Master-thesis

# What web objects contained in popular websites are delivered through hypergiants.

Fakultät IV - Elektrotechnik und Informatik
Intelligent Networks / Intelligente Netze (INET)
Research Group of Prof. Anja Feldmann, Ph.D.

Soumya Ranjan Parida
October 8, 2016

Prüfer: Prof. Anja Feldmann, Ph. D.
Betreuer: Ingmar Poese
Juhoon Kim

Die selbständige und eigenhändige Anfertigung versichere ich an Eides Statt.

Berlin, den October 8, 2016 Max Mustermann

# Zusammenfassung

Mit der Verbreitung des Internets, Hyperriesen wie Google, Content-Lieferungs-Netzwerke wie Akamai usw. spielen häufig eine wichtige Rolle in den Inhalt einer Website bereitstellt. Diese Hyperriesen nicht nur verschiedene Dienste zur Verfügung stellen, sondern auch reich an Inhalt. Flash-Medien von Youtube, Login-System von Google, Facebook, Werbung von Google Ad Sinn, populäre Social-Media-Sites wie Facebook, Twitter, LinkedIn, etc. sind weit verbreitet und beliebte Dienste in den meisten Websites heute eingebettet.

Um mit dieser Nachfrage zu bewältigen, Hyperriesen haben eine große Anzahl von skalierbaren und kostengünstigen Hosting und Content-Delivery-Infrastrukturen auf der ganzen Welt einsetzen. Diese Hosting-Infrastrukturen von wenigen großen Rechenzentren, eine große Anzahl von Cache-Speicher oder eine beliebige Kombination zusammengesetzt sein. Ein solches Szenario sowie riesige Abhängigkeit zwischen populären Webseiten und Hyperriesen eine große Menge des Verkehrsflusses von Hyperriesen führen.

Um zu wissen, wie die Einbindung von beliebten Webseiten mit Hyperriesen weiterentwickelt, ist, befasst sich diese These die folgenden Forschungsfragen .Firstly gibt es jede Anwesenheit dieser Hyperriesen in interenet Architektur? Wenn Hyperriesen vorhanden sind, dann wie viel Prozent der populären Web-Sites von 100.000 Top-Websites von alexa sind mit verschiedenen Hyperriesen verbunden? Drittens Welcher Prozentsatz von verschiedenen Objekten (text / html, img, Skript, Medien etc.) werden mit Hyperriesen verbunden?

Die vorliegende Arbeit liefert eine quantitative Forschung von Web-Verbindungen für alexa Top 100.000 Websites. Wir präsentieren Ihnen die Planung, Durchführung und Analyse von verschiedenen Arten von Objekten in verschiedenen Websites enthalten, die zu unterschiedlichen Hyperriesen verbunden sind. Die Arbeit folgt zwei Pfaden. Zum einen, um verschiedene Arten von Objekten in Webseiten zu beziffern, werden wir Homepages von oben 100.000 Websites von alexa Webseite und sammeln das Vorhandensein von Hyperriesen Infrastruktur verknüpft verschiedene Objekte wie Bilder, externe Links, Skripte, eingebettete Videos, CSS-Dateien kriechen usw. in diesen Websites. Zweitens untersuchen wir die Objekte, um den Grad der Verbindung zwischen den oberen 100.000 Websites und hypergiant zu erfahren.

Die experimentellen Ergebnisse dieser Arbeit diskutiert werden durch umfangreiche Analysen der gesammelten Daten unterstützt, die in dieser Arbeit vorgesehenen Nachweise zur Stützung des Abschlusses zur Verfügung stellen.

# Abstract

With the proliferation of the Internet,hyper giants such as Google ,Content delivery networks like Akamai etc. often play a vital role in providing the content of any website. These hyper giants not only provide different services but also rich in content. Flash medias from Youtube, login system from Google,Facebook, advertisements from Google ad sense,popular social media sites like Facebook,Twitter,LinkedIn etc. are common and popular services embedded in most of the websites today.

To cope with this demand, hyper giants have been deploying a large number of scalable and cost effective hosting and content delivery infrastructures all around the globe. These hosting infrastructures can be composed of a few large data centers ,a large number of caches or any combination. Such a scenario cause a large amount of traffic flow from hyper giants as well as huge dependency between popular websites and hyper giants.

In order to know how the involvement of popular websites with hyper giants is evolving,this thesis addresses the following research questions .Firstly,are there any presence of these hyper giants in Internet architecture ? If there are hyper giants present,then what percentage of popular web sites out of 100,000 top websites of Alexa are connected with various hyper giants ? Thirdly ,What percentage of different objects (text/html ,img ,script,media etc.) are connected with hyper giants ?

This thesis provides a quantitative research of web connectivity for Alexa's top 100,000 websites. We present the design, implementation and analyses of different types of objects contained in various websites which are connected to different hyper giants . The thesis follows two paths. Firstly ,in order to quantify different types of objects involved in websites ,we will crawl home pages of top 100,000 websites of Alexa's website and gather the presence of hyper giants infrastructure linked to different objects like images,external links,scripts ,embedded videos,css files etc. in those websites. Secondly, We examine those objects to find out the degree of connection between the top 100,000 websites and hyper giant.

The experimental results discussed in this thesis are supported by extensive analyses of data collected which provide evidence in support of the conclusion provided in this thesis.

# Contents

# 1  Introduction

The Internet has changed a lot within last decade both in technical as well as in user experience aspects.With increase in Internet users and their demand for more and richer content has led to exponential increase of Internet traffic.Social networking sites like Facebook, Twitter enable users to publish their own content and share with other users.Users also share videos in different social media sites like Youtube,Facebook etc.The highly popular on-demand video and streaming sites like Netflix etc., are also playing vital role in increasing Internet user base and traffic.Recent traffic studies [3,4] show that a large fraction of Internet traffic is originated by a small number of prominent infrastructure which can be highly distributed content delivery networks (CDNs) like Akamai or content providers like Google.Poese et al. [3] report a similar observation from the traffic of a European Tier-1 carrier. Labovitz et al. [2] infer that more than 10 % of the total Internet inter-domain traffic originates from Google, and Akamai claims to deliver more than 20 % of the total Web traffic in the Internet [5].

Traditional hosting model like CDNs are the most important technical solutions for providing high performance delivery system till now where popular contents are stored in servers of CDNs .But with the increase of content within sites ,it is not possible for the popular web sites to provide better performance to end customers by using only the traditional hosting model.Instead ,content providers now build their own global backbones, cable Internet service providers offer wholesale national transit, and transit ISPs offer CDN and cloud / content hosting services.CDNs also build highly distributed infrastructures and data centers to replicate the most popular content at different distributed cache servers and locate them at the edge of the network.It help them to provide popular content from the nearest server to customers.Hence when a user request for a popular web content ,CDN just redirect the user to most suitable server by bypassing the saturated links.

Hence due to this change in content delivering phenomenon some researchers termed this companies as hyper giants [7] which include large content providers, such as Google and yahoo, as well as highly distributed CDNs like Akamai,big cloud computing CDNs like Amazon aws etc. Most of these hyper giants are operating not only a substantial number of data centers but are also building up their own network. Some networking researchers are claiming that, due to the phenomenal growth of hyper giants, the topological structure of the Internet must be redrawn to include them, together with the Global transit and backbone networks as part of the Internet core,resulting in the topology sketched in Figure 2. This may leave the ISPs as dump pipe providers to the consumer.

Again it is important to understand that hyper giants are not usually the main operators of the network but they play vital role in delivery of the content by creating interdependency between them and the main operators by different ways and business needs. The producers of the content (popular websites) want their content to be delivered to end user in less time for which they have to rely on the main operators or hyper giants. Such a scenario cause a large amount of traffic flow from hyper giants as well as huge dependency between popular websites and hyper giants.It is this symbiosis between the two parties that motivates our work ,giving an overview on how far the reach of hyper giants in todays Internet.

Again hyper giants now not only provide rich content,they also provide different other services.For example now most of the popular websites having their own login systems or login systems which are provided through Google,Twitter,Facebook etc.In later case the authentication is verified by Google,Facebook etc.For that the popular websites need to embed third party login systems in their websites.Like this there are lot of other services provided by these hyper giants like for websites add Google adsense to advertise their

products etc.Moreover popular websites need CDNs to store web contents like HTML files,audio files,video files etc. This dependency of popular websites on hyper giants also gives us motivation to check what percentage of these web objects are delivered through hyper giants to popular websites.

A few studies have already investigated about hyper giants and their relationship with popular websites in the recent past.In 2010, Craig Labovitz, then of Arbor Networks [2],defined a new type of network entity. By placing google in this list,he characterized the hyper giant as a content provider that makes massive investments in bandwidth, storage, and computing capacity to maximize efficiencies and performance.The concept of hyper giants also aligns with Schmidts [8] assertion which talks about "gang of four" companies which are responsible for the growth and innovation of Internet. Google, Apple, Amazon, and Facebook .Bernhard et al.[9,10] also worked on identifying and mapping the content infrastructures that are hosting the most popular content.Th author also purposed a light weight automated technique discover Web content hosting and delivery infrastructures.He identified different types of infrastructures like highly distributed CDNs,CDNs,data centers and hyper giants presence in Internet.Gao et al.[6] analyzed operator interconnections from a more technical perspective.They used a methodology to quantify the type of inter-Autonomous System (AS)relationships that exist in the Internet and classify them into three groups based on the state of Border Gateway Protocol (BGP) messages: customer-to-provider, peer-to-peer, and sibling-to-sibling relationships.Shavitt and Weinsberg [11] in their paper discussed the topological trends of content providers. They create a snapshot of the AS-level graph from late 2006 until early 2011, and then analyzed the interconnection trends of the transit and content providers and their implications for the Internet ecosystem. AS graphs are built by traversing IP trace routes and resolving each IP address to its corresponding AS.Shavitt and Weinsberg also found that large content providers like Google, Yahoo!, Microsoft,Facebook, and Amazon have increased their connectivity degree during the observed period and are becoming key players in the Internet ecosystem, strengthening the idea that the Internet is becoming flatter.Palacin et [13] defined hyper giants not only content providers, they are basically content aggregators. Small companies started using high speed infrastructures to deliver their content to end users.But with increase of content these high speed infrastructures started absorbing content from the long tail, entering fully into the niche of the traditional hosting companies.

By end of this thesis we are able to give answers for some of the important research questions which can be summarized as follows:

- Identification of hosting infrastructures: We propose a lightweight and fully automated approach to discover hyper giants such as highly distributed content delivery networks,content providers etc.

- Classification of hosting infrastructures: We classify individual hosting infrastructures and their different deployment strategies based on their network.

- Web content dependency: We quantify the degree of content dependency of the popular websites on hyper giants by analyzing different web objects like text files,image files,application files delivered by hyper giants to popular web sites.

This remainder of this thesis is structured as follows.This thesis is separated into 7 chapters.

- Chapter 2 :It starts with the evolution of Internet architecture from early 2000s to current time and how the dependency of popular websites on hyper giants increases with time.

- Chapter 3 :This section will provide the overall methodology used in this thesis.

- Chapter 4 :This section focuses on the details about environment and technologies used for the prototype. The implementation details of web crawler engine, its operations and configuration management are explained.

- Chapter 5 describes the measurement details.

- Chapter 6 summarizes the results.

- Chapter 7 Conclusion will be discussed here..

- Chapter 8 This section includes the possible future work.

# 2 Background

In this chapter,we discuss how Internet architecture evolved with time.Along with this we will discuss briefly on hyper giants.In addition to this we also provide the techical background of DNS and HTTP protocol which will be used in this thesis.

## 2.1 Evolution of Internet Architecture and Rise of hyper giants

In 2010, Craig Labovitz, then of Arbor Networks [2],defined a new type of network entity he argued transcended traditional "content versus carrier" dichotomy of Internet architecture.By placing Google in this list,he characterized the hyper giant as a content provider that makes massive investments in bandwidth, storage, and computing capacity to maximize efficiencies and performance.The concept of hyper giants also aligns with Schmidt's assertion which talks about "gang of four" companies which are responsible for the growth and innovation of Internet. Google, Apple, Amazon, and Facebook [8].

The Internet architecture implemented until the early 2000s was based on a multi-tier hierarchic structure.Tier 1 ISPs were on top of the hierarchy followed by the Tier 2 regional ISPs and the Access ISPs at the lower part of the hierarchy connecting the end users. In this scheme, Tier 1 ISPs were highly connected to other ISPs and offered transit services to other ISPs in lower layers.Content was distributed through Access ISPs or, in the best cases, through ISPs located at advantageous points. Traffic flows were required to go up and then down in the hierarchy to reach end users shown in figure 1.Among the different network operators,Internet traffic was exchanged at different IP exchange points according to agreements between different layer players where the dis symmetry in traffic was compensated.
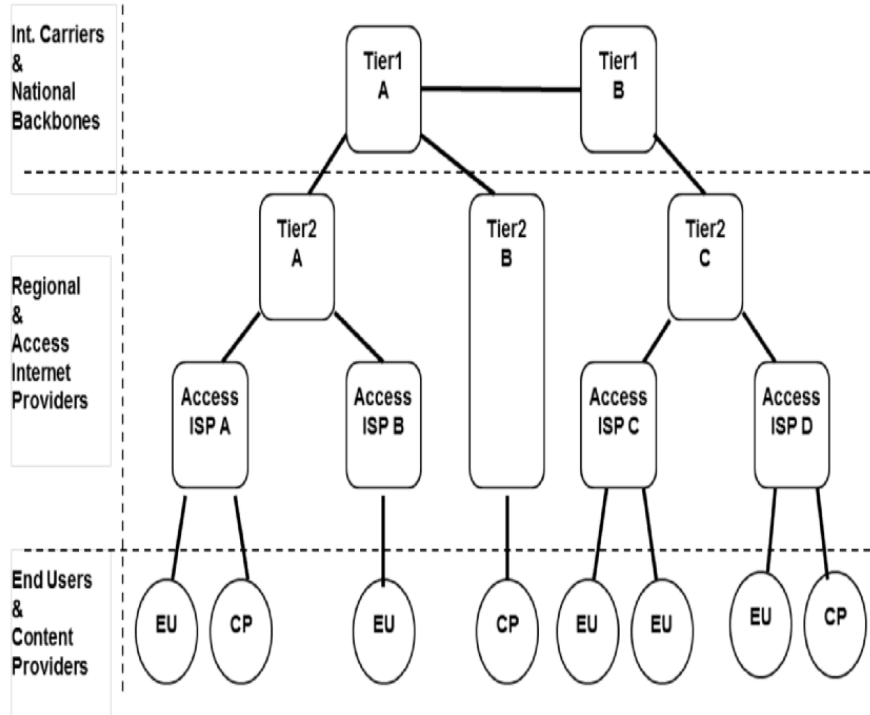


Figure 1: Traditional Hierarchic Internet Structure

But with the time,the Internet architecture has changed.Researchers found that now
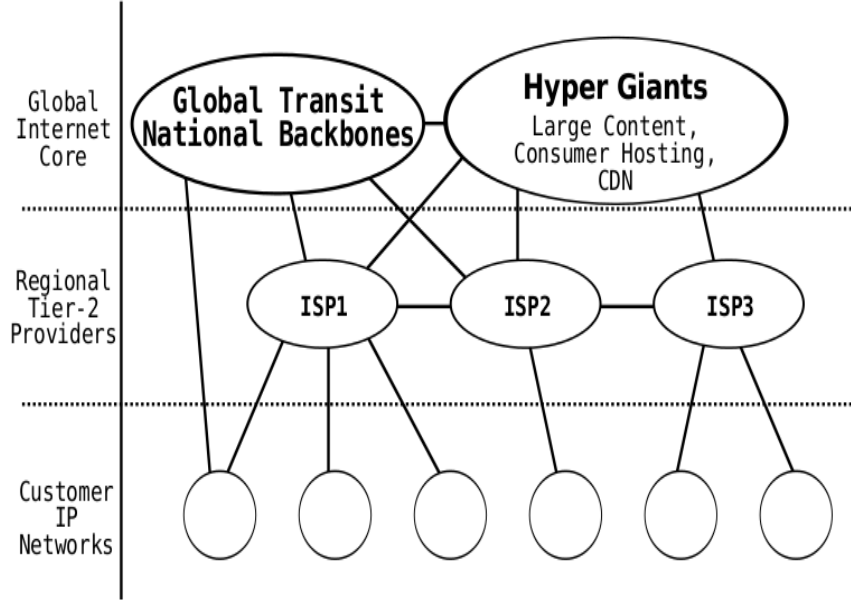
Figure 2: Modern Internet Structure

nobody has control over Internet ,instead each ISP has control over its network and depend upon the network connected with it.Even during last decade the old pyramidal structure of Internet architecture shown in figure-2 has been bypassed by big content providers, such as Google, Facebook, Amazon or Yahoo!, and content delivery network operators, such as Akamai.As a result now Internet's backbone has a flatter structure where there are few autonomous systems are playing major role in delivering content.They are connected to each other and have a big footprint by establishing small data centers all over the world.This help them to get as close as possible to the access networks used by their customers,bypassing intermediate Internet service providers.The trend towards flatter network architectures can also be found in the area of access networks.The researchers termed these infrastructure providers as "Hyper Giants" which include large content providers, such as Google and Yahoo, as well as highly distributed CDNs, like Akamai.

## 2.2   Content delivery Infrastructures

Recent traffic studies [3, 14] show that a large fraction of Internet traffic is due to content delivery and it originated by very few content delivery infrastructure (CDIs).Major CDIs include content delivery networks like Akamai,Cloudflare, content providers like Google,Microsoft,highly popular rich media sites like Youtube,Netflix and cloud computing infrastructures like Amazon aws.Most of the CDIs have a large number of servers which located across the world.CDIs cache most of the popular content of the websites at these servers.Hence we a end user request for any content,CDIs deliver the content from the server nearest to the end user.In this way CDIs reduce the load on origin servers and at the same time improve performance for the user.  CDIs follow different strategies for redirecting traffic to nearest cache server of the end user.Most of the CDIs use DNS to translate the host name of a web site request into the IP address of an server.During this

translation ,DNS takes into account the location to the end user,the location of the nearest CDIs cache server ,load of the server etc.

Independent CDIs are normally referred to as CDNs.CDNs have a large number of servers all around the world and mainly responsible for delivering content of their customers to end users.Leighton [1] proposes four main approaches to distributing content in a content-delivery architecture: (i) centralized hosting,(ii) big data center based CDNs, (iii) highly distributed CDNs, and (iv) peer-to-peer networks.In centralized hosting case ,traditional architectures sites take help of one or small number of collocation sites to host their content.These centralized hosting structures are may be enough for small content sites but not for popular websites which carries huge amount of content.Big Data Center content delivery networks have Hugh number of high-capacity data centers which are connected to major backbones.Highly popular content are cached.Hence able to increase the performance of delivery compare to centralized hosting infrastructures but still are limited in potential improvements as still they are far away from the end user.Third type of model is highly distributed content delivery networks.They have high footprint all over the world.By putting their own infrastructures inside end user's ISP,they are able to eliminating peering, connectivity, routing, and distance problems, and reducing the number of Internet components.Final approach is peer to peer networks which has very little scope in delivering the content of popular websites in today's Internet world due to serious concern of the copy right issues.

Cloud infrastructure refers to the hardware and software components, such as servers, storage, networking and virtualization software that are needed to support the computing requirements of a cloud computing model. In addition cloud computing infrastructures include a software abstraction level which virtualizes resources like servers, compute, memory, network switches, firewalls, load balancers and storage. and logically presents them to users through programmatic means.Cloud infrastructure mainly present three different types of model:infrastructure as a service (IaaS), platform as a service (PaaS) and software as a service (SaaS).Cloud infrastructures deploy a large number of data centers at certain regions of the world.In case of infrastructure as service or IaaS,cloud infrastructures give access to these data centers to their users.Users can able to access and manage remote data center infrastructures, such as compute (virtualized or bare metal), storage, networking, and networking services (e.g. firewalls).SaaS uses the web to deliver applications that are managed by a third-party vendor and whose interface is accessed on the client's side.Popular SaaS offering types include email and collaboration, customer relationship management, and health care-related applications. Paas is used for applications, and other development, while providing cloud components to software. With this technology,users can manage OSes, virtualization, servers, storage, networking, and the PaaS software itself.

Content Providers also are major player in content delivery infrastructure.Companies like Google, Facebook, Netflix etc. build their own infrastructures like data centers and interconnected them with high speed backbone networks to deliver some of their very popular services.Google connects its data centers to a large number of ISPs via IXPs and also via private peering [15].Google also now provide Google Global Cache (GGC) [16] as a service where customers can optimize network infrastructure costs associated with delivering Google and YouTube content to end users by serving this content from inside their ISP networks.Through GGC,small ISPs and which are located in areas with limited connectivity can reduced the transit cost as well as websites can deliver their content with better performance.GGC also allows an ISP to advertise through BGP the prefixes of users that each GGC server should serve.The Netflix system known as Open Connect Network[17].Netfix deploy its own infrastructure inside a lot of ISPs by partnering with

them to deliver its own content more efficiently.

## 2.3 Protocols

In this section we will discuss about the protocols used in this thesis which are domain name system (DNS) and hyper text transfer protocol (HTTP).Both protocols used in out thesis extensively to get CNAMEs of a website and to get the http header information respectively.

### 2.3.1 Domain name System (DNS)

Domain name system (DNS) is used to translate IP address to corresponding host names.Internally it is maintaining a hierarchal structure of domains.Before the invention of DNS on year 1983,a simple text file (hosts.txt) file was used to do this translation from IP address to host name.But with a growing number of host names it was difficult to keep maintain in hosts.txt file and Domain name system introduced.

The administration of domains is divided into different zones. The zone information is distributed using authoritative name servers.The top most level of DNS starts with root zone and the root zone information is served by root servers.Responsibility of specific parts of zone can be given to some other authoritative name servers which in turn divided responsibility with other authoritative name server.For, e.g.,the responsibility of .org domain is delegates to the Pub- lic Interest Registry by the root zone which in turn delegates responsibility for acm.org to the Association for Computing Machinery (ACM).At the end its site is responsible for its own zone and keep maintain its own database of authoritative name server.The information about a particular domain of a zone is kept in Resource Records (RRs) which specify the class and type of the record as well as the data describing it.Multiple RRs with the same name, type and class are called a resource record set (RRset).

To resolve a IP address into host name,the procedure starts with the end user's stub resolver queries to local name server called caching server.if caching server can not able to resolve it,it redirects the query to authoritative name server of the domain.If resolver does not know how to contact the corresponding authoritative name server of the domain,it redirects the query to root name server .The root name server again refers the resolver to the authoritative name server responsible for the domain just below the root server.This procedure continues till resolver is able to resolve the domain properly.

Figure-3 shows the DNS reply by the resolver when querying a host name served by a content infrastructure.Here the host name is www.bmw.com.The answer section contain a chain of CNAMEs which resolve into two ARecord set (RRset) with different IP addresses which can be for used for load balancing.

### 2.3.2 Hyper text transfer protocol(HTTP)

Hyper text transfer protocol (HTTP) is an application layer protocol mainly used as defector standard to transport content in world wide web.HTTP works on top of the TCP/IP protocol and follows the client server architecture via request-response communication procedure.It allows end-users to request, modify, add or delete resources identified by Uniform Resource Identifiers (URIs).

HTTP message consists of HTTP header which shows the meaning of message and HTTP body which is actual message.HTTP message can be a request message or response message.The HTTP client sends a request message to server .There are different types of methods used in HTTP request message like GET,HEAD,POST ,PUT,DELETE,CONNECT

```
; <<>> DiG 9.10.3-P4-Ubuntu <<>> www.bmw.com
;; global options: +cmd
;; Got answer:
;; ->>HEADER<<- opcode: QUERY, status: NOERROR, id: 11656
;; flags: qr rd ra; QUERY: 1, ANSWER: 4, AUTHORITY: 0, ADDITIONAL: 1

;; OPT PSEUDOSECTION:
; EDNS: version: 0, flags:; udp: 512
;; QUESTION SECTION:
;www.bmw.com.                    IN      A

;; ANSWER SECTION:
www.bmw.com.            3600    IN      CNAME   cn-www.bmw.com.edgesuite.net.
cn-www.bmw.com.edgesuite.net. 3600 IN   CNAME   a1586.b.akamai.net.
a1586.b.akamai.net.     20      IN      A       104.121.76.73
a1586.b.akamai.net.     20      IN      A       104.121.76.64

;; Query time: 22 msec
;; SERVER: 127.0.1.1#53(127.0.1.1)
;; WHEN: Wed Oct 05 21:32:17 CEST 2016
;; MSG SIZE  rcvd: 143
```

Figure 3: DNS Reply for a host using dig command line tool

etc.But in this thesis we have used extensible GET method and the HEAD method.The GET method is used to retrieve information from the given server using a given URI. Requests using GET should only retrieve data and should have no other effect on the data.Same as GET, but it transfers the status line and the header section only.The introductory line in an HTTP request shown in figure 4 consists of a method, a server-side path, and the HTTP version in use.The introductory line in an HTTP response shown in figure 4 starts out with the HTTP version in use, followed by a standardized three-digit status code and a textual status description. The status code tells the requester about the success of the query or indicates the reason of an error.Both request and response messages are followed by multiple header lines.Some header information are valid for request ,some are for response and some are valid in both the ways.Since HTTP1.1,the Host header is mandatory for request messages.The meta information encompasses information about the file type, the character set in use, preferred language etc.HTTP also allows server to set cookies in client side which help the server to track client requests.

## 2.4   Conclusion

Within last decade the Internet architecture changed vastly due to the introduction of hyper giants which can be highly distributed CDNs,cloud computing CDNs etc.Todays Internet traffic is dominated by HTTP traffic.Again to deliver the content fast ,DNS protocol is used as the load balancing mechanism by these big hyper giants.

```
HTTP Request

GET / HTTP/1.1
Host: www.example.com
User-Agent: Mozilla/5.0 [...]
Accept: text/html [...]
Accept-Language: en-us
Accept-Encoding: gzip, deflate
Connection: Keep-alive

HTTP Response

HTTP/1.1 200 OK
Accept-Ranges: bytes
Content-Type: text/html
Date: Mon, 27 Jul 2009 12:28:53 GMT
Server: Apache/2.2.14 (Win32)
Last-Modified: Wed, 22 Jul 2009 19:15:56 GMT
Content-Length: 88
<!doctype html>
<html>
[...]
```

Figure 4: HTTP Request (top) and Response (down) for example.com

# 3    Methodology

In this section we will discuss about the whole approach to identify the hyper- giants presence in today's Internet.The key idea is to collect IP addresses that DNS returns for Alexa top 100,000 websites and the website links embedded in these top 100,000 websites.We will use the content type of all the HTTP request and find out the different object type (text/html,image,video,audio etc) delivered by these hyper giants for the popular websites and the embedded links in those websites.

To achieve our goals of identifying and classifying hyper giants in Internet,we divided the whole process into 6 parts shown in figure 6.We now elaborate the steps we followed ,choices we have taken to achieve our goals.

## 3.1    Hosting Infrastructure Coverage

To achieve our goals to find out the hyper giants in Internet,we will use the SLD infrastructures which are serving popular websites.Given that Internet traffic normally consistent with the Zipf's law [22],it is very likely possible that the hyper giants are the very popular content delivery networks or content providers etc. which are responsible for the majority of the HTTP traffic flow in today's Internet.For example, Akamai claims to deliver about 20% of the total Web traffic in the Internet [5]. According to Labovitz et al. [2] ,Google serves up to 10 % of all Internet traffic.Similarly top 10 hosting infrastructures serve more than 15% and top 100 responsible for more than 40% of the traffic.In this thesis we will take the Alexa top 100,000 popular websites/
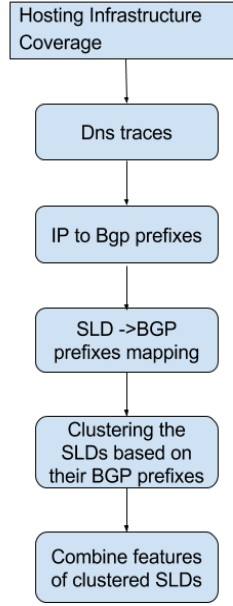
Level Of Approach.png



Figure 5: High Level Of Approach

## 3.2 DNS traces

We base our study on DNS traces collected within single vantage point within Germany.So this thesis will provide a better overview of hyper giants behavior within our vantage point location in Germany.As hyper giants are normally vast hosting infrastructures,content providers etc. ,their footprint mostly present all over the world .Hence the identification of hyper giants should not be affected much due to this limitation.But there are few other research questions arises due to this network coverage limitation .We will discuss about these research questions in "Future Work" chapter.For collecting data we use "Scrapy" framework.Scrapy framework is an open source web crawler.

## 3.3 IP to BGP Prefixes Mapping

We extract IP addresses of each website link using scrapy crawler.Along with this we also collect the ARecord names associated with the corresponding IP address.From ARecord names we collect the SLDs.The set of IP addresses for a perticular SLD shows the degree to which the corresponding SLD infrastructure is network-wise and geographically distributed.Hence natural choice for the features to consider are IP addresses,AS Numbers and the bgp prefixes.The number of bgp prefixes shows the network footprint and the number of ASNs shows how infrastructures are distributed all over the world.For example the highly distributed infrastructures have a lot of prefixes as well as high number of ASNs.Similarly small data centers will be located within a single AS , having a limited number of bgp prefixes,and a large number of IP addresses.Eventually these features are correlated and using these features we will try to find out the hyper giants presence in Internet. To determine bgp prefixes we use BGP routing information from RIPE RIS[23].The bgp prefix information for IP address 104.121.76.73 is shown in the figure

```
% This is RIPE NCC's Routing Information Service
% whois gateway to collected BGP Routing Tables, version2.0
% IPv4 or IPv6 address to origin prefix match
%
% For more information visit http://www.ripe.net/ris/riswhois.html
%
% Connected to backend ris-whois06.ripe.net

route:         104.64.0.0/10
origin:        AS35994
descr:         AKAMAI-AS - Akamai Technologies, Inc., US
lastupd-frst: 2016-10-04 11:53Z  198.32.176.70@rrc14
lastupd-last: 2016-10-07 11:29Z  196.46.25.29@rrc19
seen-at:       rrc00,rrc01,rrc03,rrc04,rrc05,rrc06,rrc07,rrc10,rrc11,rrc12,
               rrc13,rrc14,rrc15,rrc16,rrc18,rrc19,rrc20,rrc21
num-rispeers: 161
source:        RISWHOIS

route:         104.121.76.0/24
origin:        AS20940
descr:         AKAMAI-ASN1 , US
lastupd-frst: 2016-07-06 10:09Z  198.32.124.146@rrc16
lastupd-last: 2016-10-07 11:29Z  196.46.25.29@rrc19
seen-at:       rrc00,rrc01,rrc03,rrc04,rrc05,rrc06,rrc07,rrc10,rrc11,rrc12,
               rrc13,rrc14,rrc15,rrc16,rrc18,rrc19,rrc20,rrc21
num-rispeers: 163
source:        RISWHOIS
```

Figure 6: RIPE RIS bgp prefixes using whois command

7.As we can see from figure-7, the information provided are like routes which is the bgp prefixes,origin shows the corresponding AS number.From the figure this can be seen that the ip address 104.121.76.73 can be routed via two different prefixes 104.64.0.0/10 and 104.121.76.0/24.So we will consider both the prefixes for the IP address 104.121.76.73.

## 3.4   SLD to BGP Mapping

From the second stage shown in figure-6,we receive IP address and corresponding ARecords.From the ARecord,we get second level domain which gives us idea about the infrastructure involved with the website.Again from stage-3 shown in figure-6 ,we get the IP address to BGP prefixes mapping using RIPE RIS bgp prefixes.In this stage we will create the mapping of SLD infrastructures collected in stage -2 shown in figure-6 to BGP prefixes collected in stage-3 shown in figure-6.Now we have all features IP addresses,BGP prefixes,ASN numbers which are required for determining the type of infrastructure.

## 3.5   Clustering Algorithm

The clustering algorithm can be divide into two parts.In the first step ,we will try to aggregate the prefixes of a SLD into the set of parent prefixes and in the subsequent step we will compare prefixes of the SLDs with each other and cluster them if they are sharing most of the infrastructures.

### 3.5.1   Aggregate Prefixes

From the above steps,we collected SLDs with corresponding number of links served by SLD,number of IP addresses,bgp prefixes and number of ASNs.Now first we will sort the

SLDs according to their number of links in decreasing or- der.This will give us all the pop-ular SLDs which served maximum number of links.Now from the bgp prefixes,keep only parent prefixes and the prefixes which do not have any child in the prefix set.for example, googledomains.com has the prefix set ('216.239.32.0/19', '216.239.32.0/24', '216.239.34.0/24', '216.239.36.0/24'), '216.239.38.0/24').After the above step the prefix set will become ('216.239.32.0/19']) as other prefixes are subset of parent prefix 216.239.32.0/19.This procedure will be done for all the SLDs staring from first SLD sorted by decreasing order of number of links.Now starting from first and compare each SLD with the other SLDs.Between two prefix sets if child prefixes present then replace with parent prefix.This will make two prefix set with only parents.Now compare the two sets of prefixes to find out whether they are sharing same infrastructure or not.More details about the similarity procedure will be discussed in the next section.

### 3.5.2 Similarity between two prefixes set

Based on the similarity between two prefix set,we will decide if they belong to same SLD infrastructure or not.For this we define the similarity between two prefix set as follows,

$$similarity(s1, s2) = \frac{|s1 \cap s2|}{|s2|} \tag{1}$$

where s1,s2 are the bgp prefix sets.

If the similarity between two prefix sets are greater than equal to 70% then we cluster both the SLDs together.Here we assume than if s2's 70% prefixes are present in the common infrastructure between s1 and s2 set,then it shows that s2 sharing most of the infrastructure of s1.Hence we club them together.We will continue this procedure for all the SLDs starting from the first SLD sorted by their number of links. If two SLDs are matched ,they it will be removed from comparison with other SLDs and it will be mapped to the SLD with which it matched.For example 'googleusercontent.com' matched with 'google.com' with similarity 100%.Once this is matched we clubbed googleusercontent.com with google.com and will not be available for any further similarity matching with other SLDs.70% of similarity is chosen after extensively testing between bgp prefixes.But this can be taken for future work.

## 3.6 Combine features of clubbed SLDs

Now after the cluster algorithms, we will get prominent SLD infrastructures under which multiple number of SLD mapped to based on similarity comparison between the bgp prefixes they routed to.Hence now we will club the features of the child SLDs with the parent SLDs.After this stage we have all the prominent infrastructures with number of unique bgp prefixes by both parent SLD and child SLDs,unique links by both parent SLD links and child SLDs ,unique ASNs both parent SLD links and child SLDs.

## 3.7 Conclusion

In this section we discussed about the whole method we are going to follow to reach our goals.Along with this we identify some of the future works which will be discussed briefly in "Future work" section.

# 4 Implementation

This section describes how the first prototype of an crawler was built. It explains the choice of programming language as well as go into some details of the objects, the rough internal working and their interaction with each other. In the end, the usage of the program will be explained briefly.

## 4.1 Choice of programming language

Before going for implementation,it is important to know that the whole process involve two major parts.In the first part,we use a web crawler to collect popular website links as well as embedded links.And in second part we are going to analysis these crawled data to find out our result.Hence we need to choose a programming language which should be powerful in development as well as can be used for data analysis.So that any kind of dependency between two procedure can be handled easily.Again language should be free, open and usable by anyone who wishes to run this crawler implementation.Again the language should have good number of open libraries,big community who can help in necessary and should have a lot of documentation available in Internet .For this purpose we choose python as it fulfills all necessary requirement. For development we are going to use python and for data analysis we will use python,pyspark (Python version of spark).

## 4.2 Development Tools

For development,it is important to choose correct IDE which allows us to write code in less time with minimum effort.Along with this it helps in code completion, syntax highlighting, debugging and refactoring.For this purpose we have chosen "sublime editor" which is free and easy to write python code.

Now second most important aspect is to choose python web crawling tool which will suit our requirement and can be used in future.We will focus on programs that request web services from service providers and programs that scrape data from web sites. Web service applications will involve us in a new kind of programming called client-server programming; the programs we will look at will be client programs making requests from service on the Internet. Although the underlying foundation of a web-scraping program is also a client-server interaction, we will use some tools that hide the details of those interactions, and allow us to fetch web page content directly.For this we have couple of choices which are well known web crawling tools like urlib,beautiful soup,scrapy etc. But Python Scrapy is the best out there, Scrapy crawling is faster than any other platforms, since it uses asynchronous operations (on top of Twisted). Scrapy has better and faster support for parsing (x)html on top of libxml2. Scrapy is a mature framework with full unicode, redirection handling, gzipped responses, odd encodings, integrated HTTP cache etc.Again it is open source having a big community and documentation.

There are lot of different machine learning tools available but as we expect out data to be of some gigabytes,it is better to use that tool which can process data faster.Python and R are popular languages for data scientists due to the large number of modules or packages that are readily available to help them solve their data problems. But traditional uses of these tools are often limiting, as they process data on a single machine where the movement of data becomes time consuming, the analysis requires sampling (which often does not accurately represent the data), and moving from development to production environments requires extensive re-engineering. Spark provides a powerful, unified engine that is both fast (100x faster than Hadoop for large-scale data processing) and easy to

use.Hence we are going to Pyspark for our initial data analysis and then we well do further analysis using python and RStudio.

Apart from above tools we use some other tools for the implementation.The following are the list of important softwares and tools used.

- dnspython:it is a DNS toolkit for python.This is used in the code to get the A records of hosts.

- pygeoip :The libarary is used to get the ASN numbers associated with IP addresses.This library is based on Maxmind's GeoIP C API.

- urlparse :this is used to convert a relative url to an absolute url.

- Public suffix List :This is the collection of all registered host names given by all internet users.The Public Suffix List is an initiative of Mozilla, but is maintained as a community resource.It is available for use in any software, but was originally created to meet the needs of browser manufacturers.In our code we use it to get second level domain from a host name.The "effective_tld_namesdat" is the file which is free downloadable from their site.

- IPython :IPython notebook is used for pyspark code writing and execution.

- matplotlib :Python library used to for making different graphs used in this thesis.

- pygal :Pygal is also another python library which we used to make graphs.

## 4.3 Design of Crawler Engine

In this section,we will discuss the whole architecture behind the implementation.The design process involves two parts,first one is the architectural overview of crawler engine which takes 100,000 top ranked websites of Alexa as input,crawl the websites , return result file as output from the engine and second part is processing of result file using data processing engine which internally use "pyspark".Along with this,we will also discuss little bit about "Scrapy framework" which is main backbone behind the crawler engine.The input to crawler and result output file matrices also will be discussed on the below section.

### 4.3.1 Crawler Engine

In this section we will discuss the internal architecture behind the crawler engine.We will start with "Scrapy framework" [19] which is the main crawler engine work behind the whole crawler engine.Scrapy framework is based on spiders which are self-contained crawlers worked by set of given instructions.As it is one of the top ranked open source project,lots of documentation can be found.The main objective of crawler engine is to take input data in the format of text file which contain the top 100,000 websites of Alexa.Then divide this master domain list into multiple sub domain lists with equal weighted websites and then process each sub domain list using separate crawler instances.Each crawler instance work independently and store all website information like HTTP header information,rank of website in Alexa etc. in result file.

### 4.3.2 Scrapy Framework

Scrapy framework is one of top ranked open source project used for,a fast web crawling framework,used to crawl websites and extract structured data from their pages.It

provide option for both focused and broad crawling.In case of focused crawler scrapy crawls a specific domain while in case of broad crawling a large (potentially unlimited) number of domains.Hence for our implementation ,we are going to use the broad crawling.As we can see from figure-7,the main components of the framework contain scrapy engine,scheduler,downloaders,spiders,item pipeline,downloader middlewares.

- Scrapy Engine : The engine is responsible for controlling the data flow between all components of the system, and triggering events when certain actions occur.

- Scheduler : The Scheduler receives requests from the engine and enqueues them for feeding them later (also to the engine) when the engine requests them.

- Downloader : The Downloader is responsible for fetching web pages and feeding them to the engine which, in turn, feeds them to the spiders.

- Spiders : Spiders use for to parse responses and extract items from them or additional URLs (requests) to follow.

- Item Pipeline : The Item Pipeline is responsible for processing the items once they have been extracted (or scraped) by the spiders. Typical tasks include cleansing, validation and persistence (like storing the item in a database). For more information see Item Pipeline.

- Downloader middlewares : Downloader middlewares are specific hooks that sit between the Engine and the Downloader and process requests when they pass from the Engine to the Downloader, and responses that pass from Downloader to the Engine.
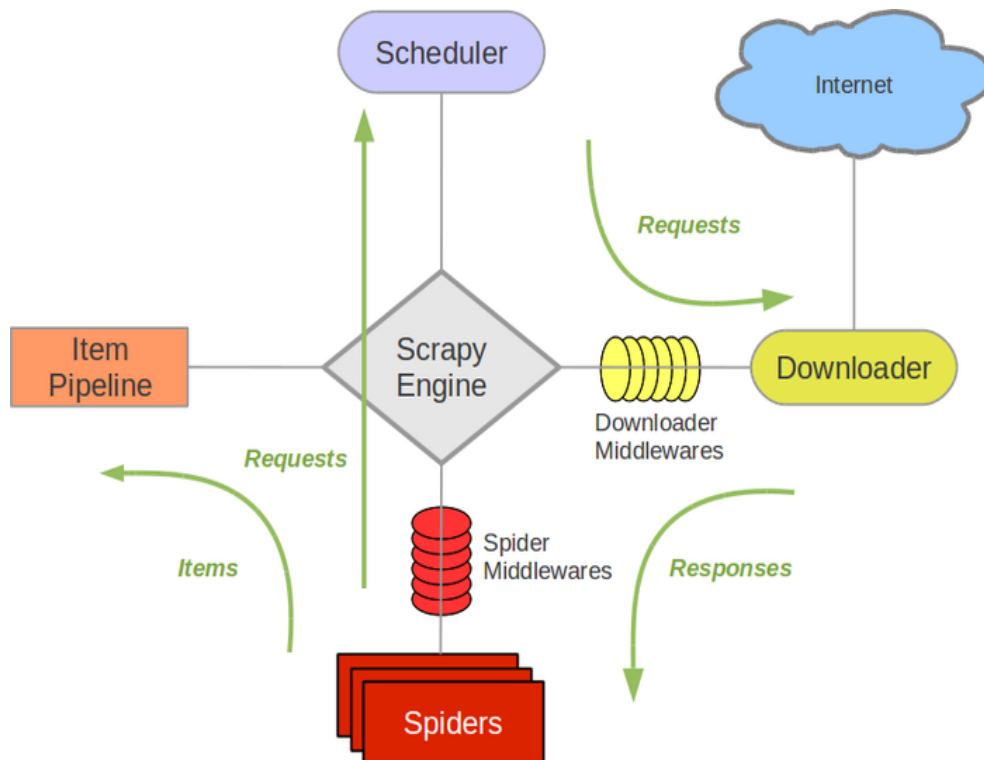


Figure 7: Scrapy Architecture

### 4.3.3 Scrapy data types

Scrapy uses some specific classes and strings which are used to crawl data.These are written in python and are open source.Scrapy also use twisted framework for multi threading.Some of the classes which we are going to use for implementation are as follows,

1. Spider :Spiders are the classes used to crawl single or multiple domains using different methods like start_requests(),parse().

   - start_requests() :This scrapy method is used to pass domains to parse method for further crawling process.
   - parse() :he parse method is in charge of processing the response and returning scraped data and/or more URLs to follow. Other Requests callbacks have the same requirements as the Spider class.

2. Twisted Framework:Twisted supports an abstraction over raw threads — using a thread as a deferred source. Thus, a deferred is returned immediately, which will receive a value when the thread finishes. Callbacks can be attached which will run in the main thread, thus alleviating the need for complex locking solutions. A prime example of such usage, which comes from Twisted's support libraries, is using this model to call into databases. The database call itself happens on a foreign thread, but the analysis of the result happens in the main thread.

3. Requests objects:Typically, Request objects are generated in the spiders and pass across the system until they reach the Downloader, which executes the request.A Request object represents an HTTP request, which is usually generated in the Spider and executed by the Downloader, and thus generating a Response.

4. Response objects :Request object returns a Response object which travels back to the spider that issued the request.A Response object represents an HTTP response, which is usually downloaded (by the Downloader) and fed to the Spiders for processing.

## 4.4 Workflow

The whole workflow is divided into two steps.In first step divided the whole list of 100,000 domains ranked by Alexa into multiple sub lists.Each sublist then will be sent as input to multiple instances of scrapy.In the second step scrapy engine crawls each website link from sub lists and store the results.We will discuss both the steps in following subsections.

### 4.4.1 PreCrawling Step

In precrawling step,main focus is to create multiple instances of scrapy spider.As scrapy is a large memory greedy tool,after testing we decided to create 20 parallel threads which work independently shown in figure 8.Each instance will take separate input file which contain 5000 domains .So we need to divide the master domain list into sublists but we also need to divide it such a way that all the crawlers will complete their crawling in almost same time.Hence we choose to create the sublists with equal weight.Weight points to rank of the websites here.

The master domain list contain the the website url and the rank of the websites.So after division the first sublist will contain the websites of 1st rank,21st rank,...99981st rank;second list will contain 2nd rank,22nd rank...99982nd rank websites and so forth for other instances.Hence each sublist will contain 5000 websites with equal division of website ranks.
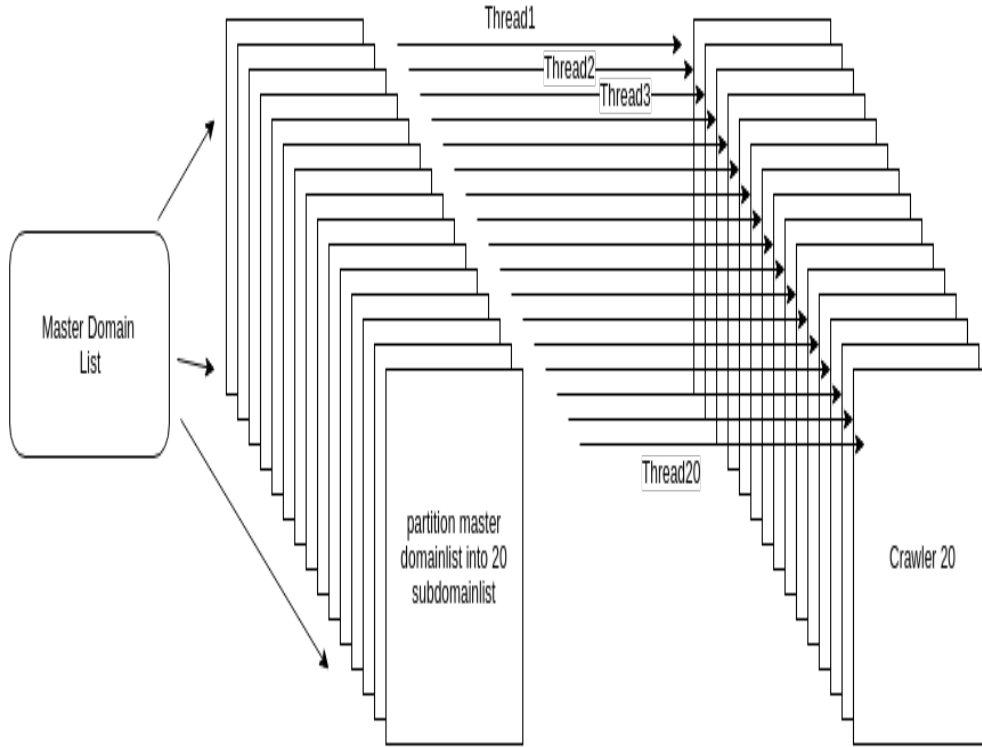
Figure 8: Pre Creawling Step

### 4.4.2 Crawling Step

In this step each crawler will take 5000 domains and parse one by one domain to the crawler engine.For each url we download the corresponding web page, extract the linked URLs, and check each url to see whether the extracted url is a fresh url which has not already been seen . With this architecture we are essentially carrying out a very large number of independent crawls of the white listed domains obtained from Alexa.Along with HTTP header information,the crawler engine also extract Arecord details ,ASN details.

### 4.5 Conclusion

In this section we talked about the softwares we are going to use to complete the whole implementation.We talked about the motivation behind choosing different languages,software tools to complete the implementation.Again we discussed about the web crawler as well internal architecture of crawler engine.In the next section we will talk about the approach to collect traces, i. e., active DNS measurements, in order to evaluate our methodology.

19

# 5 Measurement

In this section we present our approach to collect various website information,in order to evaluate our methodology.In order to achieve our goal ,we crawl a list of websites and analysis those data. This chapter explains the overview of those traces which we received after web crawling,data cleanup procedure and the final set of data after the clustering algorithm described in chapter-3 methodology section.

## 5.1 Host name Selection

To obtain a good coverage of the largest hosting infrastructures, we crawl top 100,000 ranked websites from Alexa [26].Alexa ranks websites based on Internet traffic-users of its tool bar for various web browsers like Google Chrome,Internet explorer,Firefox.After crawling we are able to generate 13919464 traces.

Moreover,websites contain a lot of embedded contents like images,videos, advertisements etc. that the browser of the user has to download from the various web servers.These web servers can be from different web content providers or from hyper giants. In our study, such embedded content has to be taken into account, as it might be served from servers other than those serving the front page of a popular host name listed in top rank websites of Alexa.To give better understanding,while crawling facebook.com,the front page is served from Facebook data centers while the logo and other meta data come from Akamai.For most of the websites,the important videos,images or other embedded contents present in the front page of the websites.So to make our study more precise,we crawl only the front pages of all the top ranked websites and the embedded links present in the front pages of those websites.

## 5.2 Data Cleanup

We perform a thorough cleanup process on the raw traces.In each crawling we gather 15 different metrics.Hence we pass each crawled link into regular expression to check the validity of the links based on number of metrics,type of matrices.

First we collect top 100,000 top ranked websites as of September 2016 from Alexa website.The scrapy crawler queries the HTTP get method to local DNS resolver for all the websites and store the results in a trace file.Each trace contain total 15 different data which are described below.

1. index :index shows the rank of the website

2. depth_level : 0 or 1. 0 shows that the website is the main page url and 1 shows the embedded links inside main page

3. httpResponseStatus :the HTTP return status code.

4. MIMEcontentType :this is included to know the type of element inside the web page.

5. content_length :content length gives the idea about the size of the element.

6. url :this is the url to be crawled by the scrapy engine.This url can be main page url or embeded links.Duplicate links are omitted for the same main page by using RFPDupeFilter.RFPDupeFilter is a scrapy class which is used to detect and filter duplicate requests [3].

7. cookies :cookies involved in a website

8. tagType :this shows the HTML element type. for example if an element is embeded in a wensite like "¡img class="desktop" title="" alt="" src="img/bg-cropped.jpg"¿",then we collect the tagtype which is img here. ARecord=This contain all the aRecord names involved in a website while resolving to the ip address.. for the the website

```
; <<>> DiG 9.10.3-P4-Ubuntu <<>> www.bmw.com
;; global options: +cmd
;; Got answer:
;; ->>HEADER<<- opcode: QUERY, status: NOERROR, id: 11656
;; flags: qr rd ra; QUERY: 1, ANSWER: 4, AUTHORITY: 0, ADDITIONAL: 1

;; OPT PSEUDOSECTION:
; EDNS: version: 0, flags:; udp: 512
;; QUESTION SECTION:
;www.bmw.com.                    IN      A

;; ANSWER SECTION:
www.bmw.com.              3600   IN      CNAME   cn-www.bmw.com.edgesuite.net.
cn-www.bmw.com.edgesuite.net. 3600 IN    CNAME   a1586.b.akamai.net.
a1586.b.akamai.net.       20     IN      A       104.121.76.73
a1586.b.akamai.net.       20     IN      A       104.121.76.64

;; Query time: 22 msec
;; SERVER: 127.0.1.1#53(127.0.1.1)
;; WHEN: Wed Oct 05 21:32:17 CEST 2016
;; MSG SIZE  rcvd: 143
```

Figure 9: dig for bmw.com

"www.bmw.com",the ARecord name are a1586.b.akamai.net. and a1586.b.akamai.net. which will be stored in trace will under ARecord column.

9. destIP :this column stores the corresponding resolved ip addresses of the website.For example,for the website bmw.com (figure-6),the destIP will store 72.247.184.130 and 72.247.184.137

10. ASN_Number=Field():This column stores the ASN number of a IP address.For this we use maxmind IP to ASN mapping file.

11. distinctASNs=This will keep the total number of ASNs involved for all the embedded website links for a particular main page url which is one of the Alexa's top ranked website.

12. ObjectCount=Field():This column stores the total number of external links embedded in a website.

13. NumberOfuniqueExternalSecondlevelSites :This columns gives the total number of unique second level domains involved in all the embedded links of a website.

14. start_time :This contains the staring time when the HTTP starts requesting for the website.

15. end_time=This contains the staring time when the HTTP ends requesting for the website.

Through scrapy we crawled total total 13919464 website links.After data cleanup, we have 13919464 clean traces that form the basis of this study.

## 5.3 Clustering Algorithm

Overall we get total 219604 unique second level domains.But from them a lot of SLDs which can be clustered into other SLDs.After clustering algorithm we are able to get 53852 unique clustered SLD infrastructures.

| | Parent.Sld | ClubbedSlds | links | ips | ASNs |
|---|---|---|---|---|---|
| 1 | cloudflare.net | 17711 | 1295505 | 29893 | 17 |
| 2 | yunjiasu-cdn.net | 6068 | 210463 | 5907 | 21 |
| 3 | us-east-1.elb.amazonaws.com | 4254 | 199109 | 10885 | 31 |
| 4 | wpengine.com | 3963 | 136400 | 4072 | 19 |
| 5 | d5nxst8fruw4z.cloudfront.net | 3331 | 75660 | 1619 | 3 |
| 6 | anycast.me | 2667 | 116024 | 2207 | 3 |
| 7 | d2t8dj4tr3q9od.cloudfront.net | 2363 | 59807 | 1316 | 8 |
| 8 | eu-west-1.elb.amazonaws.com | 1946 | 100265 | 4476 | 25 |
| 9 | jiashule.com | 1769 | 79038 | 2249 | 28 |
| 10 | kxcdn.com | 1547 | 82120 | 1235 | 7 |
| 11 | ap-northeast-1.elb.amazonaws.com | 1308 | 60343 | 3105 | 22 |
| 12 | cdntip.com | 1046 | 59457 | 1080 | 10 |
| 13 | us-west-2.elb.amazonaws.com | 1031 | 45877 | 2570 | 13 |
| 14 | alikunlun.com | 964 | 76647 | 1155 | 17 |
| 15 | cdn20.com | 792 | 34657 | 992 | 16 |
| 16 | akadns.net | 791 | 36908 | 1018 | 91 |
| 17 | windows.net | 783 | 28601 | 740 | 1 |
| 18 | scutum.jp | 768 | 21871 | 768 | 4 |
| 19 | ap-southeast-1.elb.amazonaws.com | 730 | 44386 | 1574 | 5 |
| 20 | ourwebpic.com | 722 | 113746 | 772 | 16 |

Figure 10: Top 20 clustered SLD infrastructure in decreasing order of clubbed SLD count.

The top 20 clustered SLD infrastructures in the decreasing order of clubbed SLDs are shown in figure 10.The third column contain number of links served by a whole clustered SLD infrastructure.Similarly fourth column shows number of IP addresses resolved and fifth column shows the number of ASNs managed as a whole clustered SLD infrastructure.From figure we can see that under cloudflare.net,17711 number of SLDs clubbed which is 10,68 % of all child slds clubbed.From the table also we can see three different Clustered SLD infrastructure where the main SLD having cloud flare in SLD naming pattern. cloudflare.net,d2t8dj4tr3q9od.cloudfront.net,d5nxst8fruw4z.cloudfront.net All the three SLD infrastructures are from same parent company cloudflare.net.But they shows different clustered infrastructure which shows that they have different infrastructure from each other in terms of bgp prefixes.Similarly we can five different clustered SLD infrastructure having amazon in their SLD naming pattern.Both amazon and cloudflare are big CDNs ,they also provide a lot of other services like Internet Security services,distributed domain name server services, web hosting service etc.Similarly amazon provide cloud computing services,infrastructure services etc. to their customers.This unveils that big hosting infrastructures maintain different SLD infrastructures separately which might they use for different purposes.

## 5.4   Web objects

We also collected different web objects embedded in home pages of 100,000 top ranked web sites of Alexa.After crawling we found total 285 different types of objects from 219604 unique second level domains.The main web objects which shown are text/html which is around 38% of all object types crawled.Similarly SLDs serve almost 42 % of image files.Here it is important to notice that we only get the type of object by collecting content type header field.

## 5.5   Conclusion

In this section we define our approach of selecting SLDs after clustering algorithm.This will help to find out prominent infrastructures present in today' Internet.We also talk about different metrics considered while crawling web site.

# 6 Results

In this section,the prominent hosting infrastructures will be identified first.Next these prominent infrastructures will be clustered together using the cluster algorithm described in chapter-3.Once clustered SLD infrastructures are identified,those will be analyzed further to gain insight on the their deployment and hosting strategies.Finally based on their strategies hyper giants will be determined.Once hyper giants will be determined ,the dependency between them and popular websites will be examined by taking into consideration that what type of web objects like images,videos,HTML files etc delivered through these hyper giants to popular websites.

## 6.1 Analyzing prominent infrastructures

In this section we will analyses the SLD infrastructures present in Internet based on number of URLs served by them.A small subset of the SLDs are analyzed to understand if there are common characteristics present within them.
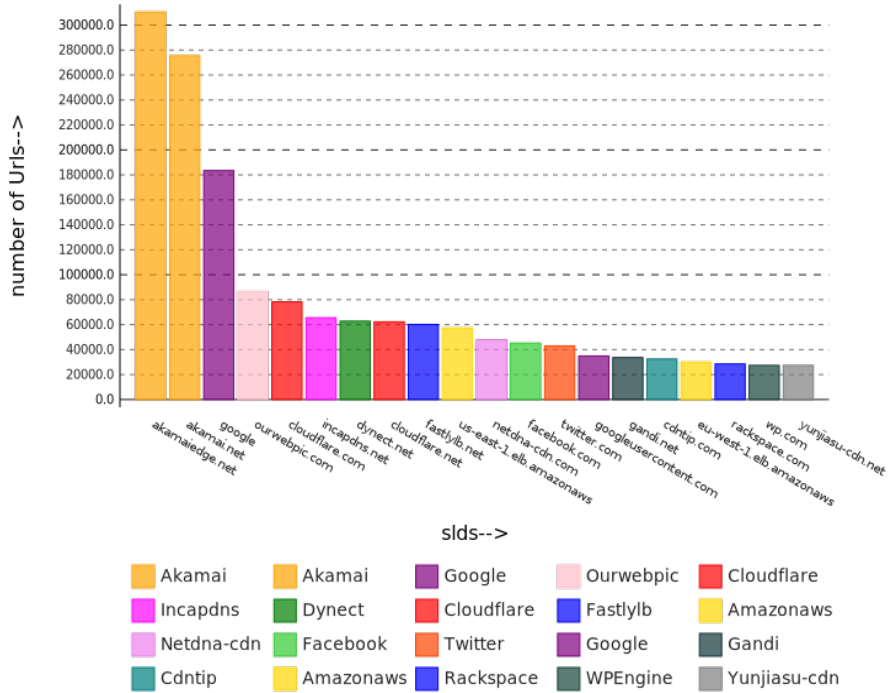


Figure 11: Top 20 SLDs

From figure-11 ,it is found that there are some SLDs which are served by same company.Like akamaiedge.net and akamai.net both are CNAMEs used by Akamai company.Similarly google.com,googleusercontent.com,googlehosted.com and googledomains.com all are used by Google.Normally companies used different names when they serve different services to the customer or different name for certain located customers.This can be seen from the CNAMEs used by amazonaws. us-east-1.elb.amazonaws.com,eu-west-1.elb.amazonaws.com are two examples of naming pattern used by the amazon to provide services to distinct located customers.But here the question arises that these names are pointing to same infrastructure or different.If they are pointing to same infrastructure they should be considered as one ,else separately.Hence it is important to identify if these SLDs

share same the infrastructure or not.This can be identified by analyzing the bgp prefixes they share.To measure this we are going to take RIPE bgp prefixes of each SLD routed to and cluster all the SLDs that use the same bgp prefixes.

The top 20 SLDs serve almost 13.13% of all the URLs crawled.These 20 SLDs contain not only CDNs like Akamai, Cloudflare etc., but also contain content providers like Google,Facebook.It also contain SLDs like amazonaws which provides cloud computing services while some other SLDs are web hosting companies like ccgslb.net.A small fraction of data set gives us multiple types of infrastructures.Classification of different infrastructures cannot be done using only number of links.To know if a infrastructure is highly distributed all across the world,the number of ASN number need to be checked.

### 6.1.1 Conclusion

From this section the following observations can be inferred.

1. There are some SLDs which are served by same parent company like akamaiedge.net and akamai.net which are served by company Akamai.Hence it is important to identify whether they can be clustered into same infrastructure or not.This cannot be analyzed with just the number of URLs instead we need to check the footprint covered by these SLDs in the world by checking their bgp prefix routes.

2. Since it is identified that there is a possibility of some SLDs getting clustered , it is important not to restrict the test sample for the top 20 of SLDs,but should be extended to all the 219604 SLDs.

In the next section,we will discuss the steps to identify the hyper giants using the clustering algorithm, as discussed in chapter 3 methodology section.

## 6.2 Identifying hyper giants

In 2010, Craig Labovitz, then of Arbor Networks,characterized the hyper giant as a content provider that makes massive investments in bandwidth, storage, and computing capacity to maximize efficiencies and performance.But as the architecture of Internet evolves,researchers found that the Internet has now become a flatter infrastructure where there are fewer autonomous systems connected to each other and they try to have a bigger footprint in terms of number of bgp prefixes than before.In this way they are able to diversify their architecture as well as able to move content to even closer to their customers.They termed this infrastructure providers as hyper giants [14].

Overall we get total 219604 unique second level domains.But from them a lot of SLDs which can be clustered into other SLDs.

After clustering algorithm we are able to get 53852 unique clustered SLD infrastructures.From there ,almost 80% of the total unclustered SLDs got clustered into first 3.73% of SLDs .It shows that these 3.73 % of top SLDs have footprint all over the world through highly distributed CDNs,data centers etc.Other SLDs share their infrastructures with these top 3.73% of clustered infrastructures.Hence there is a possibility to get the hyper giants in these 3.73% of clustered SLD infrastructures.But there are SLD infrastructures who have their own infrastructures in the form of data centers.Hence they don't share any other SLD infrastructures.Like facebook.com who has its own infrastructures in the form of data centers all over the world.In fact we also found almost 82.45% of clustered infrastructures who does not share their infrastructure with no more than another SLD infrastructure.It means there are companies who work independently by creating their own infrastructures.This can be data centers all over world.Although these 82.45% of clustered
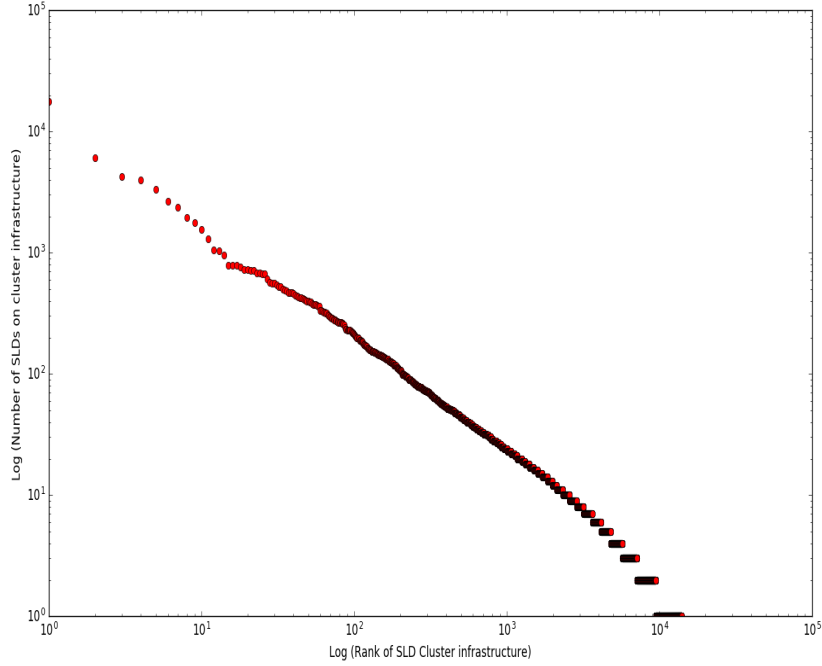
Figure 12: Number of SLDs served by different SLD infrastructure clusters.

infrastructure do not share their infrastructure with no more than single infrastructure, still some of them serve a large number of links which make them another candidate for hyper giant analysis.

Hence to get a better picture we will see clustered SLD infrastructures based on how many links they served.So we sorted all clustered SLD infrastructures in their decreasing order of links they serve and we found that top 5.95 % (=3205) clustered SLD infrastructures serve almost 80% of links and have 78.65 % of SLDs.Hence there is a possibility of getting hyper giants in this range.

### 6.2.1  Clustered SLDs

From last section we identified 5.95% (=3205) clustered SLD infrastructures as candidates for hyper giant analysis.To identify the hyper giants,two different steps will be followed.In the first step,the 5.95% clustered SLD infrastructures will be analyzed based on number of links they serve, to number of IP addresses they resolve.The big SLD cluster infrastructures will be separated from small SLD infrastructures by end of this step.In the next step we will again compare number of prefixes they resolved as a clustered SLD infrastructure to number of ASN numbers they belong to.This will give us a better idea how their whole infrastructures are distributed all over the world.After these two process we will try to identify the hyper giants.

### 6.2.2  case 1 :number of Links Vs number of IP addresses

In this section we will take the top 3205 SLD infrastructures and cluster them based on their number of links to ip addresses they served.

We used k-means clustering algorithm and number of cluster parameter 10.We found
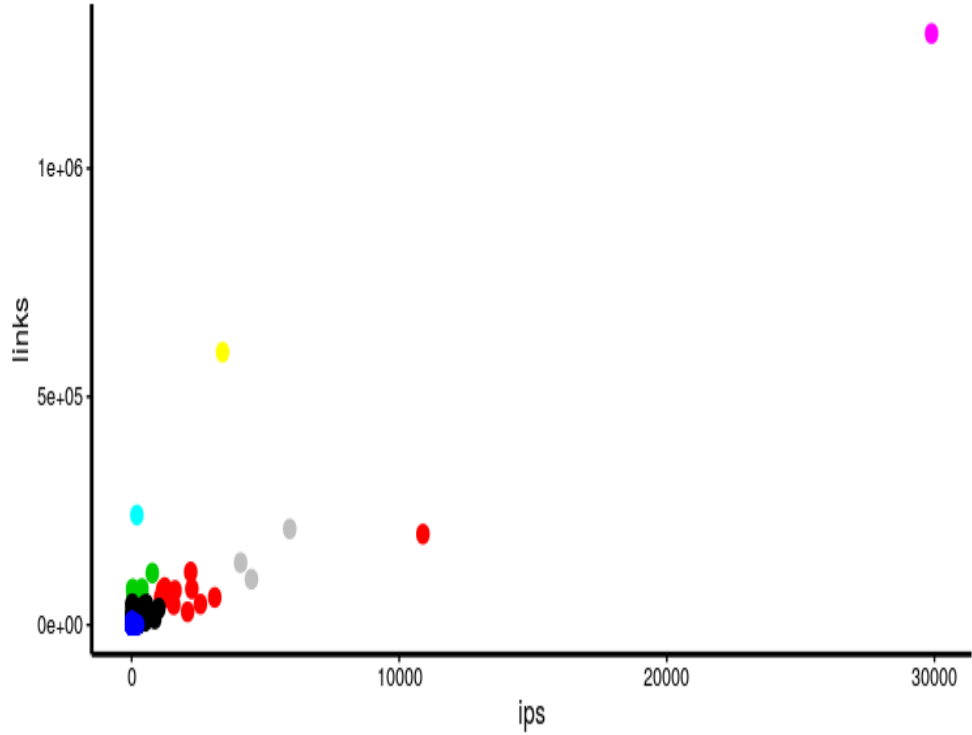
Figure 13: Clustering based on links and IP address features

6 different clusters which are clubbed total 26 SLD infrastructures and showing unique behavior.Like cloudflare.net is clustered separately as it is serving very huge number of links as well as having very high number of ip addresses.It means lots of small SLDs are serving through cloudflare.net and it has footprint all over the world.Similarly us-east-1.elb.amazonaws.com clustered separately as it has less high of links but serving a high number of ip addresses.third cluster contain google.com which is serving high number of links but not very high number of ip addresses.In this way we are able to identify total 6 clusters which resulted in a total of 26 SLD infrastructures. But it is difficult to categorize them into some specific type of infrastructure based on only links to ip address analysis. Hence these 26 SLD infrastructures will be further analyses taking into account their prefixes to their asn numbers.

### 6.2.3   case 2 :number of BGP Prefixes Vs number of ASNs

From last section we identify the SLD infrastructures which are having unique behavior. But we couldn't able to classify them .Therefore in this step we will check how they are distributed all other world.This requires these clusters to be analyzed using their corresponding ASN to prefix numbers.This is because number of ip prefixes shows the footprint of the infrastructures across the world and the number of asn numbers show the degree of distribution of infrastructures across the world.

From figure-14 ,we can cluster all the clustered SLD infrastructures into 5 parts based on their ¡number of prefixes,number of ASNs¿ analysis as below.

- very high,very high : In total 3 different clustered SLD infrastructures are clustered under this. They are yunjiasu-cdn.net,jiashule.com, us-east-1.elb.amazonaws.com.These three SLD infrastructures contain very high number of prefixes as well as they have
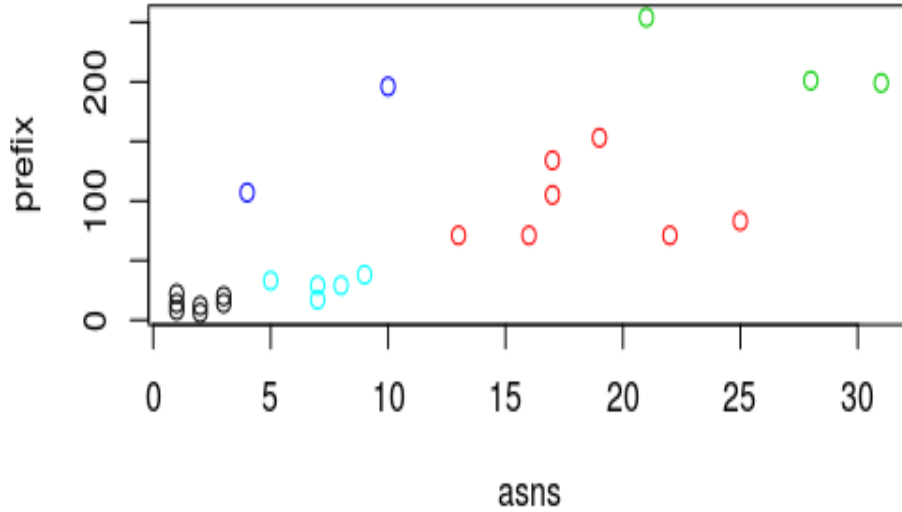
27

Figure 14: Classification of hyper giants

high number of ASN numbers,which shows that they have footprint all over the world as well as they are distributed across the world.We can classify them as highly distributed CDNs.

- high,high : Total 7 SLD infrastructures clubbed inside this. wpengine.com ,alikunlun.com,cloudflare.net,ourwebpic.com, us-west-2.elb.amazonaws.com,eu-west-1.elb.amazonaws.com and ap-northeast-1.elb.amazonaws.com.They have high number of footprint and high number of asn numbers.This shows they have presence in few of the regions.We can classify them as distributed CDNs.

- high,less : netdns-cdn.com and cdntip.com These SLD infrastructures have high number of prefixes but have less number of ASN numbers.It means they have footprint all over the region but they normally administered through very few ASN numbers.Hence these can be categorized into cloud computing infrastructures.

- medium,medium : Five different SLD infrastructures are clubbed into same cluster.They are, akamaiedge.net,ap-southeast-1.elb.amazonaws.com, kxcdn.com,incapdns.net,d2t8dj4tr3q9o All these infrastructures have few prefixes and they also not distributed which gives an evidence of multi homes data centers or web hosting companies.

- less,less : The rest of the CDN infrastructures are clubbed into same cluster which having very few number of prefixes also very few number of IP addresses.Hence they can be treated as content providers.Google and Microsoft both clustered under this.

### 6.2.4 Conclusion

From the above section,we are able to identify total 26 hyper giants which are having influence in Europe region.They are highly distributed CDNs,cloud platforms,CDNs,content producers etc.In next section we will see how they influence on number of objects delivered by them.

### 6.3 Popular websites dependency on hyper giants

In this section,first we will see different types of objects delivered through 219604 unique SLDs and compare this with web objects delivered by 26 hyper giants. Then we will examine each hyper giant separately and try to find what kind of data object they deliver.

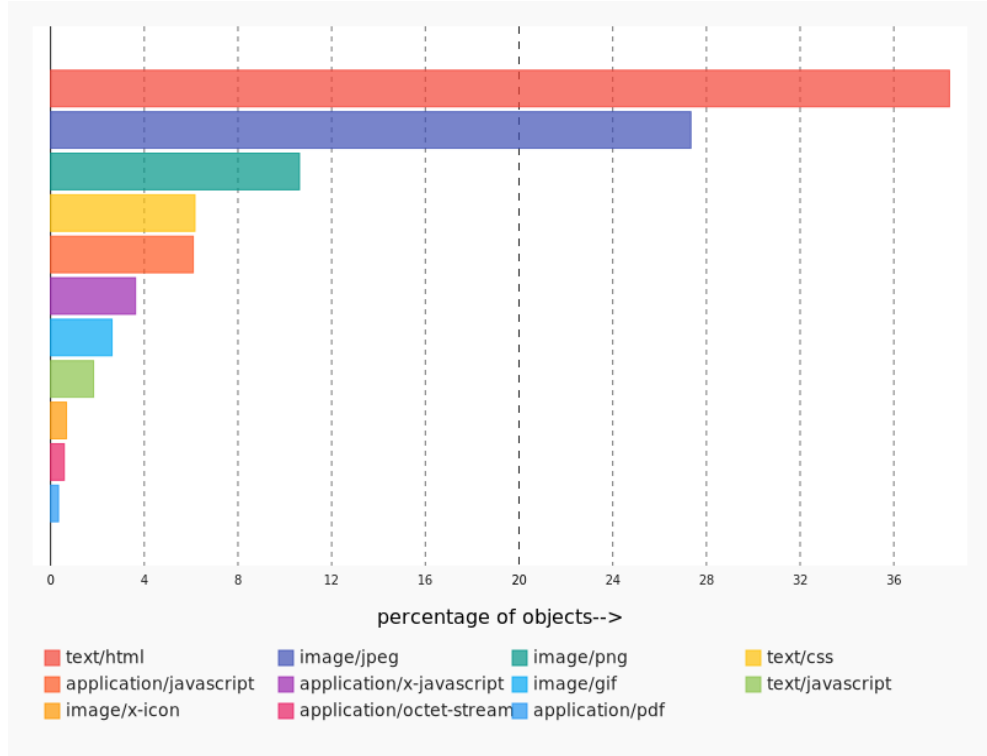### 6.3.1 Object types delivery through whole SLD infrastructures Vs hyper giants



Figure 15: top 10 object percentage used by all SLDs

Figure-15 and figure-16 shows top 10 object types in form of percentage delivered by both hype giants and slds.In case of hyper giants top 10 object types deliver almost 98.82 % of all object types.Again we can see that most of these object types are HTML files.HTML carries almost 68.27% of object type in compare to other object types.Similarly 75.04% of text object types which contain text/html,text/css,text/xml,text/js etc., are delivered through hyper giants where as only 19.29% of image files delivered through hyper giants.

In case of SLDs,top 10 object types carries almost 97.94% of all object types which is very much similar to the percentage of objects delivered by hyper giants.The highest web object type delivered through all SLDs as well as hyper giants is text/html but it can be observed that in case of all SLDs,text/html carries almost 38% of all web objects which is around 68.27% in case of delivering through hyper giants.In case of all SLDs ,object types
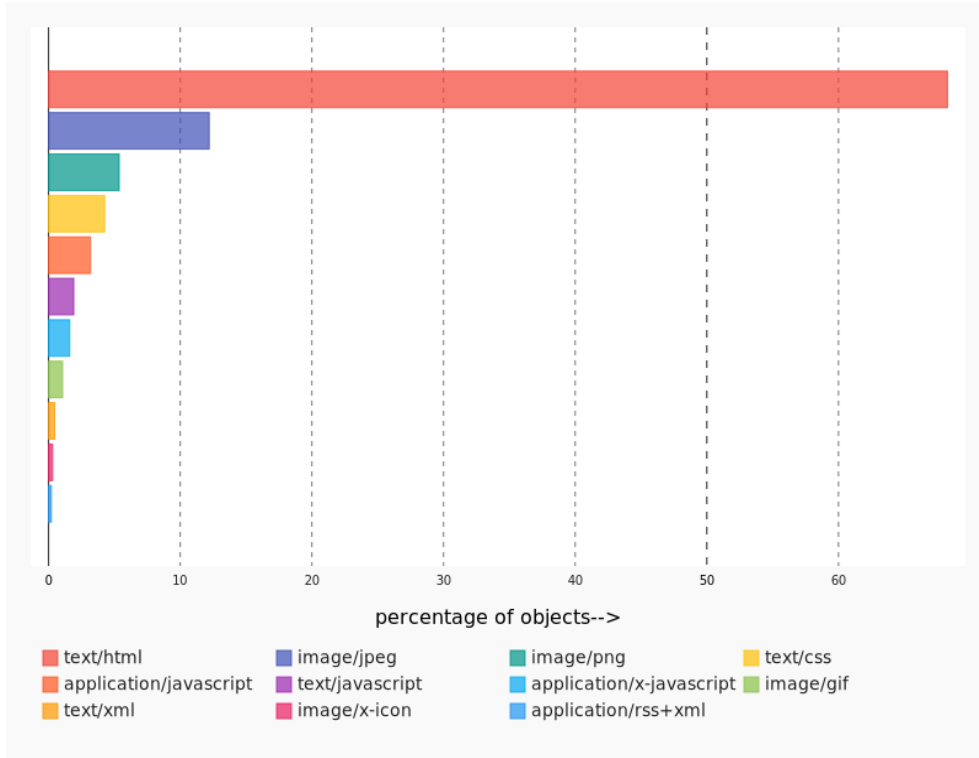
Figure 16: top 10 object percentage used by all hyper giants

are distributed properly.If we add image/jpeg,image/png,image/gif then all SLDs serve more image files than HTML files.But same is not the case for hyper giants.Similarly it can be seen that almost 46.77% of text files delivered through whole sld set which is almost 1.5 times more in case of hyper giants.But in case of image files ,SLDs serve almost 42 % of all web objects which is almost double the image files served through hyper giants.This different behavior might be because popular content websites normally store more dynamic files in CDNs where as in case of image files they store at their own servers.

### 6.3.2 Object types delivered from hyper giants to popular web sites

In this section we will see what kind of data mostly delivered through the identified hyper giants.In today's Internet ,content plays most vital role.Hence it is important to observe what kind of data mostly delivered through hyper giants .

| Hyper giant Object List | | | |
|---|---|---|---|
| Country Name | text | image | application |
| alikunlun.com | 75.99 | 20.32 | 3.65 |
| d2t8dj4tr3q9od.cloudfront.net | 49.86 | 41.10 | 8.68 |
| ap-northeast-1.elb.amazonaws.com | 90.28 | 6.23 | 3.44 |
| ap-southeast-1.elb.amazonaws.com | 87.94 | 7.36 | 4.62 |
| cdntip.com | 87.94 | 7.36 | 1.67 |
| d5nxst8fruw4z.cloudfront.net | 31.90 | 57.42 | 10.23 |
| eu-west-1.elb.amazonaws.com | 87.25 | 7.32 | 5.39 |
| fastlylb | 69.79 | 21.30 | 8.69 |
| google.com | 82.54 | 15.79 | 1.65 |
| jiashule.com | 88.85 | 8.33 | 2.80 |
| kxcdn.com | 75.62 | 17.98 | 6.34 |
| pbwstatic.com | 83.18 | 16.57 | 0.24 |
| akamaiedge.net | 73.65 | 20.35 | 5.93 |
| anycast.me | 79.64 | 14.90 | 5.37 |
| cloudflare.net | 64.86 | 11.98 | 23.12 |
| cloudflare.com | 78.15 | 15.62 | 1.70 |
| dynect.net | 94.53 | 3.74 | 3.74 |
| edgecastcdn.net | 53.23 | 38.44 | 8.17 |
| incapdns.net | 64.86 | 7.61 | 4.81 |
| netdna-cdn.com | 21.08 | 55.04 | 23.39 |
| ourwebpic.com | 86.56 | 11.58 | 1.84 |
| us-east-1.elb.amazonaws.com | 88.99 | 6.18 | 4.73 |
| wpengine.com | 80.67 | 12.42 | 6.84 |
| us-west-2.elb.amazonaws.com | 85.30 | 8.15 | 6.47 |
| windows.net | 80.52 | 12.99 | 6.41 |
| yunjiasu-cdn.net | 87.76 | 9.86 | 0.0 |

The table shows all 26 hyper giants and the percentage of text,image and application web objects delivered by them.We can see from table than most of the hyper giants deliver very high percentage of text files which contain text/html,text/css etc.But their are few exceptions .Like both the cloudfront SLDs are providing very high number of images compare to other hyper giants. Similar kind of observation can be seen for netdna-cdn which provides more image files than text files.

From last section we observed yunjiasu-cdn.net,jiashule.com, us-east-1.elb.amazonaws.com are highly distributed CDNs.It can be observed that all the three CDNs are delivering very high number of text files compare to image and application files.This might be because of their footprint all over the world and also have massively distributed CDNs.Hence they cache more of the HTML files at edge servers to provide better performance.Similarly observation can be seen from distributed CDNs .These CDNs also deliver high number of HTML files compare to image files but as they have presence in some regions the number of HTML files are not that comparable to highly distributed CDNs.Third infrastructure type we observed was cloud computing infrastructures and netdns-cdn.com,cdntip.com clustered under that.From the table it can seen that netdns-cdn.com delivers more images and very less number of HTML files.Again cdntip delivers highest number of application data.The data centers provide both images and HTML files in a very balance way.cloud front delivers 49% of HTML file and 41% of image files.Similarly akamaiedge.net provide around 20% of images.Content providers like window.net and google.com delivers very high number of links compare to number of images.This is evident as they have more content .

### 6.3.3 Conclusion

From the section we can infer that both SLDs and hyper giants deliver maximum number of text/HTML files but it also found that hyper giants delivered almost 1.5 times more HTML files than SLDs.Same kind of observation can be seen for image files where SLDs deliver double the image content than hyper giants.Again we found that highly distributed CDN generally deliver more HTML links compare to image or application files which might be because of their massive CDN distribution.Distributed CDNs provide high number of HTML files because of their presence in few regions.Cloud computing cdns provide more images than HTML file which might be because cloud computing provide scalability.Data centers delivers both HTML file and images in almost same ratio.Content providers also delivers more HTML files compare to image files which because of their rich content.

# 7 Conclusion

In this thesis,we introduce a automated process to find out the prominent infrastructures in today's Internet as well as to classify these prominent infrastructures to find out the presence of hyper-giants using DNS measurement and bgp prefixes.We presented a clustering algorithm which will help to find out which SLDs are sharing their infrastructures.The advantage of this automated approach is that it uses each to retrieve SLD and bgp prefixes ,hence this procedure can be used in future.

Along with this we measure what object types are delivered by major hyper giants.This will help researchers to get a better view to classify hyper giants based on their object type delivery.Not all popular websites provide same kind of content.Some websites are popular for delivering videos and some other are for user content.Hence with this change of content type,we provide a overview of hyper giants according to different object type they serve.

The data is collected at a single vantage point at Germany.This thesis was able to identify high distributed CDNs,cloud service providers,content providers etc. and their role which will mostly hold good for across Europe.

Furthermore our thesis is an important step towards answering some of the very crucial questions for highly distributed CDNs,distributed CDNs,content providers etc.It will give them idea to find out how other CDNs distribute their infrastructure as well as their network footprint distribute across different region which will give them a competitive advantage over their competitors in content delivery market place.Moreover it will help the research community to discover the Internet architecture changes with time.They also can able to track the hyper giants and their dependency with other popular websites.

# 8 Future Work

There are certain areas which can be investigated further in future which are not covered in the current scope of this thesis.This section will discuss about these points.

- First of all the thesis is done at single vantage point at Germany.Hence the identification of hyper giants,their role and their relationship with popular websites can be changed when the whole procedure will be done in whole world basis.Hence it will be interesting to see how the clustering algorithm works when taking all the SLDs across the world.

- Secondly while clustering the hyper giants,the k-means parameter is taken by going through very small number of observations.Hence in future this can be tested more precisely which will help to cluster the hosting infrastructures at a granular level.

- In the clustering algorithm,we club two SLDs if one SLD prefixes matches with other SLD prefixes by more than equal to 70%.This matching index is chosen after extensively testing.Hence this can be taken as future work to see what is the best matching index.

- After clustering algorithm,we have chosen to take top 5.95 % (=3205) clustered SLD infrastructures which serve almost 80% of links and have 78.65 % of SLDs .In this way we argumented to get most of the big hosting infrastructures as well as content providers for further analysis.Hence this can be taken for future work to validate the argument properly.

# 9 Appendix

# 10   List of Acronyms

| | |
|---|---|
| CDN | Content delivery network |
| SLD | Second level domain |
| ASN | Autonomous system number |
| BGP | Border gateway protocol |
| HTTP | Hypertext transfer protocol |
| HTML | Hypertext markup language |
| DNS | Domain name system |
| IP | Internet protocol |
| QoS | Quality of service |
| XML | Extensible markup language |
| ISP | Internet service provider |
| CDI | Content delivery infrastructure |
| IaaS | Infrastructure as a service |
| PaaS | Platform as a service |
| SaaS | Software as a service |
| IXP | Internet exchange point |
| GGC | Google global cache |
| CNAME | Canonical name |
| RR | Resource record |
| URIs | Uniform resource identifiers |
| URL | Universal resource locator |
| IDE | Integrated development environment |

# References

[1] T. Leighton *Improving Performance on the Internet.* Commun. ACM, 52(2):4451, 2009.

[2] C. Labovitz, S. Lekel-Johnson, D. McPherson, J. Oberheide, and F. Jahanian. *Internet Inter-Domain Traffic.* In Proc. ACM SIGCOMM, 2010.

[3] I. Poese, B. Frank, B. Ager, G. Smaragdakis, and A. Feldmann. *Improving Content Delivery using Provider-Aided Distance. Information.* In ACM IMC, 2010.

[4] G. Maier, A. Feldmann, V. Paxson, and M. Allman. *On Dominant Characteristics of Residential Broadband Internet Traffic.* In Proc. ACM IMC,2009.

[5] Nygren, R. K. Sitaraman, and J. Sun. *The Akamai Network:A Platform for High-performance Internet Applications.* Syst. Rev., 44:219, August 2010.

[6] R. Krishnan, H. Madhyastha, S. Srinivasan, S. Jain,A. Krishnamurthy, T. Anderson, and J. Gao. *Moving Beyond End-to-end Path Information to Optimize CDN Performance.*

[7] Internet evolution *https://atos.net/content/dam/global/ascent-whitepapers/ascent-whitepaper-internet-evolution.pdf*

[8] Schonfeld, E. Eric Schmidts *Gang Of Four:Google, Apple,Amazon,And Facebook.* TechCrunch.Retrieved from https://techcrunch.com/2011/05/31/schmidt-gang-four-google-apple-amazon-facebook/

[9] Bernhard Ager,Wolfgang Mhlbauer,Georgios Smaragdakis,Steve Uhlig *Web Content Cartography.*

[10] Bernhard Ager *Impact of Location on Content Delivery.* http://people.ee.ethz.ch/ bager/papers/A-ILCD-11.pdf

[11] Yuval Shavitt, Udi Weinsberg. *Topological Trends of Internet Content Providers.* SIM-PLEX 12: Proceedings of the Fourth Annual Workshop on Simplifying Complex Networks for Practitioners (2012): 13-18.

[12] Xenofontas Dimitropoulos, Dmitri Krioukov, Marina Fomenkov, Bradley Huffaker, Young Hyun, kc claffy, George Riley *AS Relationships: Inference and Validation.* SIG-COMM Computer Communications Review 37, no. 1(2007): 31-40.

[13] Manuel Palacin, Miquel Oliver, Jorge Infante, Simon Oechsner and Alex Bikfalvi *The Impact of Content Delivery Networks on the Internet Ecosystem.* Journal of Information Policy, Vol. 3 (2013), pp. 304-330

[14] A. Gerber and R. Doverspike. *Traffic Types and Growth in Backbone Networks.* OFC/NFOEC, 2011.

[15] Mike Axelrod *The Value of Content Distribution Networks and Google Global Cache.*

[16] Benjamin Frank, Ingmar Poese, Georgios Smaragdakis, Anja Feldmann, Bruce M. Maggs, Steve Uhlig, Vinay Aggarwal, Fabian Schneider *Collaboration Opportunities for Content Delivery and Network Infrastructures.* 2013

[17] Netflix Open Connect *https://openconnect.netflix.com/en/*

[18] *http://www.hpl.hp.com/research/idl/papers/ranking/adamicglottometrics.pdf.*

[19] Scrapy *http://doc.scrapy.org.*

[20] *Technology: How and Why We Crawl the Web.* Alexa. Archived from the original on April 2, 2014. Retrieved November 6, 2011.

[21] P. Mockapetris *Domain Name System.* https://www.ietf.org/rfc/rfc1034.txt

[22] K RISTOL , D., AND M ONTULLI , L. *HTTP State Management Mechanism.* RFC 2109.

[23] K RISTOL , D., AND M ONTULLI , L. *HTTP State Management Mechanism.* RFC 2965.

[24] Lada A. Adamic ,Bernardo A. *Huberman Zipfs law and the Internet* . Glottometrics 3, 2002,p 143-150

[25] RIPE NCC *RIPE Routing Information Service.* http://www.ripe.net/ris/.

[26] Alexa *http://www.alexa.com/topsites*