Outils de traitements de corpus textuels développées à Paris-Est : présentations, démonstrations, formations







CorText Manager

Application Web collaborative d'analyse et de cartographie de données hétérogènes

http://manager.cortext.net/















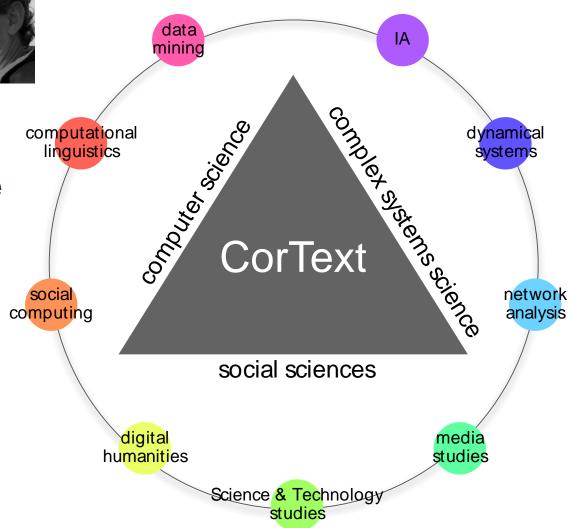


- Breucker Philippe
- Cointet Jean-Philippe
- Duloquin Clhoé
- Duong Tam-Kien
- Laurens Patricia
- Martinez Cristian
- Mazières Antoine
- Mogoutov Andreï
- Schoen Antoine
- Turenne Nicolas
- Villard Lionel

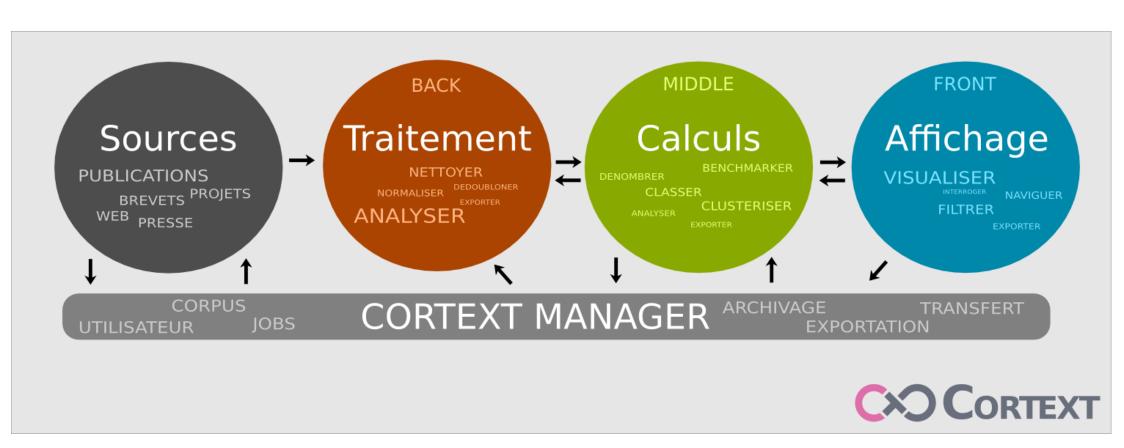




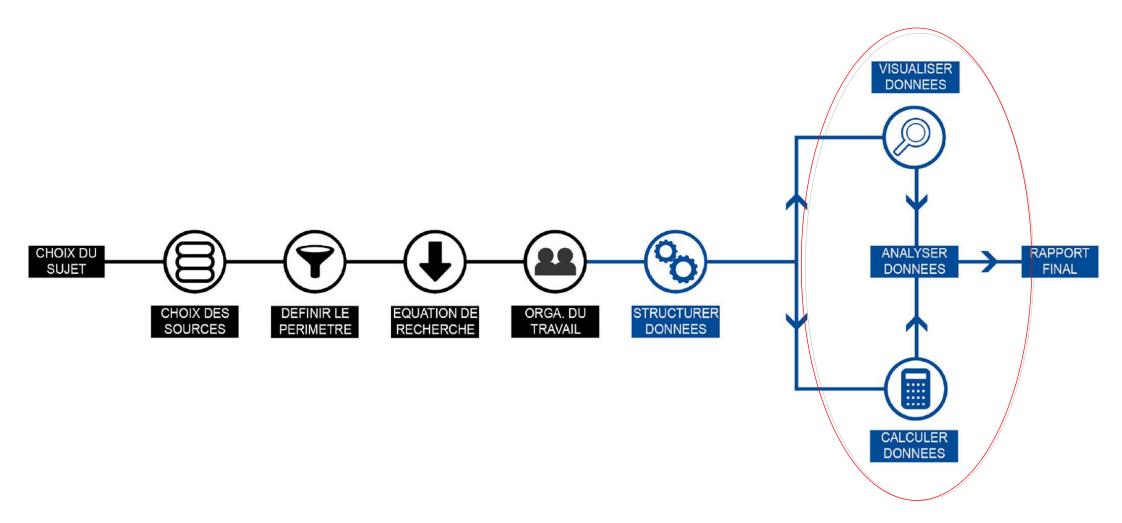




CorText Manager : principales étapes



Du sujet de l'étude aux résultats et visualisations



Des aller-retours entre données – visualisations – interprétations

scientific productions



Web Of Science IS



Microsoft Academic Search



Medline Pubmed

specific databases



rare disease database



projects database



clinical trials database

media productions (press+web)



web crawler



Factiva, press articles archive



online forums

Capacity to collect, parse and handle corpora from various arenas

Intervenir sur les échelles pour dimensionner le sujet :

- L'échelle THEMATIQUE (les sous-disciplines et champs technologiques concernés)
- L'échelle TEMPORELLE (la période)
- L'échelle GEOGRAPHIQUE (par ex. la ville, la région, le pays, le monde)

Espaces stratégiques analysés:

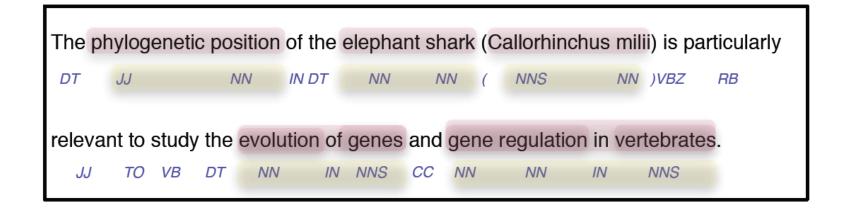
- L'espace INSTITUTIONNEL (par ex. les universités, entreprises et laboratoires publics centraux)
- La GEOGRAPHIE (par ex. les villes, les régions et les pays centraux)
- Les dynamiques SOCIO-SEMANTIQUES (association des cooccurrences de mots et des acteurs) et THEMATIQUE (disciplines, technologies...)

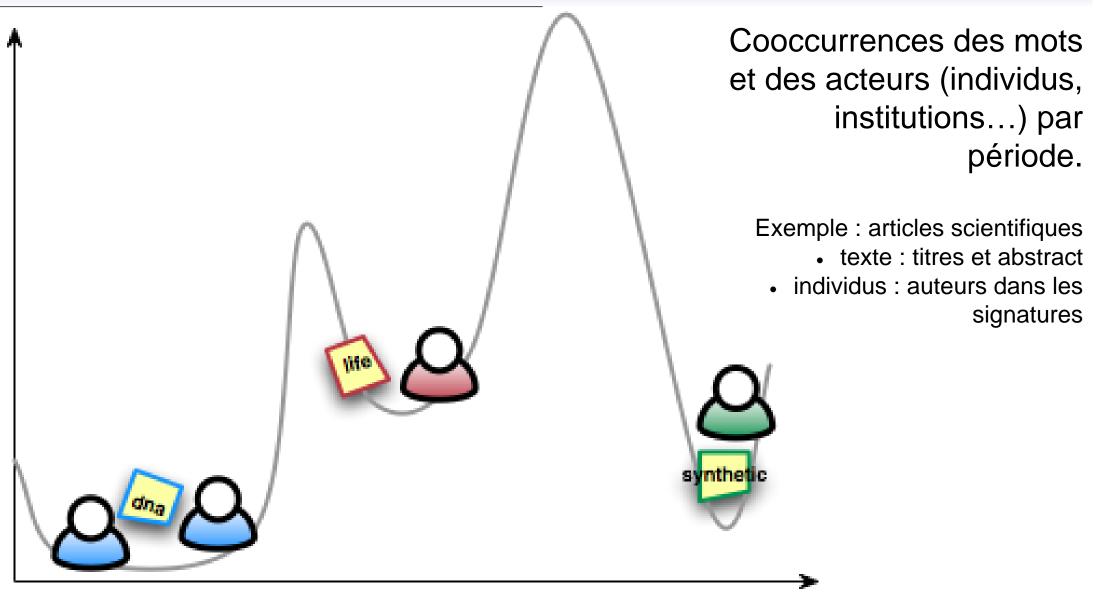
*Le choix des échelles permet notamment d'accroître la pertinence et de réduire le volume de données à extraire.

Construire le paysage socio-sémantique

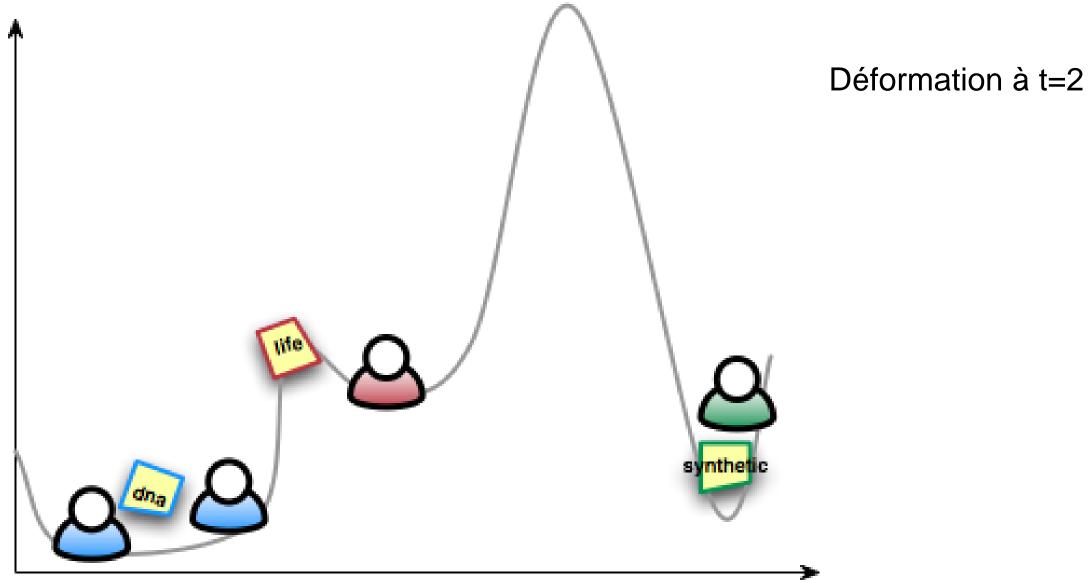
Le paysage socio-sémantique est reconstruit notamment par une EXTRACTION TERMINOLOGIQUE qui s'effectue en plusieurs étapes :

- Étiquetage morpho-syntaxique : associer aux mots des textes disponibles dans les corpus les informations grammaticales (le genre, le nombre...)
- Extraction des groupes nominaux (noms et adjectifs...)

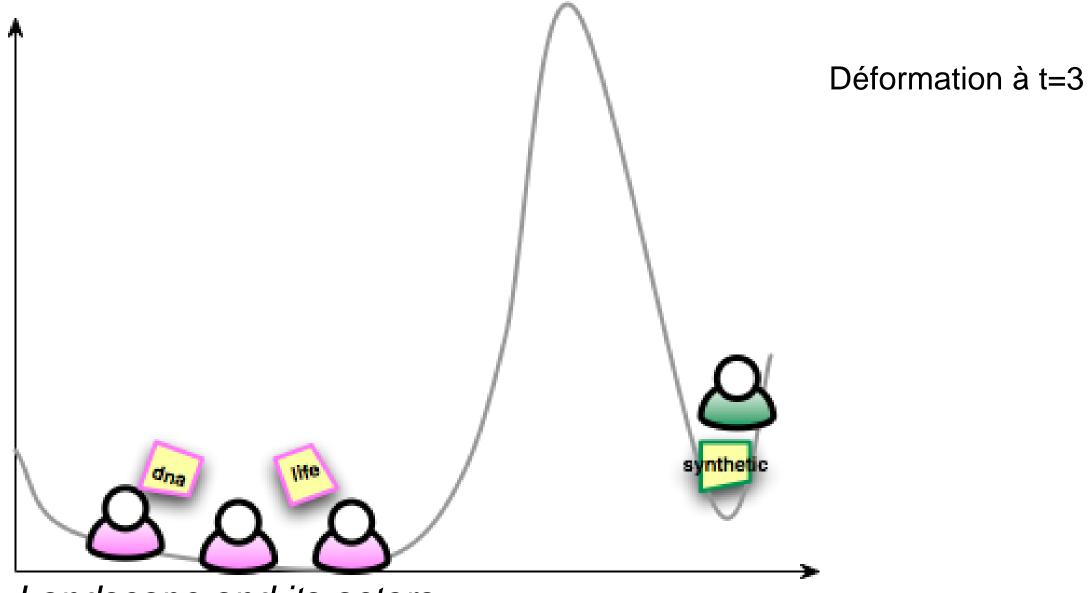




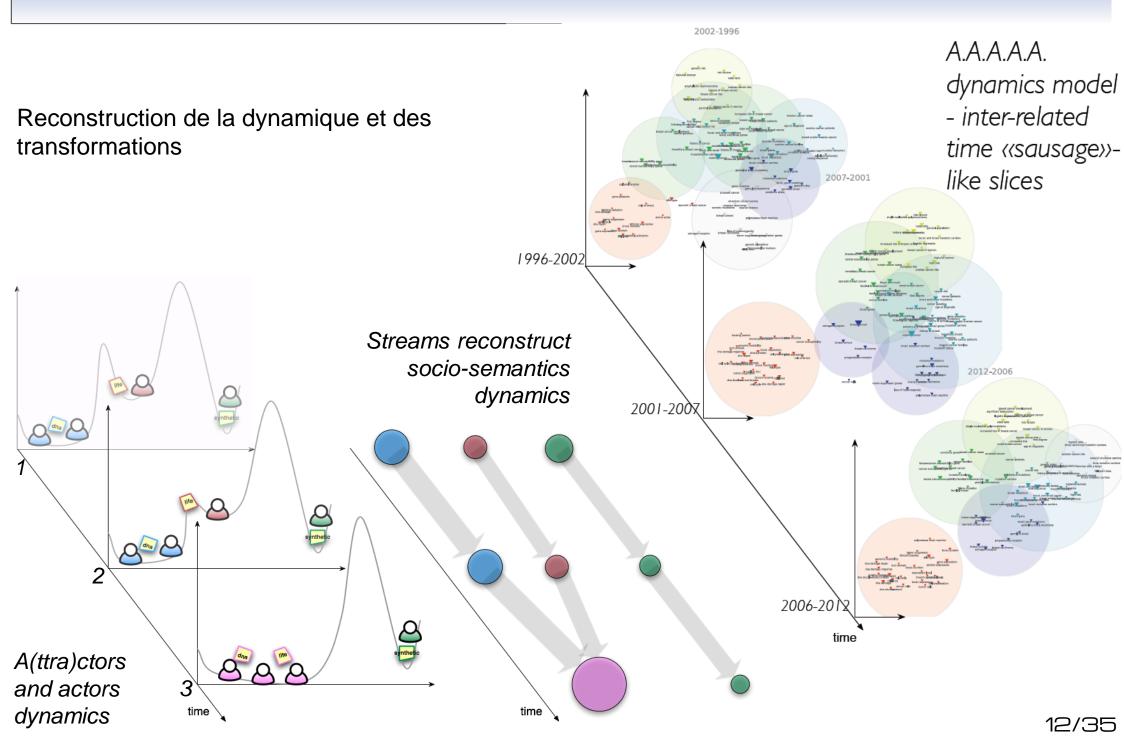
Landscape and its actors (t=1)



Landscape and its actors (t=2)

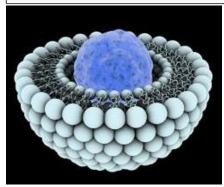


Landscape and its actors (t=3)

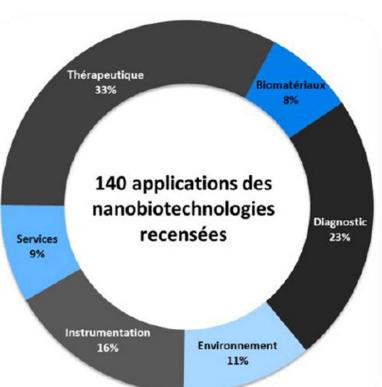


Exemples avec un corpus scientifique concernant les nanobiotechnologies en France 2004-2013

Introduction / Données tabulaires Publications / Réseaux publications



Un liposome : les billes blanches sont les lipides phosphorescants. Ils sont liés entre eux par du cholestérol (petits fils). Le corps bleu au centre de la vésicule est le principe actif médicamenteux (doxuribicine est un produit toxique contre les tumeurs hépatiques).





Bio-puces (séquençage ADN), 2008, Philippe Houdy, La révolution des nanotechnologies, futura-sciences

Figure 8 : répartition des applications en nanobiotechnologies par sous-secteur industriel. Analyse menée à partir d'une centaine de start-up & industriels présents en France.

De la recherche fondamentale aux applications

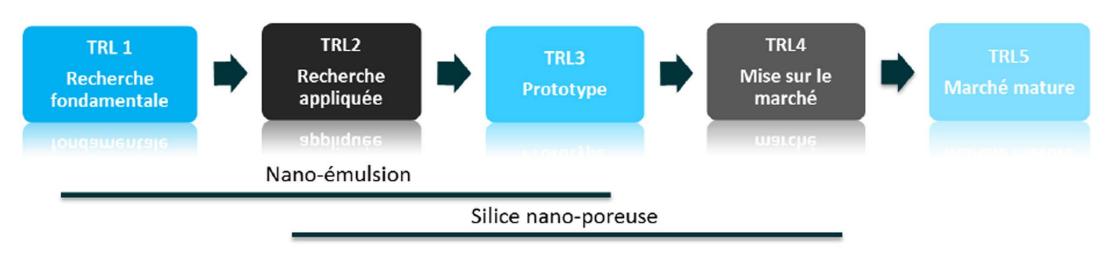
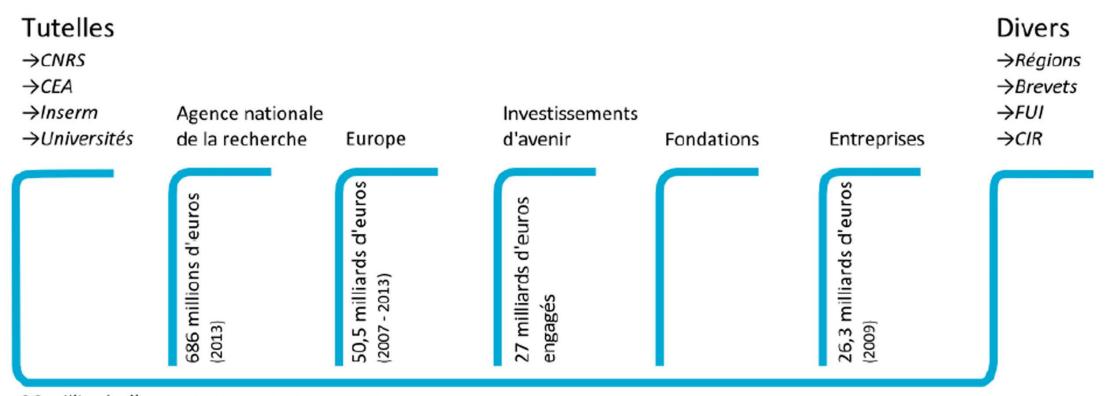


Figure 11 : niveaux de maturité pour différentes technologies.

NanoThinking, Les nanotechnologies en France, 2013

Financement du développement des nanobiotech en France



26 milliards d'euros Budget total public en 2013

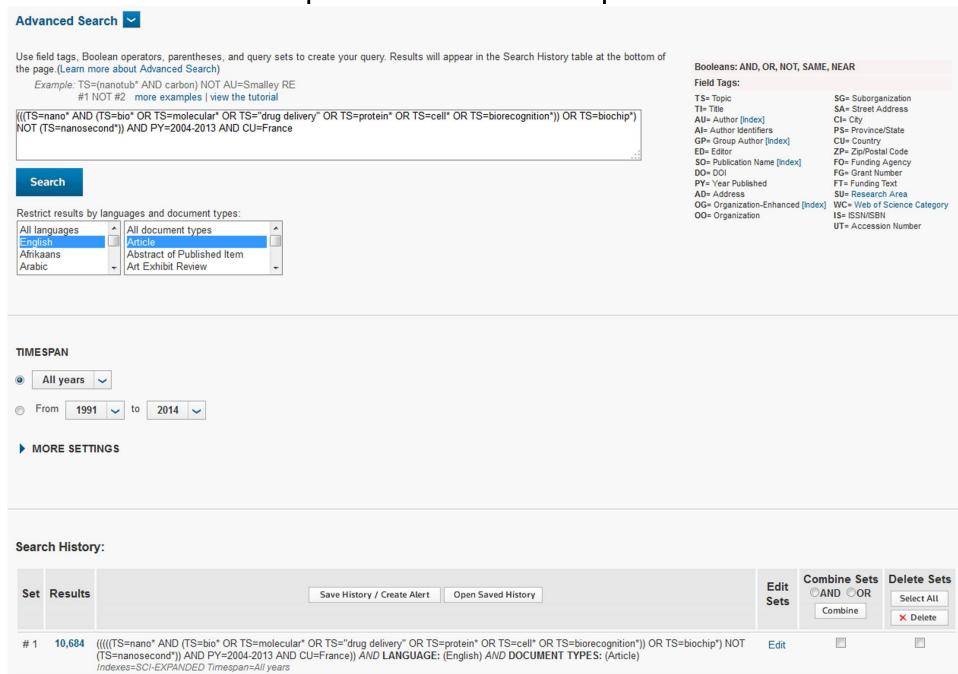
Figure 7 : mécanismes de financement de la recherche en France (source MESR).

NanoThinking, Les nanotechnologies en France, 2013

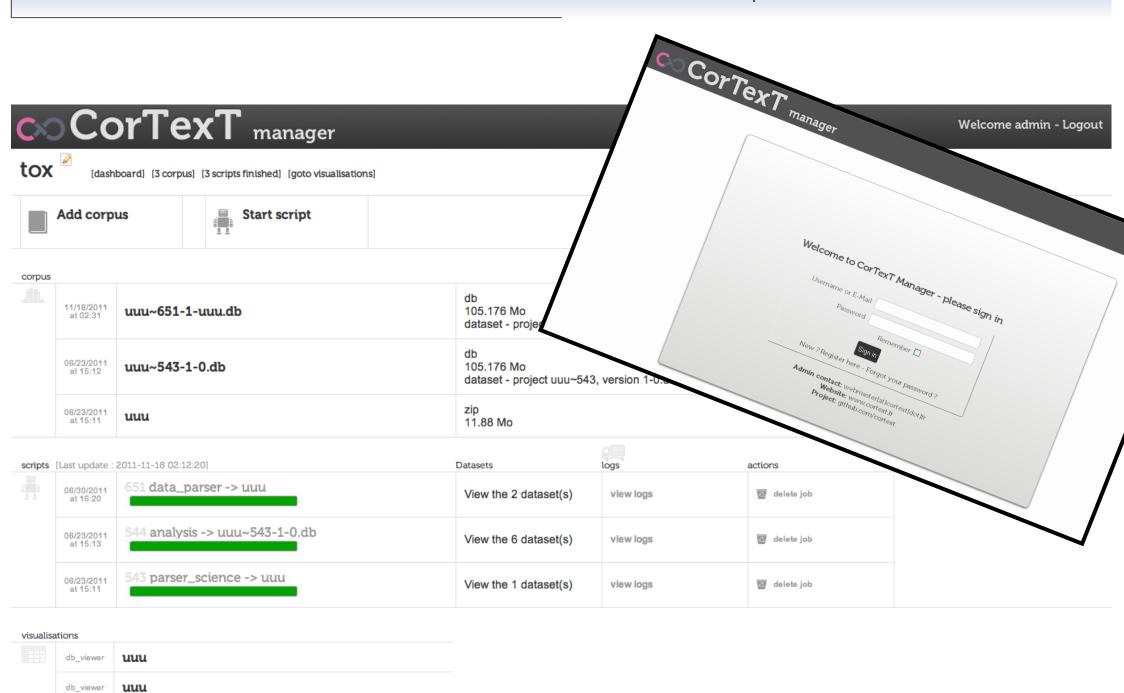
La nanobiotechnologie est donc récente et encore très orientée « recherche »!

Introduction / Données tabulaires Publications / Réseaux publications

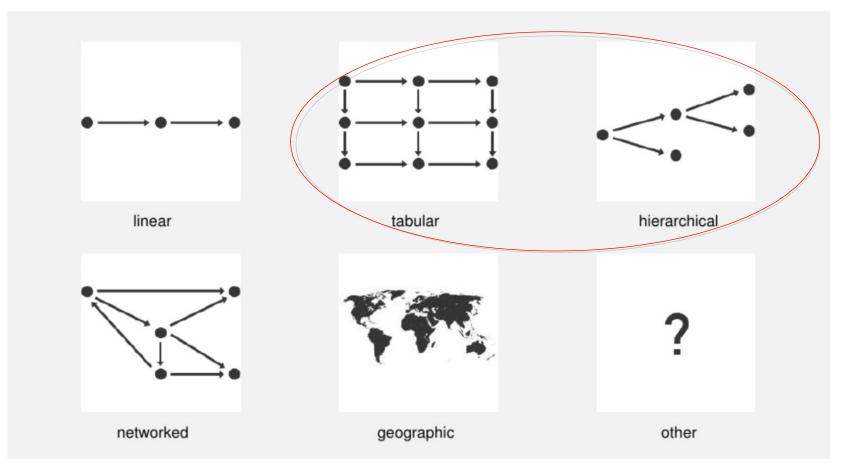
Les nanobiotech dans les publications scientifiques



CorText Manager / Paysage socio-sémantique / Exemples NanoBioTechnologies Introduction / Données tabulaires Publications / Réseaux publications



Examine the Data



Trisnadi Kurniawan, *Infographics and Data Visualisation*, http://fr.slideshare.net/trisnadi/infographics-data-visualisation, 2009

Les publications

CorText Manager / Paysage socio-sémantique / Exemples NanoBioTechnologies Introduction / Données tabulaires Publications / Réseaux publications

Analyse structurelle : statistiques descriptives simples

- caractérisations d'ensemble;
- indicateurs d'activité;
- indicateurs de spécialisation.

Ces statistiques simples s'appuient sur des données au **format tabulaire** (tableaux sql et csv simples). Elles peuvent **exprimer une hiérarchie** (ex : pays / villes / discipline / Nbre de publications).