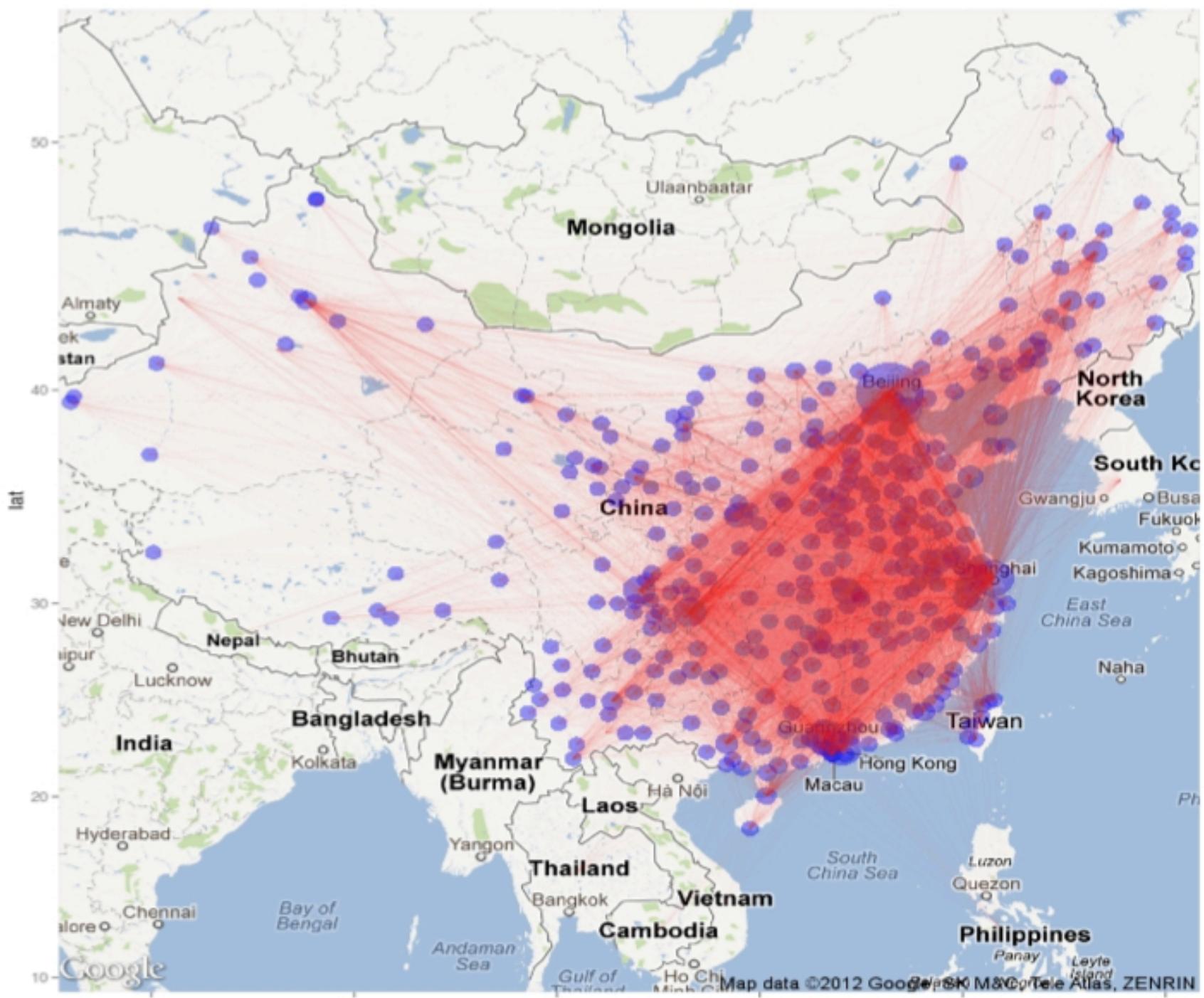


# Big Data: An Introduction

Qiang Yang  
Chair Professor, HKUST  
<http://www.cse.ust.hk/~qyang>





# Baihe.com

百合 baihe.com  
实名婚恋网开创者

9000万  
诚意会员  
实名婚恋网站开创者  
每天数千对牵手幸福

手机号注册 牵手更快

账号

密码

已经阅读并同意 [百合服务条款](#)

[免费注册](#)

第三方账号登录

- 每年上亿用户
- 5-10短信／天
- 日短信量:千万
- 2012年成功:50万对

# What Big Data Means

## Big Data: Volume, Velocity, Variety

### Social

**Facebook:** 1 Billion Users  
daily increment: 10 TB

### Search

**Google:** Daily 20 PB

**Turing Award Winner: Jim Gray: Fourth Paradigm**

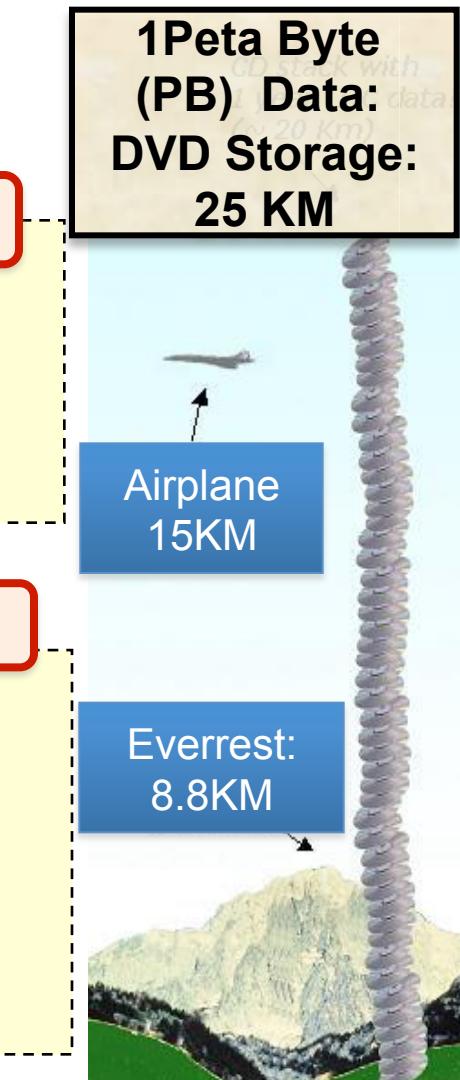
• Doubles every 18 months

• IDC: in 2020, will reach  $35 \times 10^6$  PB

1Peta Byte  
(PB) Data:  
DVD Storage:  
25 KM

Airplane  
15KM

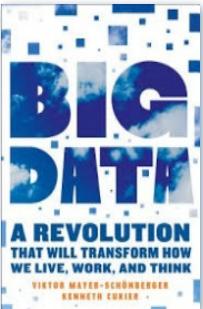
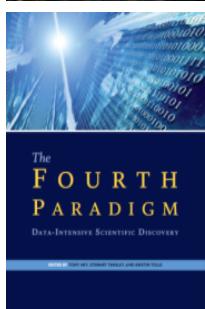
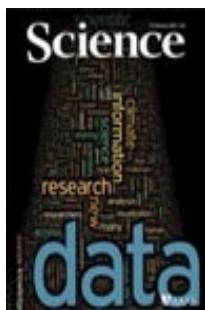
Everest:  
8.8KM



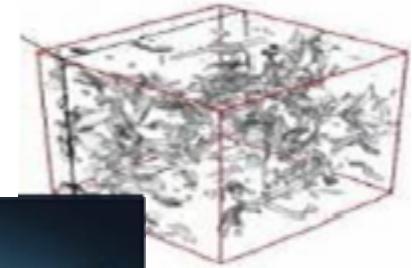
# The Fourth Paradigm

- Thousand Years Ago
  - Empirical observation
- Hundreds Years Ago - Theory
  - Newton Law, Maxwell ...
- Decades Ago - Simulation
  - Complex Phenomena

Today — Data Science  
Causal Relation → Statistical Correlation



$$\left(\frac{\dot{a}}{a}\right)^2 = \frac{4\pi G\rho}{3} - K \frac{c^2}{a^2}$$



From: <The Fourth Paradigm>

# Data Analytics Tasks

## Query Data

- Clustering (聚类)
- Association (关联)
- Segmentation (分群)
- Topic Analysis (主题)

## Prediction

- Regression (回归预测)
- Classification (分类)
- Recommendation (推荐)

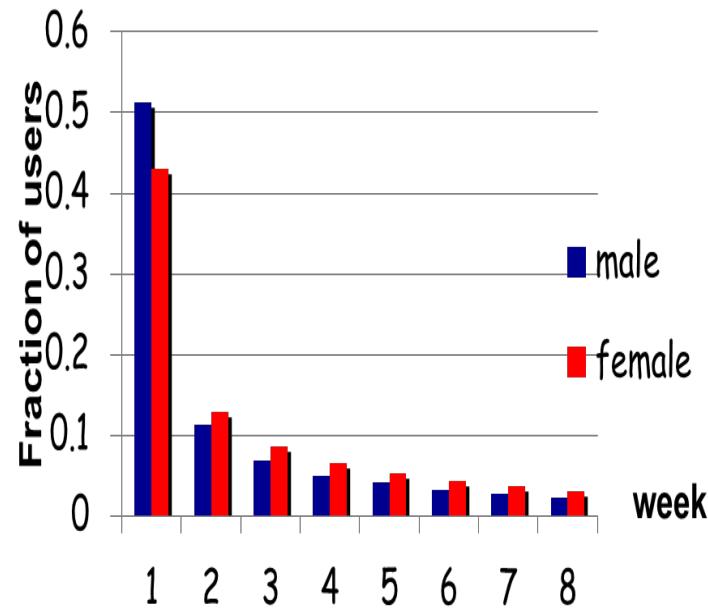
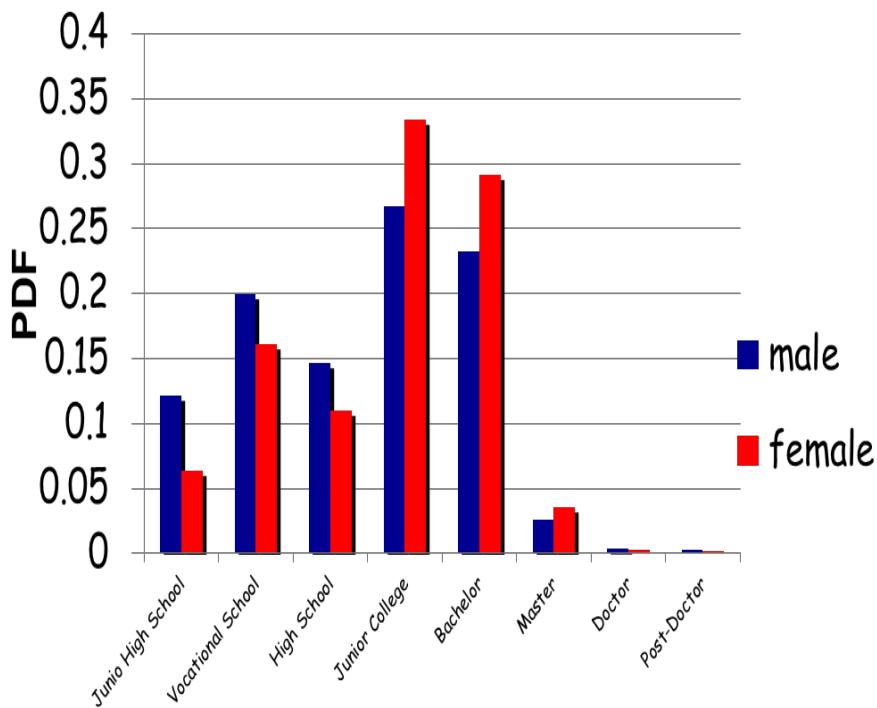
# Data Analytics Tasks

## Query Data

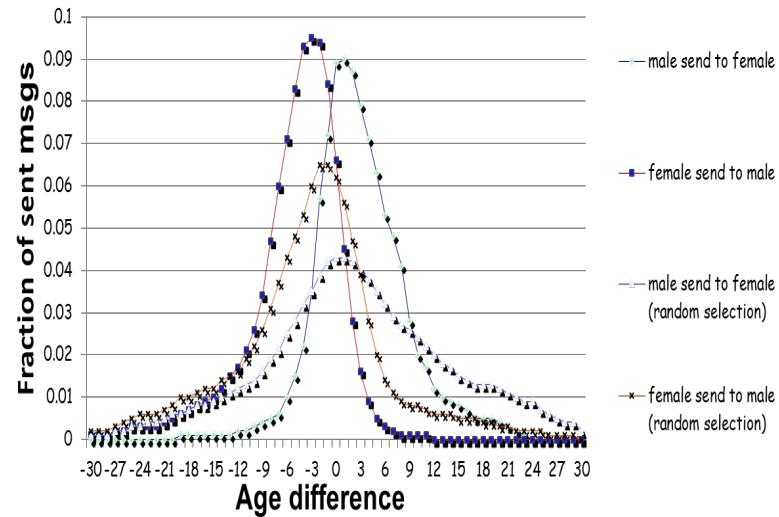
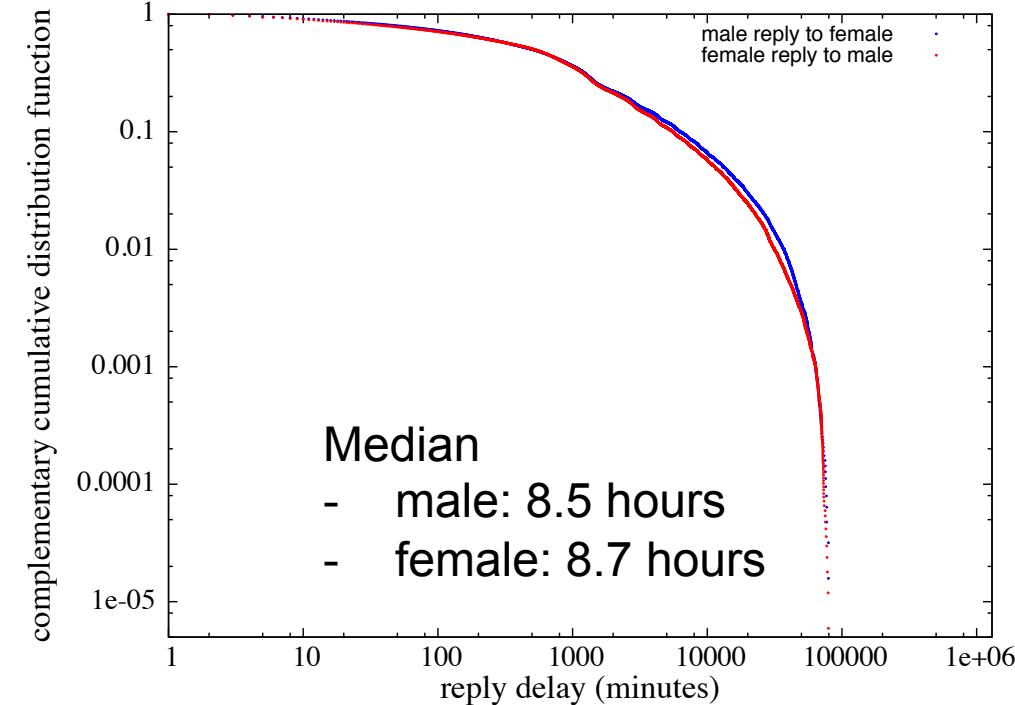
- Summary (汇总)
- Clustering (聚类)
- Association (关联)
- Segmentation (分群)
- Topic Analysis (主题)

# Who Date Online?

- Relevant to Education Level?
- How patient are they?

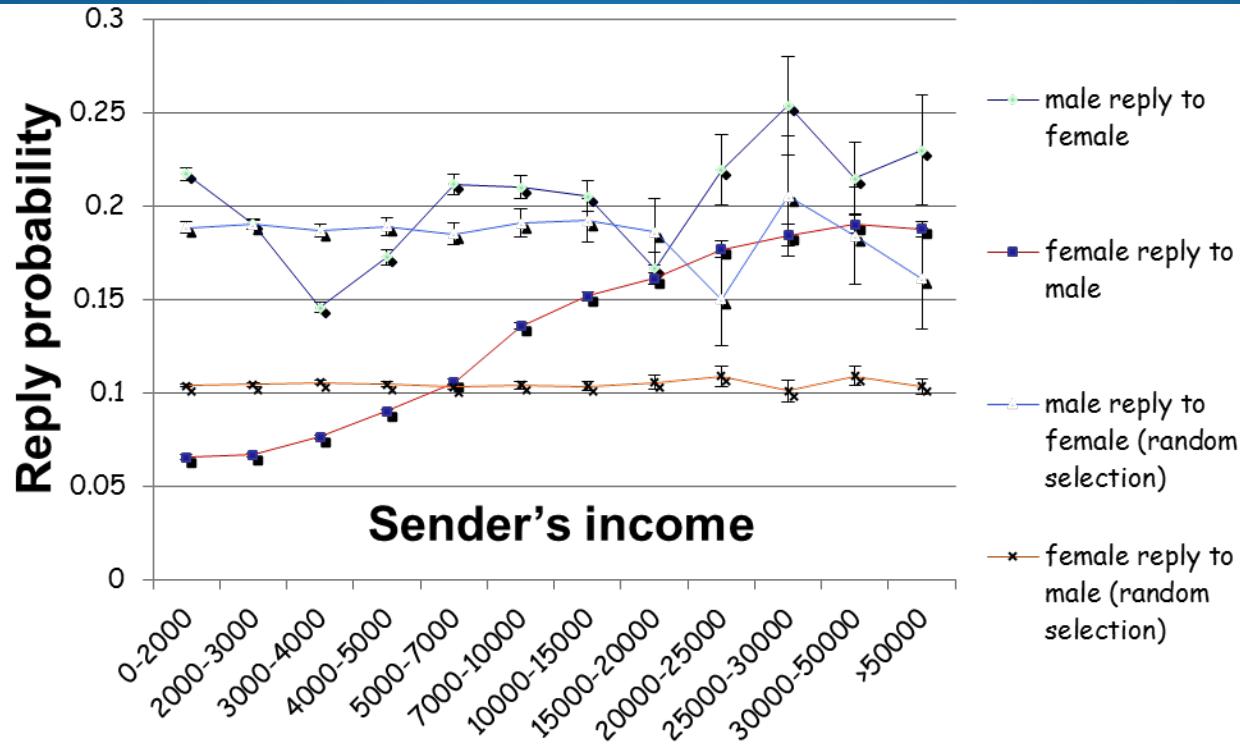


# How Long Love Waits?



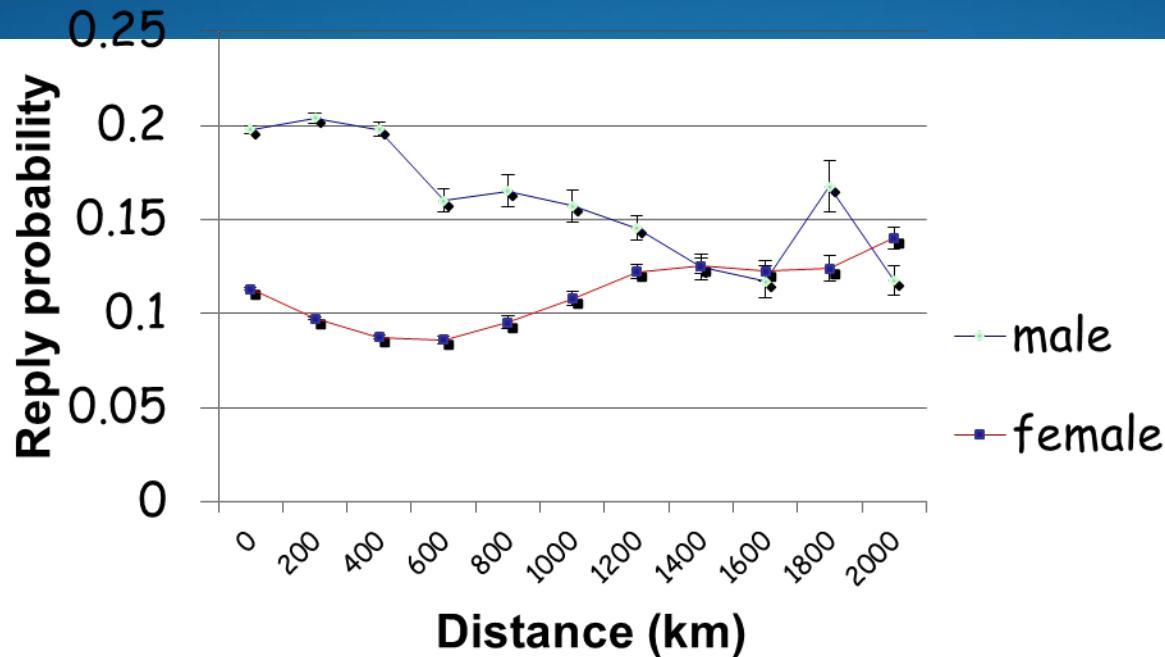
Males tend to look for younger females and females tend to look for older males

# Does \$\$ Matter in Love?



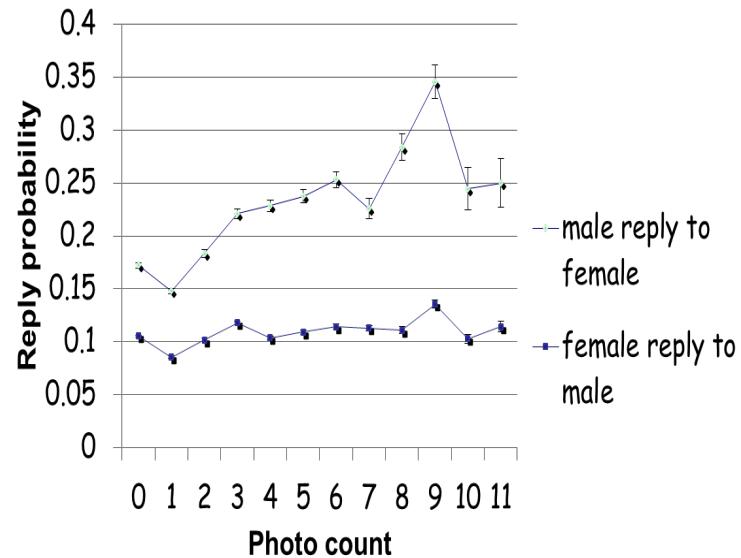
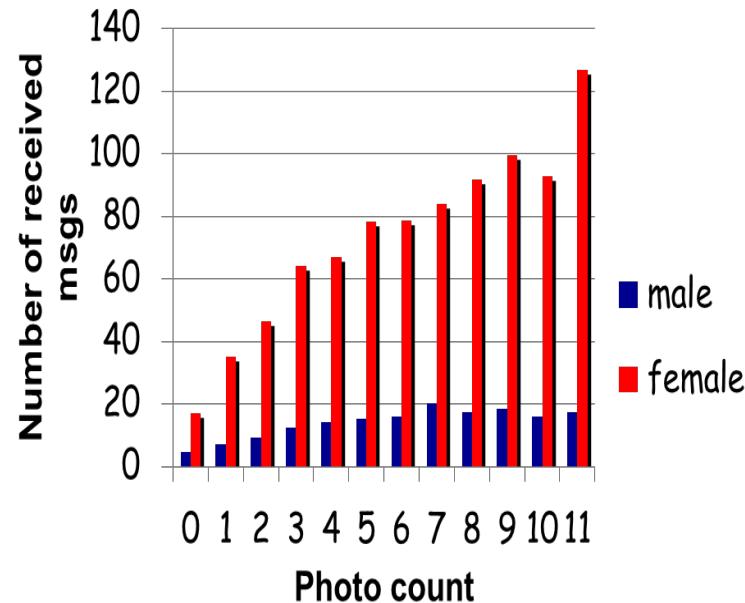
- Female reply probability increases with male senders' income,
- Males don't care much about the income of a female

# Does Distance Matter?



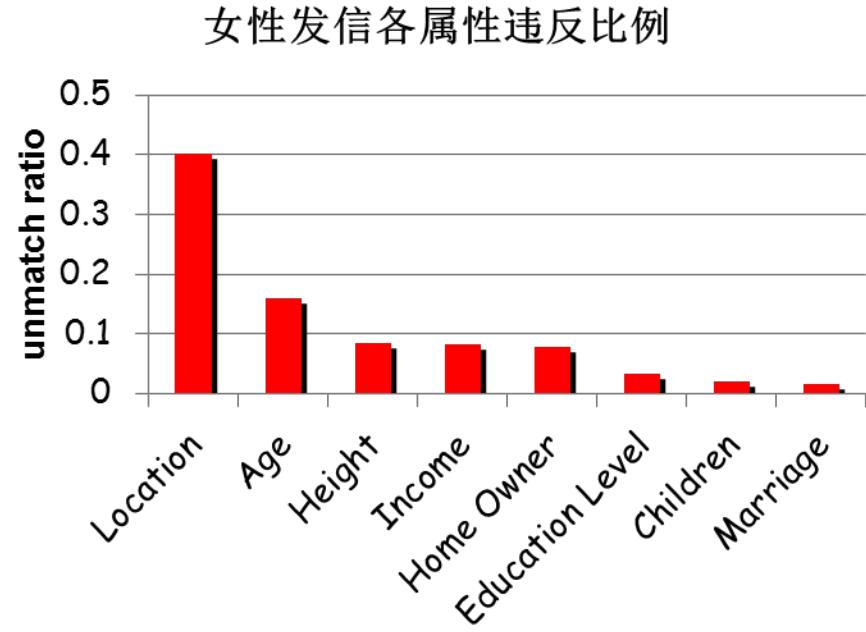
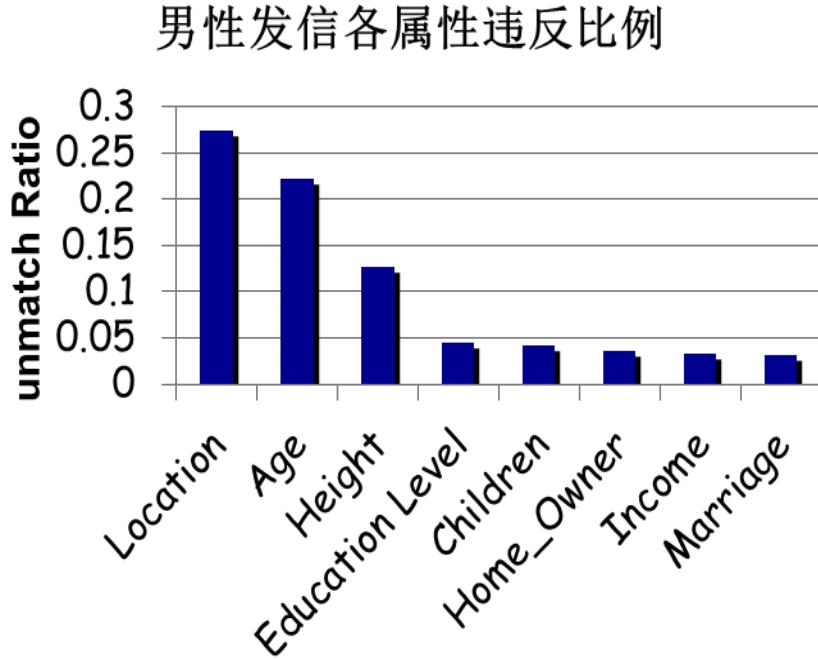
- Male reply probability decreases with distance
- Female reply probability first decreases with distance but increases in the range from 800 to 1,400km

# Are Looks important?



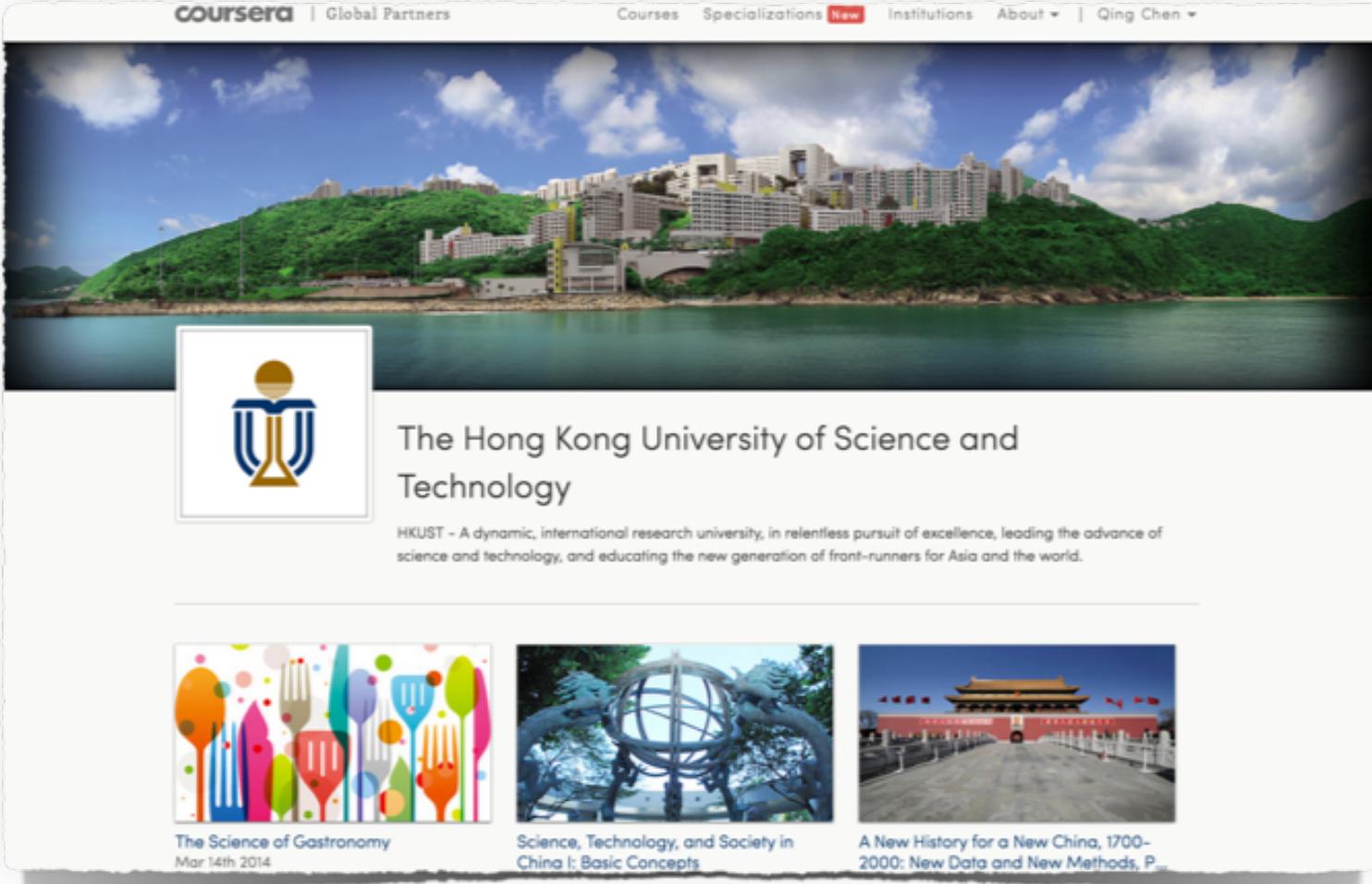
- One receives more msgs if posted more photos
- More so for females

# Does your heart follow your mind?



- Both male and females do not care about: age, location and height, contrary to their stated preferences
- Females are the most strict with marriage and children status, as well as education level of male senders.

# HKUST on Coursera



The screenshot shows the Coursera website interface for HKUST. At the top, the Coursera logo and "Global Partners" are visible, along with navigation links for "Courses", "Specializations New", "Institutions", "About", and a user profile for "Qing Chen". The main banner features a scenic view of the HKUST campus buildings perched on a hillside overlooking a body of water under a blue sky with white clouds. Below the banner, the HKUST logo (a stylized figure composed of vertical bars) is displayed next to the university's name: "The Hong Kong University of Science and Technology". A brief description follows: "HKUST - A dynamic, international research university, in relentless pursuit of excellence, leading the advance of science and technology, and educating the new generation of front-runners for Asia and the world." Three course thumbnails are shown at the bottom: "The Science of Gastronomy" (Mar 14th 2014), "Science, Technology, and Society in China I: Basic Concepts", and "A New History for a New China, 1700–2000: New Data and New Methods, P...".

**coursera | Global Partners**

Courses Specializations **New** Institutions About | Qing Chen

 The Hong Kong University of Science and Technology

HKUST – A dynamic, international research university, in relentless pursuit of excellence, leading the advance of science and technology, and educating the new generation of front-runners for Asia and the world.

---

 The Science of Gastronomy  
Mar 14th 2014

 Science, Technology, and Society in China I: Basic Concepts

 A New History for a New China, 1700–2000: New Data and New Methods, P...



# A New History for a New China, 1700-2000

第二回

在逃地主統計表

（續）

## Revolutionary Victims in Shuangcheng and Elsewhere

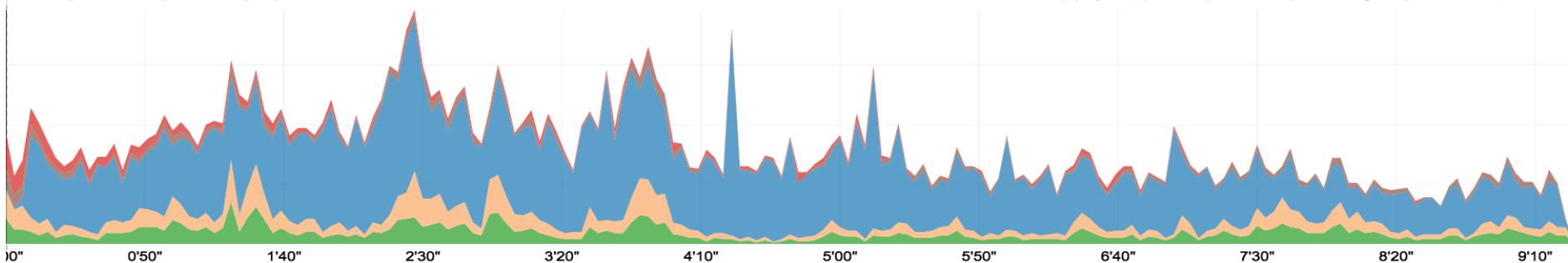
日期	人数	地点	情况
4月20日	70人	哈尔滨	哈中港
3月19日	70人	利州小富人	行踪不明
3月20日	10人	"	"
3月18日	70人	"	"
3月22日	50人	哈后石	哈尔滨
4月28日	50人	哈后石	行踪不明
6月15日	50人	哈肇安庄	"
5年2月5日	100人	哈佳村长被逼至社至回冲人	石峰要霸土地
5年3月20日	和隆村庄、哈佳村庄	哈佳巨聚村布	"
11月同	50人	哈佳四工	哈佳后品庄
每十月间	20人	哈佳三道河	哈佳后品庄
5年3月20日	20人	哈佳三道河	哈佳后品庄
5年3月20日	20人	哈佳三道河	哈佳后品庄

0:00 / 9:30

▶ 🔍

● Stacked ○ Stream ○ Expanded

● play ○ pause ● seeked ● ratechange ● stalled ● error



# Which course has more difficult homework assignments?

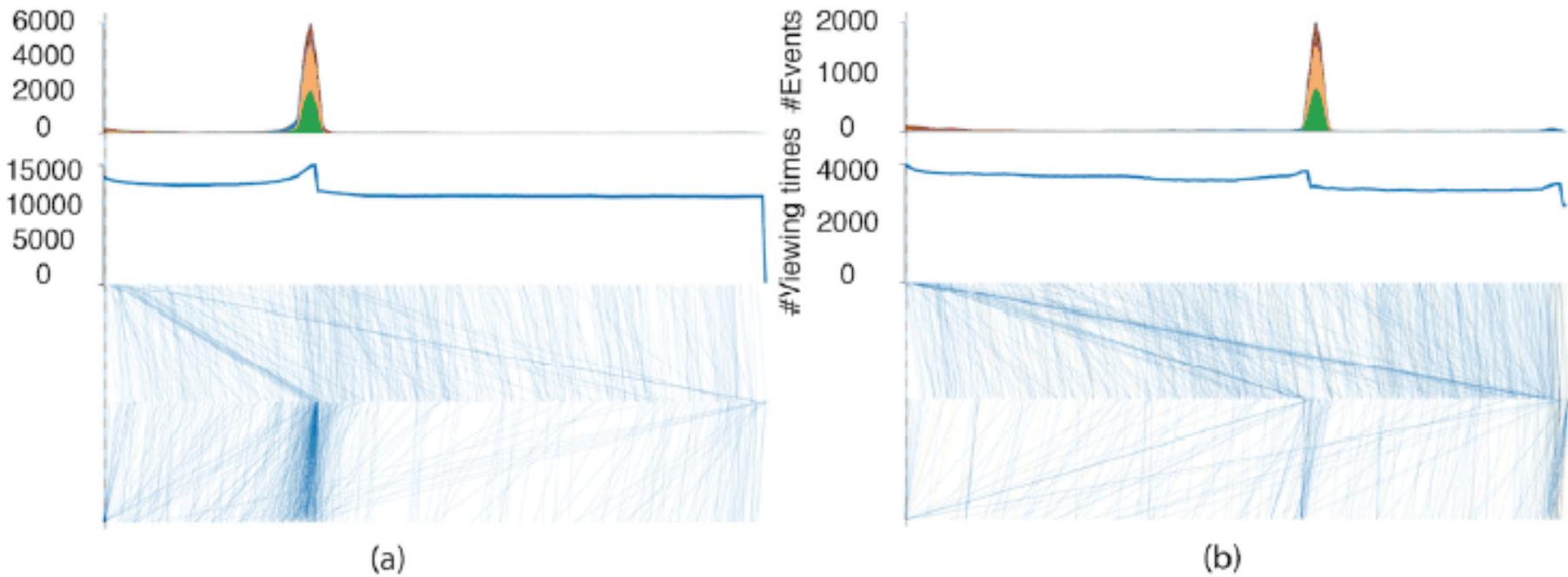


Figure 3: Comparison between the Content-based views of two videos with an in-video question. We can see that: there are a consid-

# US and Chinese Students Behaves Differently?

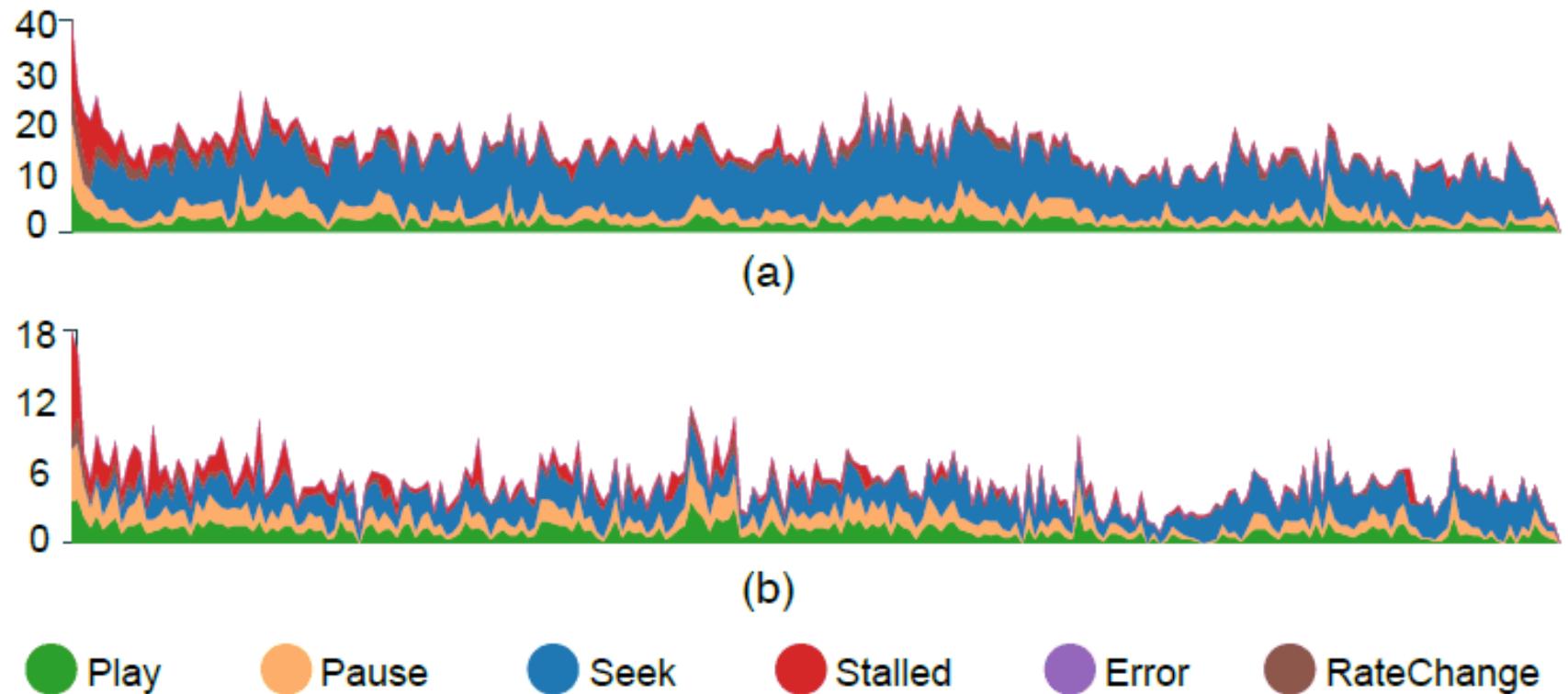
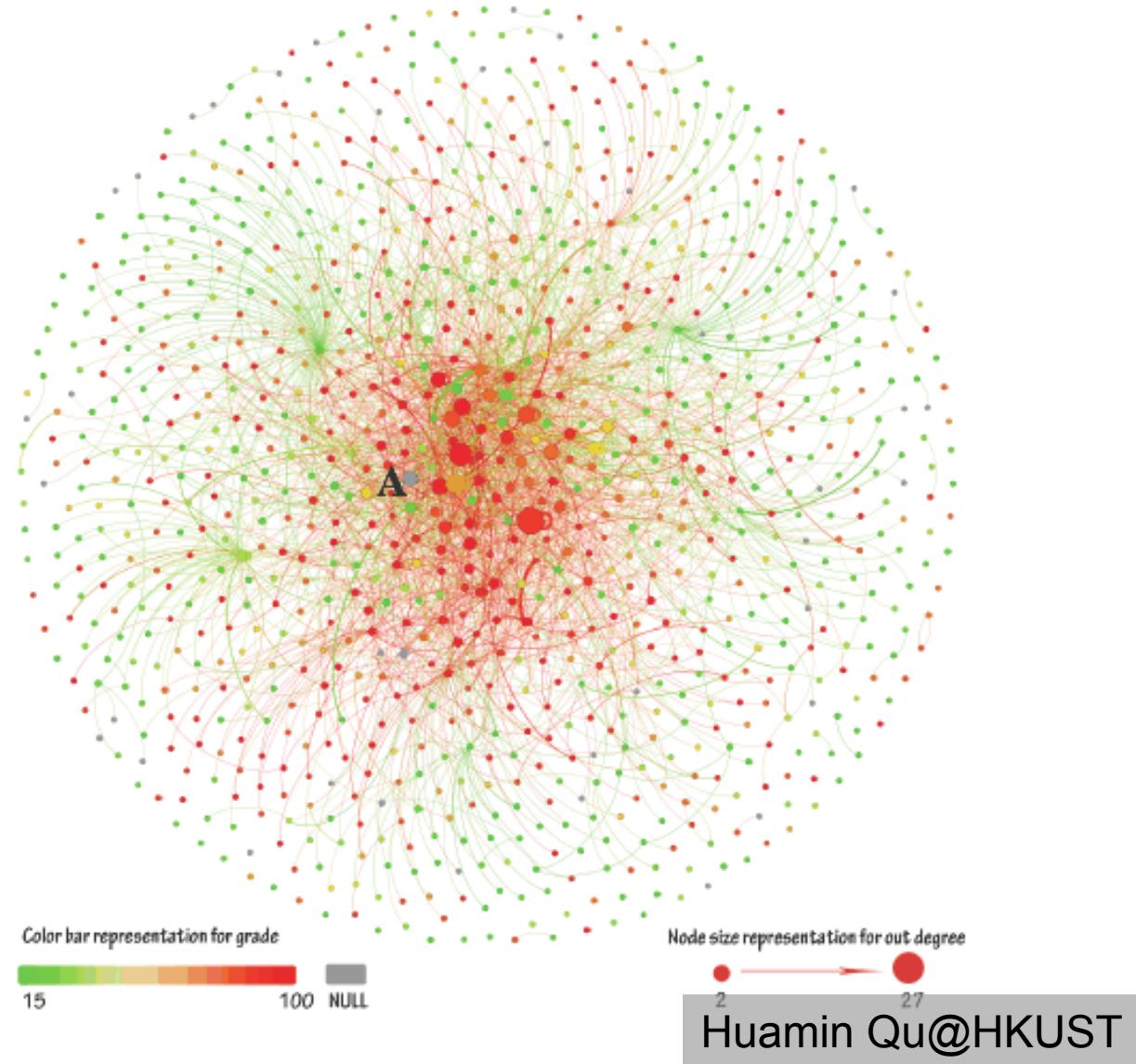
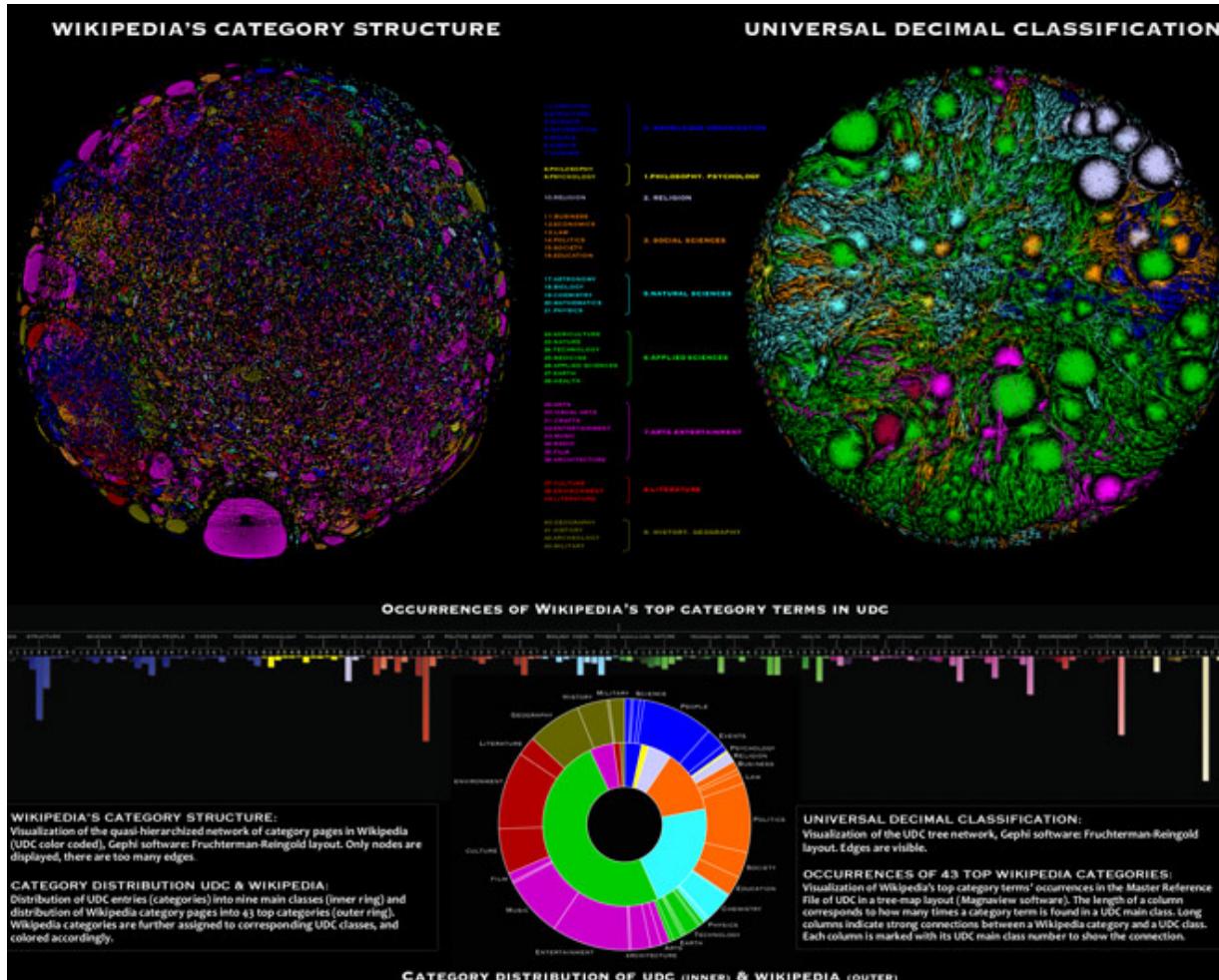


Figure 2: The Event Graphs showing the clickstream data of the same course during the same time period but for learners different countries. a) Learners from the U.S; b) Learners from China.

# Are Good Students More Vocal?

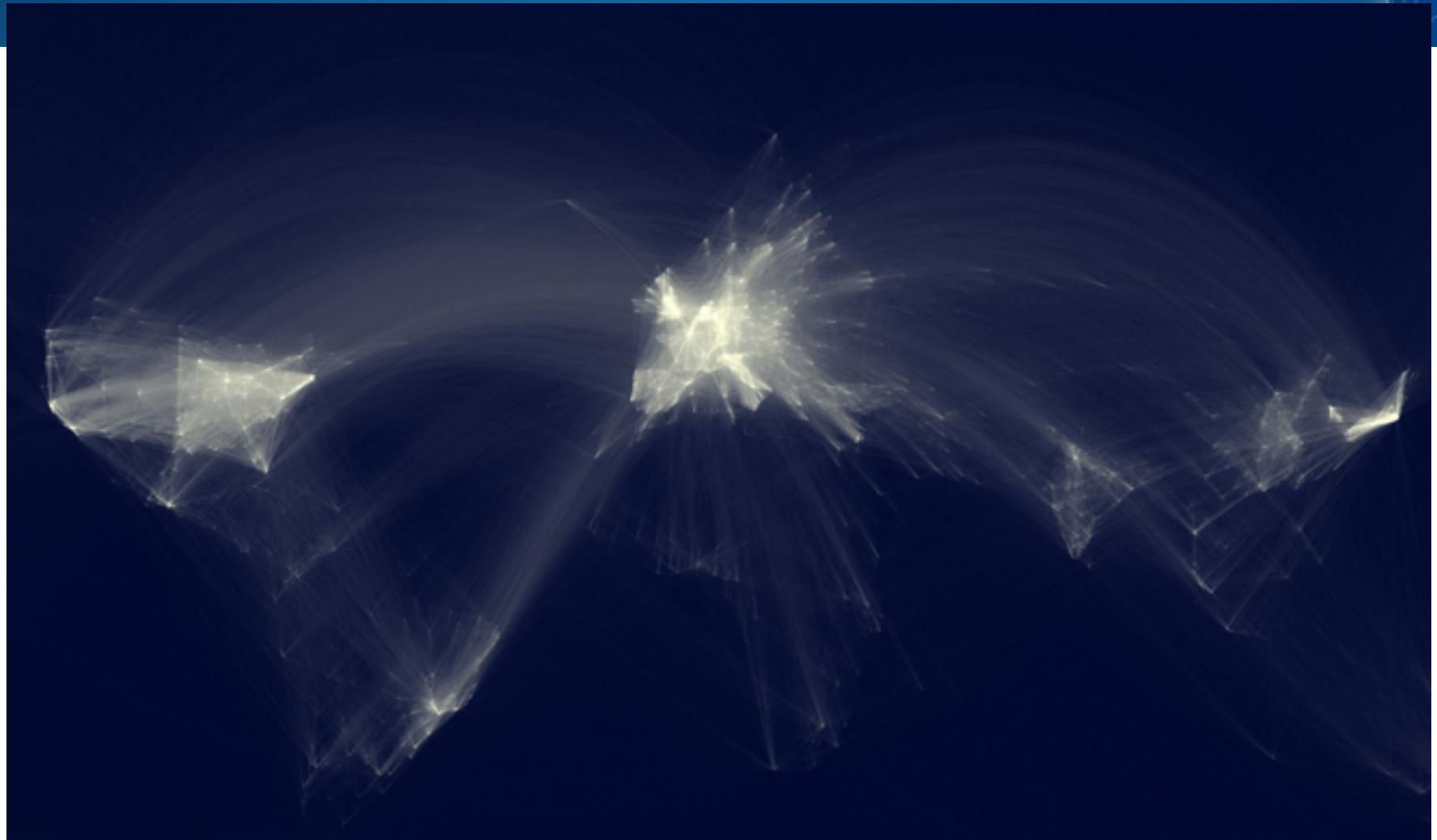


# Are Concepts Different in Wikipedia and Library?



Almila Akdag Salah/Cheng Gao/Krzystof Susecki/Andrea Scharnhorst/  
Den Haag/[Knowledge Space Lab \[high-resolution version\]](#)

# Where Do Scientists Collaborate Most Often?



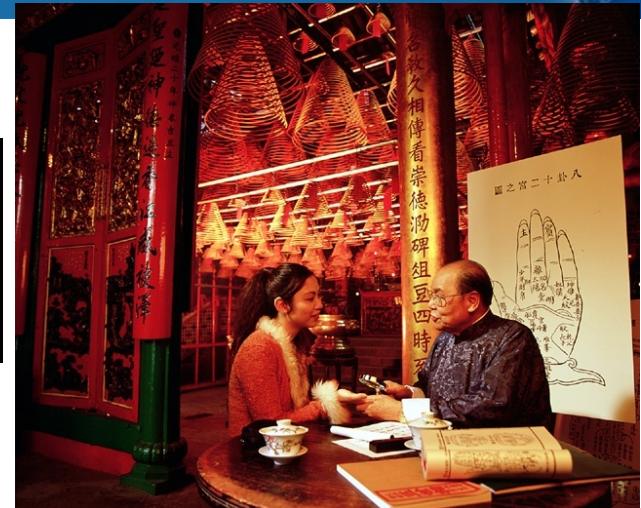
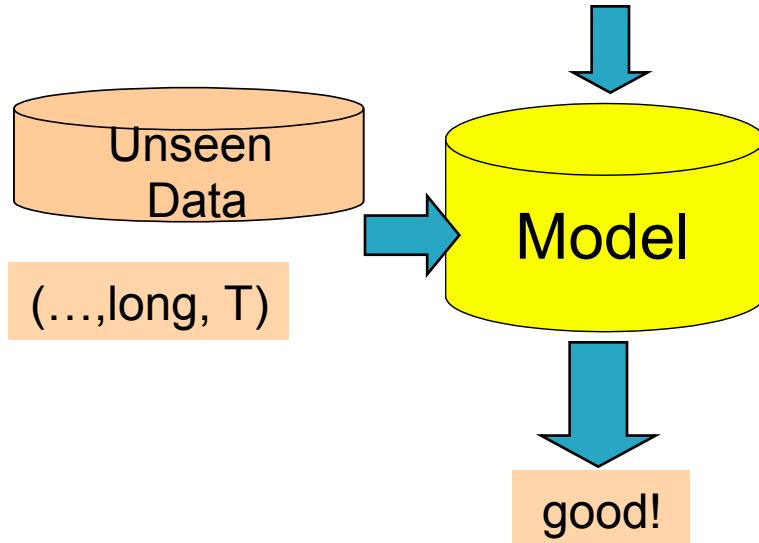
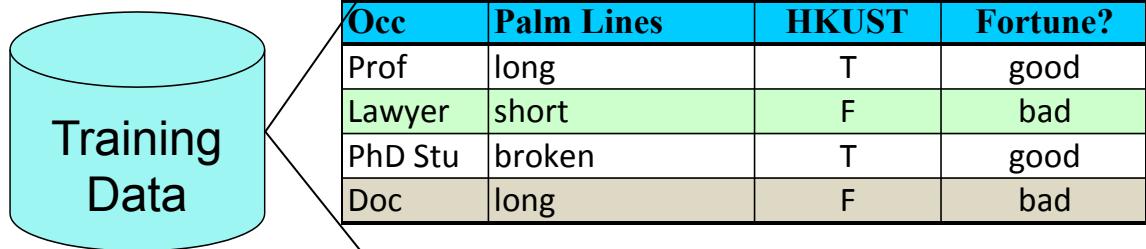
*Image: Olivier H. Beauchesne/Science-Metrix [[high-resolution version](#)]*

# Data Analytics

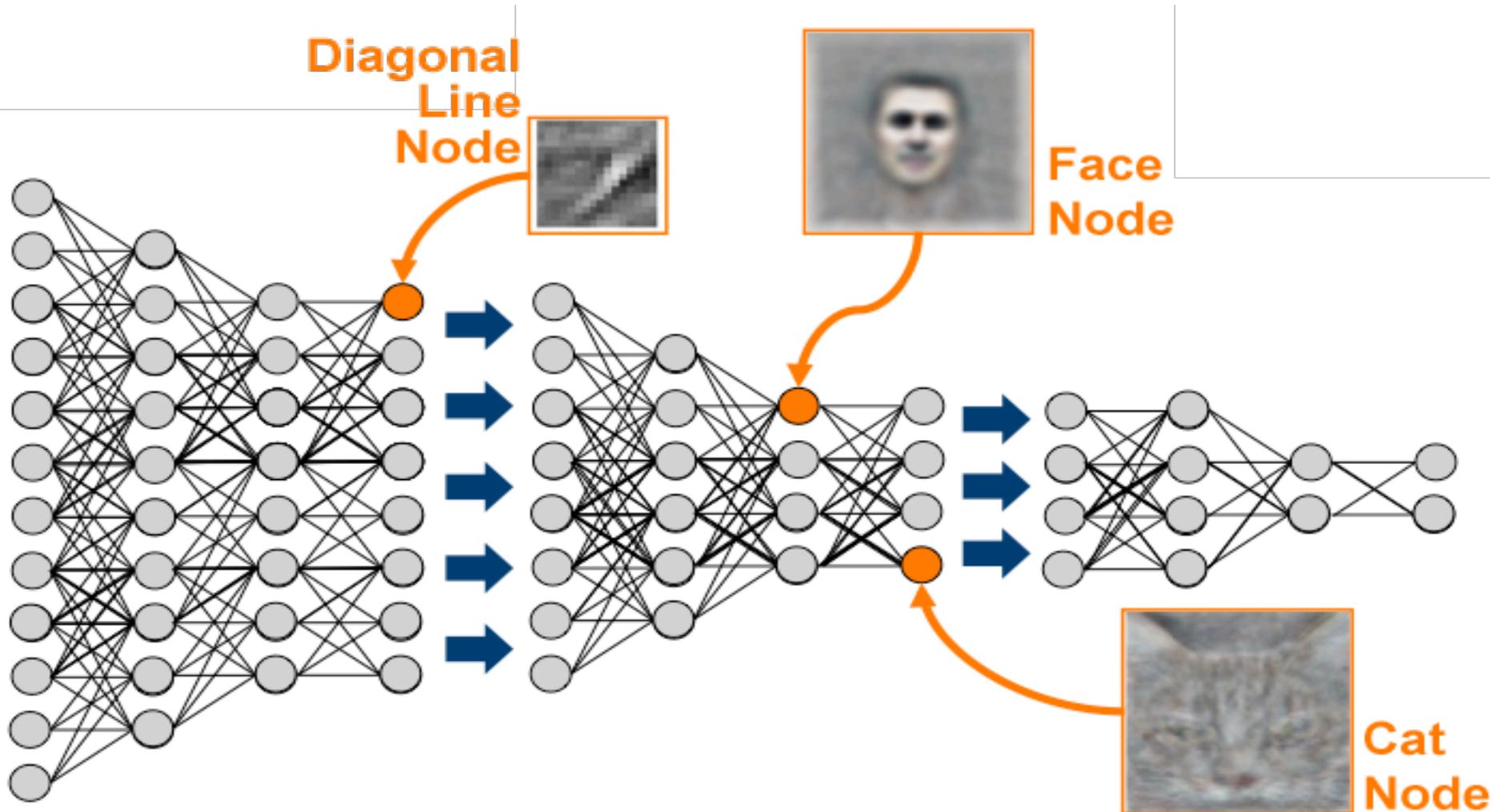
## Prediction

- Regression
- Classification
- Recommendation

# What's Common between a Data Scientist and a Fortune Teller?



# Deep Learning



## Traditional Hand-crafted Features

Macro-level Information, Human generated

Sampling

Selected Features

Core Feature Set( $10^2 \sim 10^3$ )

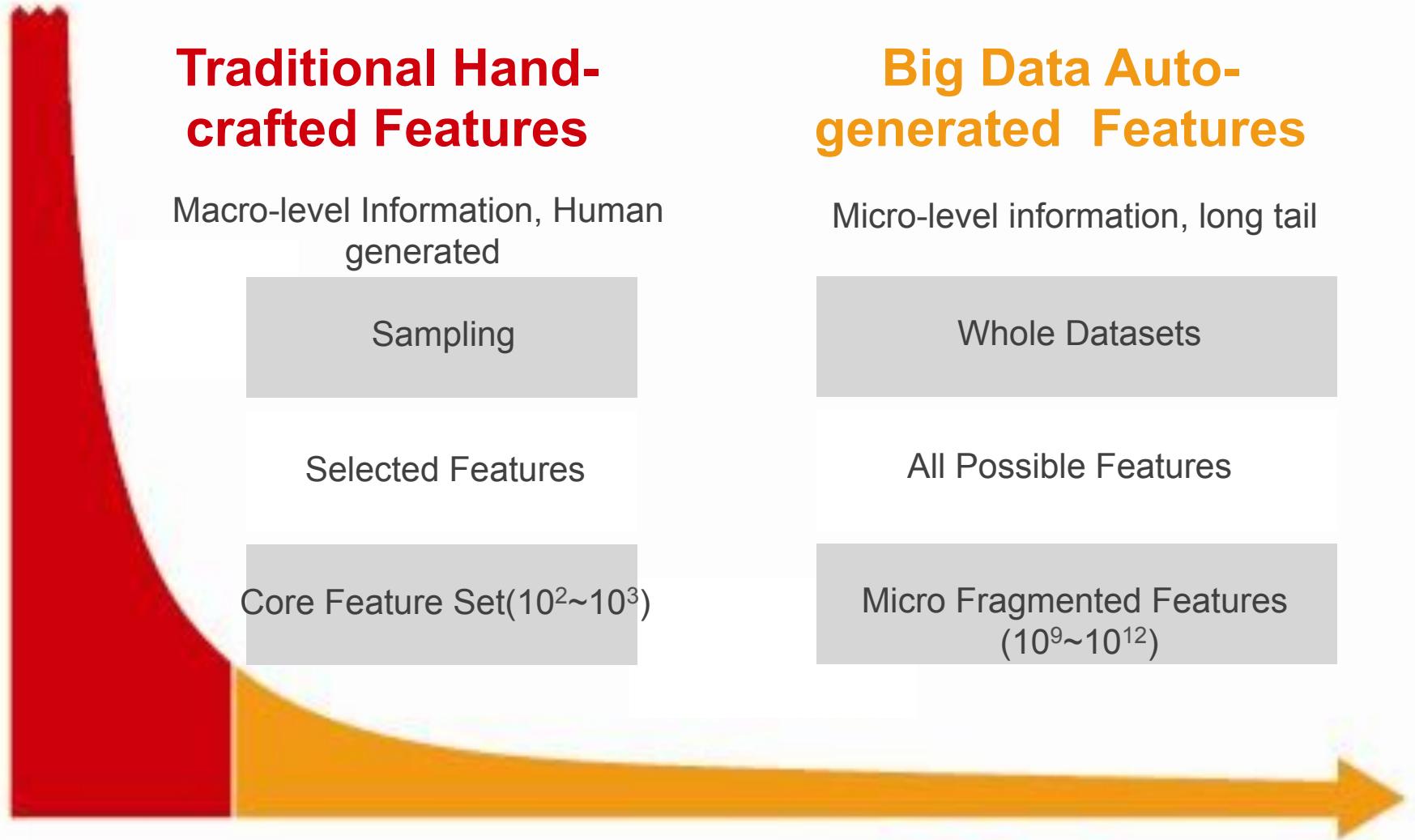
## Big Data Auto-generated Features

Micro-level information, long tail

Whole Datasets

All Possible Features

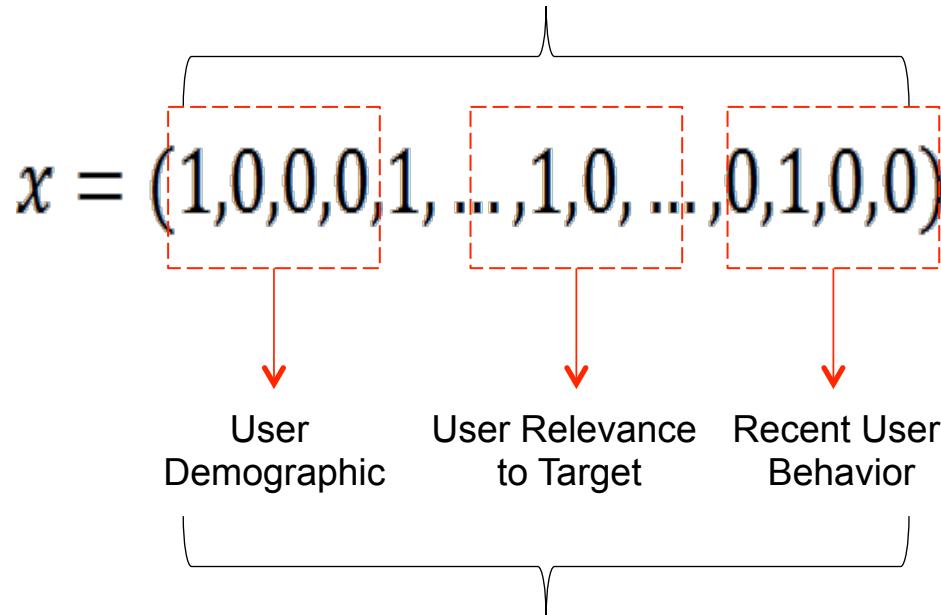
Micro Fragmented Features ( $10^9 \sim 10^{12}$ )



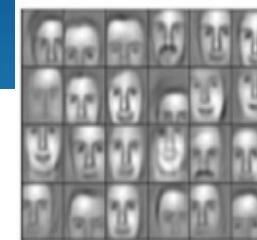
# Feature Engineering



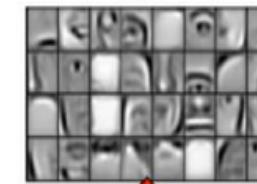
Over 100 Billion Features



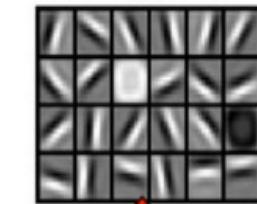
**Human Managed  
Vs.  
Automatically Generated  
and Maintained**



object models



object parts  
(combination  
of edges)

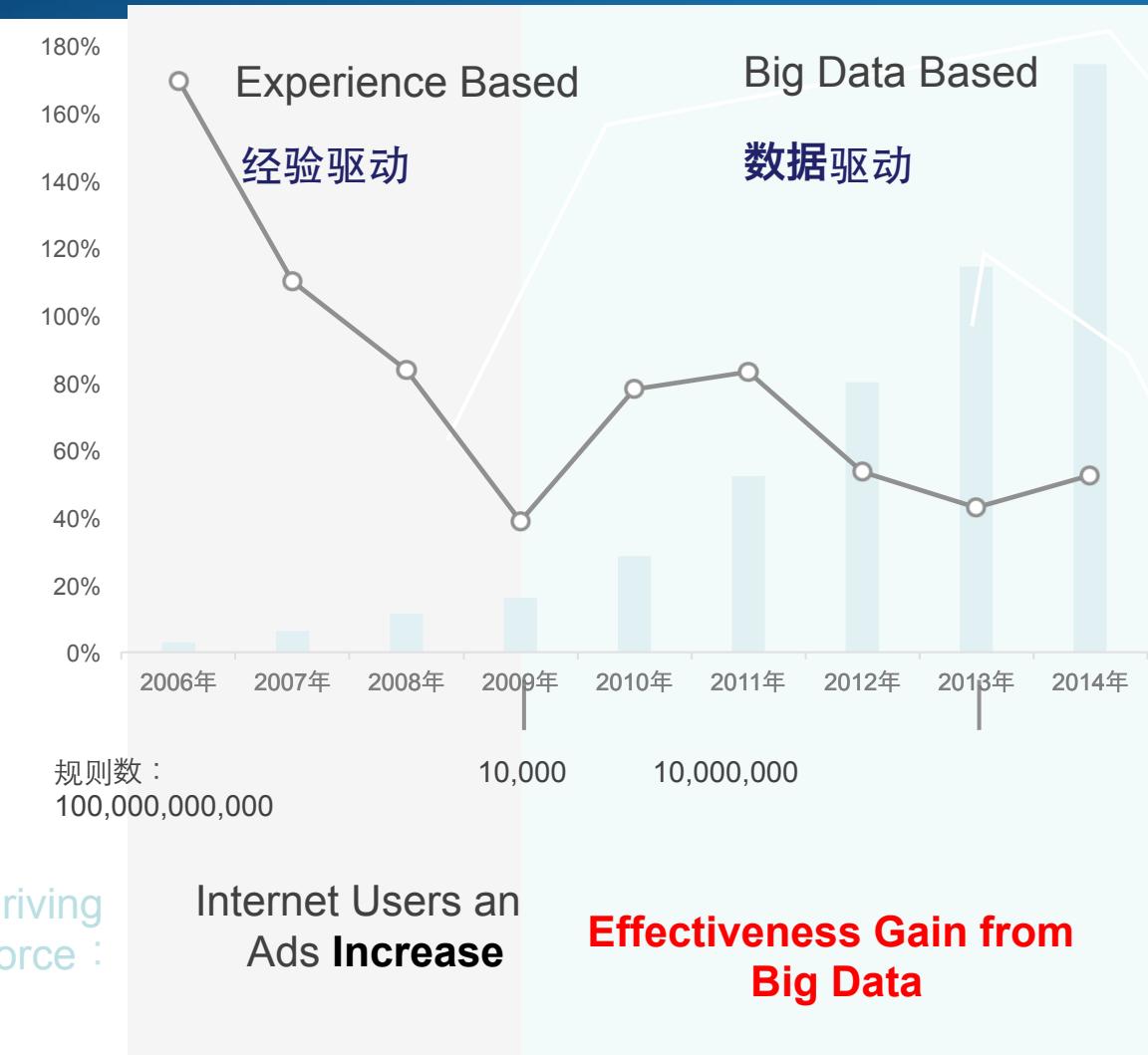


edges



pixels

# Baidu : Big Data → New Growth



2009 : 1 server online,  
#features: 10,000,  
income increase = 40%

2013 : 20,000 server  
cluster, # features:  
100,000,000,000,  
income increase = 800%



**Target:** Millions of credit card transactions per day,

**Action:** Installment recommendation via SMS

Income: **+61.7%**

Note : same ratio of SMS

# Genetic Big Data

## ATCG vs 0101



DNA: Data Source

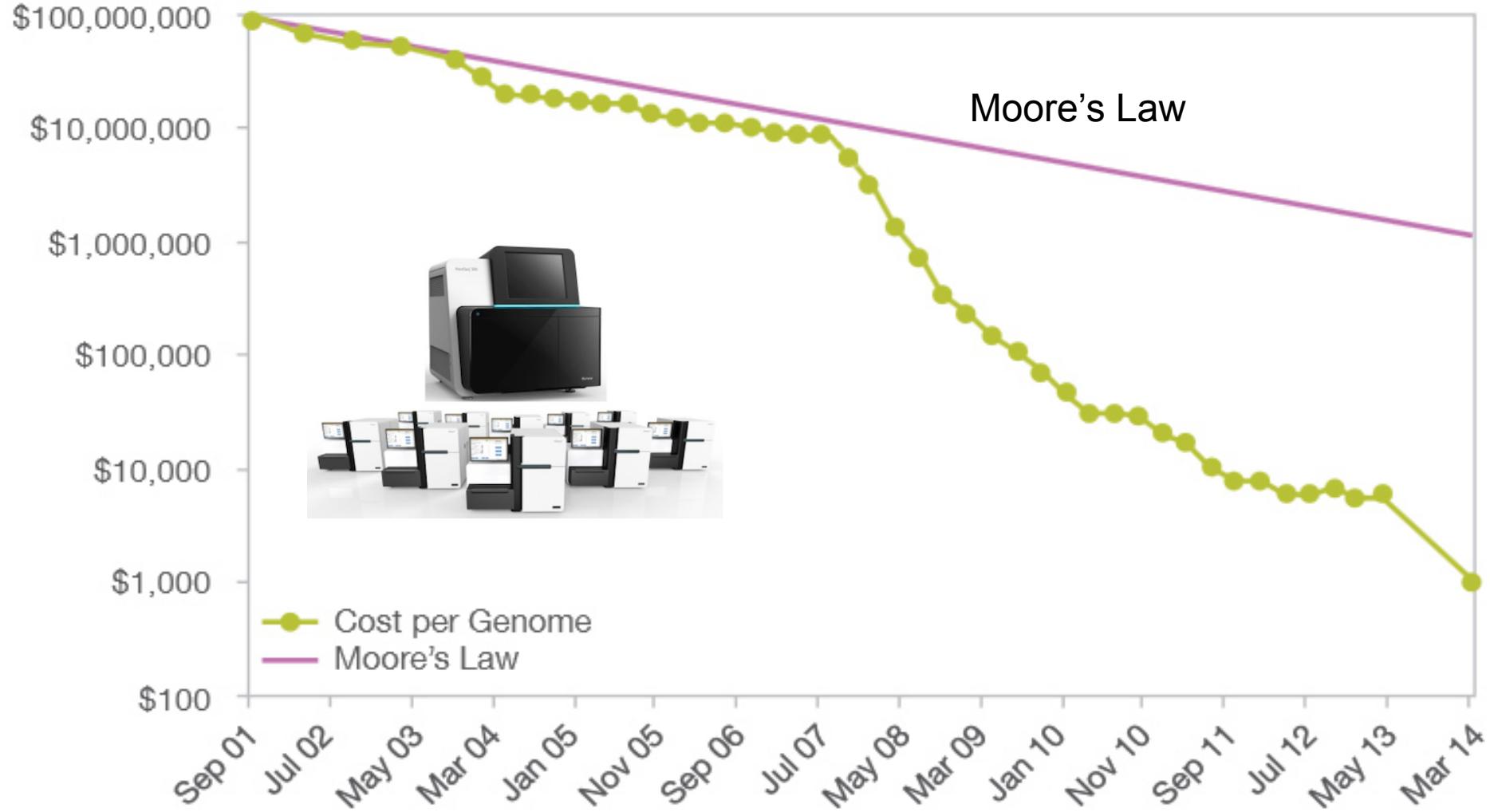
DNA: Most Effective Data Storage

0.000,000,000,001,5g DNA  
3,000,000,000 Bytes

Distance between: 3.4 Å (0.34nm)

Under certain conditions,  
DNA data can be stored for  
several million years!

# We can now read our genome at low cost, but understanding it is enormously difficult



But Understanding it is the most difficult!



Life's Equation!

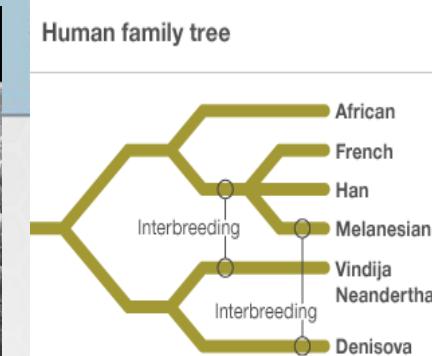
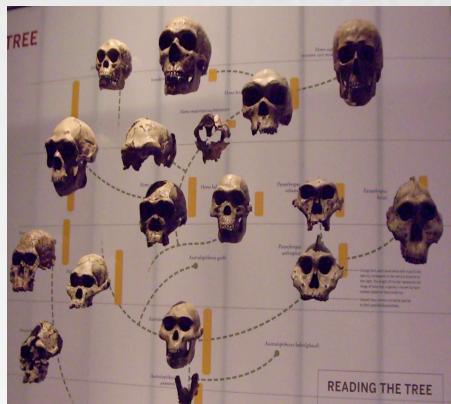
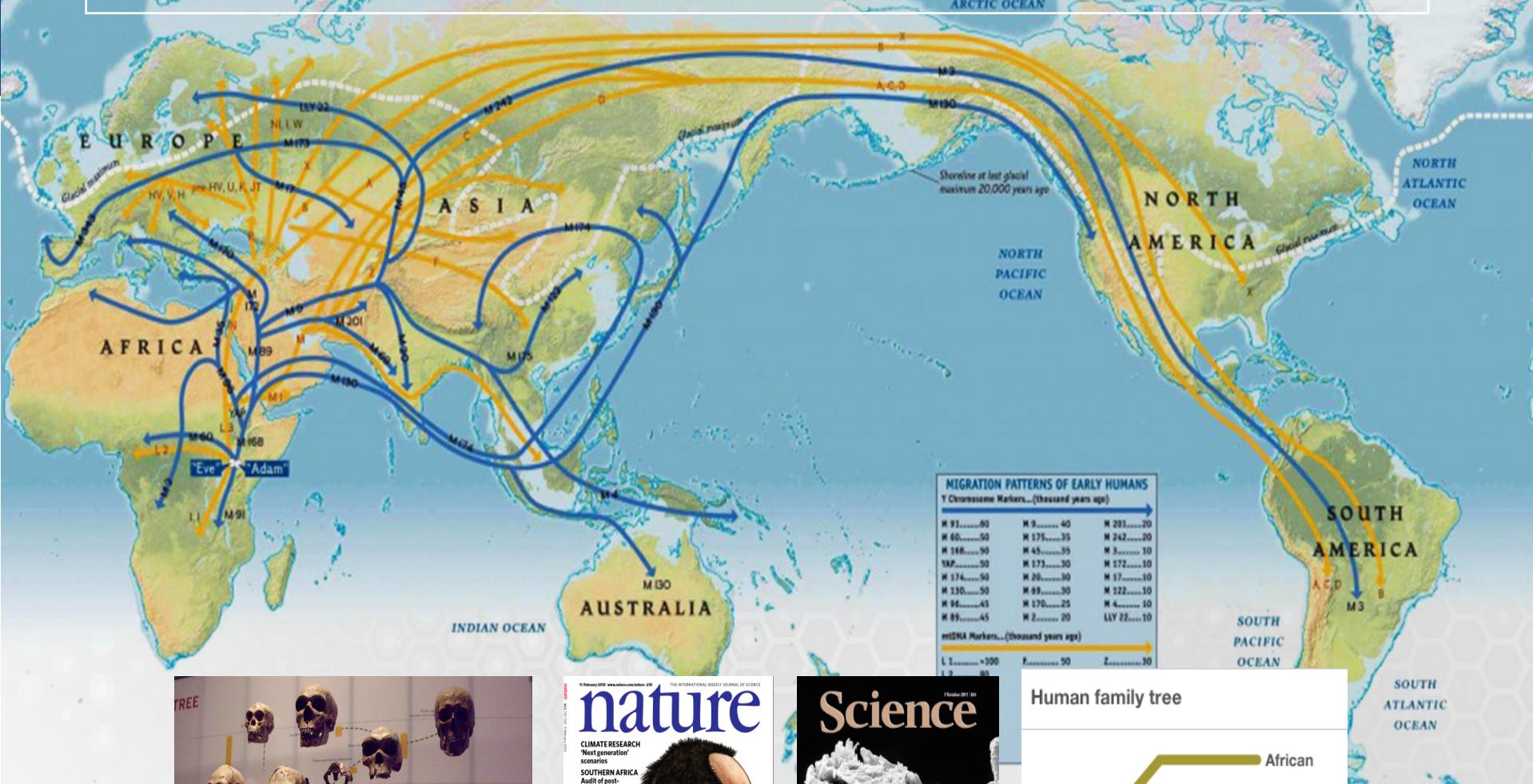




BGI Shenzhen:  
Found relation  
between genes and  
environment

Tibet Adaptive to  
High Altitude

# Understand How Humans Migrate

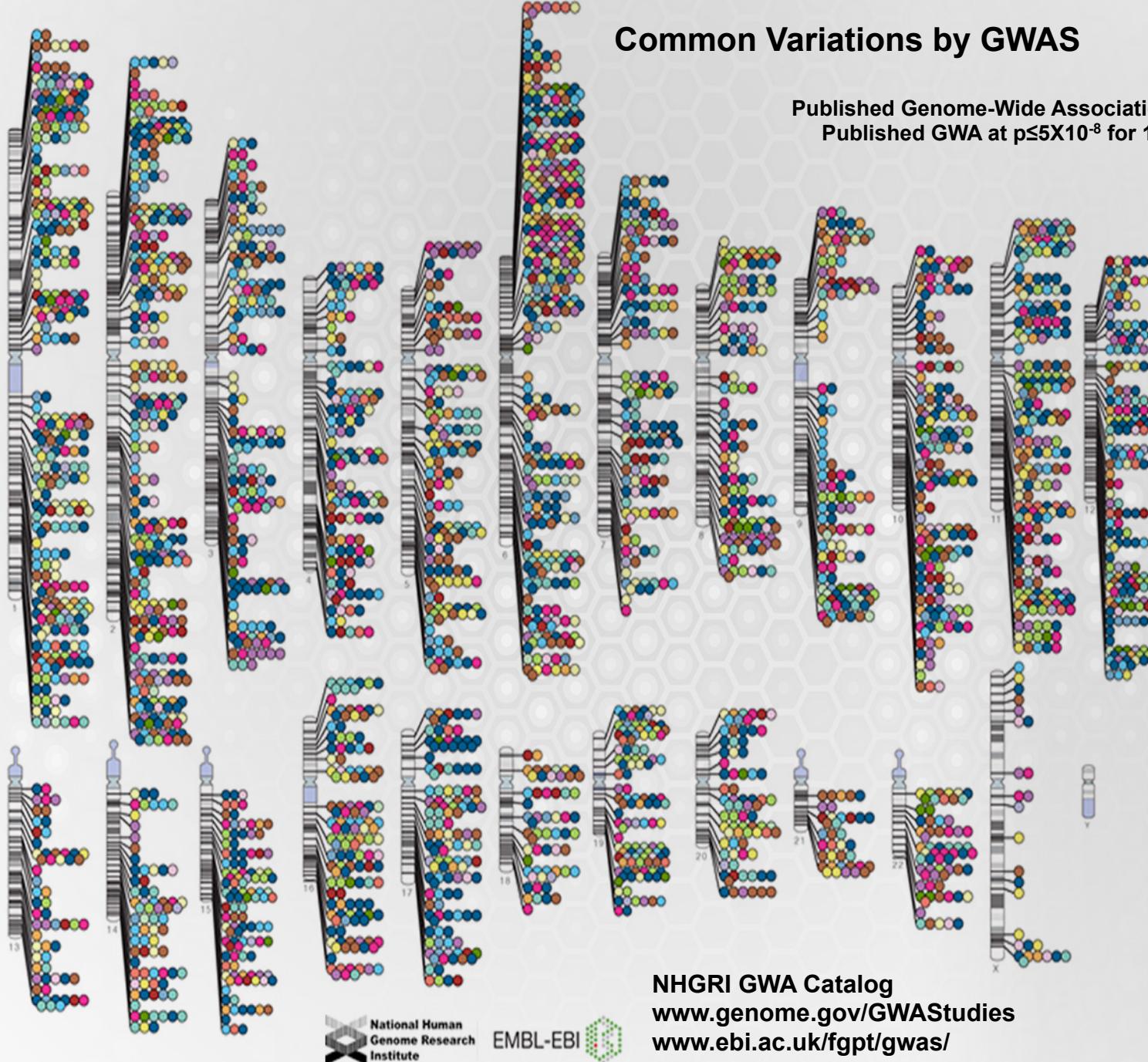


Source: Nature

# Common Variations by GWAS

Published Genome-Wide Associations through 12/2012  
Published GWA at  $p \leq 5 \times 10^{-8}$  for 17 trait categories

SNPs and  
GWAS



NHGRI GWA Catalog  
[www.genome.gov/GWASStudies](http://www.genome.gov/GWASStudies)  
[www.ebi.ac.uk/fgpt/gwas/](http://www.ebi.ac.uk/fgpt/gwas/)



## ARTICLE

doi:10.1038/nature11450

# A metagenome-wide association study of gut microbiota in type 2 diabetes

Junjie Qin<sup>1\*</sup>, Yingrui Li<sup>1\*</sup>, Zhiming Cai<sup>2</sup>, Shenghui Li<sup>1\*</sup>, Jianfeng Zhu<sup>1\*</sup>, Fan Zhang<sup>3\*</sup>, Suisha Liang<sup>1</sup>, Wenwei Zhang<sup>1</sup>, Yuanlin Guan<sup>1</sup>, Dongjian Shen<sup>1</sup>, Yingqiang Peng<sup>4</sup>, Dongya Zhang<sup>5</sup>, Zhiye He<sup>6</sup>, Wenzian Wu<sup>7</sup>, Youwen Qin<sup>1</sup>, Wenbin Xue<sup>1</sup>, Junhua Li<sup>1</sup>, Lingchuan Han<sup>1</sup>, Donghai Lu<sup>1</sup>, Peixian Wu<sup>1</sup>, Yili Dai<sup>1</sup>, Xiaojuan Sun<sup>2</sup>, Zesong Li<sup>1</sup>, Alfa Tang<sup>1</sup>, Shikong Zhong<sup>4</sup>, Xiaoping Li<sup>1</sup>, Weining Chen<sup>1</sup>, Ran Xu<sup>1</sup>, Mingbang Wang<sup>1</sup>, Qiang Feng<sup>1</sup>, Meihua Gong<sup>1</sup>, Jing Yu<sup>1</sup>, Yanyan Zhang<sup>1</sup>, Ming Zhang<sup>1</sup>, Torben Hansen<sup>5</sup>, Gaston Sanchez<sup>6</sup>, Jeroen Raes<sup>7,8</sup>, Gwen Falony<sup>7,8</sup>, Shujiro Okuda<sup>7,8</sup>, Mathieu Almeida<sup>9</sup>, Emmanuelle LeChatelier<sup>9</sup>, Pierre Renault<sup>9</sup>, Nicolas Pons<sup>9</sup>, Jean-Michel Batto<sup>9</sup>, Zhaoxi Zhang<sup>1</sup>, Hua Chen<sup>1</sup>, Ruifu Yang<sup>1,10</sup>, Weimou Zheng<sup>1</sup>, Songgang Li<sup>1</sup>, Huanning Yang<sup>1</sup>, Jian Wang<sup>1</sup>, S. Dusko Ehrlich<sup>7</sup>, Rasmus Nielsen<sup>5</sup>, Olaf Pedersen<sup>5,11,12</sup>, Karsten Kristiansen<sup>1,12</sup> & Jun Wang<sup>1,5,13</sup>

# Cure for Diabetes?

Thus, we need  
big data!

Life Big Data

碳基

Power

Endurance

VO2max

Recovery

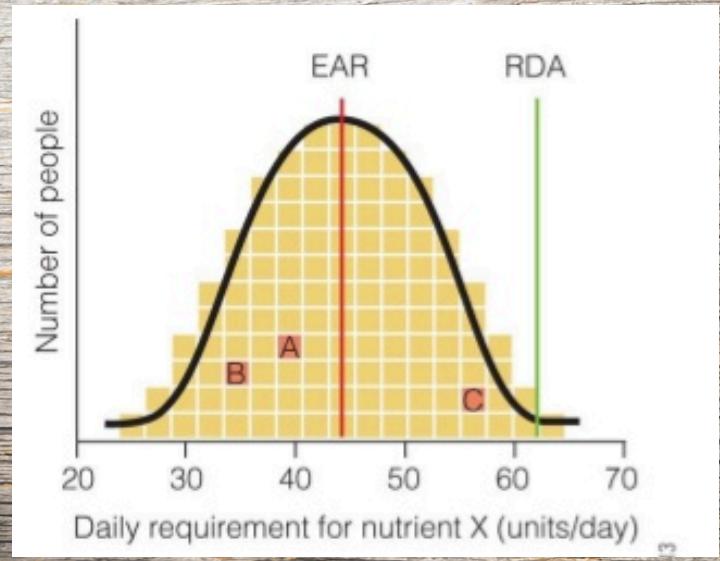
Sickness Risks



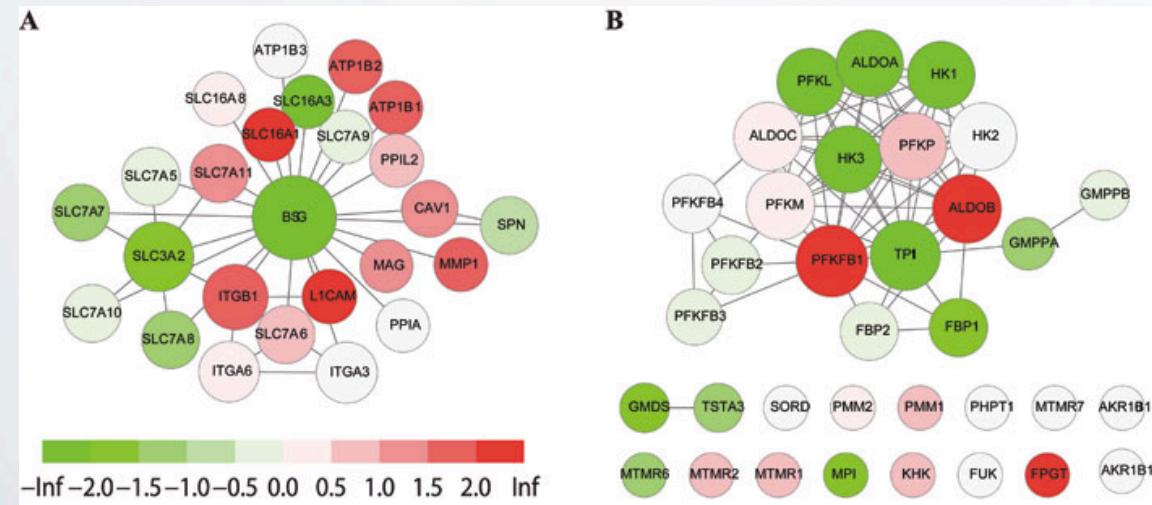
Gene+ Life Style



# Gene + Life Style



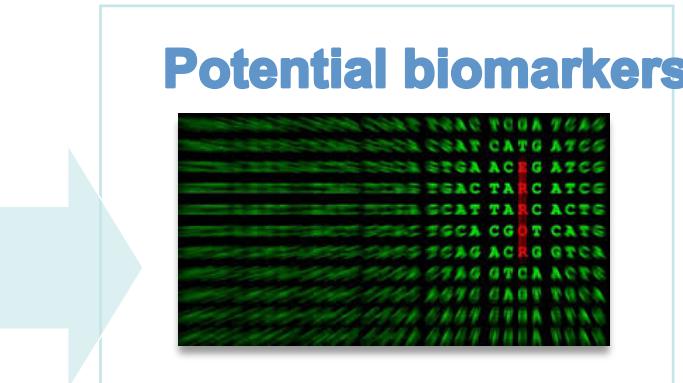
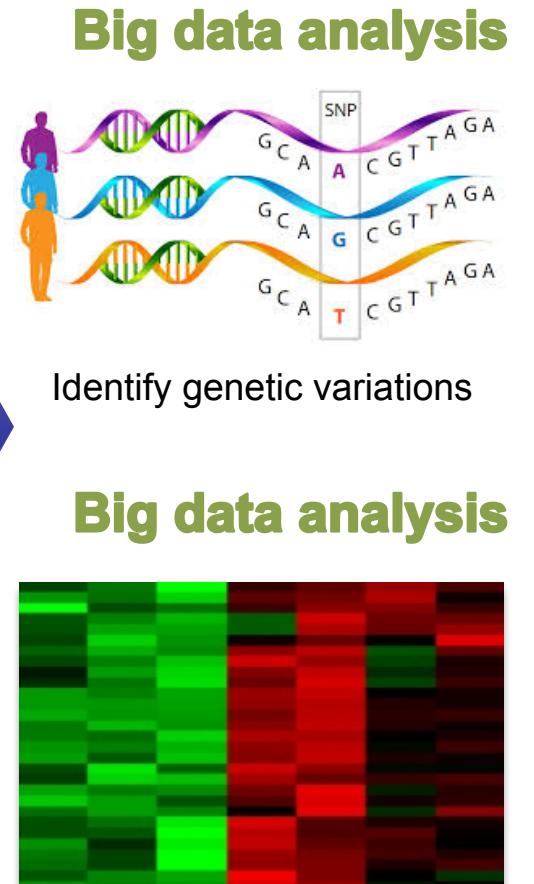
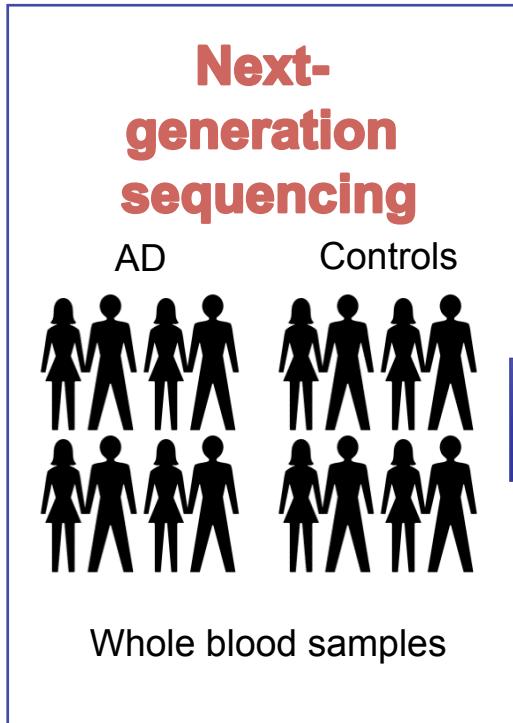
# Gene + Environment



中研网  
CIRN.com

Wang Jun@BGI Shenzhen

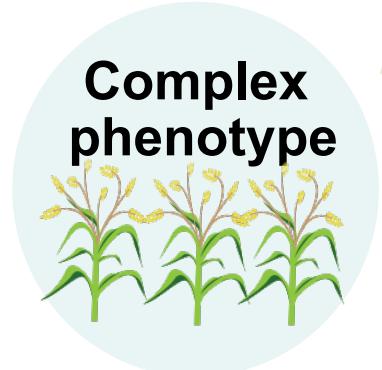
# Next-generation sequencing and big data analysis of human diseases (Alzheimer's disease)



- Quantitative analysis
- Differential expression of genes

# Project Millet: HKUST and BGI Shenzhen

**Input:** Very high dimension and low sample size labeled data ( $N \approx 2000$ ,  $D \approx 240K$ )



Machine Learning

Accelerate Hybridization Breeding

**Task:** Train an accurate phenotype predictor using genetic data.

E.g. Facilitate understand biology process

markers underlying specific phenotype.

# How to Educate a Data Scientist?

- Find structure in data
  - Can use statistical learning tools
  - Can generate new features
    - Combinational features
    - From 10 to 10 Billion
    - Samples and Features
- Build good models
  - Model must fit type of data
    - Text, signal, images, speech, social media, tables, etc.
  - Knows how to trade off between model complexity and data quantity

# How to Educate a Data Scientist?

- Focus on **questions**, not methods
  - Given the data, what questions are interesting?
  - For questions, what tools are available?
- Knows how to **interpret results**
  - Explain black box models to non-computer experts

# Conclusions: Big Data + Science = Data Science

- **Data Scientists**
  - Collect Data
  - Understand the data
  - Generate Features
  - Build and test models
  - Explain the results
  - Find new knowledge!