

01

De la rareté à l'abondance

## Hier:

Moyens de production rares et onéreux (35mm, caméra) Moyens de diffusion rares et centralisés (chaine unique, cinémas)



# Aujourd'hui:

Moyens de production facile et peu onéreux (smartphones) Moyens de diffusion nombreux et décentralisé (internet, chaines , )





# Une Abondance qui pose des enjeux de :

Traçabilité, provenance Mise en contexte, vérification de l'information Importance de la donnée d'accompagnement :





02

Du document à la donnée

### Document vs donnée

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Fusce mollis neque in ante vulputate, quis accumsan dui euismod. Nunc lobortis aliquet orci, ut iaculis nunc feugiat id. Interdum et malesuada fames ac ante ipsum primis in faucibus.

Un **document** est un ensemble logique, fini d'**informations** dont les limites peuvent être définies par des caractéristiques physiques.

Le fait que « Le premier mot du texte soit Lorem » est une donnée.

Le fait que « ce texte soit extrait d'un texte de Cicéron » est une donnée, qu'on appelle métadonnée (donnée sur la donnée)



## Hier:

Prima du document:

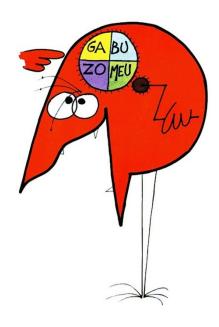


#### On observe:

- → Cloisonnement des collections
- → Des systèmes documentaires tournés vers peu d'usages
- → Opacité des collections

# Aujourd'hui:

Prima de la donnée:



## Objectifs:

- → Décloisonner les collections
- → Libérer de nouveaux usages
- → Transparence des données
- → Résilience des systèmes documentaires



# Exemple: Numériser une collection audiovisuelle

Un processus au coût maîtrisé, industrialisé une équation de la problématique l'importance aussi de la donnée documentaire

Une collection sans matériel d'accompagnement, sans données :

→ moins de valeurs

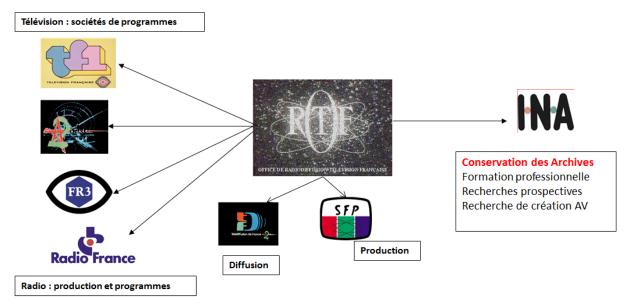
Car le temps de description documentaire nécessite d'énormes ressources.



# Historiquement les collections de l'Ina ce sont deux fonds principaux :

#### 1 / LE FONDS DES ARCHIVES PROFESSIONNELLES

L'éclatement de l'ORTF en 1974 : Archives de l'ORTF + TF1 + Antenne 2 + FR3 + Radio France : Le début des archives Patrimoniales



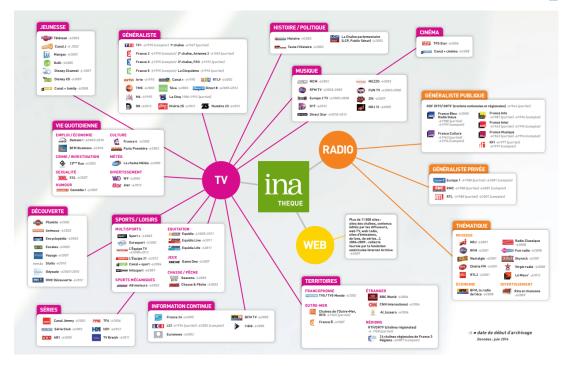
- → 2 000 000 d'heures
- → usage : une vente des archives pour réutilisation dans de nouveau programmes



## 2 / LE FONDS DU DEPOT LEGAL

Création du dépôt légal de l'audiovisuel français en 1992 :

100 chaines tv, 30 chaines radio et 11 000 sites web média captés 24h/24



- → 14 174 000 heures au titre du Dépôt légal
- → usage : consultation par les chercheurs



#### Deux Fonds distincts

- 2 logiques différentes
- 2 logiciels différents

- → peu d'interopérabilité
- → outil documentaire tourné vers le matériel et adapté à un usage précis



POURQUOI FAIRE SIMPLE QUAND ON PEUT FAIRE COMPLIQUÉ?!

(recherche ou vente)

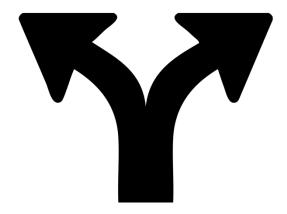


# Les problèmes posés

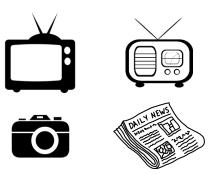




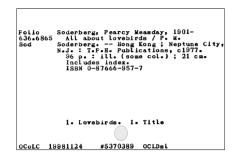
Stock vs flux



Evolution en parallèle



## Différents médias



# Description rigide



# Projet en cours à l'Ina : Le lac de données

Des données + des référentiels + Une interface de recherche suivant le profil d'utilisateur



## Objectifs:

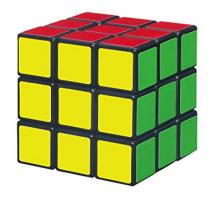
- → Réunir les données des différentes bases
- → Créer des liens entre collections
- → Favoriser de nouveaux usages



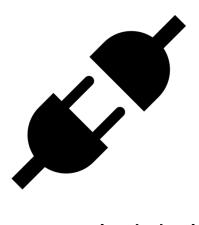
# Les points importants



Souplesse



Cohérence



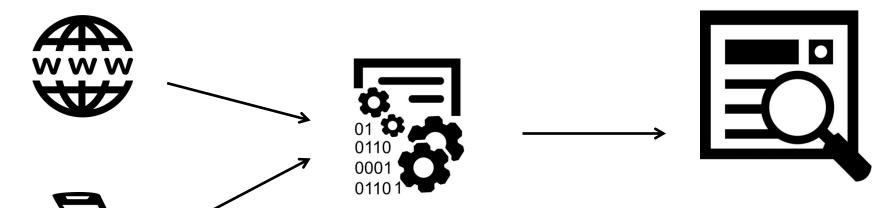
Interopérabilité



REMISE EN CAUSE DU TRAVAIL DE CATALOGAGE ET D'INDEXATION ?



# L'indexation automatique du contenu ou la remise en cause du catalogage



Article dans la revue : *EEE Intelligent System* journal

mars 2009 "the unreasonable effectiveness of data"



Peter Norvig

Responsable de la recherche chez Google



## La folksonomie ou la remise en cause du référentiel





Article sur le blog : *Clay Shirky's Writings About the Internet*Printemps 2005

"Ontology is Overrated: Categories, Links, and Tags"

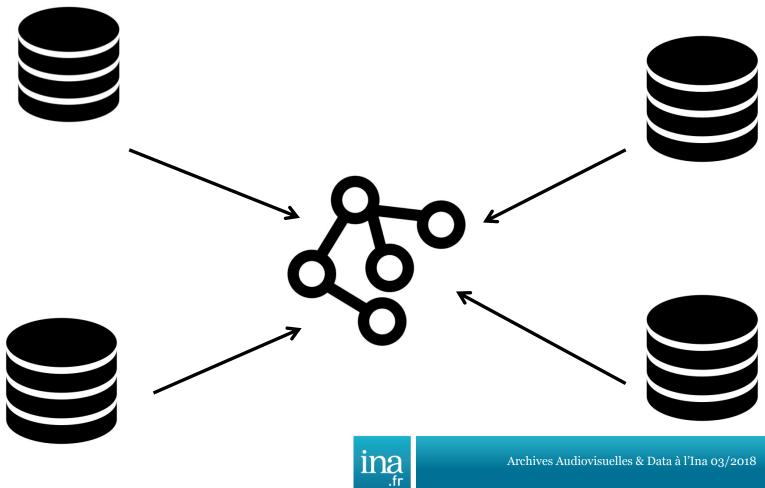


Clay Shirky



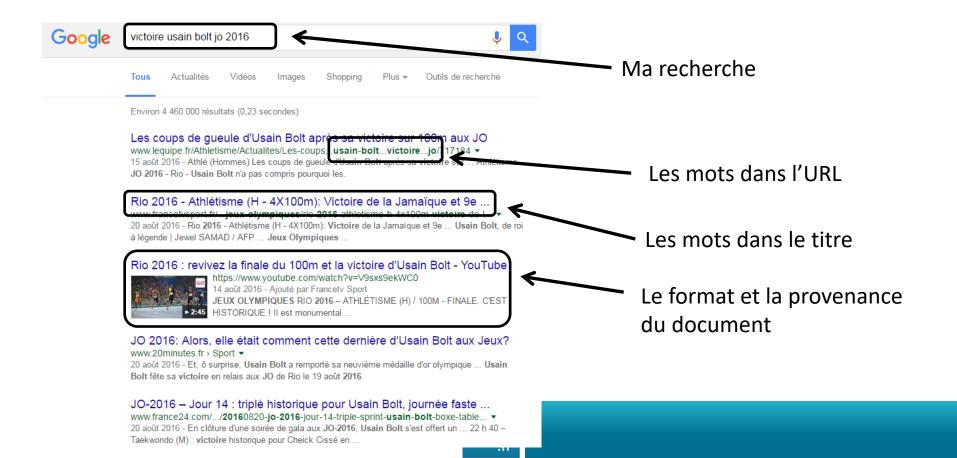
## Mettre en cohérence les données

ou la nécessité de disposer de référentiels structurés



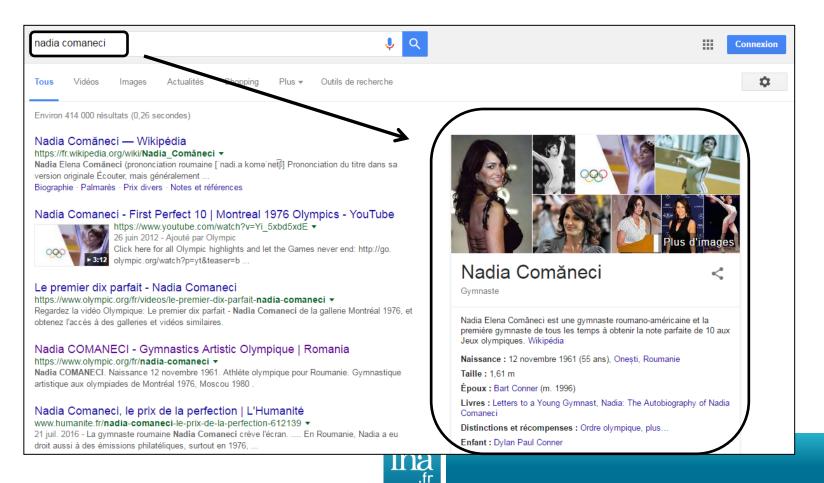
## Hiérarchiser les données

# ou la nécessité de disposer de métadonnées pertinentes



## Organiser les données

# ou la nécessité de disposer de données structurées



03

Typologie des données

# Qu'est-ce-qu'une donnée?



Connaissances assimilées et/ou connaissances partagées

Synthèse ou déduction obtenue à partir de plusieurs informations

ina

Un ensemble organisé de données OU résultat d'un traitement sur un ensemble de données

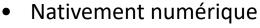
DATA

Des faits, signaux et symboles formant une unité indépendante les unes des autres et non traitée

# Typologie des données



Contenu



 Issu d'un processus de numérisation (OCR, Speech to text, reconnaissance de formes)



Métadonnées

Informations sur le document qui peut-être d'ordre :

- descriptive
- administrative (technique, traçabilité, juridique)
- structure



Référentiels

Informations de référence permettant d'assurer une cohérence entre les données



Données d'usage et contributions

Toutes les données laissées par les utilisateurs : logs, traces, contributions sur les réseaux sociaux en rapport avec les données

1100 1010 0101

Données générées /secondaires

Données obtenues à partir de l'exploitation des autres types de données : données quantitatives, extraction, machine learning, inférence...



Archives Audiovisuelles & Data à l'Ina 03/2018

#### De la notice à la donnée

### Anatomie d'un changement de granularité

#### Seule la description est numérique



**Notice** 

La description d'un document forme lui-même un document, un tout cohérent et fini



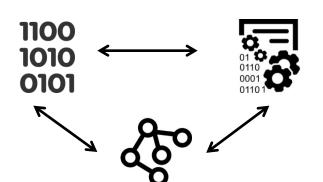




Métadonnée

Chaque information est indépendante Elargissement de la description (technique, juridique, structure...)

Contenu nativement numérique ou numérisé



Donnée

Fin de la hiérarchie entre les données (description, contenu, référentiels) Diversification des données

(données d'usages / contributions)



# L'importance des données d'usage

#### Intérêt dans la recherche de documents :

Le nombre de vue -> une information importante pour le consommateur de vidéo





Les crêpes (avec beaucoup d'alcool) de Raymond Oliver | Archive INA

147 043 vues















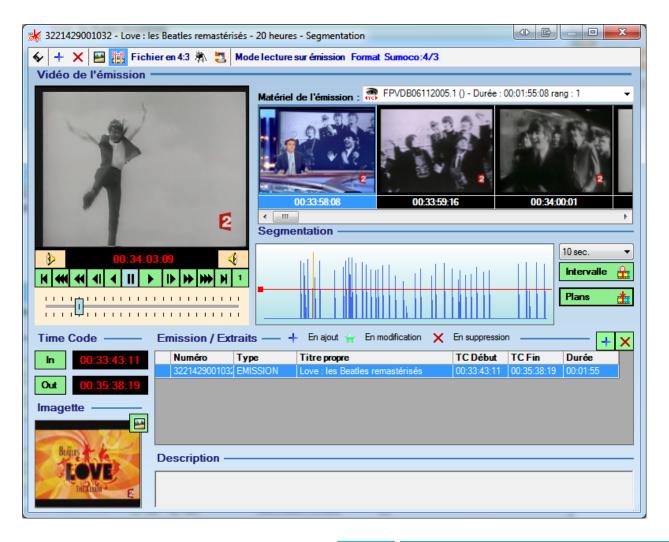
# Des modèles économiques basés sur l'utilisation de ces données d'usage :

Algorithme + données usages = Bizness juteux





# Les données de contenu (signifiant / signifié)





Les données d'accompagnement

Les métadonnées

Les référentiels (thésaurus, listes contrôlées,...)

Les données générées automatiquement

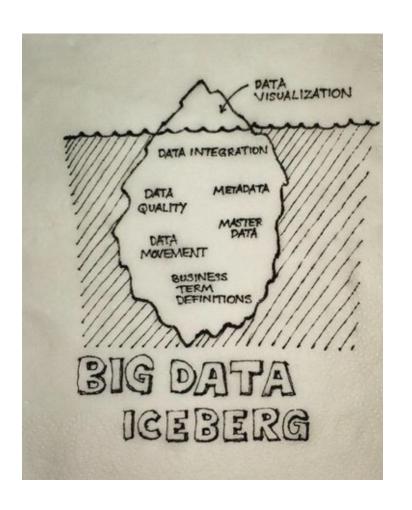


# Notion de Big Data

## Explosion quantitative de la donnée numérique







## Essai de définition

Le Big data désigne la capacité technologique à traiter de très grandes masses de données avec des infrastructures matérielles standards.

Par extension, on désigne avec cet expression un mouvement général se caractérisant par un intérêt particulier du monde numérique pour la donnée.





# Objectif du département Recherche et Innovation de l'Ina

→ Nouveau outils pour les usages existants

Fournir des outils

Assister, préparer le travail de documentation, d'annotation, d'indexation, de description

Proposer des annotations, des index, des repères

→ Proposer de nouveaux usages / services

Recherche par similarité

Réalité augmenté



# Le Département Recherche et Innovation de l'Ina



# Objectif des technologies 1/2

#### → Détecter

Image : présence de visages, de textes, de logos, d'objets divers

Audio : présence de parole, de musique, d'applaudissements,

de sons divers

#### → Reconnaitre

Image : le visage de qui ? Qu'est il décrit ? Quel tableau ?

Audio: qui parle? Qu'est-il dit? Quelle chanson?

Image + audio : qui?



## Objectif des technologies 2/2

#### → Organiser

Image: regrouper les visages

Audio: musiques par genre

#### → Structurer temporellement

Découpage des Jt en « Plateaux » et « Reportages »

Chapitrage

Détection des sections d'un morceau de musique (couplet -refrain)



# Démonstrations sur l'Espace Recherche http://recherche.ina.fr





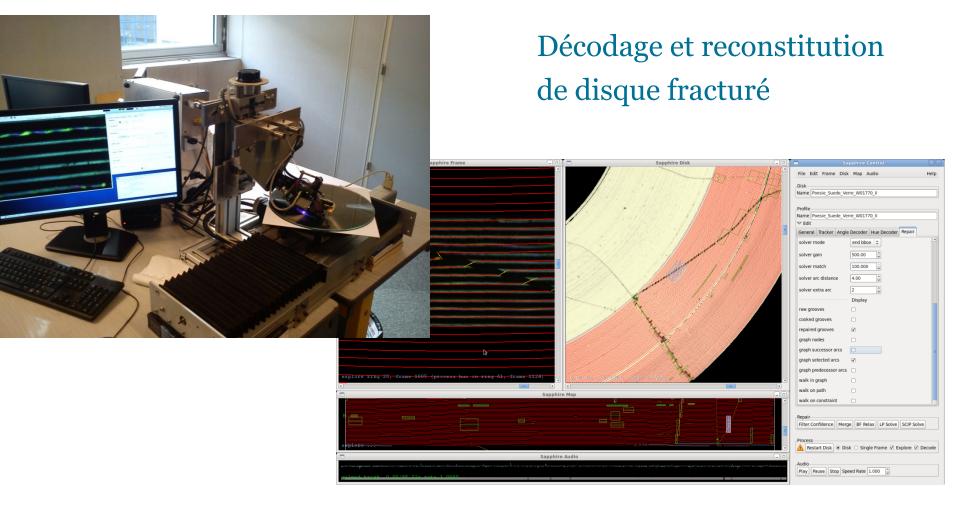
Exemple du projet SAPHIR, à l'Ina : récupération de données sur supports défectueux

20 000 disques concernant la mémoire radio de 1930 à 1960 Lecture avec contact impossible





## Lecture de l'angle du sillon par la couleur réfléchie et



## SIGNATURE : signature de vidéos Empreinte numérique, fingerprinting



#### **DIGINPIX**

## Reconnaissance de peintures dans des archives vidéo

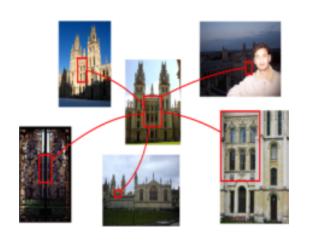








## Recherche rapide de parties d'images (bâtiments, logos...)

































## Détection de « concepts » visuels

#### Reconnaître Londres, Prague et Paris

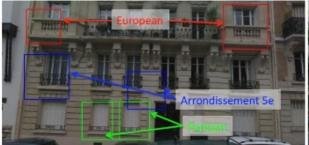


Source: What makes Paris look like Paris? [Doersch et al., 2012]

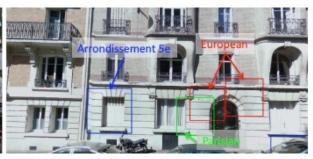


Prague, Czech Republic

London, England







### DESIGN PRINT Automatisation de la segmentation



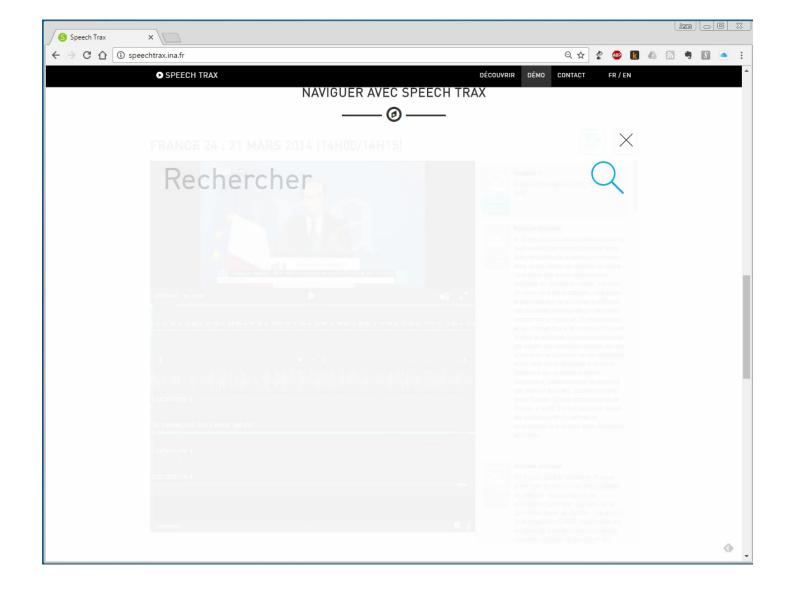
La segmentation des fichiers numériques (pose de TC in/out) est une opération essentiellement manuelle, laborieuse et chronophage. L'Ina a mis au point un système de segmentation automatique qu'il a testé sur les journaux télévisés de 1989. Le projet a consisté à imaginer une interface qui permette de valider, corriger et capitaliser ces repérages temporels automatiques.



### SPEECH TRAX Identification de voix célèbres dans des archives vidéo

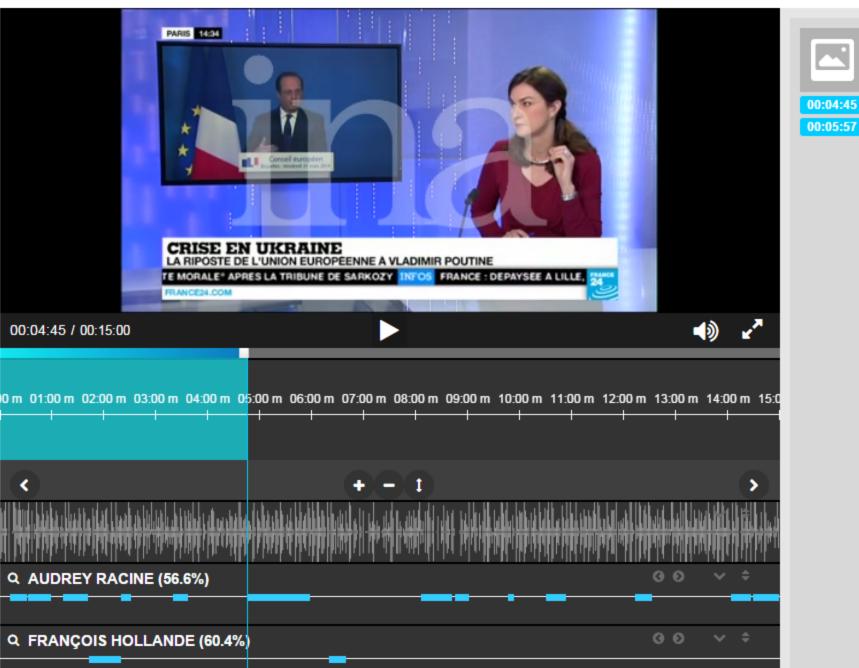








## FRANCE 24: 21 MARS 2014 (14H30/14H45)



audrey raci la peur juste parlera dans focus mais e annoncé la f avec la Russ toutes les m déclaration a français de l actuellemen et au bord d a également France était avions de ch renforcer sa baltes des s pas vraimen président VI la dernière r Crimée à la d'application plus tôt sans Parlement ru

> l'unanimité le les faits, tou

mort. Dans I échanges de l'opposition e

socialistes a

#### **UTILISATION:**

→ Pour les documentalistes :

Proposer une première documentation brute et imparfaite des contenus indexés par l'Ina pour amorcer le travail documentaire

→ Pour les clients et usagers de l'Ina :

Proposer des parcours d'archives obliques reposant sur les interventions orales de personnalités du paysage audiovisuel français



## DATASET Ouvrir les données Ina pour la Recherche





→ Mise à disposition de la communauté scientifique et technologique d'un corpus de documents audiovisuels issus de ses collections, de fiches documentaires et de métadonnées associées à ces documents.

Ce corpus est destiné à la mise au point, l'expérimentation et l'évaluation d'outils de recherche et d'analyse de contenus multimédias dans un strict cadre de recherche scientifique.

http://dataset.ina.fr





#### La mutation du métier de documentaliste audiovisuel

→ L'éditorialisation prends une place majeure

Le documentaliste n'est plus qu'un médiateur, il enrichit le document brut par des données.

- exemple de data + archives (tour de France) :





#### La mutation du métier de documentaliste audiovisuel

→ Automatisation de certaines taches chronophages

Indexation automatique : lancer des calculs par lot

- → Bâtiments
- → Personnalités
- → Logos
- → Segmentation
- → Lien donnée







#### Merci de votre attention

Délégation Ina Nord (Hauts de France / Normandie) 8 Allée de la Filature 59000 Lille

Clément Mouly Responsable Documentaire cmouly@ina.fr

