

UNIVERSITÀ DEGLI STUDI DI PISA



FACOLTÀ DI SCIENZE MATEMATICHE FISICHE E NATURALI

CORSO DI LAUREA IN INFORMATICA

RELAZIONE DI TIROCINIO

svolto presso il

DIPARTIMENTO DI INFORMATICA

Creazione di un sistema di reputazione per domini Internet

STUDENTE Antonino Lorefice

TUTORE ACCADEMICO Prof. Luca Deri

Anno Accademico 2012/13

Abstract

Per contrastare la diffusione delle minacce informatiche alcune comunità virtuali mettono gratuitamente a disposizione degli utenti un insieme di servizi per la prevenzione dei rischi derivanti dalla navigazione web e dallo scambio di messaggi di posta elettronica.

Sono disponibili inoltre dei servizi che hanno la funzionalità di distinguere fra le varie tipologie di contenuto dei siti web, allo scopo di interfacciarsi con le applicazioni di filtraggio del traffico.

Entrambe le tipologie di servizi, che si possono definire di reputazione e di categorizzazione, sono offerti in maniera molto differente e quindi inutilizzabile in modo automatico.

L'attività di tirocinio ha avuto lo scopo di analizzare ed omogeneizzare tutti questi servizi e di renderli fruibili attraverso un'unica interfaccia.

In particolare è stato creato un sistema per il salvataggio e per l'aggiornamento automatico di tutte quelle informazioni scaricabili ed un sistema unificato di interrogazione di tutti quei servizi disponibili via web.

Per la validazione del sistema sviluppato sono stati effettuati dei test, riguardo a delle liste di siti web significative, sia dal punto di vista della pericolosità sia da quello della tipologia di contenuto.

I test hanno evidenziato l'affidabilità del sistema come strumento da utilizzare dalle applicazioni di monitoraggio del traffico di rete

che si occupano della sicurezza e del filtraggio dei contenuti. Il sistema costruito è stato reso disponibile tramite una pagina web, che fornisce un'interfaccia di accesso ai suoi dati, sia alle applicazioni che alle persone umane.

Indice

1	Introduzione	5
1.1	Struttura della relazione	6
1.2	Motivazioni	8
1.3	Reputazione di domini	9
1.4	Categorizzazione di siti web	10
1.5	Obiettivi del tirocinio	11
2	Stato dell'arte	13
2.1	Categorie di malware	13
2.2	Antivirus e firewall	15
2.2.1	Limiti degli antivirus	16
2.2.2	Firewall	17
2.2.3	Filtraggio dei contenuti	19
2.2.4	Proxy web	19
2.2.5	Parental Control	20
2.3	Servizi di reputazione commerciali	21
2.3.1	TrustedSource	21
2.3.2	Commtouch	21
2.3.3	Brightcloud	22
3	Servizi Analizzati	23
3.1	DMOZ	23
3.2	Blocksi	25
3.3	Google Safe Browsing	26
3.4	Alexa	28
3.5	URL.BlackList.com	29
3.6	I-BlockList	29
3.7	Spamhaus	30
3.8	SURBL	31
3.9	SORBS	32
3.10	The Abusive Hosts Blocking List	33
3.11	APEWS.ORG	33
3.12	inps.de-DNSBL	33
3.13	mailspike	34
3.14	DNS-BH – Malware Domain Blocklist	34
3.15	Malware Domain List	35
3.16	Zeus Tracker	35
3.17	SpyEye Tracker	35
3.18	Palevo Tracker	36
3.19	Norton Safe Web	36
3.20	AVG Threat Labs	37

4	Architettura ed implementazione del sistema	39
4.1	Architettura del sistema	39
4.2	Sottosistema di categorizzazione	41
4.2.1	Schema di categorie adottato	42
4.2.2	Funzionamento del sistema di categorizzazione	44
4.3	Sottosistema di reputazione	45
4.3.1	Funzionamento del sistema di reputazione	46
4.4	Dettagli di implementazione	47
5	Validazione	49
5.1	Liste note di siti web e domini	49
5.2	Risultati ottenuti	50
5.2.1	Validazione del sottosistema di categorizzazione	50
5.2.2	Validazione del sottosistema di reputazione	54
5.3	Confronto con i sistemi commerciali	55
5.4	Casi d'uso	60
5.4.1	Caso d'uso della pagina web	60
5.4.2	Caso d'uso dell'interfaccia di accesso JSON	62
5.5	Contatto e download del codice	62
6	Conclusioni	64
6.1	Sviluppi futuri	65

Capitolo 1

1 Introduzione

Internet mette a disposizione tutta una serie di servizi allo scopo di fornire informazioni sulla pericolosità di particolari domini e sulla tipologia di contenuto ospitato in particolari siti web.

Questi servizi sono molto variegati e differiscono per tipologia di informazioni offerte, per disponibilità del servizio, per gli strumenti che mettono a disposizione e per gli aspetti legati all'uso del loro sistema.

I servizi possono occuparsi della reputazione dei domini, della categorizzazione dei siti web o di entrambi. Possono essere offerti gratuitamente ed illimitatamente o gratuitamente solo in parte oppure esclusivamente a pagamento. Possono essere fruibili come database di dati scaricabili, come strumenti consultabili via web, come API remote o come DNSBL[1], delle liste interrogabili tramite richieste DNS.

Inoltre ogni servizio ha un proprio sistema di categorie e di classificazione delle possibili minacce, che rende inutilizzabili tutte le informazioni globalmente.

Tutte queste informazioni non sono utilizzabili in modo automatico, per esempio se si volesse costruire un antivirus, che blocchi la navigazione verso quei siti elencati come pericolosi, bisognerebbe

controllare singolarmente tutte le liste dei domini e degli indirizzi pericolosi ed andare ad interrogare tutti quei servizi via web che forniscono informazioni sulla sicurezza.

Lo scopo del tirocinio è stato quello di sviluppare e rendere disponibile a tutti gratuitamente un sistema di raccolta automatizzato di tutte le informazioni scaricabili e di uno strumento automatico per l'utilizzo dei servizi remoti disponibili.

L'utilizzatore del sistema otterrà informazioni circa la tipologia del contenuto di un sito web ospitato da un certo dominio e sulla sua pericolosità, non dovendo utilizzare singolarmente i vari servizi ed avendo a disposizione un interfaccia di accesso standard.

1.1 Struttura della relazione

In questo paragrafo verrà illustrata brevemente la struttura della relazione, descrivendo i contenuti per ogni capitolo.

Capitolo 1

Il primo capitolo introduce l'argomento del tirocinio e le motivazioni che hanno dato vita al progetto. Descrive inoltre le tipologie di servizi disponibili per la reputazione dei domini e per la categorizzazione dei siti web, le loro caratteristiche, ed infine gli obiettivi dell'attività svolta.

Capitolo 2

Nel secondo capitolo verrà analizzato lo stato dell'arte inerente alle minacce della rete Internet e le principali contromisure adottate. Verranno infine presentati dei prodotti commerciali simili al sistema che si vuole sviluppare.

Capitolo 3

Nel terzo capitolo verranno descritti i servizi di reputazione e di categorizzazione analizzati ed utilizzati per l'implementazione del progetto di tirocinio.

Capitolo 4

Nel quarto capitolo verrà approfondita l'architettura del sistema realizzato, analizzando le caratteristiche di ogni elemento sviluppato e le motivazioni delle scelte progettuali.

Capitolo 5

Nel quinto capitolo verrà descritta la validazione del sistema realizzato, analizzando i risultati in termini di qualità e di prestazioni, facendo le dovute considerazioni per la parte di categorizzazione e per quella di reputazione. Verranno inoltre illustrati gli utilizzi del sistema sviluppato.

Capitolo 6

Nel sesto ed ultimo capitolo verranno trattate le conclusioni, evidenziando gli obiettivi raggiunti, gli ambiti in cui è stato utilizzato il sistema costruito ed i possibili sviluppi futuri.

1.2 Motivazioni

Attualmente la maggior parte delle tecnologie che si occupano della sicurezza e del monitoraggio del traffico di rete adottano un approccio basato sull'analisi del traffico ricevuto, sia a livello di header che a livello di contenuto[2].

Tali tecnologie però non garantiscono l'assoluta sicurezza e correttezza delle informazioni ed inoltre il controllo in tempo reale dei contenuti pesa notevolmente dal punto di vista computazionale.

È possibile però utilizzare un approccio alternativo, basato sulla reputazione e sulla categorizzazione dei domini, in modo tale da attivare i controlli in tempo reale solo quando ritenuto necessario.

A questo proposito sono disponibili un insieme di servizi che hanno lo scopo di istruire gli utenti di Internet riguardo alla tipologia di contenuto di un sito web ed alla pericolosità rappresentata da un particolare dominio.

Questi servizi sono per la maggior parte gratuiti e liberi da qualsiasi licenza d'uso, ma utilizzabili singolarmente, ognuno con le sue

caratteristiche ed il suo livello di dettaglio di informazioni fornite.

Esistono degli aggregatori di tali risorse, che comunque non comprendono globalmente tutte le risorse in tutte le tipologie disponibili e che non possono essere utilizzati in maniera automatica.

A onor del vero esistono dei sistemi a pagamento che in qualche modo aggregano tali risorse e le integrano con un motore proprio che ha lo scopo di reputare e categorizzare i domini che gli vengono sottomessi.

Da questa considerazione nasce l'esigenza di costruire e rendere disponibile gratuitamente, un sistema di reputazione e di categorizzazione, che integri tutti i servizi disponibili in un unico strumento che fornisca un'interfaccia di accesso standard per l'interazione con terze parti.

Il sistema ricavato potrà essere utilizzato dalle applicazioni di monitoraggio del traffico di rete che si occupano della sicurezza e del filtraggio dei contenuti.

1.3 Reputazione di domini

I sistemi di reputazione dei domini hanno lo scopo di fornire indicazioni circa la pericolosità di un certo dominio.

La loro pericolosità deriva dal potenziale danno che possono arrecare ad altri computer in Internet a causa della loro attività. Non per forza tali computer devono essere consapevoli del danno arrea-

to, possono semplicemente essere vittima loro stessi di attacchi che li hanno fatti diventare una minaccia per gli altri.

Per contrastare tale attività malevole alcune comunità virtuali o anche delle aziende specializzate nella lotta al malware, mettono a disposizione dei servizi, che possono essere delle semplici liste di domini e/o indirizzi IP di computer coinvolti in attività di malware. Le liste possono essere messe a disposizione via web e sono interrogabili tramite richieste al sito che le ospita o anche tramite delle API remote messe a disposizione degli sviluppatori.

Inoltre le liste possono anche essere messe a disposizione come DNSBL, un meccanismo che pubblica una lista di indirizzi IP interrogabile tramite richieste DNS. Tale meccanismo è principalmente utilizzato per la pubblicazione di indirizzi IP in qualche modo legati all'attività di spam e la maggior parte dei mail server possono essere configurati per rifiutare o contrassegnare i messaggi provenienti da host presenti nella lista.

1.4 Categorizzazione di siti web

I sistemi di categorizzazione si occupano di catalogare i siti web in base al contenuto, utilizzando un sistema di categorie. Non esiste una catalogazione standard ed in genere ogni sistema ne adotta una propria, in base agli utilizzi a cui è destinato.

La necessità di categorizzare un sito web scaturisce dall'esigenza di conoscere la tipologia di contenuto di un sito web, prima di accederne i contenuti.

Tale esigenza è quella che hanno i motori di ricerca[3] per fornire i risultati che più si avvicinano ai criteri di ricerca impostati. Anche alcune applicazioni di monitoraggio del traffico quali ad esempio quelle di filtraggio dei contenuti utilizzano tali informazioni.

I servizi di categorizzazione sono meno numerosi di quelli di reputazione e anche loro vengono resi disponibili tramite liste, servizi utilizzabili via web o API destinate agli sviluppatori.

Esistono anche un numero ristretto di DNSBL che si occupano della categorizzazione dei domini e degli indirizzi IP. Questi funzionano sempre come normali server DNS ma forniscono risposte diverse in base alla categoria del dominio o dell'indirizzo IP richiesto.

1.5 Obiettivi del tirocinio

L'attività di tirocinio ha avuto lo scopo di ricercare ed analizzare gli strumenti di reputazione e di categorizzazione disponibili su Internet. Dopo lo studio è stata effettuata una selezione di quelli che fossero utilizzabili in modo automatico.

L'obiettivo del tirocinio è stato quindi quello di omogeneizzare tutti questi servizi e di creare uno strumento totalmente gratuito, utilizzabile da tutti, che fosse estendibile e che avesse un tempo

di risposta accettabile, in considerazione del fatto che dovrà essere utilizzato durante la navigazione e quindi non potrà rallentarla eccessivamente.

Il sistema sviluppato è stato integrato in una pagina web che ha la possibilità di interagire con gli utenti umani e con le applicazioni tramite un'interfaccia di accesso standard JSON.

Capitolo 2

2 Stato dell'arte

In questo capitolo verranno analizzati gli studi e le ricerche svolte nell'ambito dell'attività di tirocinio. Il capitolo inizia con un'introduzione sulle principali minacce diffuse su Internet, illustrando le metodologie di contrasto più comuni. Di seguito vengono illustrati i principali servizi commerciali di reputazione e categorizzazione a cui il progetto si è ispirato.

2.1 Categorie di malware

Nella terminologia informatica il termine malware[4] indica un qualsiasi software realizzato con lo scopo di arrecare danni ad altri computer. La loro diffusione è in continuo aumento a causa del proliferare dei dispositivi connessi a Internet ed al diffondersi della cultura informatica.

Esistono molte tipologie di malware ed alcuni di questi sono il risultato di una composizione e rientrano pertanto in più tipologie.

Per questo vengono presentati i più conosciuti :

- Virus: programmi che fanno parte di altri programmi o che si trovano in particolari sezione del disco fisso. Si diffondono tra computer tramite lo spostamento di file effettuato dagli utenti.
- Worm: questi[5] non hanno bisogno di infestare altri programmi per diffondersi, perchè modificano il sistema operativo ospitante in maniera tale da essere eseguiti automaticamente. Si diffondono principalmente tramite Internet utilizzando tecniche di ingegneria sociale o approfittando dei difetti di alcuni programmi. Il loro scopo è quello di rallentare il sistema facendogli eseguire operazioni inutili e dannose.
- Trojan horse: software[6] che oltre ad avere funzionalità lecite che ne favoriscono la diffusione tra gli utenti, contengono delle parti dannose che vengono eseguite a loro insaputa. Il nome deriva dal fatto che non hanno la capacità di autoriprodursi e per diffondersi devono essere consapevolmente inviati alla vittima.
- Backdoor: letteralmente porta sul retro, sono dei programmi che consentono l'accesso senza autorizzazione nei sistemi in cui sono in esecuzione.
- Spyware: software[7] che raccolgono informazioni del sistema su cui sono installati. Tali informazioni spaziano da quelle che descrivono il comportamento dell'utente fino alle password.

- Dialer: si occupano di gestire la connessione a Internet tramite la normale linea telefonica. Quelli malware dirottano la connessione su numeri a tariffazione speciali.
- Hijacker: programmi che si appropriano degli strumenti di navigazione e provocano l'apertura di pagine web indesiderate.
- Rootkit: programmi utilizzati per mascherare all'utente ed ai programmi di antivirus Trojan horse e Spyware.

Il mezzo di comunicazione utilizzato per la trasmissione del malware è principalmente Internet, in particolare la navigazione web e la posta elettronica.

2.2 Antivirus e firewall

Una soluzione per la protezione dei dispositivi dal malware è l'installazione di un software chiamato antivirus[8], che ha il compito di prevenire, rilevare e rimuovere eventuale malware rintracciato.

Uno dei principali metodi di funzionamento degli antivirus è quello che si basa sulla ricerca nel disco fisso del dispositivo, di programmi con caratteristiche, chiamate firme, tipiche dei malware[9]. Il successo di questa tecnica è condizionato dal continuo aggiornamento delle firme che l'antivirus è in grado di riconoscere.

Questa tipologia di tecnica è detta analisi statica[10], in contrapposizione all'altra tecnica utilizzata, l'analisi dinamica[11], basata sullo studio dell'esecuzione dei programmi sospettati di essere malware[12].

Queste ultime tecniche sono raramente utilizzate dai software antivirus commerciali; vengono utilizzate dalle aziende produttrici di software antivirus, come supporto allo studio dei malware. Il principale limite di tali tecniche è costituito dall'alto overhead. Tuttavia alcuni antivirus implementano tale tecnica col pericolo di rallentare talmente il sistema da indurre a chi lo usa di disabilitare i controlli.

2.2.1 Limiti degli antivirus

Un normale antivirus è in grado contrastare soltanto il malware presente nel proprio database, quindi i nuovi virus non vengono riconosciuti ed occorre aggiornare costantemente le firme. Inoltre l'antivirus riesce a rintracciare il malware solamente quando è già entrato nel sistema e lo ha infettato.

Un altro limite è dato dal fatto che gli antivirus sono dei grandi consumatori di risorse del computer e rallentano in maniera importante il sistema. Inoltre non sono in grado di proteggere i dispositivi dalle minacce derivanti da attività dannose e illegali svolte da altri computer.

Le principali sono:

- lo spamming[13], cioè l'invio di messaggi di posta elettronica indesiderata
- il phishing[14], cioè l'invio di messaggi di posta elettronica che imitano quelli inviati dalle banche ed hanno lo scopo di rubare le credenziali di accesso dei servizi bancari online.

Queste sopra non compromettono l'integrità del sistema, ma rappresentato comunque degli abusi.

- DoS che è la sigla di Denial of service[15], un malfunzionamento causato da un attacco informatico che ha lo scopo di esaurire le risorse di un sistema informatico che fornisce un servizio, ad esempio un sito web, fino a renderlo non in grado di erogare il servizio

2.2.2 Firewall

Dato che un antivirus da solo, per quanto affidabile ed efficiente, non è una protezione totale contro la totalità dei malware esistente al mondo, un'ulteriore protezione è il firewall[16].

Un firewall può essere configurato per permettere di bloccare i malware, anche non conosciuti, prima che vengano a contatto con il

computer. Permette di bloccare anche quelli già presenti all'interno evitando così che possano infettare la rete a cui si è collegati.

Un firewall è quindi uno strumento aggiuntivo che impedisce ai malware di infettare la macchina, prima che possa essere individuato dall'antivirus.

La sua funzionalità principale è quella di filtrare tutti i pacchetti entranti ed uscenti, da e verso un computer o una rete, in base a delle regole che aumentano la sicurezza del sistema. Può effettuare sui pacchetti azioni di controllo, modifica e monitoraggio.

Può essere realizzato con un computer con due schede di rete, una per i pacchetti in input e l'altra per quelli in output, dotato di uno specifico software. Oppure può essere una funzionalità logica implementata in un apparato di rete.

In tutti i casi apre il pacchetto IP e legge le informazioni dell'header, ed in alcuni casi quelle del payload. La tipologia deep packet inspection[17] effettua controlli fino al livello applicativo dei pacchetti, per esempio riconoscendo e bloccando i dati appartenenti a malware noti.

Esistono anche i firewall personali[18], che sono software che permettono di filtrare i pacchetti che entrano ed escono dal calcolatore su cui sono installati[19], utilizzando in tal caso una sola scheda di rete.

In questi il principio di funzionamento differisce, in quanto le regole che definiscono il traffico permesso non vengono impostate

in base all'indirizzo IP sorgente, quello di destinazione e la porta attraverso la quale viene erogato il servizio, ma in base alla specifica applicazione.

2.2.3 Filtraggio dei contenuti

Alcuni firewall hanno la funzionalità di filtrare il traffico che arriva da Internet sulla base di criteri non riguardanti la sicurezza, ma volti a limitare l'utilizzo della rete sulla base dei protocolli o per quanto riguarda il web, a determinate categorie di siti.

Ad esempio siti con contenuti non adatti ai minori, non pertinenti all'attività lavorativa o in base alla tipologia di informazione trattata.

Il firewall può anche essere uno strumento di censura per esempio per limitare la diffusione della conoscenza e della libertà di stampa[20].

2.2.4 Proxy web

Un proxy web è un programma che si interpone tra un client ed un server http, filtrando le richieste in entrambe le direzioni.

Oltre che per migliorare le prestazioni e ridurre il consumo di banda, può essere utilizzato per monitorare il traffico effettuato, per limitare l'ampiezza di banda utilizzata dal client oppure per bloccare le pagine web in transito, per esempio bloccando quelle il cui contenuto viola determinate regole.

2.2.5 Parental Control

È un servizio[21] per il controllo delle pagine web accedute in base a prefissati criteri. Possibili utilizzi sono per la censura ai bambini dei contenuti considerati pericolosi e violenti e per limitare l'accesso in ambienti lavorativi ai contenuti non inerenti all'attività lavorativa.

Sono possibili due approcci complementari, uno di black list, in cui la navigazione è consentita verso tutti i siti non contenuti nella black list, ed uno di white list, in cui invece è consentita la navigazione solo verso quei siti contenuti nella white list. Le liste possono essere organizzate anche per categorie di contenuti ospitati dai siti web.

Il servizio può essere fornito da un software, di solito a pagamento, installato sul computer, da configurare per consentire l'accesso e proibirlo a certe categorie di siti. La maggior parte classifica le categorie dei siti in base ad un controllo in tempo reale delle pagine visitate, categorizzandoli in base alle parole trovate. Altri si basano su delle liste di siti web suddivise per categorie.

Alcuni sistemi operativi odierni integrano un software di parental control. Anche alcuni fornitori di accesso a Internet (ISP) danno la possibilità di attivare dei blocchi di contenuti non adatti ai minori. Inoltre è possibile utilizzare dei servizi di DNS che garantiscono l'accesso a siti con contenuti per minori, il più famoso e gratuito è

FamilyShield di OpenDNS.

Sono disponibili anche dei browser dedicati ai minori, che garantiscono l'accesso ai siti a loro adatti.

2.3 Servizi di reputazione commerciali

In questo paragrafo verranno illustrati e descritti alcuni dei servizi di reputazione commerciali che sono stati analizzati. Tali sistemi sono utili per bloccare efficacemente gli attacchi basati sulla rete, inviati tramite messaggi di posta elettronica ed altri protocolli.

2.3.1 TrustedSource

TrustedSource è un sistema di reputazione di Internet, di proprietà di McAfee, un'azienda che si occupa di sicurezza informatica e che produce software antivirus.

Fornisce reputazione di indirizzi IP, di url e di domini analizzando il traffico in tempo reale, dando indicazioni sul grado di pericolosità e sulla categorizzazione dei contenuti.

2.3.2 Commtouch

Commtouch è una società di sicurezza informatica, che fornisce servizi per il filtraggio dei contenuti Internet e la protezione da malware e spam. Per identificare le nuove minacce, analizza automaticamente miliardi di transazioni Internet e sulla base di questi modelli, identifica nuovo spam e attacchi malware.

2.3.3 Brightcloud

È un servizio offerto da Webroot, una società di sicurezza informatica. Lo strumento categorizza e reputa le risorse di Internet e permette di implementare soluzioni di sicurezza personalizzate.

Capitolo 3

3 Servizi Analizzati

In questo capitolo verranno illustrati i servizi di reputazione e di categorizzazione analizzati ed utilizzati per l'implementazione del sistema.

I servizi scelti sono stati confrontati con tanti altri non menzionati e gli sono stati preferiti per la loro completezza e per la mole di informazione che offrivano. Sono stati selezionati inoltre quelli aggiornati più frequentemente.

Verranno descritte le loro funzionalità, le loro caratteristiche ed il loro utilizzo all'interno del sistema.

3.1 DMOZ

L' Open Directory Project[22], anche conosciuta come DMOZ, è una web directory costruita e mantenuta da una comunità virtuale di editori volontari che viene utilizzata dai motori di ricerca per comprendere la categoria di un sito web.

Chiunque può segnalare un sito indicando la categoria più ap-

propriata e la correttezza del servizio è favorita dal fatto che ogni segnalazione viene analizzata da una persona umana.

Siccome queste persone lo fanno nel loro tempo libero, non ci sono tempi certi di risposta ed inoltre le segnalazioni possono essere bocciate.

I dati di ODP sono resi disponibili gratuitamente all'indirizzo <http://rdf.dmoz.org/rdf/> a condizione di inserire dei crediti nei siti web che li utilizzano.

Utilizza uno schema di categorizzazione gerarchico, indirizzi IP ed url con contenuti simili sono raggruppate nella medesima categoria che a sua volta può contenere delle sottocategorie. Al momento categorizza più di cinque milioni di risorse web in più di un milione di categorie e le macrocategorie sono: Adult, Arts, Business, Computers, Games, Health, Home, Kids and Teens, News, Recreation, Reference, Regional, Science, Shopping, Society, Sports e World.

La categoria Adult non è presente nei link della home page, ma è raggiungibile andando all'indirizzo <http://www.dmoz.org/adult> ed è disponibile in un separato file scaricabile.

La categoria Kids and Teens contiene risorse web appropriate per persone sotto i diciotto anni di età e viene resa disponibile anch'essa in un file separato.

Inoltre mentre tutte le categorie principali, al loro interno sono organizzate in sottocategorie per argomento, nella categoria Regional sono organizzate per regione geografica.

Nuove versioni dei file vengono fornite in genere settimanalmente; i dati ODP danno vita al nucleo di molti dei più grandi motori di ricerca, tra cui Netscape Search, AOL Search ed Alexa. Google Directory usava le informazioni ODP, fino al suo oscuramento nel luglio del 2011.

Per la costruzione del sistema sono state scaricate solo le informazioni riguardanti i domini ed ignorate quelle riguardanti le url.

3.2 Blocksì

Blocksì è un piccolo team di lavoro sloveno il cui obiettivo principale è quello di fornire protezione verso i contenuti web illegali, immorali e rischiosi.

A tale proposito offre una estensione per browser, Blocksì - Web filtering and parental control, che permette di filtrare il traffico web.

Sono disponibili inoltre delle API che consentono agli sviluppatori di applicazioni di interrogare il motore Blocksì ottenendo informazioni sulla tipologia di contenuto e sulla pericolosità di un sito web.

Al momento categorizza i contenuti web in 79 categorie e possiede un database con più di 76 milioni di siti web. Lo schema di categorizzazione è organizzato in sette macrocategorie che sono: Potentially Liable, Controversial, Bandwidth Consuming, Security Risk, General Business (Business), General Interest (Personal) e

Unrated. Tutte ad eccezione di Unrated hanno delle sottocategorie che specificano meglio la categoria del sito web della risorsa web sottomessa.

Blocksi non è solo un motore di categorizzazione ma anche di reputazione, in quanto offre indicazioni sulla sicurezza del dominio a cui la risorsa web che si è sottomessa appartiene (macrocategoria Security Risk).

Durante i test effettuati, nessun tipo di blocco è stato attivato sia utilizzando il servizio in modo non automatico sia in modo automatico. Si suppone quindi che almeno per il momento non sia attivo nessun tipo di blocco e che si possa utilizzare il servizio illimitatamente.

Si è riscontrato tuttavia una sorta di registrazione da parte del motore, delle richieste di categorizzazione, in quanto alcune risorse web, nella prima sottomissione al servizio davano una risposta diversa da tutte le successive sottomissioni.

3.3 Google Safe Browsing

Google Safe Browsing è un servizio offerto da Google che fornisce liste di risorse web che ospitano malware. I browser Google Chrome, Apple Safari e Mozilla Firefox utilizzano le liste di tale servizio per bloccare i siti pericolosi.

All'indirizzo `http://www.google.com/safebrowsing/ diagnostic? site=miosito.it` sono disponibili diverse informazioni, come lo sta-

to attuale del sito sottomesso, cosa è successo dopo la visita di Google al sito, se ha assunto la funzione di intermediario per la distribuzione di malware e se il sito stesso ha ospitato malware.

Le informazioni rilevanti ai fini del sistema di reputazione sono quelle sullo stato attuale del sito, le possibili stringhe ritornate sono:

- Site is listed as suspicious - visiting this web site may harm your computer.
- This site is not currently listed as suspicious.

L'utilizzo automatico del servizio ha prodotto il blocco da parte di Google, tramite l'introduzione di captcha. Il servizio è stato comunque riattivato dopo qualche ora.

Google fornisce inoltre una API pubblica per il servizio, che utilizza un cookie che aiuta Google a conoscere tutti i siti che sono stati visitati. Lo strumento ritorna una risposta http con il risultato della richiesta: se il dominio ospita qualche minaccia, il contenuto della risposta sarà malware altrimenti la risposta non avrà contenuto.

L'utilizzo intensivo dello strumento non ha causato blocchi del servizio.

3.4 Alexa

Alexa è un motore di ricerca[24] con un servizio di web directory che si occupa anche di raccogliere statistiche sul traffico di Internet.

Il motore classifica i siti basandosi sulle visite effettuate dagli utenti delle barre degli strumenti integrate nei browser.

Inoltre Alexa Top Sites, è un servizio che fornisce l'accesso alle liste dei siti web più visitati globalmente e singolarmente per ogni nazione.

Alexa non esporta i dati riguardanti la categorizzazione dei dati e le liste dei siti web più visitati, ma le rende disponibili via web.

Le informazioni di categorizzazione dei siti web, sono reperibili all'url `http://www.alexa.com/siteinfo/dominio`, dove dominio è il dominio di cui si vogliono conoscere le informazioni, nonostante funzioni anche con le url, si ottengono però informazioni riguardanti il dominio a cui la url appartiene. Il motore di categorizzazione usa lo stesso schema di categorie adottato da DMOZ.

È un servizio a pagamento e gli utilizzi automatici vengono rilevati e bloccati, quindi è utilizzabile per singole richieste di categorizzazione.

Le lista dei siti web mondiali più visitati è disponibile all'indirizzo `http://www.alexa.com/topsites`, e quelle delle singole nazioni all'indirizzo `http://www.alexa.com/topsites/countries/sigla`, dove sigla è la sigla del paese preso in considerazione, ad esempio

IT per l'Italia. Il servizio è utilizzabile in maniera automatica, in quanto durante i test non si è rilevato nessun blocco.

3.5 URL.BlackList.com

Si tratta di un servizio commerciale che fornisce una lista di domini e di url organizzati in cartelle, il cui nome rappresenta la categoria delle url e dei domini contenuti dentro ognuna di essa.

La maggior parte delle informazioni sono raccolte da vari siti gratuiti, quindi la lista è notevolmente più grande di altre liste che è possibile trovare. La lista viene generata automaticamente in genere ogni giorno e non contiene soltanto le risorse web che hanno una cattiva reputazione, ma contiene anche molte altre categorie di siti. Secondo la documentazione del sito la lista è scaricabile gratuitamente una sola volta, ma i test hanno evidenziato che il blocco del download dopo un po di tempo, circa 48 ore, viene disattivato. Quindi il sistema tenta di scaricare la lista, se ci riesce costruisce il database dal file scaricato altrimenti rimane tutto invariato.

Per la costruzione del sistema sono state scaricate solo le informazioni riguardanti i domini ed ignorate quelle riguardanti le url.

3.6 I-BlockList

È un servizio in parte gratuito dedicato alla raccolta ed alla distribuzione di liste proprie di indirizzi IP e di altri servizi.

Gli intervalli di indirizzi IP sono raggruppati per appartenenza ad una stessa tipologia di azienda, di organizzazione, di servizio offerto o di attività dannosa svolta.

Le categorie disponibili coprono poche tipologie di contenuti di siti web, mentre le varie attività dannose sono tutte elencate separatamente.

3.7 Spamhaus

Spamhaus è organizzazione internazionale no-profit la cui missione è quella di tenere traccia delle operazioni e delle fonti di spam su Internet, con lo scopo di fornire gratuitamente protezione in tempo reale della rete.

Mantiene un certo numero di basi di dati di intelligence di sicurezza e di basi di dati per il blocco dello spam e dei malware in tempo reale (DNSBL).

I DNSBL di Spamhaus sono utilizzati dalla maggior parte dei provider di posta elettronica, dalle organizzazioni governative e dalle università.

Per soddisfare la domanda per i suoi DNSBL, Spamhaus dispone di una delle più grandi infrastrutture DNS del mondo.

L'utilizzo gratuito non consente di effettuare un elevato volume di query e l'utilizzo del servizio per scopi commerciali.

Le liste che vengono rese disponibili sono:

- La Spamhaus Block List (SBL) che elenca gli indirizzi IP che sono fonte di spam o che forniscono loro servizi.
- La Exploits Block List (XBL) è una banca dati di indirizzi IP di worm, virus, motori di spam, PC e server infettati da trojan horse.
- La Domain Block List (DBL) è una lista di nomi di dominio che fanno attività di spam o che ospitano malware.

Le liste sono interrogabili in tempo reale dai sistemi di posta, attraverso Internet.

3.8 SURBL

SURBL è un insieme di liste di siti web che appaiono nei corpi dei messaggi indesiderati.

Queste liste sono raccolte da altri siti web in un DNSBL e possono essere utilizzate per bloccare le connessioni verso questi siti o per bloccare i messaggi di posta elettronica che nel loro corpo contengono un sito listato.

Alla pagina web <http://www.surbl.org/surbl-analysis> è possibile controllare lo stato di qualsiasi dominio o indirizzo IP. Non è possibile usare questo form per test automatizzati in quanto è necessario immettere un captcha.

I dati SURBL sono forniti agli utenti di tutto il mondo attraverso i server DNS pubblici o attraverso un servizio di feed di dati.

Il primo (query DNS) è completamente gratuito e soggetto a determinate restrizioni di uso, mentre il secondo è un servizio a pagamento.

Per gli utenti individuali, le piccole organizzazioni di beneficenza o non-profit, piccole imprese o altri enti che hanno meno di mille utenti o che sottomettono meno di 250,000 richieste al giorno, il servizio di query SURBL è completamente gratuito.

3.9 SORBS

SORBS (Spam and Open Relay Blocking System) è un DNSBL di server di posta elettronica sospettati di invio o di inoltro di messaggi di spam, di host che sono stati attaccati e dirottati o infestati da trojan.

Il servizio è gratuito ed oltre ai messaggi di spam, consente di bloccare attacchi di phishing, altre forme dannose di posta elettronica, i server che sono stati attaccati e dirottati e quelli infestati da trojan horse.

L'utilizzo del DNSBL non deve superare le 10 richieste DNS al secondo pena il possibile blocco del servizio senza preavviso. Alla pagina web <http://www.sorbs.net/lookup.shtml> è disponibile uno strumento di lookup per controllare indirizzi o nomi di host. Per procedere è però necessario inserire un codice che compare in

un'immagine.

3.10 The Abusive Hosts Blocking List

AHBL (The Abusive Hosts Blocking List) è un database di host noti per la loro attività dannosa su internet come spam, attacchi denial of service e molto altro.

I dati forniti sono in parte dati estratti da varie fonti su Internet e in parte rilevati da strumenti che analizzano gli host. È stato sviluppato per l'impiego in servizi di posta elettronica e nei sistemi di filtraggio.

3.11 APEWS.ORG

APEWS è il successore di SPEWS un servizio anonimo che manteneva un DNSBL di intervalli di indirizzi IP appartenenti a fornitori di servizi internet (ISP) che ospitavano spammer e che mostravano poca attenzione nella prevenzione dei pericoli.

APEWS identifica spammer noti e le operazioni di spam, elencandoli appena iniziano ed anche prima. Il servizio è criticato in quanto blocca tutto quanto un ISP.

3.12 inps.de-DNSBL

È un sito web che pubblica una lista di indirizzi IP da cui hanno ricevuto messaggi di posta elettronica che hanno classificato come spam.

Sono disponibili delle API per segnalare delle fonti di spam e un DNSBL altamente affidabile interrogabile gratuitamente.

3.13 mailspike

È un servizio che permette l'identificazione e quindi il blocco degli spammer noti.

Tutti gli indirizzi IP elencati dal servizio vengono costantemente monitorati e le liste vengono aggiornate frequentemente.

Il servizio consiste di due insiemi di dati complementari:

- dati basati sul comportamento nel tempo di un indirizzo IP
- indirizzi IP che hanno partecipato ad invio di spam distribuito

Il servizio può essere utilizzato tramite query DNS che non devono superare le 100,000 al giorno.

3.14 DNS-BH – Malware Domain Blocklist

Il progetto DNS-BH crea e mantiene una lista di domini che sono noti per essere utilizzati per propagare malware e spyware. La lista è disponibile in vari formati ed è interrogabile tramite richieste DNS al DNSBL.

La lista viene fornita gratuitamente per uso non commerciale, come strumento della lotta al malware. Qualsiasi uso commerciale è severamente vietato senza preventiva autorizzazione.

3.15 Malware Domain List

Malware Domain List è un progetto non commerciale di distribuzione di una lista di domini da considerare pericolosi. La lista fornisce informazioni sulla tipologia di pericolo rappresentato e può essere utilizzata gratuitamente da chiunque. Il servizio è disponibile anche via web e per ogni minaccia fornisce il dominio, l'indirizzo IP, la tipologia di minaccia riscontrata e la data di riscontro.

3.16 ZeuS Tracker

ZeuS Tracker è un servizio globale che cattura e rintraccia host che ospitano zeus, che sono trojan che rubano le credenziali di accesso di vari servizi online, come social network, conti bancari online, account ftp, account di posta elettronica e altri (in generale phishing).

L'obiettivo principale è quello di fornire la possibilità di bloccare host zeus noti, rendendo disponibili liste di domini e di indirizzi IP.

3.17 SpyEye Tracker

Spyeye Tracker è un progetto molto simile a ZeuS Tracker con la leggerezza differenza che SpyEye tiene traccia e monitorizza i Spyeye maliziosi, dei malware che dovevano essere i successori di ZeuS.

3.18 Palevo Tracker

Palevo Tracker offre tre diverse liste che possono essere utilizzate per bloccare l'accesso a reti note infettate dal worm palevo, un malware che è in grado di autoreplicarsi ed in grado di diffondersi senza legarsi ad altri eseguibili.

Palevo si diffonde utilizzando la messaggistica istantanea, le reti peer-to-peer e le unità rimovibili.

3.19 Norton Safe Web

Norton Safe Web[26] è un servizio commerciale sviluppato da Symantec Corporation che ha lo scopo di aiutare gli utenti ad identificare siti web dannosi.

Le informazioni fornite sono basate su analisi automatizzate e feedback degli utenti. Quando avviene un drive-by download da un sito web, l'url sospetta viene segnalata automaticamente a Norton Safe Web per l'analisi. Il sito segnalato è classificato come pericoloso se l'analisi conferma che il download è dannoso. Per garantire che il rating rispecchi fedelmente lo stato attuale di un sito, Norton Safe Web esegue frequente rianalisi. I siti non sicuri che hanno più probabilità di essere stati ripuliti vengono rianalizzati spesso, mentre quelli che potrebbero richiedere più tempo per rimuovere le minacce vengono rianalizzati meno frequentemente.

Una versione limitata, standalone di Safe Web è disponibile come

freeware, inoltre è disponibile come plugin per i browser e via web.

La versione via web fornisce informazioni sulla reputazione del sito, il numero ed i tipi di minacce rilevate, il numero di osservazioni pervenute da parte dei clienti di Norton e il traffico del sito.

3.20 AVG Threat Labs

AVG Threat Labs[27] è un portale di informazione online che raccoglie le minacce del web che AVG ottiene quotidianamente dai suoi 100 milioni di utenti che usano il suo software AVG Secure Search.

Gli utenti che hanno installato il software sono protetti in tempo reale dai siti considerati minacciosi; AVG Secure Search avverte l'utente prima che visiti pagine Web pericolose.

Lo strumento effettua quindi una protezione preventiva, ogni pagina web viene controllata prima di essere richiesta.

La funzionalità è disponibile via web e permette di ottenere direttamente dal portale informazioni riguardo la sicurezza di una url di un sito web.

Oltre alle informazione sulla sicurezza vengono fornite statistiche ed analisi del sito.

Non è consentito l'uso del servizio tramite qualsiasi mezzo diverso dall'interfaccia che è fornita da AVG Technologies, a meno che non sia stato specificatamente autorizzato a farlo in un accordo separato.

L'utente accetta espressamente di non accedere (o tentare di aver accesso) al servizio tramite mezzi automatici (incluso l'utilizzo di script o crawler web) a meno che non sia stato specificatamente autorizzato ad agire così in un separato contratto.

Capitolo 4

4 Architettura ed implementazione del sistema

In questo capitolo viene descritta l'architettura del sistema e la metodologia di utilizzo dei servizi descritti nel capitolo precedente. Verranno così esposte le scelte progettuali impiegate.

4.1 Architettura del sistema

Il sistema è suddiviso principalmente in tre parti:

- un sottosistema che si occupa del salvataggio e dell'elaborazione delle informazioni di reputazione e di categorizzazione, dai servizi descritti nella sezione 3. Le informazioni riguardanti i domini Internet vengono salvate in dei database SQLite3, mentre quelle riguardanti gli indirizzi IP vengono salvate in degli alberi, detti patricia tree, che sfruttando le caratteristiche degli indirizzi IP, li ordinano, li gestiscono e li ricercano velocemente.
- un sottosistema di categorizzazione, che data una certa risorsa web, quale può essere un dominio, un indirizzo IP o una url,

restituisce la categoria del contenuto ospitato. La categoria è riferita al dominio preso in considerazione anche se la risorsa sottomessa era una url. Il sottosistema ha a disposizione una cache delle ultime richieste di categorizzazione servite che viene invalidata ogni 24 ore, un database di tutte le richieste di categorizzazione servite e i database delle informazioni scaricate dai servizi descritti nella sezione 3. L' utilità della cache e del database delle richieste di categorizzazione servite è quella di diminuire il tempo di risposta del sistema. Il database delle richieste di categorizzazione servite è stato previsto anche per l'eventuale controllo di correttezza delle richieste servite.

- un sottosistema di reputazione, che restituisce informazioni riguardanti la pericolosità di una risorsa web. Il risultato sarà anch'esso riferito al dominio. Il sottosistema di reputazione ha anch'esso a disposizione una cache delle richieste di reputazione servite, che viene invalidata ogni 5 minuti. La cache ha lo scopo di far diminuire il tempo di risposta del sistema.

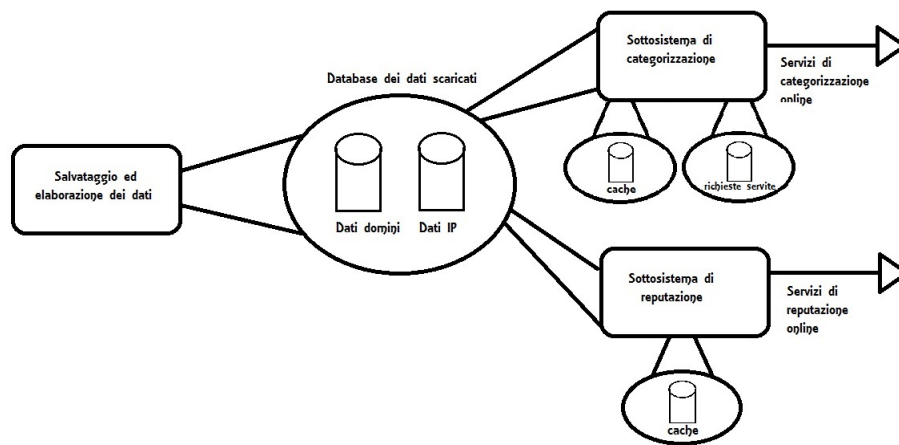


Fig. 4.1: Architettura del sistema

L'architettura del sistema rende indipendente l'aggiornamento dei tre sottosistemi e delle componenti al loro interno. Erano possibili altri schemi architetturali ma è stata preferita questa per favorire l'estensione e la modifica delle varie componenti.

4.2 Sottosistema di categorizzazione

Per la categorizzazione delle risorse web è stata realizzato uno schema di categorie che riuscisse a catalogare tutti i possibili contenuti in maniera corretta e non troppo dettagliata. Tale schema è stato il risultato dello studio e della sintesi degli altri schemi di categorie adottati dai servizi di categorizzazione analizzati.

Tutti gli schemi adottati dai motori di categorizzazione analizzati non rispondevano alle esigenze del sistema implementato, alcuni non avevano un numero sufficiente di categorie per catalogare tutte le risorse web esistenti. Altri avevano un elevato numero di categorie ma che coprivano solamente alcune tipologie di contenuti. Infine altri coprivano esaustivamente ogni tipologia di contenuto in maniera però troppo dettagliata. L'elevato livello di dettaglio rendeva impossibile il mappaggio dei dati degli altri sistemi di categorizzazione.

4.2.1 Schema di categorie adottato

Lo schema di categorie adottato cataloga i contenuti dei siti web per macrocategoria. Ogni macrocategoria copre un determinato settore che non sconfinava in altri settori. Si è cercato di fare un minimo comune denominatore di tutte gli schemi di catalogazione analizzati con lo scopo di ottenere il medesimo risultato da ognuno di essi.

Di seguito l'elenco delle categorie e la loro descrizione:

1. Adult and Controversial: siti con contenuti vietati ai minori, come pornografia, nudità, sessuologia, violenza, droga, alcool e comportamenti pericolosi
2. Arts and Entertainment: siti con contenuti artistici e di intrattenimento, come televisione, film, musica, opere d'arte e fotografia

3. Business and Economy: siti con contenuti di economia e di affari, come lavoro, aziende e finanza
4. Computers and Internet: siti con contenuti di informatica e di servizi Internet, come mail, portali, motori di ricerca e software
5. Education and Learning: siti con contenuti inerenti l'istruzione e l'apprendimento
6. Games: siti con contenuti inerenti giochi e passatempi
7. Health and Medicine: siti con contenuti sulla salute e sulla medicina
8. Home and Family: siti con contenuti sulla casa e sulla famiglia, come giardinaggio, cucina e fai da te
9. Kids and Teens: siti con contenuti adatti ai bambini ed ai ragazzi
10. News and Media: siti con contenuti di informazione on line
11. Politics and Society: siti con contenuti di politica e riguardanti la società, come governo, pubblica amministrazione e religione
12. Recreation and Sports: siti con contenuti sul tempo libero e gli sport
13. Shopping: siti con contenuti inerenti gli acquisti on line e la pubblicità

14. Social Network: siti di social network e di pagine personali

4.2.2 Funzionamento del sistema di categorizzazione

Il sistema si divide in due componenti, una che scarica e salva i dati dai database gratuiti disponibili online ed una che serve le richieste di categorizzazione.

Il primo si esegue quotidianamente e scarica dai database online i dati. Le informazioni riguardanti i domini vengono salvati in dei database sqlite3, le informazioni riguardanti gli indirizzi IP vengono salvati in una struttura dati patricia tree. I database dei domini sono tre e vengono interrogati nell'ordine in cui sono di seguito descritti.

Il primo contiene i 500 siti più visitati in Italia, servizio offerto da Alexa, categorizzati con il motore online di Blocksì. Questo database ha lo scopo di servire velocemente le richieste di categorizzazione dei siti più visitati. Il secondo contiene i domini scaricati da DMOZ ed il terzo quelli scaricati da URL.BlackList.

La struttura dati patricia tree contiene gli indirizzi IP scaricati da I-Blocklist.

Il secondo componente si occupa di servire le richieste di categorizzazione pervenute al sistema dall'interfaccia di accesso JSON, andando a ricercare nei dati scaricati dal primo componente. In caso di assenza di risultati, viene interrogato il motore di categorizzazione online di Blocksì.

Il sistema mantiene una cache ed un database sqlite3 delle richieste servite. La cache ha validità di 24 ore ed ha lo scopo di abbassare il tempo di risposta del sistema. Il database delle richieste servite ha anch'esso il medesimo compito, ma anche quello di strumento di controllo del funzionamento del sistema tramite la verifica della correttezza delle richieste servite.

Inoltre vengono registrate in un separato database le richieste che il sistema non è stato in grado di servire allo scopo di essere eventualmente servite.

In caso che il sistema non riesca a categorizzare la risorsa sottomessa, è possibile usufruire del motore di categorizzazione di Alexa che però fa aumentare il tempo di risposta del sistema a causa della sua natura remota. Il servizio però non è gratuito e se utilizzato in modo intensivo viene sospeso.

4.3 Sottosistema di reputazione

I sistemi di reputazione analizzati avevano ognuno uno schema di reputazione diverso. Alcuni si limitavano ad indicare se la risorsa fosse malware o meno, altri fornivano indicazioni dettagliate riguardo la tipologia di malware o altri parametri riguardanti l'attività malevola svolta (numero di attacchi perpetrati, data ultimo attacco, numero di computer infestati ecc.). Altri ancora fornivano vari gradi di pericolosità della risorsa web sottomessa al sistema.

Invece nel sistema si è scelto di non fornire indicazioni sulla tipolo-

gia di malware rintracciata e sul suo grado di pericolosità, indicando semplicemente la situazione come malware.

4.3.1 Funzionamento del sistema di reputazione

Il sistema si divide in due componenti, una che scarica e salva i dati dai database gratuiti disponibili online ed una che serve le richieste di reputazione.

Il primo si esegue ogni 2 ore e scarica dai database online i dati. Le informazioni riguardanti i domini vengono salvati in dei database sqlite3, le informazioni riguardanti gli indirizzi IP vengono salvati in delle strutture dati patricia tree. I database dei domini sono cinque, che contengono i domini reperiti dal sito Malware Domain List, da URL.BlackList, da Malware Domain Blocklist, da ZeuS Tracker e da SpyEye Tracker.

Le struttura dati patricia tree contengono gli indirizzi IP scaricati da I-Blocklist, da ZeuS Tracker, da SpyEye Tracker, da Malc0de, da Spamhaus e da Emerging Threats.

Il secondo componente si occupa di servire le richieste di reputazione pervenute al sistema dall'interfaccia di accesso JSON, andando a ricercare nei dati scaricati dal primo componente. In caso di assenza di risultati, vengono interrogati in ordine il motore online BlocksI, l'API di Google Safe Browsing e i DNSBL di SURBL e di Spamhaus.

Il sistema mantiene una cache che ha validità di 5 secondi ed ha lo scopo di abbassare il tempo di risposta del sistema. Si è scelto

di non mantenere un database delle richieste di reputazione servite, in quanto la reputazione di una risorsa web può cambiare da un momento all'altro.

L'assenza di risposta dal sistema non necessariamente indica l'assoluta sicurezza della risorsa web sottomessa.

Per maggiore sicurezza è possibile in caso consultare dei motori di reputazione online offerti da AVG Threat Labs e da Norton Safe Web.

Il motore AVG reputa le risorse in base alle sue liste mentre quello Norton, scansiona e analizza i siti. Il loro utilizzo aumenta notevolmente il tempo di risposta del sistema.

4.4 Dettagli di implementazione

Per l'implementazione del sistema si è scelto di utilizzare il linguaggio Python-2.7 perchè facile da usare e portabile su tutte le piattaforme purché dotate dell'interprete. Nonostante sia un linguaggio interpretato è performante grazie al fatto che il codice viene compilato in un bytecode molto efficiente che permette di raggiungere prestazioni vicine ai linguaggi in codice nativo. Ha una implementazione efficiente di molte strutture dati e funzioni e gestisce la memoria con un meccanismo di garbage collection.

Per i database si è scelto di utilizzare il modulo di python per sqlite3, il quale ha il vantaggio di implementare un DBMS SQL che non necessita di appoggiarsi a un server. Il database durante i test

è risultato essere più leggero e veloce di altri DBMS testati.

Per le strutture dati patricia tree è stata utilizzata una libreria c esistente che sfrutta le caratteristiche degli indirizzi IP per memorizzarli, ordinarli e ricercarli velocemente.

Capitolo 5

5 Validazione

In questo capitolo si vuole validare e testare il sistema sviluppato, analizzando i risultati ottenuti e paragonando il sistema a quelli commerciali di reputazione e di categorizzazione menzionati nello stato dell'arte. Per poter visualizzare i dati restituiti dal sistema sono stati effettuati dei test automatici con delle liste note di siti web e domini ed i risultati salvati in dei file di testo.

5.1 Liste note di siti web e domini

La validazione del sistema è stata realizzata sottomettendo al sistema di categorizzazione la lista dei siti più visitati in Italia e nel mondo, entrambe servizi di Alexa, e la lista dei siti più visitati nel mondo fino a luglio 2011 secondo Google adplanner. La scelta di queste liste ha lo scopo di valutare i risultati forniti dal sistema implementato.

Per quanto riguarda il sottosistema di reputazione la validazione è intrinseca nello stesso sottosistema, in quanto tutte le liste di risorse

web considerate dannose che sono state reperite, sono stata inserite nel sistema. Sono stati comunque effettuati dei test con le liste sopracitate ed altre per fare delle valutazioni riguardo il tempo di risposta del sistema.

5.2 Risultati ottenuti

In questo paragrafo verranno illustrati i risultati dei test effettuati, utilizzando le liste sopracitate, sul sottosistema di categorizzazione e su quello di reputazione. Verranno inoltre valutati i risultati sul sottosistema di categorizzazione includendo il motore Alexa e quelli del sottosistema di reputazione includendo i motori AVG e Norton. Questi ultimi due non si limitano a controllare la presenza delle risorse nelle loro liste ma vanno ad analizzare le risorse direttamente, alla ricerca di malware. Le validazioni inoltre sono state eseguite con il database del sistema vuoto e la cache disattivata.

5.2.1 Validazione del sottosistema di categorizzazione

La lista dei 500 siti web più visitati secondo la toolbar di Alexa è stata categorizzata per circa il 90% ed i risultati sono stati compatibili con quelli del motore di Alexa. Le categorizzazioni sono state tutte servite utilizzando il database costruito con tali domini e categorizzati con il motore remoto Blocksì. Questo database viene consultato come prima risorsa per diminuire il tempo di risposta del sistema per i domini più noti.

Il tempo di risposta del sistema è risultato paragonabile a quello del motore di categorizzazione di Blocks i e a quello di Alexa.

Aggiungendo al sistema il servizio offerto da Alexa (a pagamento) si ha come risultato quello di categorizzare alcuni dei domini che il sistema non aveva categorizzato al costo di un tempo di risposta più elevato, a causa della natura remota del motore Alexa.

Top site Italia	Categorizzazione del sistema	Categorizzazione Alexa
google.it	Computers and Internet	World,Italiano,Computer,Internet,Ricerca
facebook.com	Social Network	Computers,Internet,Social Networking
google.com	Computers and Internet	Computers,Internet,Search Engines
youtube.com	Arts and Entertainment	Arts,Video,Community Video
yahoo.com	Computers and Internet	Computers,Internet,Web Portals
wikipedia.org	Education and Learning	Computers,Open Source,Open Content
libero.it	Computers and Internet	not categorized
ebay.it	Shopping	World,Italiano,Acquisti Online,Aste
repubblica.it	News and Media	World,Italiano,Notizie,Quotidiani
corriere.it	News and Media	World,Italiano,Notizie,Quotidiani

Fig. 5.1: Categorizzazione top site Italia con il sistema e con Alexa

La lista dei 500 siti più visitati nel mondo è stata categorizzata anch'essa per il 90% circa, tutti i domini categorizzati sono compatibili con quelli del motore di Alexa. Le richieste di categorizzazione sono state servite in maggior parte dal database creato con i dati scaricati da DMOZ e da URL.BlackList.com, in parte servite dal

database creato con i 500 siti più visitati in Italia categorizzati con Blocksì ed in minor parte dal motore remoto Blocksì.

Il tempo di risposta è stato ancora una volta paragonabile a quello dei motori Blocksì ed Alexa, seppur leggermente maggiore, causato dal fatto che per la categorizzazione di alcuni domini è stato necessario l'utilizzo del motore remoto Blocksì.

Ancora una volta aggiungendo al sistema il servizio offerto da Alexa (a pagamento) si ha come risultato quello di categorizzare alcuni dei domini che il sistema non aveva categorizzato al costo di un tempo di risposta più elevato.

Top site Italia	Categorizzazione del sistema	Categorizzazione Alexa
google.com	Computers and Internet	Computers,Internet,Search Engines
facebook.com	Social Network	Computers,Internet,Social Networking
youtube.com	Arts and Entertainment	Arts,Video,Community Video
yahoo.com	Computers and Internet	Computers,Internet,On the Web,Web Portals
baidu.com	Computers and Internet	World,Chinese Simplified CN
wikipedia.org	Education and Learning	Computers,Open Source,Open Content
qq.com	Computers and Internet	World,Chinese Simplified CN
linkedin.com	Business and Economy	Computers,Internet,Social Networking
live.com	Computers and Internet	World,Vietnamese,Tin học,Internet
twitter.com	Social Network	Computers,Internet,Social Networking

Fig. 5.2: Categorizzazione top site mondo con il sistema e con Alexa

La lista dei 1000 siti più visitati nel mondo fino a luglio 2011 è

stata categorizzata anch'essa per il 90%, tutti i domini categorizzati sono compatibili con quelli forniti da Google adplanner. Le richieste di categorizzazione sono state in maggior parte servite dal database creato con i dati scaricati da DMOZ e da URL.BlackList.com, in parte servite dal motore remoto Blocksì ed in minor parte con il database dei siti più visitati in Italia categorizzati da Blocksì.

Il tempo di risposta è stato ancora una volta paragonabile a quello dei motori Blocksì ed Alexa, ma leggermente superiore, a causa dal fatto che per la categorizzazione di alcuni domini è stato necessario l'utilizzo del motore remoto Blocksì.

Aggiungendo al sistema il servizio offerto da Alexa (a pagamento) si ha come risultato quello di categorizzare alcuni dei domini che il sistema non aveva categorizzato al costo di un tempo di risposta più elevato.

Top site Google adplanner	Categorizzazione del sistema	Categorizzazione Google adplanner
facebook.com	Social Network	Social Networking
youtube.com	Arts and Entertainment	Video-Sharing
yahoo.com	Computers and Internet	Search
live.com	Computers and Internet	Portal
msn.com	Computers and Internet	News
wikipedia.org	Education and Learning	Reference
blogspot.com	Politics and Society	Blogging
baidu.com	Computers and Internet	Search
qq.com	Computers and Internet	Instant Messaging
adobe.com	Computers and Internet	Software

Fig. 5.3: Categorizzazione top site Google adplanner con il sistema

5.2.2 Validazione del sottosistema di reputazione

La validazione del sottosistema di reputazione utilizzando le liste sopracitate ha evidenziato l'assenza totale di minacce fra i domini elencati in esse. Il tempo di risposta è paragonabile al sottosistema di categorizzazione nonostante integri al suo interno un numero maggiore di servizi via web.

Per la validazione del sistema dal punto di vista della correttezza è stata costruita una lista random di domini malevoli. La totalità delle minacce è stata intercettata ed il sistema ha servito le richieste con un tempo di risposta minore di quello dei test precedenti, grazie alla presenza di tutti i domini sottomessi all'interno dei database

costruiti.

dominio o indirizzo IP malevolo	servizio che ha intercettato la minaccia
ianfette.org	Blocksi (API)
barakair.com	I-BlockList (patricia tree indirizzi IP)
18dd.net	SURBL (DNSBL)
fullsfinesr.info	Blocksi (API)
servegame.org	URL.BlackList.com (db domini)
twonext.com	I-BlockList (patricia tree indirizzi IP)
109.73.106.6	I-BlockList (patricia tree indirizzi IP)
74.208.85.228	I-BlockList (patricia tree indirizzi IP)
ygl.a.ru	URL.BlackList.com (db domini)
spark29.ru	I-BlockList (patricia tree indirizzi IP)

Fig. 5.4: Reputazione domini malevoli random con il sistema

Integrando al sistema i servizi di reputazione via web, di Norton Safe Web e di AVG Threat Labs, si aumenta l'affidabilità ma il tempo di risposta è notevolmente maggiore a causa sia della natura remota dei servizi, sia della tipologia di analisi effettuata, non basata solamente su delle liste.

5.3 Confronto con i sistemi commerciali

I test delle liste sopracitate con i sistemi commerciali descritti nello stato dell'arte hanno prodotto dei risultati del tutto paragonabili in termini di correttezza e di tempo di risposta a quelli del sistema

sviluppato, ad eccezione del sistema di reputazione di Commtouch (fig. 5.6) che non intercetta la maggior parte delle minacce. Di seguito le tabelle con degli esempi di test effettuati, le prime due colonne di dati riguardano la categorizzazione, le ultime due la reputazione.

dominio	Categoria Trusted Source	ip o dominio	Reputazione Trusted Source
google.com	Search Engines	depenam.com	High Risk (Web)
facebook.com	Social Networking	barakair.com	High Risk (Web)
youtube.com	Streaming Media	18dd.net	High Risk (Web)
yahoo.com	Portal Sites	fullsfinesr.info	High Risk (Web)
baidu.com	Search Engines	servegame.org	Unverified (Web)
wikipedia.org	Education/Reference	twonext.com	High Risk (Web)
qq.com	Portal Sites	109.73.106.6	High Risk (Web & Mail)
linkedin.com	Professional Networking	74.208.85.228	Minimal Risk (Web & Mail)
live.com	Search Engines	ygl.ru	High Risk (Web)
twitter.com	Social Networking	spark29.ru	High Risk (Web)
msn.com	Portal Sites	artvideo3d.ru	High Risk (Web)
blogspot.com	Blogs/Wiki	4yoursecret.co.tv	High Risk (Web)
adobe.com	Software/Hardware	dewell.ru	High Risk (Web)

Fig. 5.5: Categorizzazione e reputazione Trusted Source

dominio	Categoria Commtouch	ip	Reputazione Commtouch
google.com	Search Engines & Portals	37.230.212.0	Unknown
facebook.com	Social Networking	5.34.242.0	Unknown
youtube.com	Entertainment	193.106.173.198	Unknown
yahoo.com	Search Engines & Portals	176.9.36.151	No Risk
baidu.com	Search Engines & Portals	109.68.190.148	Unknown
wikipedia.org	Education	31.186.3.99	Unknown
qq.com	Search Engines & Portals	109.73.106.6	Unknown
linkedin.com	Social Networking	74.208.85.228	Unknown
live.com	Web-based Email	49.212.32.154	Unknown
twitter.com	Social Networking	178.32.54.90	Unknown
msn.com	Search Engines & Portals	116.254.188.24	High Risk
blogspot.com	Personal Sites	16.54.12.15	Unknown
adobe.com	Computers & Technology	140.113.207.143	Unknown

Fig. 5.6: Categorizzazione e reputazione Commtouch

dominio	Categoria Brightcloud	ip o dominio	Reputazione Brightcloud
google.com	Search Engines	depenam.com	High Risk
facebook.com	Social Network	barakair.com	High Risk
youtube.com	Streaming Media	18dd.net	High Risk
yahoo.com	Internet Portals	fullsfinesr.info	High Risk
baidu.com	Search Engines	servegame.org	Trustworthy
wikipedia.org	Reference and Research	twonext.com	Moderate Risk
qq.com	Internet Portals	109.73.106.6	Trustworthy
linkedin.com	Social Network	74.208.85.228	Trustworthy
live.com	Web based email	ygl.ru	High Risk
twitter.com	Social Network	spark29.ru	Low Risk
msn.com	Internet Portals	artvideo3d.ru	High Risk
blogspot.com	Personal sites and Blogs	4yoursecret.co.tv	Trustworthy
adobe.com	Business and Economy	dewell.ru	High Risk

Fig. 5.7: Categorizzazione e reputazione Brightcloud

I test effettuati sul sistema di categorizzazione hanno evidenziato dei risultati che differiscono da quelli dei sistemi commerciali. Tali differenze scaturiscono principalmente dal fatto che nel sistema è stato adottato uno schema di categorie differente da quello adottato nei sistemi testati.

Inoltre le differenze riscontrate nel sistema non sono da considerarsi degli errori, in quanto alcuni contenuti potrebbero essere categorizzati correttamente in più categorie.

dominio o IP	sistema	Trusted Source	Commtouch	Brightcloud
google.com	Computers and Internet	Search Engines	Search Engines	Search Engines
facebook.com	Social Network	Social Networking	Social Networking	Social Network
youtube.com	Arts and Entertainment	Streaming Media	Entertainment	Streaming Media
yahoo.com	Computers and Internet	Portal Sites	Search Engines	Internet Portals
baidu.com	Computers and Internet	Search Engines	Search Engines	Search Engines
wikipedia.org	Education and Learning	Education/Reference	Education	Reference
qq.com	Computers and Internet	Portal Sites	Search Engines	Internet Portals
linkedin.com	Business and Economy	Professional Networking	Social Networking	Social Network
live.com	Computers and Internet	Search Engines	Web-based Email	Web based email
twitter.com	Social Network	Social Networking	Social Networking	Social Network

Fig. 5.8: Confronto del sistema di categorizzazione con quelli commerciali

I test effettuati sul sistema di reputazione hanno evidenziato dei risultati ugualmente o maggiormente affidabili di quelli dei sistemi commerciali. Le risorse minacciose sono state tutte intercettate dal sistema e catalogate come “Spam and Virus”. I sistemi commerciali testati che hanno intercettato la maggior parte delle minacce hanno dei sistemi di catalogazione con dettaglio riguardo al grado di pericolosità.

Trusted Source ha rintracciato circa il 90% delle minacce che il sistema ha rintracciato nei test effettuati. Commtouch non fornisce nessuna informazioni riguardo i domini e ha rintracciato il 10% circa delle minacce rintracciate dal sistema. Brightcloud ha invece rintracciato il 70% circa delle minacce.

dominio o IP	sistema	Trusted Source	CommTouch	Brightcloud
depenam.com	Spam and Virus	High Risk (Web)	Nessuna informazione	High Risk
barakair.com	Spam and Virus	High Risk (Web)	Nessuna informazione	High Risk
18dd.net	Spam and Virus	High Risk (Web)	Nessuna informazione	High Risk
fullsfinesr.info	Spam and Virus	High Risk (Web)	Nessuna informazione	High Risk
servegame.org	Spam and Virus	Unverified (Web)	Nessuna informazione	Trustworthy
twonext.com	Spam and Virus	High Risk (Web)	Nessuna informazione	Moderate Risk
109.73.106.6	Spam and Virus	High Risk (Web & Mail)	Unknown	Trustworthy
74.208.85.228	Spam and Virus	Minimal Risk (Web & Mail)	Unknown	Trustworthy
ygl.a.ru	Spam and Virus	High Risk (Web)	Nessuna informazione	High Risk
spark29.ru	Spam and Virus	High Risk (Web)	Nessuna informazione	Low Risk

Fig. 5.9: Confronto del sistema di reputazione con quelli commerciali

5.4 Casi d'uso

La validazione del sistema ha evidenziato l'affidabilità dei risultati forniti, infatti i risultati sono paragonabili ed in alcuni casi migliori di quelli di alcuni sistemi commerciali analizzati. Per l'utilizzo del sistema sono stati implementati dei casi d'uso. In questo paragrafo vengono descritti anche allo scopo di mostrare gli ambiti di utilizzo, l'utilità e le funzionalità del sistema.

5.4.1 Caso d'uso della pagina web

Allo scopo di rendere i dati dei sistemi accessibili agli utenti, è

stata costruita una pagina web, che permette di interrogare separatamente il sottosistema di reputazione e quello di categorizzazione.

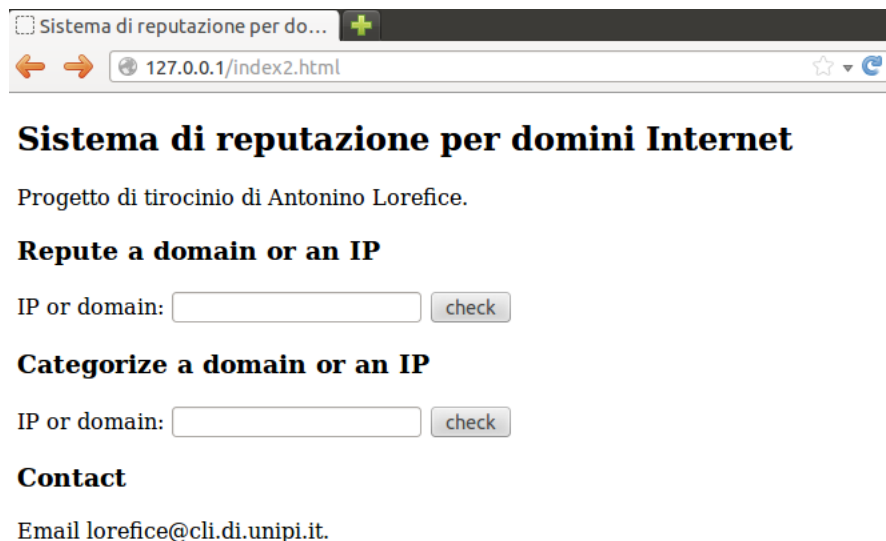


Fig. 5.10: Pagina web per l'utilizzo del sistema



Fig. 5.11: Esempio di reputazione del dominio barakair2.com



Fig. 5.12: Esempio di categorizzazione del dominio `repubblica.it`

5.4.2 Caso d'uso dell'interfaccia di accesso JSON

Il sistema mette a disposizione una interfaccia di accesso standard JSON per la comunicazione con le applicazioni. Tale interfaccia è stata integrata in Ntop[28], un applicazione per l'analisi ed il monitoraggio del traffico di rete dove viene utilizzata per la categorizzazione del traffico http.

5.5 Contatto e download del codice

Il codice sviluppato è disponibile gratuitamente all'indirizzo `http://www.cli.di.unipi.it/~lorefice/`. Per eseguirlo è necessario l'interprete Python 2.7 e l'installazione dei moduli contenuti nel file `"packages_you_need_to_install"`. Il codice è composto da degli script, uno per l'aggiornamento dei database di categorizzazione, da eseguire quotidianamente ed uno per l'aggiornamento dei database di reputazione, da eseguire ogni due ore. Il sistema di reputazione è

implementato nello script `repute_domain.py` e quello di categorizzazione nello script `categorize_domain.py`. Inoltre nel codice sono presenti degli eseguibili `c` per la gestione degli alberi utilizzati per la memorizzazione degli indirizzi IP.

Per qualsiasi informazione e chiarimento, contattare: lorefice@cli.di.unipi.it.

Capitolo 6

6 Conclusioni

Lo scopo del tirocinio è stato quello di costruire un sistema di reputazione libero e gratuito, che raccogliesse ed omogeneizzasse tutte le informazioni ed i servizi disponibili gratuitamente su Internet. L'obiettivo era quello di creare uno strumento che si interfacciasse con le applicazioni di analisi e di monitoraggio del traffico e che avesse un tempo di risposta tale da non rallentare eccessivamente la navigazione.

Si considerano raggiunti gli obiettivi iniziali del tirocinio, il sistema costruito reputa e categorizza correttamente le risorse con tempi di risposta paragonabili ai sistemi commerciali.

Il sistema è stato reso disponibile tramite interfaccia web agli utenti umani ed alle applicazioni. L'applicazione di analisi e monitoraggio del traffico Ntop lo utilizza nella parte che si occupa della categorizzazione del traffico http.

L'ambiente e le condizioni in cui è stato svolto il tirocinio sono state ottimali, il tutore si è sempre dimostrato disponibile ed

interessato, seguendo attivamente tutta l'attività svolta.

Il lavoro di sviluppo è stato svolto autonomamente sotto la costante supervisione del tutore per quanto riguarda le decisioni strategiche prese.

6.1 Sviluppi futuri

Le aziende che sviluppano software per la sicurezza informatica si stanno orientando sempre di più verso una gestione dei controlli di sicurezza da effettuare, basata sulla reputazione che ha un determinato dominio[30][29].

Per rispondere a tale esigenza il sistema di reputazione sviluppato potrebbe essere esteso e integrato in un sistema più complesso che fornisca diversi livelli di pericolosità dei domini e che in base alle informazioni raccolte valuti la possibilità di adottare misure di sicurezza maggiori, il tutto mantenendo un approccio aperto e gratuito.

Un'altra possibile implementazione sarebbe l'integrazione del sistema in netfilter, il meccanismo che implementa il firewall di linux. Netfilter è estendibile tramite plugin standard, un esempio è quello della temporizzazione delle regole, che permette ad esempio l'utilizzo di un protocollo di comunicazione solo in determinate ore della giornata. Queste regole sono sia di ingresso che di uscita del traffico dal pc.

Netfilter permette la realizzazione di uno stateful firewall, ovve-

ro un meccanismo che tiene traccia dell'appartenenza dei pacchetti alle comunicazioni e permette al sistema di ricordarsi le regole di trasmissione e di non doverle applicare per ogni pacchetto.

Il sistema sviluppato potrebbe essere utilizzato per la categorizzazione delle risorse web filtrate dal firewall, tenendo presente che le regole devono essere scritte basandosi allo schema di categorie adottato nel sistema.

Riferimenti bibliografici

- [1] Jaeyeon Jung and Emil Sit. 2004. An empirical study of spam traffic and the use of DNS black lists. In Proceedings of the 4th ACM SIGCOMM conference on Internet measurement (IMC '04). ACM, New York, NY, USA, 370-375. DOI=10.1145/1028788.1028838 <http://doi.acm.org/10.1145/1028788.1028838>
- [2] Loris Degioanni, Mario Baldi, Fulvio Risso, Gianluca Varenni. WinPcap: una libreria Open Source per l'analisi di rete.
- [3] Junghoo Cho and Sourashis Roy. 2004. Impact of search engines on page popularity. In Proceedings of the 13th international conference on World Wide Web (WWW '04). ACM, New York, NY, USA, 20-29. DOI=10.1145/988672.988676 <http://doi.acm.org/10.1145/988672.988676>
- [4] Konrad Rieck, Thorsten Holz, Carsten Willems, Patrick Dussel, and Pavel Laskov. 2008. Learning and Classification of Malware Behavior. In Proceedings of the 5th international conference on Detection of Intrusions and Malware, and Vulnerability Assessment (DIMVA '08), Diego Zamboni (Ed.). Springer-Verlag, Berlin, Heidelberg, 108-125. DOI=10.1007/978-3-540-70542-0_6 http://dx.doi.org/10.1007/978-3-540-70542-0_6
- [5] Smith, B. A Storm (Worm) Is Brewing.

- [6] antivirus.com. Trojan Horse or Trojan: It's Not All a Myth. <http://www.antivirus.com/security-software/definition/trojan-horse/index.html>
- [7] A Moshchuk, T Bragin, SD Gribble, HM Levy. A Crawler-based Study of Spyware in the Web.
- [8] Orathai Sukwong, Hyong Kim, and James Hoe. 2011. Commercial Antivirus Software Effectiveness: An Empirical Study. *Computer* 44, 3 (March 2011), 63-70. DOI=10.1109/MC.2010.187 <http://dx.doi.org/10.1109/MC.2010.187>
- [9] Saverio Verrascina, Daniele Gozzi, Mirco Marchetti. Architettura collaborativa per la rilevazione e l'analisi di malware.
- [10] Aubrey-Derrick Schmidt, Rainer Bye, Hans-Gunther Schmidt, Jan Clausen, Osman Kiraz, Kamer A. Yüksel, Seyit A. Camtepe, and Sahin Albayrak. 2009. Static analysis of executables for collaborative malware detection on android. In *Proceedings of the 2009 IEEE international conference on Communications (ICC'09)*. IEEE Press, Piscataway, NJ, USA, 631-635.
- [11] Ulrich Bayer, Andreas Moser, Christopher Kruegel, Engin Kirda. Dynamic Analysis of Malicious Code. *Journal in Computer Virology* August 2006, Volume 2, Issue 1, pp 67-77
- [12] Manuel Egele, Christopher Kruegel, Engin Kirda. Dynamic Spyware Analysis.

- [13] Amir Lev Commtouch Software Ltd. THE MARRIAGE OF SPAM AND MALWARE: IMPLICATIONS FOR SMTP MALWARE DEFENCE.
- [14] Sujata Garera, Niels Provos, Monica Chew, and Aviel D. Rubin. 2007. A framework for detection and measurement of phishing attacks. In Proceedings of the 2007 ACM workshop on Recurring malcode (WORM '07). ACM, New York, NY, USA, 1-8. DOI=10.1145/1314389.1314391 <http://doi.acm.org/10.1145/1314389.1314391>
- [15] Roger M. Needham. 1993. Denial of service. In Proceedings of the 1st ACM conference on Computer and communications security (CCS '93). ACM, New York, NY, USA, 151-153. DOI=10.1145/168588.168607 <http://doi.acm.org/10.1145/168588.168607>
- [16] Networking e sicurezza reti. http://www.8volante.com/sicurezza_reti.php
- [17] Ido Dubrawsky. Firewall Evolution - Deep Packet Inspection.
- [18] Almut Herzog, Nahid Shahmehri. Usability and Security of Personal Firewalls.
- [19] Informazioni generali sui firewall. <http://support.mozilla.org/it/kb/Informazioni%20generali%20sui%20firewall>
- [20] The Great Firewall: ecco come la Cina censura Internet. <http://www.terrefertili.net/2008/08/the-great-firewall-ecco-come-la-cina-censura-internet/>

- [21] Filtro famiglia. http://it.wikipedia.org/wiki/Filtro_famiglia
- [22] Dmoz.org: cos'è e come funziona la directory più autorevole del web? <http://www.newcomweb.it/blog/trucchi-e-curiosita-dal-web/article/dmoz-org-cos-e-e-come-funziona-la>
- [23] Google Safe Browsing. http://en.wikipedia.org/wiki/Google_Safe_Browsing
- [24] Alexa Internet. http://en.wikipedia.org/wiki/Alexa_Internet
- [25] The Spamhaus Project. http://en.wikipedia.org/wiki/The_Spamhaus_Project
- [26] Norton Safe Web. http://en.wikipedia.org/wiki/Norton_Safe_Web
- [27] AVG Threat Labs: controllo siti internet preventivo. <http://www.pctuner.net/news/14368/AVG-Threat-Labs-controllo-siti-internet-preventivo/>
- [28] L. Deri and S. Suin. 2000. Effective traffic measurement using ntop. *Comm. Mag.* 38, 5 (May 2000), 138-143. DOI=10.1109/35.841838 <http://dx.doi.org/10.1109/35.841838>
- [29] Stefanie Hoffman. Reputation Scoring: A Step Ahead Of Malware. <http://blog.fortinet.com/Reputation-Scoring--A-Step-Ahead-Of-Malware-/>
- [30] The Importance of Client Reputation. http://www.fortinet.com/resource_center/whitepapers/importance_client_reputation.html