

La chaîne documentaire

Cours « Document et Web Sémantique »

Nicolas Malandain, Nicolas Delestre

- 1 Le document
 - Description
 - Contenu & Forme
 - Diffusion
 - Interprétation

- 2 La chaîne de traitements
 - Problématique
 - Reconstruction document papier
 - Interprétation

- 3 Le sens du document
 - Interprétation
 - Indexation

- 4 Conclusion

Quelques définitions

Document

- Larousse
« pièce écrite servant d'**information**, de preuve »
- Wikipedia
« Un document renvoie à un ensemble formé par un **support** et une **information** (le contenu), celle-ci enregistrée de manière persistante. » (2021)

Données Suite de caractères représentant une valeur

Information Données contextualisées de manière formelle ou informelle

Connaissance Information qui permet de valider une information et/ou d'inférer une nouvelle information

Évolution des supports

À voir : « L'Odyssée de l'écriture »

La pierre les peintures rupestres

La tablette d'argile en Mésopotamie 4^e millénaire avant JC
recto-verso, lecture haut-bas, droite-gauche (forme la plus ancienne d'écriture)

Le Papyrus Égypte et proche orient 3^e millénaire avant JC

- craint le pliage, l'humidité et le feu
- apparition du *Volumen* : rouleau formé de deux axes verticaux, que l'on déroule.
- longueur jusqu'à 15 mètres
- nécessite les deux mains pour le dérouler et lire
- ne permet pas de relire quand on écrit, pas de marquage possible
- très utilisé pour la lecture à haute voix

Le parchemin Pergame (aujourd'hui Bergama, Turquie) 2^e millénaire avant JC

- fabriqué à partir de peau (mouton/chèvre)
- souple
- écriture des deux côtés
- apparition du *Codex* : assemblage de parchemins
- la bible est copiée par les chrétiens sur codex

Le papier Chine 2^e siècle avant JC

- monopole de 6 siècles
- se répand à partir de 751
- arrive en Europe au 12^e siècle
- support souple et léger, il se conserve mieux
- au 15^e siècle apparition de l'imprimerie avec Gutenberg

Dématérialisation de nos jours le document se dit électronique ou numérique et n'a plus de support spécifique.

Caractérisations

Un document se caractérise par :

- Son contenu
- Sa forme
- Sa diffusion
- Son utilisation
- Son interprétation

Contenu 1 / 2

Informel, semi-formel et formel

Alan Turing

☛ Pour les articles homonymes, voir Turing.

Alan Mathison Turing, né le 23 juin 1912 à Londres et mort le 7 juin 1954 à Wilmslow, est un mathématicien et cryptologue britannique, auteur de travaux qui fondent scientifiquement l'informatique.

Pour résoudre le problème fondamental de la décidabilité en arithmétique, il présente en 1936 une expérience de pensée que l'on nommera ensuite machine de Turing et des concepts de programme et de programmation, qui prendront tout leur sens avec la diffusion des ordinateurs, dans la seconde moitié du xx^e siècle. Son modèle a contribué à établir la thèse de Church, qui définit le concept mathématique intuitif de fonction calculable.

Durant la Seconde Guerre mondiale, il joue un rôle majeur dans la cryptanalyse de la machine Enigma utilisée par les armées allemandes. Ce travail secret ne sera connu du public que dans les années 1970. Après la guerre, il travaille sur un des tout premiers ordinateurs, puis contribue au débat sur la possibilité de l'intelligence artificielle, en proposant le test de Turing. Vers la fin de sa vie, il s'intéresse à des modèles de morphogénèse du vivant conduisant aux « structures de Turing ».

Poursuivi en justice en 1952 pour homosexualité, il choisit, pour éviter la prison, la castration chimique par prise d'œstrogènes. Il est retrouvé mort par empoisonnement au cyanure le 8 juin 1954 dans la chambre de sa maison à Wilmslow. La reine Elisabeth II le reconnaît comme héros de guerre et le gracie à titre posthume en 2013.

230	A. M. TURING	[Nov. 12,
ON COMPUTABLE NUMBERS, WITH AN APPLICATION TO THE ENTSCHEIDUNGSPROBLEM		
By A. M. TURING.		
[Received 28 May, 1936.—Read 12 November, 1936.]		
<p>The "computable" numbers may be described briefly as the real numbers whose expressions as a decimal are calculable by finite means. Although the subject of this paper is ostensibly the computable numbers, it is almost equally easy to define and investigate computable functions of an integral variable or a real or computable variable, computable predicates, and so forth. The fundamental problems involved are, however, the same in each case, and I have chosen the computable numbers for explicit treatment as involving the least cumbersome technique. I hope shortly to give an account of the relations of the computable numbers,</p>		

Naissance	23 juin 1912 Maida Vale (Londres) (Royaume-Uni)
Décès	7 juin 1954 (à 41 ans) Wilmslow (Cheshire) (Royaume-Uni)
Domicile	Wilmslow (Cheshire, Angleterre)
Nationalité	 Britannique
Domaines	Informatique, mathématiques, logique, cryptanalyse
Institutions	Université de Manchester National Physical Laboratory Université de Cambridge
Diplôme	Université de Manchester Université de Princeton
Renommé pour	Problème de l'arrêt Machine de Turing Cryptanalyse d'Enigma

Contenu 2 / 2

Avant le document numérique

- contenu principalement unique : texte / image / vidéo / voix / ...
- accès séquentiel ou direct
- quelques mélanges : livre avec texte et image
- duplication pouvant entraîner des pertes

À l'ère du document numérique

- contenu totalement hétérogène : texte + image + vidéo + voix ...
- accès séquentiel / direct / hypertextuel
- duplication 100% conforme à l'original

Forme

Avant le document numérique

- support physique
- humainement déchiffrable
- stockage volumineux et relativement pérenne

À l'ère du document numérique

- le support est lié au stockage, mais non à l'utilisation
- stockage peu volumineux
- quasiment humainement indéchiffrable
- le support est pérenne mais les documents non (problème de format)

Accélération de la diffusion

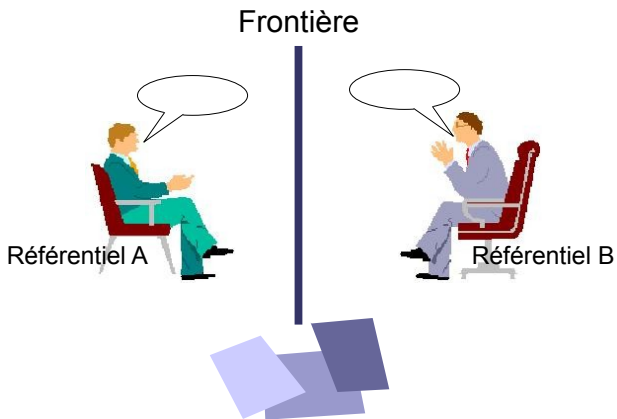
Changement du support

- poids
- transport
- copie
- disparition du support

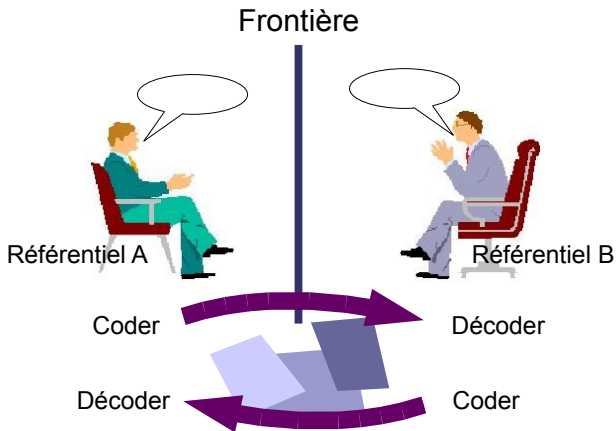
Changement de la publication

- rédaction : des “élites” au quidam
- publication : des éditeurs au quidam
- moins de sélection
- disponibilité immédiate

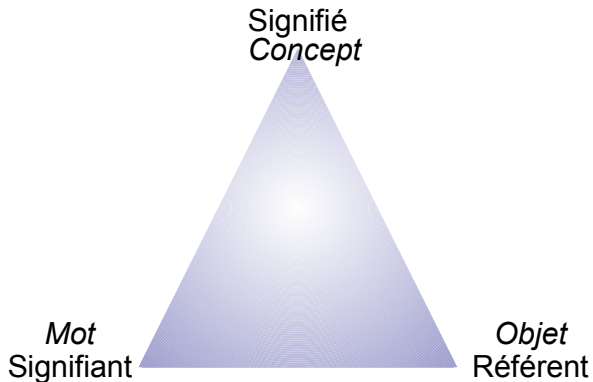
Le mécanisme de la communication



Le mécanisme de la communication



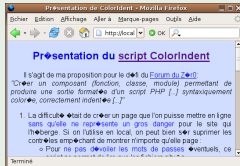
Le triangle sémiotique



L'interprétation d'un document 1 / 2

Elle dépend

- écriture/lecture asynchrone
- écriture/lecture référentiels différents
 - homme ↔ homme / homme ↔ machine / machine ↔ machine
- codages multiples
- mises en forme



http://si2ra.ouvaton.org/si2ra/pub/mode_d039emp/installation/un_encodage_correcte.php

自己PR

私はコンサルティング会社に勤めていた新卒。システムエンジニアとして、多種多様なプロジェクトに関わってきました。これらの経験を通じて、現在の私は、どのような環境においても豊く貢献できるようになりました。また、外向的な性格を活かしたコミュニケーション能力と、チームワークを大切にすることを大切にしています。

これまで担当したプロジェクトでは、プログラマーとしての職務経験のみならず、プロジェクトマネージャーの業務アプリケーション開発担当として責任あるポストも経験しており、プログラムスキルに加えて、基礎的なマーケティング能力も身に付けております。

私は、仕事においても、趣味においても、新しいことに意欲的にチャレンジし、常に目標を持って取り組む性格です。今後、貴社へ入社することになれば、積極的にプロジェクトマネージャーとして活躍したいと考えております。貴社で活躍いただける機会、今まで勇に付けてきた経歴やコミュニケーション能力を全力で、幅広い分野で貴社に貢献していく所存です。

私はフランス人ですが、母が日本人であるおかげで、フランス文化と日本文化の基本的な違いは心得ていますので、プロジェクトの進め方やチーム管理にもなるものと思います。また、英語とフランス語の両方を活かせる貴社の取引先と関係がある貴社での仕事を希望しております。

ご採用いただければ、必ず精進に努める成果が出せるよう全力を尽くしますので、どうぞよろしくお願ひ致します。

http://www.tartoinjapan.com/fr/blog_post-157-cv_lettre_motivation.html

Pierre Laurence

L'interprétation d'un document 2 / 2

Elle nécessite

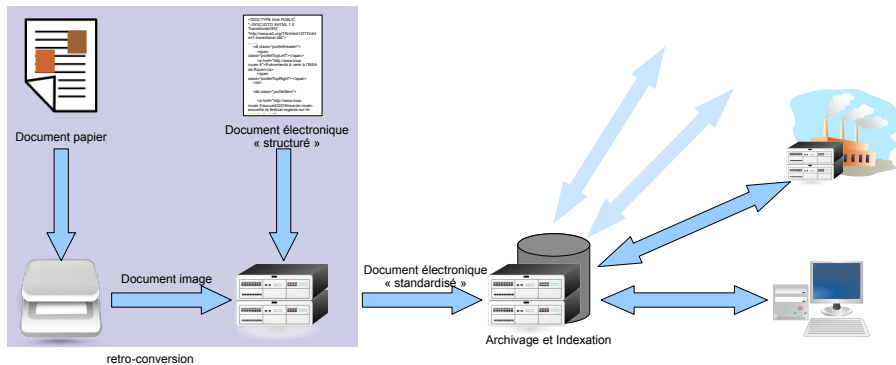
- L'ajout de données sur (ou au sein d') un document : métadonnées
 - Par exemple : l'auteur du document, la date de parution
- Ces métadonnées peuvent être intrinsèques ou externes

Une chaîne de traitements, pourquoi ?

Quels sont les problèmes liés aux documents ?

- Pérennité
 - la masse de documents papiers, audio et vidéo (analogiques) existante
 - la masse d'information papier/audio/vidéo en perpétuelle augmentation
 - la masse de documents électroniques en perpétuelle augmentation
 - la perte importante de documents
 - le stockage
- Accès à l'information
 - accès distant
 - mise en relation de documents
 - recherche d'informations : une aiguille dans une meule de foin
 - langues
 - méta-informations
 - mises à jour
 - format

Présentation d'une chaîne de traitements générale



Exemple de chaînes de traitements

- Projets à grande échelle : BnF, Google (Google books)
- GED pour les entreprises : Nuxeo
- GED pour les particuliers : Paperwork
<https://github.com/jflesch/paperwork>

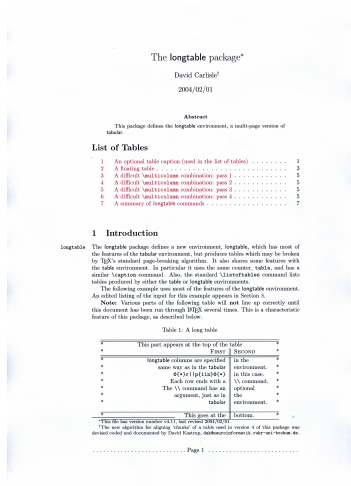
Numérisation

3 “formats”

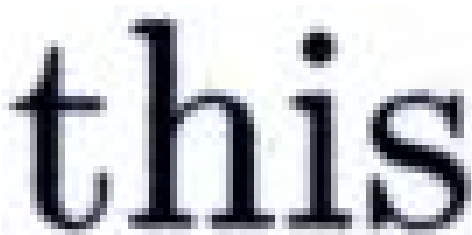
- format papier
- format électronique “image”
- format électronique “structuré”

Le format image

issu de l'utilisation d'un scanner



- information de plus bas niveau
- format bitmap décrivant chaque pixel
- seul le pixel est “connu” de l'ordinateur
- volumineux



Difficultés 1 / 2

Variabilité du contexte

- mise en page : journal, magazine, livre, article scientifique
- culture : langue, sens de lecture
- ...

Variabilité de la présentation

- police de caractère
- taille, style
- alignement
- ...

Difficultés 2 / 2

Structuration

- caractères, mots, paragraphes, ...
- structure logique typée
 - titres, sous titres, ...
 - colonnes
 - tableaux
 - types d'images (photo, schéma, ...)
 - fil de lecture



NÉGOCIER UN RABAI

VINCENT GAUVIN, de Bazacoan, a fait ce qu'il est devenu d'appeler une affaire. Il nous raconte comment.

« Il y a peu, je décide de m'acheter un ordinateur IBM PC. Mais, déçu par le prix de moyen, je me tourne vers le marché de l'occasion. Mon choix se porte vers un Amstrat PC 1512 800 monochrome. J'ai constaté avec surprise que les prix pratiqués étaient de 70 à 15 % inférieurs à ceux du marché neuf. J'ai donc décidé d'acheter ce modèle. Mais, sachant que certains hypermarchés vendent du PC, je me rends dans l'un d'eux (Centrom de Bénes), pour de surprise, les prix y sont presque les mêmes que chez un distributeur classique. C'est alors que j'ai l'idée de demander au chef de rayon si je pourrais éventuellement acquies un modèle de démonstration, qui est en très état, en bénéficiant d'un rabais important. Il me répond qu'il n'est pas question pour lui de se séparer de ce bien tant qu'il a une bonne affaire. Mais il ajoute qu'étant donné la difficulté pour eux d'écouler ce matériel, en partie à cause de l'incompétence des vendeurs, il me contacte un rabais de 5 500 F sur un modèle neuf.

C'est ainsi que j'ai acheté un PC 1512 800, avec monnaie contre plusieurs jeux, Kakomai, Evolution Sunset et Superbase, pour 6 500 F avec une garantie d'un an ! Il me semble que les autres chaînes d'hypermarchés ont les mêmes problèmes, tout au moins en province. Peut-être est-ce intéressant à avoir les lecteurs de S.V.M. ? »

UNE MACRO POUR SYMPHONY, de Lucien Bossourol de Taverny.

« Votre exemple d'appel à Print Graph à partir de la feuille de calcul de Symphony peut être amélioré pour permettre le basculement rigide de la feuille à cet utilitaire d'impression. L'incorporation de votre exemple vient de ce qu'après le détachement de MS-DOS, un mot et appel à la machine entraîne l'erreur : 'DOS déjà attaché'. Voici la solution :
G CASEROLEUR PGRAPH
SERVICES ARDOS
PGRAPH SERVICES
...ARDOS C ; SYM PGRAPH
En notant la macro en deux champs nommés. J'erreur entraîne le saut sur la seconde ligne et permet le rappel de PGraph. La commande sortie de PGraph appelle alors la feuille de calcul.

Quelques précisions sur le choix des polices de caractères sur Macintosh, d'origine, qui a pu être de signer.

« Pour alléger le système du Macintosh, vous conseillez au

LES FONCTIONS CACHÉES DES LOGICIELS

utilisateurs d'éliminer les fontes de grande taille. Ceci n'est valable que si on se soucie pas

police New York, en impression de qualité supérieure donc, sur imprimante L.

NEW YORK CORPS 12, avec corps 24 dans le système

NEW YORK CORPS 12, sans corps 24 dans le système

drop de la qualité d'impression. En effet, lorsqu'on utilise une fonte de taille 12, par exemple, il est utile d'avoir la même fonte en 24, car lors de l'impression en qualité supérieure, le système du Macintosh recherche cette dernière ; si elle n'existe pas, l'ordinateur imprime avec la même fonte de points que celle qui figure à l'écran. Dans ces conditions, le résultat peut être très médiocre, surtout dans le cas de polices de caractères au dessin complexe (cascades, gothiques, etc.). En revanche, si le corps 24 est disponible, le système calcule à partir de ce dernier le dessin des caractères à imprimer, et le résultat est nettement meilleur.

Voici un exemple obtenu avec la

LA GUERRE DES SALONS

A LA SUITE DE NOTRE ARTICLE « La guerre des salons est déclarée » (S.V.M. Actualités n°47), Max Hérmau, président du SICOB, nous écrit : « Mes Poyes a effectivement passé quatre ans dans l'équipe du SICOB, mais le SICOB (qui) a fondé avec quelques collègues en 1950 existe depuis 30 ans. Il est donc prévisible et légitime que Mes Poyes affirme être sa son origine. Elle a par ailleurs été à l'origine des positionnements réussis : de la manifestation PNRD de septembre dont elle est aujourd'hui le chef de file jusqu'à l'été 1987 (convention de la Convention informelle et l'Exposition S.V.M.) qui effectivement ne peuvent être considérées comme des réussites : ce sont manifestations qui ont été placées sous la responsabilité de Mes Poyes. Cela a contribué au SICOB à coexister avec le S.V.M. et le SYSTEC, à soutenir la Convention informelle. Pour l'avenir, le SICOB continue à servir la profession pour acquies et au nom de la quelle il est organisé, tout parti de ses ressources, corrigent ce qui doit être amélioré. S'adaptant sans cesse depuis 30 ans aux besoins de ses exposants, le SICOB n'a jamais fait à sa mission... »

DÉBORDÉ PAR LE SUCCÈS

LE PROGRAMME LE CIEL, proposé gratuitement dans S.V.M., est en train de faire une victime : son auteur, Jean-Jacques Tigney, qui croit sous les demandes d'envoi.

« En un mois, nous écrits, j'ai reçu environ 500 demandes, dont une certaine sans les timbres pour le retour. Toutes les demandes accompagnées de timbres ou de diverses contri-

butions spontanées ont été satisfaites. Mon je ne puis continuer indéfiniment à ce rythme, sous peine de perdre la santé. Je ne souhaite faire savoir que je ne suis plus en mesure de continuer à diffuser moi-même le programme. Ce sera écrit, et je ne puis par un ou plusieurs organismes spécialisés dans la distribution des logiciels du domaine public, tractions en cours, surveiller la presse autonome et informatique. Les personnes qui ont pas joint de timbres pour le retour de leur demande sont priées de le faire (tim : 70 F ; rapide : 5,00 F ; échanger : échanger réponse (internationaux), et de rappeler leurs nom et adresse, afin que je retourne leur demande.

Interpréter un document

Extraire du sens
Mais dans quel but ?

L'interprétation d'un document n'a de sens qu'en fonction du besoin

La "compréhension" en fonction de l'objectif












- **indexation**
- recherche d'information particulière (extraction)
- correction orthographique / grammaticale
- veille technologique
- résumé automatique
- traduction automatique
- ...

Documents à indexer

Contenu d'un document électronique

- texte
- image
- son
- video

une indexation principalement textuelle

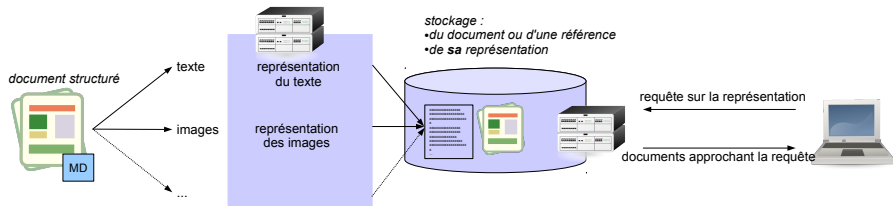
Describes without errors	Describes with minor errors	Somewhat related to the image	Unrelated to the image
 <p>A person riding a motorcycle on a dirt road.</p>	 <p>Two dogs play in the grass.</p>	 <p>A skateboarder does a trick on a ramp.</p>	 <p>A dog is jumping to catch a frisbee.</p>
 <p>A group of young people playing a game of frisbee.</p>	 <p>Two hockey players are fighting over the puck.</p>	 <p>A little girl in a pink hat is blowing bubbles.</p>	 <p>A refrigerator filled with lots of food and drinks.</p>
 <p>A herd of elephants walking across a dry grass field.</p>	 <p>A close up of a cat laying on a couch.</p>	 <p>A red motorcycle parked on the side of the road.</p>	 <p>A yellow school bus parked in a parking lot.</p>

exemple de résultats de l'algorithme de Google

Indexation de document

Objectif

- sélectionner les informations caractéristiques dans le document
- extraire ces informations pour en faire une **représentation**
 - représentation structurée : métadonnées
 - représentation non structurée



Deux types d'indexation

- 1 manuelle
- 2 automatique

Indexation manuelle

Caractéristiques

- Elle permet de fixer des valeurs à la représentation structurée :
 - Ces valeurs peuvent être structurée (utilisation de thésaurus) ou pas (mots libres)
- Elle n'est réalisable que lorsque le nombre de documents est faible
- Elle peut varier en fonction de l'indexeur

Paramètres variables

- personne (connaissance, subjectivité, ...)
- temps
- technique
- objectif
- langue
- ...

Indexation automatique

Caractéristiques

- Elle permet d'indexer de gros volumes de documents
- Elle permet de renseigner la représentation structurée et la représentation non structurée
 - représentations structurée : à partir des métadonnées associées aux documents (facile si ce sont les mêmes, difficile sinon) ou à partir d'une analyse du texte
 - représentation non structurée : à partir d'une analyse du texte

Indexation plein texte

Principe

- Extraire automatiquement des informations caractérisant le texte

Analyse du texte

- segmentation (tokenization) : découpage
- analyse lexicale : association des tokens à des lemmes + propriétés
- analyse syntaxique : détermination de la classe morpho-syntaxique dans le corpus
- ...

Extraire les informations pertinentes

Qu'est ce qu'un mot "pertinent" :

- fréquence importante ?
- mot seul ? co-occurrences ?
- structure logique : résumé, titre, ...
- ...

⇒ pondération des mots

⇒ association des mots de poids important au document

Plan du cours

- Du document papier au document numérique
- Rappels :
 - XML le métalangage de représentation des documents
 - DTD le langage (non XML) de base de description de grammaires XML
- XPath le langage (non XML) d'identification/de désignation d'une partie d'un document XML
- XSLT le langage (XML) de transformation de documents XML
- XSD le langage (XML) de description de grammaires XML

Plan du cours

- La recherche d'information avec des représentations de documents non structurées
- Les métadonnées ou la recherche d'information avec des représentations structurées
- RDF et SPARQL où comment représenter et interroger des métadonnées
- RDFS et OWL où comment définir des schémas de métadonnées
- Prolog où comment raisonner/inférer sur des connaissances

