

L'Internet: un nouveau mode de communication et de diffusion de l'information

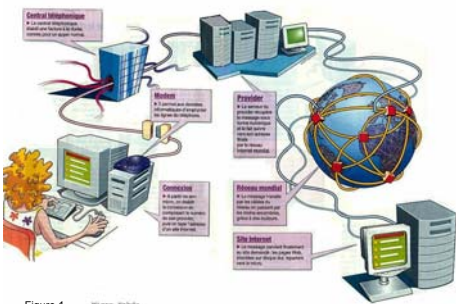


Figure 1

Internet est un vaste réseau formé de centaines de réseaux dans le monde entier. Tous les ordinateurs connectés à ce réseau communiquent entre eux grâce à un protocole commun, le protocole TCP/IP et utilisent l'architecture client/serveur.

Les composantes d'Internet les plus utilisées:

E-Mail:

L'E-Mail ou messagerie électronique est un service qui permet d'envoyer et de recevoir des messages textuels ou des fichiers sur Internet.

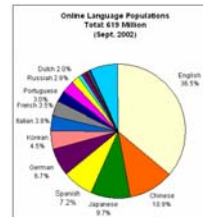
FTP:

FTP ou « File Transfer Protocol » est un protocole qui permet le chargement et le téléchargement de fichiers via Internet

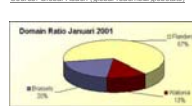
World Wide Web:

Le WWW, basé sur le principe de l'hypertexte, permet de naviguer de manière conviviale parmi les ressources du web. Le web fournit une interface unique d'accès à différents types de ressources, à des fichiers de formats divers (web, image, vidéo, etc...) ainsi qu'aux différents protocoles d'accès (FTP, Telnet, e-mail...)

Figure 2



Source: Global Reach (global-reach.biologista.it)



Source: Stratemet

L'Internet moyen est un homme âgé entre 15 et 50 ans, il est majoritairement anglophone et navigue avec Internet Explorer 5.x sous Windows 98.

67% des Américains utilisent Internet contre seulement 32% des belges.

Les activités principales sur Internet sont: surfer sur la toile: 82% et utiliser le courrier électronique: 61.5%.

La Belgique compte 3.2 millions d'internautes réguliers en 2002, à savoir une augmentation de 16% par rapport à 2001 (Belgian Internet Mapping de InSites Consulting).

Le nombre de Wallons connectés reste plus bas que la moyenne nationale.

Le web: des sites de rencontre, chat, hot ... mais pas uniquement



Chaque image, page, dossier sur le Web possède sa propre adresse. Cette adresse est appelée URL (Uniform Resource Locator).

Elle identifie l'endroit exact sur le réseau où se situe une ressource ainsi que le protocole nécessaire (http, ftp, telnet, etc) pour pouvoir accéder à l'information.



L'adresse Internet se compose en général de la façon suivante:



HTML

HTML (Hyper Text Markup Language) est un langage de description de documents qui permet de concevoir les pages web et de décrire les liens hypertextes.

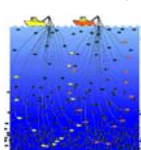


Source: Wamadoo.fr

L'HTML offre la possibilité de faire des liens à l'intérieur d'un document, vers un autre document html, vers des images, des vidéos, des fichiers audio...

Le logiciel navigateur utilise les informations contenues dans le langage HTML pour restituer à l'écran l'aspect des pages web.

L'information sur le web peut être scindée en deux catégories:



Web visible ou « surface web » qui regroupe les pages statiques

Web invisible ou « deep web » regroupe l'ensemble des pages web générées dynamiquement à la demande de l'utilisateur.

Figure 5

Le web invisible représente environ 550 milliards de documents uniques soit 260 fois le web visible. Les 60 sites les plus importants représentent à eux seuls plus de 40 fois le volume du web visible.

Ce sont les pages accessibles par mot de passe ou exigeant une autorisation préalable, celles pour lesquelles il faut remplir un formulaire pour afficher les pages, plus important encore c'est toute l'information résidant dans les bases de données traditionnelles -qu'elles soient bibliographiques ou textuelles.

Le contenu du Web Invisible émanant de professionnels est de qualité supérieure au web visible.

Le Web Invisible est la catégorie augmentant le plus rapidement sur le net.

95% du Web Invisible est accessible gratuitement sans restriction ni inscription.

Distribution de l'information sur le web

Le Web Visible

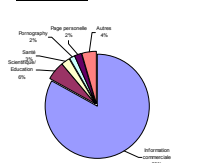


Figure 7

Le Web Invisible

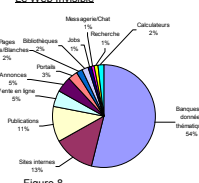


Figure 8

Les pages du Web Visible sont principalement des pages commerciales.

Vu la facilité de publication sur le Web Visible et le fait qu'aucun contrôle de qualité ne s'applique à ces pages, le contenu informationnel est très pauvre.

Par exemple: pages personnelles, photos de vacances, pages publicitaires...

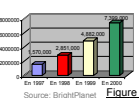
Le contenu du Web Invisible est principalement constitué de bases de données scientifiques à haute valeur ajoutée, le plus souvent l'information est validée par des comités de scientifiques.

Publications d'articles scientifiques en ligne: Pre-print Loos Alamos en Physique...

Bases de données scientifiques: PubMed, la base de données biomédicales...

Trouver ce que l'on cherche et même ce que l'on ne cherchait pas...

Evolution du nb de sites web



Source: BrightPlanet

Le nombre de sites web étant en croissance continue, il devient de plus en plus difficile de trouver l'information pertinente sur Internet, information qui est le plus souvent diluée parmi une quantité importante d'informations inutiles et non désirées.

Les outils permettant d'effectuer des recherches sur Internet sont apparus en grand nombre pour répondre à cette demande. Parmi ces outils on peut discerner les services suivants:

- Les annuaires de recherche:** Sites présentant des listes de sites classés de façon thématique
 - Points forts:** - simple d'utilisation
- référencement des sites incontournables
- classement parmi des catégories structurées
 - Points faibles:** - non exhaustif
- mise à jour souvent lente
- Les moteurs de recherche:** Outil basé sur l'utilisation d'un robot (programme) qui parcourt toutes les pages et les indexe dans une base de données
 - Points forts:** - relativement exhaustif
- recherche par mot-clés
- différentes options de recherche (images,...)
 - Points faibles:** - nombre de réponses souvent très élevé
- beaucoup de parasites
- Les méta-moteurs:** Outils permettant de consulter différents moteurs de recherche en une seule interrogation.
 - Points forts:** - simple d'utilisation
- référencement des sites incontournables
- classement parmi des catégories structurées
 - Points faibles:** - non exhaustif
- mise à jour souvent lente

Principes de fonctionnement des moteurs de recherche

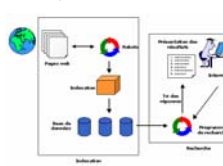


Figure 11

Principaux moteurs de recherche:

- GOOGLE: www.google.com
- MSN Search: www.msn.com
- HotBot: www.hotbot.com
- AltaVista: www.altavista.com
- AltaVista: www.altavista.com
- AltaVista: www.altavista.com

Annuaire/répertoire:

- Yahoo: www.yahoo.com
- Cherche: www.cherche.com
- répertoire de sites en chimie
- InfoBot: www.infobot.com
- Annuaire téléphonique belge

Méta-moteur de recherche:

- Proton: www.proton.com
- Doopile: www.doopile.com
- Metacrawler: www.metacrawler.com
- Fastie (SearchOnline.info): www.fastie.com
- Excite: www.excite.com
- Webcrawler: www.webcrawler.com

Taux d'utilisation des moteurs de recherche



Figure 12

Google est actuellement le moteur de recherche le plus complet et le plus utilisé en Europe et dans le monde. 73% des pages web répertoriées sont indexées, 1.75% de la base de données concerne des formats autres que le format html (pdf, gif...). 25% des pages ne sont pas encore indexées et seulement 0.15% de la base de données est réindexée journalièrement.

Taux d'utilisation des moteurs de recherche

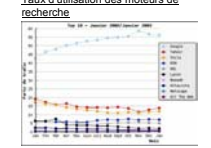


Figure 13

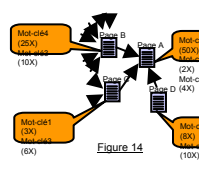


Figure 14

Google classe les sites par ordre de pertinence par rapport à une requête exprimée par l'utilisateur, ce classement est effectué grâce à un algorithme complexe appelé PageRank qui fait intervenir le nombre de citations reçues par cette page et procède à l'analyse du contenu de la page qui contient le lien pour juger de l'importance de celle-ci et les mot-clés qu'elle contient.