

# hurence

**Big Data – get its magical power**

[Laurence.Hubert@hurence.com](mailto:Laurence.Hubert@hurence.com), CEO & CTO

<http://www.hurence.com>



hurence

your Big Data expert





# Hurence : le 'pure player' Big Data



## Une équipe de spécialistes du Big Data

Une équipe d'architectes, de data scientists, de développeurs seniors passionnés par le Big Data.

Une équipe innovante qui participe à la plupart des initiatives Big Data en France (projets d'investissements dans le cadre du Grand Emprunt).

## Une approche indépendante

Une stratégie d'indépendance vis-à-vis des fournisseurs de solutions Big Data.

Une volonté de promouvoir l'intégralité des offres Big Data pertinentes pour nos clients et d'entretenir des relations privilégiées avec les plus grands acteurs (IBM, EMC, Oracle, Microsoft, HP, Dell, SAS, Teradata, Cloudera, Hortonworks..)

## Une expérience importante

Plus de 30 clients « Big Data » en France.

La première société à se positionner sur une offre pure Big Data, autour de Hadoop et de son écosystème, en France.



## Consulting & Services

### ARCHITECTURE BIG DATA

Une offre de conseil en sélection d'outils et architectures Big Data, sur tout type de socle.

### INSTALLATION ET GESTION

Une offre de dimensionnement, provisionnement, installation et gestion d'infrastructures matérielles et logicielles Big Data.



## Formation

### FORMATIONS BIG DATA

Une gamme de formation pour tous les profils et toutes les technologies Big Data open source.



## Technologies

### LOGICIEL BIG DATA

Une offre logicielle basée sur l'écosystème Hadoop, dont plusieurs produits natifs Hadoop & HBase et des composants d'extraction Big Data.

### SERVEUR D'ENRICHISSEMENT

Une offre logicielle dédiée analyse de logs dont un serveur d'enrichissement de données de logs.



## Service

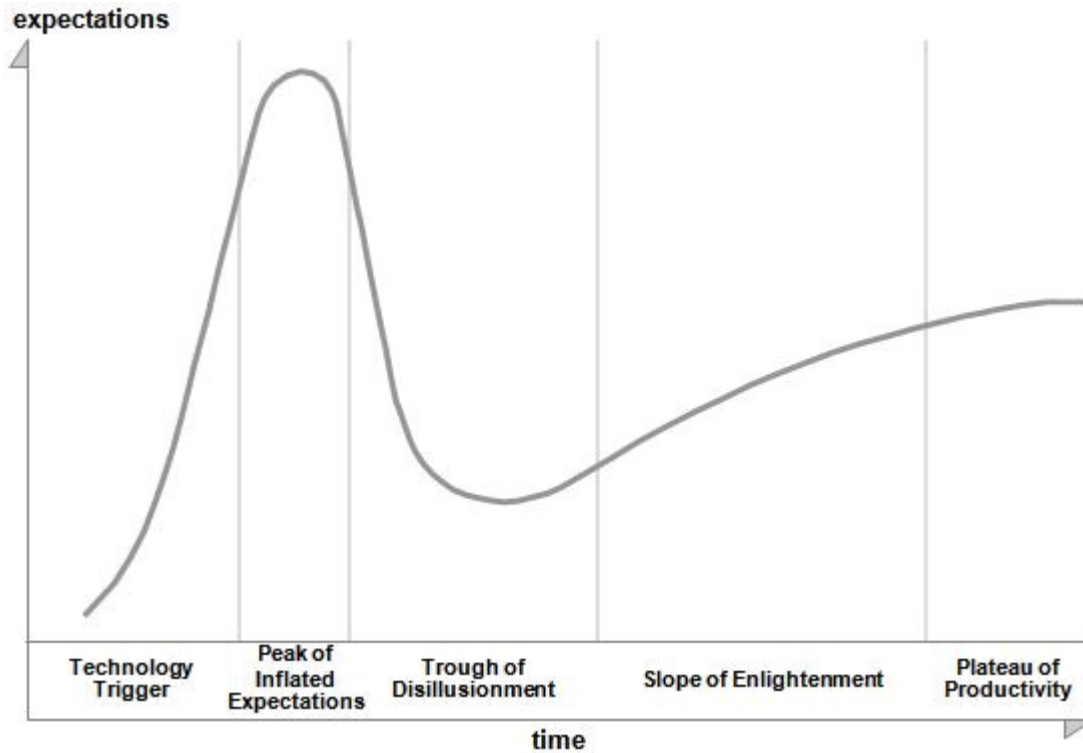
### EXPRESS DATA SERVICE

Une offre de traitement de données sur nos clusters Hadoop.



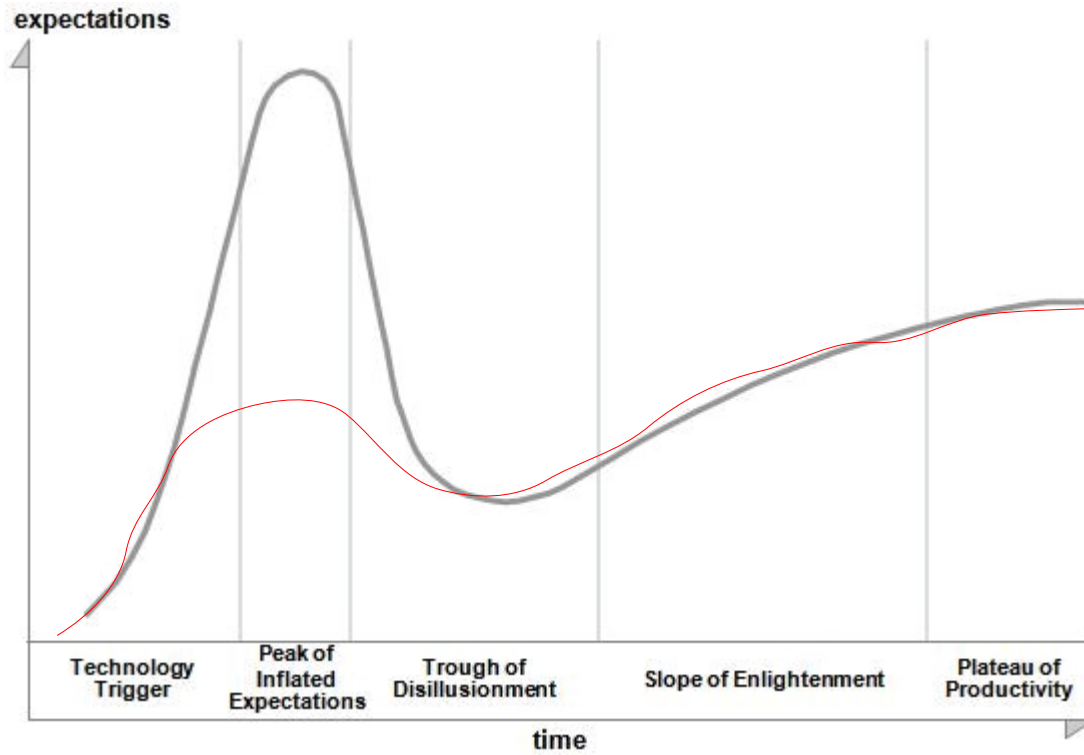
# Gartner Hype Cycle

---





# French curve





## Big Data challenges and opportunities

---

- The Big Data « me too » phenomenon
- The Big Data «NIH»
- Big Data is data !
- Big Data is hard !
- Big is Big !
- Big Data is sometimes poor data !
- Big Data is Hadoop but not only Hadoop
- Big Data is not just data... it can be beautiful
- Big Data is ... ambition



## The Big Data « me too » phenomenon

---

- Vendors and consultants tend to use the same « Big Data » technologies and use cases... if not this is research or free POCs !
  - Customer Relationship Management through analysis of social networks, Advertising and campaign management, predictive churn etc.
  - Security and Fraud Management

**Big Data cloning :)**

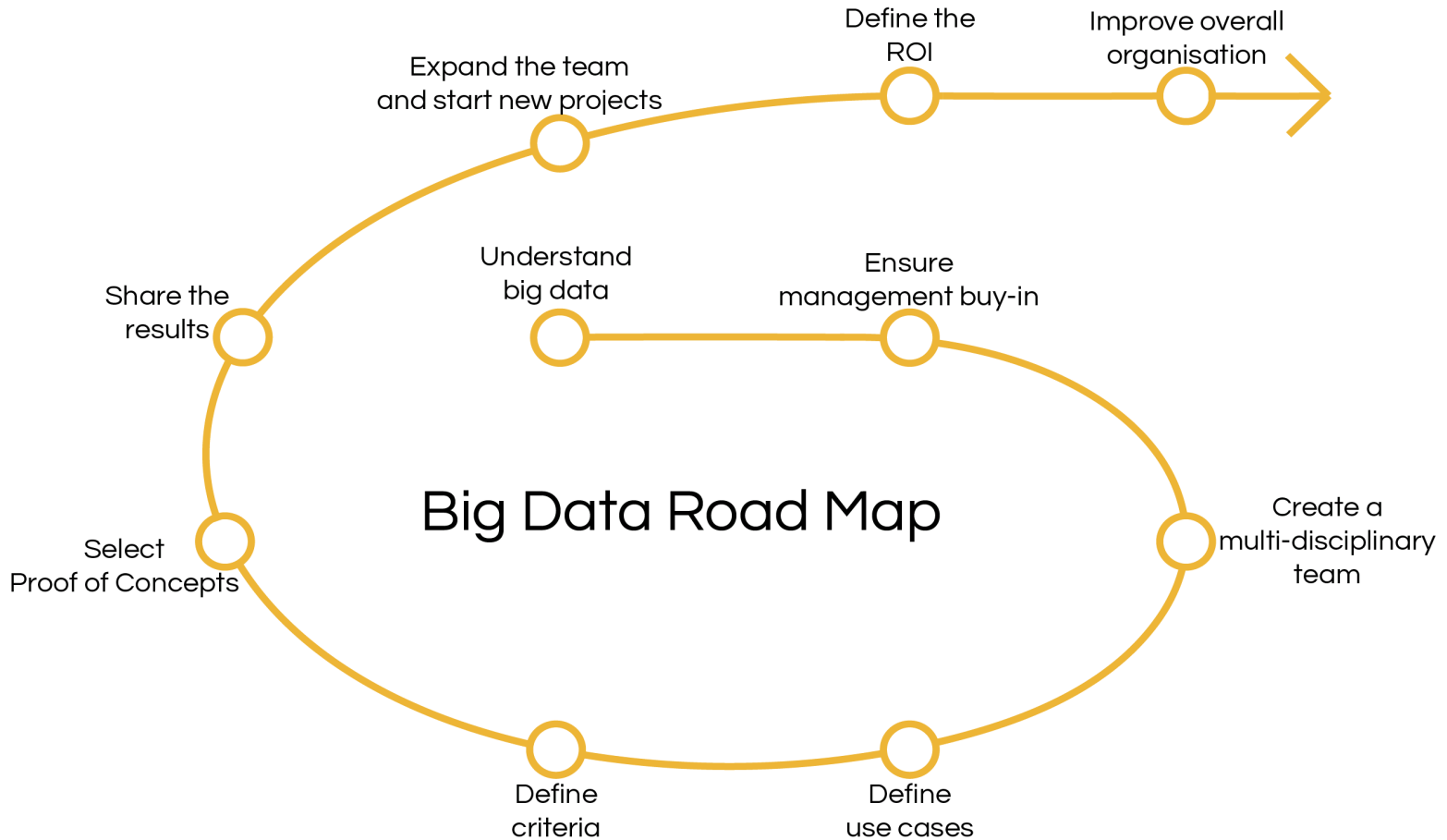


- Managers engage in Big Data because they think they have to...
  - First consulting action is to train them
  - Second consulting action is to help them build an ambitious and innovative Big Data plan with respect to « their » business!





# A Big Data Roadmap...



<http://www.bigdata-startups.com/the-big-data-roadmap/>



hurence

your Big Data expert



## The Big Data «NIH»

---

- Big Data was marketed as a « rare resource » which it is
  - trend : a lot of French managers are trying to turn their teams into « Big Data specialists »
  - but... the maturity of Big Data tools does not yet allow this
  - and... not sure it will ever allow this...
- Big Data is « private » data sometimes => analysis is done in-house
- Not every engineer can turn into a Big Data scientist : on 100 engineers we train, less than 20 have the « Big Data spirit » !







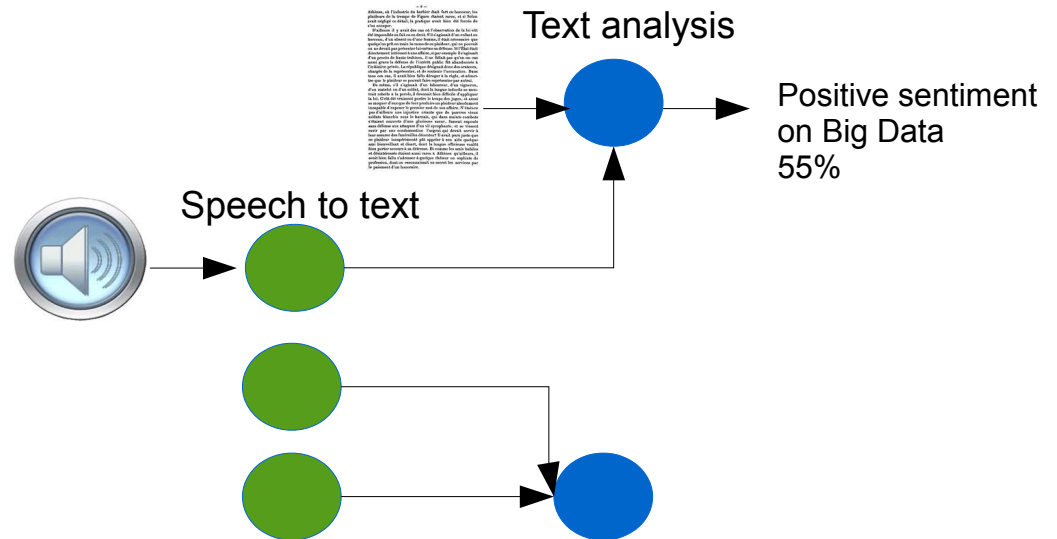
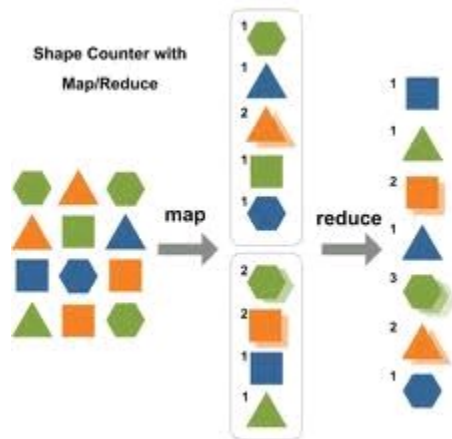
- Siloed data, confidential data
  - Sometimes HUGE political and managerial problems to access data
- Data is power !
- Data is money !





# Big Data is hard

- Thinking parallel is inherently hard
  - **Map Reduce** is « fairly » easy but current implementation in Hadoop has limitations (latency, re-entrance, long-life jobs)
  - **MPI** (Message Passing Interface) in other words **stream processing** is much more difficult to setup and program (IBM streams, Storms)





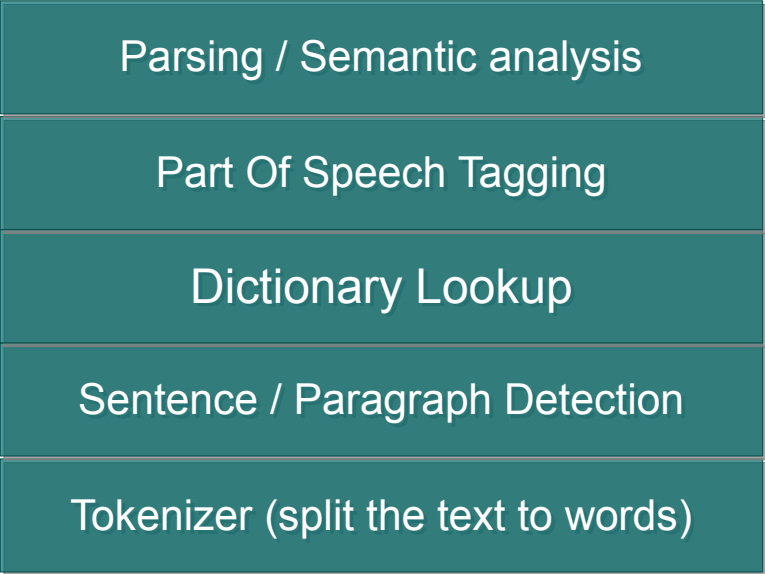
# Big Data is hard

- Advanced sentiment analysis or text mining or web crawling
  - Not just counting positive or negative words !
  - Involves linguistic knowledge and resources (domain specific)
  - Text Analytics is a domain for specialists

## Mix of ML and traditional tools

Big Data Big Brother ! Super!

— 4 —  
 Athènes, où l'industrie du barbiere étoit fort en honneur, les  
 plébeux de la tresse de Figure étoient rares, et si Sotus  
 avoit voulu en offrir, la pratique avoit bien dû faire de  
 l'un accuser.  
 Et comme il y avoit des cas où l'observation de la loi étoit  
 impossible en fait ou en droit. S'il s'agissoit d'un enfant ou  
 barbare, d'un idiot ou d'une femme, il étoit nécessaire que  
 quelqu'un prît en main le soin de le plébeux, qui ne pouvoit  
 en avoir pu prendre lui-même en défense. Si l'État étoit  
 démocratiquement gouverné, que seroit-ce les juges, d'un  
 parti de haine barbare, il ne s'étoit pas qu'en un cas  
 sans genre le défenseur de l'intérêt public fit allusion à  
 l'industrie privée. Le républicain étoit donc un orateur,  
 chargé de la représentation et de toutes les fonctions. Mais  
 tous ces cas, il avoit bien fallu déroger à la règle, et admettre  
 que le plébeux ne pouvoit faire autrement que se défendre.  
 De même, s'il s'agissoit d'un laboureur, d'un vigneron,  
 d'un marchand ou d'un ouvrier, dont la langue indécise ne seroit  
 venue à la parole, il seroit bien difficile d'appliquer  
 la loi. C'est où étoient parés le temps des juges, et aussi  
 se seroit d'un cas que de leur produire un plébeux absolu  
 impossible d'en parler le premier mot de son affaire. N'est-ce  
 pas d'ailleurs une injustice évidente que de punir un vieillard  
 indigent, sans la barbare, que dans toutes les langues  
 étoient couvertes d'une glorieuse robe, sans être exposés  
 sans défense aux attaques d'un vil esclave, et un vieillard  
 n'est pas un condamné. Les juges qui devoient servir à  
 leur donner des fonctions d'honneur. Et avoit pour justice  
 ou plébeux inépuisable qui étoit à son aide quelque  
 soit. Surtout, et d'ailleurs, dont la langue indécise étoit  
 bien partie secourue à sa défense. Et comme les sans loi  
 et d'ailleurs étoient sans ressource à Athènes, qu'il étoit, il  
 avoit bien fallu s'adresser à quelque orateur ou quelque de  
 profession, dont un reconnaissance au secret les services par  
 le paiement d'un honoraire.

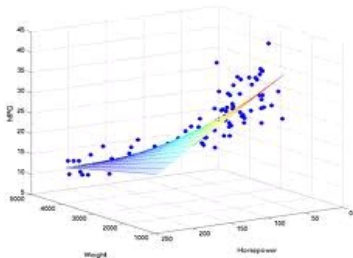




# Big Data is hard

---

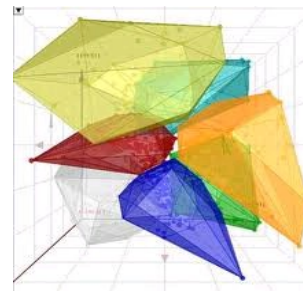
- Advanced Data Mining (beyond traditional BI)
  - Not just drag and dropping data with nice tools !
  - Need to define the right variables
  - Your variables model somehow your hypothesis about the world – so you need to have an hypothesis !
  - Big Data is just here to compute !
  - This is also a domain !



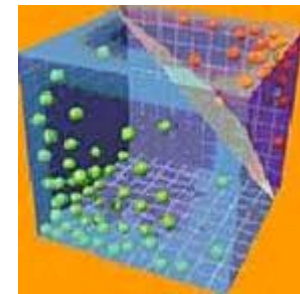
**Linear Regressions**

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

**Bayesian models**



**K-Means**



**Support Vector Machines**

Etc.



**hurence**

your Big Data expert





## Big Data is Big

---

- The 4Vs are for marketing : Big is Big
- Big Data technologies like Hadoop are overkilling for small and medium data
  - Hadoop YARN (2.0) should improve this
  - Hadoop focused on making Hadoop « consumable » through SQL (Impala/Hive/Stinger) but is slow on improving on the real-time and small/medium data for Map Reduce jobs and to propose other parallel paradigms (MPI).
  - Still big fans of Hadoop but a bit deceived on how Hadoop has evolved over the last 2 years to prevent « advanced users » desillusions.
- Avoid the 5 VMs nodes Hadoop cluster and the 20 mega tests !

US



MY HADOOP IS  
**BIGGER**  
THAN YOURS...

France





Big Data is sometimes poor data



# GetMore

## Enrichissement classique

- IP → **latitude, longitude, ville, densité de population, altitude moyenne,...**
- Ville française → densité, altitude, ...
- Villes françaises → **distance** qui les sépare
- Date (+ heure), latitude, longitude → **météo**



Où?  
**Geolocalisation**  
**Localisation**



Qui?  
**Profiling** visiteur



Qui?  
**Profiling** produit

## Enrichissement custom

- Association d'un groupe de mots à plusieurs **caractéristiques**
- Url : <http://www.myhealth/node/6/> → produits, fitness, femme,...
- Produits : chips allégée → chips, régime, snack,...



Contexte?  
**Météo**





## Big Data is not only Hadoop (NoHadoop ? After NoSQL !)

- Diamonds in the open source domain : Elastic Search is one of them
- You can index the world with Elastic Search and even geolocalize it !

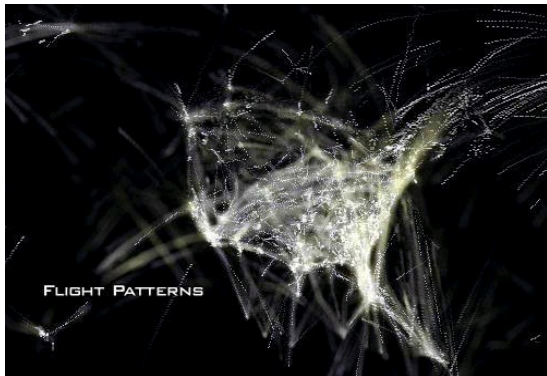




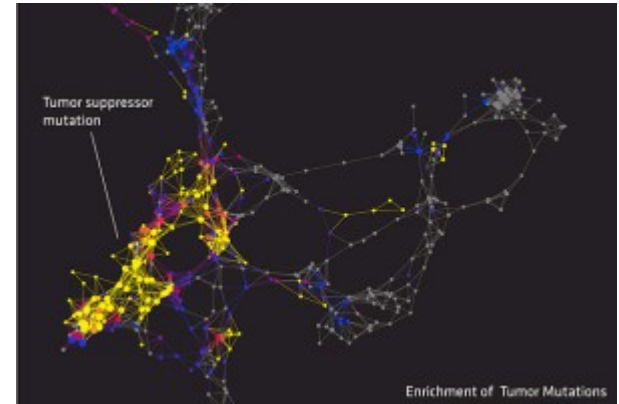
## Big Data can be beautiful

---

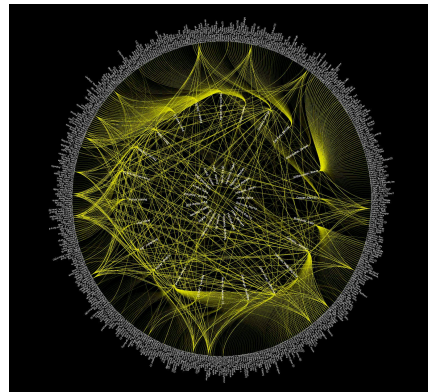
- <http://www.bigdata-startups.com/the-beauty-is-in-the-big-data/>



US Flight patterns by FAA



Tumor mutations



Connections between Oscar winners  
<http://www.pitchinteractive.com>







## Big Data won't make you an innovator but...

---

- If you are an innovator, this is your technology !
  - Big Data players dreamed « Google Search » , « Google Earth » and « Google Map » and figured out the tools they needed for that...
  - They made the work for us ... now our challenge is « just » to have ambition !!!
- The magic powers you get ?
  - **The possibility to store, describe and analyse the world !**





---

## Thank you!

- Laurence Hubert
- Email : [Laurence.Hubert@hurence.com](mailto:Laurence.Hubert@hurence.com)
- Web : <http://www.hurence.com>
- Twitter : @hurence