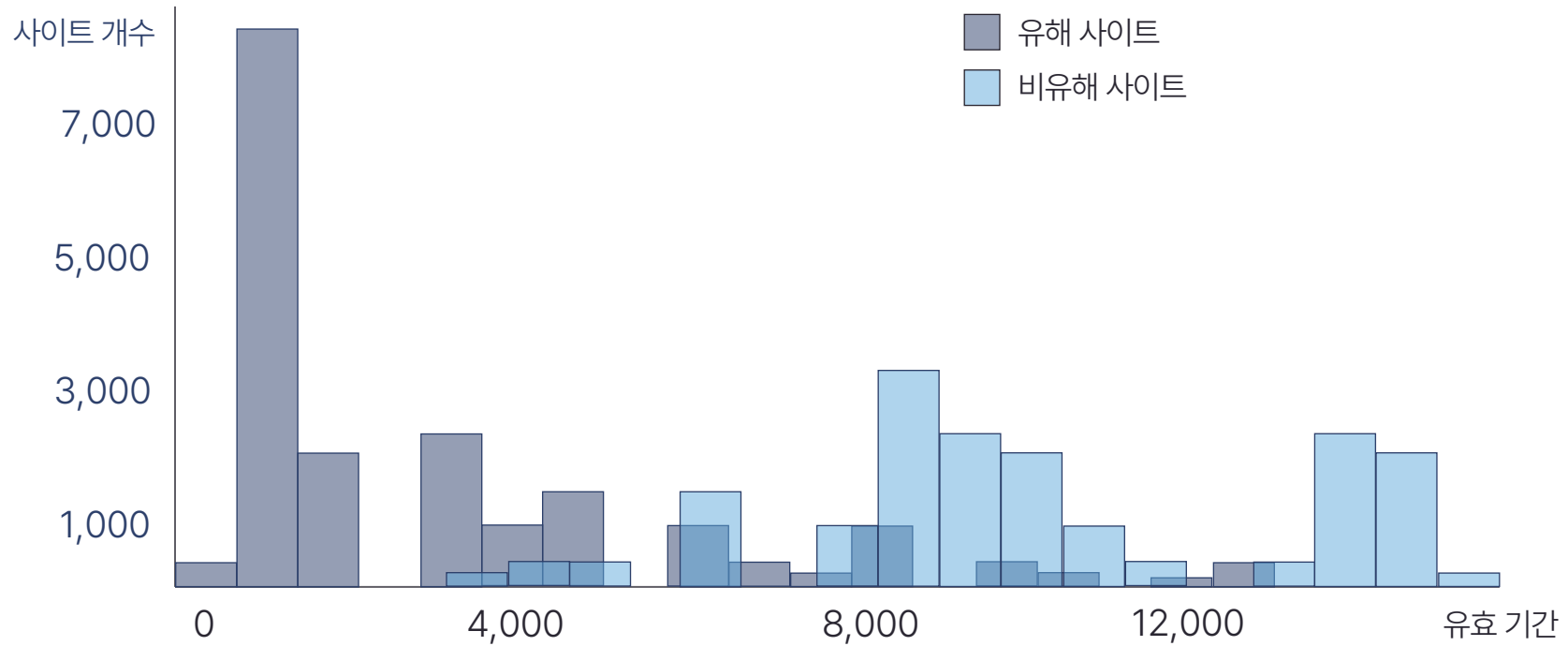
피쳐 선정

24가지

연결 상태	로딩 시간	주소 길이	생성 기간	최종 주소	리다이렉트
HTTPS	Meta	A	IP	Title	Div
Popup Count	Js_Len	Js_density	Link	Country	Label
City	Img	Img_density	Tld	SSL	Registrar

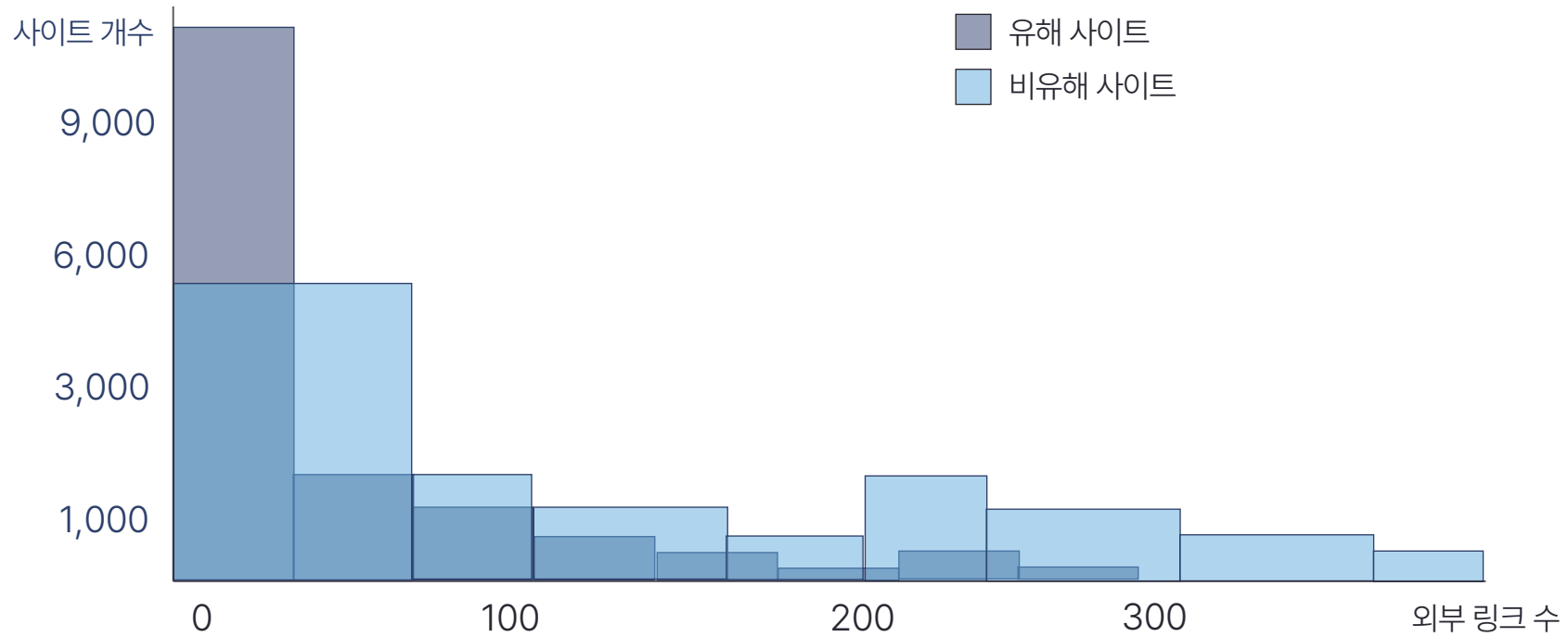
유효 기간 분포 비교

(단위: 개)



외부 링크 수 분포 비교

(단위: 개)



IP	Js_Len	외부 링크 수	도메인 생성일	도메인 마감일	Div 수
213.166.69.26	7391	0	2022-02-13 15:58:13	2025-02-13 15:58:13	25
104.21.17.249	2074	83	2023-08-30 10:09:23	2024-08-30 10:09:23	14
167.114.35.174	9789	5	2021-07-12 16:05:39	2024-07-12 16:05:39	118
104.21.59.9	5108	57	2023-08-21 00:10:21	2024-08-21 00:10:21	27
104.21.61.222	22634	18	2005-02-23 17:07:11	2025-02-23 05:00:00	617
172.67.172.38	79591	11	2002-04-19 22:42:04	2027-04-19 22:42:04	1135
172.67.165.12	76516	11	2020-10-30 09:36:57	2025-10-30 09:36:57	1138
172.67.163.235	5135	57	2023-08-21 00:10:08	2024-08-21 00:10:08	27
104.21.65.109	5010	59	2023-08-21 00:10:34	2024-08-21 00:10:34	30
104.21.19.16	1162	2	2024-01-19 04:24:53	2025-01-19 04:24:53	333

팝업 수	이미지 수	로딩 시간	리다이렉트	Js_density	Img_density
0	7	0.00297904	0	0.103187346	0.28
0	4	0.002000093	0	0.036336878	0.285714286
1	21	0.003004789	0	0.101218049	0.177966102
0	2	0.002022743	0	0.088093267	0.074074074
0	25	0.001998186	0	0.157867939	0.040518639
0	229	0.002015591	0	0.099017181	0.201762115
0	248	0.001999855	0	0.094648587	0.217926186
0	2	0.001999617	0	0.088368411	0.074074074
0	2	0.00199914	0	0.068043841	0.066666667
1	47	0.002000332	0	0.017107859	0.141141141

Numeric

전처리

Int

Float



별도의 전처리 없이
머신러닝 학습 데이터로
바로 이용 가능

연결 상태	Meta	A	Title	Tld	Registrar
False	[]	[]		ca	0
True	['lbp8V0h90drtN...']	['#site-navigation'...]	Wildblaster Casino	com	NAMECHEAP INC
True	['X_uHNqglEIAWnl...']	['https://adappgeo'...]	M11HKB : Raja...	net	NAMECHEAP INC
True	['width=device-wid...']	['#content', '#'...]	A Beginner's Guide t...	com	Dynadot Inc
True	['3p0Q0q3gFDaG'...]	['https://basahter'...]	PAPASLOT : Link...	org	NameCheap, Inc.
True	['width=device-wid...']	['#fl-main-content'...]	Home – Canadian...	ca	GoDaddy Domains Canada, Inc
True	['text/html; charset...']	['https://casinome'...]	Casinos Online...	com	Key-Systems GmbH
True	['width=device-wid...']	['https://casinome'...]	Casinos Online...	com	Key-Systems GmbH
True	['text/html; charset...']	['https://casinosc'...]	Online Casino...	com	NAMECHEAP INC
True	['Ps9Z0Mp9ayNFm...']	['https://basahter'...]	AGEN123 : Link Raja...	org	NAMECHEAP INC

TF-IDF / One-Hot 인코딩

String

Konolpy | Kss



문자열 속 문장 중
단어만 추출,
단어를 숫자형으로
변환



머신러닝, XGBoost

데이터 학습 과정에서 과적합 방지
오차 데이터만 반복 학습을 통해 정확도 향상
각 카테고리 성인, 도박, 웹툰, 스트리밍에 따라
우수한 정확도 제공

카테고리 별
차단율 향상

데이터 학습
과적합 방지

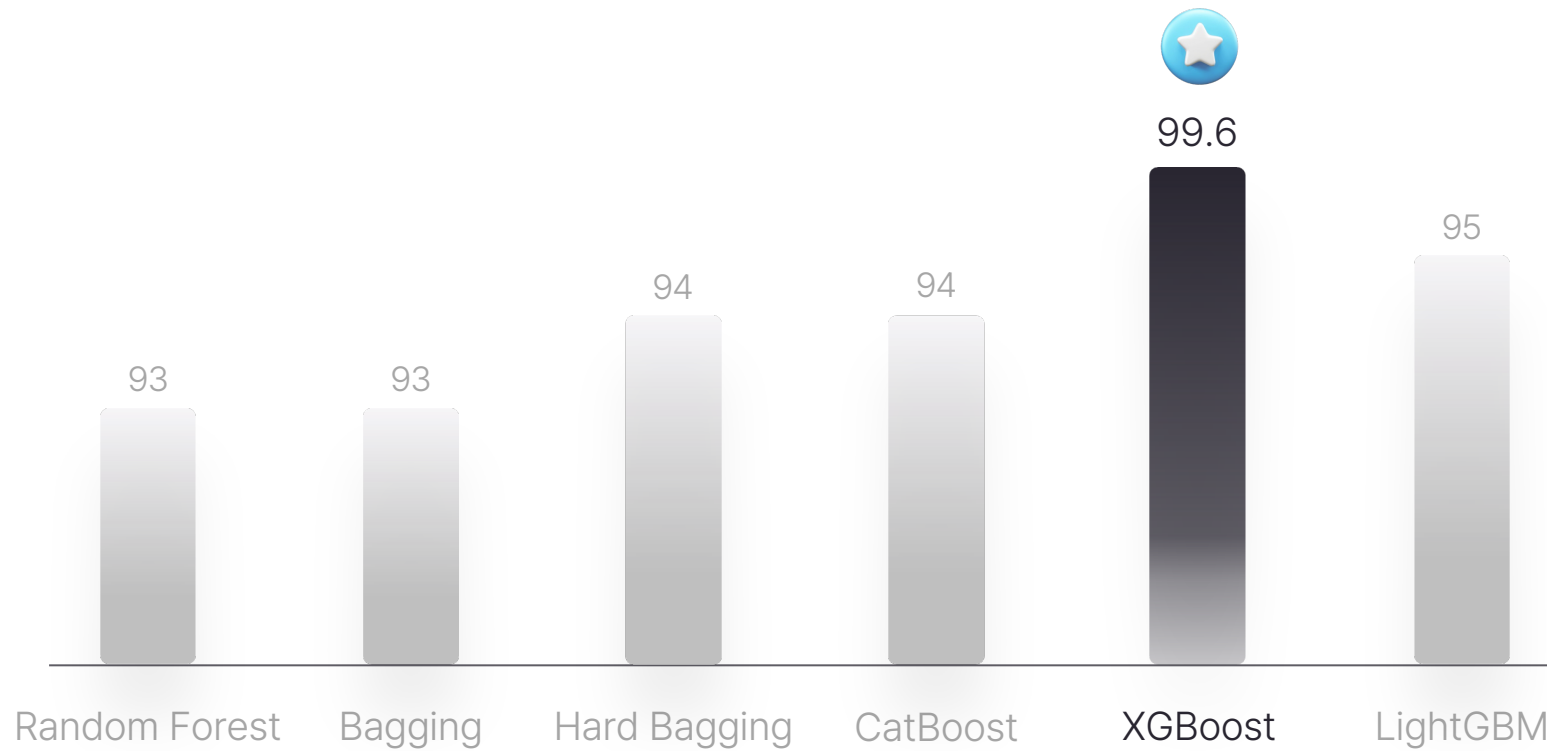
오차 학습
정확도 향상

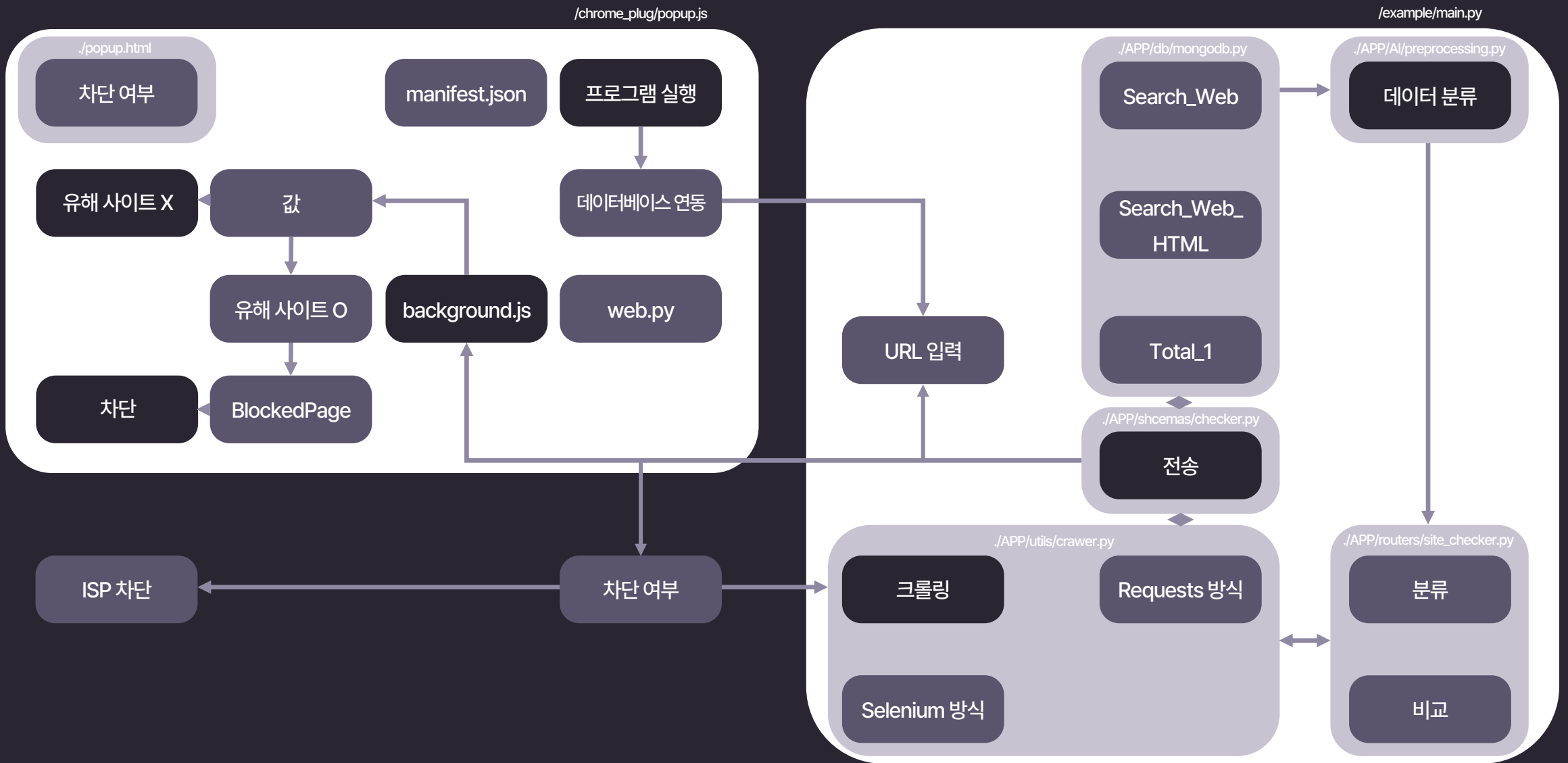
타 머신러닝보다
우수한 정확도

XGBoost

머신러닝 학습 정확도

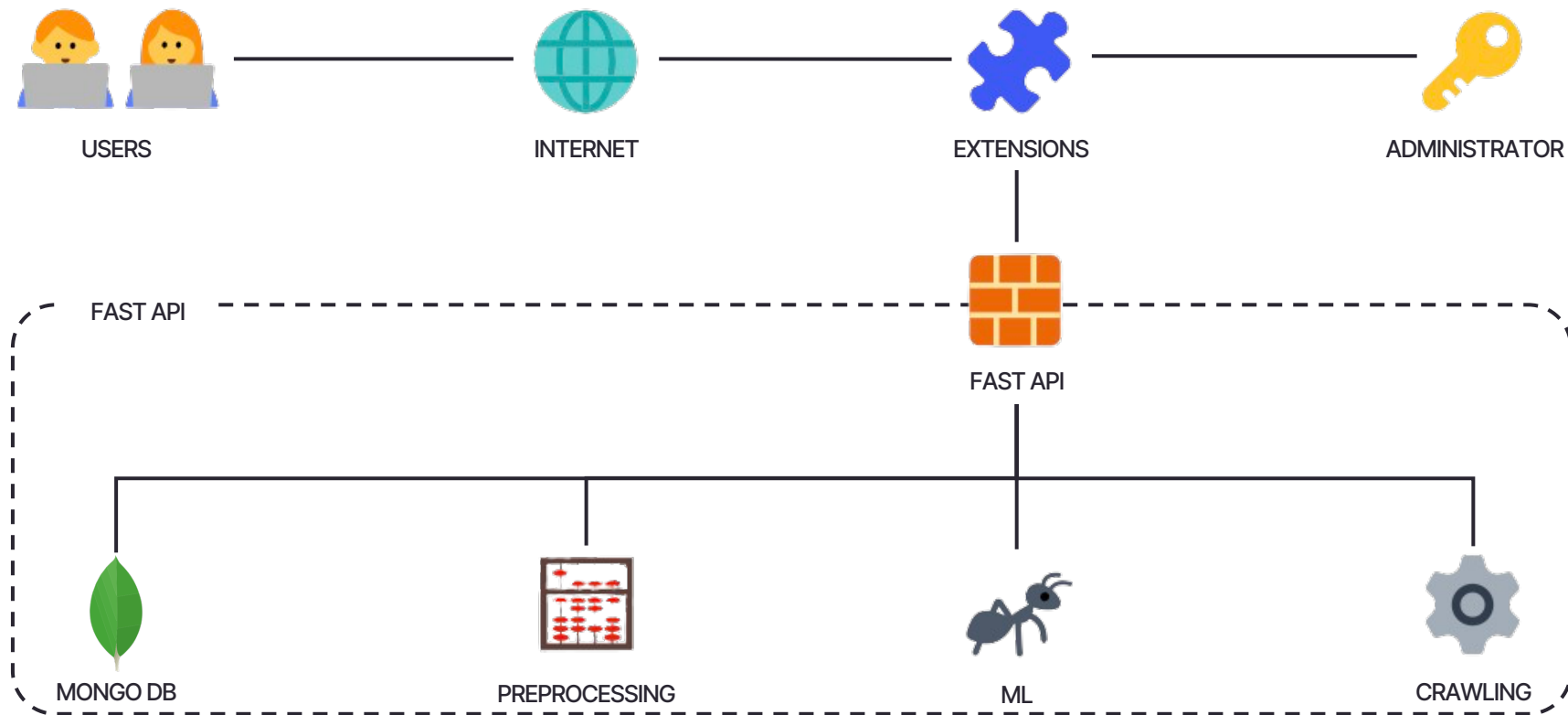
(단위: %)





불법·유해 사이트 차단 프로그램

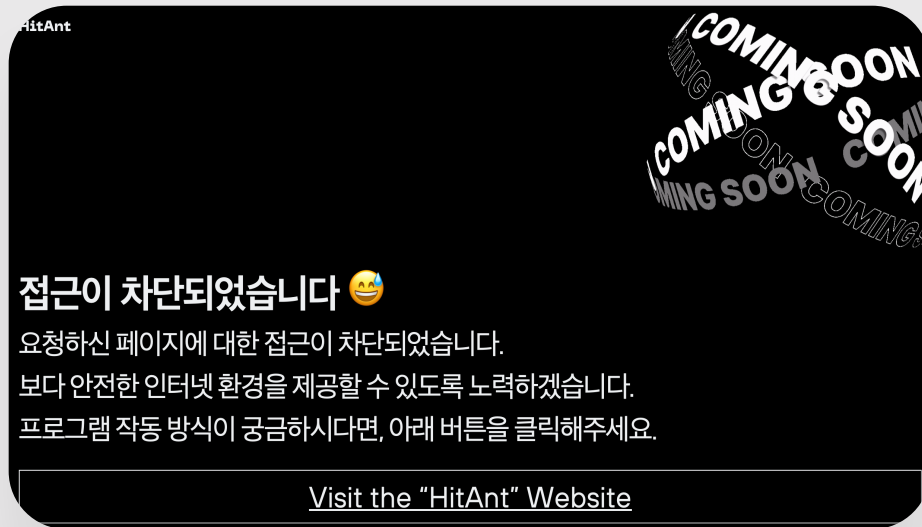
Flow Chart





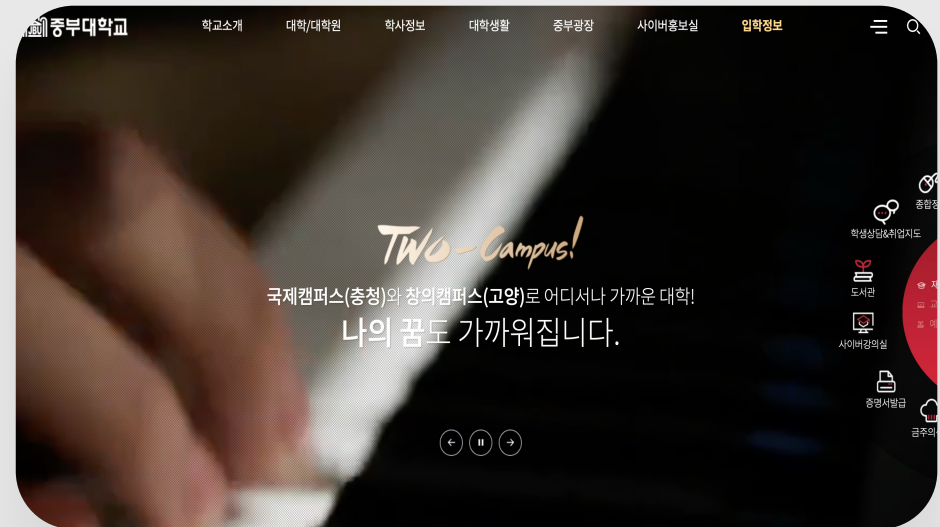
불법·유해 사이트 판정

차단 페이지로 리다이렉트



비유해 사이트 판정

정상적으로 이용 가능

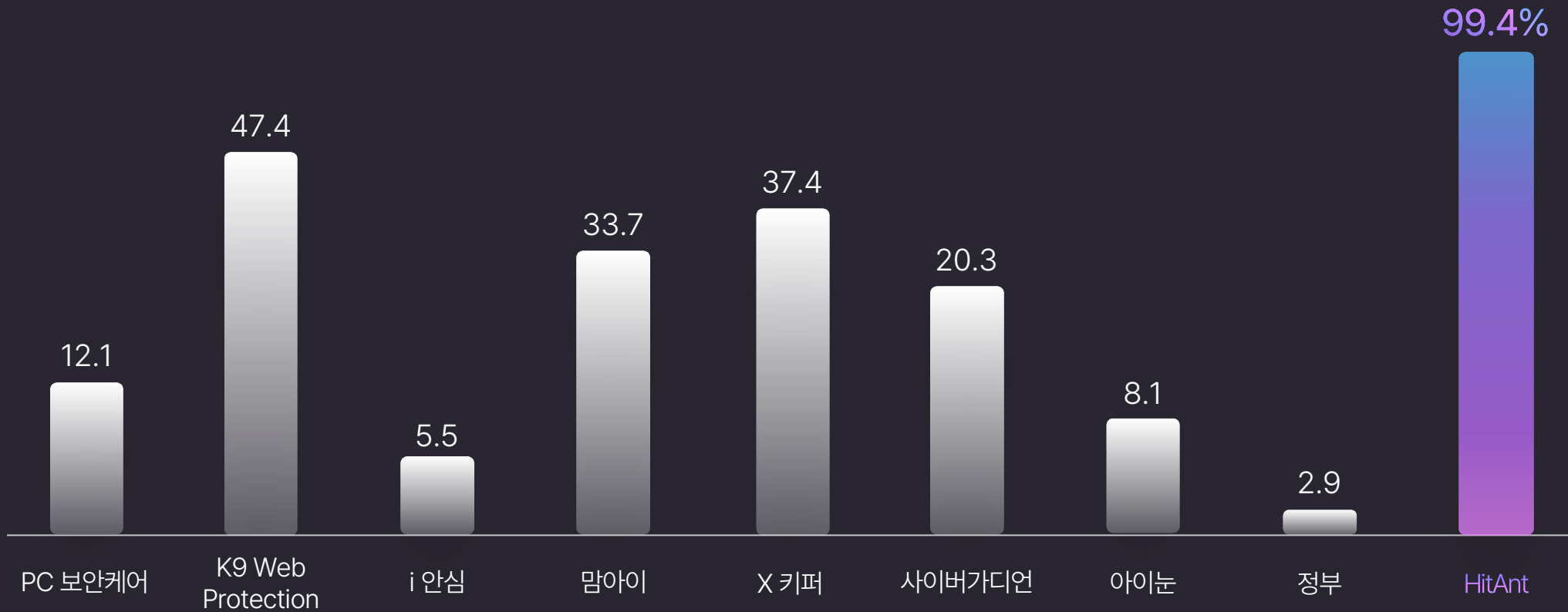


3

프로그램 검증 및 배포

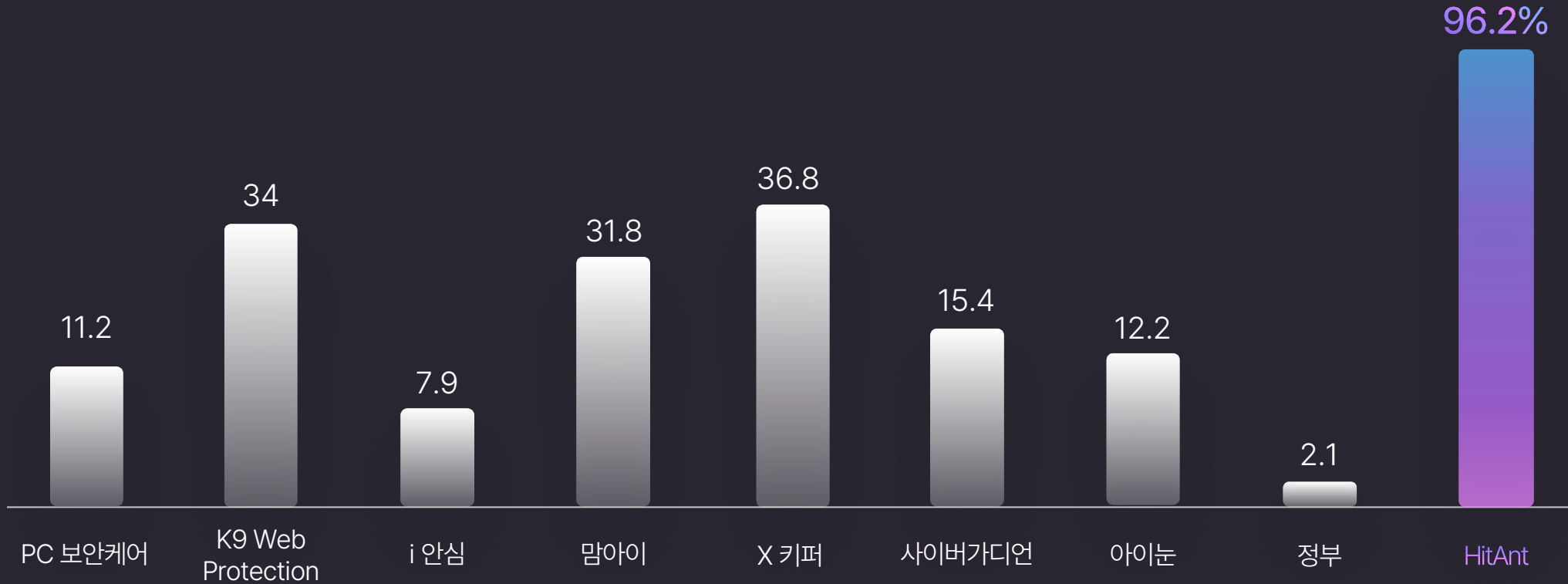
현 차단 프로그램의 차단율

(단위: %)



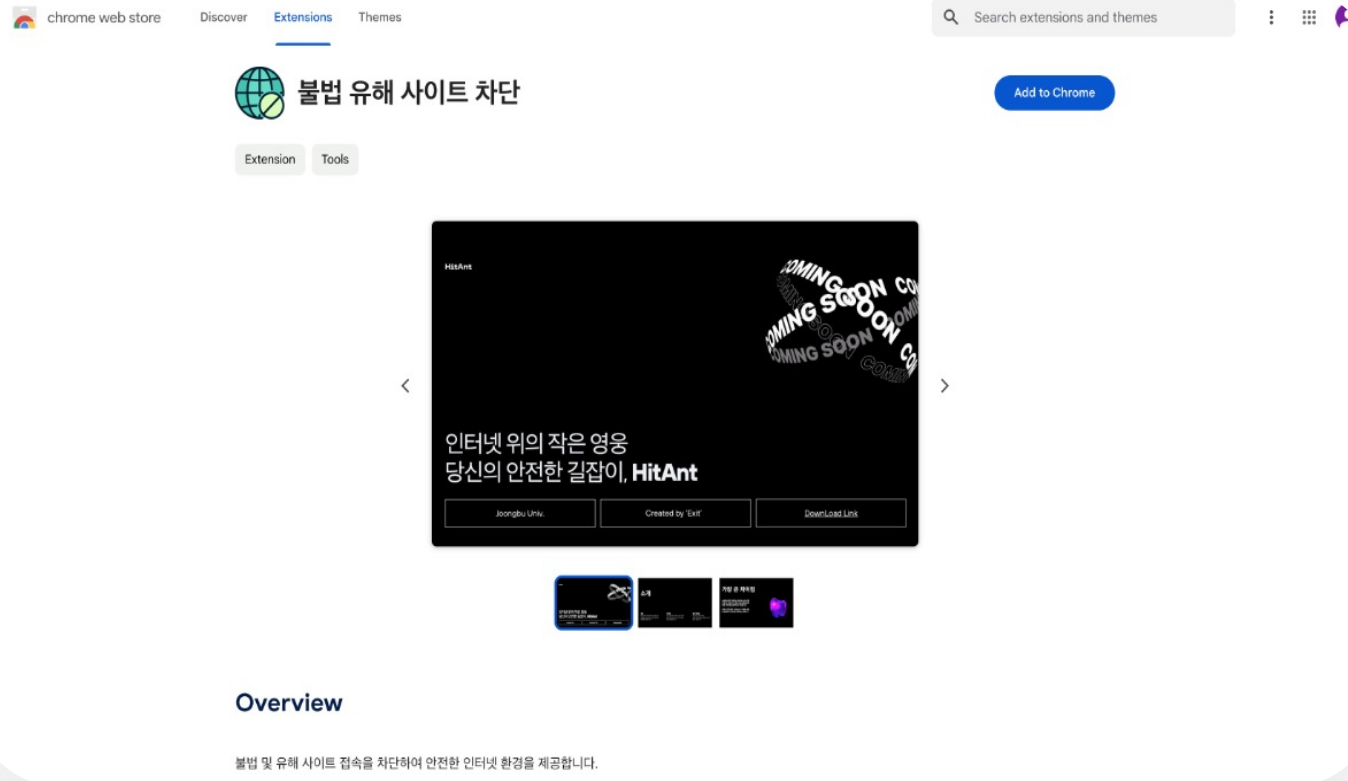
현 차단 프로그램의 차단율

KISA 데이터 자료 (단위: %)



불법 · 유해 사이트 차단 프로그램 배포

크롬 확장자 플러그인 프로그램



불법 유해 사이트 차단

실행

접속 중인 사이트의 리다이렉션 확인

리다이렉션이 감지되지 않았습니다.

차단하고 싶은 URL 추가하기

URL 입력...

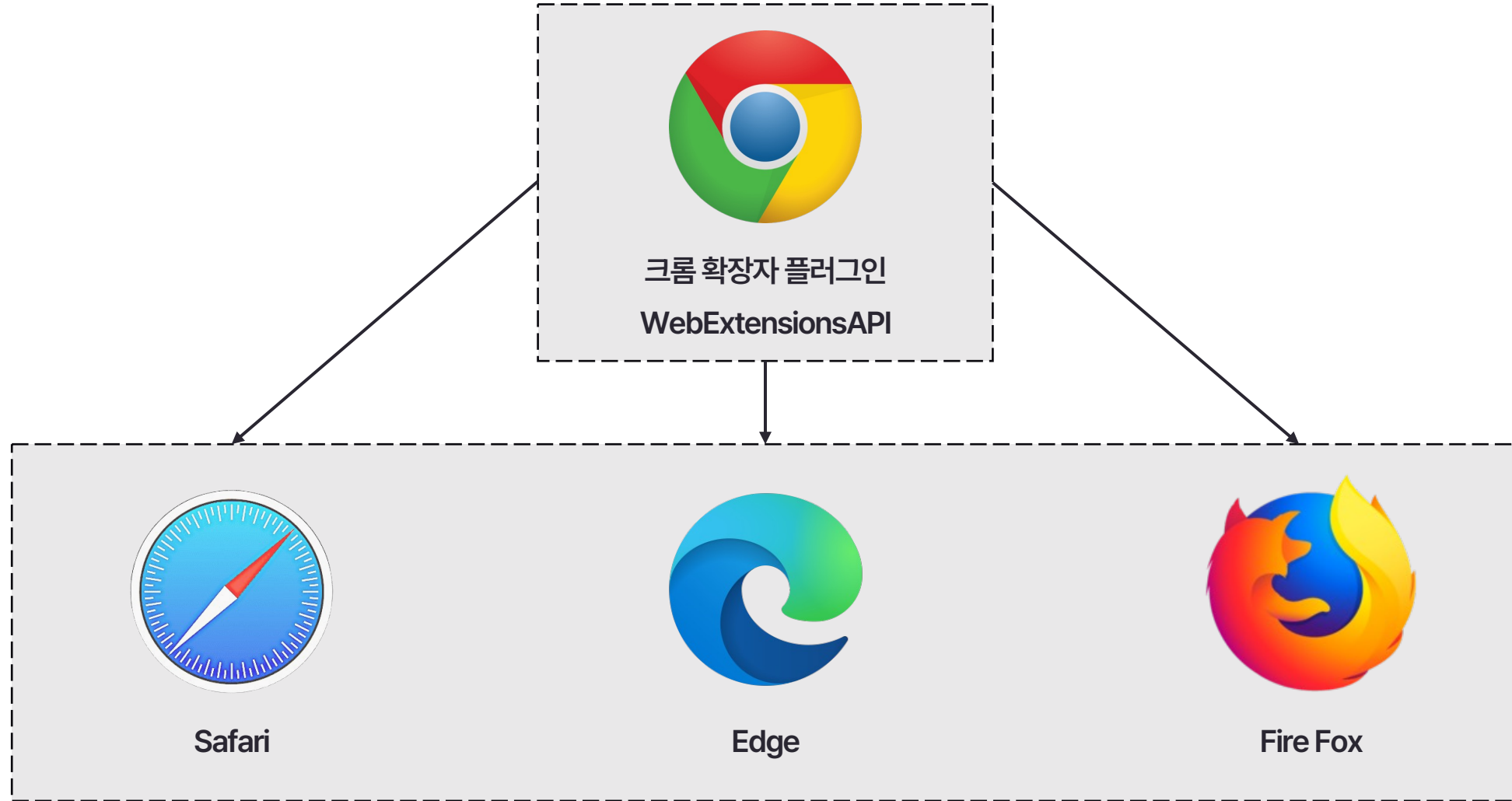
차단하지 않고 싶은 URL 제거하기

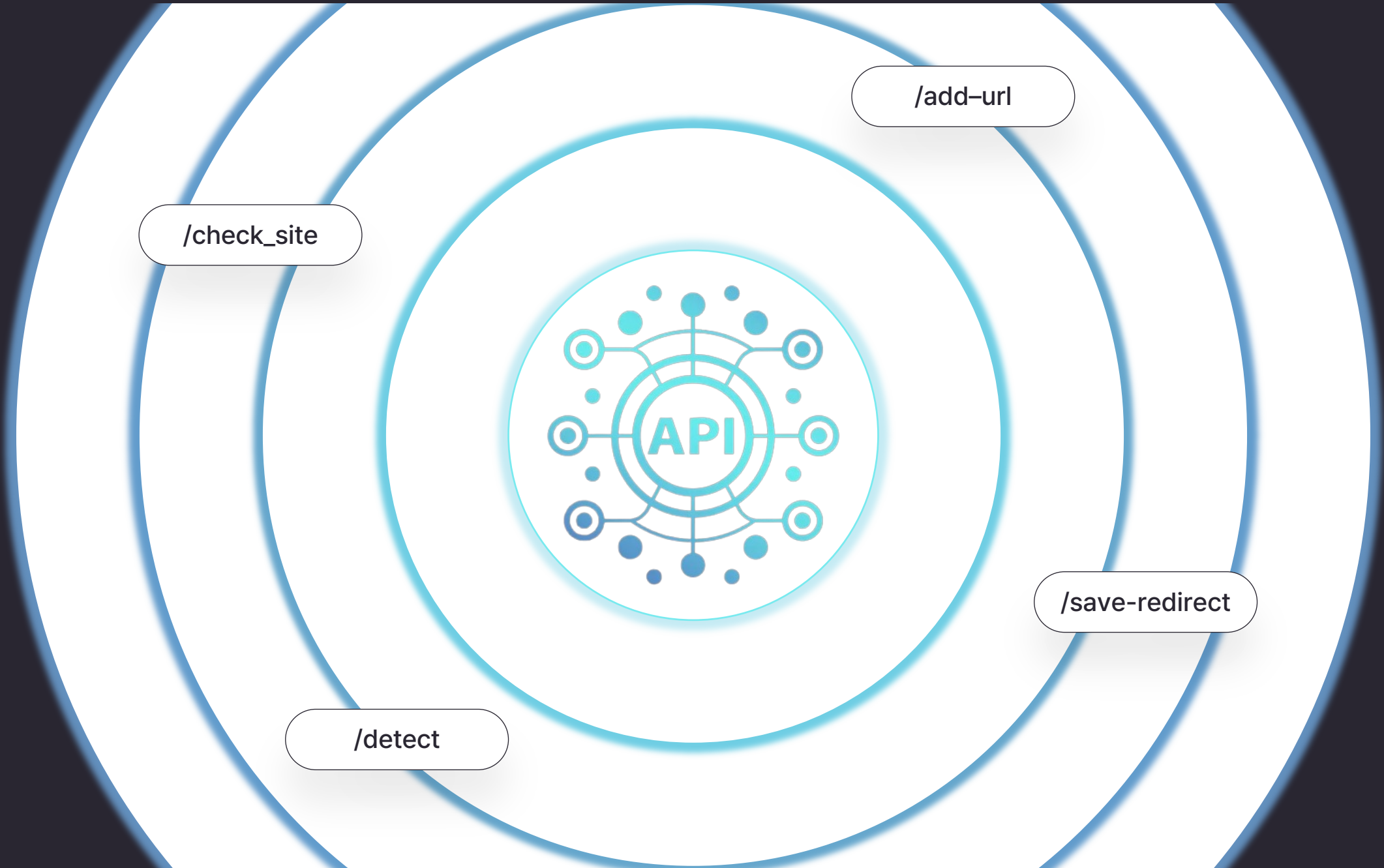
URL 입력...

[크롬 확장자 스토어 링크](#)

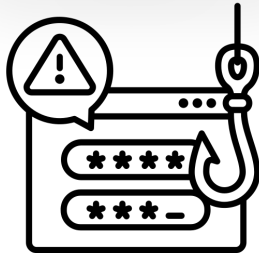
4

마무리





/detect



피싱 사이트 탐지 [전]

피싱 사이트 탐지

현재 안전한 사이트입니다!

피싱 사이트 탐지 [후]

피싱 사이트 탐지

현재 피싱 사이트에 접속하셨습니다!

구글 API 활용한 피싱 탐지






Google Cloud Platform

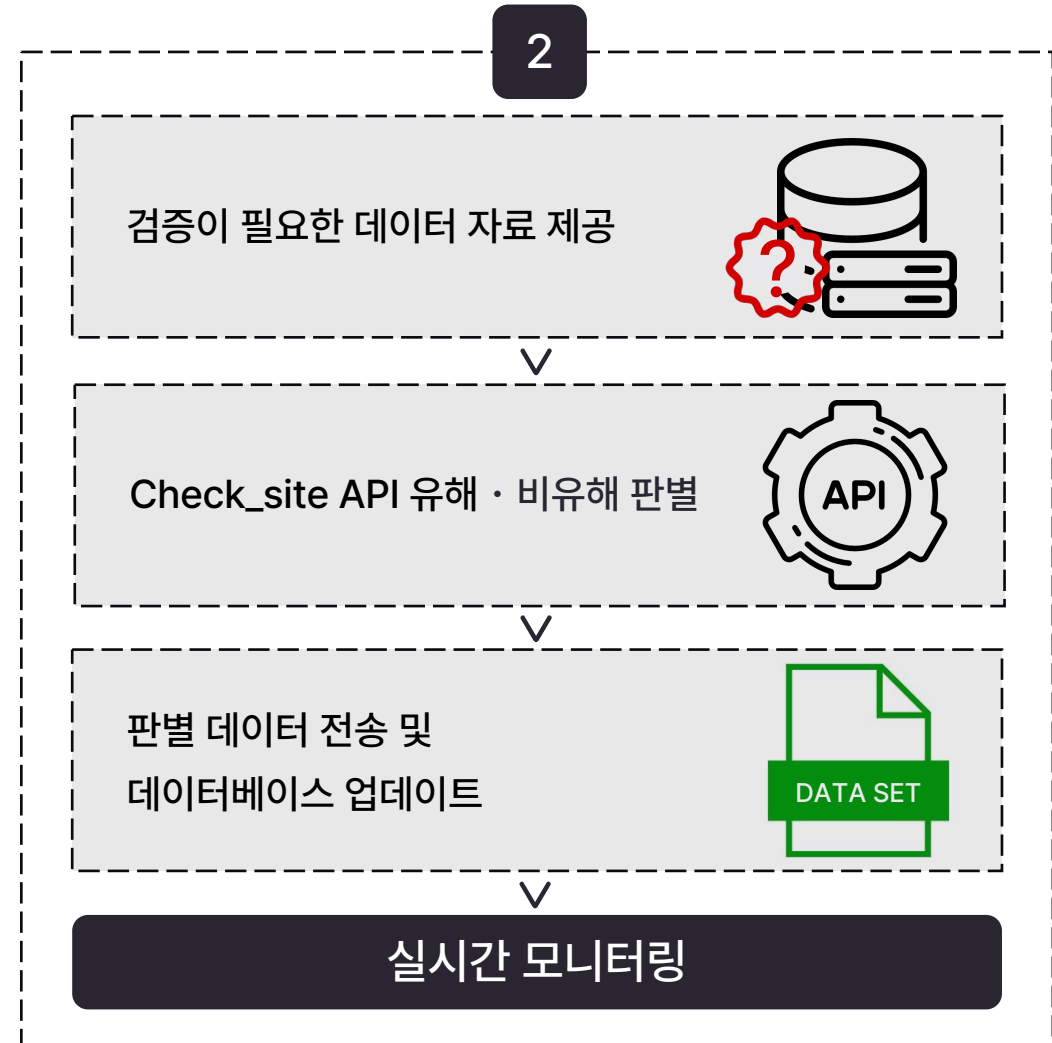
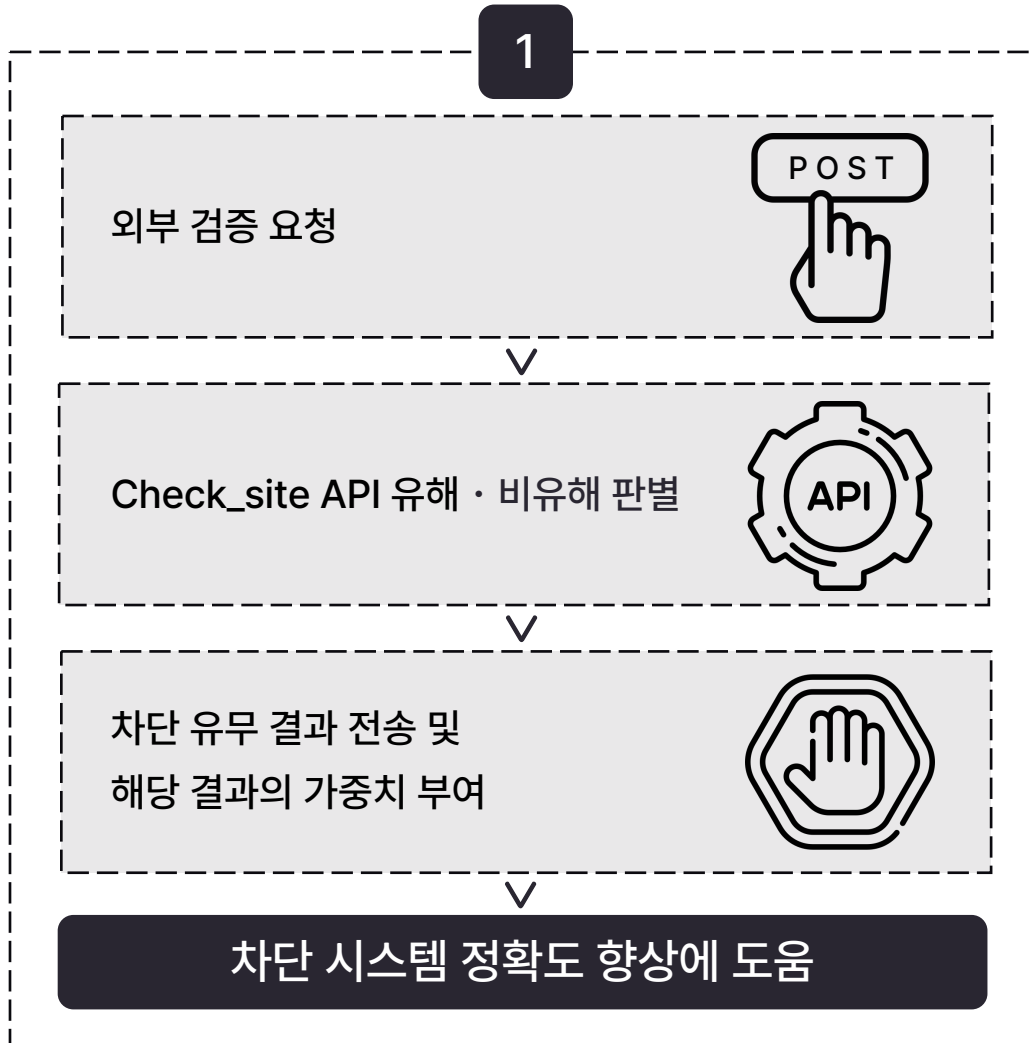
정규 표현식을 이용한 피싱 URL 탐지

/Reg[ex]+/

머신러닝 비교 전 실행



<p>/check_site</p> <p>SAFE</p>  <p>DANGEROUS</p>	<p>도메인 정규 표현식 검증</p> <p>/Reg[ex]+/</p>	<p>화이트 / 블랙 리스트 검증</p> 
<p>TLD 검증</p> <p>.ORG .COM .INFO</p> <p>.CO .NET .CO.UK</p>	<p>HitAnt 머신러닝 검증</p> 	<p>검증 결과 형 변환</p> <p>✓ TRUE → 0</p> <p>✗ FALSE → 1</p>



감사합니다.