



# 빅데이터 이야기

## @ 더블에스텍

2015. 8. 4.

김 성 수

[sungsoo@etri.re.kr](mailto:sungsoo@etri.re.kr)

Data Management Research Section

**ETRI**

# About Me

2

## Sung-Soo Kim

*Senior Researcher*

Electronics and Telecommunications Research Institute (ETRI),

South Korea

[sungsoo@etri.re.kr](mailto:sungsoo@etri.re.kr)

[sungsookim@kaist.ac.kr](mailto:sungsookim@kaist.ac.kr)

<http://sungsoo.github.com>

[about me](#)



## WORK EXPERIENCE

### Electronics and Telecommunications Research Institute (ETRI)

September 2000 – Present

*Senior Researcher*

Research on multiscreen services, photo-realistic rendering (global illumination), real-time rendering, geometry compression algorithms, spatio-temporal data modeling/indexing and route determination algorithms.

### Korea Advanced Institute of Science and Technology (KAIST)

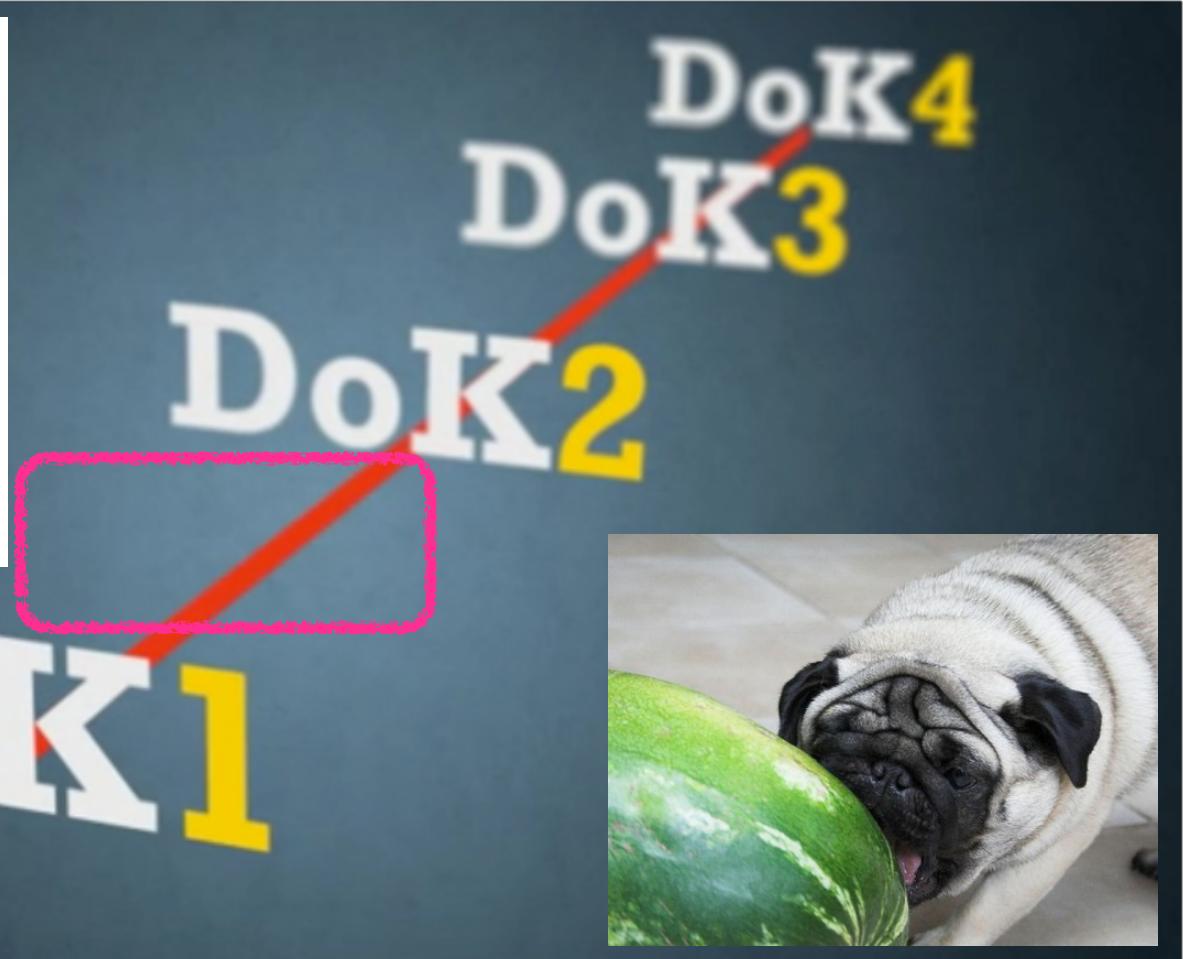
- August 2000

*Researcher*

Research on reuse based software engineering for large information systems.

# Depth of Knowledge (DoK) Level of the Seminar

3





Electronics and Telecommunications  
Research Institute

# Part I: Big Data 101

*Big Data, NoSQL, MapReduce Framework and Hadoop*

04 August 2015

Sung-Soo Kim

[sungsoo@etri.re.kr](mailto:sungsoo@etri.re.kr)

Data Management Research Section

ETRI



# The Fourth Paradigm

6



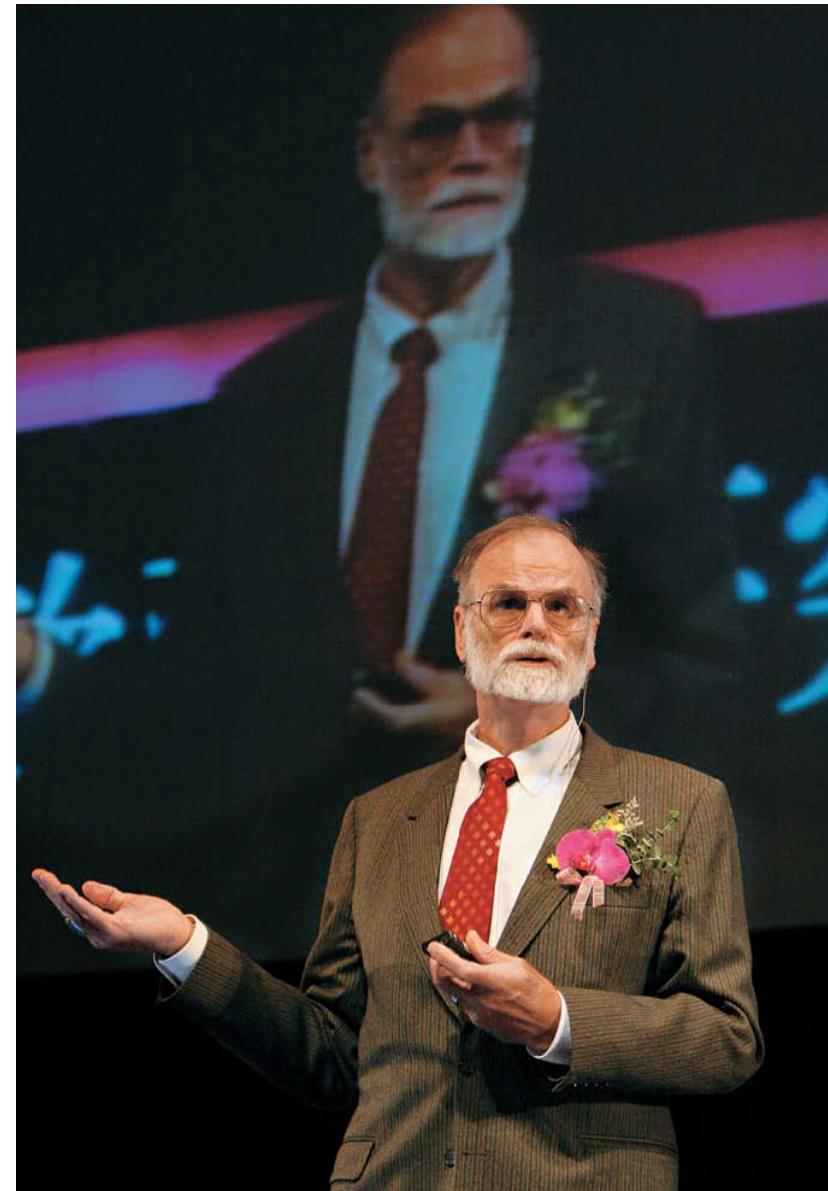
*The*  
**FOURTH  
PARADIGM**

DATA-INTENSIVE SCIENTIFIC DISCOVERY

EDITED BY TONY HEY, STEWART TANSLEY, AND KRISTIN TOLLE

# Jim Gray (Computer Scientist, Born Jan. 12. 1944 ~ Disappeared Jan. 28. 2007)

7



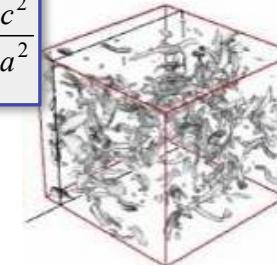
# Science Paradigms

8

# Science Paradigms

- Thousand years ago:  
**science was empirical**  
describing natural phenomena
- Last few hundred years:  
**theoretical branch**  
using models, generalizations
- Last few decades:  
**a computational branch**  
simulating complex phenomena
- Today:  
**data exploration (eScience)**  
unify theory, experiment, and simulation
  - Data captured by instruments  
Or generated by simulator
  - Processed by software
  - Information/Knowledge stored in computer
  - Scientist analyzes database / files  
using data management and statistics

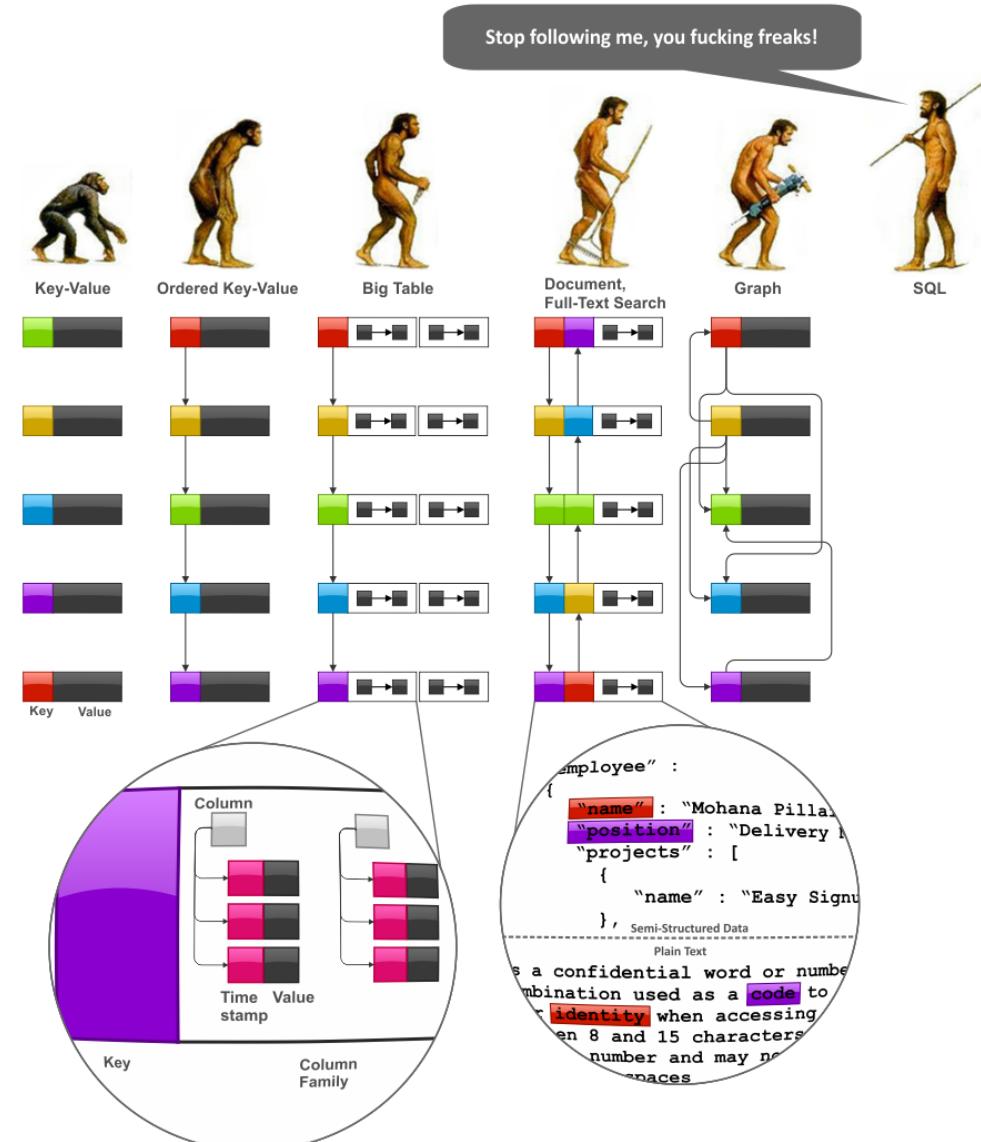
$$\left(\frac{\dot{a}}{a}\right)^2 = \frac{4\pi G\rho}{3} - K \frac{c^2}{a^2}$$



# NoSQL

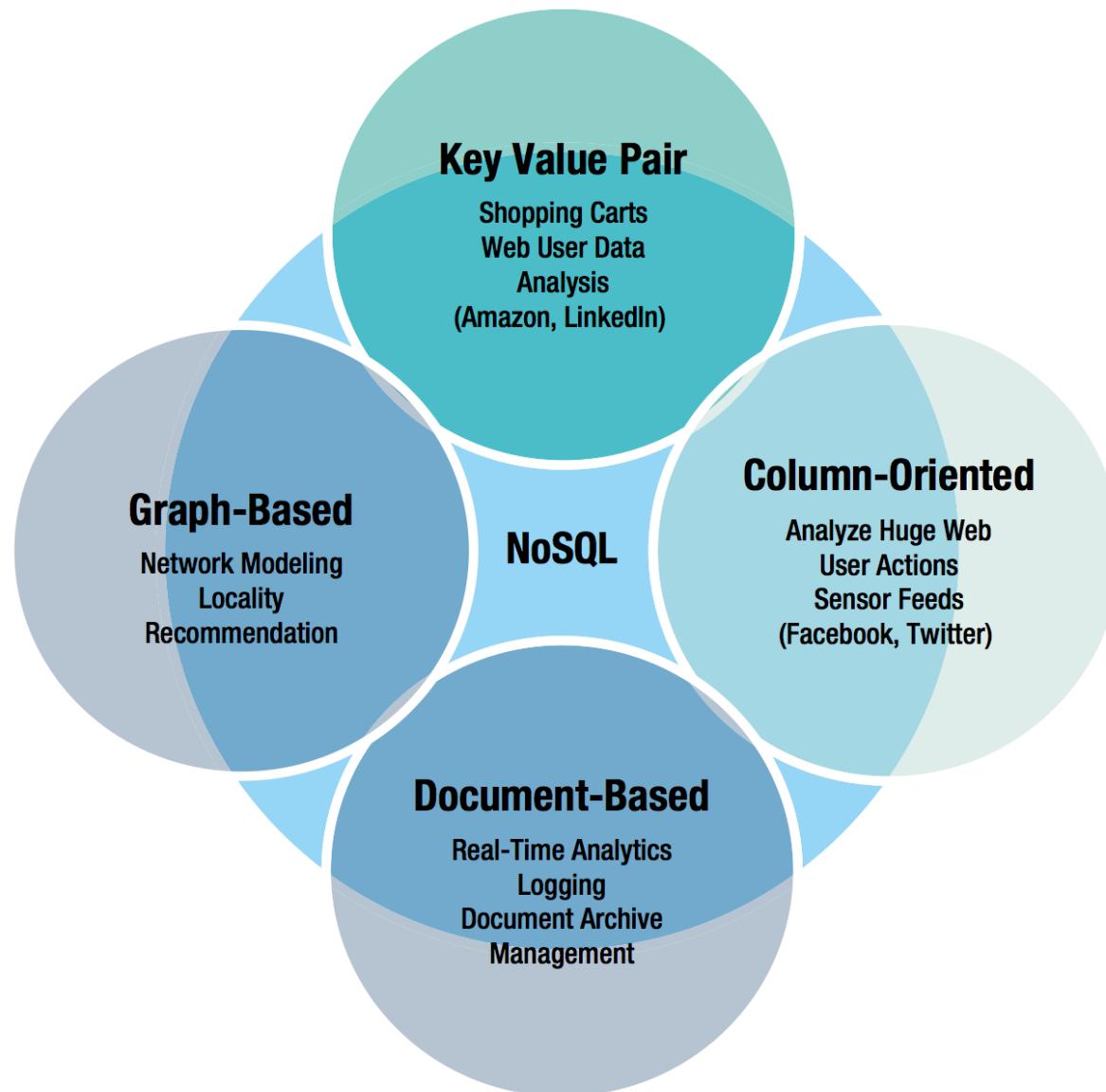
9

- “NoSQL”=“Not Only SQL”
- NoSQL Data Modeling Techniques
  - Key-Value
  - Ordered Key-Value
  - Big Table (Column-Oriented)
  - Document, Full-Text Search
  - Graph Database



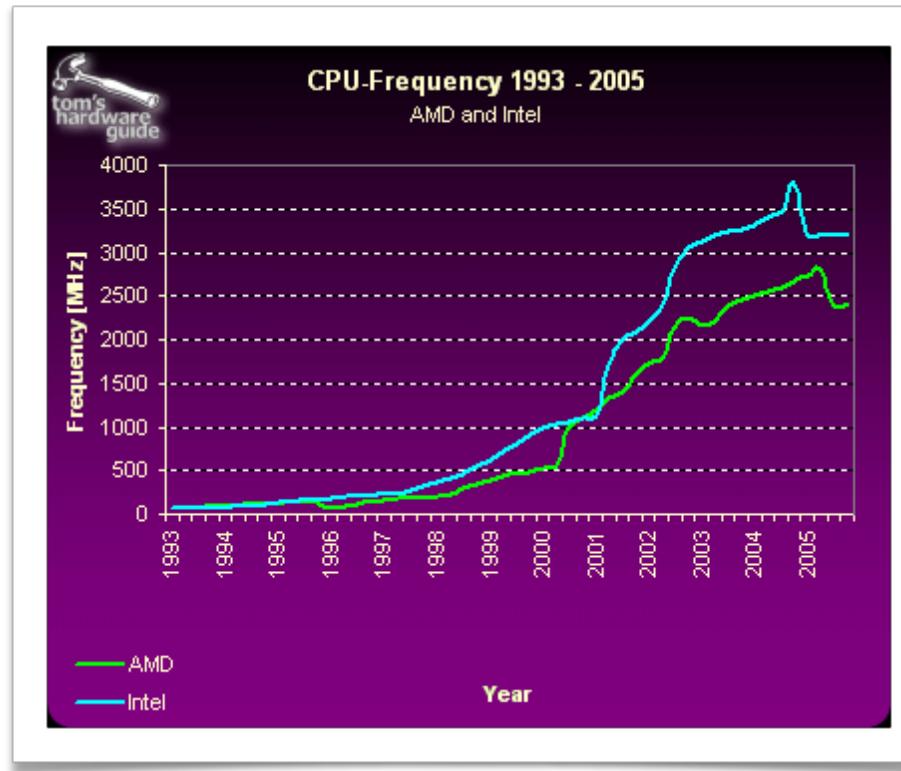
# NoSQL Database Typical Business Scenarios

10



# Computer Speedup

11



## Moore's Law

"The density of transistors on a chip *doubles every 18 months*, for the same cost" (1965)

# Distributed Problems

12

- Rendering multiple frames on high-quality animation



# Distributed Problems

13

- Simulating several hundred or thousand characters



**Happy Feet © Kingdom Feature Productions; Lord of the Rings © New Line Cinema**

# Distributed Problems

14

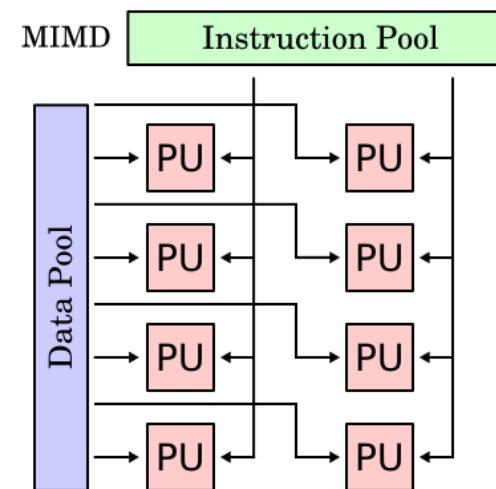
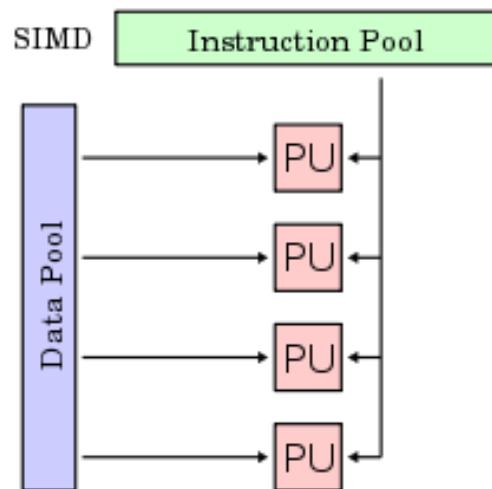
- **Indexing the web (Google)**
- **Simulating an Internet-sized network of networking experiments (PlanetLab)**
- **Speeding up content delivery (Akamai)**

**What is the *key attribute* that all these examples have in common?**

# Parallel vs. Distributed

15

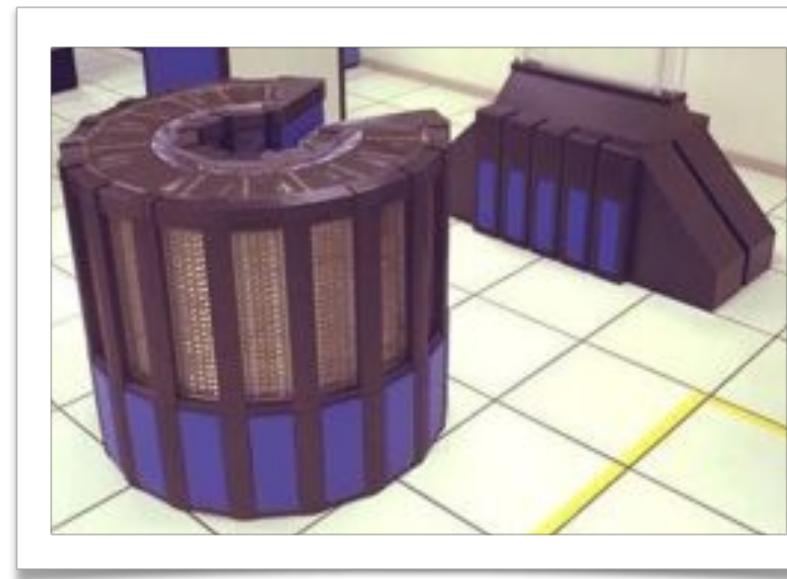
- Parallel computing can mean:
  - Vector processing of data (SIMD)
  - Multiple CPUs in a single computer (MIMD)
- Distributed computing is multiple CPUs across many computers (MIMD)



# A Brief History ... 1975 - 1985

16

- Parallel computing was favored in the early years
- Primarily vector-based at first
- Gradually more thread-based parallelism was introduced

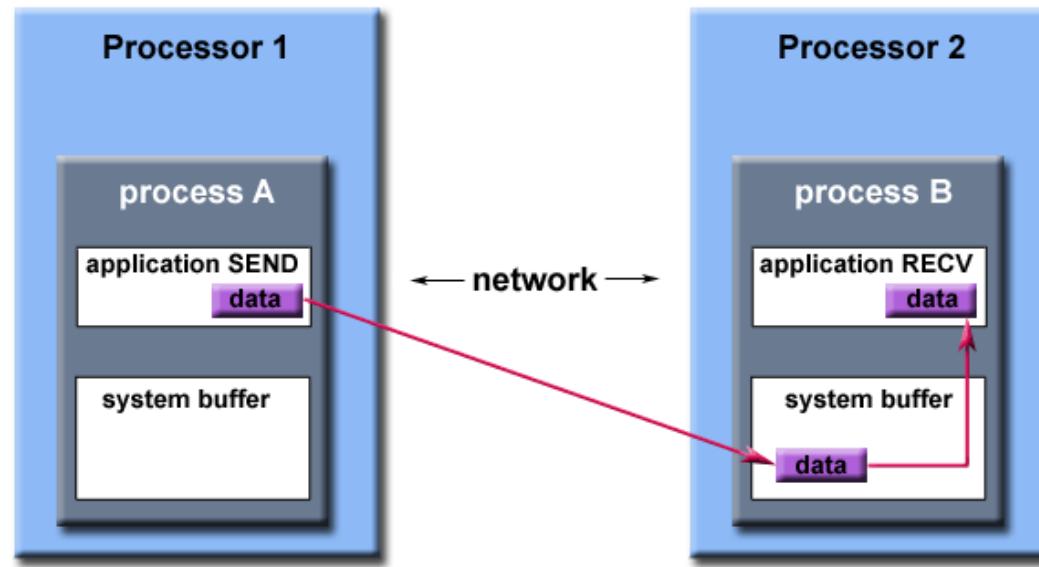


Cray 2 supercomputer ([Wikipedia](#))

# A Brief History ... 1985 - 1995

17

- “Massively parallel architectures” start rising in prominence
- Message Passing Interface (MPI) and other libraries developed
- Bandwidth was a big problem



Path of a message buffered at the receiving process

# A Brief History ... 1995 - Today

18

- Cluster/grid architecture increasingly dominant
- Web-wide cluster software
- Companies like Google take this to the extreme (10,000 node clusters)



**Detect Neonatal Patient Symptoms Sooner  
Up to 24 Hours**

**Continuously correlate data  
Thousands of events each  
second**

**Signal Processing and Data Cleansing  
Heart Rate Variability**



**"Helps detect life threatening  
conditions up to 24 hours  
sooner"**



# University of Ontario Institute of Technology (UOIT)





# University of Ontario Institute of Technology (UOIT)



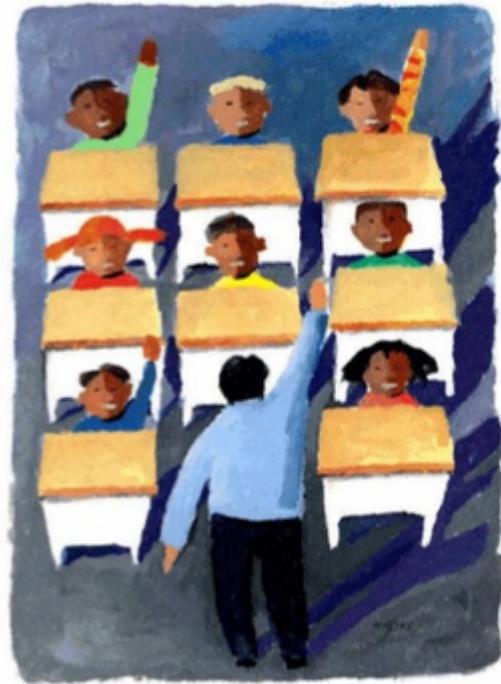


Electronics and Telecommunications  
Research Institute

# MapReduce and Hadoop Ecosystem

# Understanding MapReduce

23



MapReduce의 이해  
- 나는 반장이다 -

# Word Counting Problem

24



선생님 말씀

“반장. 이 책에 각 단어들이  
몇 번씩 나왔는지 세어줄래?”

# Oh My God!

25

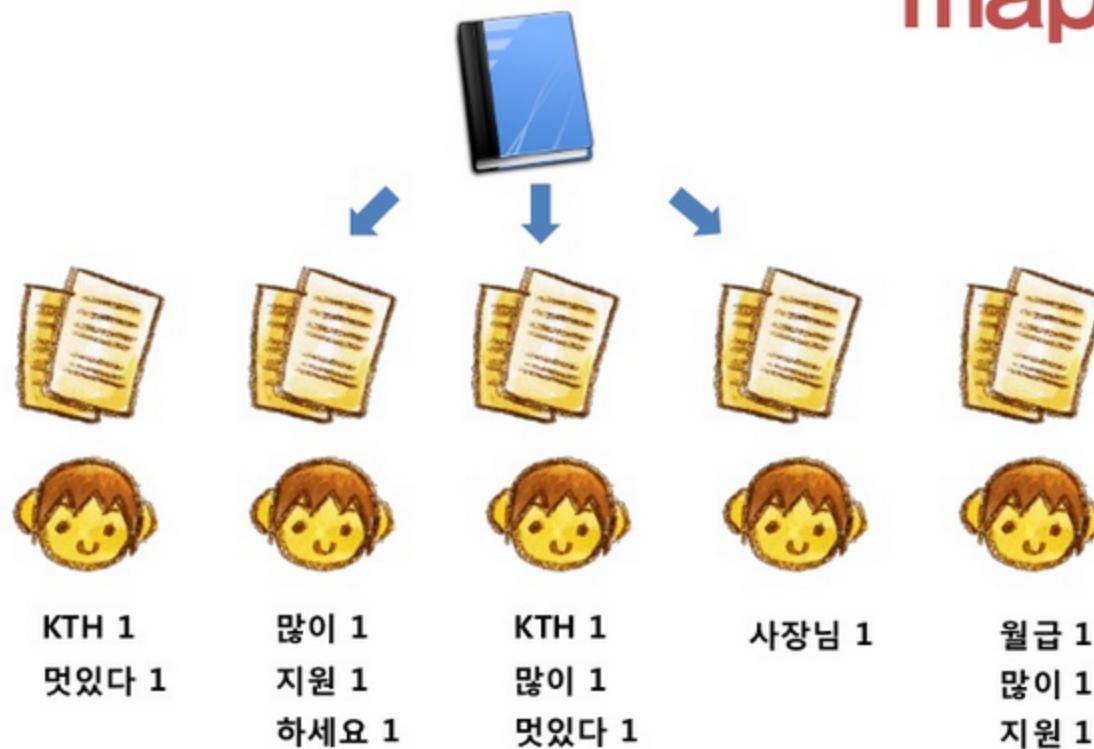


한국어

# Problem Decomposition

26

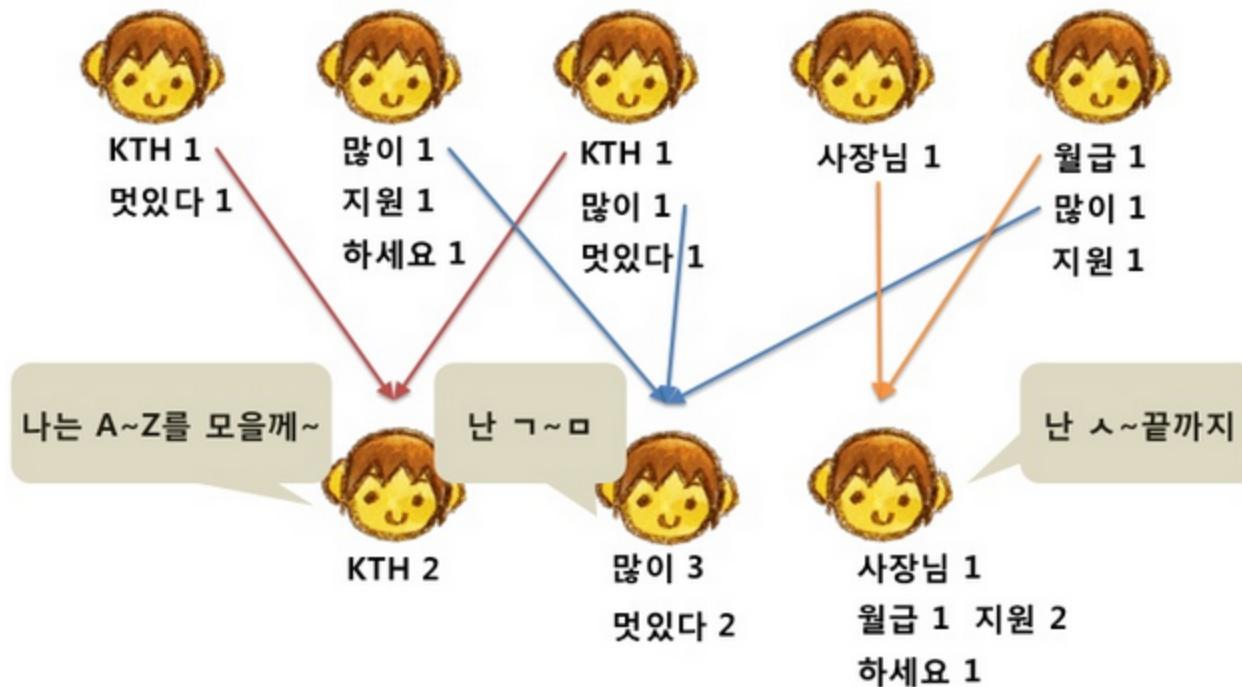
애들을 풀어 일을 나눠 시킵니다.  
**map**



# Reduction

27

몇 명이 결과를 모아 계산합니다.  
**reduce**



# What is MapReduce?

28

# MapReduce

큰 작업을

잘게 나누기(map)

종류별로 모으기(reduce)

로 처리하는 방식을 말합니다.

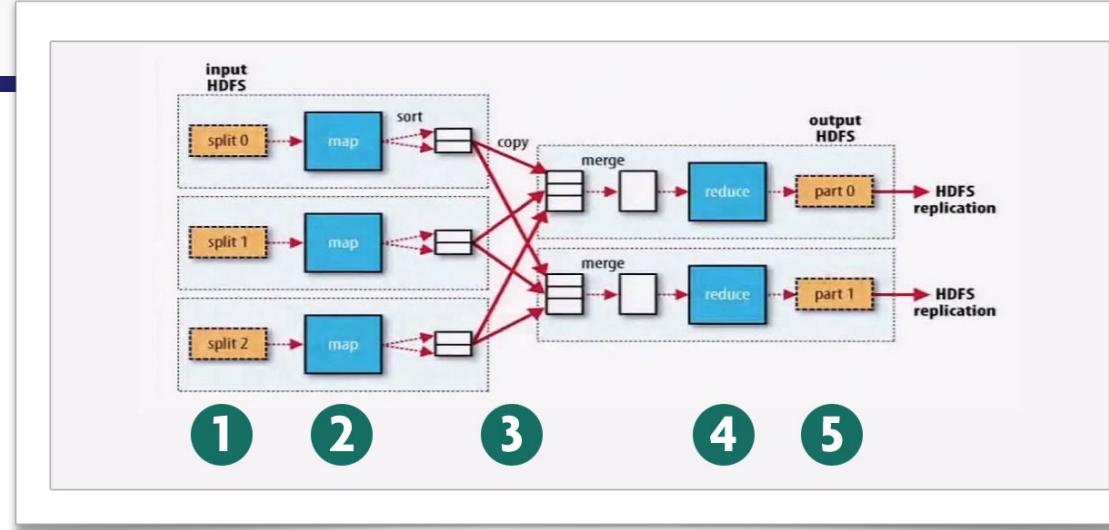
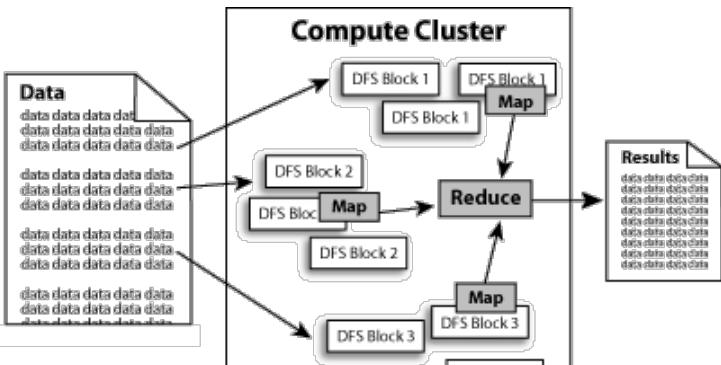
# MapReduce Framework

29

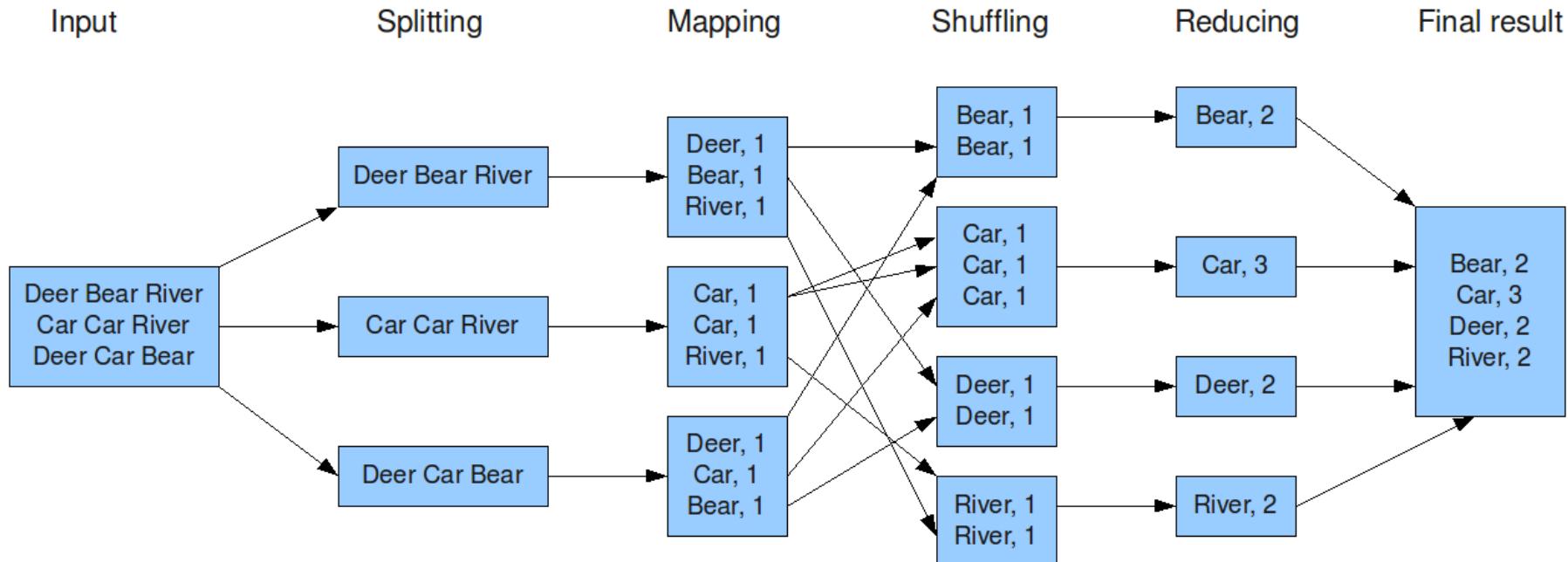
- Originally from Google, open source Hadoop
  - No data model, data stored in files [GFS, HDFS; Hadoop Distributed File System]
  - User provides specific functions
    - map(), reduce(), reader(), writer(), combiner()
  - System provides data processing “glue”, fault-tolerance, scalability
- Map and Reduce Functions
  - Map: *Divide problem into subproblems*
    - map(item) → 0 or more <key, value> pairs
  - Reduce: *Do work on subproblems, combine results*
    - reduce(key, list-of-values) → 0 or more records

# WordCount Example

30



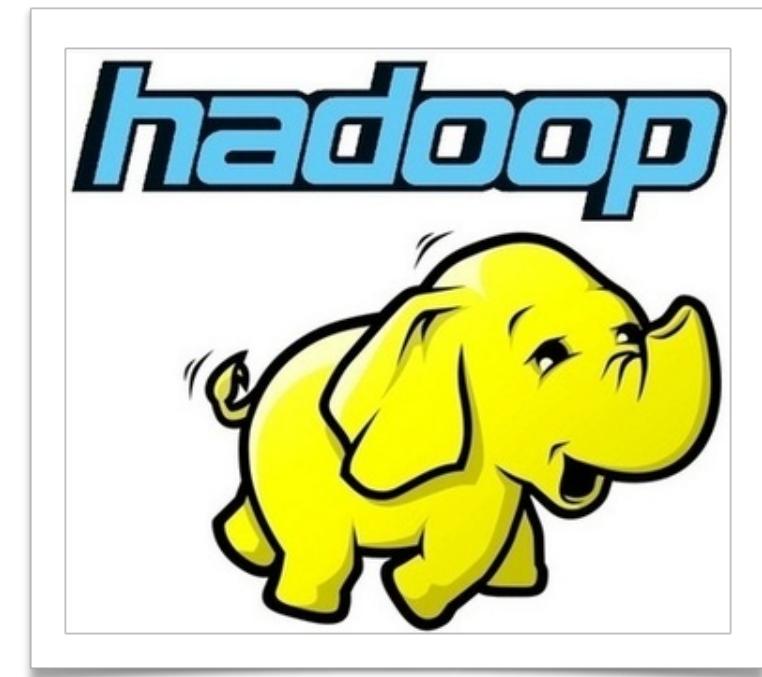
The overall MapReduce word count process



# What is Apache Hadoop?

31

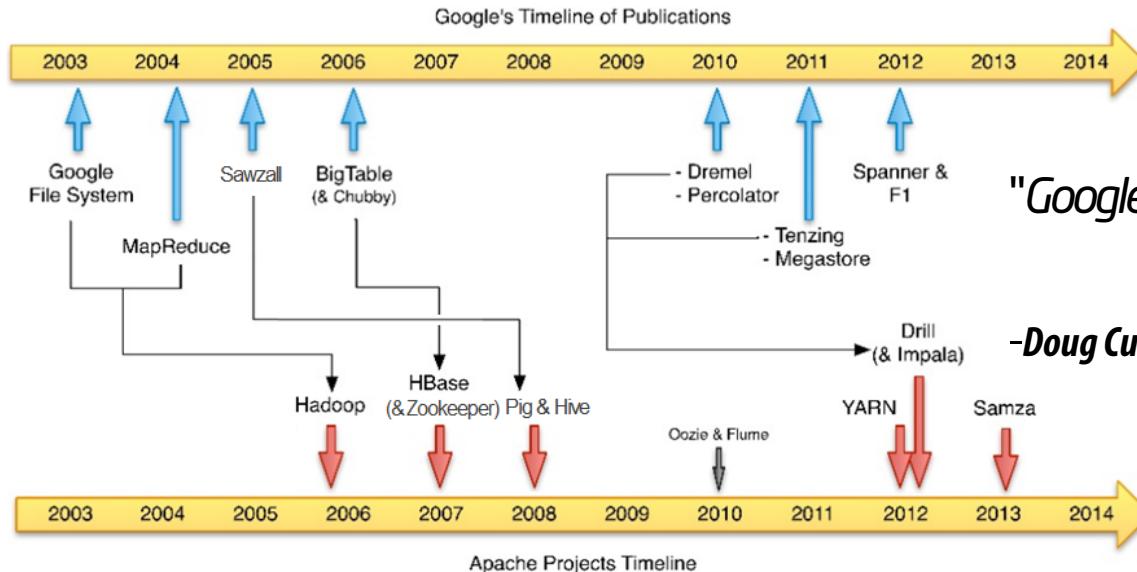
- Open source project for scalable data storage and processing
  - Inspired from Google's GFS, BigTable, MapReduce
  - Harnesses the power of *commodity servers*
  - Provides a distributed and *fault-tolerant solution*
- “Core” Hadoop consists of two main parts
  - HDFS (storage)
  - MapReduce (processing)



# Brief History of Hadoop

32

- **Created by Doug Cutting**
- **Originated Apache Nutch [2002]**
  - Nutch: web crawler engine, a part of Lucene project
  - Lucene: text search engine
- **NDFS(Nutch Distributed File System, 2004)**
- **MapReduce (2005)**
- **Official start of Apache Hadoop project (Feb 2006)**

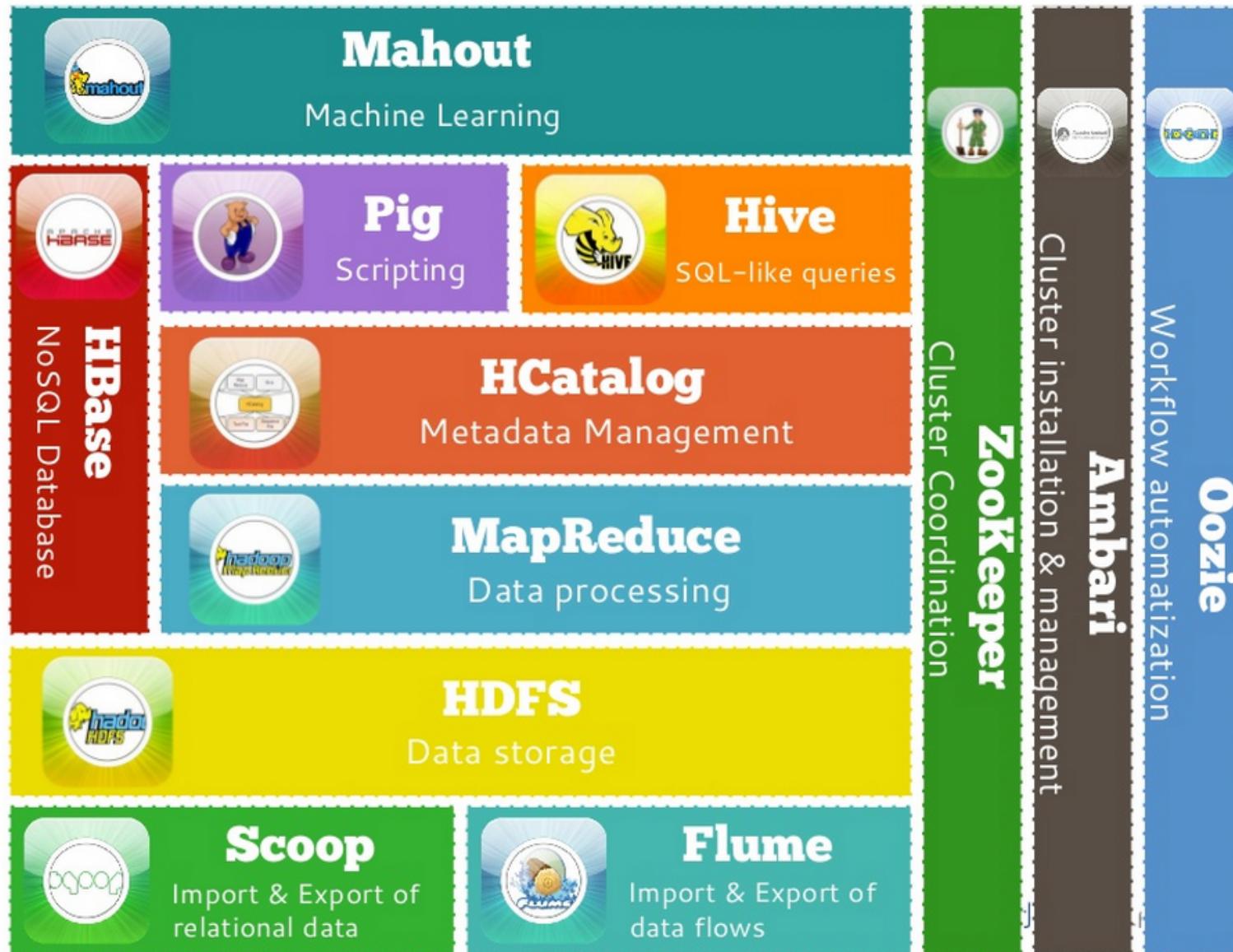


*"Google is living a few years in the future  
and sending the rest of us messages."*

-**Doug Cutting**@ O'Reilly Strata Conference in London, 2013

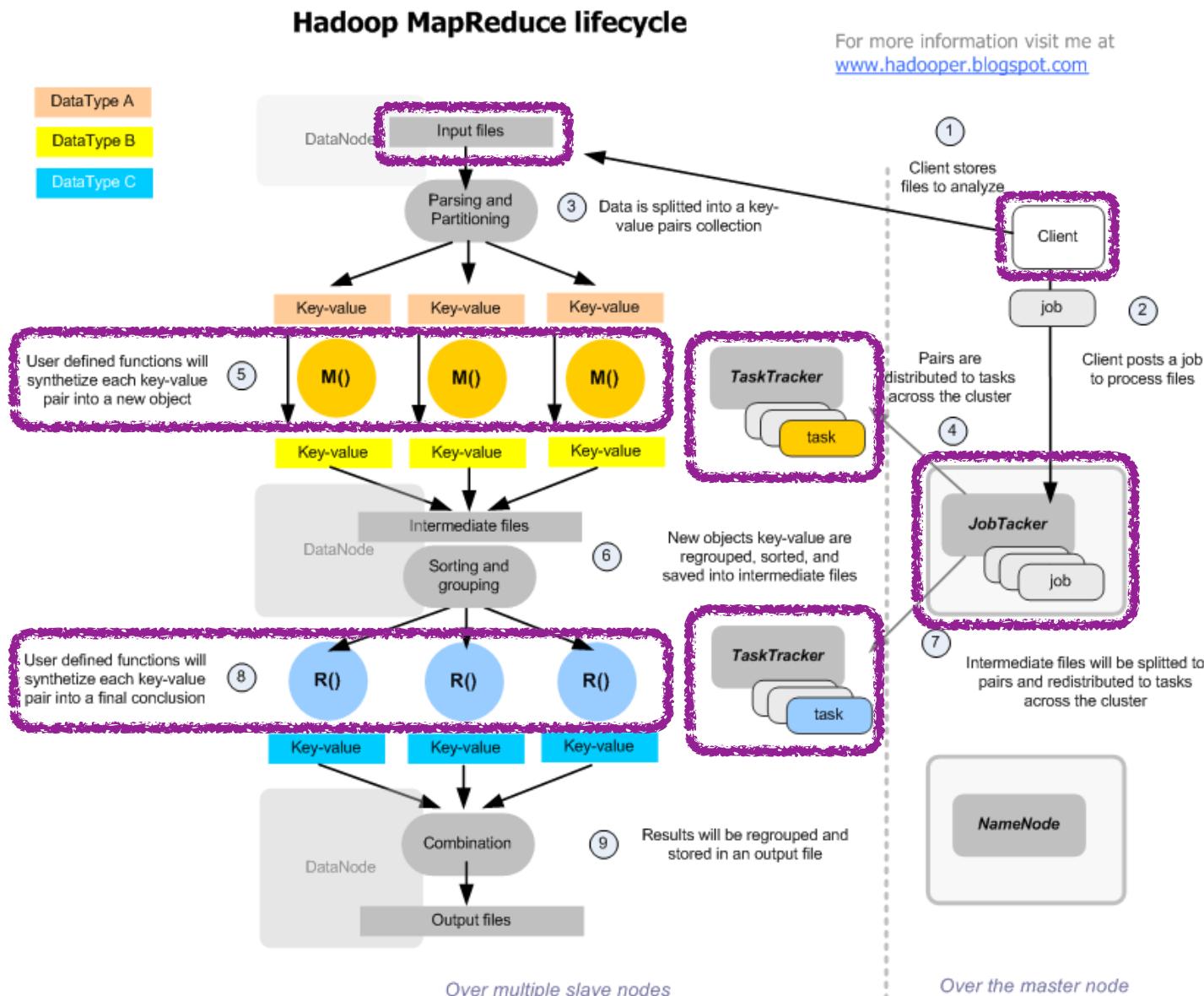
# Hadoop Ecosystem

33



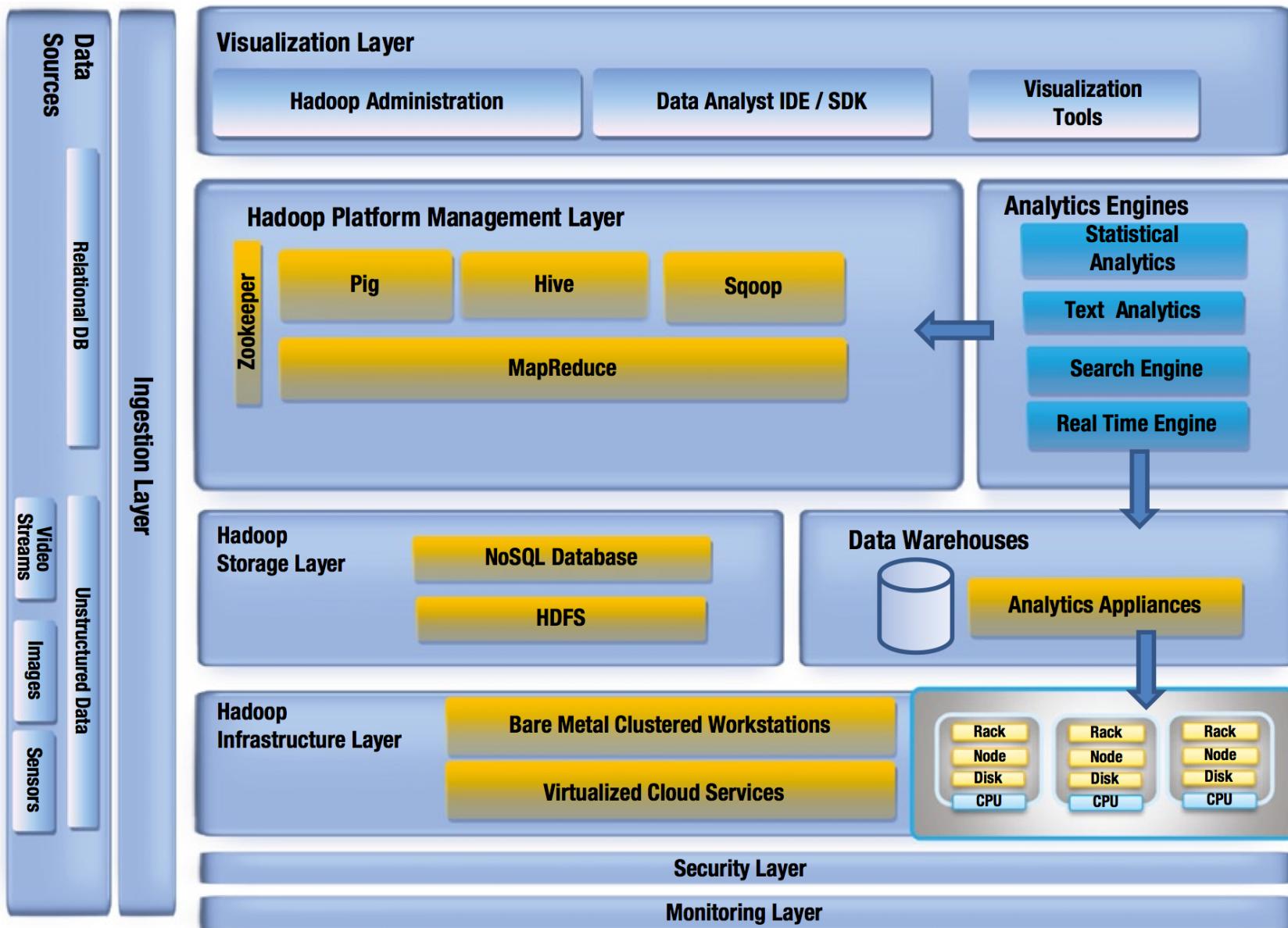
# Hadoop MapReduce Lifecycle (MRv1)

34



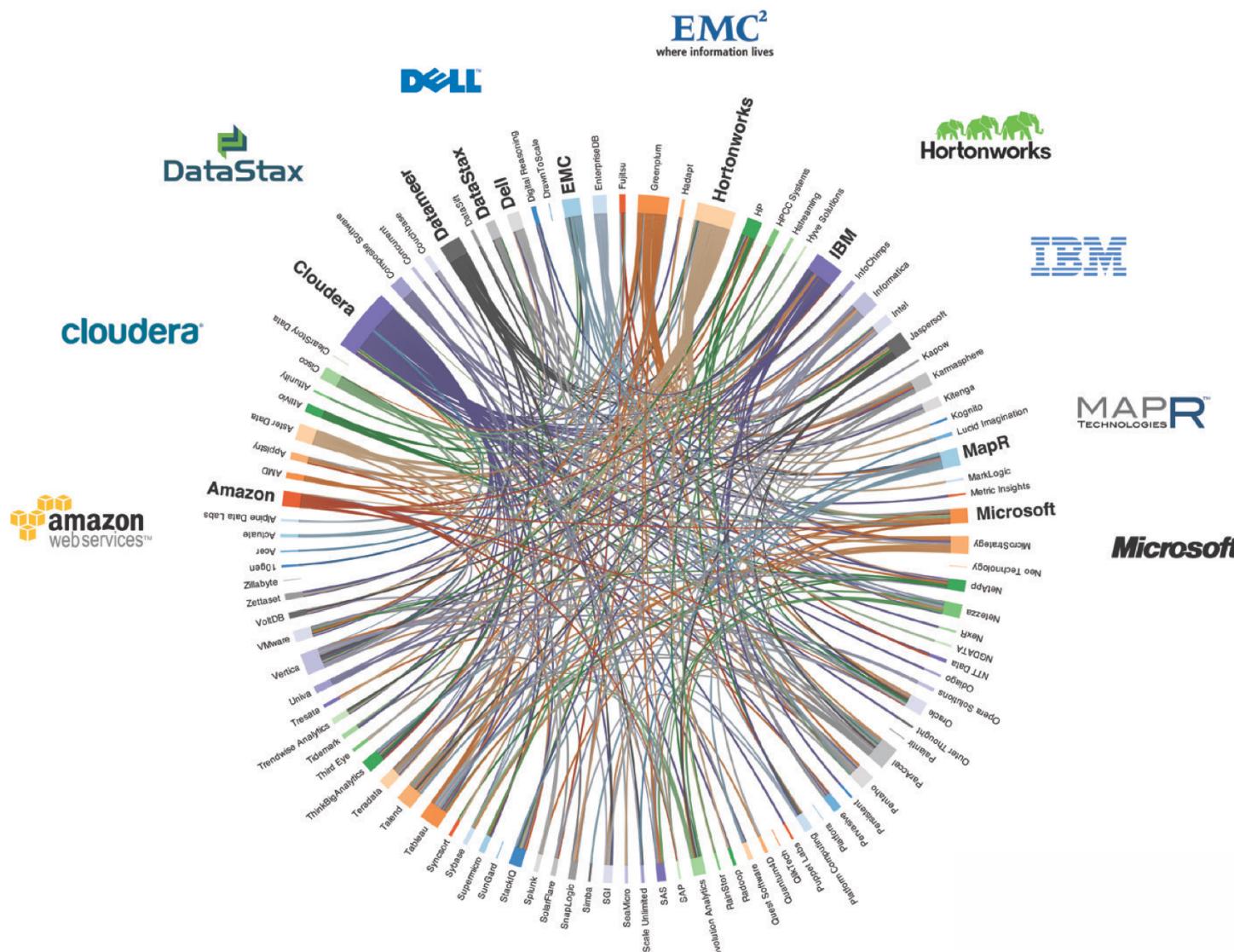
# The Big Data Architecture

35



# Hadoop Ecosystem Components as Visualized by Datameer

36



## Hadoop ecosystem components as visualized by Datameer.

# Summary: Part I

37

- **Big Data**
  - IBM's 4Vs : *Volume, Velocity, Variety, Veracity*
  - **Big data at rest**
  - **Big data in motion**
- **NoSQL**
- **MapReduce Framework**
- **Hadoop Ecosystem**
  
- **Let's take a break! 15 min**

A close-up photograph showing the lower halves of several people's bodies. They are all wearing business attire, including various shades of suits, shirts, and ties. Their hands are stacked in a circular pattern, palm up, at waist level. The hands belong to different ethnicities, suggesting a diverse team. The background is plain and light-colored.

**Collaboration is everything!**

# Brainstorming

39

- Key Requirements
  - Issue #1: Raw Data Gathering/Construction
    - 날씨, 행사, 매출(지역, 고객층), 카드사, 통신사
  - Issue #2: Business Logic, Decision 기준
    - Prototype 개발, 데이터 구성, 가중치 알고리즘 (sampling, ranking)
    - Action: customized (국내 현실에 맞는 모델이 필요)
    - Measure: 매출 이력 정보 활용
- Business Models
  - Big Data for CRM
    - 프렌차이즈: 고객 성향 분석
  - Twitter API

# Part II: Predictive Policing 101

## *Preventing Crime with Data and Analytics*



04 August 2015

Sung-Soo Kim

[sungsoo@etri.re.kr](mailto:sungsoo@etri.re.kr)

Data Management Research Section

**ETRI**





# References

42

1. Craig D. Uchida et. al, Data-Driven Crime Prevention - New Tools for Community Involvement and Crime Control, Justice & Security Strategies, Inc., 2013.
2. Walter L. Perry et.al, Predictive Policing: The Role of Crime Forecasting in Law Enforcement Operations, Research Report of RAND Safety and Justice Program, 2013.
3. Jennifer Bachner, Predictive Policing; Preventing Crime with Data and Analytics, Improving Performance Series, IBM Center for The Business of Government, 2013.
4. Ryan Gale, An Application of Risk Terrain Modeling to Residential Burglary, TCNJ Journal, 2013.
5. James Byrne, A History of Crime Analysis, Technology and the Criminal Justice Program (44.203 Course Note), UMass Lowell, Fall, 2012.
6. Mohler, G.O., M.B. Short, P.J. Brantingham, F.P. Schoenberg, and G.E. Tita, Self-exciting point process modeling of crime. Journal of the American Statistical Association 106(493):100-108, 2011.
7. Erik Lewis et. al, Self-exciting point process models of civilian deaths in Iraq, Security Journal, 2011.
8. Leslie Kennedy, Joel Caplan and Eric Piza, Risk Clusters, Hotspots, and Spatial Intelligence: Risk Terrain Modeling as an Algorithm for Police Resource Allocation Strategies, Journal of Quantitative Criminology 27, pp. 339–362, 2011.
9. Joel Caplan, Leslie Kennedy, and Joel Miller, Risk Terrain Modeling: Brokering Criminological Theory and GIS Methods for Crime Forecasting, Justice Quarterly 28:2, pp. 360–381, 2011.
10. Graham Farrell and William Sousa, Repeat Victimization and Hot Spots: The Overlap and its Implications for Crime Control and Problem-Oriented Policing, Crime Prevention Studies 12, pp. 221–240, 2001.
11. Mike Egesdal et. al, Statistical and Stochastic Modeling of Gang Rivalries in Los Angeles, SIAM Journal, 2010.
12. AUSTIN C. ALLEMAN, GEOGRAPHIC PROFILING THROUGH SIX-DIMENSIONAL NONPARAMETRIC DENSITY ESTIMATION, Technical Report, DEPARTMENT OF MATHEMATICS AND COMPUTER SCIENCE, SANTA CLARA UNIVERSITY, 2012.
13. P Jeffrey Brantingham, Prey selection among Los Angeles car thieves, Crime Science 2013.
14. George Mohler, MODELING AND ESTIMATION OF MULTI-SOURCE CLUSTERING IN CRIME AND SECURITY DATA, Annals of Applied Statistics, 2013.
15. George Mohler, ESTIMATING THE HISTORICAL AND FUTURE PROBABILITIES OF LARGE TERRORIST EVENTS, Annals of Applied Statistics, 2013.

# Outline

43

- **Why is the Predictive Policing Important?**
- **Background**
- **Mathematical Frameworks**
- **Data Used in Predictive Policing**
- **Predictive Methodologies for Predictive Policing**
- **Summary**

# Why is the Predictive Policing Important?

44



# Why is the Predictive Policing Important?

45

- **What is Crime Analysis?**

- **Administrative, Tactical and Strategic Problem Solving**
- **Forecasting and predictive analysis**



*August Vollmer*  
Source: LAPD Web Site

**1909**

- **Principles of problem-oriented policing**
- **Evidence-based policing**
- **Real-time crime analysis**
- **Intelligence-led policing and other proven policing models**



*Orlando W. Wilson*

**1950**

# History

## ■ *Origin of Crime Analysis*

- Crime Mapping: 1829, Italian geographer and French Statistician
  - Crime Analysis: 1842, London's Metropolitan Police

## ■ *First Data Miner*

- John Graunt [Amateur Scientist]
  - Bill of Mortality (1st: 21 December 1592)
    - essentially random set of information
    - the study of the patterns causes and effects of disease
    - *Plague*
      - caused by person-to-person contact
      - tended to increase during the first year of the reign of a new king

ss section  
ve  
agation

22 Hertz

$\lambda$  - wavelength 26.71

a - amplitude 80.42

v - velocity 5.02

T - time 1.28

f - frequency 0.22

$\rho$  = density  
 $p$  = pressure  
 $\mu$  = viscosity

$$\rho \left( \frac{\partial \mathbf{v}}{\partial t} + (\mathbf{v} \cdot \nabla) \mathbf{v} \right) = -\nabla p + \mu \nabla^2 \mathbf{v} + \mathbf{f}$$
$$\nabla \cdot \mathbf{v} = 0,$$

# Clustering Patterns

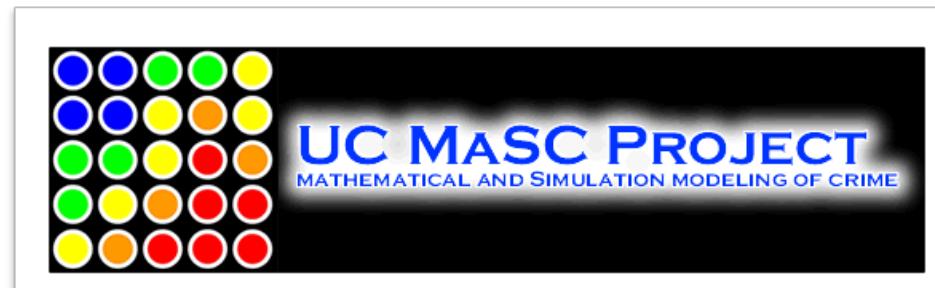
48

- Loma Prieta Earthquake (1989)

- Earthquake analysis
  - Clustering patterns

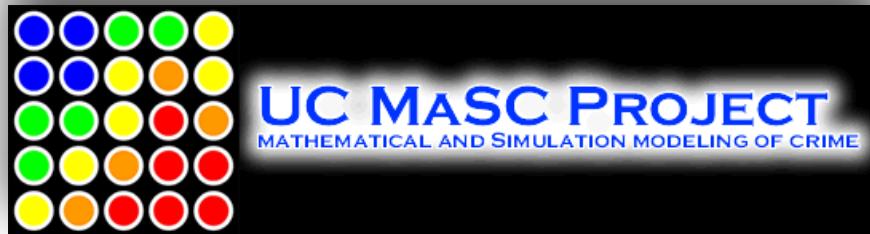
- MASC Project

- Clustering patterns are also seen in *crime data*
  - Future events *nearby in space and time* after shocks of crime
  - Finding and clustering the *patterns* from chaotic human behaviours
  - *Self-exciting point process model similar* to that used in earthquake analysis
  - Pilot Project : Los Angeles
    - Data: 13 million past crimes
    - Duration: past 80 years



# Mathematical and Simulation Modeling of Crime

49



***Jeff Brantingham***  
**UCLA Anthropology**



***George Mohler***  
**Santa Clara Mathematics**



***George Tita***  
**UCI Criminology,  
Law and Society**

# Mathematical Framework [7]

50

## ■ The Hawkes Process Model (1974)

- any given event, or collection events can be causally linked to a background Poisson process and foreground self-exciting process.

### *Self-excitation*

*the distribution of crimes following an initial event*

$$\lambda(t) = \mu + k_0 \sum g(t - t_k)$$

*background rate of events*  
*stationary Poisson process*

*the density of prior events*  
*necessary to trigger excitation*

*how much excitation is generated*  
*by a collection of prior events*

SERRA MEDICAL CENTER

$$\lambda(t) = u + \lambda(t-t_1)$$

17:352

BURGLARY

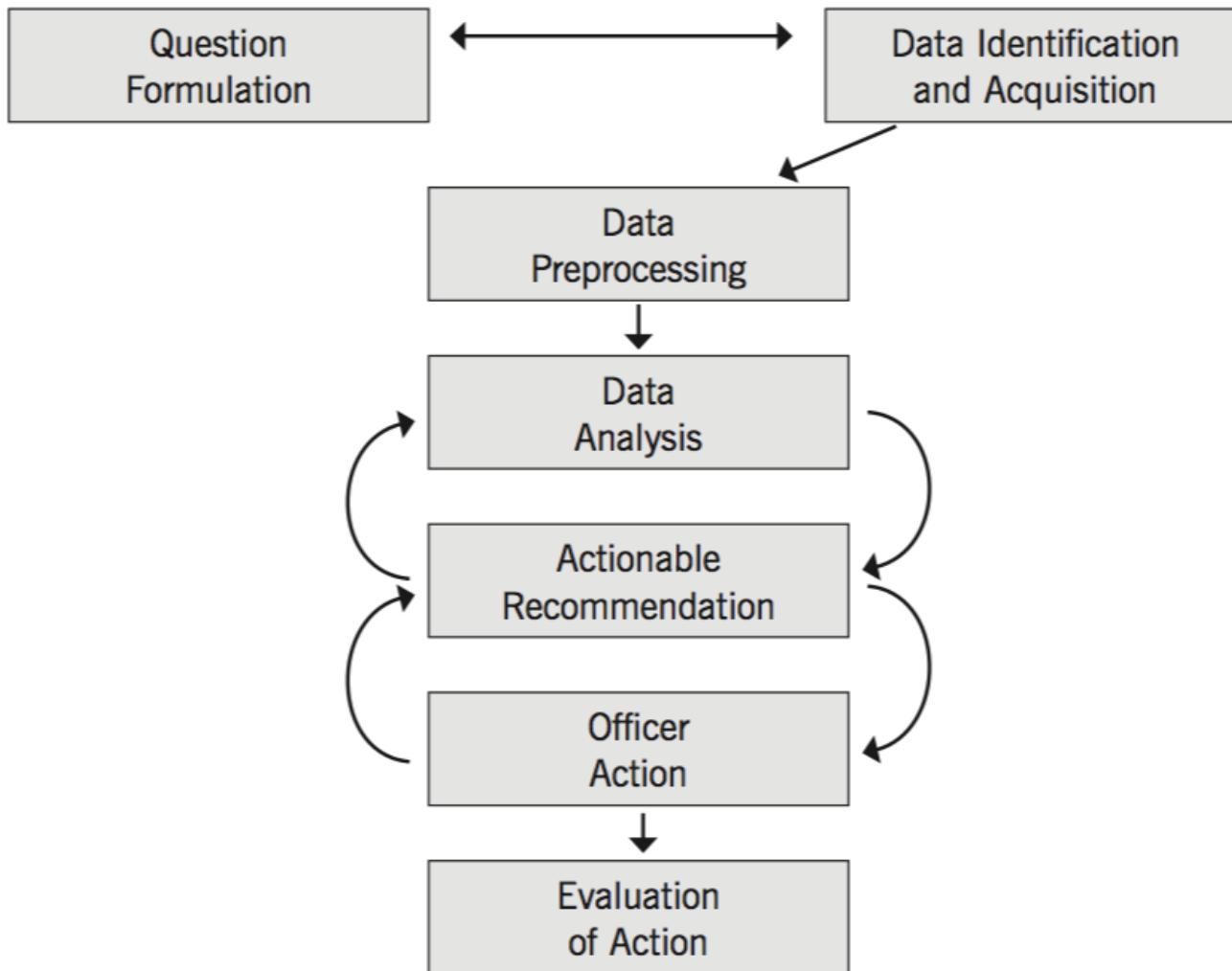
CASES

CASES



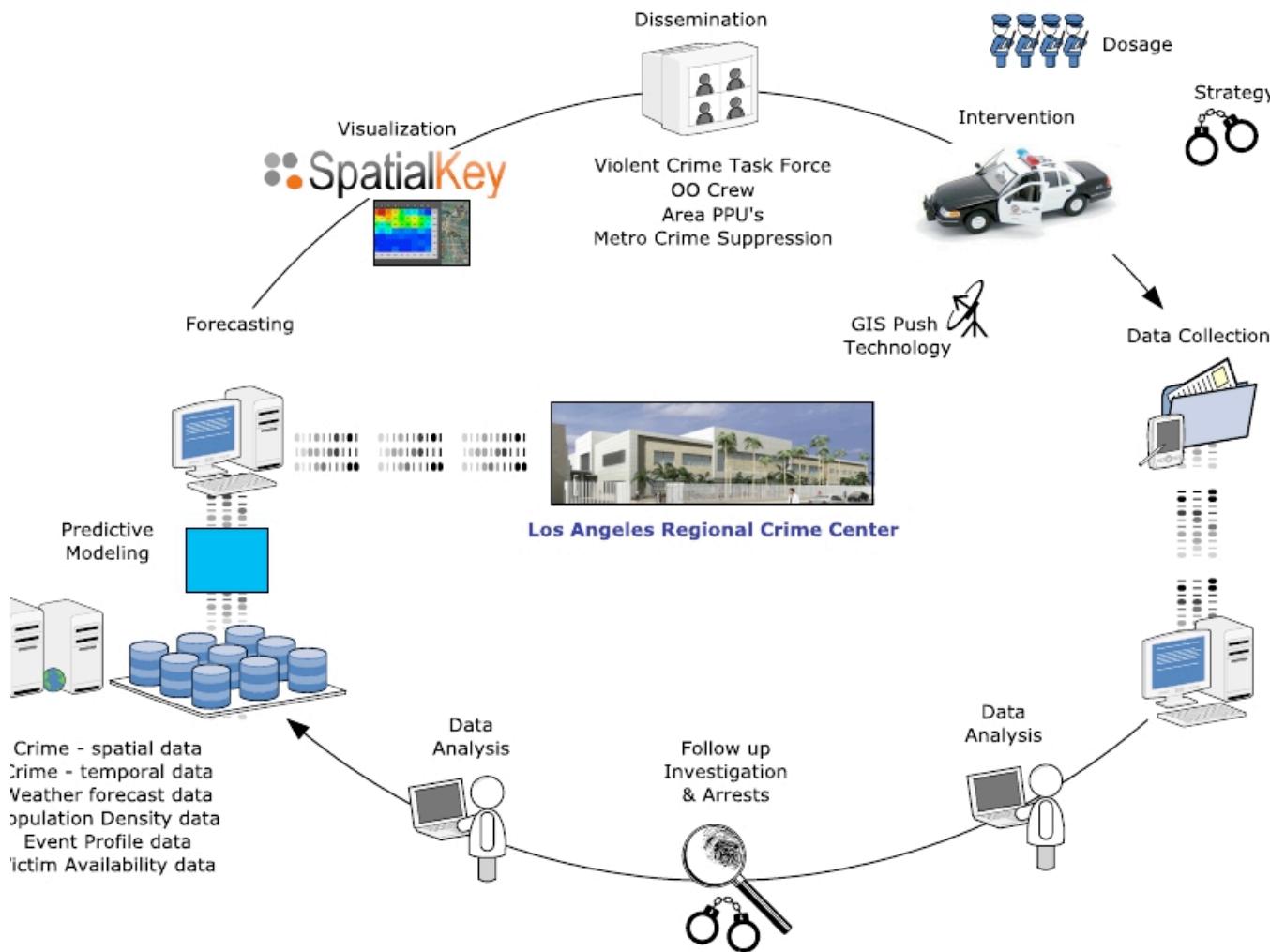
# Operational Challenges of Predictive Policing

52



# Case Study: Los Angeles Regional Crime Center

53



# Data Used in Predictive Policing

54

## ■ Three Categories

- **Spatial**
- **Temporal**
- **Social Network**

Spatial Variables	Temporal Variables	Social Network Variables
<p><b>Indicators of Areas with Potential Victims/Targets</b></p> <ul style="list-style-type: none"><li>• Shopping malls</li><li>• Property values</li><li>• Hotels</li><li>• Area demographics</li><li>• Population density</li><li>• Residential instability</li></ul> <p><b>Indicators of Escape Routes</b></p> <ul style="list-style-type: none"><li>• Highways</li><li>• Bridges</li><li>• Tunnels</li><li>• Public transportation</li><li>• Railways</li><li>• Dense foliage</li></ul> <p><b>Indicators of Criminal Residences</b></p> <ul style="list-style-type: none"><li>• Bars and liquor stores</li><li>• Adult retail stores</li><li>• Fast food restaurants</li><li>• Bus stops</li><li>• Public health information</li><li>• Areas with physical decay</li></ul>	<ul style="list-style-type: none"><li>• Payday schedules</li><li>• Time of day</li><li>• Weekend vs. weekday</li><li>• Seasonal weather (e.g., hot versus cold weather)</li><li>• Weather disasters</li><li>• Moon phases</li><li>• Traffic patterns</li><li>• Sporting and entertainment events</li></ul>	<ul style="list-style-type: none"><li>• Kinship</li><li>• Friendship</li><li>• Affiliation with an organization</li><li>• Financial transaction</li><li>• Offender/victim</li></ul>

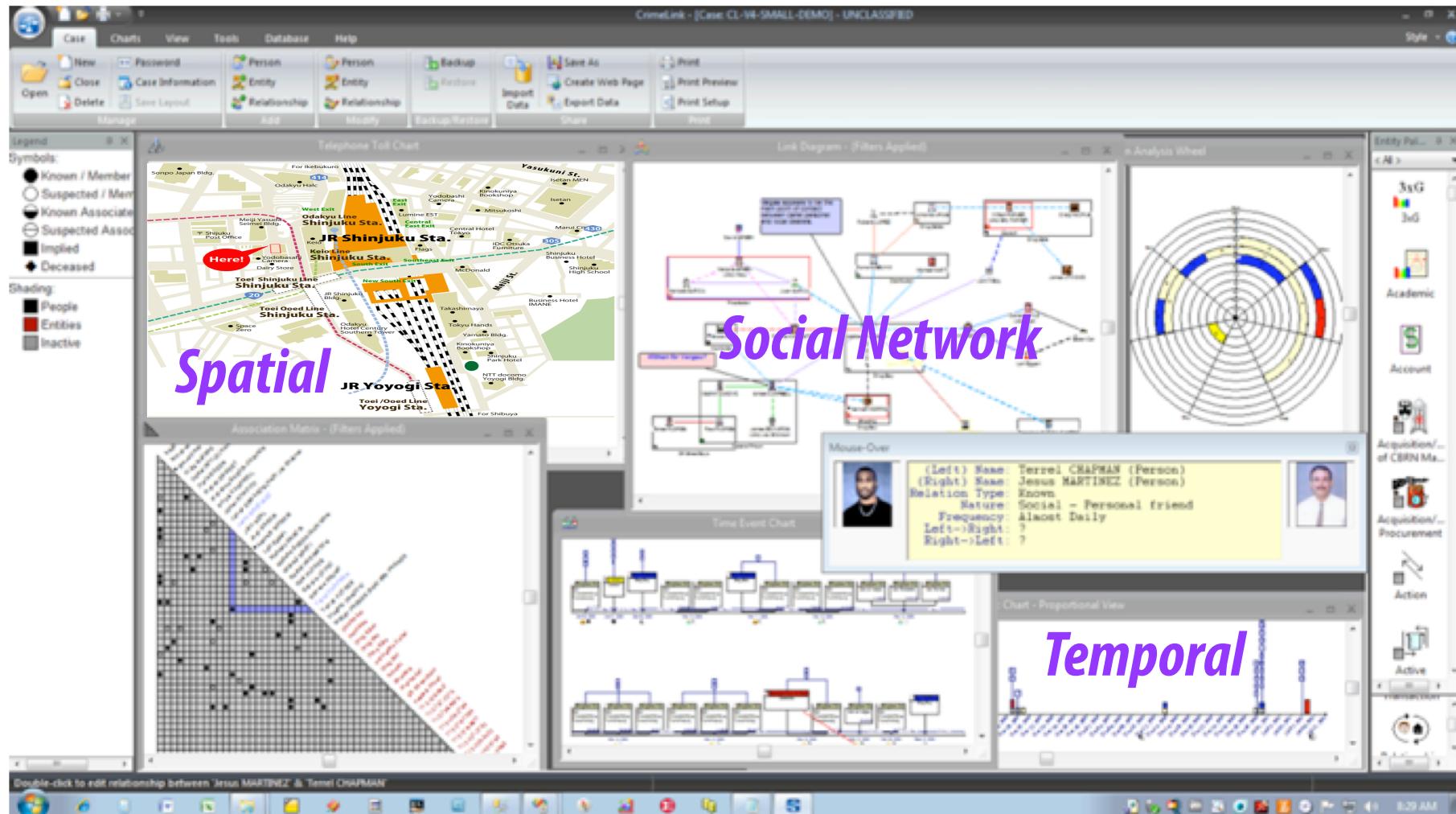
# The Transition from Data to Knowledge in a Police Agency

55

Data		Information		Knowledge		Result
Analysis	Individual Incident Reports in a records management system	Six of the reports are related in a series of robberies	Communication	Robbery series is prime topic of discussion in next detective's meeting	Strategy & Action	Robbery offender is apprehended
	Statistics showing number of officers per capita throughout the state	Your police department has 20% fewer officers per capita than average		Chief has this information in mind when preparing his budget proposal		Agency is granted additional officers by city
	Crime volume of current year compared to past years; individual records in RMS; jurisdictional information	Auto theft is up 20%, with most of the increase in Police Beat 5 on the midnight shift, probably influenced by new sports arena		Officers internalize this information and consider it when patrolling Beat 5		Auto theft is reduced

# Full Scale Analysis

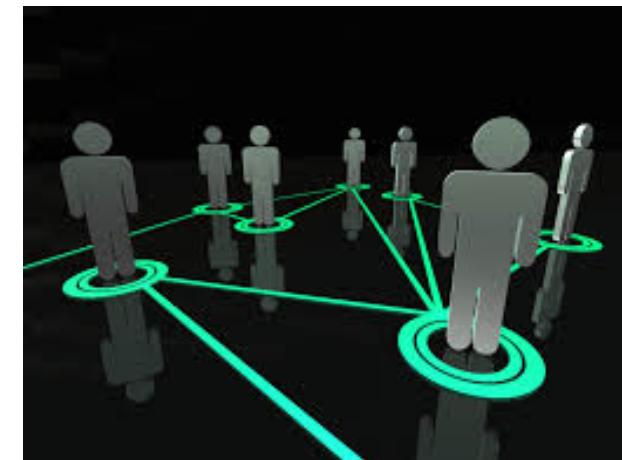
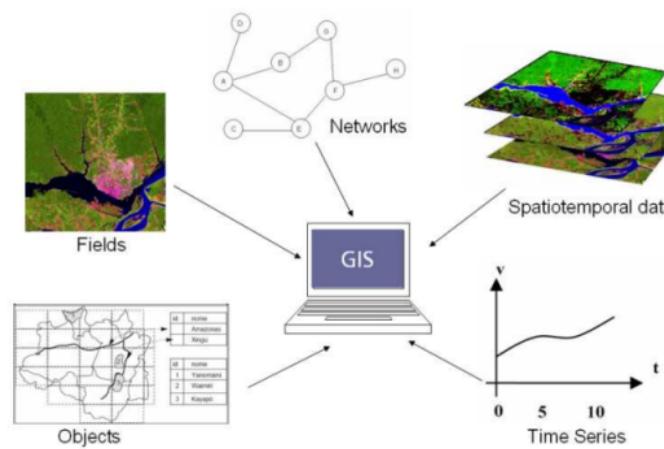
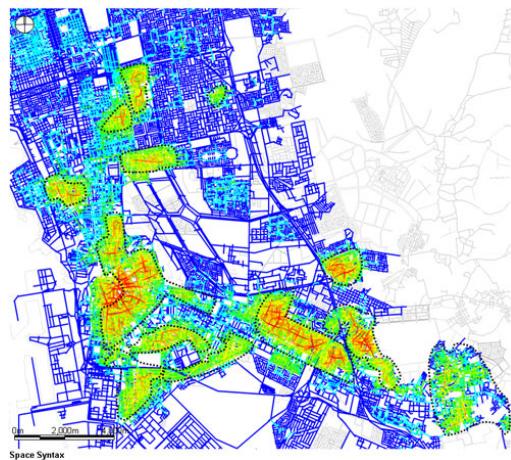
56



# Predictive Methodologies

57

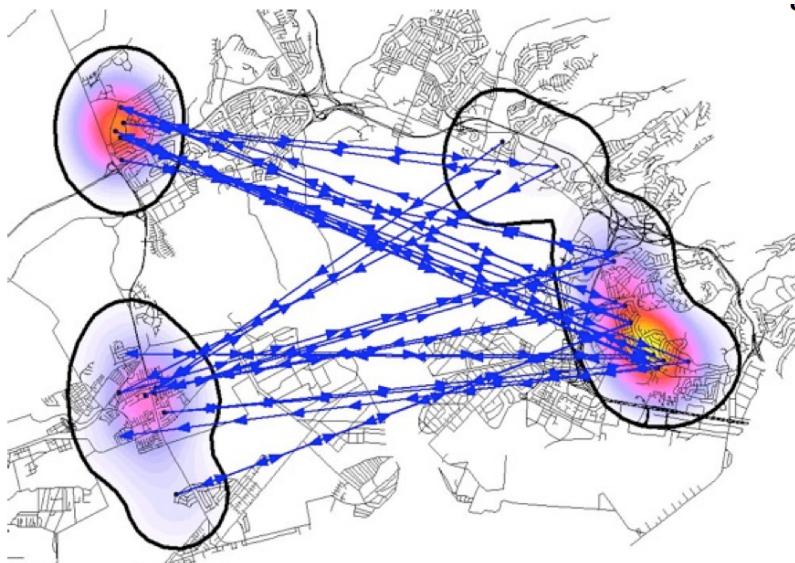
- **Analysis of space**
- **Analysis of time and space**
- **Analysis of social networks**



# Predictive Methodology One: Analysis of Space

58

- Point (or offense) locations
- Hierarchical clusters
- Partitioned clusters
- Fuzzy clusters
- Density mapping
- Risk-terrain modeling (RTM) clusters



# Predictive Methodology One: Analysis of Space

59

- Point (or offense) locations
  - *Theory of repeat victimization, 500x500 feet [PredPol Software]*
- Hierarchical clusters
  - use a nearest-neighbor technique [display the clusters: *ellipses, convex hulls*]
- Risk-terrain modeling (RTM) clusters



Gun shootings example



Gun shootings example

# Predictive Methodology Two: Analysis of Time and Space

60



# Predictive Methodology Two: Analysis of Time and Space

61

- CrimStat III : a software program (sociologist + National Institute of Justice)
  - spatial-temporal moving average (STMA)
    - the *average time and location* for a subset of incidents
  - correlated walk analysis (CWA) : temporal and spatial relationships between incidents
    - computing the correlation between *intervals* [time, distance, direction between two events]





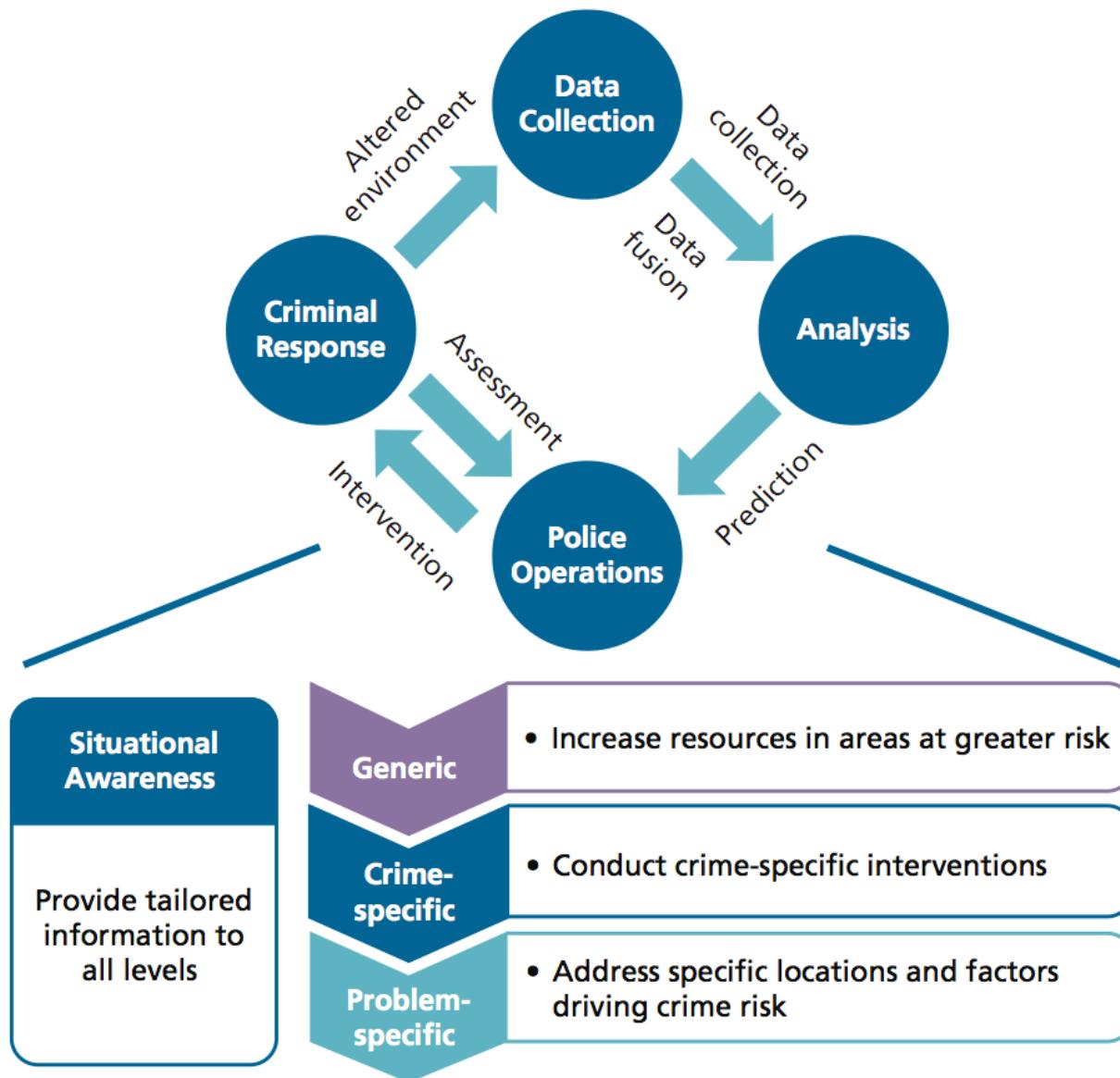
# Predictive Methodology Three: Analysis of Social Networks

63

- **Social Network Analysis (SNA) : *cutting edge* of crime analysis**
  - detect *persons of interest*, as opposed to *locations of interest*
  - identify individuals that are central to criminal organizations (eg., gangs and drug networks)
- **Building blocks of a social network: *relationships* between two actors**
- **In crime-fighting applications,**
  - SNA is frequently used to identify central nodes [high level of connectivity]
  - **Measures of *centrality***
    - ***degree*** : the number of links possessed by a node  
node's level of connectedness
    - ***closeness*** : the total distance from a node to all other nodes in the network  
ease of obtaining information from the network
    - ***betweenness*** : the number of instances a given node appears in the shortest path between other nodes  
relevance to the passage of information within the network

# Review: Prediction-Led Policing Business Process [2]

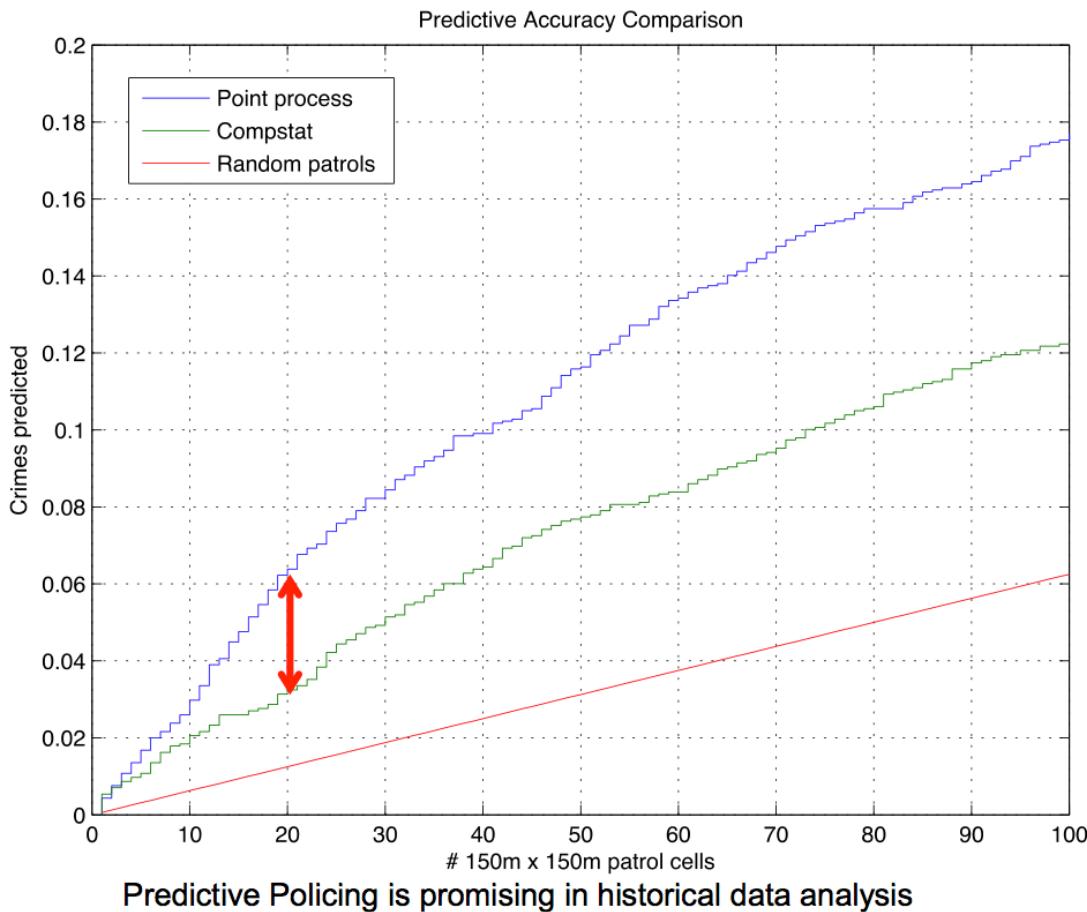
64



# Places on the Frontier of Predictive Policing in the United States

65

- *Los Angeles, California*
- **Santa Cruz, California**
- **Baltimore Country, Maryland**
- **Richmond, Virginia**
- **Memphis, Tennessee**



# Summary: Part II

66

- **Predictive Policing Concept**
- ***Mathematical Frameworks for Predictive Policing***
- ***Data Used in Predictive Policing***
- ***Predictive Methods for Predictive Policing***
- ***Predictive Policing Business Process***

The background is a reproduction of Raphael's fresco "The School of Athens" from the Vatican Palace. It depicts a gathering of ancient Greek philosophers in a grand, light-colored stone building with columns and arches. Numerous figures in classical robes are shown engaged in discussion or study, with some sitting on the floor and others standing. Large statues of figures like Apollo and Athena are visible on the left and right sides. The ceiling features a perspective drawing of a vaulted ceiling with a grid pattern and small windows showing a blue sky.

**Thank you**