

# Data Marketplace Initiative

Weicheng Huang

National Center for High-performance Computing

National Applied Research Laboratories,

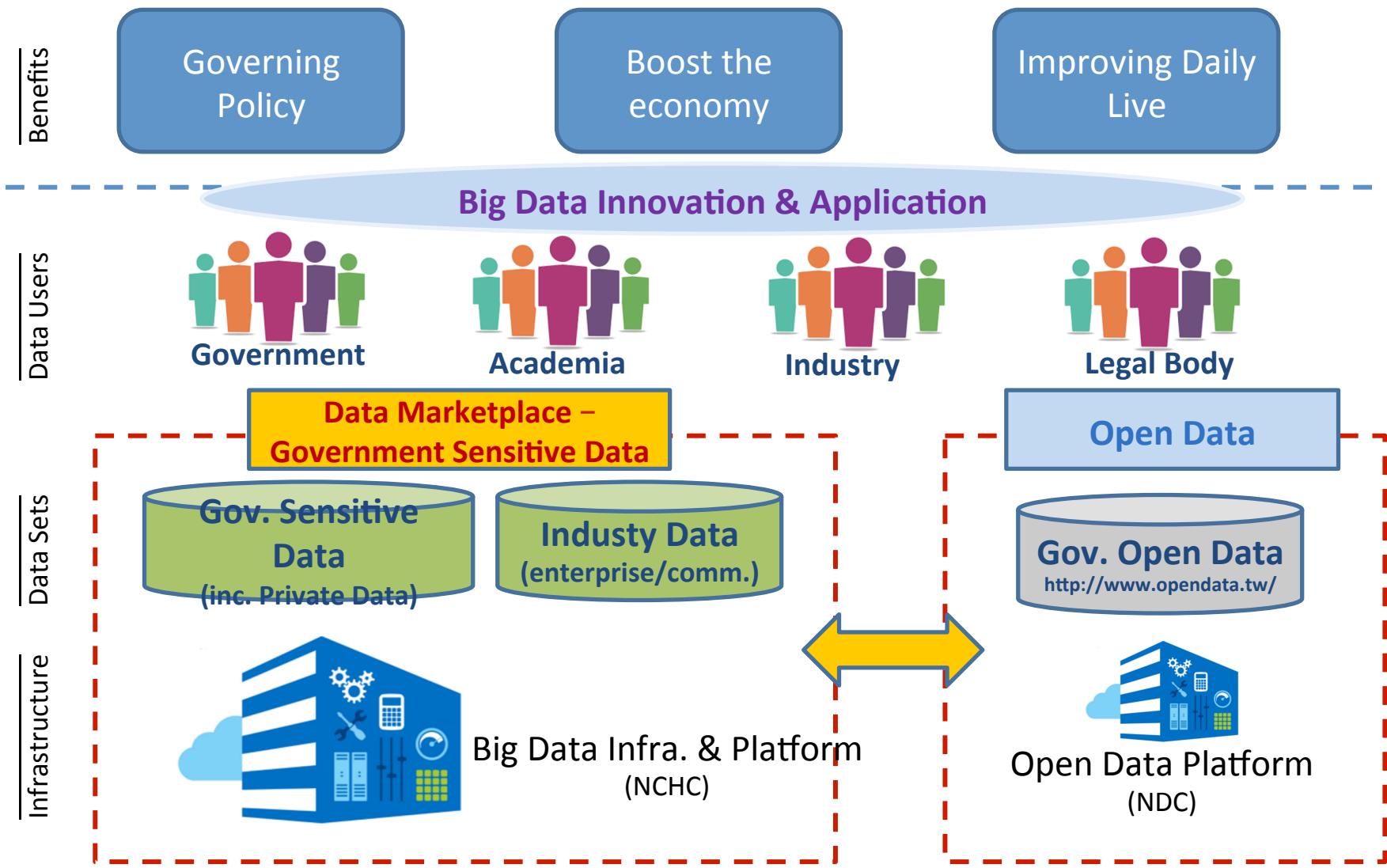
(NCHC), Taiwan

# Contents

- Goal
- Activity
  - Legal obstacles
  - mechanism
  - Tool/Software
  - Platform & Service

# Vision of Data Mart

- Data Economy requires the application of “sensitive” data.
- Secured Data Platform shared by the government and industry, can boost the data value



# Data Marketplace



## 1 Data Cloud & Data Integration

- Data backup** : hosting important/sensitive data w/o utilizing cold data
- Large scale computing** : Big Data Analysis Platform, e.g. Hadoop, ... etc.

## 2 Tech. & efficiency

- Variety of Data Analysis tools** : big data toolkits, data processing tool, ... etc. Tool/tech. for Big Data
- Applications** : promoting Big Data tools to gov. agencies

## 3 Data Analysis & Vis.

- Data Processing, promote data value** : data cleaning, catalogue, aggregation, ... etc. work with the Data Analysis.
- On-line visual analysis tool & env.** : selection and adoption of tools/software

## Benefiting Gov. Agencies

### Security & low cost

- Reduce the cost of hardware/software
- Shared data center & data management

### High performance platform on demand

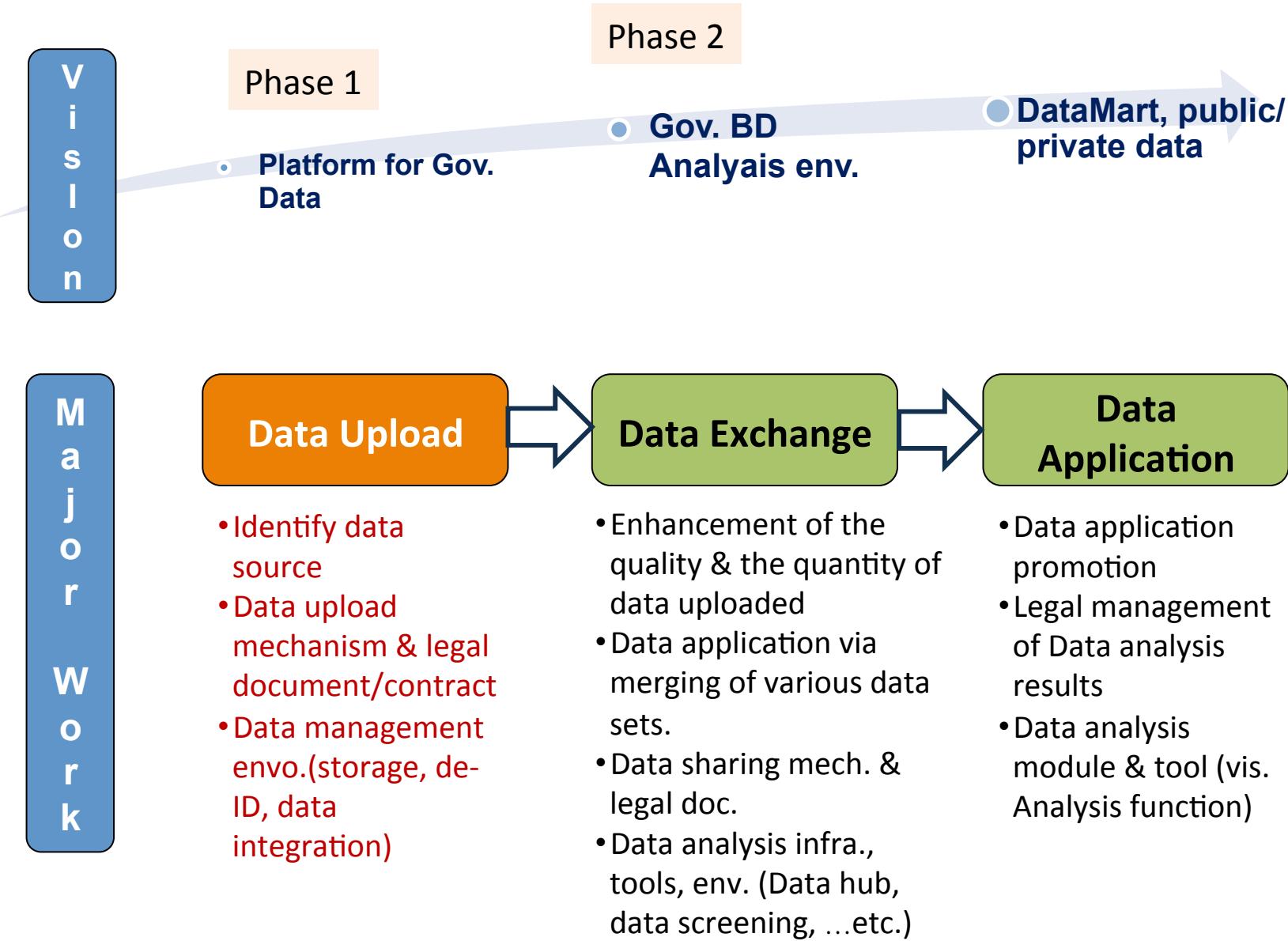
- Shared big data analysis tools, w/updated version
- Reduced IT efforts, training time & cost

4

### Result deploy

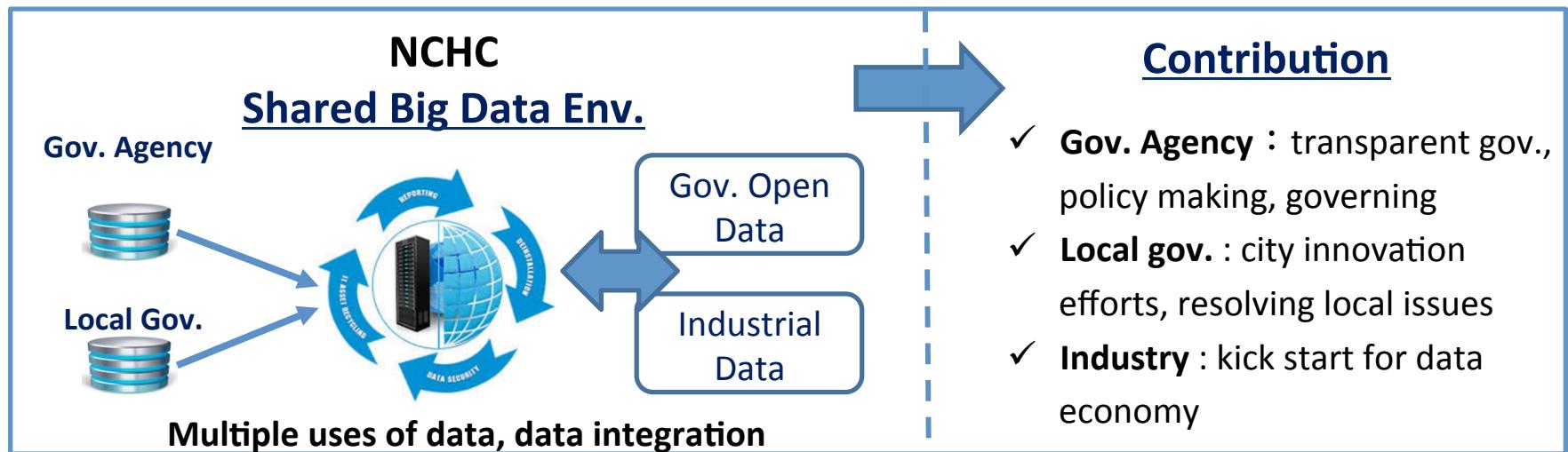
- Data processing mechanism is taken cared of
- Data analysis via web page

# 資料市集三階段發展重點



# Shared Data Env., for Gov. data application

goal	Non-open data from gov. agencies, co-located to NCHC's shared env., Provide to gov. and academia application	
merit	Env.	<ul style="list-style-type: none"><li>Secure &amp; flexible data analysis env. &amp; tools w/lowered hardware &amp; software investment and data management cost</li><li>Simplified gov. data integration for better data application efficiency</li></ul>
	Ap.	<ul style="list-style-type: none"><li>Identify core data sets that meet the demand from gov. agencies (Top-down research topics) &amp; public concerns (such as education, society security, safety, ... etc.)</li><li>Gov. agencies identify research topics, academia answer the call for proposal, MOST support the research projects</li></ul>



# Conclusion – Legal Issues

- Goal : Recommendation of legal adjustment for big data application
- Based on :
  - Study of International cases
  - Study of Domestic law & regulation
- Resolution
  - de-ID
  - Case by case investigation for data from gov. agencies, No single solution that fits all
    - Purpose of data collection, nature of the data, data policy of gov. agency
  - Merging of Gov. data and civilian data
  - Single Gov. entity or single platform
  - One-stop service

# Mechanism

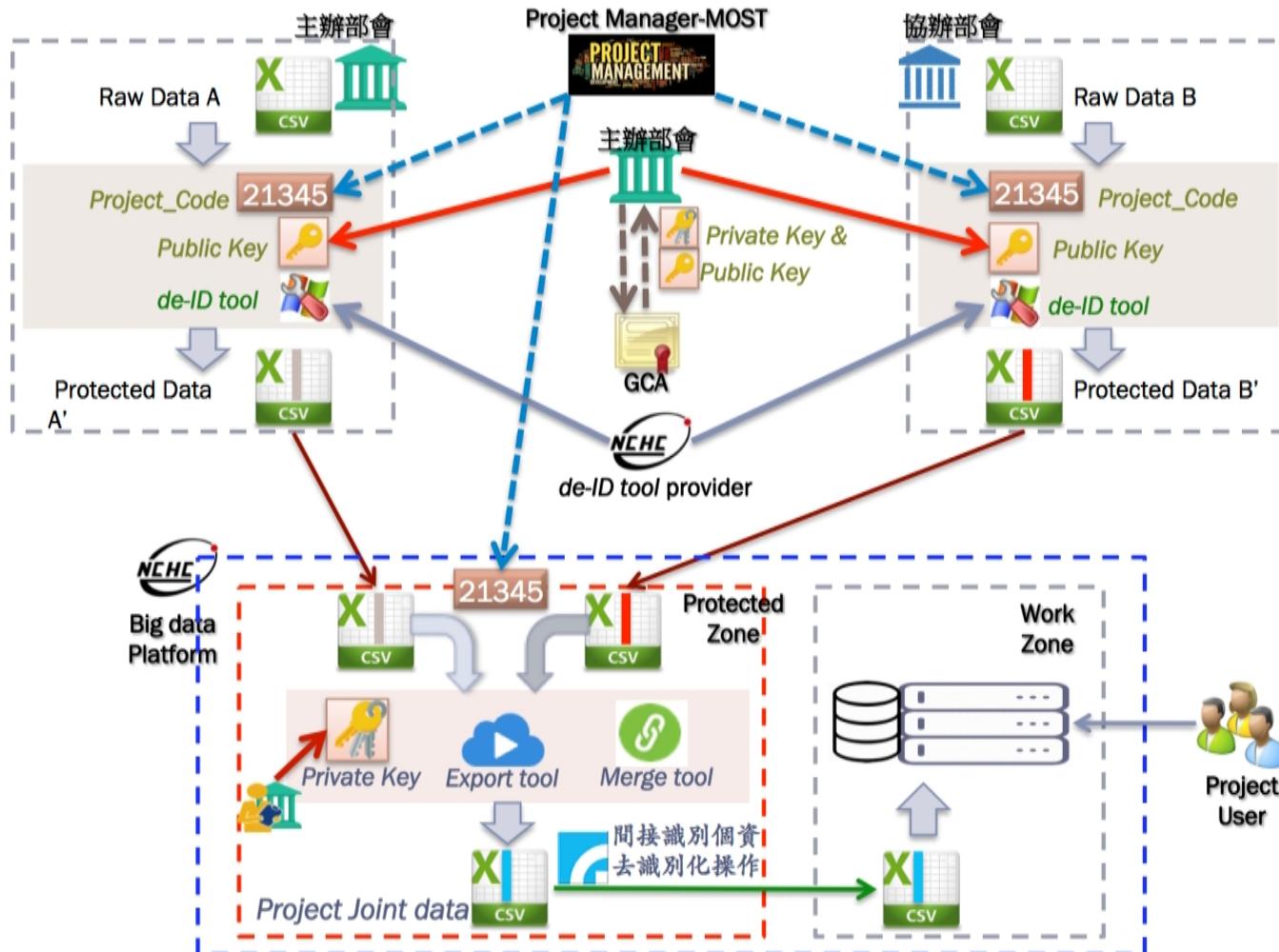
- Goal : Recommendation on the mechanism for the project
  - Surrounding the development around specific application, with 4 major arenas
    - Disaster Mitigation, Security, Traffic, Medical and Health application
  - Delivered by integrated project or data competition
- Based on
  - Study of International cases related to common data analysis environment as well as cases for data marketplace
  - Study of the status of Domestic data applications

# Tools & Software

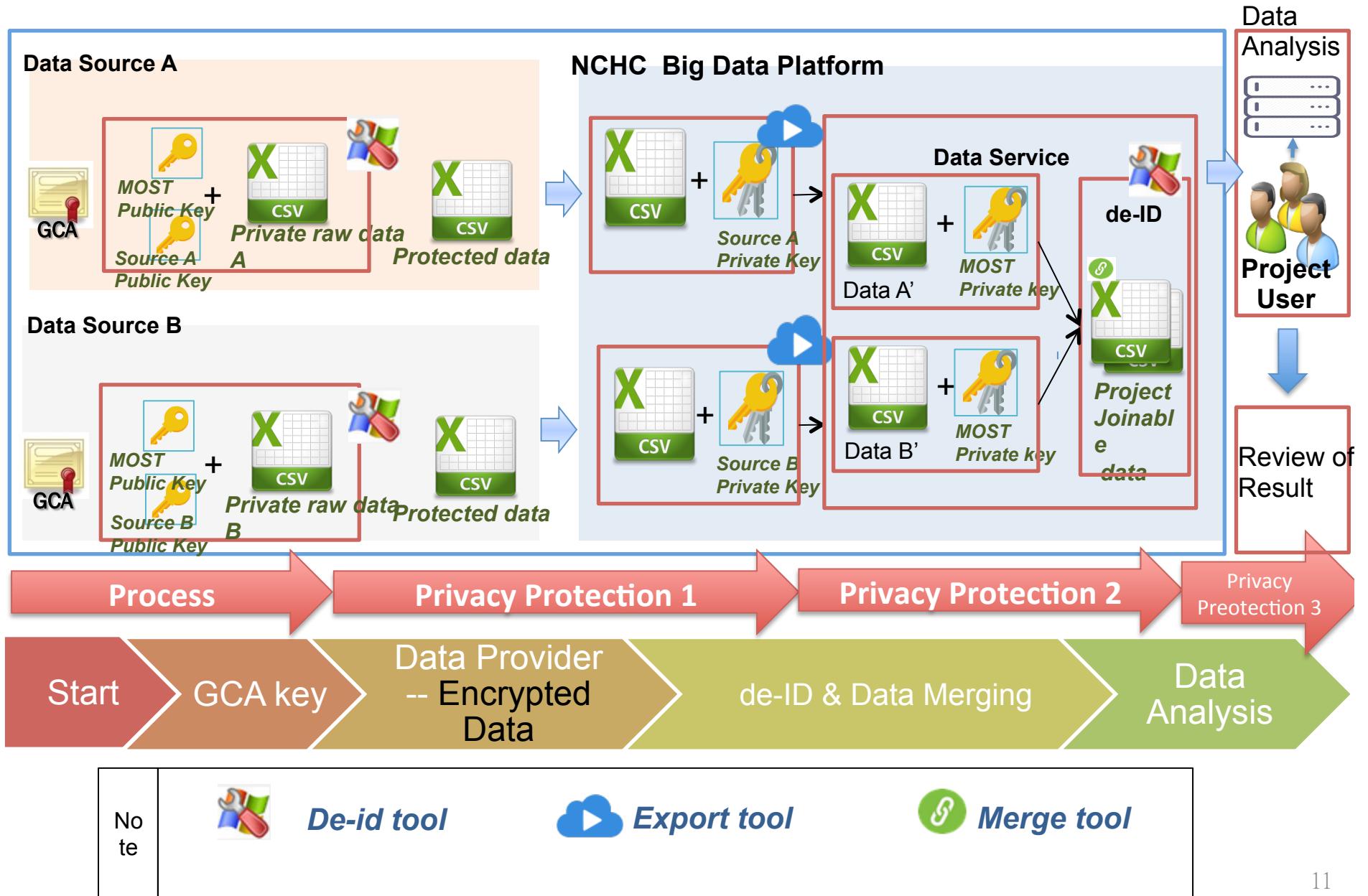
- Goal : provide sensitive data from Gov. agencies for academia to analyze
- Tools :
  - Mechanism for de-Identification
    - Randomizing of direct IDentification information
    - Anonymizing in-direct IDentification information
  - Data integration across different agencies\

# Tools & Software

- Mechanism for Randomizing the direct ID



# Tools & Software



# Tools & Software

- Tool for Randomizing the direct ID
  - Provided to the data owner



# Tools & Software

- Example of Randomized direct ID

## Original Data File

身份證字號,地址

B391987589,40646臺中市北屯區 北屯路 雙 30號至 56巷  
A399274755,40678臺中市北屯區 陳平一街 單 133號至 139號  
G399741197,30348新竹縣湖口鄉 長富路 2 段 全  
D399078461,10843臺北市萬華區 峨眉街 雙 120號以上  
G397129977,51442彰化縣溪湖鎮 民生街 全  
F399293757,26144宜蘭縣頭城鎮 武營路 全  
H399225240,32060桃園市中壢區 內定二街 全  
B392762833,32656桃園市楊梅區 上陰影巒 連 62號至 81之 1號  
B399510515,33456桃園市八德區 高城三街 全  
F393098274,41447臺中市烏日區 自強街 全  
F395266496,24841新北市五股區 成泰路 3 段 單 579號以上  
H392706057,30266新竹縣竹北市 新寮街 單 73號以上  
H391994284,30642新竹縣關西鎮 新德街 全  
B395922966,26546宜蘭縣羅東鎮 新興街 全  
H394767850,43252臺中市大肚區 華山路 連 57號至 160號  
G393113115,40644臺中市北屯區 軍榮二街 全  
E399933808,33064桃園市桃園區 南豐二街 全  
D398245299,40859臺中市南屯區 文山二街 全  
D397441495,24890新北市新莊區 五工五路 雙 8號以下  
C396056651,33372桃園市龜山區 楓樹三街 全  
C391552012,50243彰化縣芬園鄉 大彰路 2 段 連 934號以下  
A391468140,23542新北市中和區 德光路 77巷全  
H398472772,32053桃園市中壢區 中正路 3 段 220巷連 501號以上  
D397629166,42350臺中市東勢區 東關路友誼二巷 全  
C391916212,35244苗栗縣三灣鄉 石馬店 全  
C391581757,33849桃園市蘆竹區 蘆竹街 全  
G391673147,30352新竹縣湖口鄉 大同路 全  
C391532012,35059苗栗縣竹南鎮 口增圍 全  
H398687746,33862桃園市蘆竹區 大福街 全  
C397015449,50862彰化縣和美鎮 愛永街 全  
A391615434,11570臺北市南港區 南港路 3 段 雙 54號至 80號  
D394984537,42076臺中市豐原區 豐東路 250巷全

## de-ID Data File

deId 身份證字號,地址

L1nGYbEPs2ppfigahp/03zNLvHYx0LQYwZn5y1fsSDChpy/DdXlzbKn3GpjjerIQxeAeKKFSpY7SyXK04MdWNd  
etEQ4m0q/lfuidYsCxSv6b7TudQq/5v7LGgH8fYFw9CMszzaTOU5kexYh6M/J5CKgUpmOx0UXBbFXIL0ghynN  
bTKpWe2s3xv748yTN/opqKd42aSHjd5EfWfVJMWVrjjusCNl1mQoJB6VJGRhb7mwT8RbzbsfQfmB3VcSwQpc9  
m9sUZuLjbo004fffnGCvoDiaryURCMnuOtoDCASP8U+T0/S4cSVBVKOSHOnANqc4aj/+cNqCetrVPe6ufBfA  
h6LTJjnQDCmUPEgVEBzBzroj4zLNe9g8aBpaoSrpaJQ9dDw4ivzaOA GddZ1qQFFOSFWcNVCZGoekN2/3jtVgc  
Z5DOC0dyKZiK1XCifGMmEcoHkpXLsR63Ve/BRETFT+hJo/WC/zHkCfxZfkalr0ILTLMH8efH+k1JnVtARu/7f  
LL+vVvDZX1qKQ1lxWRvzAhYBVGtZjc8WGRdzfHeVEXMzMTxzcB1RHxb6/9wG91HmzDosKOSFenuSDndVOX0QY  
LkWjj9csVqu0A17zce4iu054G6GP33IN+4H+DMFOkBhbb1kuA+9vhD9VZ1u2cLRy+CsjKMya9rjp01mf+vWHY  
03DPQguGp1Nz06dM6zYGS2vmydijnc/iILJjzXYuM2K3mayzqcR/eGbauQ1zTkVU/v8Xu0By0uzqsXJPNY50h  
JLGZw5900JAQyOLK+gMR5cbVPnE66jNmqB09LDNURXmsBSManHdbivlx3Z0Ah iEfzFRDwWREB8aLSdLyerRXfv  
eG8ofvBBbRIw75umKmnngPB1ec68Bg4s02zKx3jmtRX7JhuCa3T5g6+T2960wNRPOA8j02RZ+jpjbgDulJaL8  
jaBd0CvtmxHATh34bLnCzxK1f0VY+y+2t1PM0uNM1iCPBzUGzdjntLCvluxflr7qtFc1NRRQUCcd2U4Ft8H/  
WHK/BUTmSF2/U/w71sPQdYQ1ud9SjW5bNyuk6yKs3rb5jQfe9LqnkvN49td1nN+HzstWw1zZJRNdEc1Fn6DE  
d+KbhrV2n60j+41+dKlrGZ6oF44LK1EcSrOpIXQy8DpN1rPCNeDKxC8ApYHj1M2gsSXY6ZeVmQsV1Iex/mYAh  
WSEhucy2YCia+dWa414vIDfen2y3s+IG0Aw7nM1Rb8oMrMo10d1zK+gRFnRc0491Jy2F1MWv6SH169VxTKqKV9  
XuCr0U1PWOXtrfcdgBAJ1Mp02u8xS+4cMJ1sbZ178mPxxy+H1VWn8XfHsVbVHHDDMfl3cKLRcn0VLmzXhkFGy  
R4A1ZZ5Q0seIMfV0negdAmrgWoIPN/UR4tZYxTzTfgM5Xn6zSNWzRG1uk9MgjdB9QU87YdAQ6DyNPsCn10tlm19  
FicewsBALbgPeX0B3hFGuzN/6th9mUN1arWbjq5IBAXBHfoSxeiSJMujiANE9uyHgtuNeLamNh6g/uT0ZIEgB6  
bUZKwkAmGsEQZff7XYXG1uZ0kBNicwHjaWUBHYiwZDOW2FSDw2ram4PNOiYQFELyviuEQVg6h6Xkg1Gt1GdpLM  
c0R/G6g+1bx1toDqJ1nQuKBn6znLjBdZ+PxkgqzwFy/+4mlemp8oS0rlb0+tq2HsuuvL9ZtqV3KT44KackhY2tS  
h1m1c510ArJLhVR580kHJ4R5gUTt07CF3XPtRv+vRGQLVQ/vHQ225+6sfuE5YLdgYjATr643Q0eqiyy0BmTZ39  
LpHU7bYzqFwd2mqYRBTkxBiG5suww7f/50hBZ6WfpMtutJ4Zyxfn+hGdx+sKMU19Y1R2uEcK3SyQ1Drndhhw  
SL9FzgNUzySfmCygsmkupRoe7y+4zBiX7bsM+qGMXVDtb8KaEbaSLyXhV3C110RbTnQ95190CxQn08gTMKVMeu  
L6RI0ekIMv9UrQEoqr0FT193CW2QEG1LcmfC0soxtQxb1E4jCwrfEkdiEHLCevGzPROWtGjw2tGYKeStIjAlfv  
LnleG iGNEARtLtaec7ae4cCYMaXet/4zELYY98HAHQo4K221WpTZVbzazS702gQDolTzCs3fcB38pH/wyBVULLH  
SfRNzZXwood1VL7jhCFAzAP1zkiuNhp2JgKTZucwsCIcmaxD/Z///4Px1uCIW1iUSf10TMslBDfBwzNBvs60tqN  
m4ztprXDxUTYvrqptRAzUqSb03h7JYCie20RCF1hRTBf7mf7yPdEnL3KZLG2p6zltz+D0qRbjzQCI1UDTI/W  
B3kLQ6o2Mhb2po8ix8AC0IC1Ua1KF3yN1Z1/7C4g8/ITb2g3R0zwicBKeWgljTek/MrKTpxwB1VxxbObCVw17N  
b+rTEDRKr6LAtvDHdoy0DgQOFszgbyp7dK9Tx6RhXcQ+17J55qnewCPdMBs6ttgWCbCobSYbpFwL7Htkh4Caf  
eEaxnKP15+aNWNDbreN8chPqsPb2bx0CaxF7wJMVUajjU36nBh2HKKuhBubSw9Vd1bwso+Ka1/SQ7pneksIQkM  
El1twZzDOPJ1G1ljfJ19Xo2N/SCp2+Zoh5Nzds9/vzGwGbjeB6FEo0Dz7oNRT5ZE216NEX50QCxP90x7MaBGH  
KWkEkWaPp9g+uufEY4N8f2HbKkpUQq9CYdtkjG6+5QVKBVqY5iX2Z5RHCCqnw90hWnRz2EQ8wiHe6hWc6AI5dF

# Tools & Software

- Example of Randomized direct ID
  - the rainbow table

de-ID & encrypted ID

Original ID

9zH15g iGTFtjm81LcgHFVHTrXWYE4JyoRGb9YXE3tJNioZeXGP8U5hMcnOja16+keOA5f2BvjGGis=, B391987589  
mga+sGfMDXCjKEP1Emcd1duPXfvI+3XTLpds6uKPNL3i6yXhOZU0rEx6y+vLRtRN1NxX1AzAwKbjw=, A399274755  
CRGWsEximwKnS98YLgc0LjKSvpukoPJ2vYANcZ/gi1uRXHEA/Cy96H5omBHtTynR0R5pL6n4riAwY=, G399741197  
+Zkga0dVA75i+u6UNF/0AaEiG+maFoHP1Y6kqWRaewx1QQWW1CW4DqA9lz4+PTUXMsA0gKwyftdzI=, D399078461  
1E3yq7vbFMU1+/Nrppm6+YfA0icnFAI2zZy1WKUa254Pg2pmeeqahXPfxRX5XER7iMeVNi2XzbWII=, G397129977  
PY6Gha3+mLtWYBuPitP+w6isGS19rtS/W+uKYBHiaZ3sDg5s1QV2vBQXU+8yOLk4sT04Fono6C61A=, F399293757  
WrIZAg+KiCuXNDQIf0iaHWxIku5DcuFxMtpgoD+Cjp1G0e6T6sIJZ1jXBw7E0IZovr7QL5Z/J8UY=, H399225240  
TAYxa71mo0zZW1BwNxcNd8/eE7MGVVU9ZRS+is17dTLTj78gQA/Mxgy0/VC60CgvKYIk05mh8n/4tI=, B392762833  
i9VDG284SqtA5JGY53G6b1puEkEia18go4Q0bK1BGw7e02F7x6TqZeE1eoMWFwHEPoZ4cKVk0Ihg=, B399510515  
le1h/8EE8Q/Dqp26siWrp17iG6jSeCI+F0kz9rD11m3fs02pab5bKDeQYswHXcX21WWzNvvj1Db0=, F393098274  
NmuyW6k/sd9hCxwHI5ZIC5okN54tp0QbbXNYJ7c16McIZPBM59Tbmikvw6fxzgY1eTAo5iONwmLhk=, F395266496  
9XXb+jqzqpEYg1ZUj0oNz1Bbrt66A4SGgBKrsAYSbnbrTPGOKikTC+qeTjxTERXOsCukGjKoNw03w=, H392706057  
RQZn6M15DIvgcy6tWX40I2mNZ+HrEf6T79Z3afZ10neCMBTIU8eH09MrWarhUfMfoKuboBx0htUjk=, H391994284  
A9dmhrYb0E2DMiQJLHRzSVxzu1hS0Ch07+yoYgHvwGEP+6atUFIEkJMq72uCKM7pnZh4qNpc9kSt0=, B395922966  
kygpd4MBspF10shru0nsLwQ1goMbxB1ihJ6SKY20hyfrnzAODLCCQJqpqR37PG77k3+PJBS9Yzd=, H394767850  
/6XtDnujKHj9e1ZG2B9bU3S0hmjB6oa+t0Z1I+ephBuecjtvCj0XEIGXbeFjUyBwt1WfqgLbAcw0=, G393113115  
auGfczcCNVR7aMzaeHsg4e0MrxN22q1QtUzJHyNzivXBHXYLxLqE1MwEPDbyN1BOPBDVY7xLLCR0g=, E399933808  
J7MaZpNbNjPF3VU9Ky/WyewxwvUswHMkTNTS019nCdJvbERiCC6VX41ZyhKTsx/HX0184eo+BXtQ=, D398245299  
FvKzUNrcdRUiDnZEfvzvBZCZ1QVWrQ9CeBy2+mMNYpoe/4lqm4bwq2g1N8rq2n3JgWmlBzrHC769j8=, D397441495  
bwEIWkadPfVb+P2jiiYiyLnRpZT5QCdtCFxUVgY5D/k0VR1jIZXBLVDj1BQ1qzjY1xYUpnExD8W+I=, C396056651  
3PWZD+3TBGLgYyrsLQNGgezh0uXqOEuk5XQa0MQPPSaPuK+HabJY8iy0Cg6SR7pXbf1C8h9fpGao=, C391552012  
fOM2TrWC1Mbbant1tPoGZXn9eERct16hyB0Wm8/agP03V0LpCgwSbbxmuc+J4yiaba2+HWFizWd24=, A391468140  
EjQZjmRsUs6uPFmzhBWWrPkom2c15cHaFivnW+su/TapEtQgorG5TixUvp9vWBQ1zVozs5u+XSkgY=, H398472772  
Dyfu124Qw0SzZRYCvLM6yt5X9XYn7cp/ofwnkySNTEtTKY4NAsoiBkD1dYwA36pe1mB0qad9bjLAG=, D397629166  
FouoqUWwaPkRo9ugVtA5U1zDrevv8b50w8/9LPzoQj1g81UB+b+P7JLBMVyi/AcjmloqbMXWR1HE=, C391916212  
5E01f4Epng9JxQuyG8yyi+0GRkt9XsEjXOFCQ/te8Dgiig0Tu4f4hwezfStTo0WJPjyzwG15a8og=, C391581757  
PfVQ0qSGFvMqJ6wqa2TgtYd5jXDxTcsFMOSX1xpOZyH4ZtBndHB9305aeAa2Red6L01Y1iAf6IXV4=, G391673147  
1Nx6js1RKMqHPax8GPchgo0aQERzRg9ru4RsxE12fdhNqKq7hdsj93TPu6T1VEWtqnxykhnqGNW4=, C391532012  
8ZABkgpYgB0UnRRLeJFErfY4KHPyoawYgkJaQEsulKIOnXg7m0OL10R/5yZwCemfbxZ1P628KnRS0=, H398687746  
8Wm7mpj28URoRoVzc80/OPhRtCB6V0911GBMh439aUh0BTu2xJsozVFcFRqf/whAHA/j2a781GKow=, C397015449  
9d5CoP6f2EF0bo14u1fb++05ahNDiWgo16Cb4H1/MYXwnLJ110LJZzX9dH3kMPnR0po5Tg26n/Ges=, A391615434  
QzJMeEXdXE33ahR7m9ex7WvVooVpoJZ2I1DTb1o9KbVeFLofzvN62I/iExy8ZTfUJ2BCzXugVA1XE=, D394984537  
6vavG4CeI31Zo35xZj2txcUn9XTbcL VFuQTG4mrRY/W220Fr3Sys1kZ5FAQsM4W/374CxyCxdh3pc=, B395293120

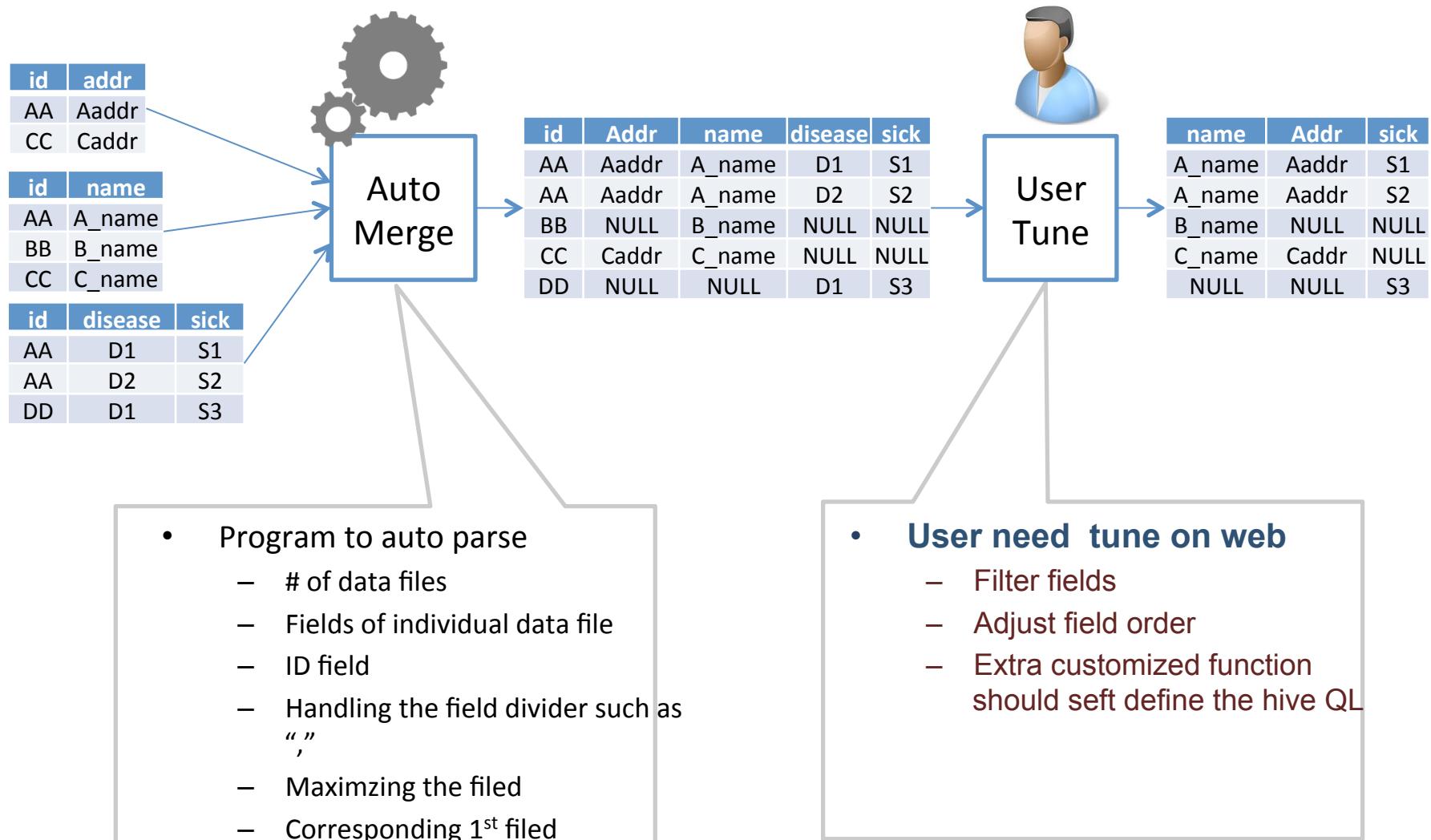
# Tools & Software

- Performance of Randomizing the direct ID
  - Parallelization
    - Multi core : processed by multi-servers simultaneously
    - MapReduce : parallelization based on Hadoop cluster
  - Performance
    - 24 millions records
    - Non-Parallel : > 100 hrs
    - Multi-core (4 nodes) : 35hrs
    - MapReduce (13 nodes) : 17.25 hrs

# Tools & Software

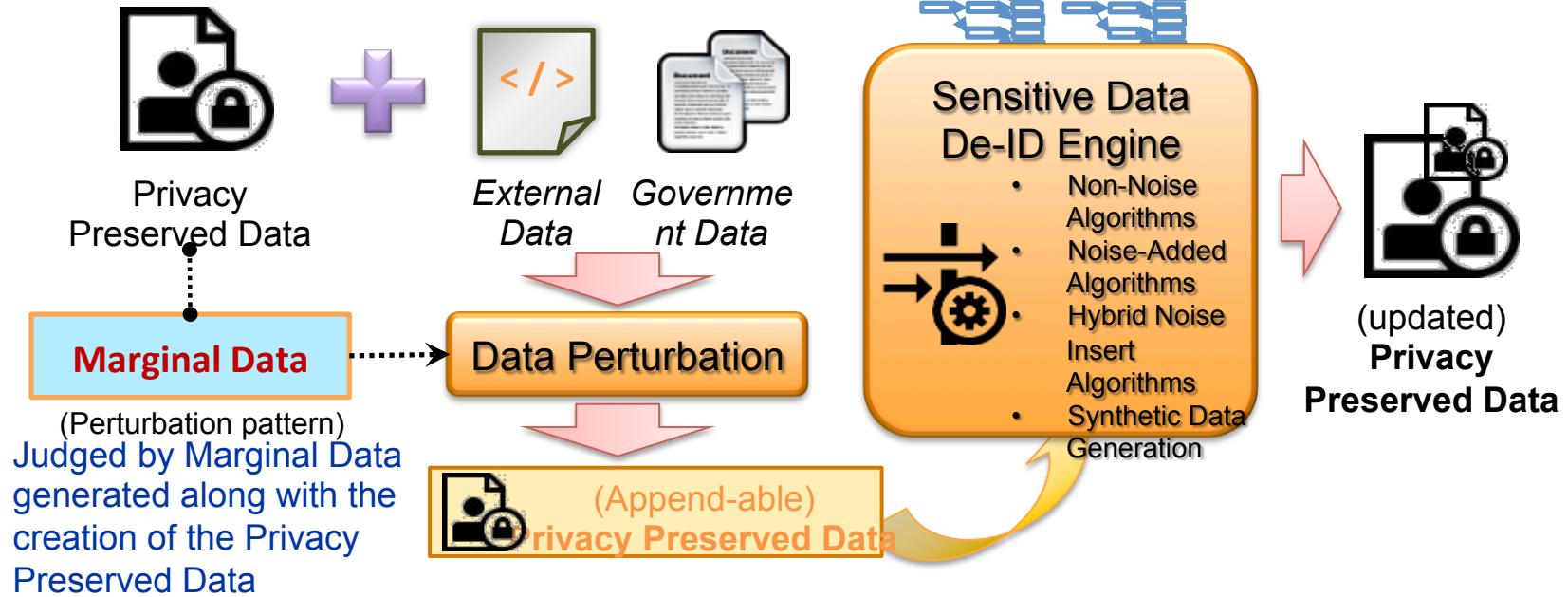
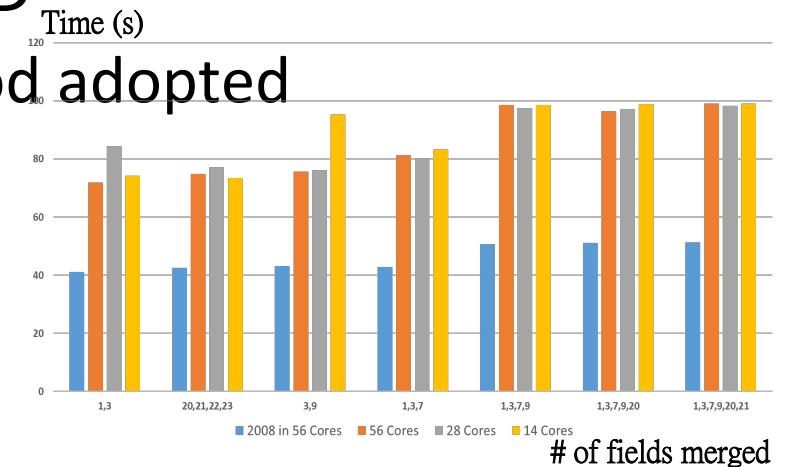
- Merging of data from various Gov. agencies

Design by ELT ( Extract – Load – Transform )

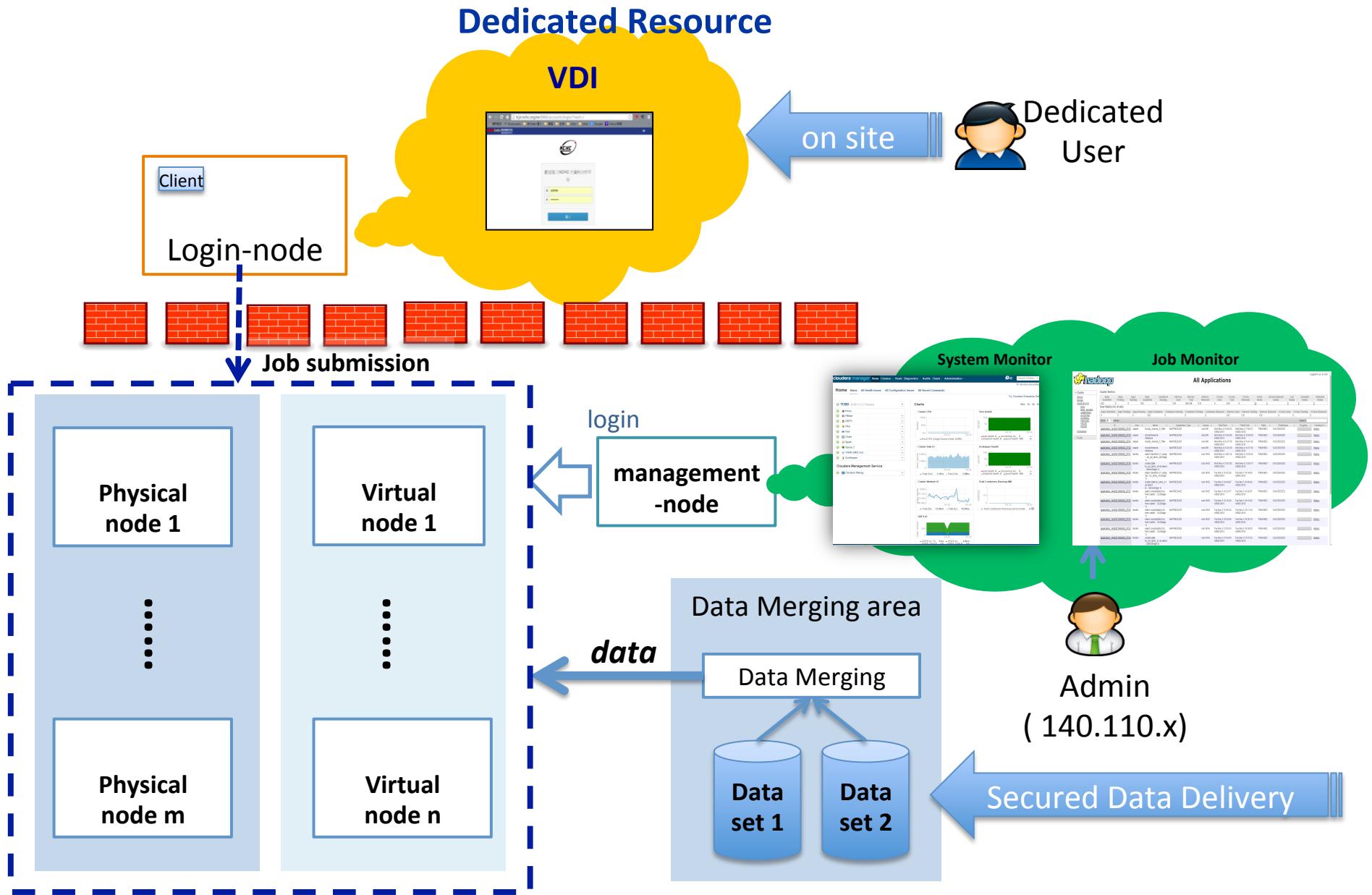


# Tools & Software

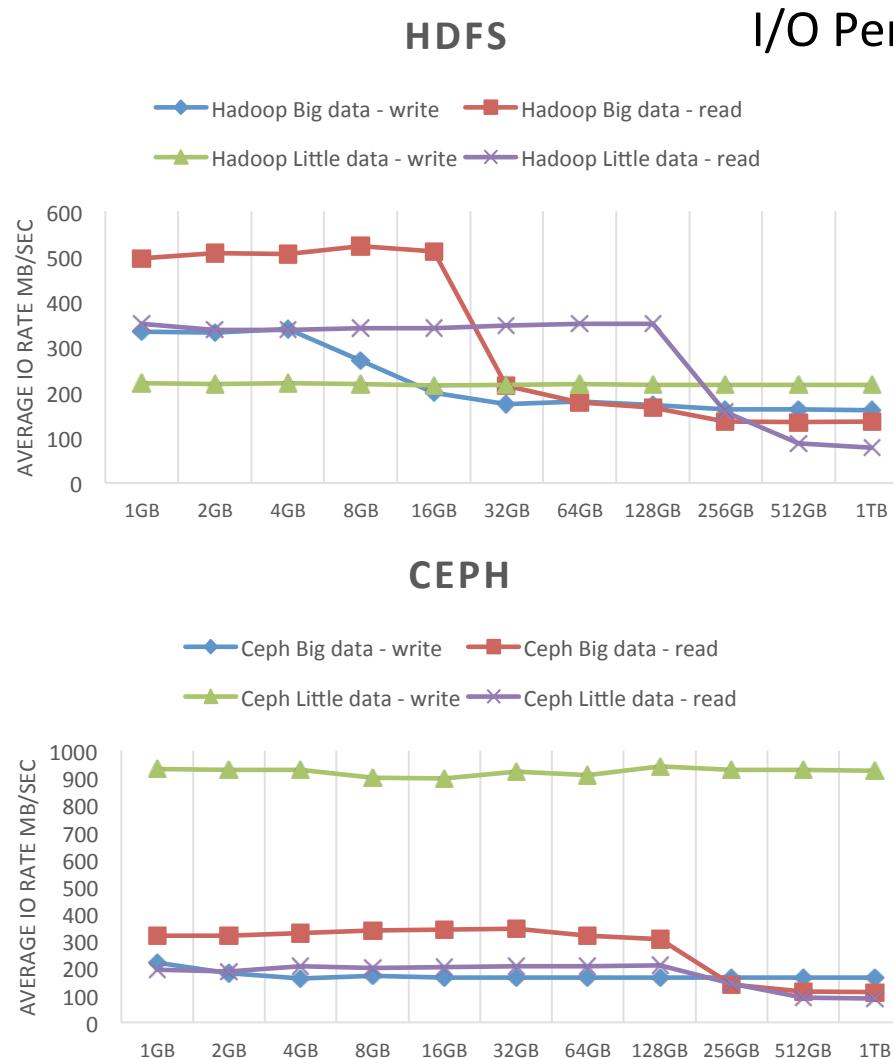
- de-IDentification of in-direct ID
  - Parallelized K-anonymity method adopted
- Performance
  - 1.4 billions records in 4 ~ 5 hr
    - w/5 fields managed and verified



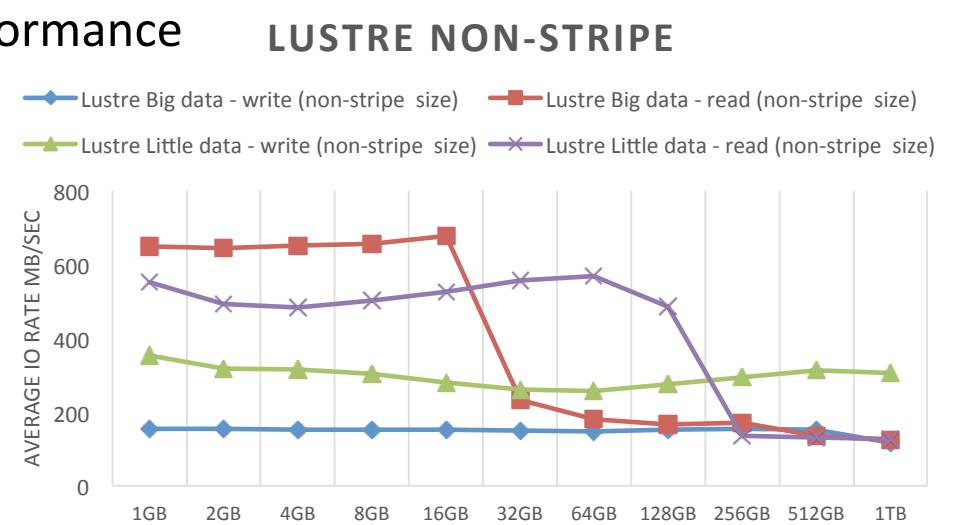
# Platform & Service



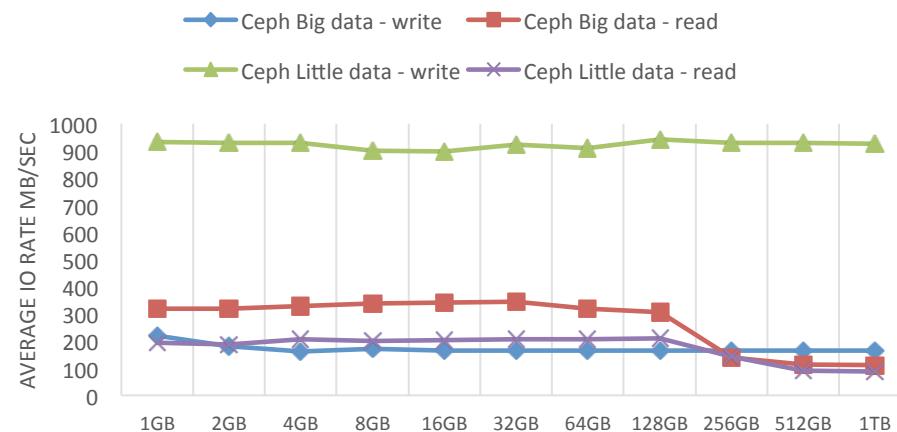
# Performance Comparison of File System



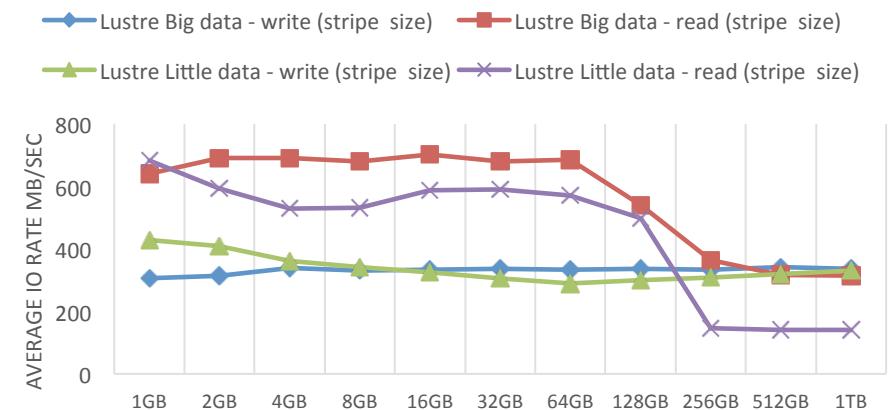
## I/O Performance



## CEPH

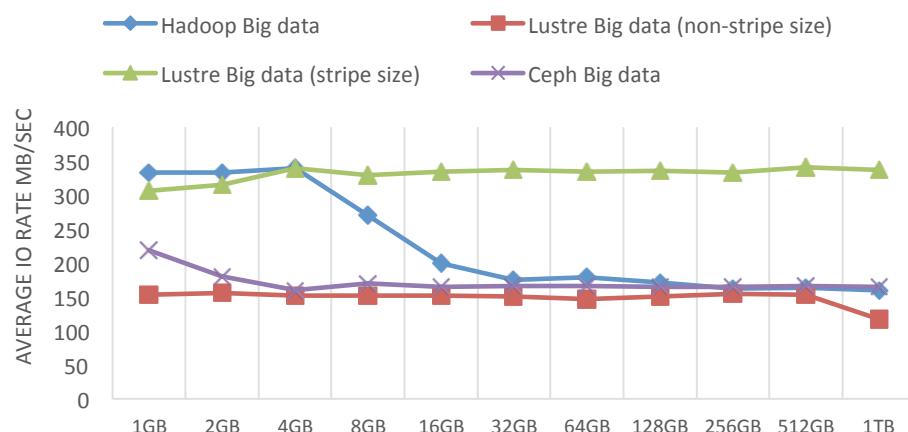


## LUSTRE STRIPE

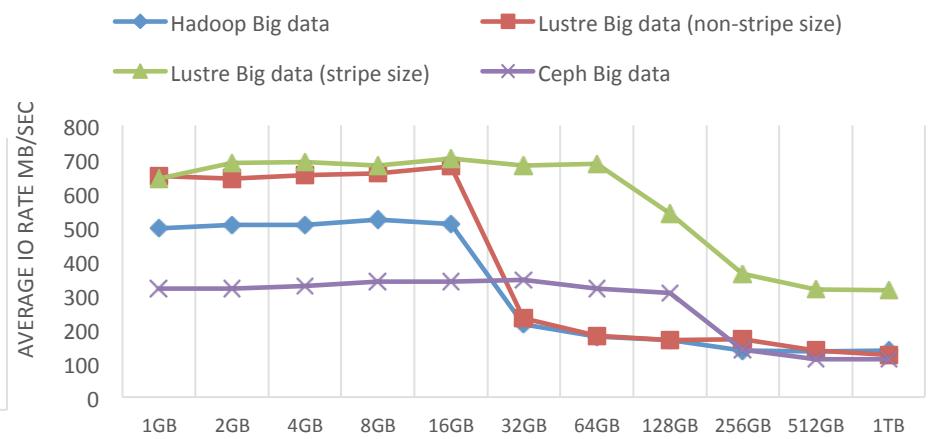


# Performance Comparison of File System

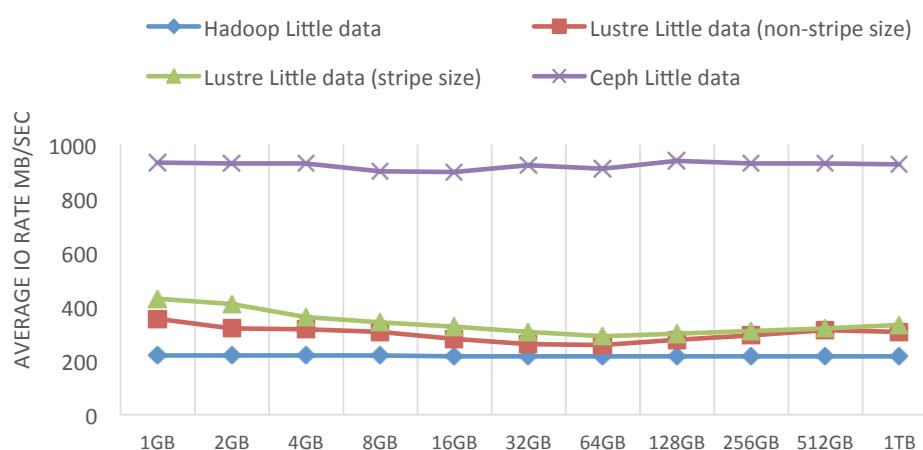
## BIG DATA WRITE



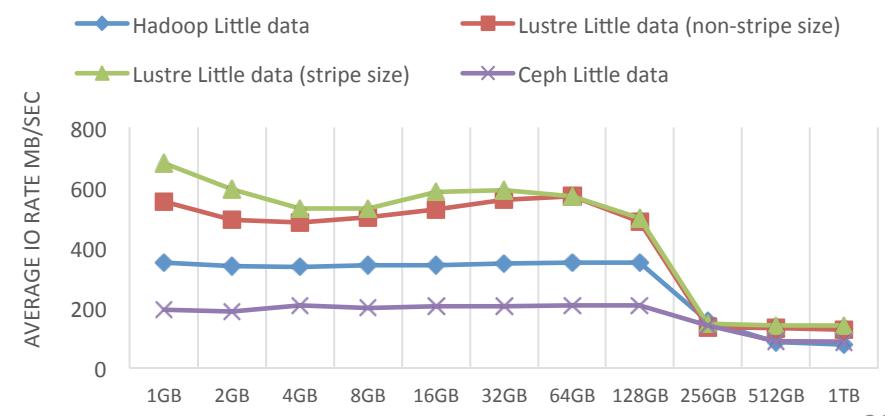
## BIG DATA READ



## LITTLE DATA WRITE

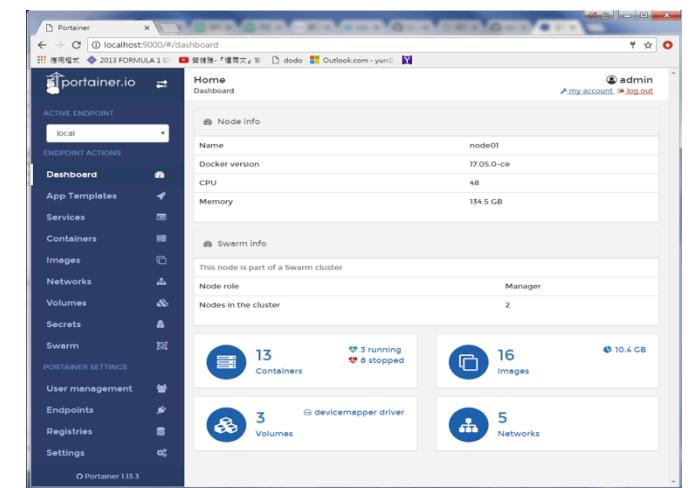
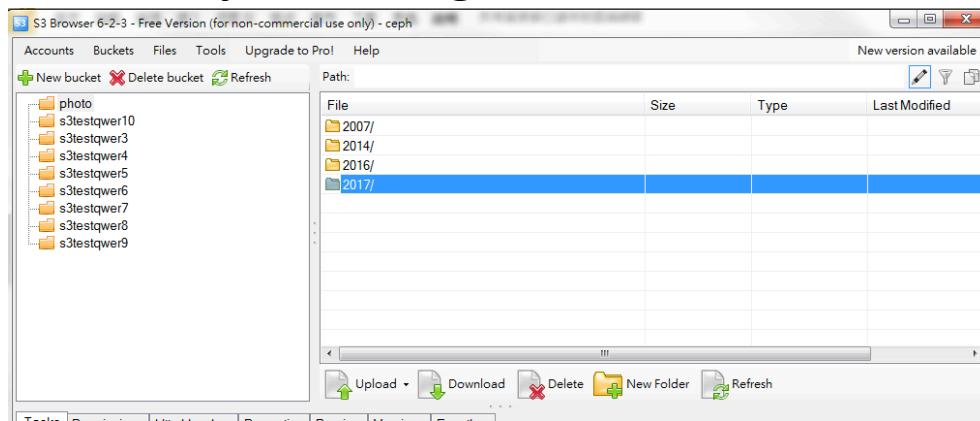
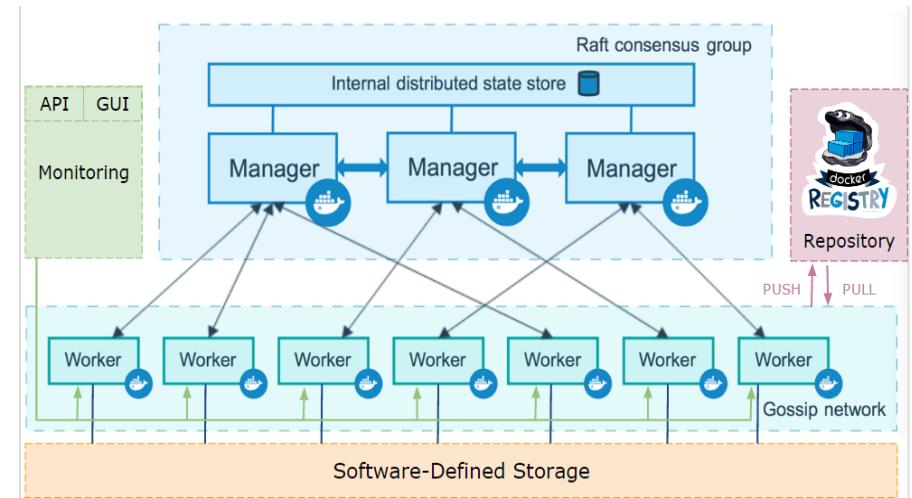


## LITTLE DATA READ



# Platform & Service

- Application Repository
  - Virtualized Data Analysis Application Sharing
  - Container based
    - Swarm as Manager
    - Jupyter for Interface
    - Example : GATE image
  - Ceph File System
    - Block Device for DataMart
    - Object Storage : DataMart Intererfaing



# Data Marketplace Initiative



The screenshot displays the homepage and a detailed view of a specific dataset from the Data Marketplace.

**Homepage Features:**

- Big Data Analytics:** A large image showing two people working on a computer screen with the text "big data analytics".
- Search Bar:** A search bar with placeholder text "搜尋資料" and categories: 熱門標籤 (Hot Tags), 主計處 (Statistics Bureau), 公共資訊 (Public Information), and 環保 (Environmental).
- Statistics:** Three circular icons showing 24.4k 資料集 (Datasets), 22 組織 (Organizations), and 10 群組 (Groups).
- Navigation:** Links for 登入 (Login) and 註冊 (Register) at the top right, and 資料集 (Datasets), 組織 (Organizations), 群組 (Groups), and 關於 (About) in the top center.
- Content Panels:** Two panels on the left:
  - Acceptance-Rejection Method:** Describes a sampling method from a uniform distribution.
  - Crime survey:** States that no datasets are available.
  - Police survey:** States that no datasets are available.

**Detailed Dataset View:**

**Path:** 首頁 / 組織 / 湖泊與環境監測 / LASS環境感測資料庫\_2017Q1

**Header:** LASS環境感測資料庫\_2017Q1

**Statistics:** 0 位業者 (Businesses), 0 組織 (Organizations).

**Thumbnail:** A small image of a landscape with a small structure in the foreground.

**Category:** 湖泊與環境監測 (Lakes and Environmental Monitoring)

**Description:** 湖泊監測與環境品質監測 (Lake Monitoring and Environmental Quality Monitoring)

**Links:** 取得更多 (Get More), 探索 (Explore) for each item.

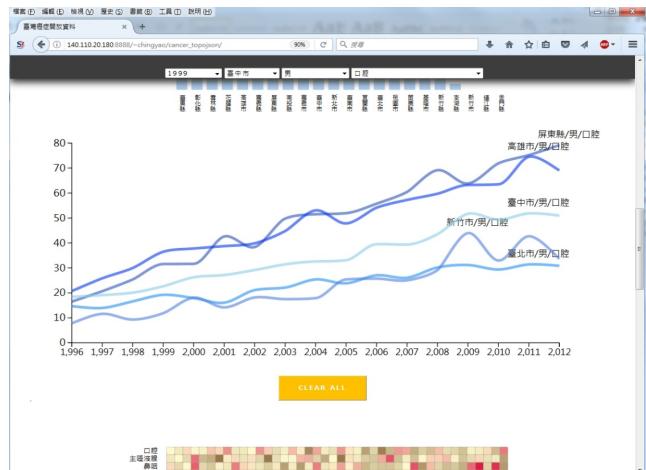
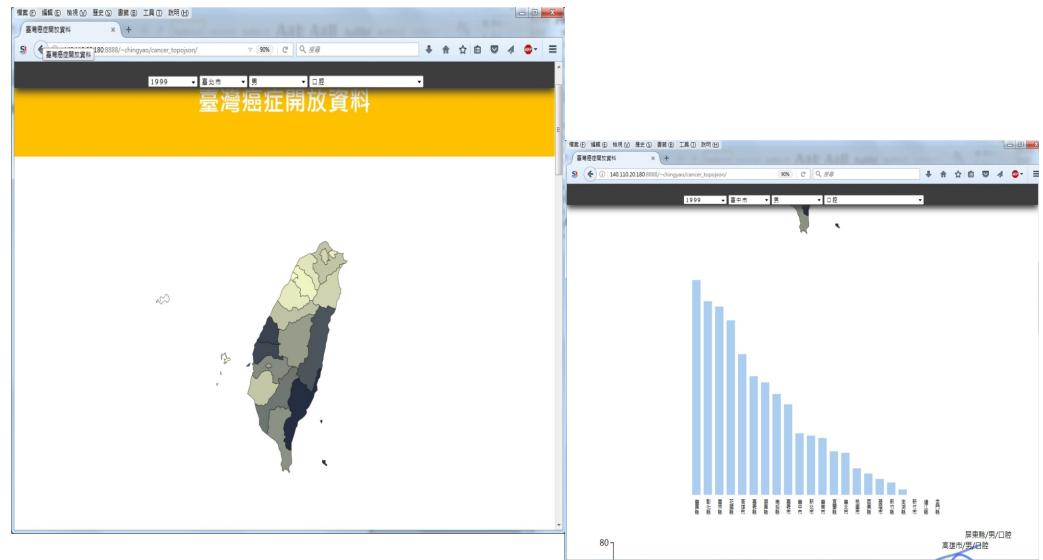
**Items:** A list of 11 JSON files from January 1, 2017, to January 11, 2017, each with a preview icon and a "探索" (Explore) button.

**Heatmap:** A large heatmap visualization at the bottom right, likely representing spatial data from the dataset.

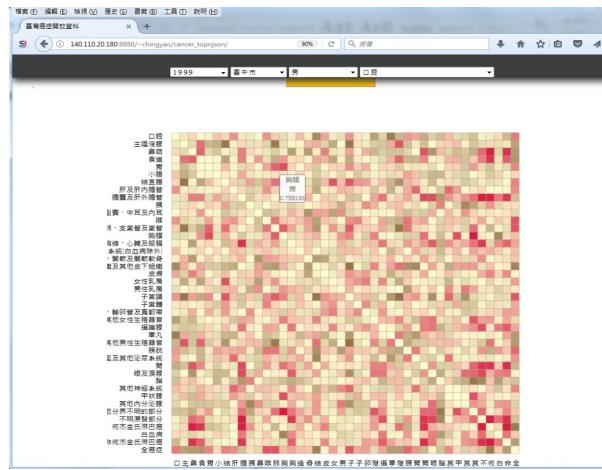
Most data : public available  
Restricted sensitive data : ~ 5TB, 5 billions records

# Data Application Example

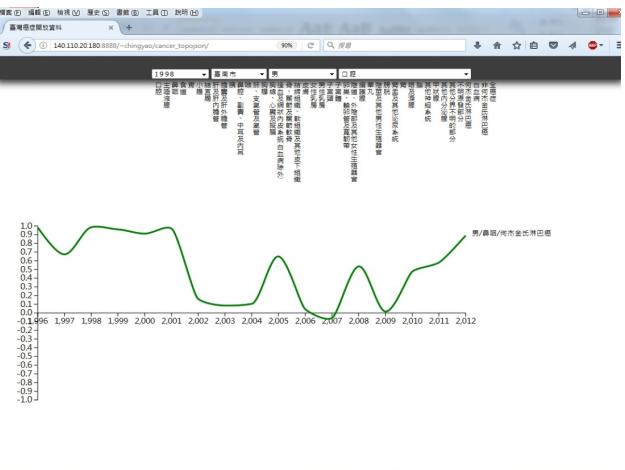
- Cancer Open Data
  - 1996 ~ 2012
  - Pierson Analysis



Geographical Distribution of Cancer



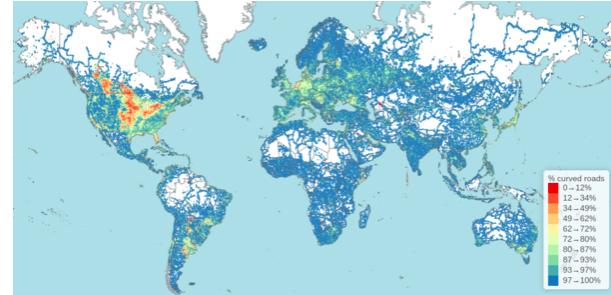
Pierson Analysis : relation between cancers



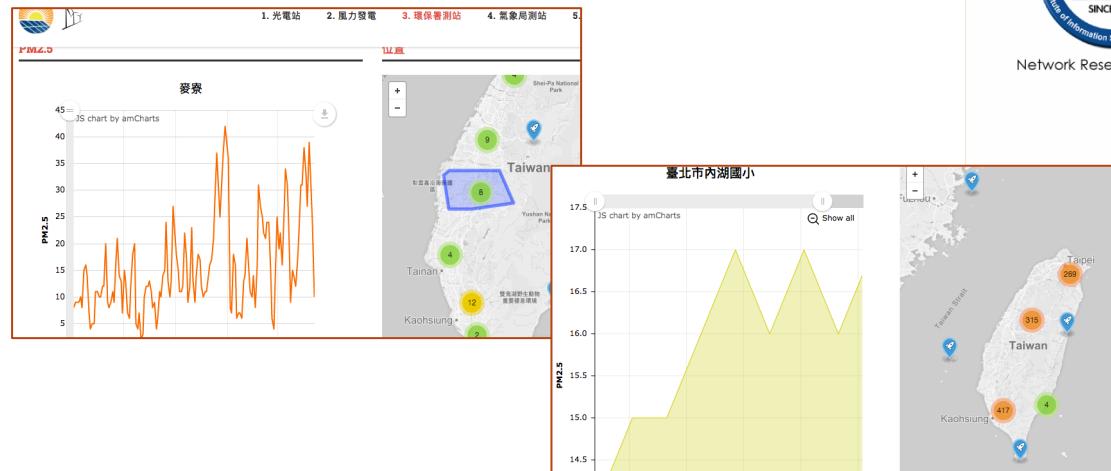
Pierson Analysis : 1995 - 2012

# Data Application Example

- OpenStreetMap(OSM)
  - The only mirror site in Asia
  - Integrated with Overpass turbo, Umap, Planet
- Location Aware Sensor System (LASS)
  - AirBox for PM2.5



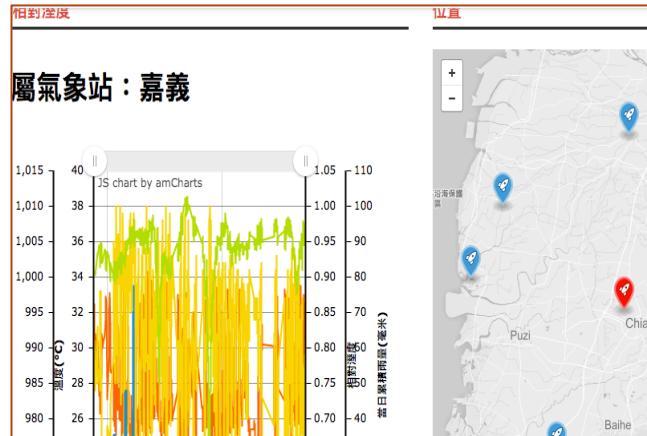
\* Image source: <http://www.openstreetmap.org>



: <http://research.sinica.edu.tw/pm25-air-box/>

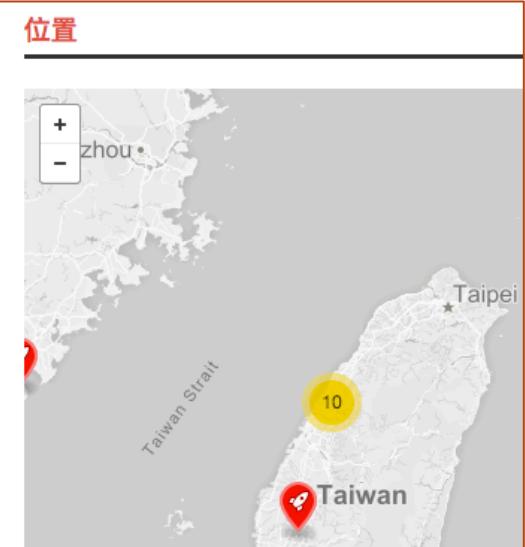
# Data Application Example

- Renewable Energy
  - Solar & Wind



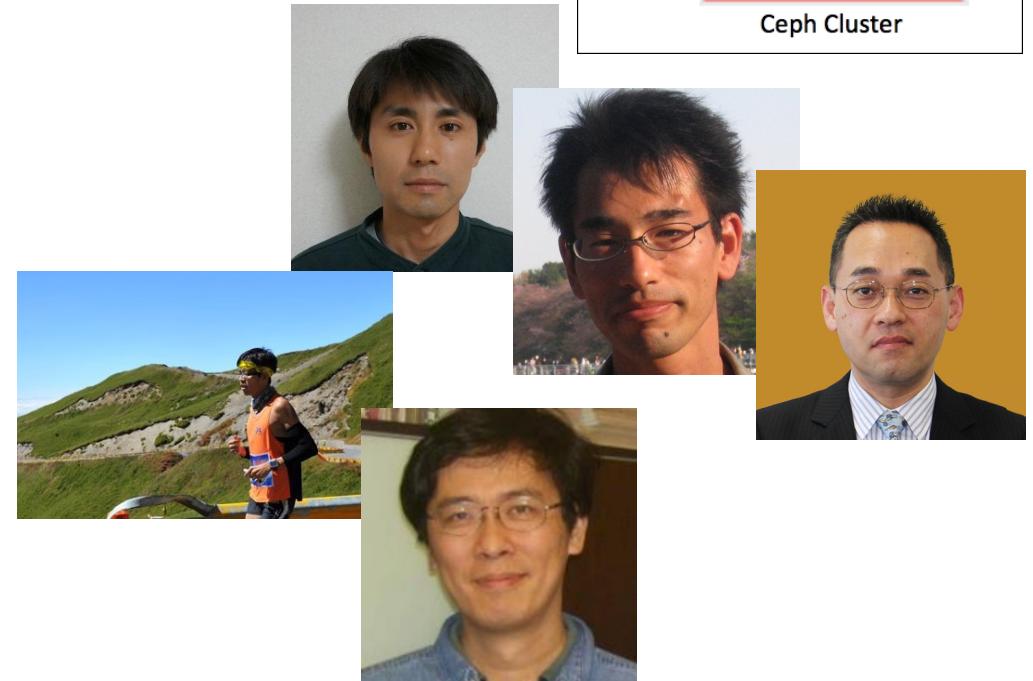
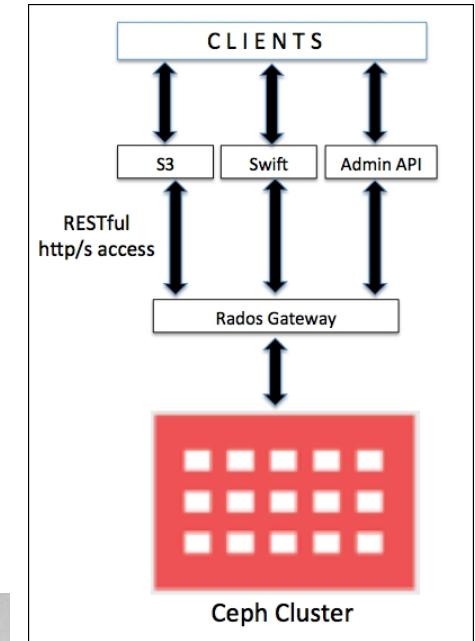
Solar energy

Wind energy

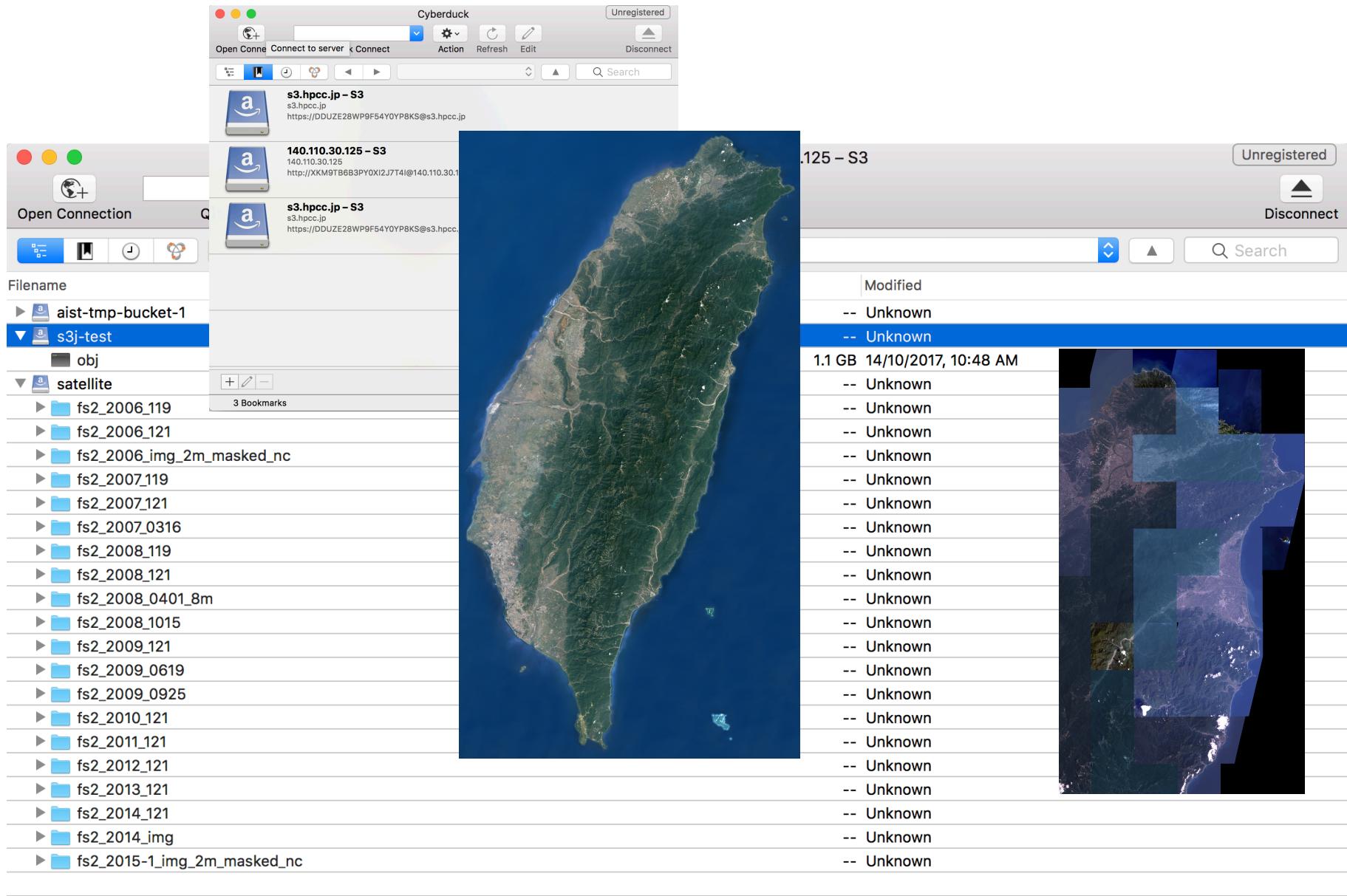


# Out Reach

- Data Cloud initiative
  - International Collaboration
    - AIST Japan
      - Yusuke Tanimura
      - Ryousei Takano
      - Yoshio Tanaka
    - NCHC Taiwan
      - Max Yu
      - Weicheng Huang
    - Ceph, S3

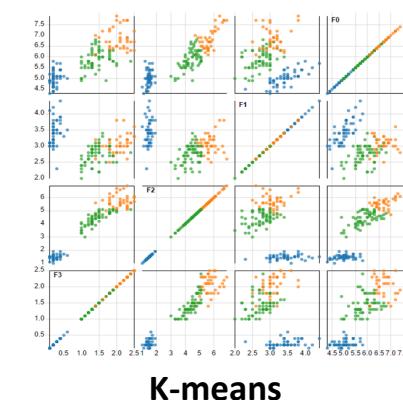
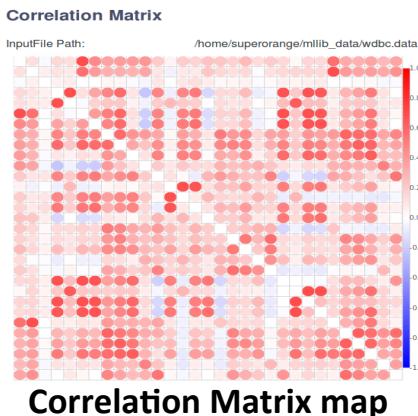


# Out Reach

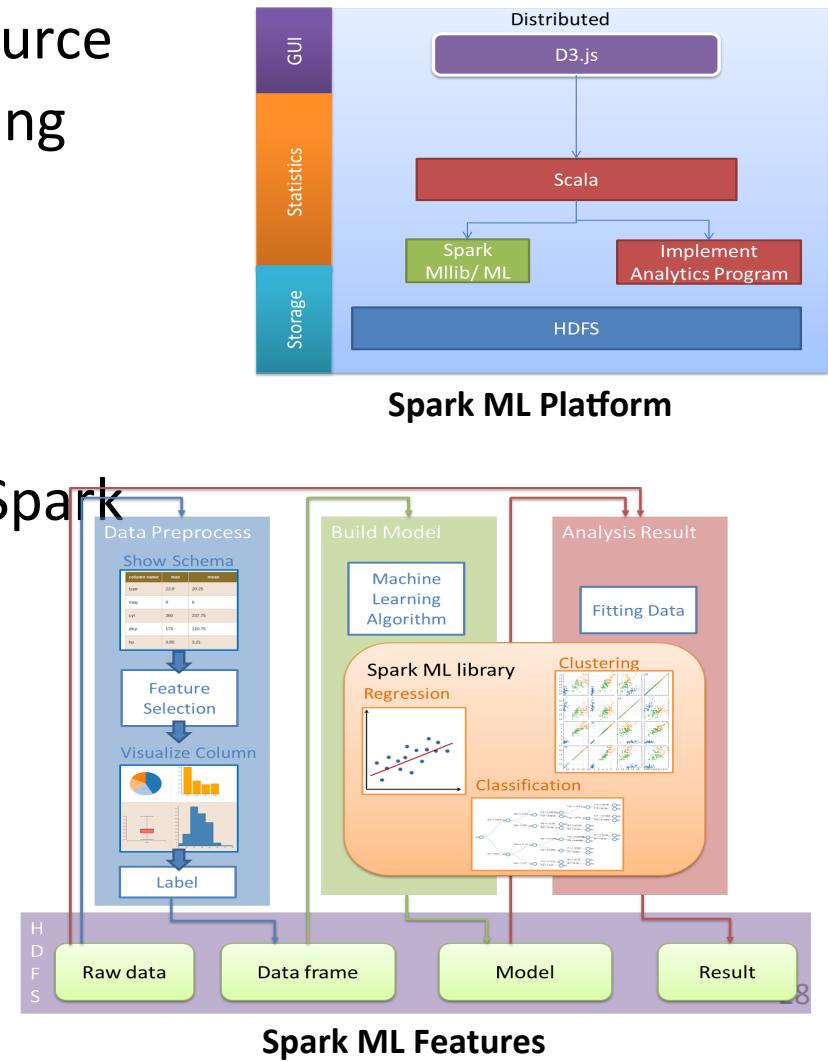


# Spark for Machine Learning

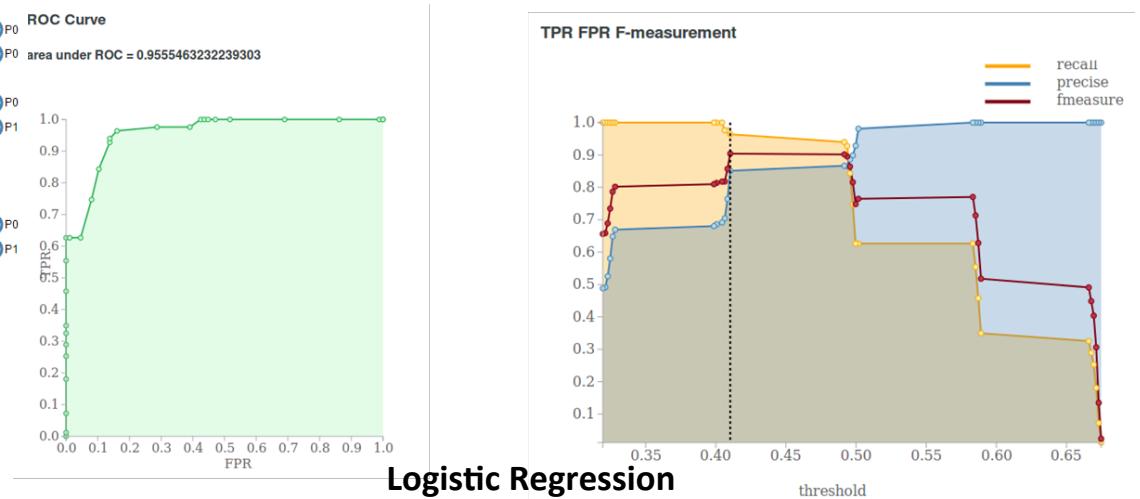
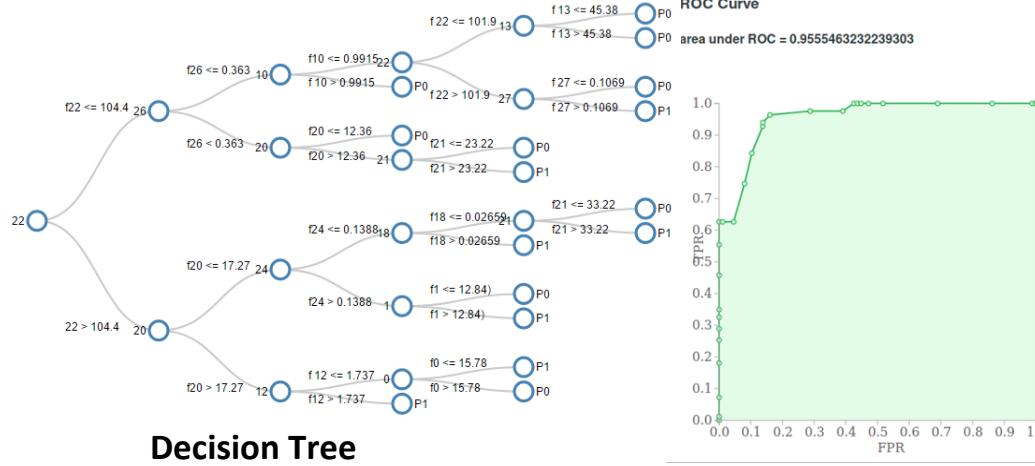
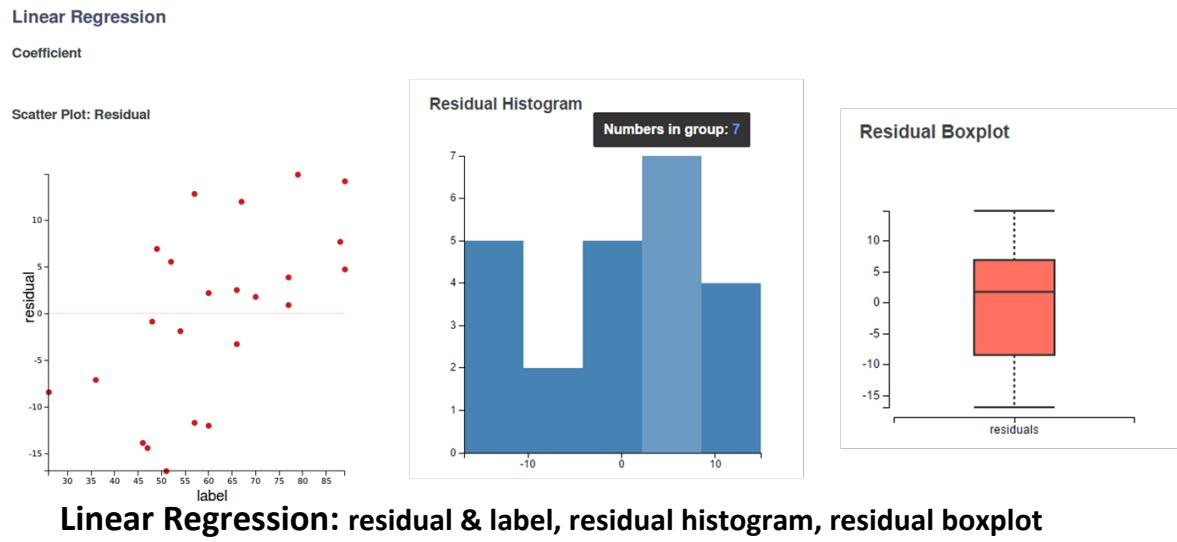
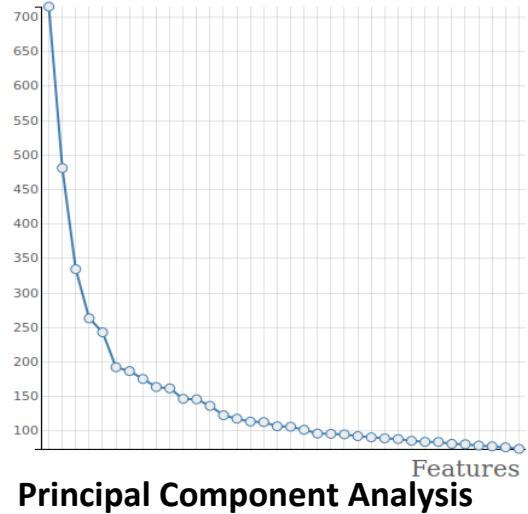
- Data Preprocessing, Analysis w/UI
  - Analysis capability via Open Source
  - ML process w/parallel processing
  - Features
    - Data Visualization : 15
    - Data Pre processing : 5
    - ML model : 4
  - To lower the barrier for using Spark
  - Better performance



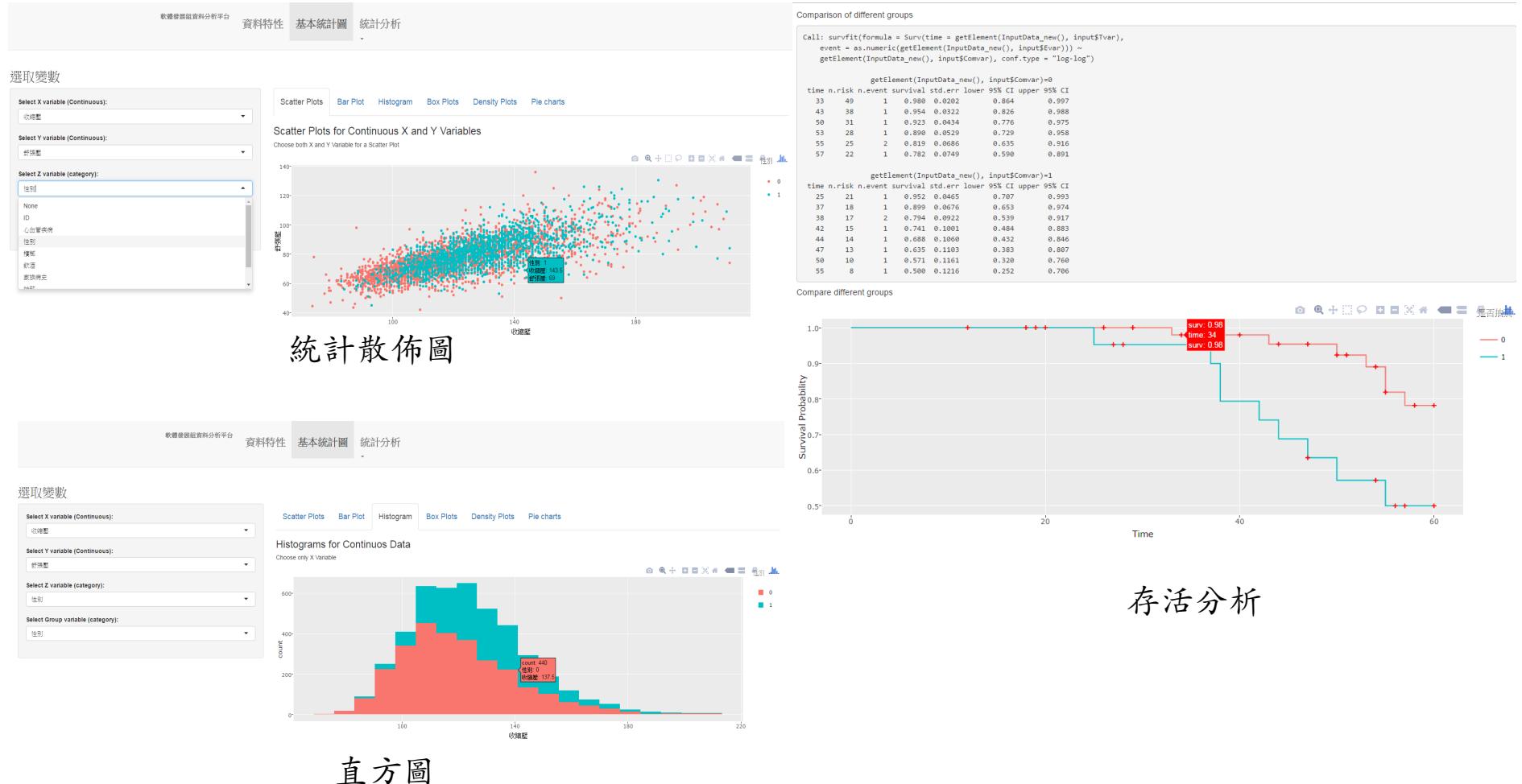
K-means



# Spark – ML w/Visualization



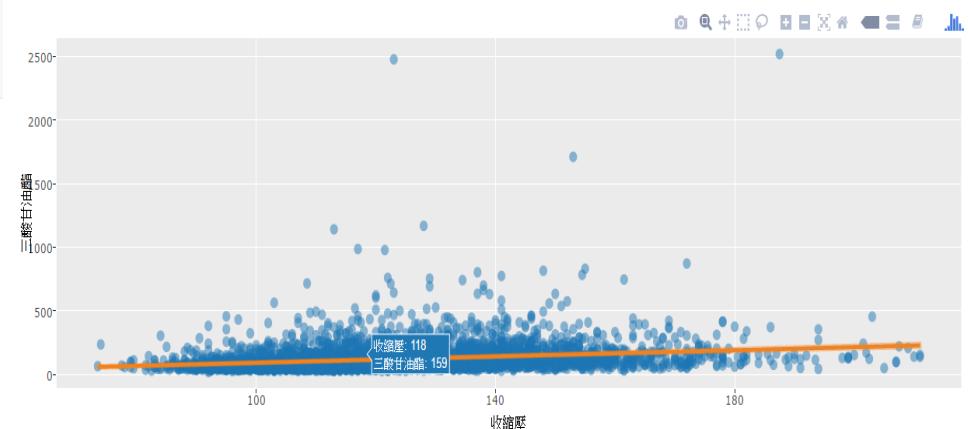
# Web-based Data Analysis Platform



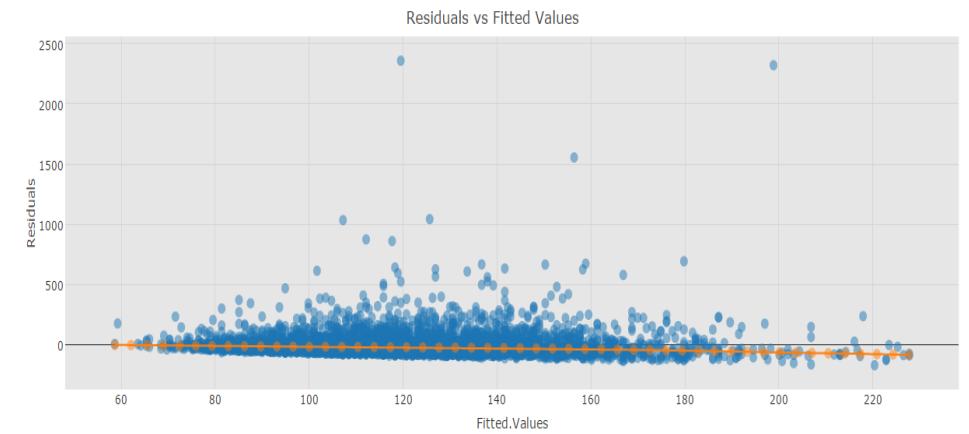
# Web-based Data Analysis Platform



配適曲線圖(當獨立變數僅有一個時才會呈現)



殘差與配適值圖



## 迴歸分析

# Web-based Data Analysis Platform

- Data Cleaning : 3
- Data Integration
- Statistics tool kit : 8
- Statistics chats : 5
- Text mining : 4
- API w/R & Python

