



BIG DATA et DONNÉES SEO

Vincent Heuschling
vhe@affini-tech.com
@vhe74

Agenda

- **Affini-Tech**
- **SEO ?**
- **Application**
- **Généralisation**



Société

3 Piliers

Méthodes projets
Outils de reporting
& Data-visualisation

**Business
&
Analyses**

BigData
Hadoop
NoSQL
Cloud

Technos

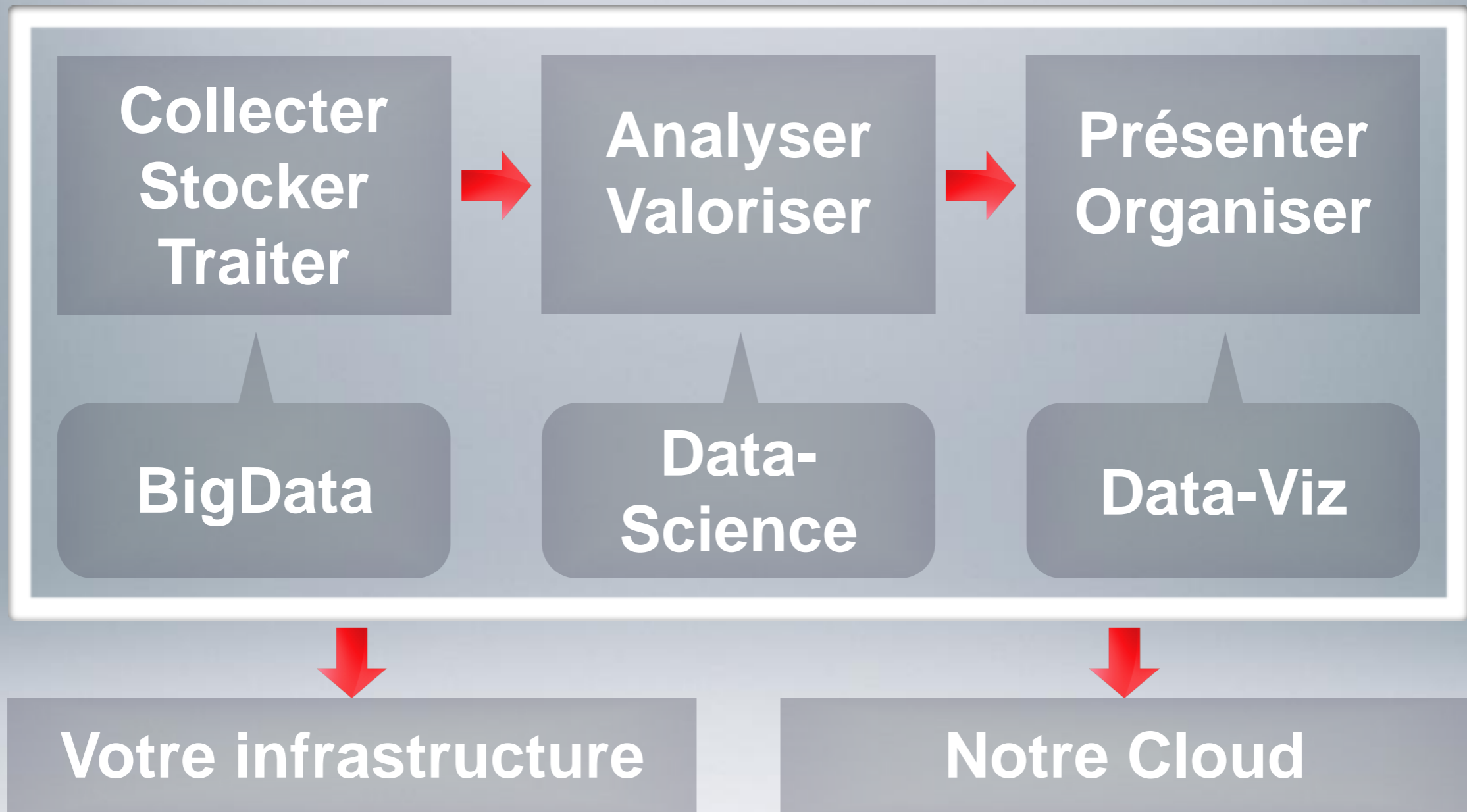
Sciences

Modélisation
Statistiques (R)
Machine Learning

Intégration, Mise en Oeuvre, Conseil et Formation

Une démarche intégrée de bout en bout





Partenaires sectoriels



**Mktg &
Ventes**

Finance

Production

Stats

Apps

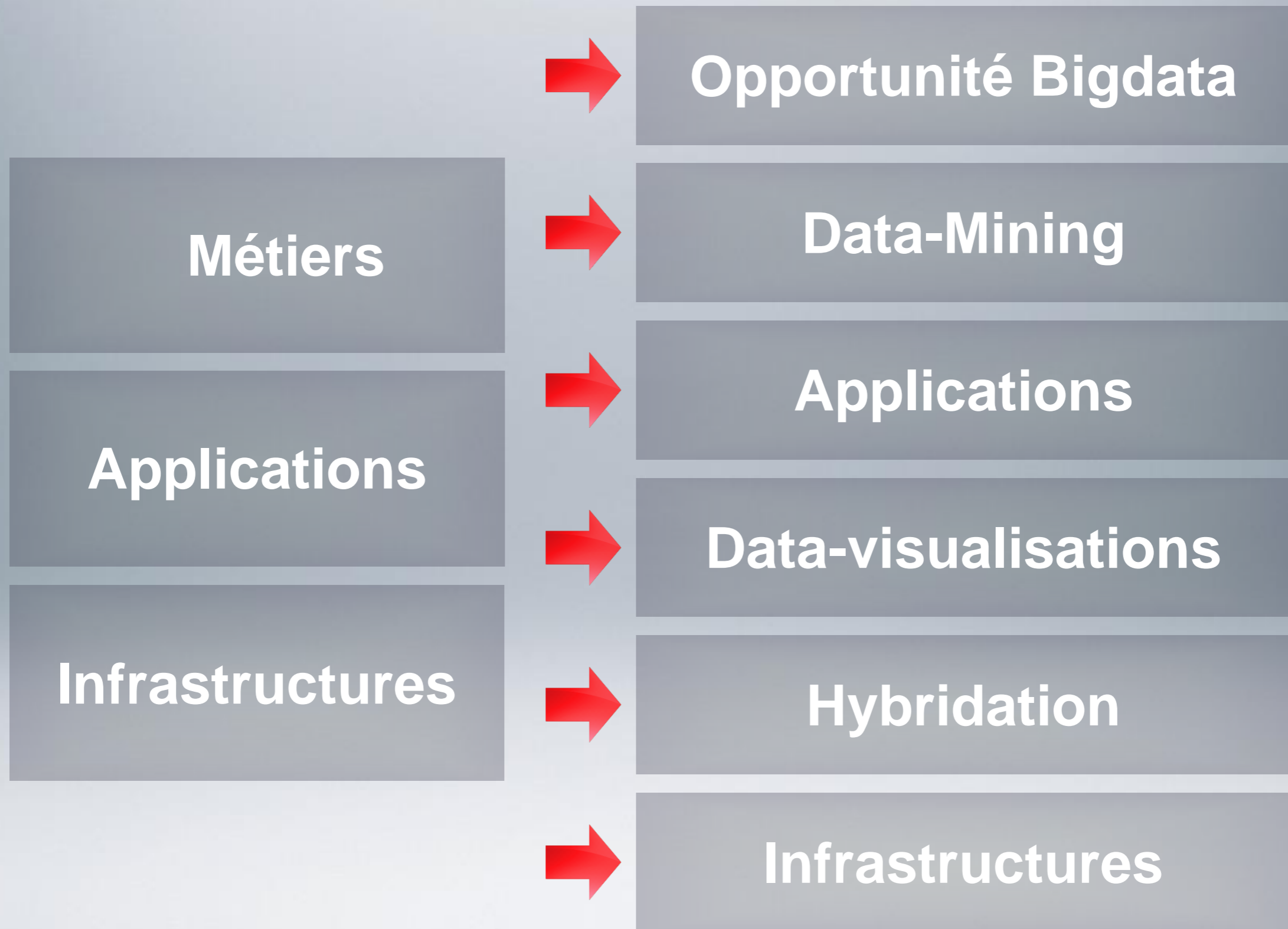
Data-Viz

Infrastructures

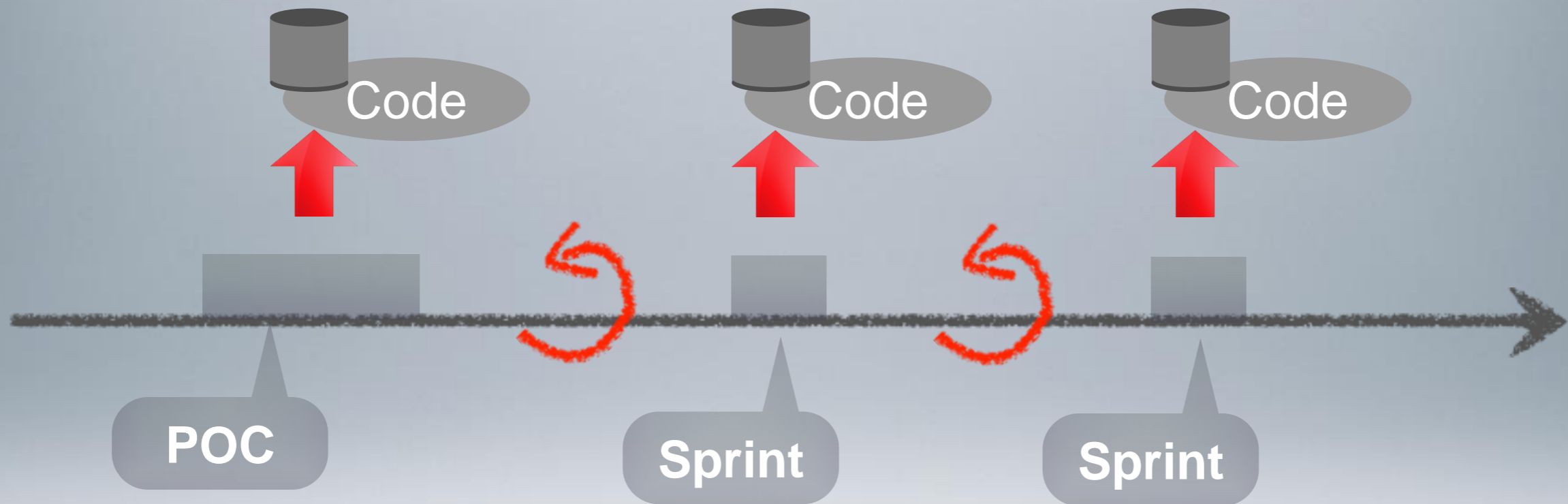


Partenaires technologiques





Agile Data



COLLECTER | STOCKER | ANALYSER | PARTAGER





SEO ?

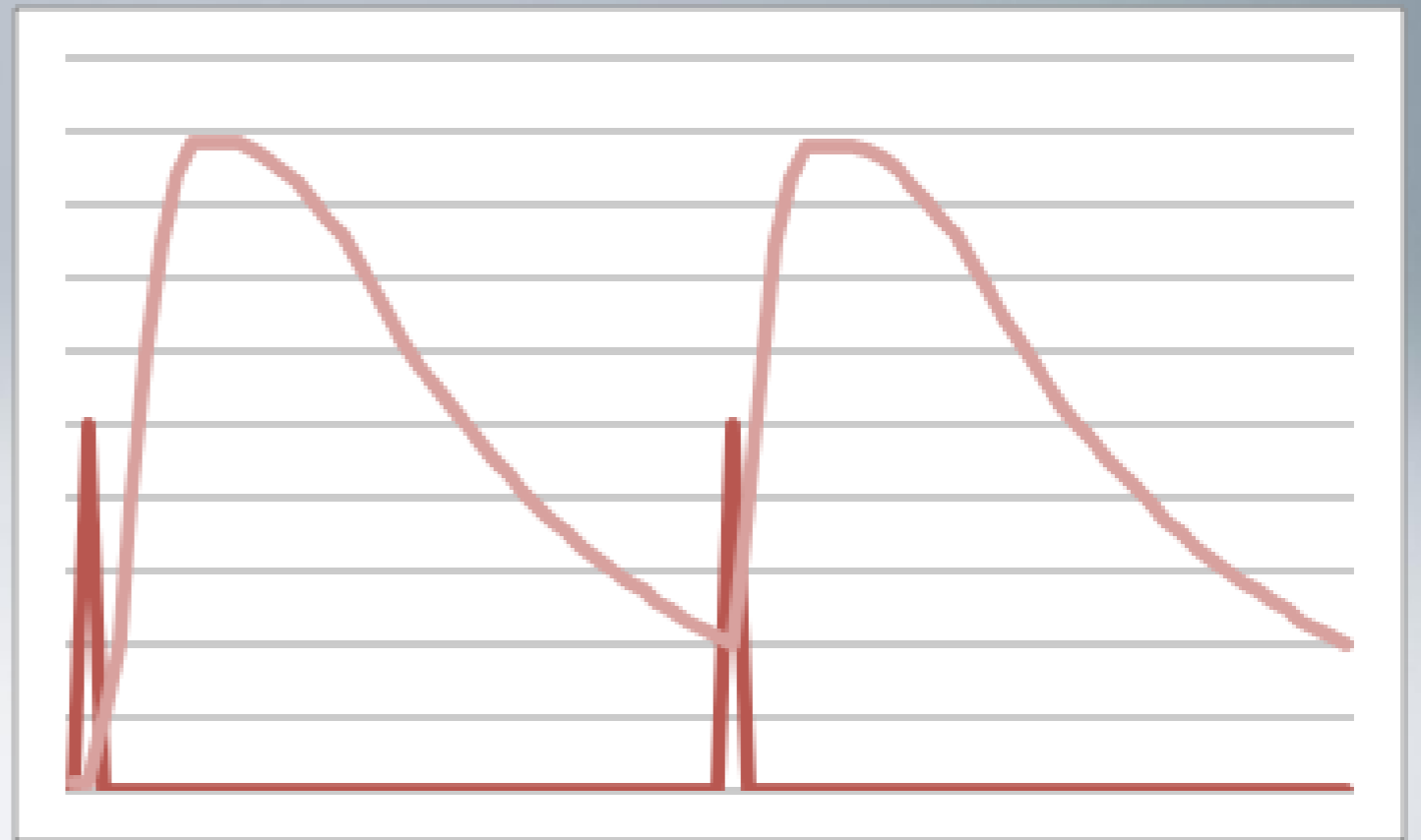
Obtenir les meilleures positions dans la page de réponse de Google.

- ▶ Définir quels éléments du site sont à forte valeur
- ▶ Les promouvoir vers les moteurs de recherches (linking, etc...)
- ▶ Mesurer et étudier le positionnement du site sur des recherches vis à vis de sa concurrence

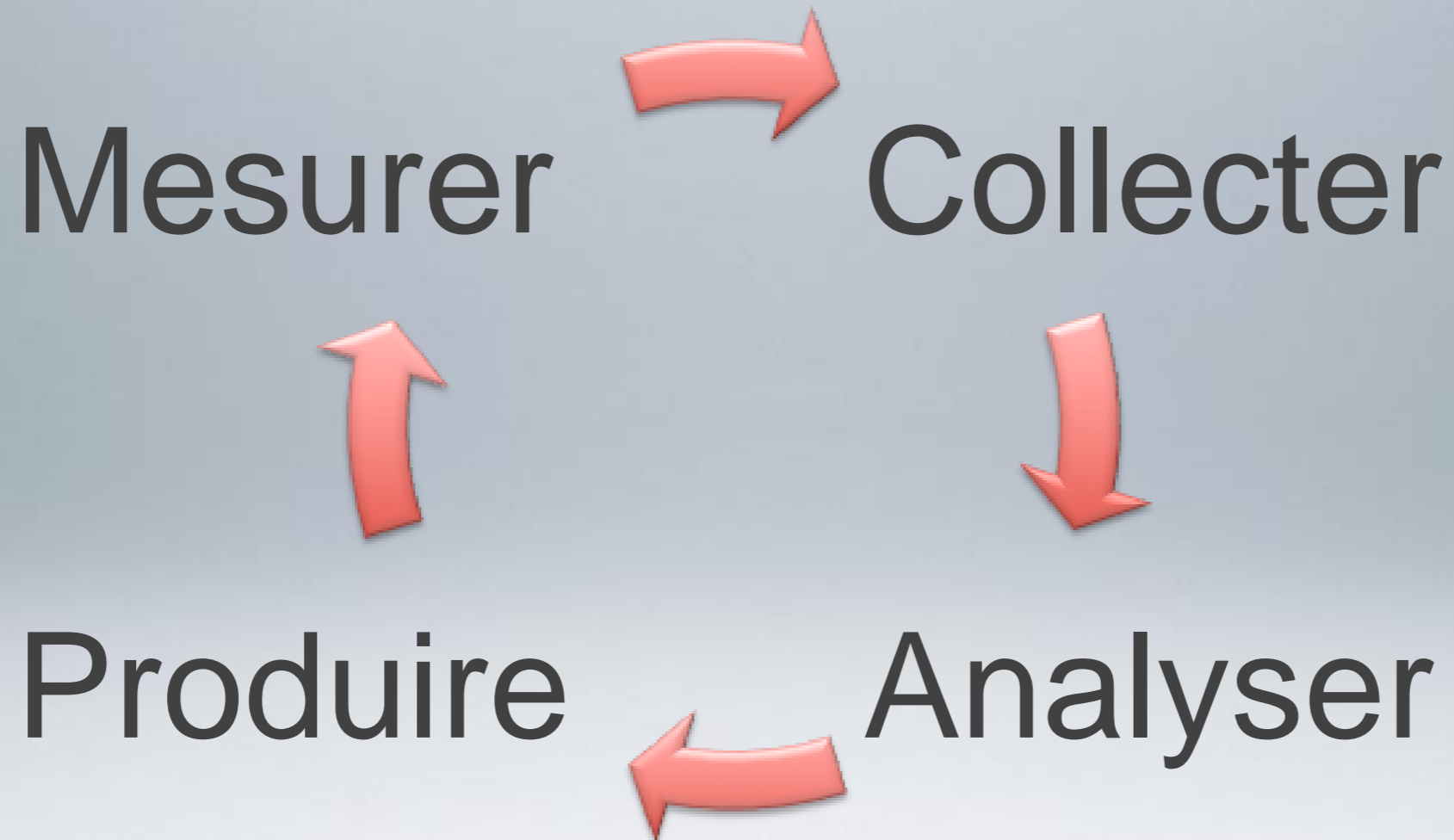


CRAWL et VISITES

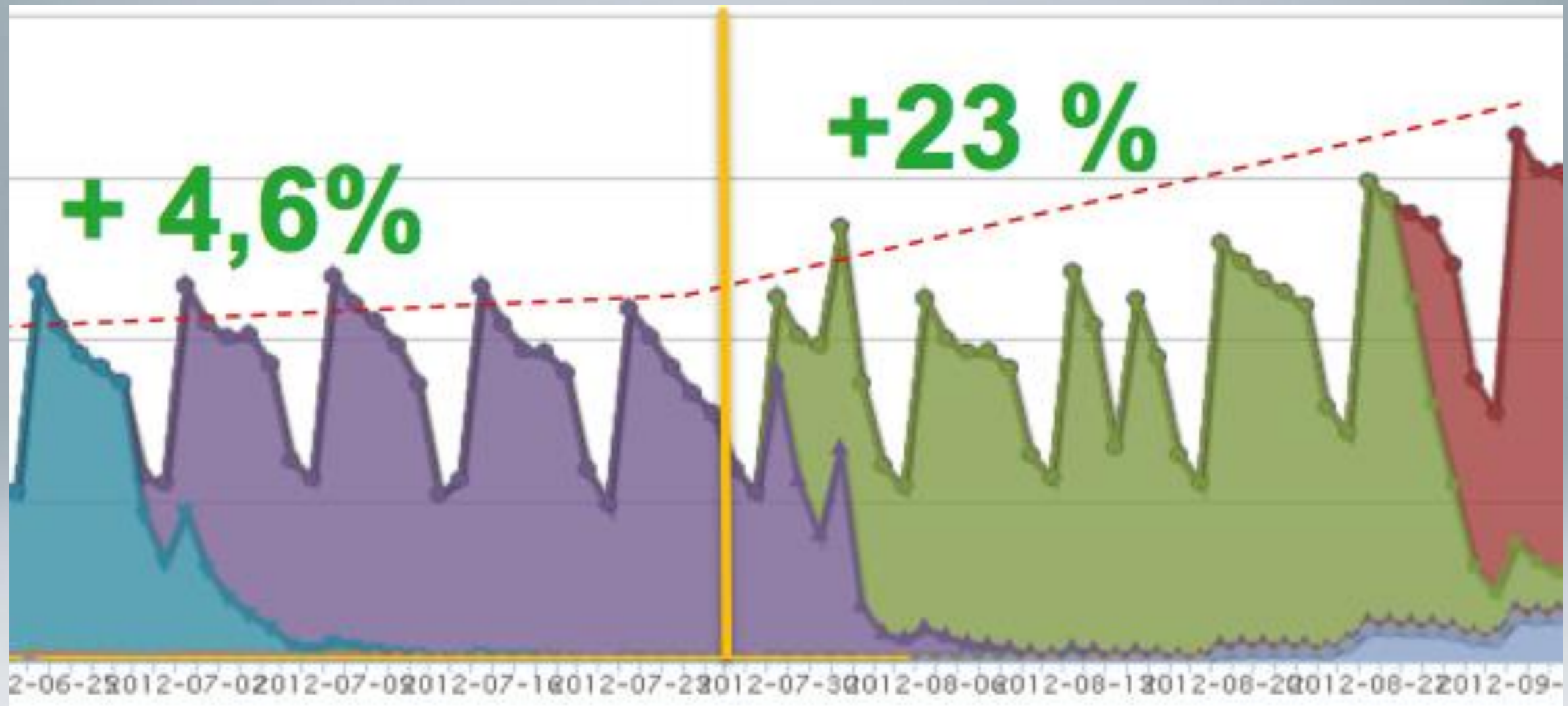
- ▶ Organiser le contenu des pages (Pagerank)
- ▶ Faire Crawler les pages par Google
- ▶ Augmentation directe du trafic



Cercle VERTUEUX de la DATA



RésUltats





Application

ANNUAIRE

- **2000 Professions**
- **40000 Communes**
- **100 M de requêtes par mois**
- **Small data : SEO = env 100 Go /an**

Combien ?

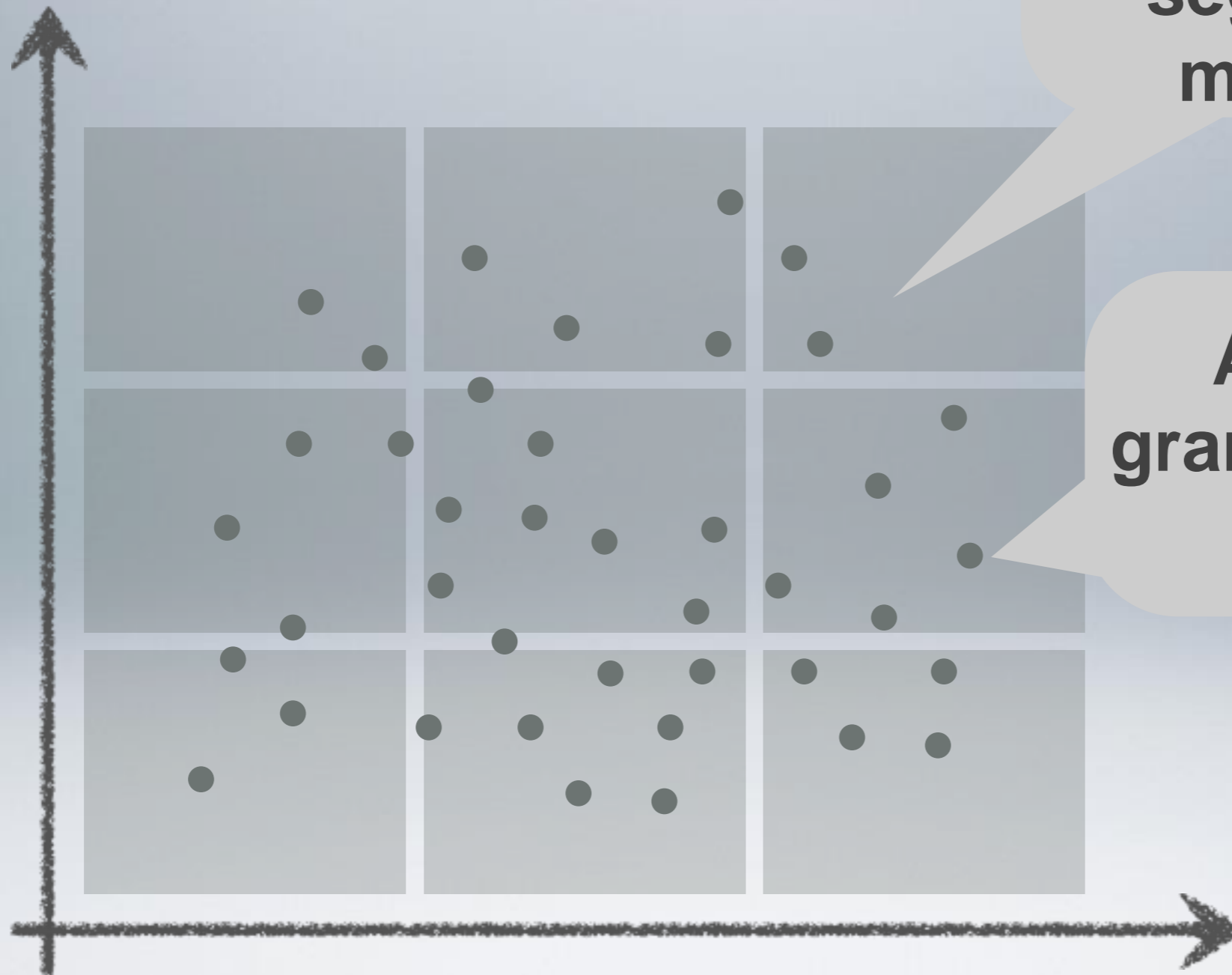
- ▶ 10 visites SEO (hors marque)
- ▶ 30 visites SEO (marque)
- ▶ 90 visites non SEO
- ▶ 20 crawl
- ▶ x7 à x10 au total (pages + ressources)
- ▶ Nécessité de filtrer à la source



TROUVER 400K NOUVELLES URLS À PROMOUVOIR PARMIS 84M ?

- ▶ Similarités et Classifications
- ▶ Recommandation & intelligence collective
- ▶ OpenData

Professions



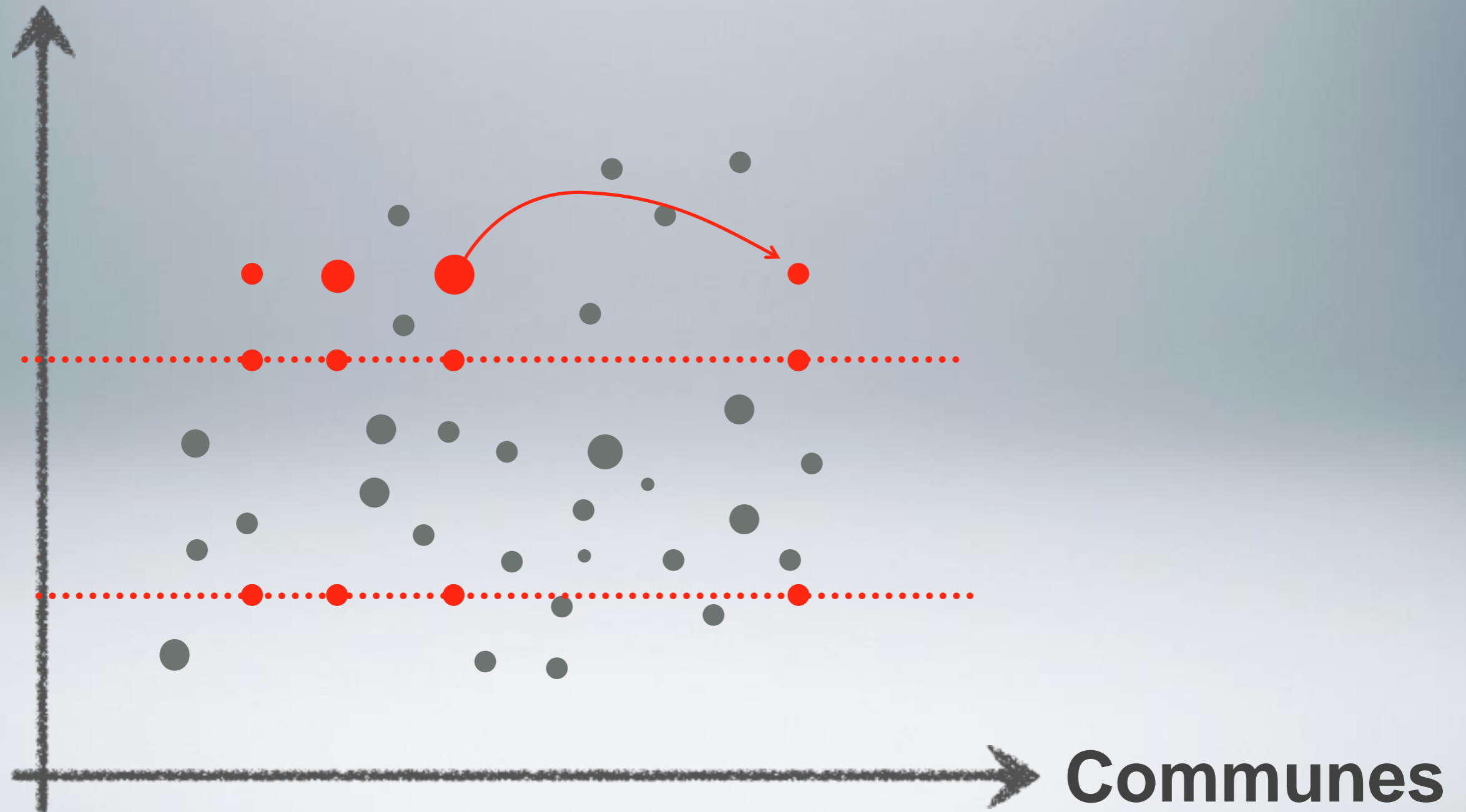
Grandes catégories, segments et moyennes

Analyse à la granularité la plus fine

Communes

RECOMMENDATIONS

Professions



AUGMENTER LA DONNÉE

Data

+ Insee

+ Opendata

a	b	c
~~~	~~~	~~~
~~~	~~~	~~~
~~~	~~~	~~~
~~~	~~~	~~~
~~~	~~~	~~~
~~~	~~~	~~~
~~~	~~~	~~~



a	b	c	m	n	n
~~~	~~~	~~~	~	~	~
~~~	~~~	~~~	~	~	~
~~~	~~~	~~~	~	~	~
~~~	~~~	~~~	~	~	~
~~~	~~~	~~~	~	~	~
~~~	~~~	~~~	~	~	~
~~~	~~~	~~~	~	~	~



a	b	c	m	n	n	x	y	z
~~~	~~~	~~~	~	~	~	~	~	~
~~~	~~~	~~~	~	~	~	~	~	~
~~~	~~~	~~~	~	~	~	~	~	~
~~~	~~~	~~~	~	~	~	~	~	~
~~~	~~~	~~~	~	~	~	~	~	~
~~~	~~~	~~~	~	~	~	~	~	~
~~~	~~~	~~~	~	~	~	~	~	~

- ▶ Une information plus riche
- ▶ Un ciblage plus pertinent

# OUTILS

Collector

Traiter

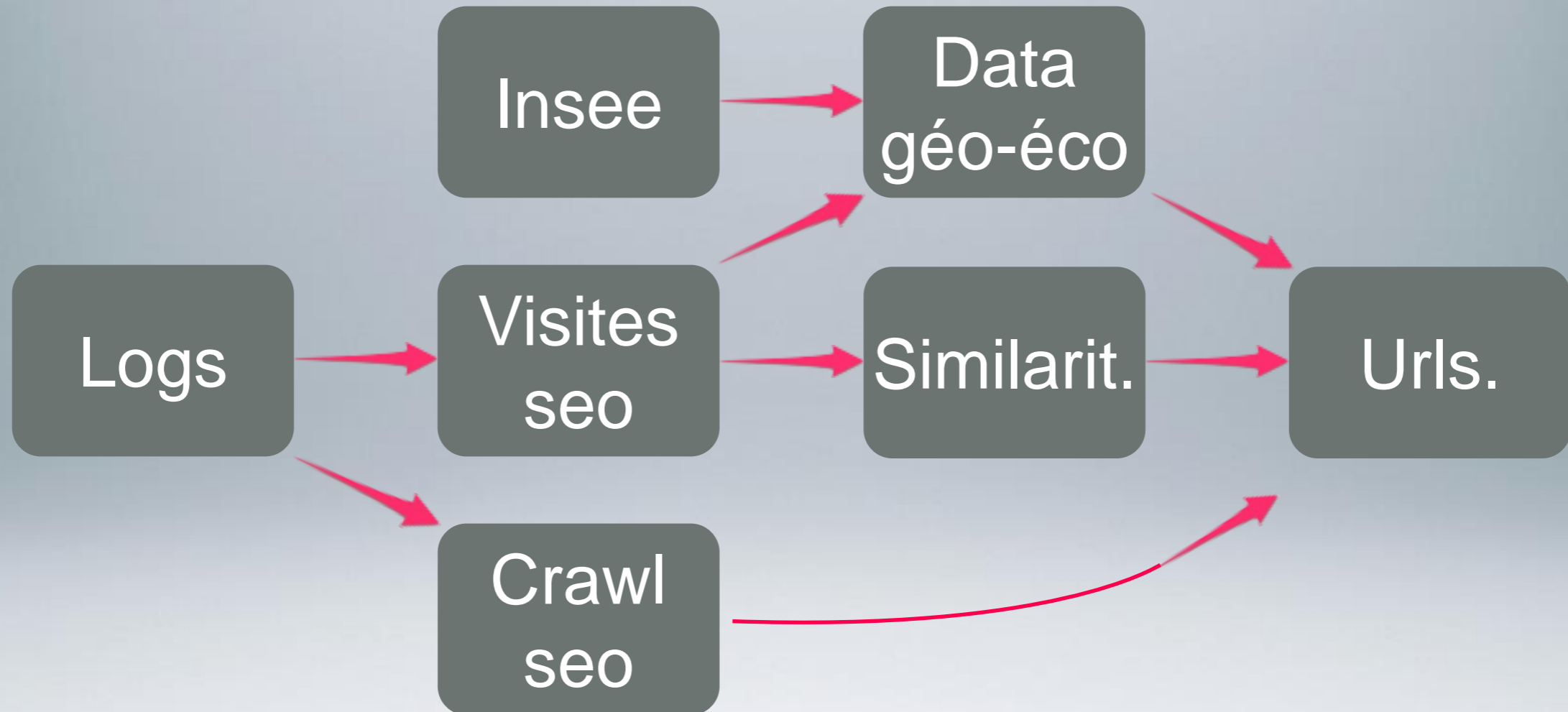
Analyser



# TYPES de REQUETES

- ▶ Analyse par zones de géographique
- ▶ Données socio-économiques
- ▶ Recherches de similarités
- ▶ Analyse au niveau Url (granularité fine)
- ▶ Impact du Crawl sur les visites

# Pipeline

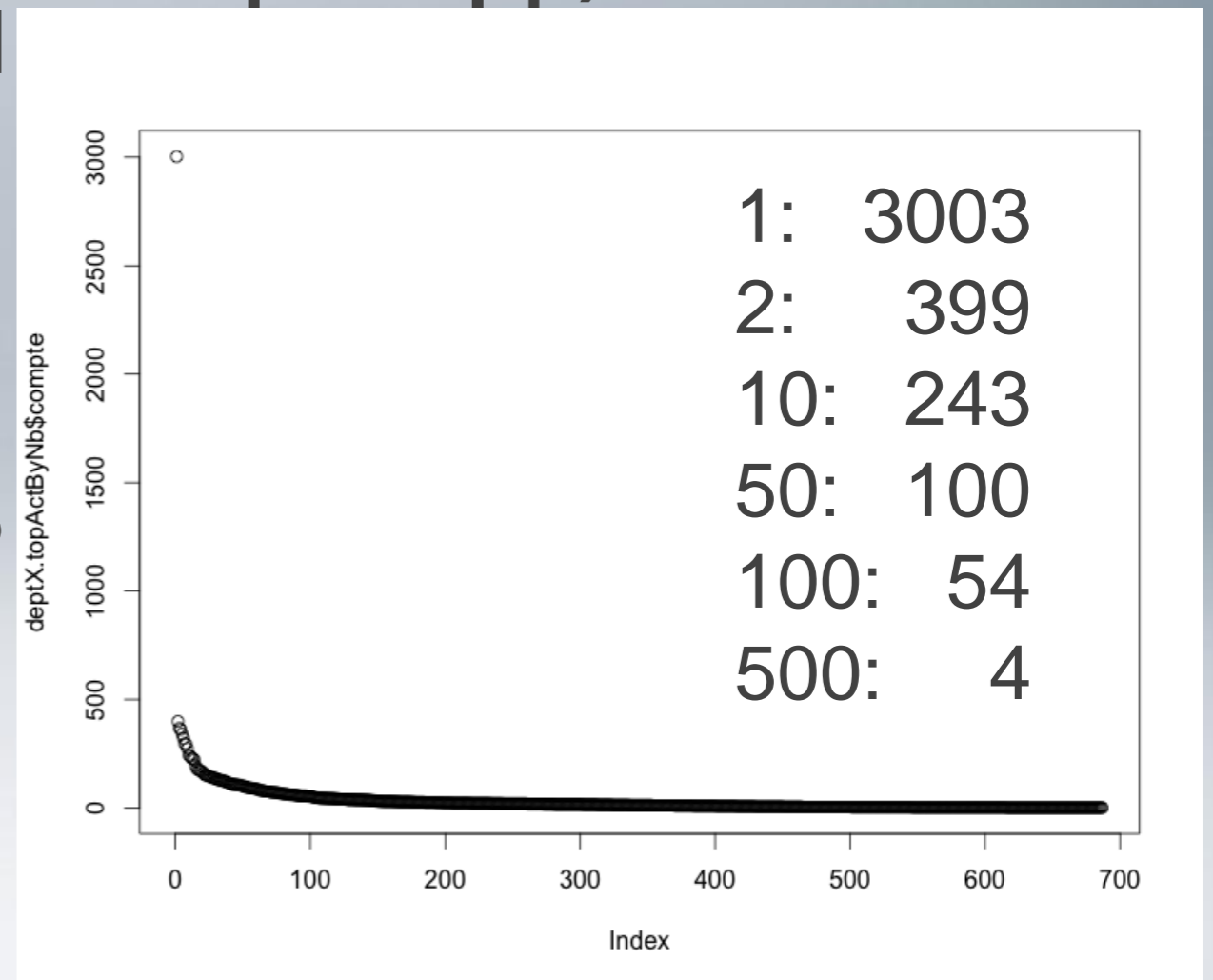




# Long TAIL

- Très grosses concentrations d'activités sur certains éléments

- ▶ Map/Reduce !!
- ▶ PIG Skewed joins



# HEATMAPS

Départements  
(96)

			2			11	1		
			9	32		9	21		
			1	3		1			
			48	9			42		
150	27	61	29	77	15	67	73	58	
1352	782	1199	757	2733	244	2255	2016	2187	
45	35	115	274	100	9	36	105	112	
		50	8			117	22	17	
13	20	21		14		158	46	9	
			8	5		2	7		
			2	2		6			
							1	1	

Volume de  
visites

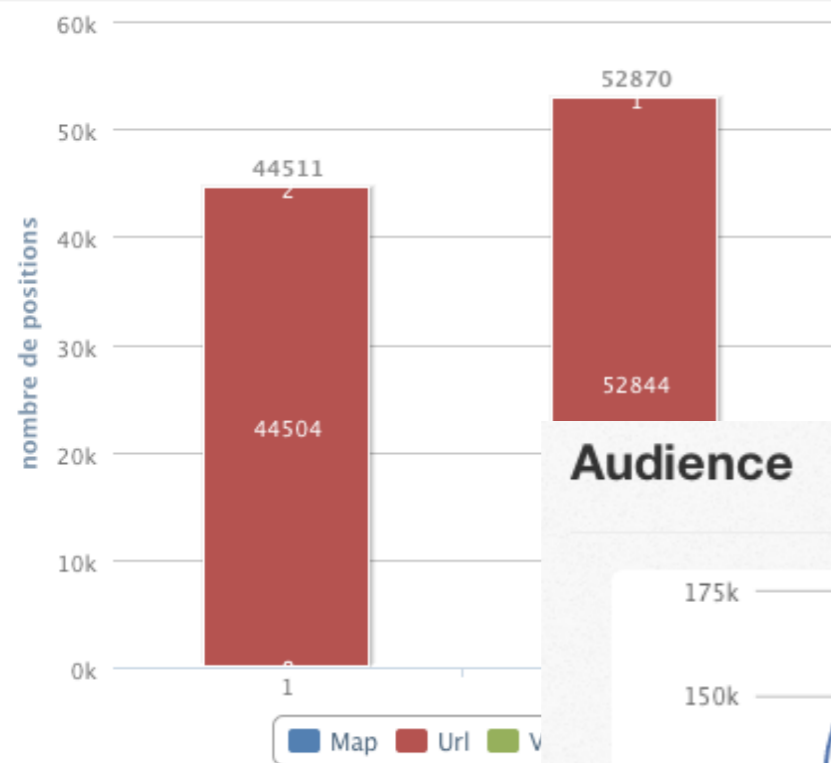
Activités (10500)



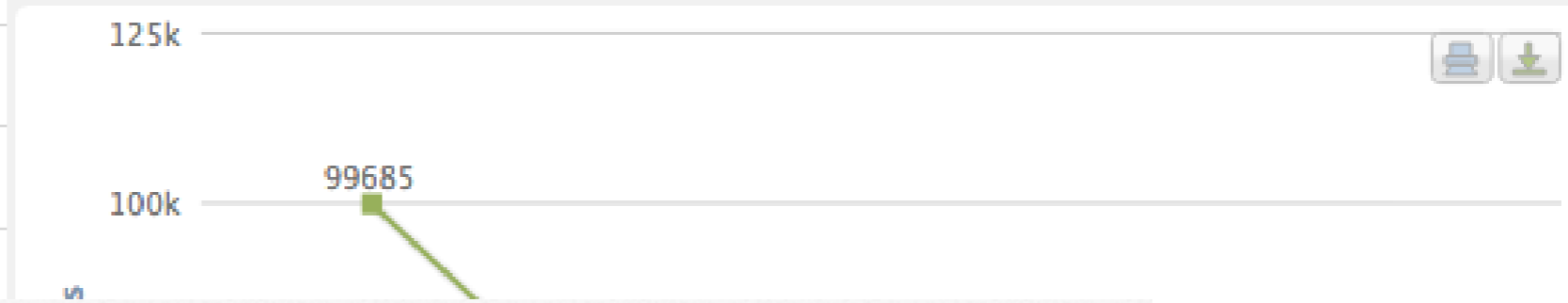
# MESurer : Rankings

- ▶ Collecte des réponses Google (30x par recherche = dizaines de millions par mois)
- ▶ Forte croissance de la volumétrie
- ▶ Classifier et Segmenter par produit, par thématiques.

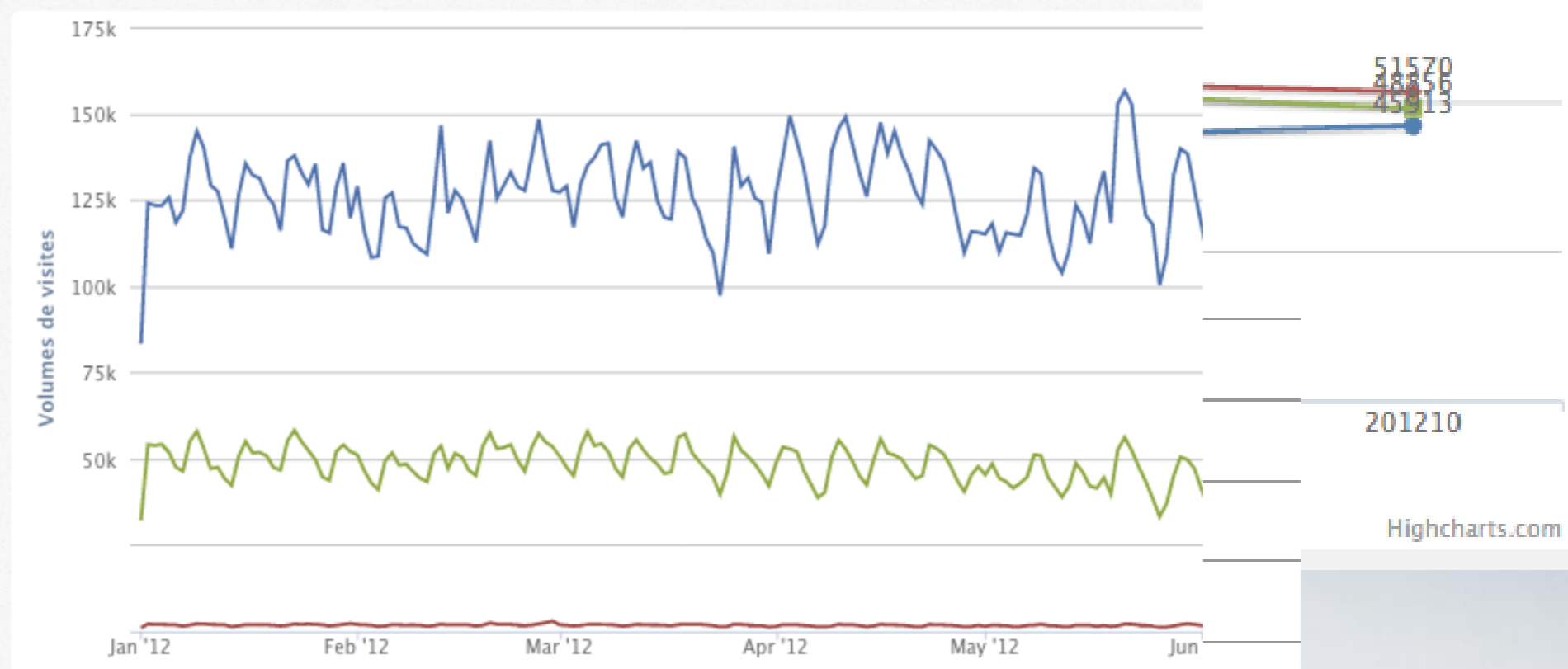
# distribution des positions en pages 1



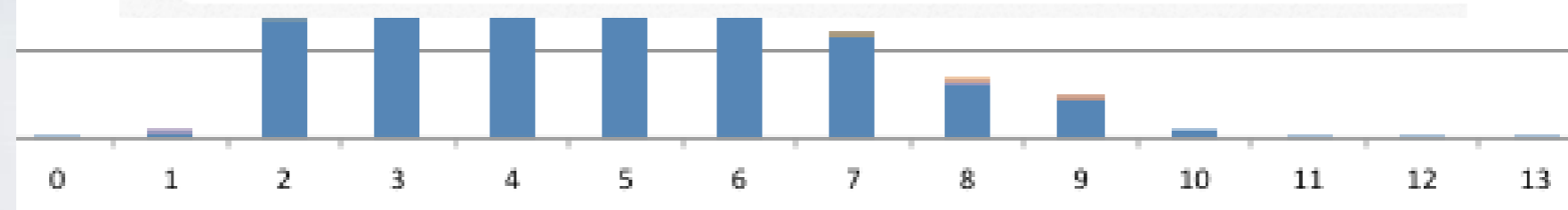
Evolution du nombre de positions par classe dans le temps :



## Audience



	1
Map	2
Url	44504
Video	5
Image	0
Adresse	0



51570  
48856  
45813  
201210

Highcharts.com



# OUTILS

Collector

Traiter

Stocker

Visualiser



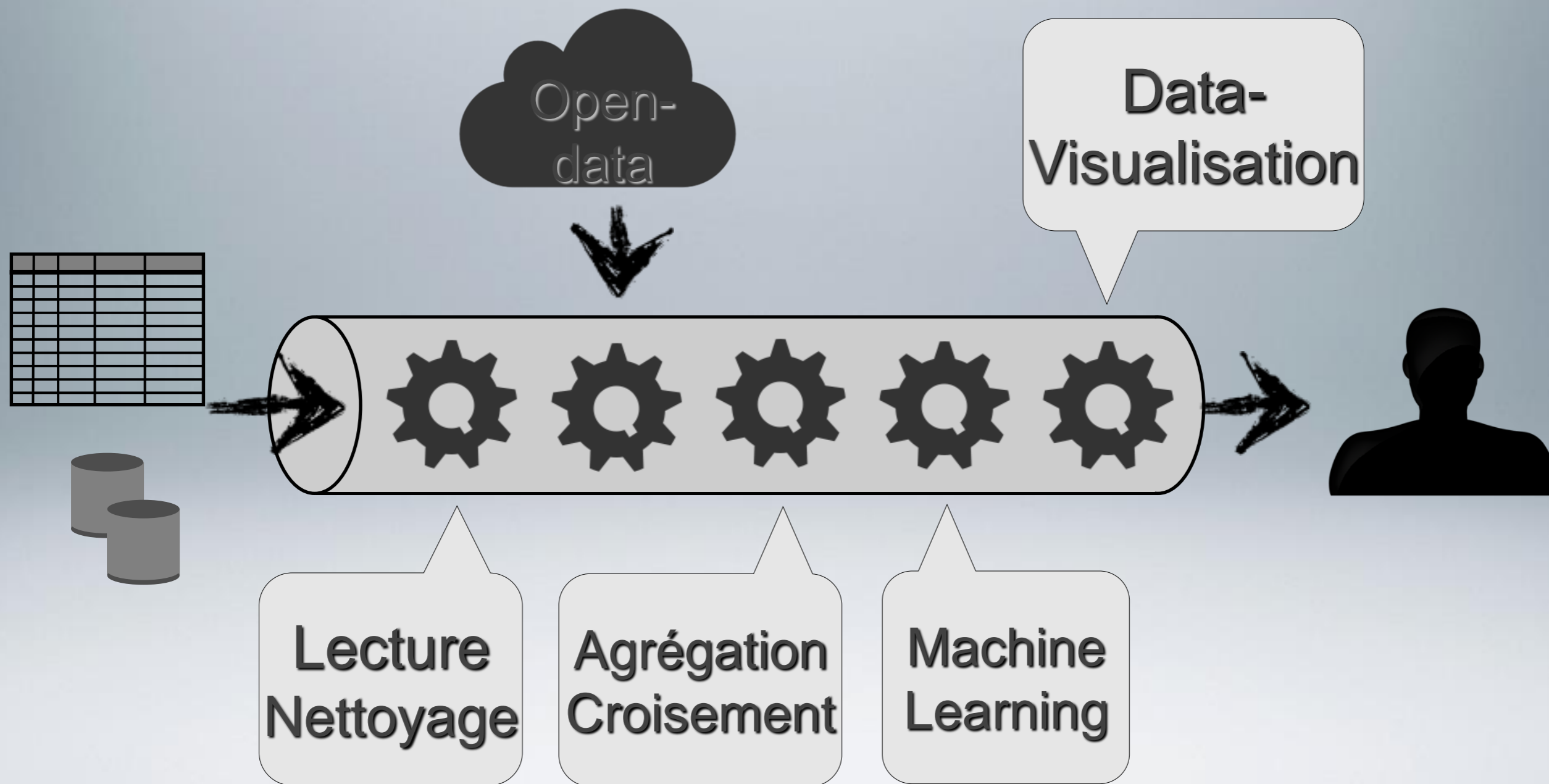
Analyser



# GENERALISATION

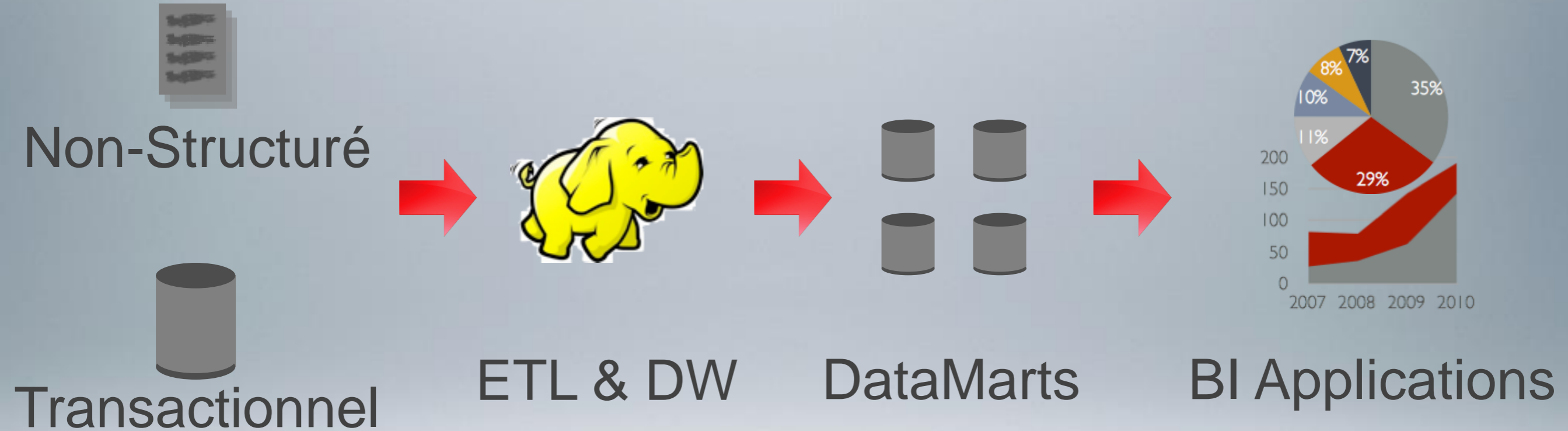
- **Applicable à toute transaction**
- **Des données brutes**
- **Augmenter la donnée**
- **Similarités et Classifications**
- **Recommandations**

# DAta-PIPELINE

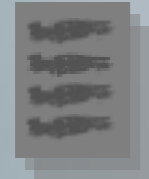




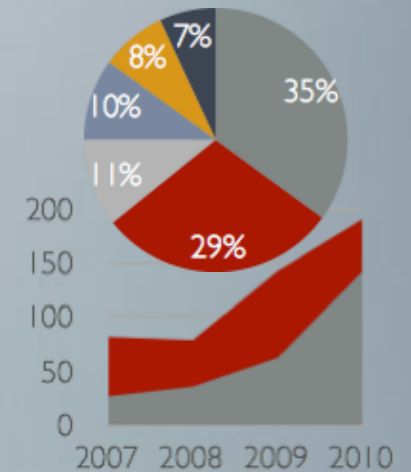
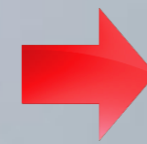
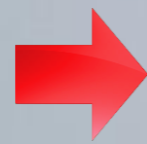
# **hadoop** : ETL & DW



# **hadoop** : EDW



Non-Structuré



  
Transactionnel

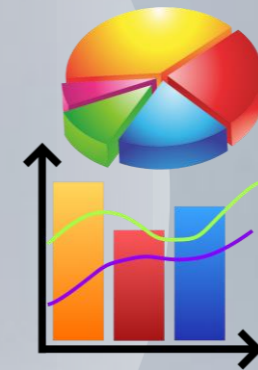
ETL & DW & DataMarts

BI Applications

# Applications & Machine Learning



Opendata




Visualisations  
Tableau & JS



## Plateformes



- 
- A close-up portrait of Marissa Mayer, CEO of Yahoo, with blonde hair and blue eyes, smiling slightly. She is wearing a blue top. The background is a solid yellow color.
- **"With data collection, 'the sooner the better' is always the best answer"**
  - **Marissa Mayer, Yahoo CEO**

# Merci !

- Vincent Heuschling
- Gsm : 06 61 88 76 71
- Email : [vhe@affini-tech.com](mailto:vhe@affini-tech.com)
- Web : <http://www.affini-tech.com>
- Twitter : @affinitech & @vhe74

