

## Search Engine Components

- a. **Crawler** – software robot that receives information about unvisited sites from the index, searches those sites retrieving key words and other significant information concerning the site such as URL's to other sites, and returns that information to the index.
- b. **Query Processor** – provides
  - i. a user front-end to receive a query request, i.e., key words
  - ii. locate the requested information in the database
  - iii. report back to the user which web pages contain those key words; the information reported is limited to that information that one of the may crawlers have located.
- c. **Index** – consists of a
  - i. Database of
    - searched sites which contains URL's, key words, and other significant information concerning each individual site.
    - URL's of identified but unvisited sites
  - ii. Index manager which receives information from a crawler and classifies it into the proper form for inclusion into the database

List the three most significant reasons that most web pages not indexed.

- a. web page that is dynamically created, e.g., Amazon.com pages designed for each individual purchaser reflecting that individuals purchase history
- b. web page consists no text, e.g., only pictures or graphic images
- c. no external web page contains a URL referencing the web page
- d. no crawler has yet reached the web page
- e. web page is password protected

Describe how Google determines the ranking of pages returned by the search engine.

PageRank – count the links to a page; the more links there are to a page, the more relevant it must be. If page A links to page B, then we may consider it to be a vote by A for B. If many sites vote for B, it must be of great interest. If A has a high rank and it votes for B, then B assumes greater rank than if it was only voted for by lower ranking sites.

## Internet Infrastructure

### a) URL

In [computing](#), a **Uniform Resource Identifier (URI)** is a compact [string](#) of [characters](#) used to [identify](#) or name a [resource](#) on the [Internet](#). The main purpose of this identification is to enable interaction with representations of the resource over a network, typically the [World Wide Web](#), using specific [protocols](#). URIs are defined in schemes defining a specific [syntax](#) and associated protocols.

<http://python.ecs.csun.edu/compsci/faculty/putnam.html>

i.e.,

server software://host name/directory/subdirectory/web page

### b) IP Address

An [Internet Protocol \(IP\) address](#) is a numerical identification ([logical address](#)) that is assigned to devices participating in a [computer network](#) utilizing the [Internet Protocol](#) for communication between its nodes.<sup>[1]</sup> Although IP addresses are stored as [binary numbers](#), they are usually displayed in [human-readable](#) notations, such as 208.77.188.166 (for [IPv4](#)), and 2001:db8:0:1234:0:567:1:1 (for [IPv6](#)). The role of the IP address has been characterized as follows: "A [name](#) indicates what we seek. An address indicates where it is. A route indicates how to get there."<sup>[2]</sup>

### c) Domain Name System/Server

- The **Domain Name System (DNS)** is a hierarchical naming system for computers, services, or any resource participating in the [Internet](#). It associates various information with [domain names](#) assigned to such participants. Most importantly, it translates human meaningful domain names to the numerical (binary) identifiers associated with networking equipment for the purpose of locating and addressing these devices world-wide. An often used analogy to explain the Domain Name System is that it serves as the "[phone book](#)" for the Internet by translating human-friendly computer [hostnames](#) into [IP addresses](#). For example, [www.example.com](http://www.example.com) translates to *208.77.188.166*.
- The Domain Name System makes it possible to assign [domain names](#) to groups of Internet users in a meaningful way, independent of each user's physical location. Because of this, [World-Wide Web \(WWW\) hyperlinks](#) and Internet contact information can remain consistent and constant even if the current Internet routing arrangements change or the participant uses a mobile device. Internet domain names are easier to remember than IP addresses such as *208.77.188.166* ([IPv4](#)) or *2001:db8:1f70::999:de8:7648:6e8* ([IPv6](#)). People take advantage of this when they recite meaningful [URLs](#) and [e-mail addresses](#) without having to know how the machine will actually locate them.

- The Domain Name System distributes the responsibility for assigning domain names and mapping them to [Internet Protocol](#) (IP) networks by designating [authoritative name servers](#) for each domain to keep track of their own changes, avoiding the need for a central register to be continually consulted and updated.
- In general, the Domain Name System also stores other types of information, such as the list of [mail servers](#) that accept [email](#) for a given Internet domain. By providing a world-wide, distributed [keyword](#)-based redirection service, the Domain Name System is an essential component of the functionality of the [Internet](#).
- Other identifiers such as RFID tags, UPC codes, International characters in email addresses and host names, and a variety of other identifiers could all potentially utilize DNS <sup>[1]</sup>.
- The Domain Name System also defines the technical underpinnings of the functionality of this database service. For this purpose it defines the [DNS protocol](#), a detailed specification of the data structures and communication exchanges used in DNS, as part of the [Internet Protocol Suite](#) (TCP/IP). The context of the DNS within the Internet protocols may be seen in the following diagram. The DNS protocol was developed and defined in the early 1980s and published by the [Internet Engineering Task Force](#) (cf. History).

#### d) TCP/IP

The **Transmission Control Protocol (TCP)** is one of the core protocols of the [Internet Protocol Suite](#). TCP is so central that the entire suite is often referred to as "TCP/IP." Whereas IP handles lower-level transmissions from computer to computer as a message makes its way across the Internet, TCP operates at a higher level, concerned only with the two *end systems*, for example a Web browser and a Web server. In particular, TCP provides reliable, ordered delivery of a stream of bytes from one program on one computer to another program on another computer. Besides the Web, other common applications of TCP include [e-mail](#) and [file transfer](#). Among its management tasks, TCP controls message size, the rate at which messages are exchanged, and network traffic congestion.

TCP provides a communication service at an intermediate level between an application program and the [Internet Protocol](#) (IP). That is, when an [application program](#) desires to send a large chunk of data across the Internet using IP, instead of breaking the data into IP-sized pieces and issuing a series of IP requests, the software can issue a single request to TCP and let TCP handle the IP details.

IP works by exchanging pieces of information called [packets](#). A packet is a sequence of [bytes](#) and consists of a *header* followed by a *body*. The header describes the packet's destination and, optionally, the [routers](#) to use for forwarding—generally in the right direction—until it arrives at its final destination. The body contains the data which IP is

transmitting. When IP is transmitting data on behalf of TCP, the content of the IP packet body is TCP payload.

Due to network congestion, traffic load balancing, or other unpredictable network behavior, IP packets can be lost or delivered out of order. TCP detects these problems, requests retransmission of lost packets, rearranges out-of-order packets, and even helps minimize network congestion to reduce the occurrence of the other problems. Once the TCP receiver has finally reassembled a perfect copy of the data originally transmitted, it passes that datagram to the application program. Thus, TCP abstracts the application's communication from the underlying networking details.

The **Internet Protocol (IP)** is a [protocol](#) used for communicating data across a [packet-switched internetwork](#) using the [Internet Protocol Suite](#) (TCP/IP).

IP is the primary protocol in the [Internet Layer](#) of the [Internet Protocol Suite](#) and has the task of delivering datagrams (packets) from the source host to the destination host solely based on their addresses. For this purpose the Internet Protocol defines addressing methods and structures for datagram [encapsulation](#). The first major version of addressing structure, now referred to as [Internet Protocol Version 4 \(IPv4\)](#) is still the dominant protocol of the Internet, although the successor, [Internet Protocol Version 6 \(IPv6\)](#) is actively deployed worldwide.

Data from an [upper layer protocol](#) is encapsulated as [packets/datagrams](#) (the terms are basically synonymous in IP). [Circuit](#) setup is not needed before a host may send packets to another host that it has previously not communicated with (a characteristic of [packet-switched](#) networks), thus IP is a [connectionless protocol](#). This is in contrast to [Public Switched Telephone Networks](#) that require the setup of a circuit before a phone call may go through (*connection-oriented* protocol).

Because of the abstraction provided by encapsulation, IP can be used over a [heterogeneous network](#), i.e., a network connecting computers may consist of a combination of [Ethernet](#), [ATM](#), [FDDI](#), [Wi-Fi](#), [token ring](#), or others. Each link layer implementation may have its own method of addressing (or possibly the complete lack of it), with a corresponding need to resolve IP addresses to data link addresses. This address resolution is handled by the [Address Resolution Protocol](#) (ARP) for [IPv4](#) and [Neighbor Discovery Protocol](#) (NDP) for [IPv6](#).