



The Web IR - Challenges

Vasudeva Varma

IIIT Hyderabad

The web and its challenges

- Unusual and diverse documents
- Unusual and diverse users, queries, information needs
- Beyond terms, exploit ideas from social networks
 - link analysis, clickstreams ...



SURFACE WEB

Google
Bing Wikipedia

DEEP WEB

Contains 90% of the information on the Internet, but is not accessible by Surface Web crawlers.

Academic Information
Medical Records
Legal Documents
Scientific Reports
Subscription Information

Multilingual Databases
Financial Records
Government Resources
Competitor Websites
Organization-specific Repositories

Social Media

(DARK WEB)

A part of the Deep Web accessible only through certain browsers such as Tor designed to ensure anonymity. Deep Web Technologies has zero involvement with the Dark Web.

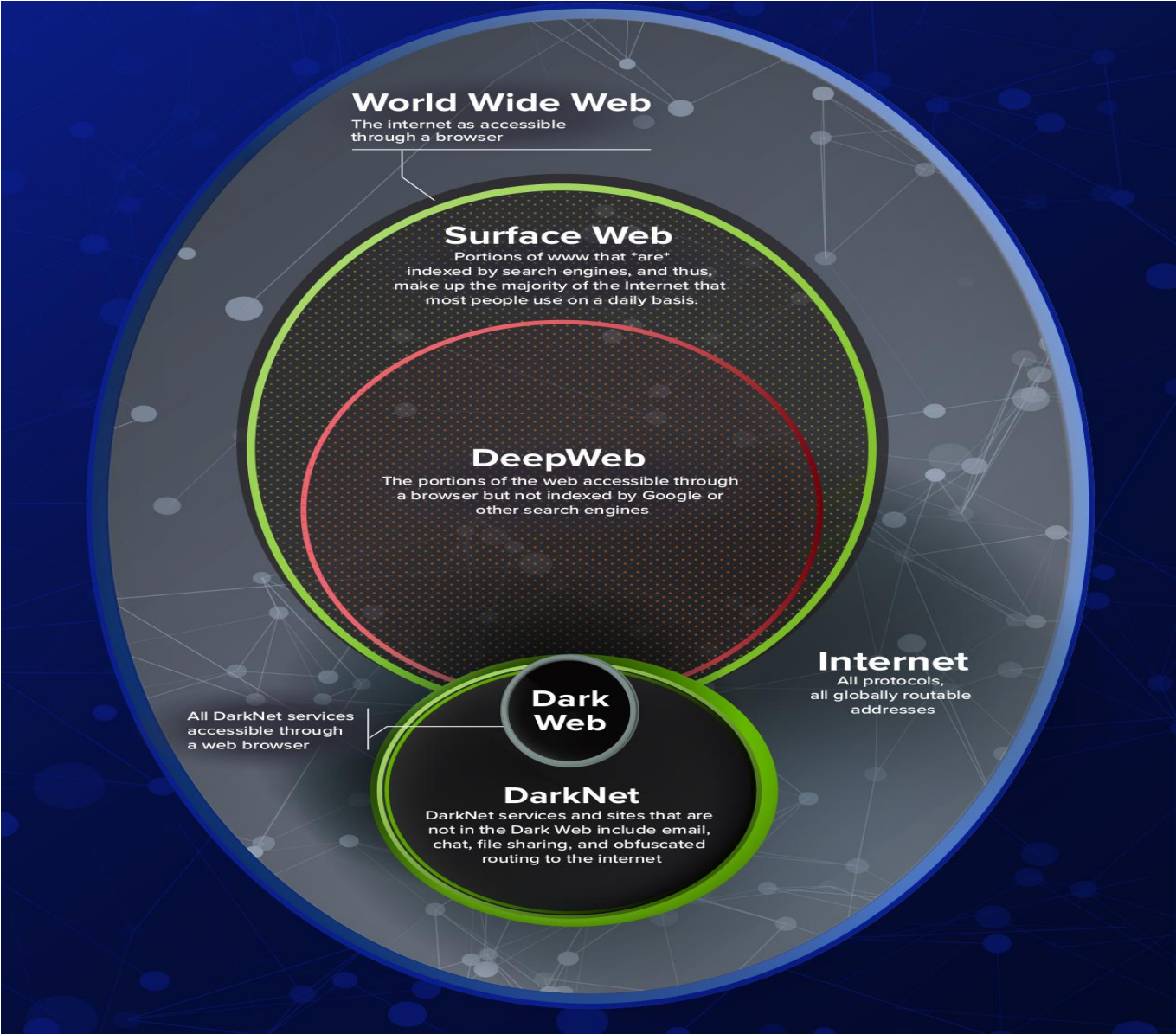
Illegal Information
TOR-Encrypted sites

Political Protests

Drug Trafficking sites
Private Communications

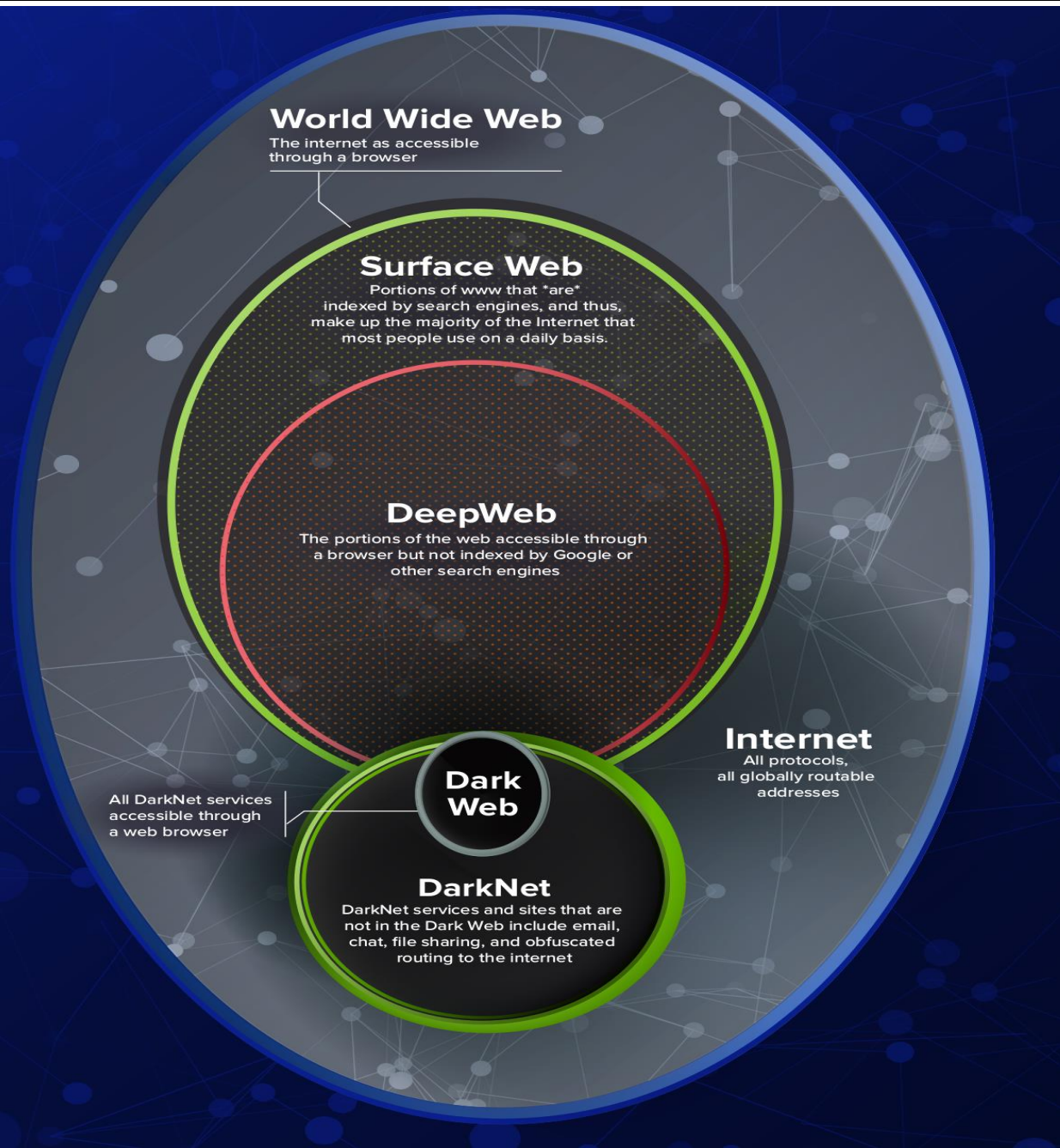
Surface web, deep web, dark web...





The web: size

- What is being measured?
 - Number of hosts
 - Number of (static) html pages
 - Volume of data
- Number of hosts – netcraft survey
 - http://news.netcraft.com/archives/web_server_survey.html
 - Gives monthly report on how many web servers are out there
- Number of pages – numerous estimates
 - For a Web engine: how big its index is
- <https://www.internetlivestats.com/>



*What is
the size of
DEEP web?*

Crawling...

Crawling overview

- Types of crawlers
- Functionality
 - Start the crawl: Seed URLs
 - ...
 - End the crawl: Halting Criteria
- Policies
- Architecture of a Web crawler

Types of crawlers

- Search engine crawlers
- Enterprise crawlers
- Monitoring crawlers
 - Copy right violations
 - DRM crawlers
 - Malware detection
 - Web analytics
- Document feeds (RSS/Atom or commercial feeds)

Functionality of the crawlers

- Start with seed URLs
 - Selection of seed URLs is important
 - Quality (avoid spam/objectionable/non-hub pages)
 - Importance (popularity/trustworthiness/reliability)
 - Potential yield documents
 - Web graph helps pick right seed URLs
- Survive
 - Avoid crawler traps
 - Causes infinite number of requests being made
 - Infinitely deep directory structures
 - Follow the rules and behave well (adhere to ***policies***)
- End when time comes
 - Some crawlers are designed to go on forever
 - Some stop when a particular criteria is met (after reaching depth K, after crawling N pages or time T, after Index reaches K Units)

Policies

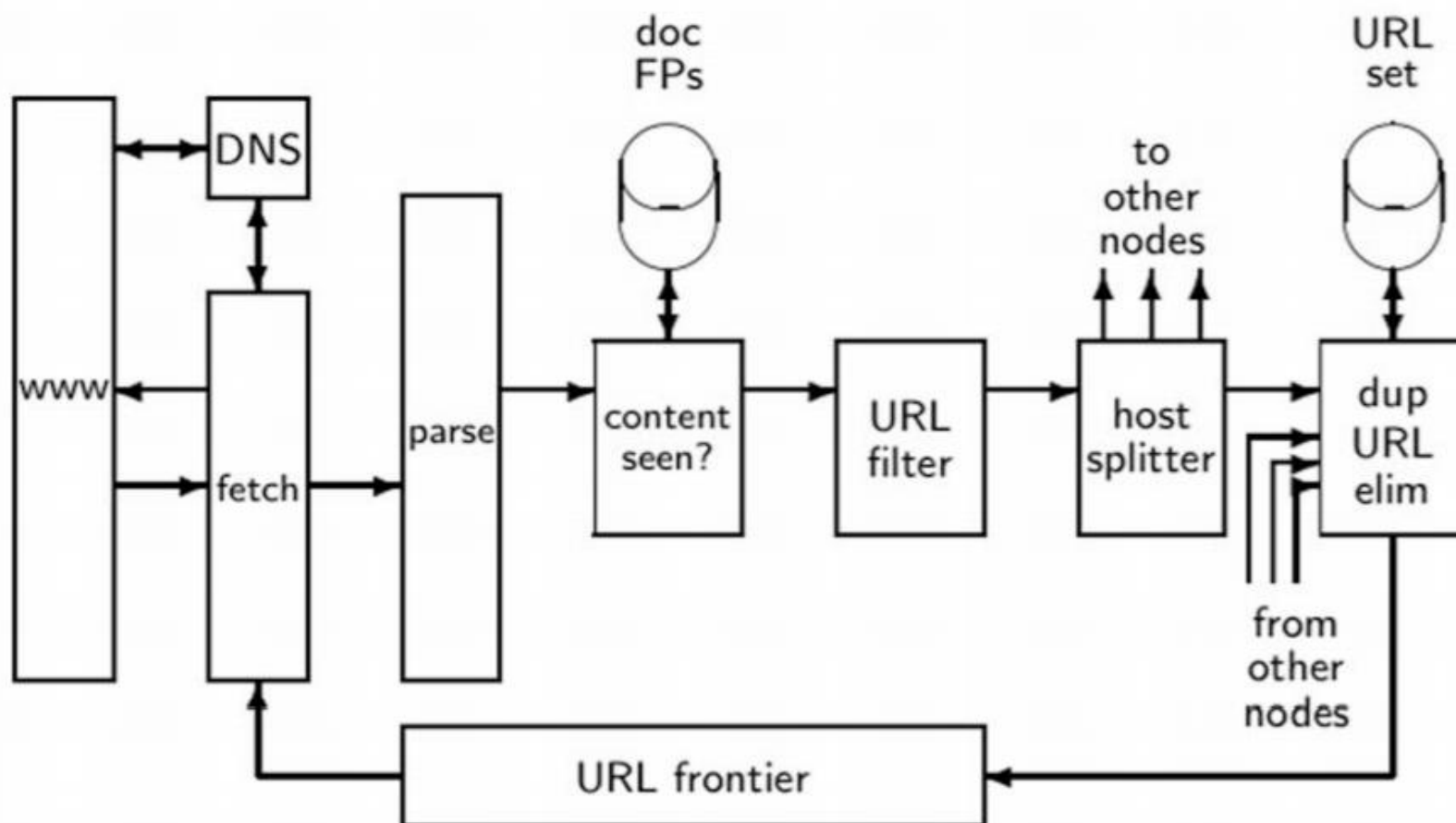
The behavior of a Web crawler is the outcome of a combination of policies

- Selection policy: states which pages to download
 - Prioritization: predict high yield pages from web graph
- Revisit policy: states when to check for changes to the pages
 - goal: high avg Freshness and low avg age
 - Two policies: Uniform policy or proportional policy
- Politeness policy: states how to avoid overloading Web sites
 - Robots.txt
 - Sitemap – organize the site to control crawling its parts
 - Meta tag: `<META NAME "ROBOTS" CONTENT="NOINDEX, NOFOLLOW">`
- Parallelization policy: states how to coordinate distributed Web crawlers

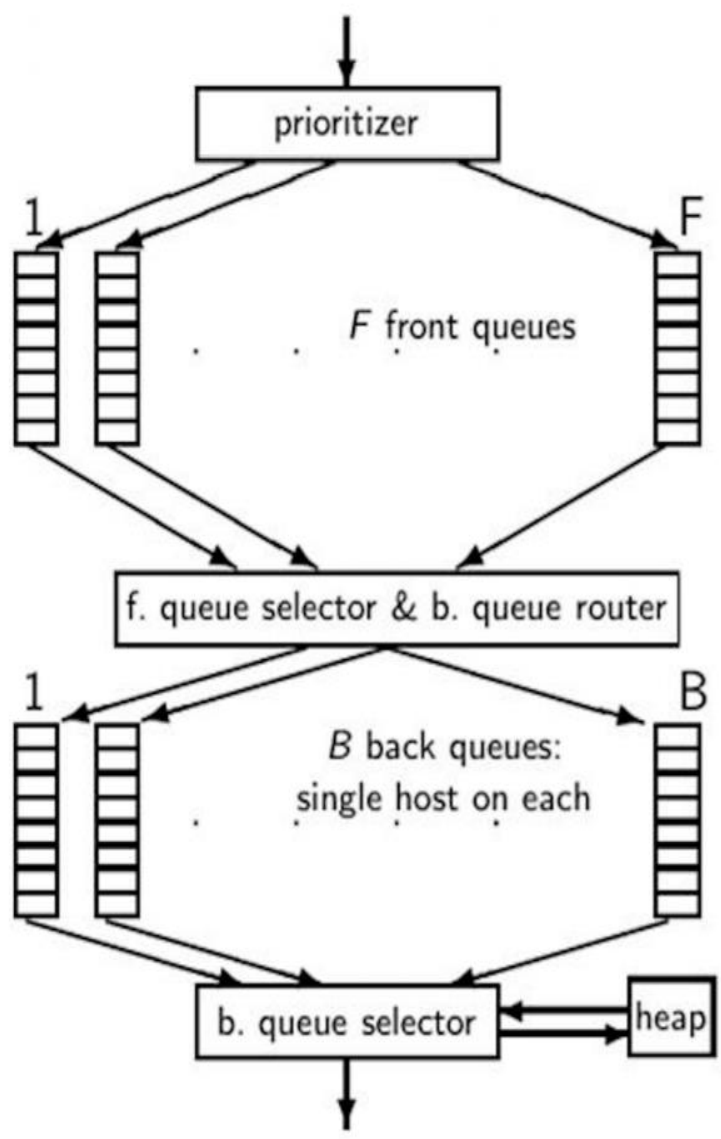
Challenge: How to handle politeness and priority simultaneously?

- An ideal crawler needs to frequently visit high priority pages but also needs to be polite.
- To achieve this we implement a URL frontier with the following goals.
 - only one connection is open at a time to any host
 - a waiting time of a few seconds occurs between successive requests to a host
 - high-priority pages are crawled preferentially

Distributed Crawler Architecture



URL Frontier Design



Examples of Popular Web Crawlers

- GooleBot
 - BingBot
 - MSNbot
 - Slug
 - Yahoo!Slurp
- Nutch
 - Scrapy
 - DataParkSearch
 - Grub
 - Heritrix

Further reading

- Introductory article on deep web [link](#)
- Section on crawling the deep web, from [this](#) university's guide on deep web:
- Optional Reading: [Paper on crawling the deep web](#)
- Victor Lavrenko short [videos on web crawling](#)
- Implementation [view](#)