

# Mind and Data Science

Herman Y Sutarto

East Continent Research Center

ECF Extension Course Filsafat  
Universitas Katolik Parahyangan





# DARE TO SLACK

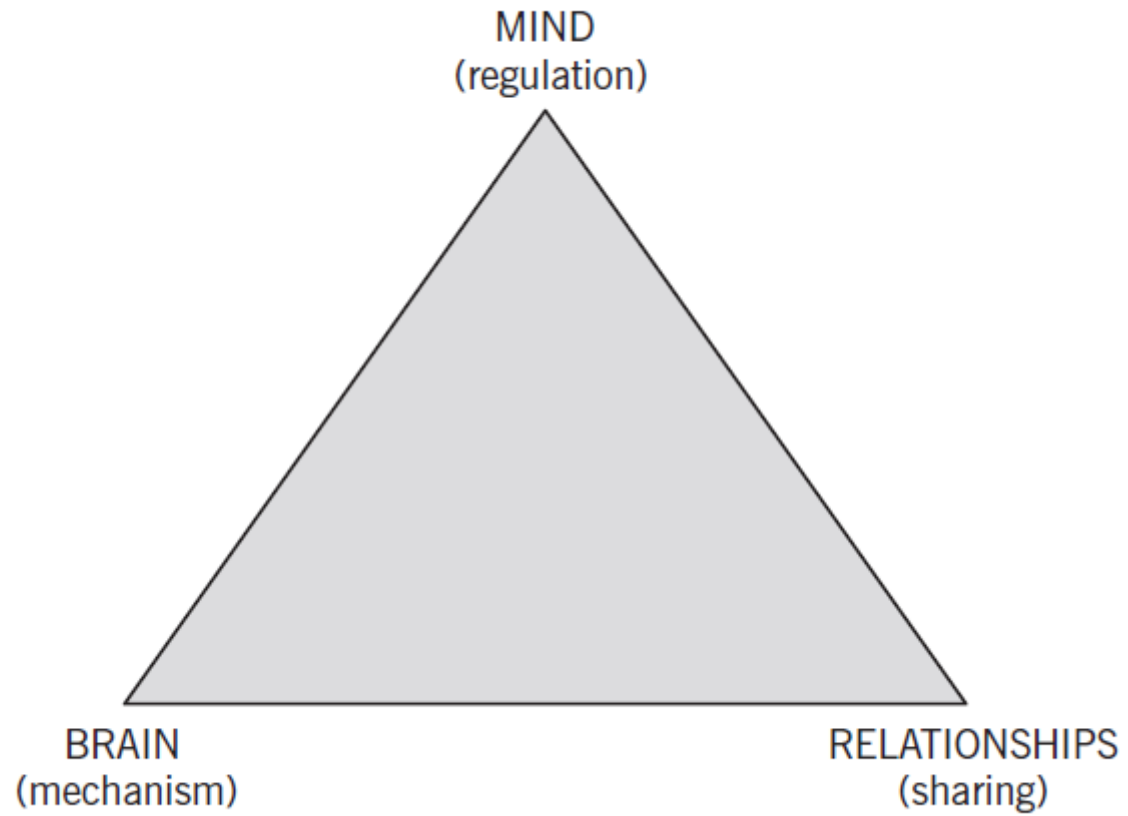
WHEN BIRDS FLY IN THE RIGHT FORMATION, THEY NEED ONLY EXERT HALF THE EFFORT.  
EVEN IN NATURE, TEAMWORK RESULTS IN COLLECTIVE LAZINESS.

[www.despair.com](http://www.despair.com)

# Video

# Emergence

- Novel behavior
- Properties of the whole
- **Cannot be predicted** from properties of the components that make up the system



The mind is an embodied and relational process that regulates the flow of energy and information.

The emergence of consciousness may be intimately related to the development of memory.

Memory is not a static thing, but an active set of processes.

Our internal experiences are constructive processes.

Experiences can shape not only what energy and information enters the mind, but also how the mind processes that information.

Our social experiences can directly shape our neural architecture.

Interpersonal experiences appear to have a direct effect on the development of explicit memory.

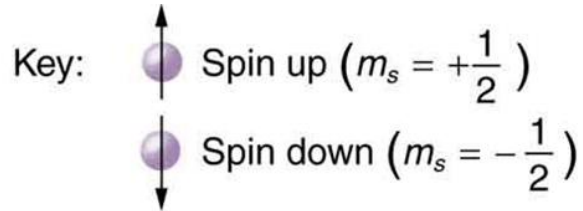
Relationships and the embodied brain are really part of one larger system.

Early experience shapes the regulation of synaptic growth and survival.

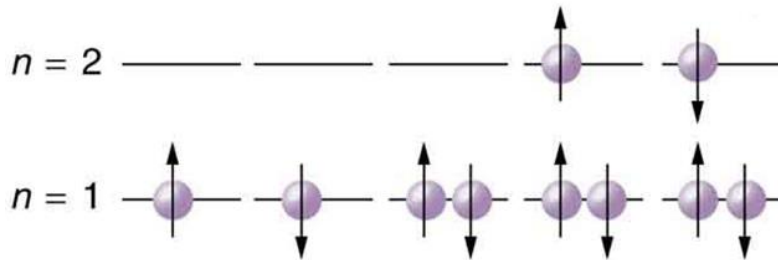
Interpersonal experiences continue to influence how our minds function throughout life.

*Consciousness* is the experience of being aware, the internal state of knowing that something is happening in the present moment.

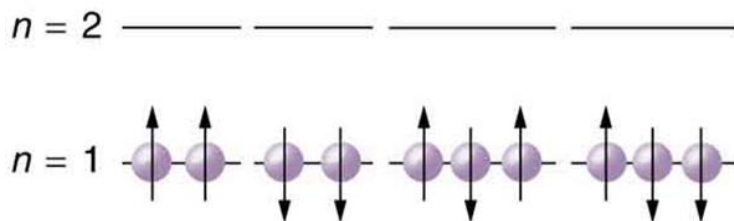
# The Pauli Exclusion Principle



Allowed



Not allowed



It is impossible for two electrons of a poly-electron atom to have the same values of the four quantum numbers:  $n$ , the principal quantum number,  $\ell$ , the angular momentum quantum number,  $m_\ell$ , the magnetic quantum number, and  $m_s$ , the spin quantum number. For example, if two electrons reside in the same orbital, and if their  $n$ ,  $\ell$ , and  $m_\ell$  values are the same, then their  $m_s$  must be different, and thus the electrons must have opposite half-integer spin projections of  $1/2$  and  $-1/2$ .

# Periodic Table

Group →	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	
↓ Period																			
1	1 H																		2 He
2	3 Li	4 Be											5 B	6 C	7 N	8 O	9 F		10 Ne
3	11 Na	12 Mg											13 Al	14 Si	15 P	16 S	17 Cl		18 Ar
4	19 K	20 Ca	21 Sc	22 Ti	23 V	24 Cr	25 Mn	26 Fe	27 Co	28 Ni	29 Cu	30 Zn	31 Ga	32 Ge	33 As	34 Se	35 Br		36 Kr
5	37 Rb	38 Sr	39 Y	40 Zr	41 Nb	42 Mo	43 Tc	44 Ru	45 Rh	46 Pd	47 Ag	48 Cd	49 In	50 Sn	51 Sb	52 Te	53 I		54 Xe
6	55 Cs	56 Ba	57 La *	72 Hf	73 Ta	74 W	75 Re	76 Os	77 Ir	78 Pt	79 Au	80 Hg	81 Tl	82 Pb	83 Bi	84 Po	85 At		86 Rn
7	87 Fr	88 Ra	89 Ac *	104 Rf *	105 Db	106 Sg	107 Bh	108 Hs	109 Mt	110 Ds	111 Rg	112 Cn	113 Nh	114 Fl	115 Mc	116 Lv	117 Ts		118 Og
				* 58 Ce	59 Pr	60 Nd	61 Pm	62 Sm	63 Eu	64 Gd	65 Tb	66 Dy	67 Ho	68 Er	69 Tm	70 Yb	71 Lu		
				* 90 Th	91 Pa	92 U	93 Np	94 Pu	95 Am	96 Cm	97 Bk	98 Cf	99 Es	100 Fm	101 Md	102 No	103 Lr		



## Morowitz's Three Claims

- (i) PEP is a *nondynamical* principle, but it influences the dynamical behavior of electrons.
- (ii) PEP has *nothing to say about the behavior of individual electrons*.
- (iii) PEP is *unrelated to the other laws of physics*.



# How does the [mind] work ?

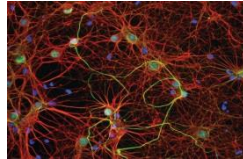
Molecules/  
Genes



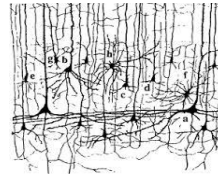
Neurons/Synapses



Neuronal circuits



Systems of neurons Brain regions



Data & state  
Representation  
+ algorithms



Behavior/Cognition



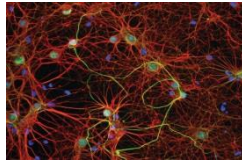
Molecules/  
Genes



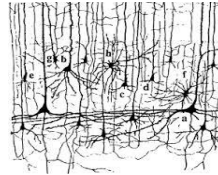
Neurons/Synapses



Neuronal circuits



Systems of neurons Brain regions



Data & state  
Representation  
+ algorithms



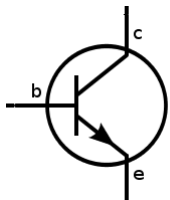
# How does the [mind] work ?

"You, your joys and your sorrows, your memories and your ambitions, your sense of personal identity and free will, are in fact ~~no more than~~ the behavior of a vast assembly of nerve cells and their associated molecules"

-- Francis Crick (Co-discover of the structure of DNA)

# How does my favorite app work?

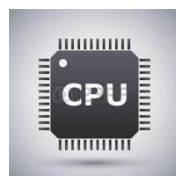
Transistor



IC



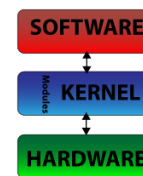
CPU/GPU



Motherboard



Kernel

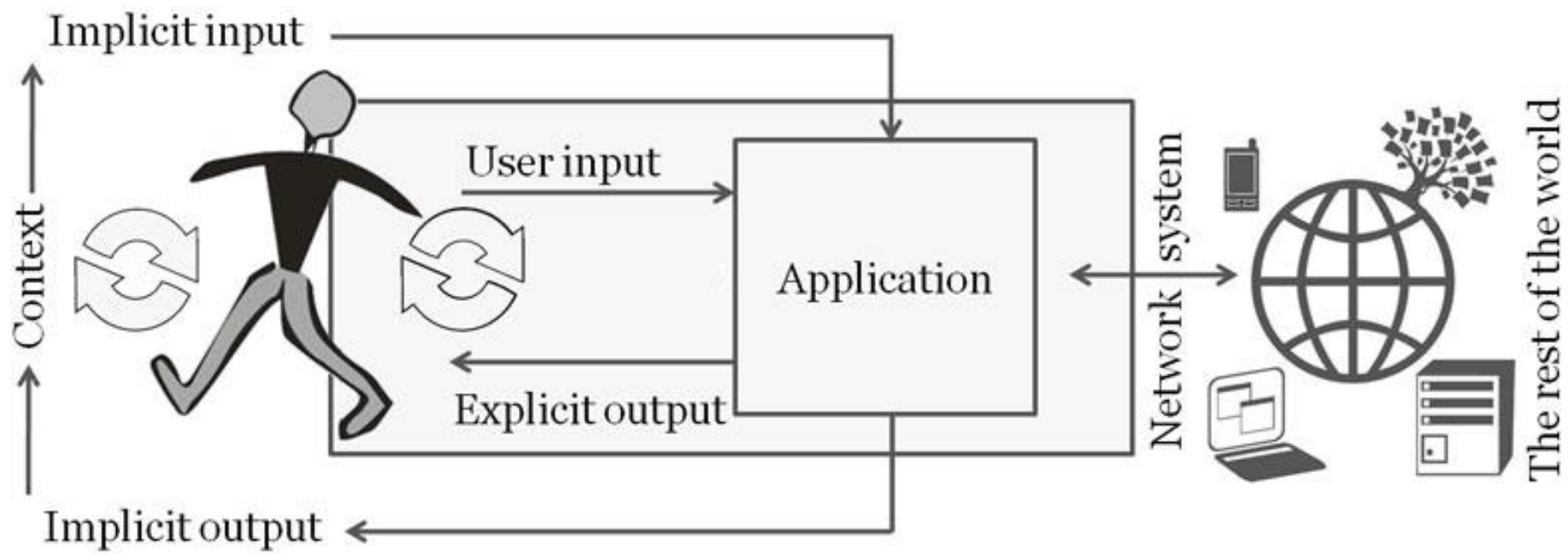


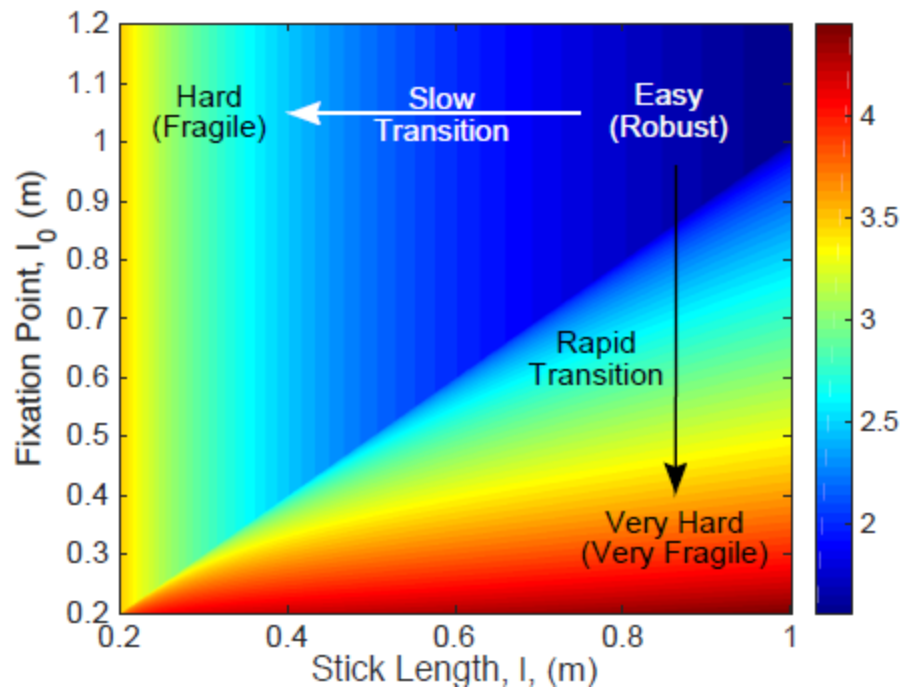
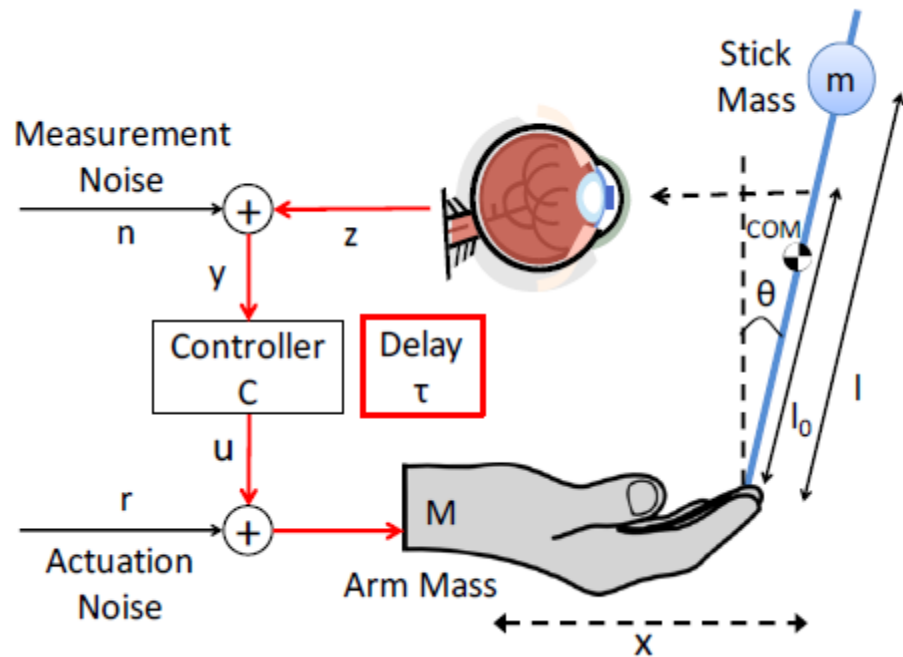
Algorithm application



User Interface







# Sigmund Freud


## The Human Mind



# Influ-Venn-Za

Who can catch which flu?

April 22nd 2013 - suspected mutation of an avian virus emerged outside Shanghai, China. Human fatality rate is unknown but out of 130 people infected, 40 have died (30%).

 Influenza Type A is divided into H & N strains (i.e. H1N1) referring to different combinations of:

**H = hemagglutinin**  
(binds to cells)

**N = neuraminidase**  
(surface enzyme)

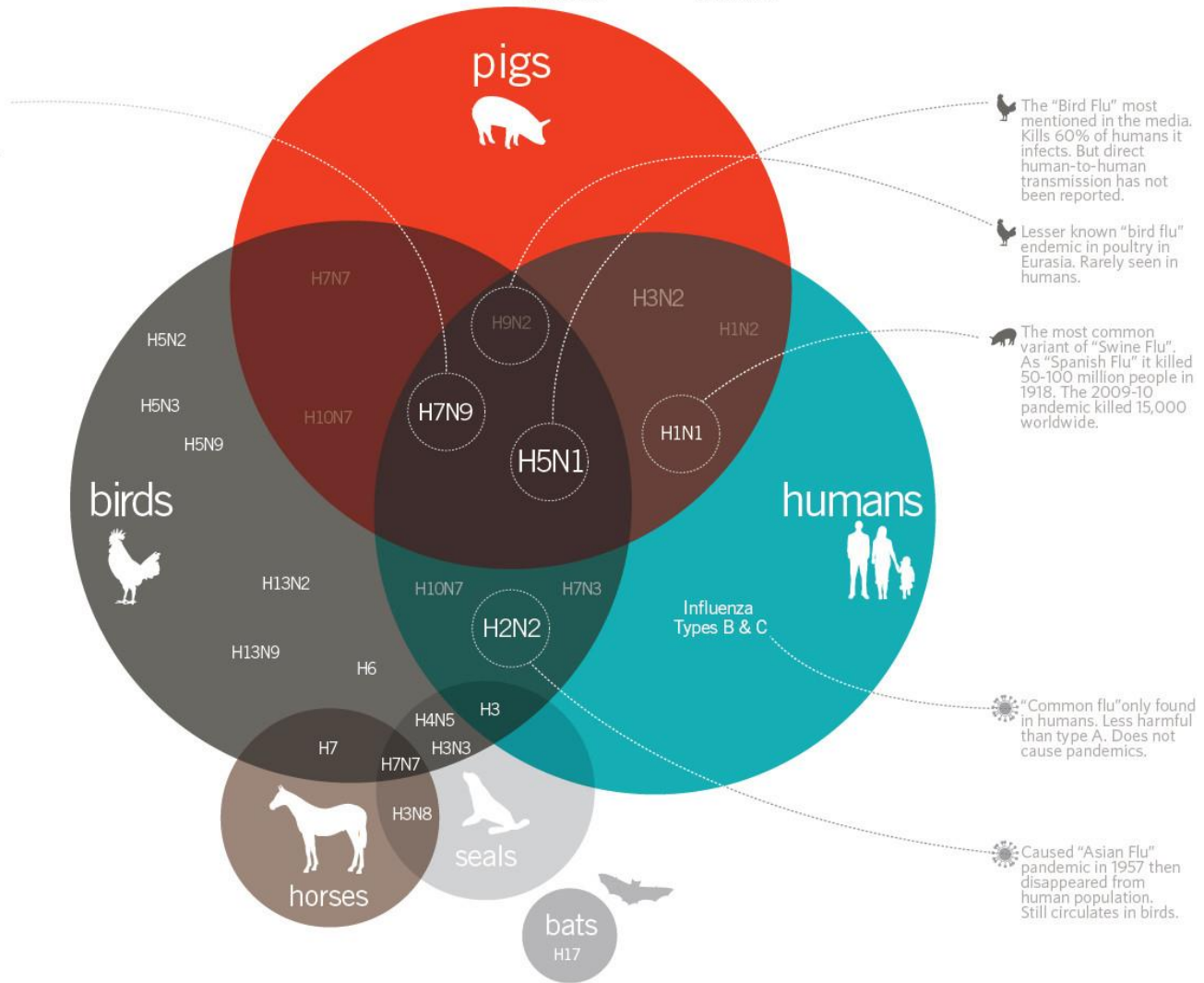
text SIZE  
= human fatality rate

LIGHT TEXT  
= rarely infects humans

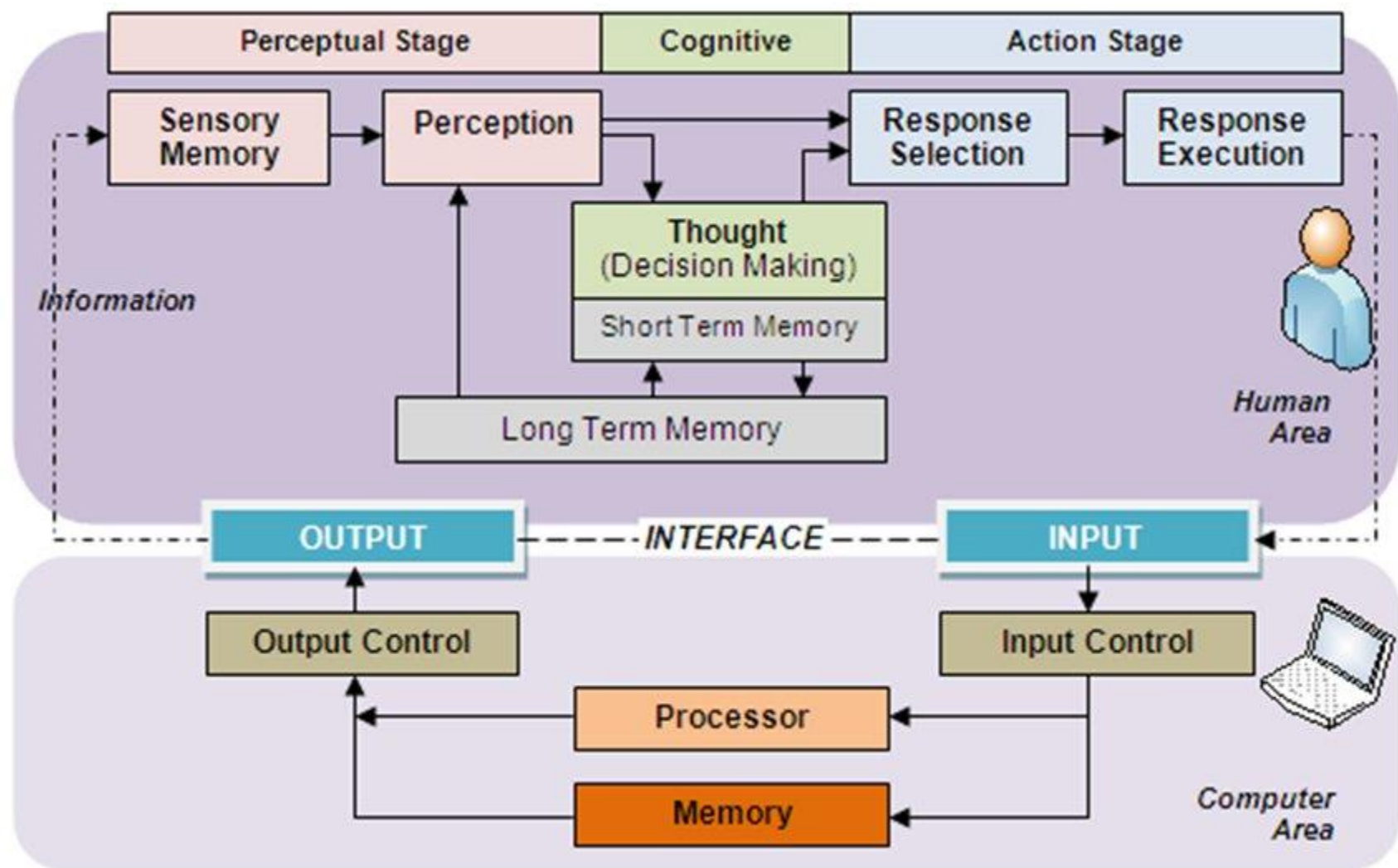
Concept & Design: David McCandless  
Research: Ella Hollowood // Additional design: Phillipa Thomas  
Version 1.22 / August 2013

 Pigs often a source of flu pandemics as they can be infected by bird, human & swine flus.

Worse-case, they act as a bridge for newly evolved virus strains to cross from birds to humans.














Culture is the widening of  
the mind and of the spirit.

Jawaharlal Nehru

© 2015

Hundreds of  
billions of dollars  
are spent every  
year to control the  
public mind.

*~Noam Chomsky*



# Data Science

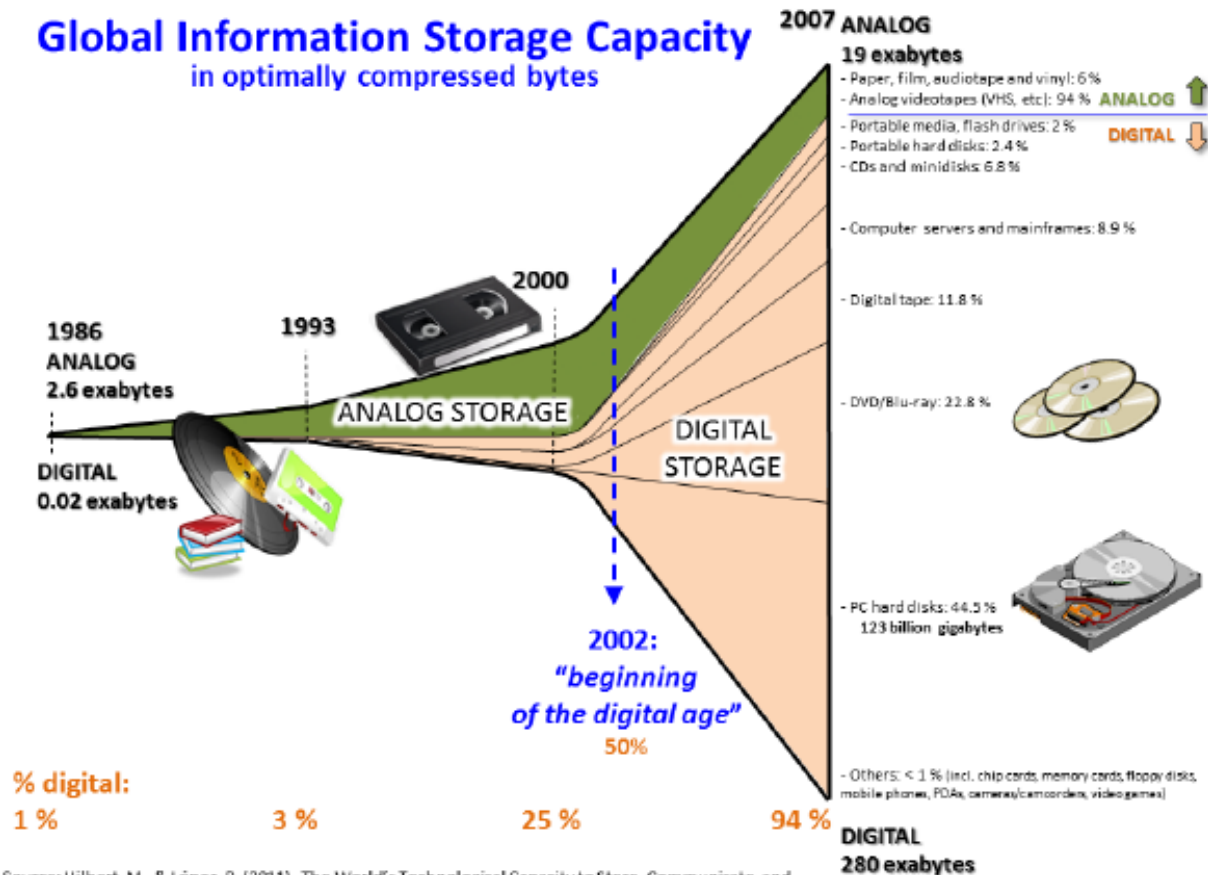
# The (Long) History of Data Processing

- (...) Manual Data Processing
  - (1832) Punch cards
  - (1936) Turing Machine
  - (1944) First IBM Computers (SSEC)
  - (1945) Von Neumann Architecture
  - (1950s) Sort & Search Algorithms
  - (1950s) Heuristics Methods
  - (1936) Pattern Recognition
  - (1951) The First Neural Network Machine
  - (1955) Concerns on data explosion by Fremont Rider, Wesleyan University Librarian
    - (1960s) DBMS
    - (1960s) Data Analysis Methods (Bayesian, Time Series, Stochastic, ...)
      - (1968) Knuth – The Art of Computer Programming
      - (1970s) Relational DBMS
        - (1974) SQL
        - (1975) First PC (MITS Altair 8800)
  - (1990s) Data Mining / KDD
  - (1990s) Complex-Event Processing
  - (1990s) Data Stream Processing
  - (1990s) Social Network Analysis
  - (1998) The Term “Big Data” was first coined by John Mashey
  - (1999) Internet of Things
  - (2001) Volume, Velocity, Variety by Doug Laney
  - (2001) “Data Science” by William Cleveland
  - (2004) MapReduce
  - (2009) No-SQL
  - (2011) Global Information Storage Capacity grows at 25% annually by Martin Hilbert, Priscila López



# How Big is Big?

## Global Information Storage Capacity in optimally compressed bytes



Source: Hilbert, M., & López, P. (2011). The World's Technological Capacity to Store, Communicate, and Compute Information. *Science*, 332(6025), 60–65. <http://www.martinhilbert.net/WorldInfoCapacity.html>

As of 2012, every day 2.5 exabytes ( $2.5 \times 10^{18}$ ) of data are generated  
[IBM, "What is Big Data"]



# VOLUME

SCALE OF DATA

**40 ZETTABYTES**  
(43 TRILLION GIGABYTES)



OF DATA WILL BE CREATED BY 2020 AN INCREASE OF 300 TIMES FROM 2005

**2.5 QUINTILLION BYTES**  
(2.3 TRILLION GIGABYTES)  
OF DATA ARE CREATED EACH DAY

**6 BILLION PEOPLE**  
HAVE CELLPHONES



WORLD POPULATION  
7 BILLION

MOST COMPANIES  
IN THE U.S. HAVE AT LEAST

**100 TERABYTES**  
(100,000 GIGABYTES)  
OF DATA STORAGE

# VELOCITY

ANALYSIS OF  
STREAMING DATA

THE NEW YORK STOCK  
EXCHANGE CAPTURES



**1 TB OF TRADE  
INFORMATION**  
DURING EACH SESSION

BY 2016, IT IS PROJECTED THERE  
WILL BE



**100 BILLION  
NETWORK  
CONNECTIONS**



ALMOST 2.5 CONNECTIONS  
PER PERSON ON EARTH

MODERN  
CARS HAVE  
CLOSE TO



THAT MONITOR  
ITEMS SUCH AS FUEL LEVEL AND TIRE  
PRESSURE

# VARIETY

DIFFERENT FORMS  
OF DATA

AS OF 2011, THE GLOBAL SIZE OF  
DATA IN HEALTHCARE WAS  
ESTIMATED TO BE



**150 EXABYTES**  
(161 BILLION GIGABYTES)



**30 BILLION PIECES OF  
CONTENT**  
ARE SHARED ON FACEBOOK  
EVERY DAY



**>4 BILLION HOURS OF VIDEO**  
ARE WATCHED ON YOUTUBE  
EACH MONTH



**400 MILLION TWEETS**  
ARE SENT PER DAY BY ABOUT  
200 MILLION MONTHLY ACTIVE  
USERS

BY 2014, IT'S ANTICIPATED  
THERE WILL BE



**420 MILLION  
WEARABLE, WIRELESS  
HEALTH MONITORS**

# VERACITY

UNCERTAINTY OF  
DATA

**1 IN 3**  
**BUSINESS LEADERS**



DON'T TRUST THE  
INFORMATION THEY USE TO  
MAKE DECISIONS



**OF  
RESPONDENTS**



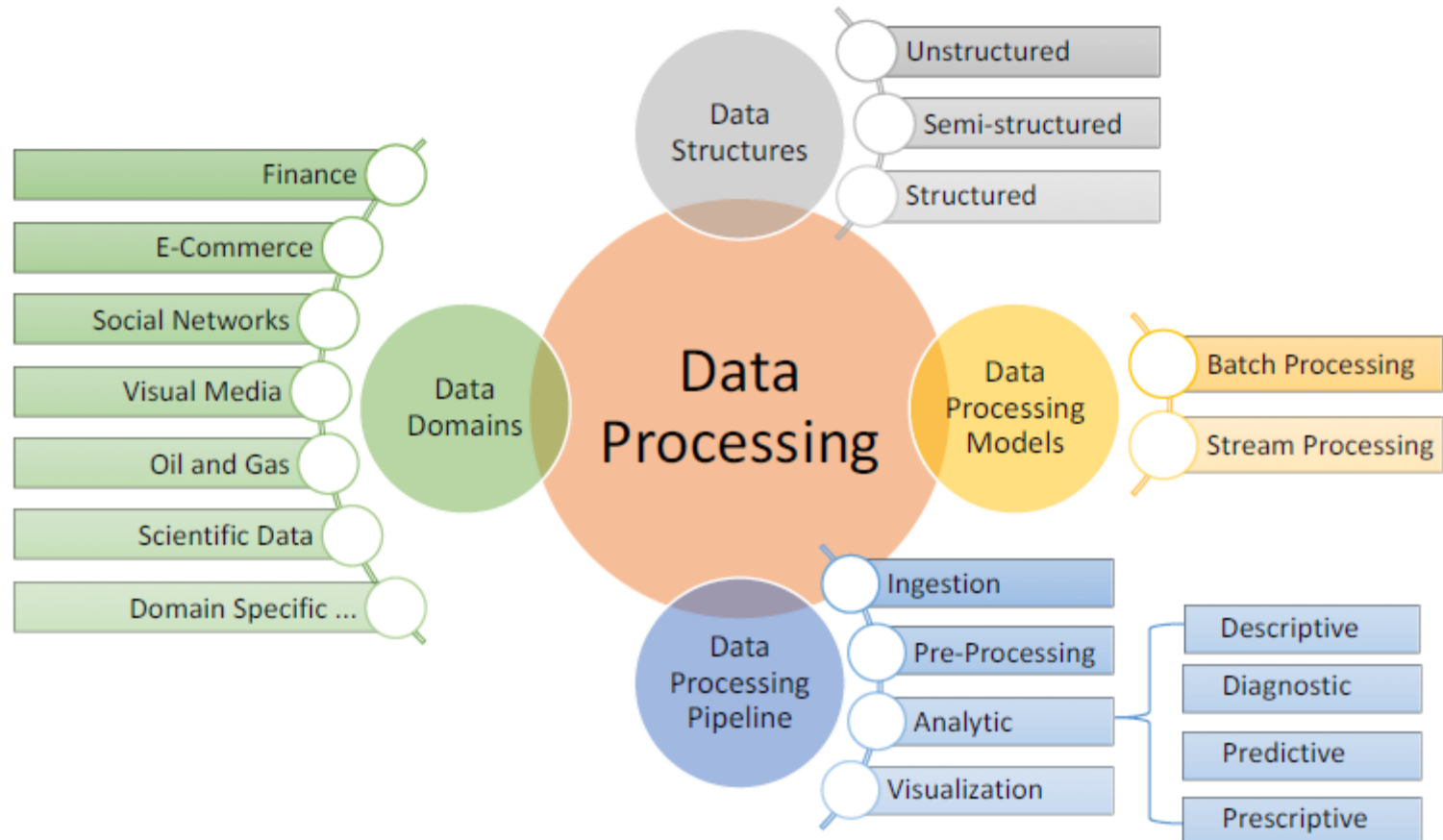
IN ONE SURVEY WERE UNSURE OF  
HOW MUCH OF THEIR DATA WAS  
INACCURATE

POOR DATA QUALITY  
COSTS THE US ECONOMY  
AROUND

**-\$3.1 TRILLION A YEAR**



# The Taxonomy of Data Processing





# Data Structures

Unstructured Data



Around 80-90% of all potentially usable business information may originate in unstructured form [Merill Lynch, 1998]

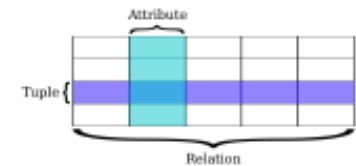


Semi-structured

Structured Data



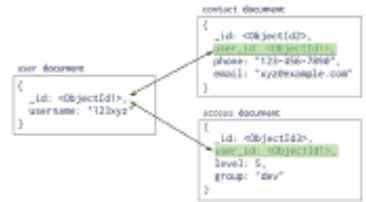
Relational



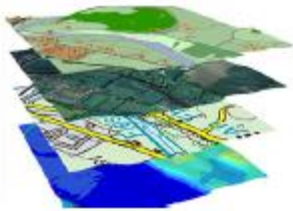
Graph DB → Network Model DB  
e.g., CODASYL (1969)



Document DB → Hierarchical Model DB  
e.g., IBM IMS (1969)



# Other Data Structures



Spatial / Geospatial Data



Spatio-temporal Data



e.g., Moving Objects



Biological Data

1000 Genomes Project → >200 Terabytes  
<https://aws.amazon.com/1000genomes/>

Million Human Genomes project → ???

# Data Stream

- ... is an ordered sequence of instances that in many applications can be read only once or a small number of times using limited computing and storage capabilities\*
- Data Stream Processing Applications:
  - IoT applications
  - Live datamart
  - Pattern mining on live data
  - ...



# Data Stream Characteristics

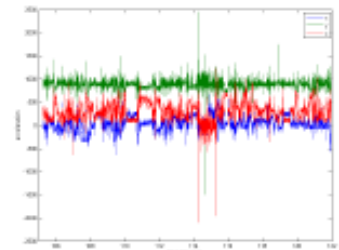
- Continuous flow of data
- Infinite length
  - not just **BIG** but **“UNLIMITED”**
  - impractical to store all data
- Examples:



Call Detail Records (CDR)



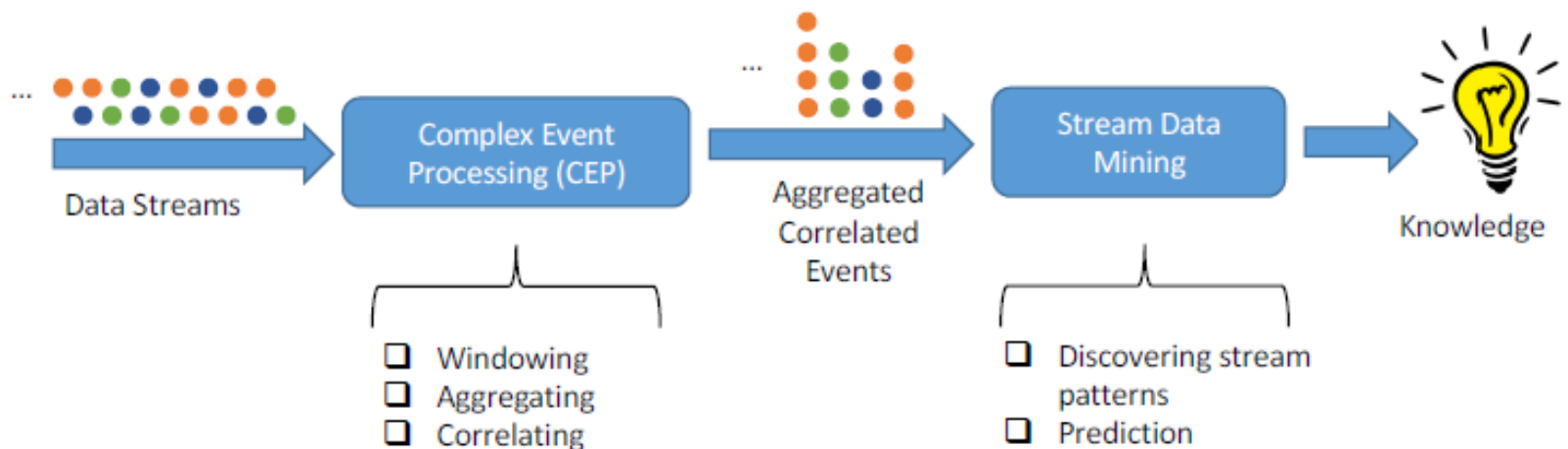
Tweets



Sensor Data

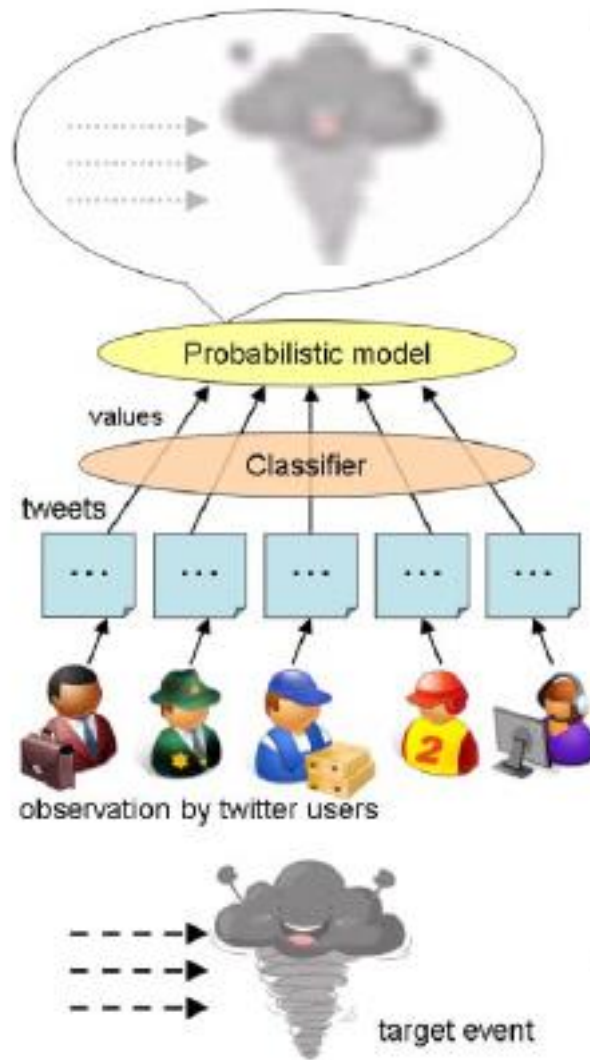
# Data Stream Processing

An example of a typical data stream processing flow

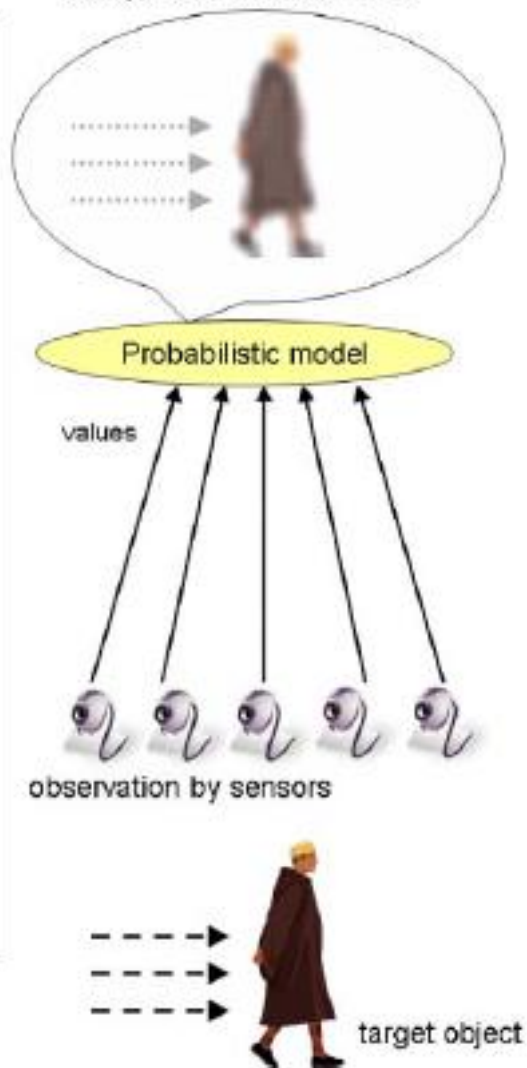


# Digital vs Analog World

Event detection from twitter

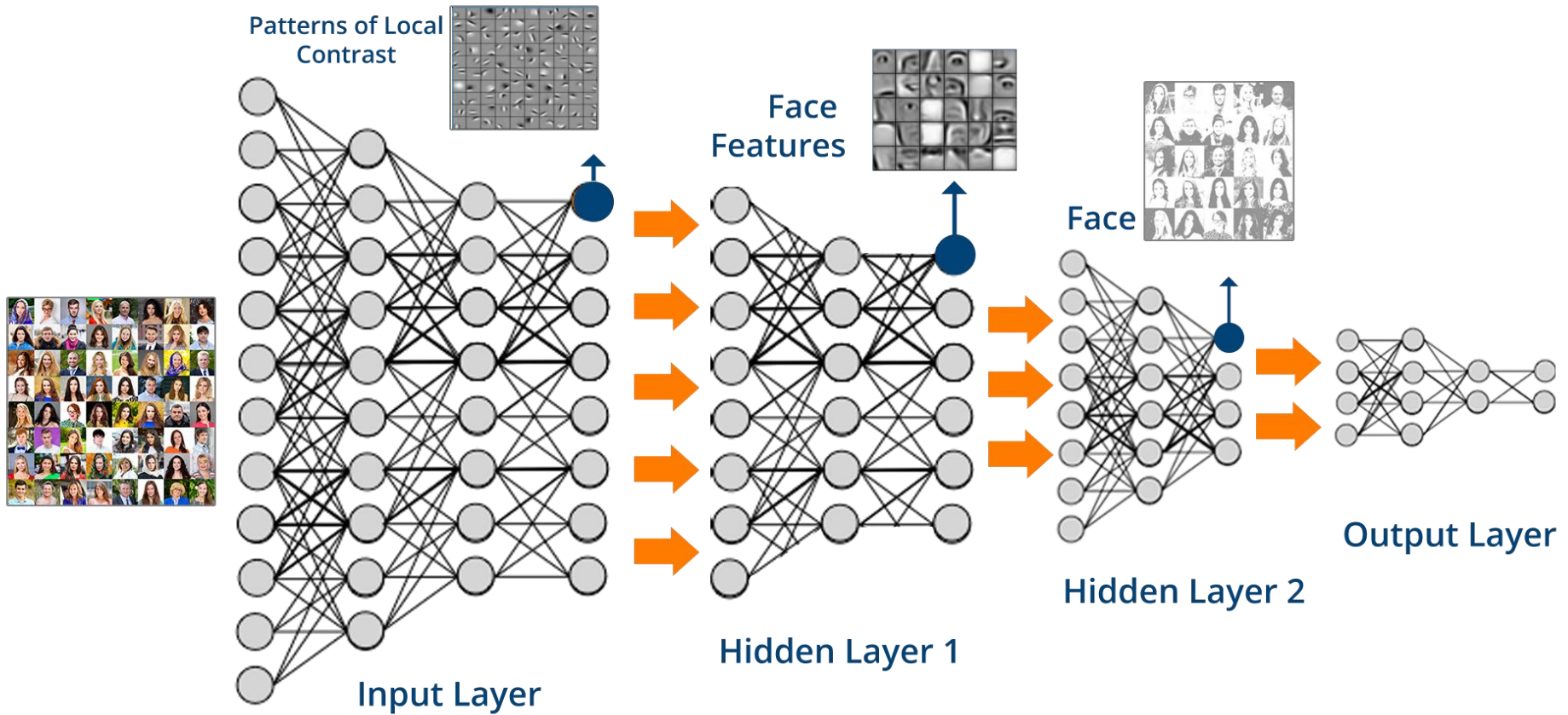


Object detection in ubiquitous environment





# Deep Learning



# Big Data

## Correlation vs Causation

Is correlation more important than causation?

Correlations play an important role as heuristic devices [but] have to be further analyzed [. . .] to assign them a meaning”

The correlations may not tell us precisely why something is happening, but they alert us that it is happening.



# Big Data

## Correlation vs Causation

**Thomas Kuhn:**

Anomalies, by definition, For such discoveries to occur, establishing that there is something that does not match our expectations is not enough. We have also to find out what it is. This process does not arise directly from data or numbers, but rather from a change in how we look at them, and it involves a reassessment of our beliefs and methodologies.

# Big Data

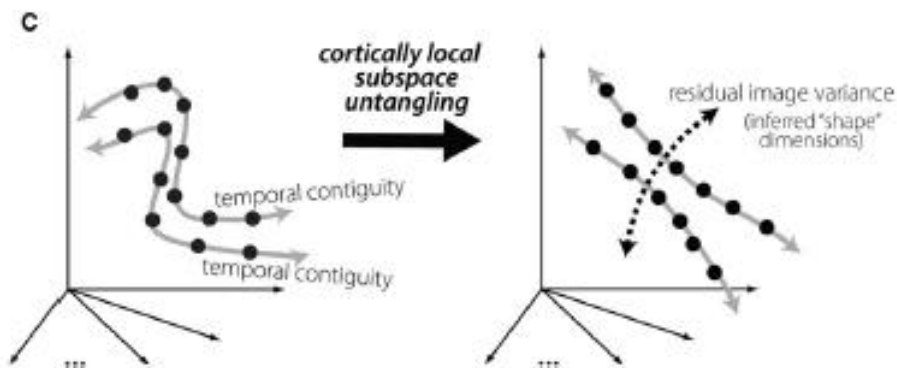
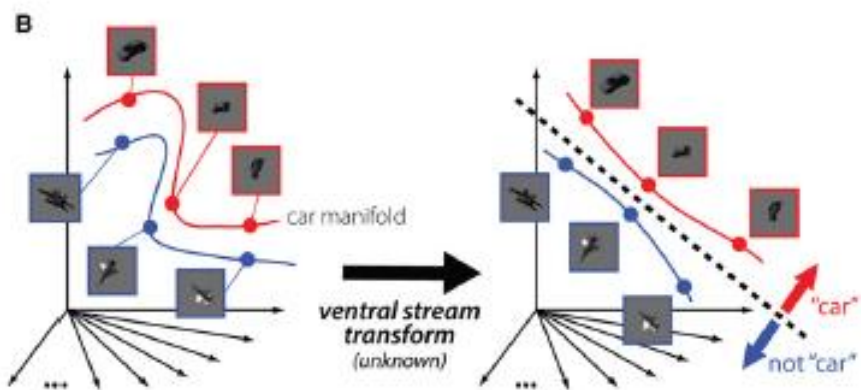
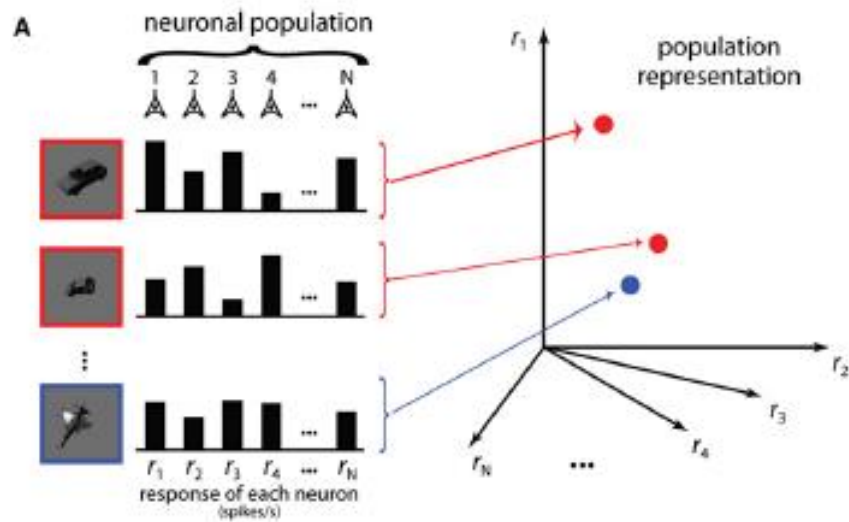
## Correlation vs Causation

“Big Data, distributed computing and sophisticated data analysis all played a crucial role in the discovery of the Higgs boson [. . .] But the discovery of the Higgs boson was not data-driven.”

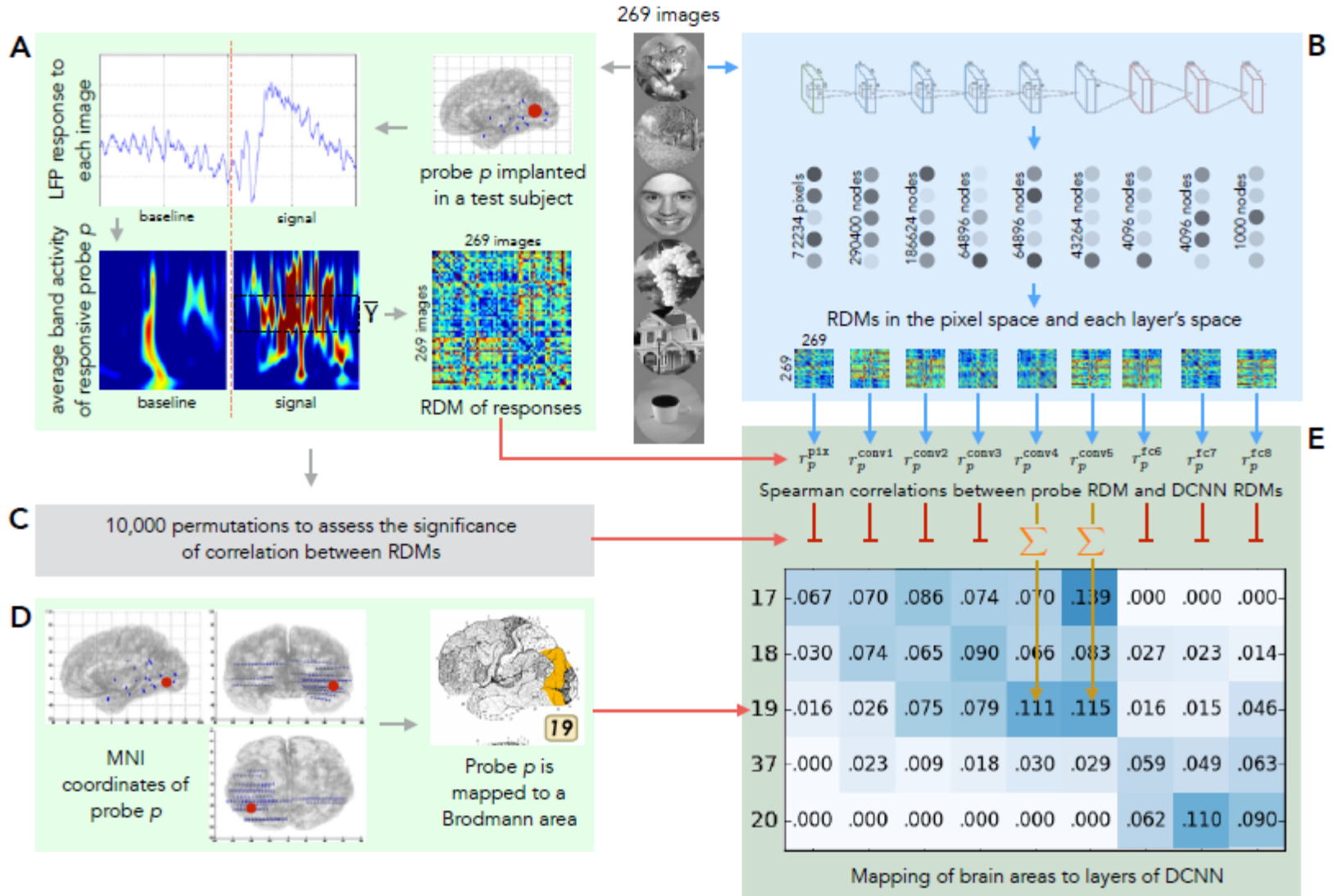
“The data-driven approach constitutes a novel tool for scientific research. Yet this does not imply that it will supersede cognitive and methodological procedures. . . .”

# Some Applications

How does the Brain solve visual  
objective recognition



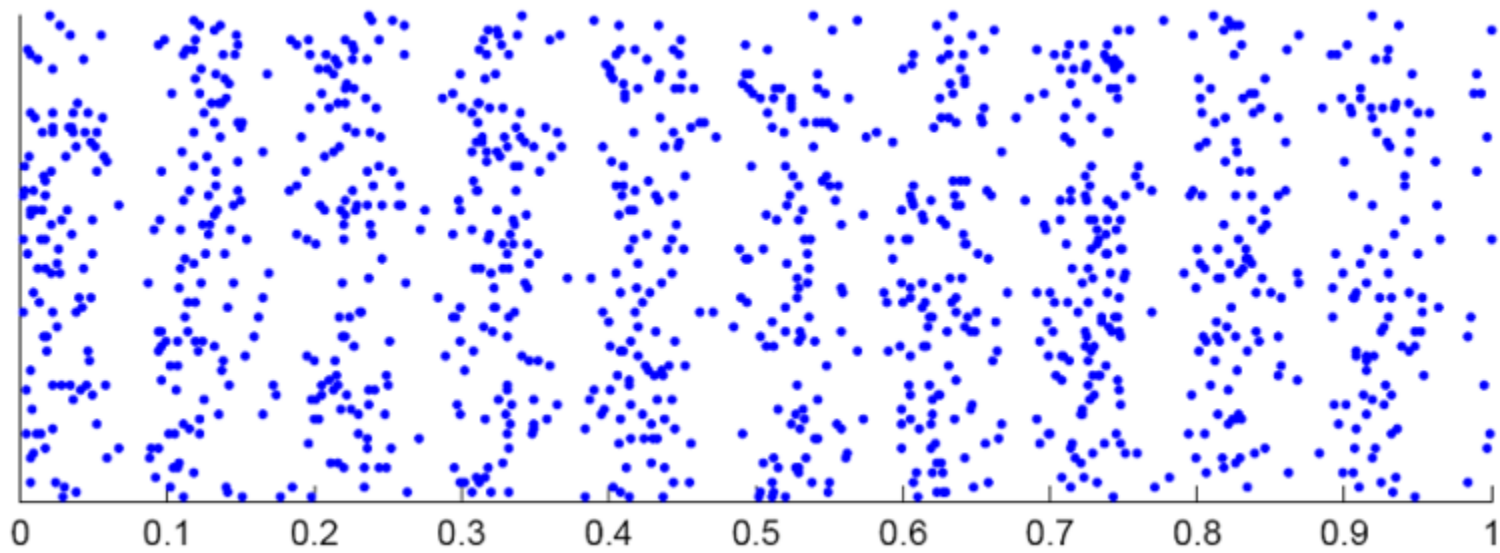
# Natural images vs artificial vision system



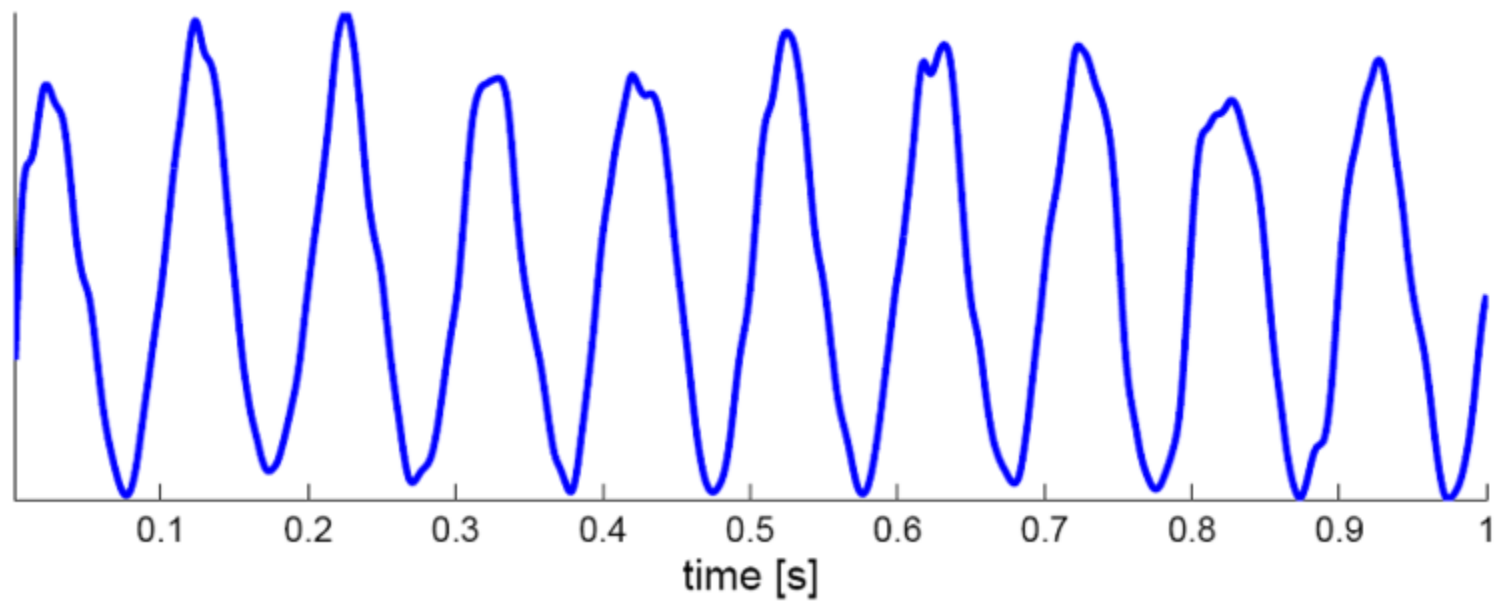
# The phenomenon of synchronization

- Circadian rhythms
- Electrical generators
- Heart, intestinal muscles
- Menstrual cycles
- Fireflies
- Applause (esp in Eastern Europe)

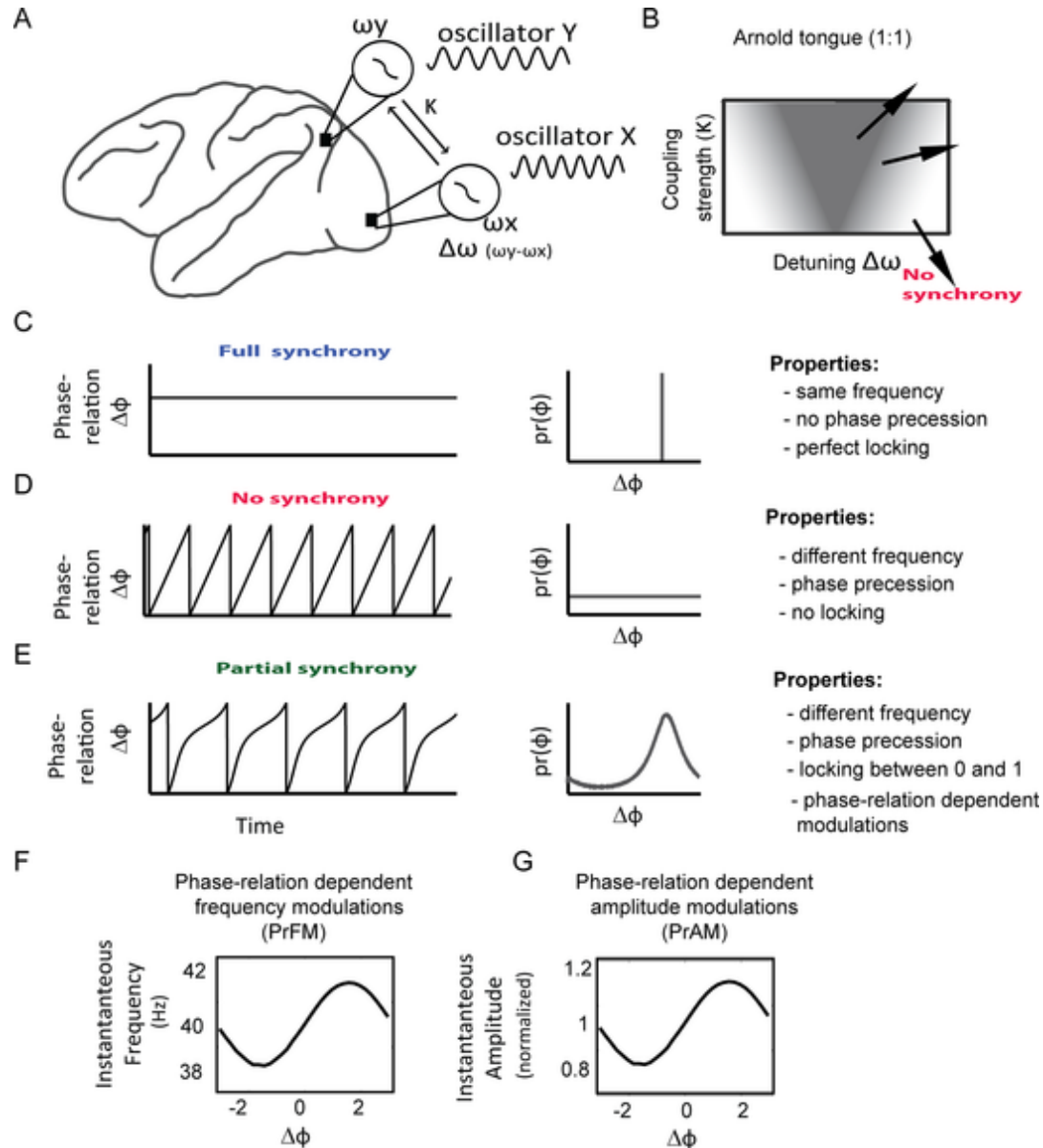
**neuronal spiking**



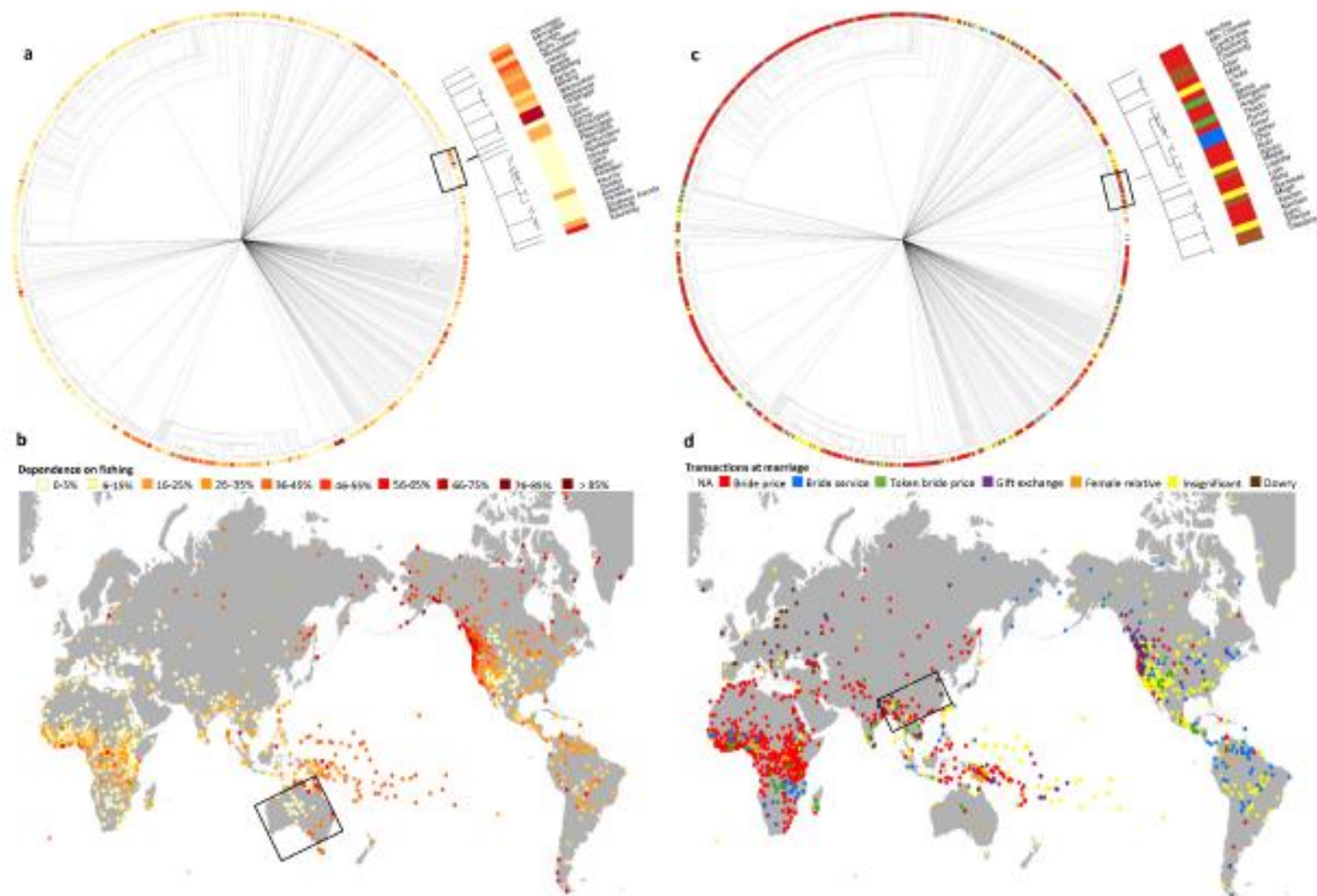
**local field potential**



The dynamic emergence of **coherent physiological activity**, such as phase-locked high-frequency electromagnetic oscillations,



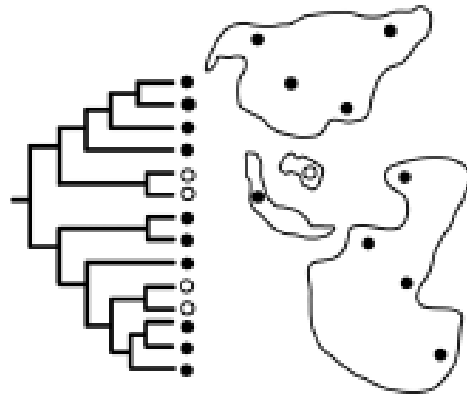




**Fig 1. D-PLACE links cultural information to language classifications and phylogenies (a, c) and to geographic locations and environmental features (b, d). This allows users to consider the relative influence of cultural ancestry, spatial proximity, and environment on diverse cultural practices. For example, panels a and b illustrate variation among societies in their dependence on fishing relative to other subsistence activities, based on data from the Ethnographic Atlas (EA) [11–15] and the Binford Hunter-Gatherer dataset [16, 17]. Panels c and d highlight diversity in the most common economic transaction at marriage, based on data from the EA. In addition to providing global results, D-PLACE allows users to focus a search on a particular geographic region or linguistic family. Here, results for societies speaking Pama-Nyungan languages (a, b) or Sino-Tibetan languages (c, d) are magnified and outlined in black boxes on the global tree and map.**

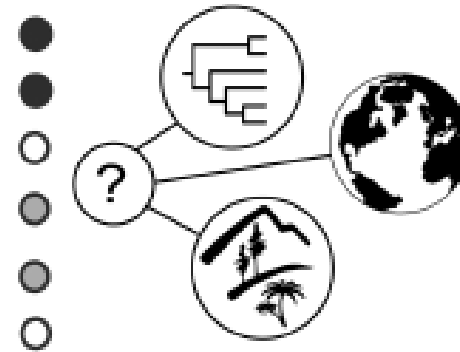
### a. Exploratory

How are features distributed across societies?



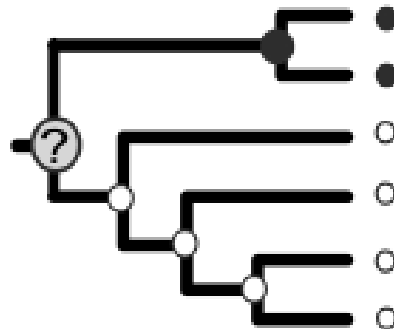
### b. Regression Analysis

What predicts patterns of cultural diversity?



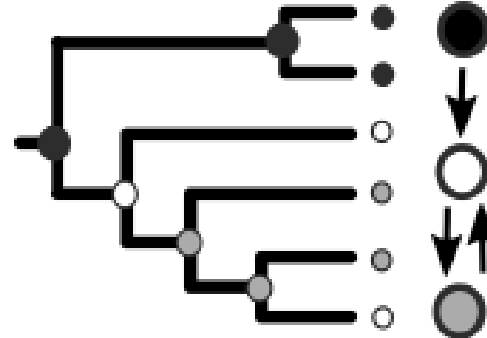
### c. Ancestral States

What was the earlier form of a feature?



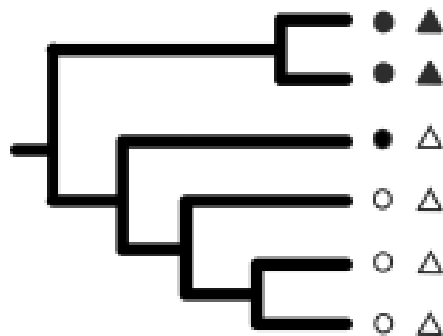
### d. Transformation

How do cultural features change form?



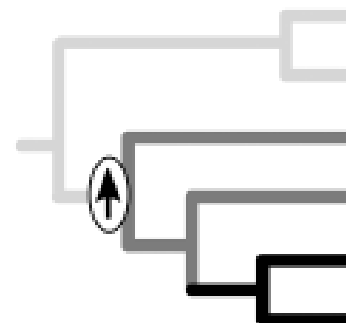
### e. Correlated Evolution

Do features change together?

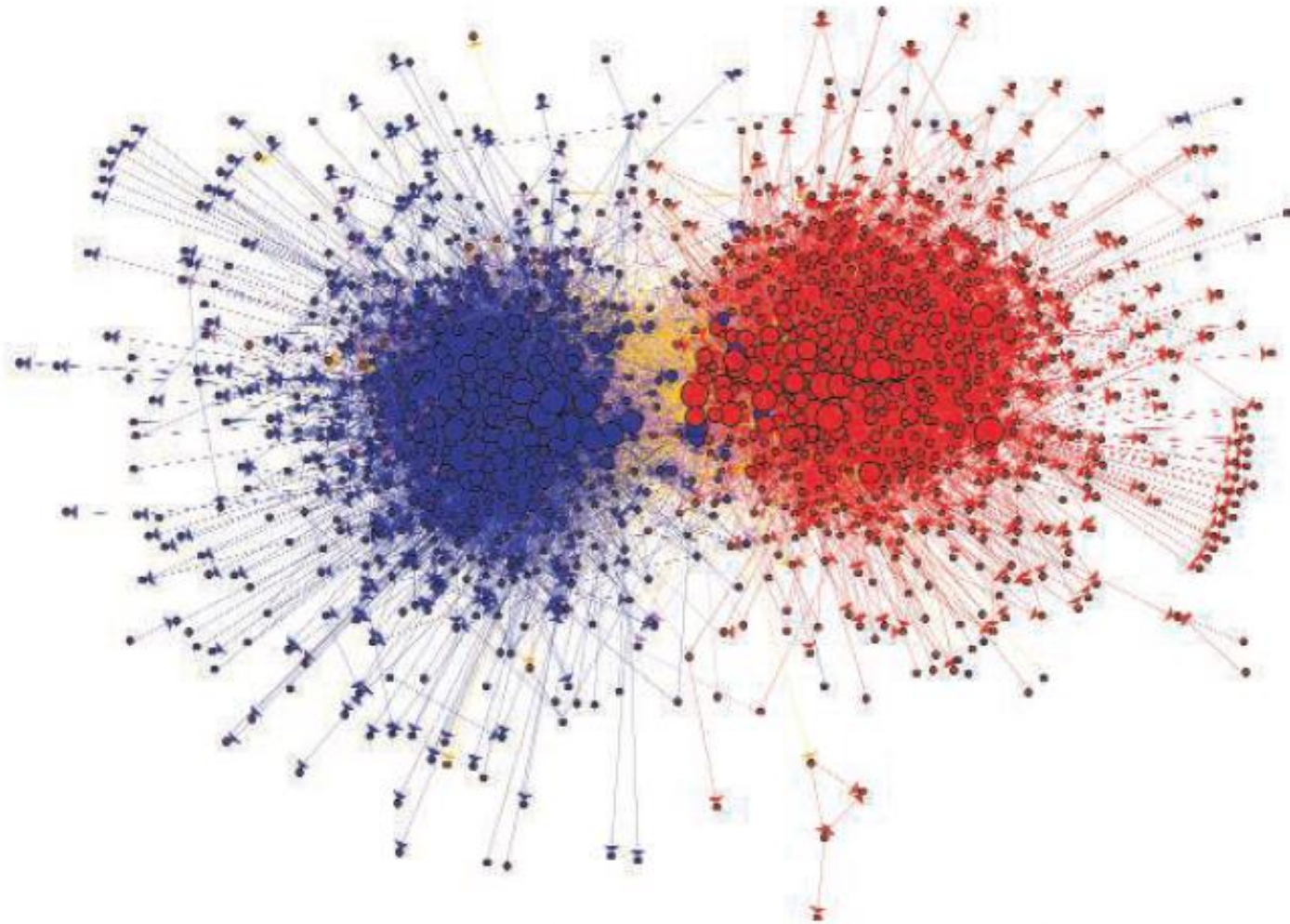


### f. Mode and Tempo

How and when do features diversify?

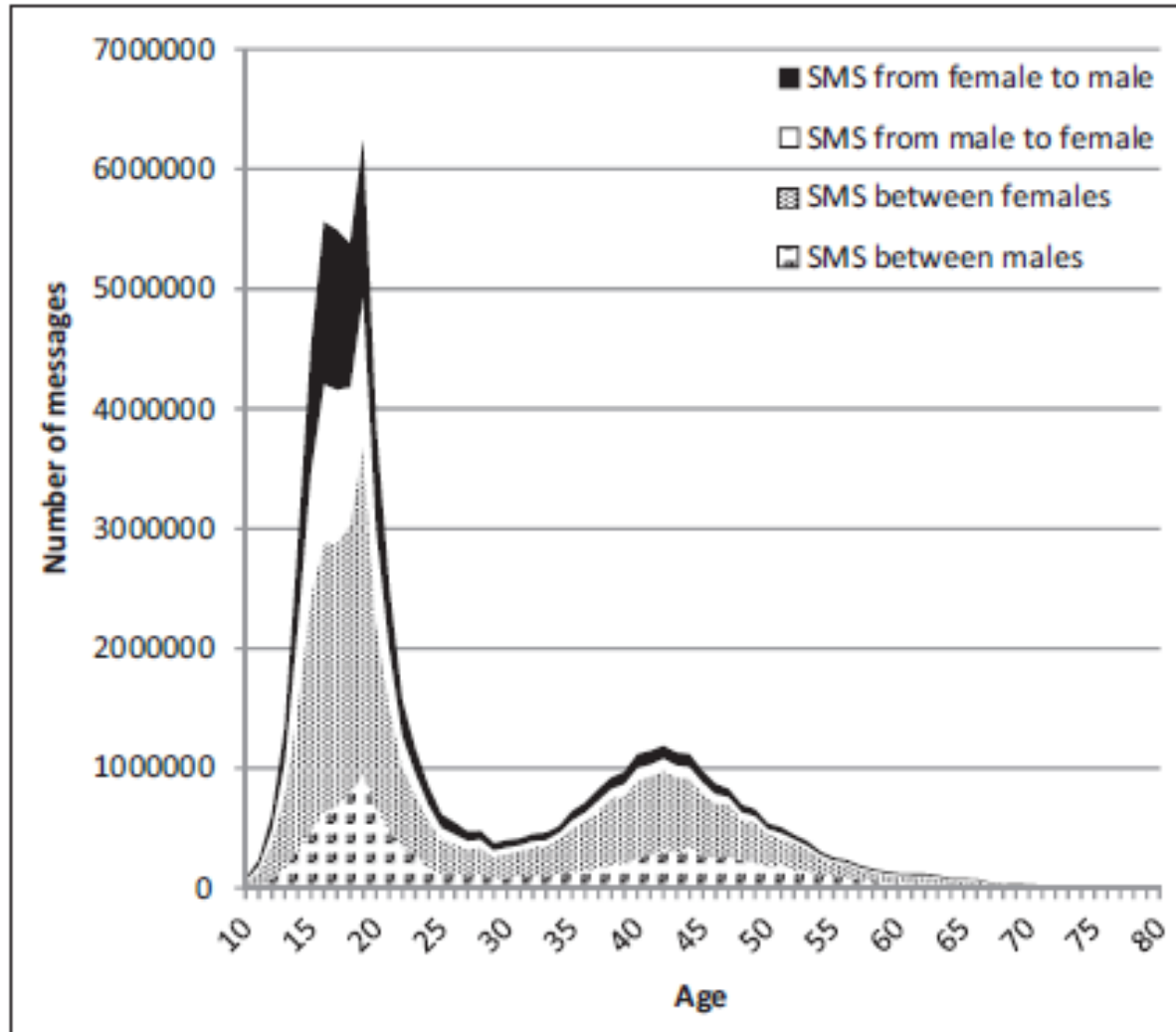


# Political Blogs



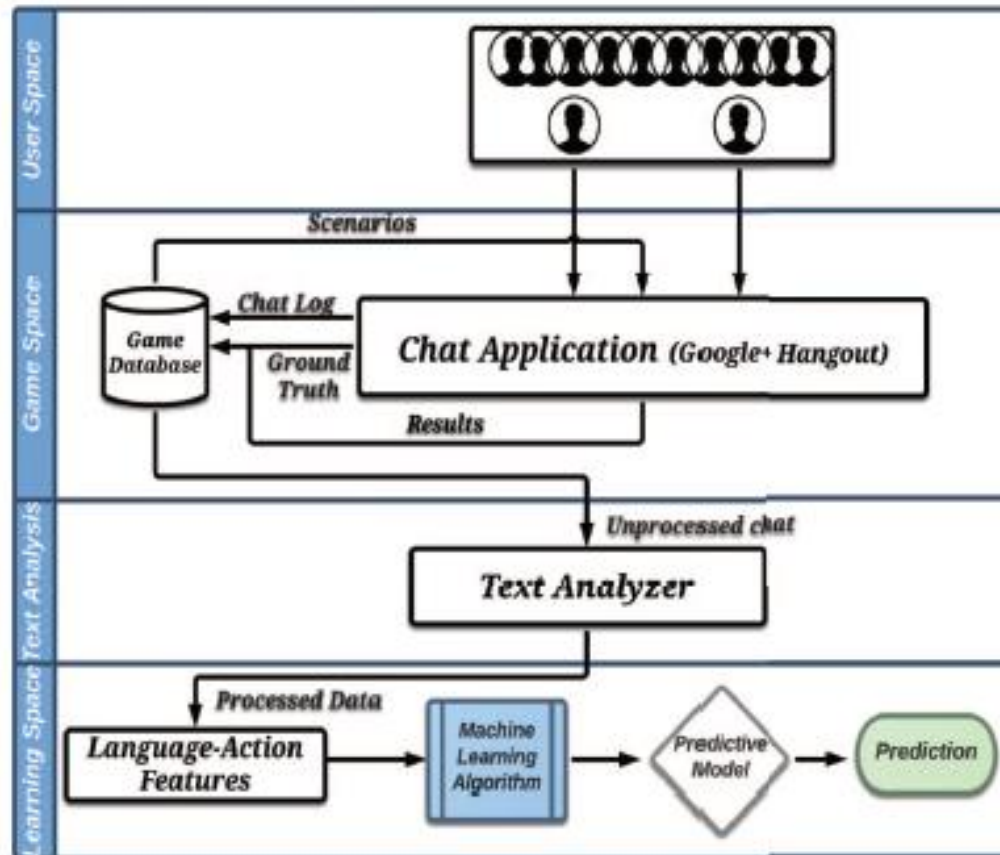
Data from the blogosphere. Shown is a link structure within a community of political blogs (from 2004), where red nodes indicate conservative blogs, and blue liberal. Orange links go from liberal to conservative, and purple ones from conservative to liberal. The size of each blog reflects the number of other blogs that link to it. [Reproduced from (8) with permission from the Association for Computing Machinery]

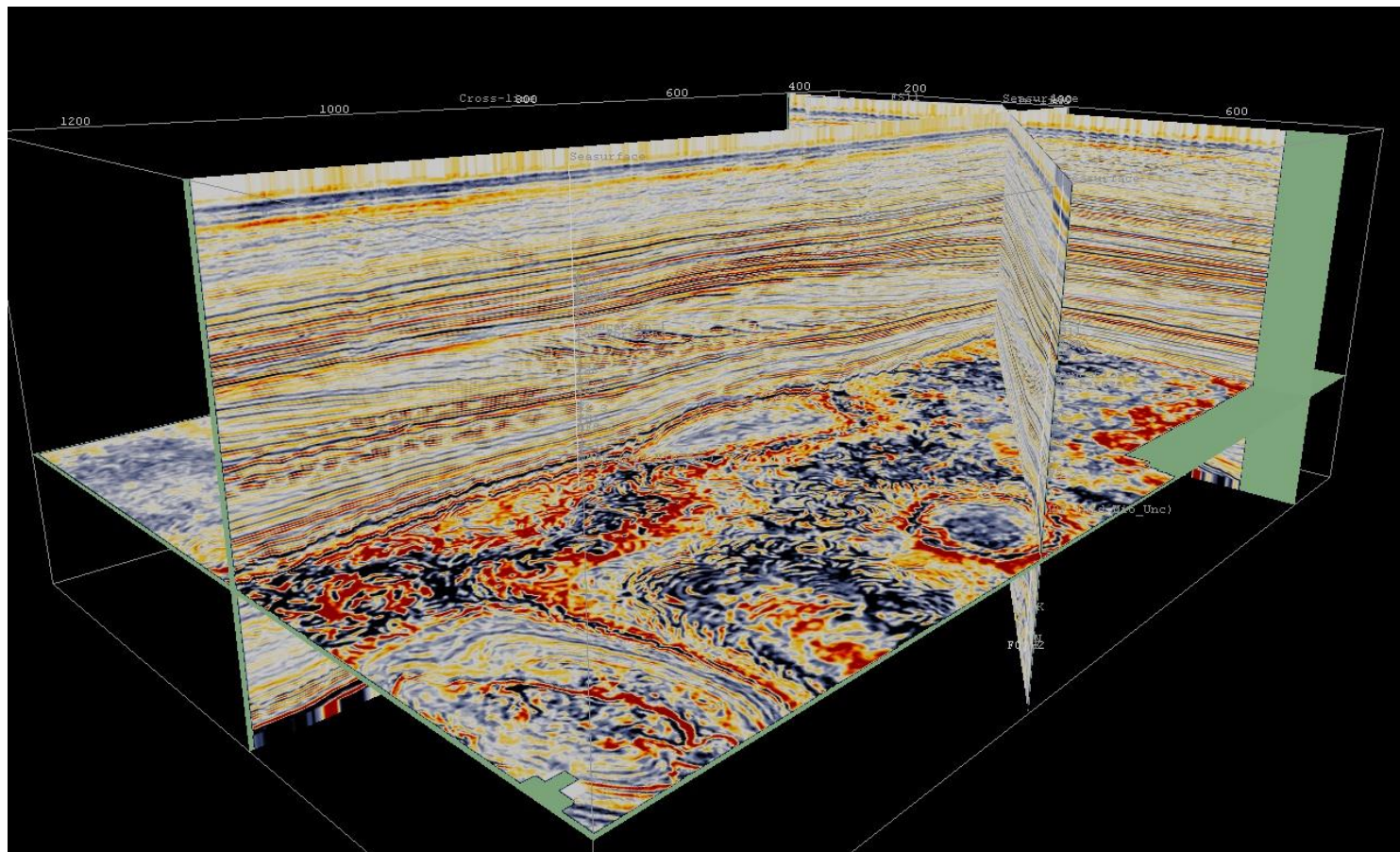
# Texting



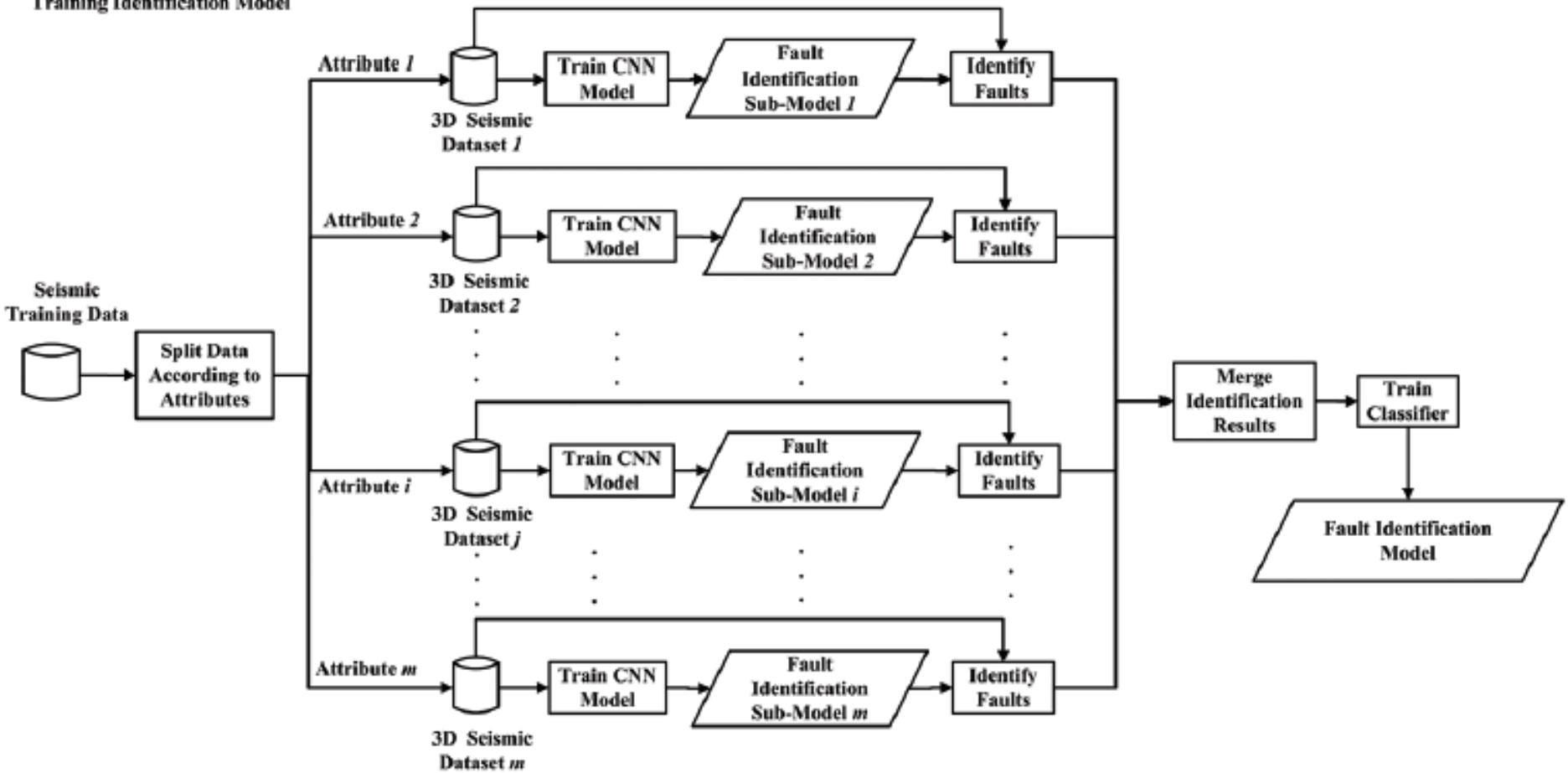


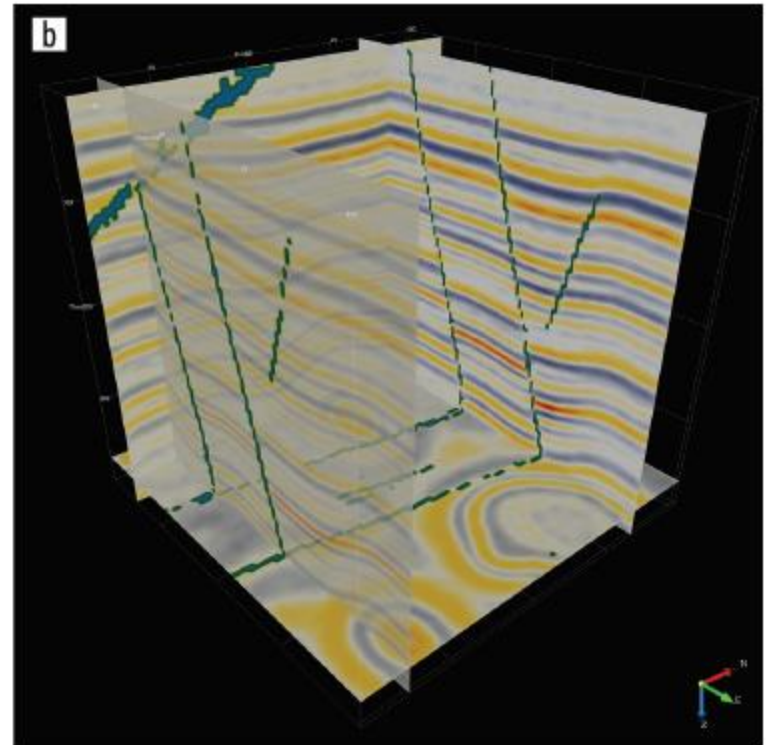
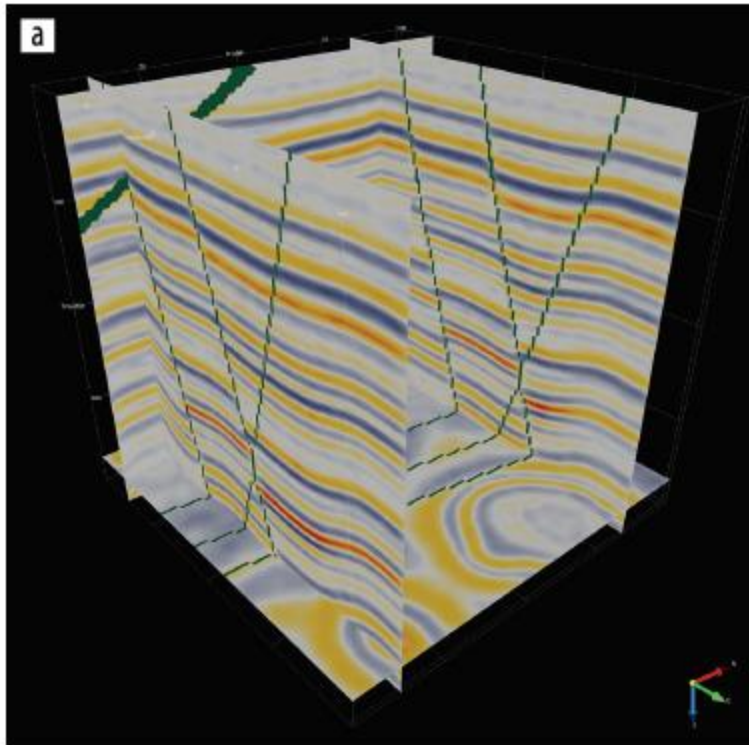
# Automated Detection of Deceptive Language-Action Cues





### Training Identification Model





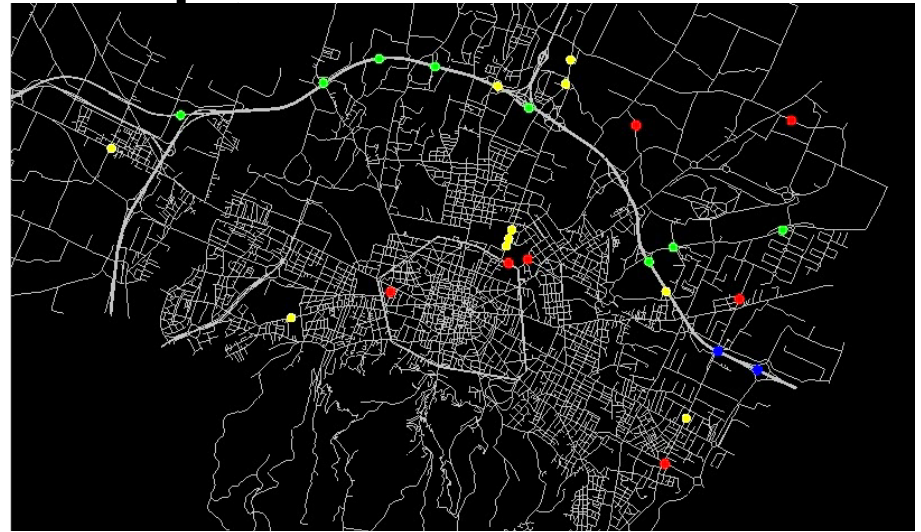


# URBAN TRAFFIC

## Heterogeneity

- Spatial and Temporal
- Congestion Level
- Topology
- Modes of transport
- Sensing equipment

## Sparse multi-sensor



## Develop travel time field and route choice information

