



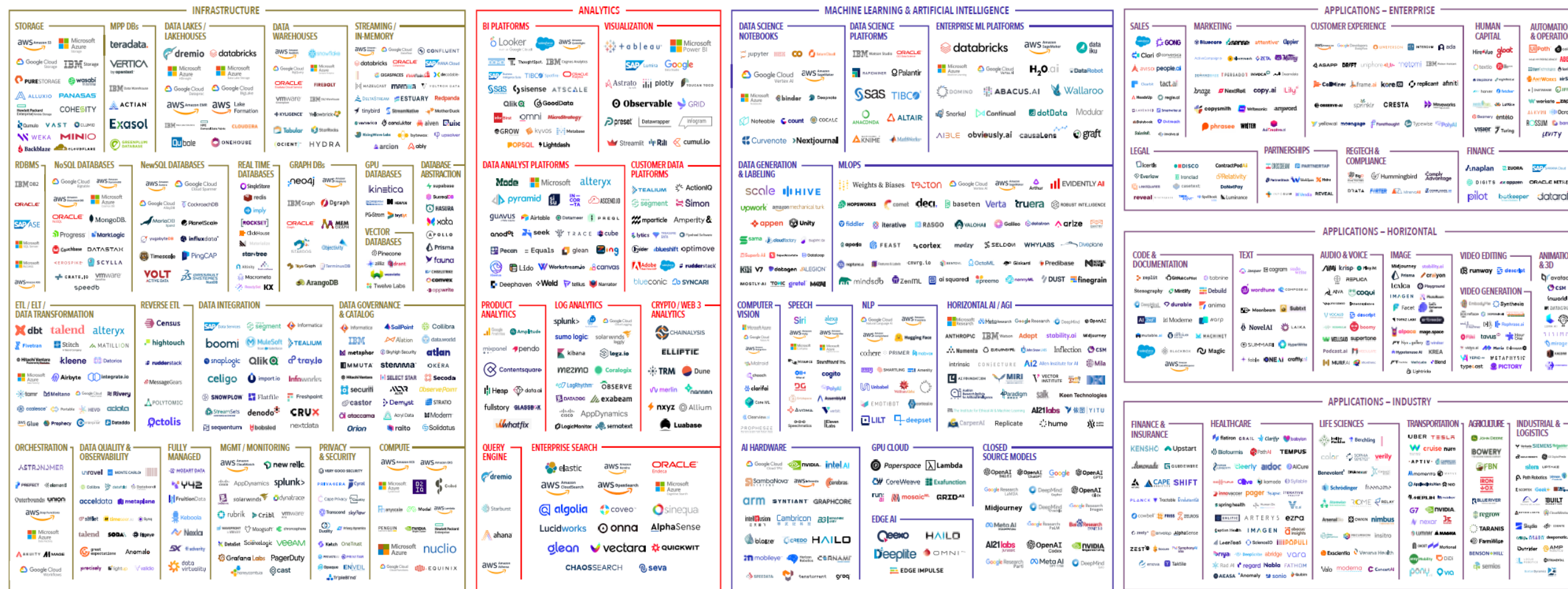
Longbing Cao | [www.datasciences.org](http://www.datasciences.org)

# Declaration

- Please be noted not every cited material has been properly acknowledged in this presentation
- The author acknowledges the wisdom and contributions of all relevant contributors
- More information about the content discussed in this presentation is in the book: Data Science Thinking



THE 2023 MAD (MACHINE LEARNING, ARTIFICIAL INTELLIGENCE & DATA) LANDSCAPE



## — OPEN SOURCE INFRASTRUCTURE

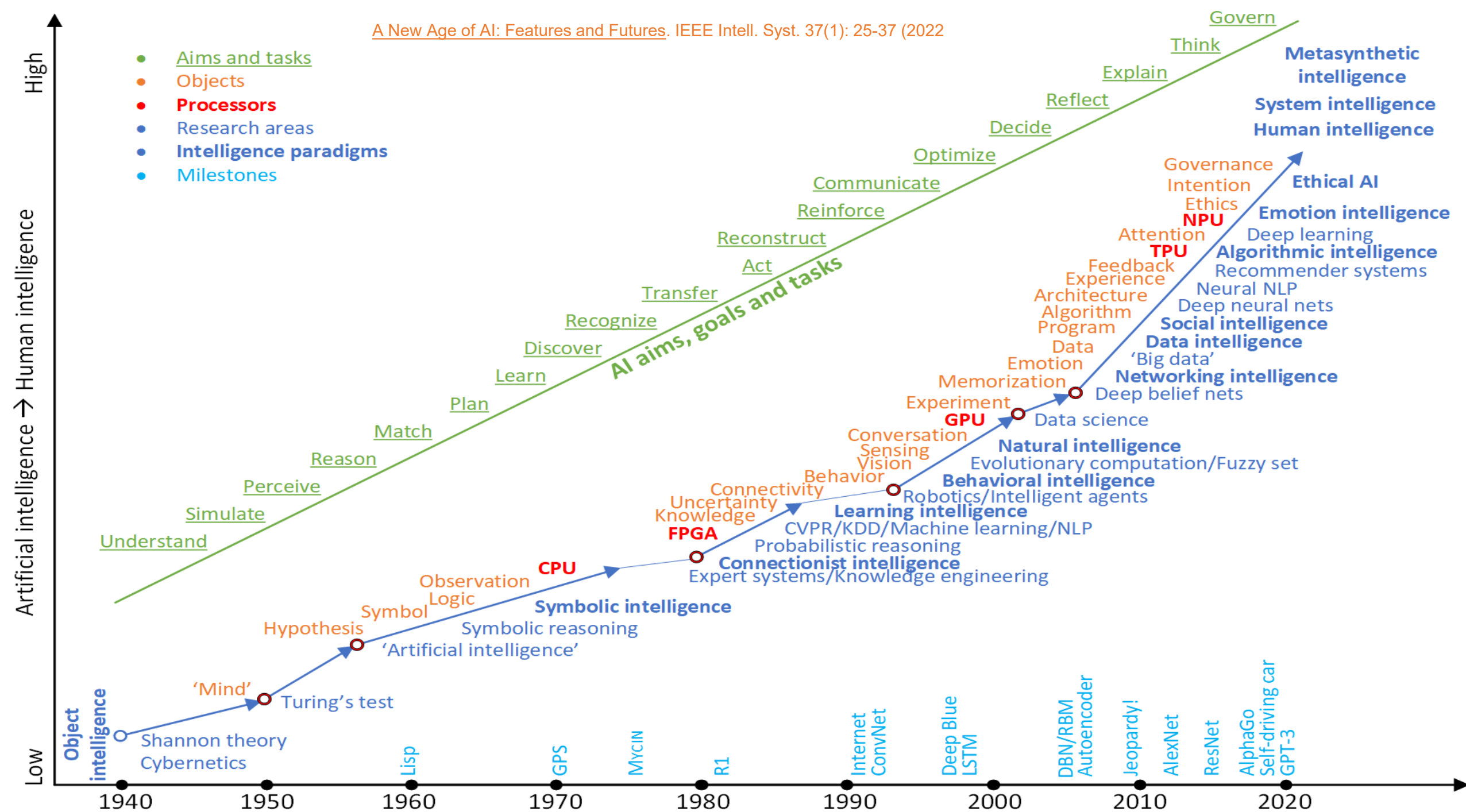


## - DATA SOURCES & APIs



## — DATA & AI CONSULTING

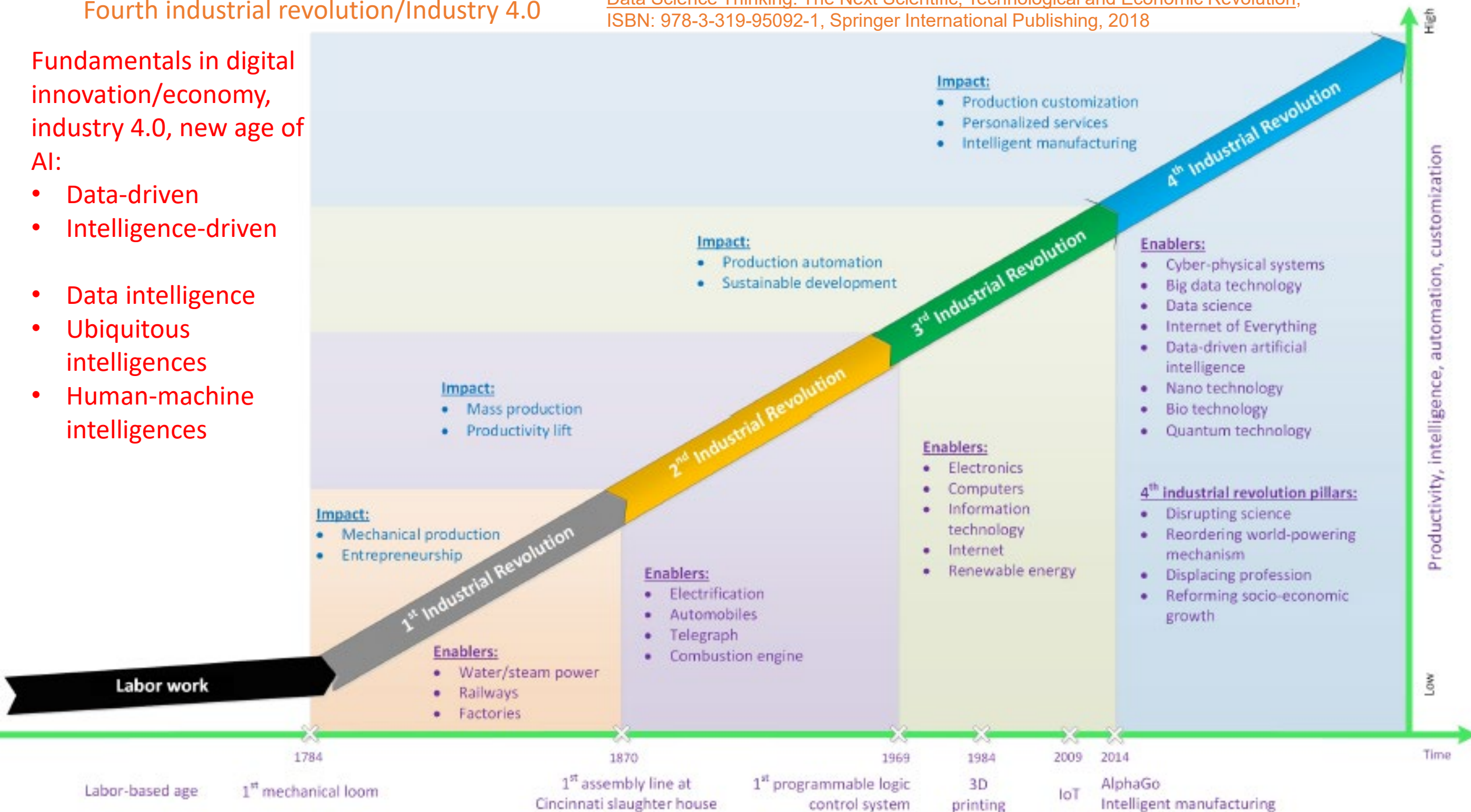




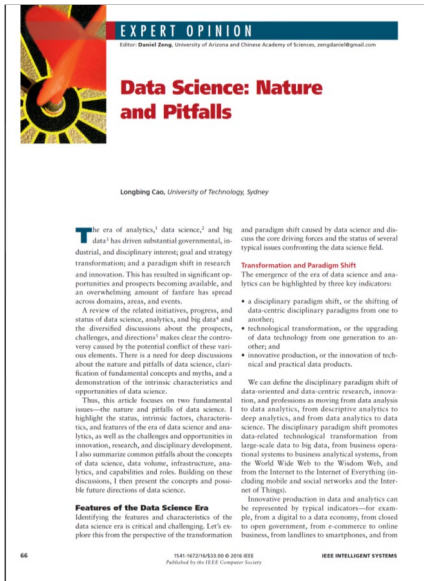


Fundamentals in digital innovation/economy, industry 4.0, new age of AI:

- Data-driven
- Intelligence-driven
- Data intelligence
- Ubiquitous intelligences
- Human-machine intelligences



# Data Science Concepts



L. Cao. IEEE Intelligent Systems, 2016



L. Cao. Communications of the ACM, 2017



L. Cao. ACM Computing Survey, 2017



L. Cao. IEEE Intelligent Systems, 2019



# 50 Years of data science: an immature discipline

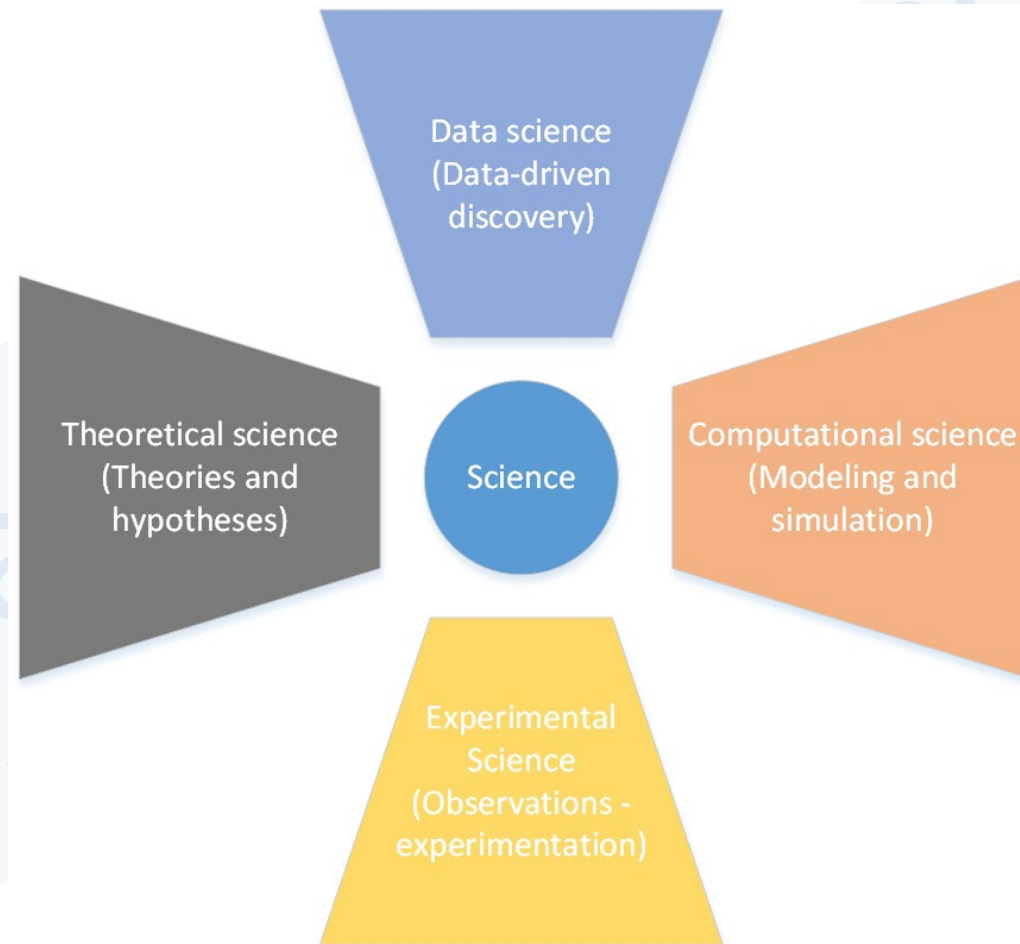
D. Donoho, “50 Years of Data Science,” 2015;  
<http://courses.csail.mit.edu/18.337/2015/docs/50YearsDataScience.pdf>



<https://link.springer.com/content/pdf/10.1007%2F978-3-319-95092-1.pdf>



# Scientific paradigms – the era of data



Big data	Refers to data that are too large and/or complex to be effectively and/or efficiently handled by traditional data-related theories, technologies and tools.
----------	---

# What is data science?

**Definition 2.2 (Data Science<sup>1</sup>).** Data science is the science of data, or data science is the study of data.

**Definition 2.3 (Data Science<sup>2</sup>).** Data science is a new trans-disciplinary field that builds on and synthesizes a number of relevant disciplines and bodies of knowledge, such as statistics, informatics, computing, communication, management and sociology, to study data and its domain employing data science thinking.

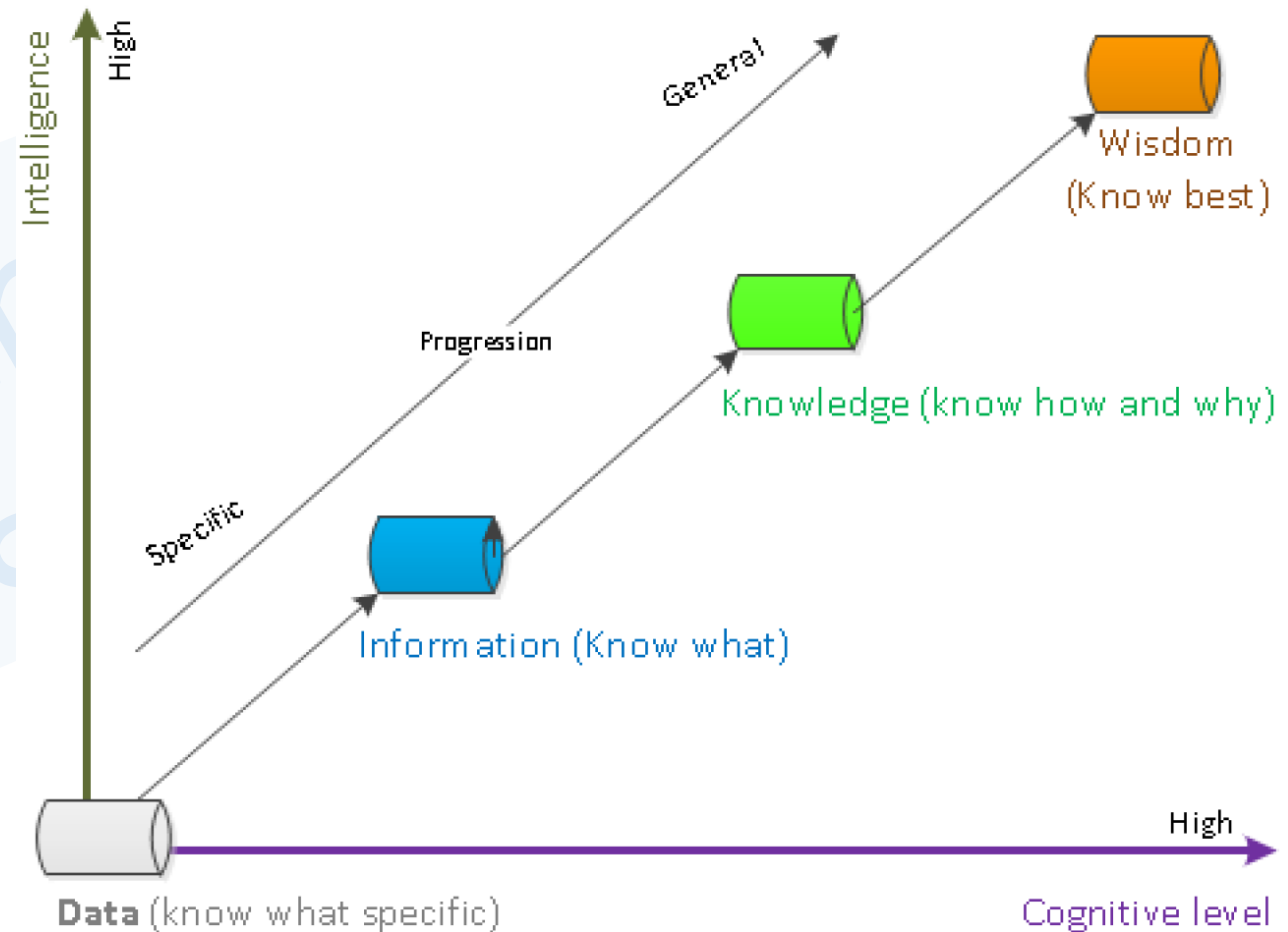
$$\text{data science} \stackrel{\text{def}}{=} \{ \text{statistics} \cap \text{informatics} \cap \text{computing} \cap \text{communication} \\ \cap \text{sociology} \cap \text{management} \mid \text{data} \cap \text{domain} \cap \text{thinking} \} (2.1)$$

where “|” means “conditional on.”

# Data science: converting data to intelligence/wisdom

- Data intelligence
- Actionable intelligence
  - unique and valuable understanding, thinking, insights and expertise that can enable significantly better and smarter planning, decision-making and outcomes

Data science thinking, Springer, 2018

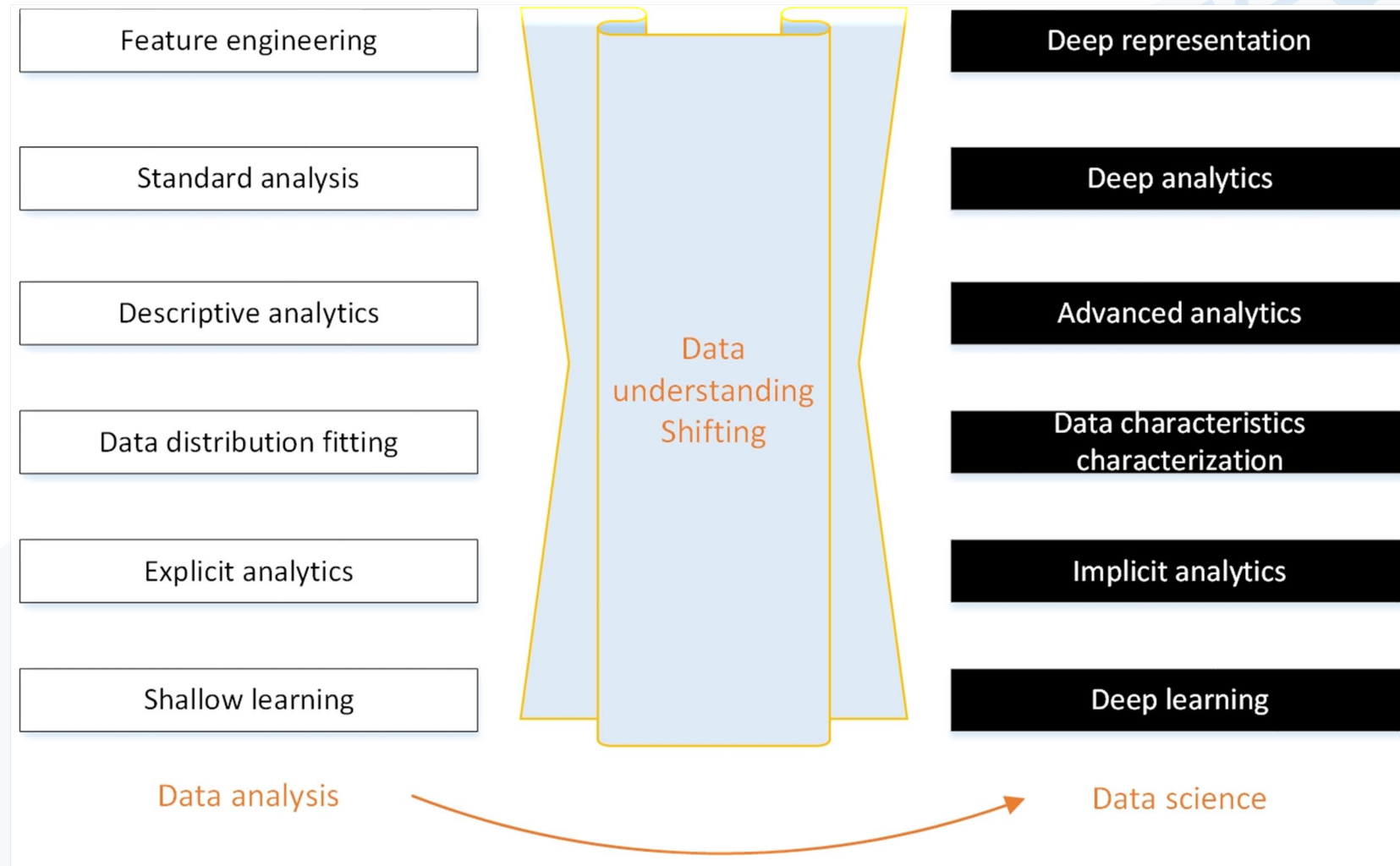


**Fig. 2.1** Data-to-information-to-knowledge-to-intelligence-to-wisdom cognitive progression.

<sup>2</sup>Note: X-axis: the increase in cognitive level; Y-axis: the increase in intelligence



# Paradigm shift: Well-developed data analysis → Immature data science

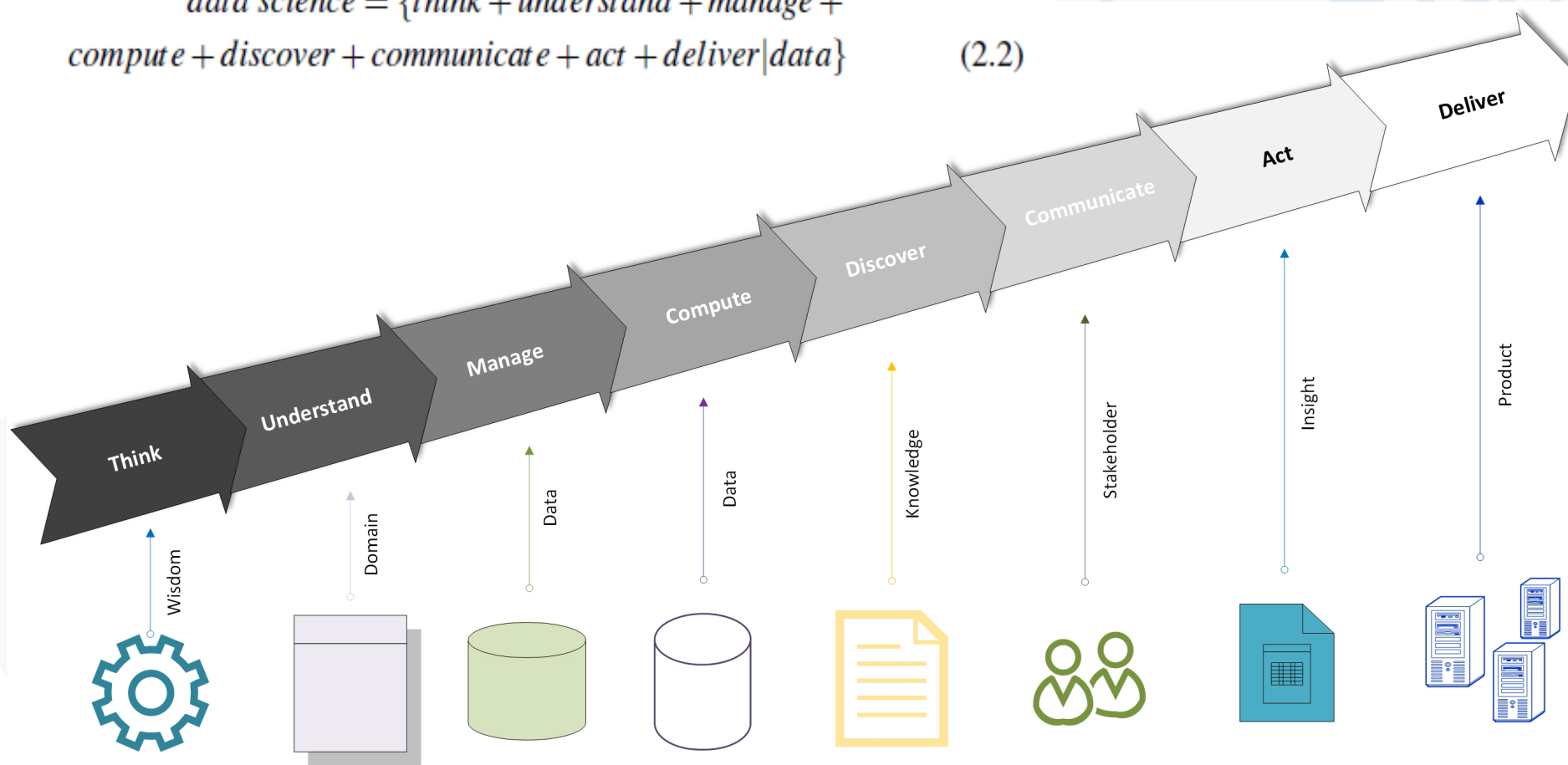


L. Cao. Data science thinking, Springer, 2018

# Data science processes

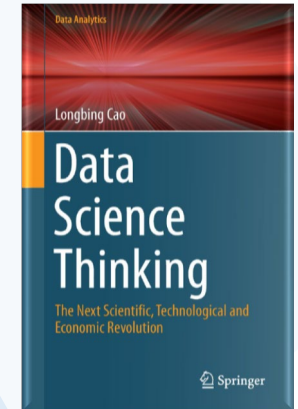
**Definition 2.4 (Data Science<sup>3</sup>).** From the *process* perspective, *data science* is a systematic approach to “thinking with wisdom”, “understanding the domain”, “managing data”, “computing with data”, “discovering knowledge”, “communicating with stakeholders”, “acting on insights”, and “delivering products”.

$$\text{data science} \stackrel{\text{def}}{=} \{ \text{think} + \text{understand} + \text{manage} + \text{compute} + \text{discover} + \text{communicate} + \text{act} + \text{deliver} | \text{data} \} \quad (2.2)$$

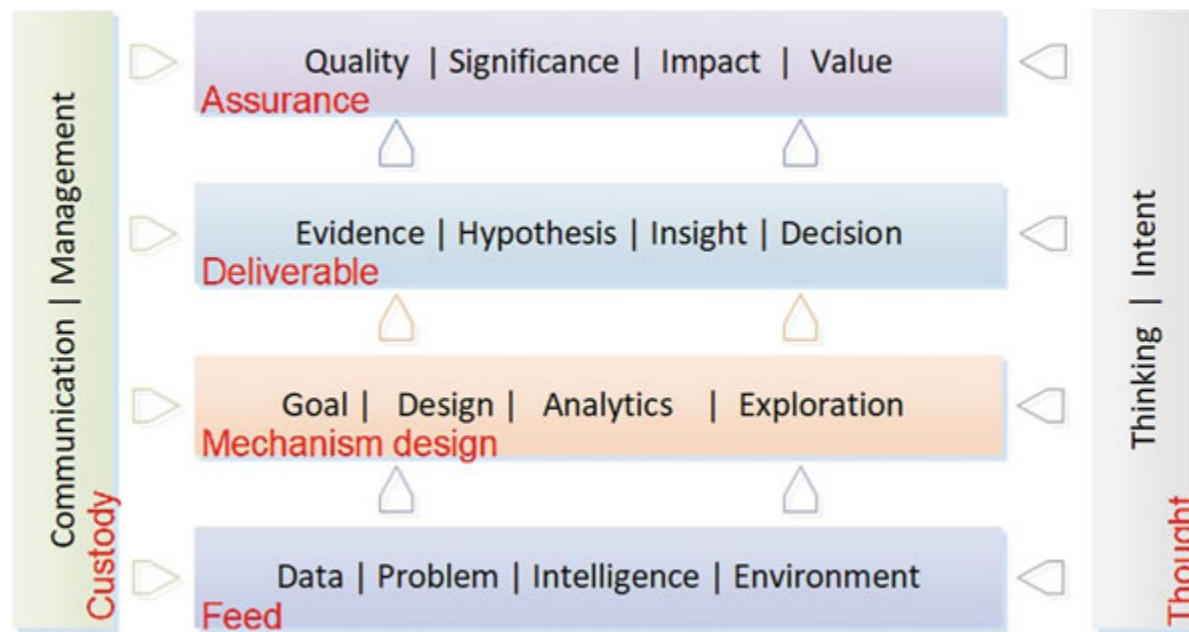


- Not necessary 80:20 rule
- End-to-end solution
- Thinking, design, and actionability
- Embed data processing and feature engineering in learning

# Data science/analytical thinking



**Definition 3.1 (Data Science Thinking)** Data science thinking refers to the perspective on the methodologies, process, structure, and traits and habits of the mind in handling data problems and systems.

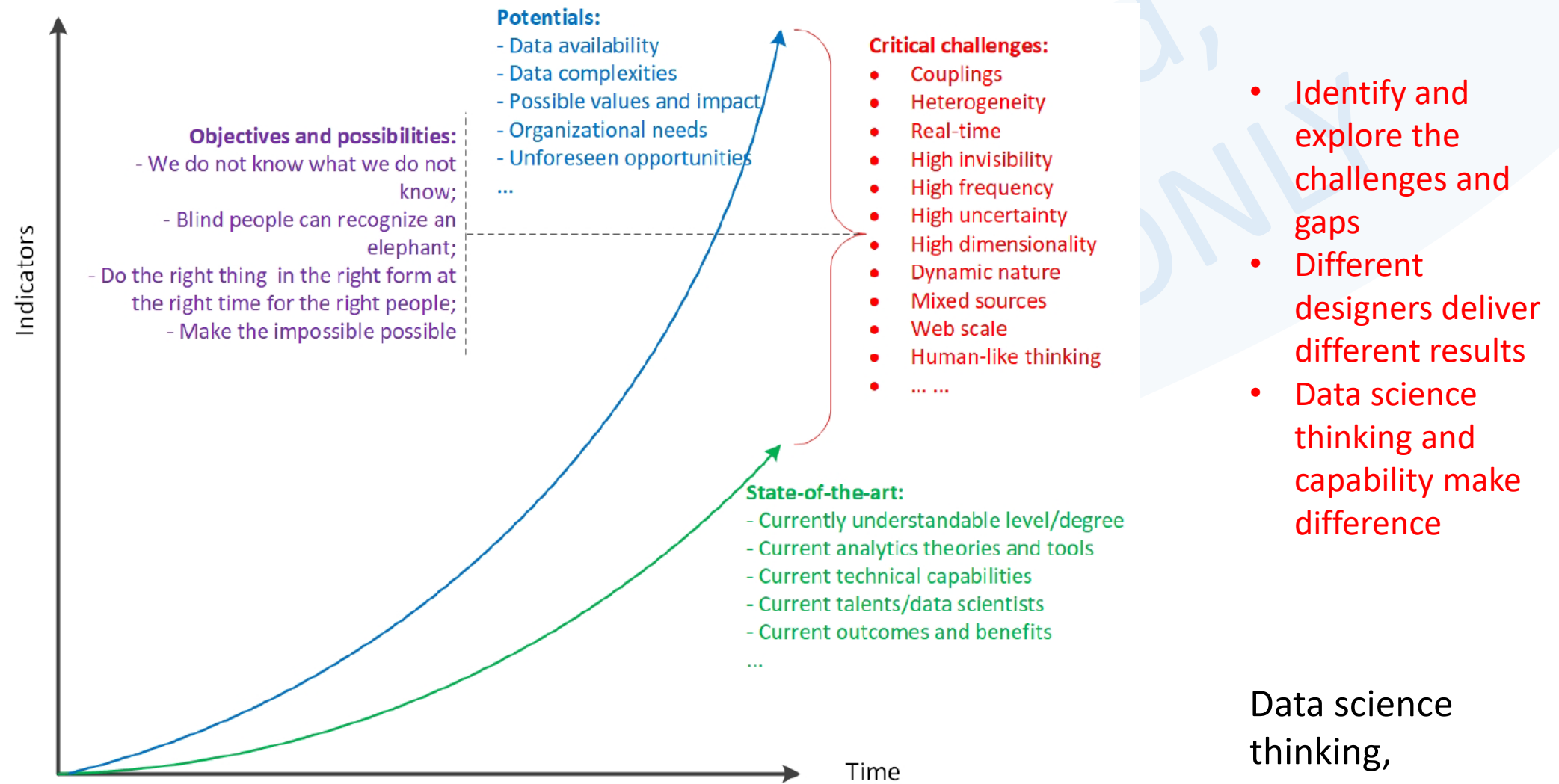


**Fig. 3.5** Data science overview and structure

- Scientific thinking
- Data-oriented critical and creative thinking
- Individualized → holistic → systematic design
- From methodologies, processes, structures, designs to evaluation, deliverables
- Design thinking in analytics and learning paradigms and methods
- The best about what, how and why – insights and strategies



# Gaps: complexities vs. capabilities/capacity



Data science  
thinking,  
Springer, 2018

**Fig. 5.1** Growing disciplinary gaps between data potential and disciplinary capabilities.

# X complexities

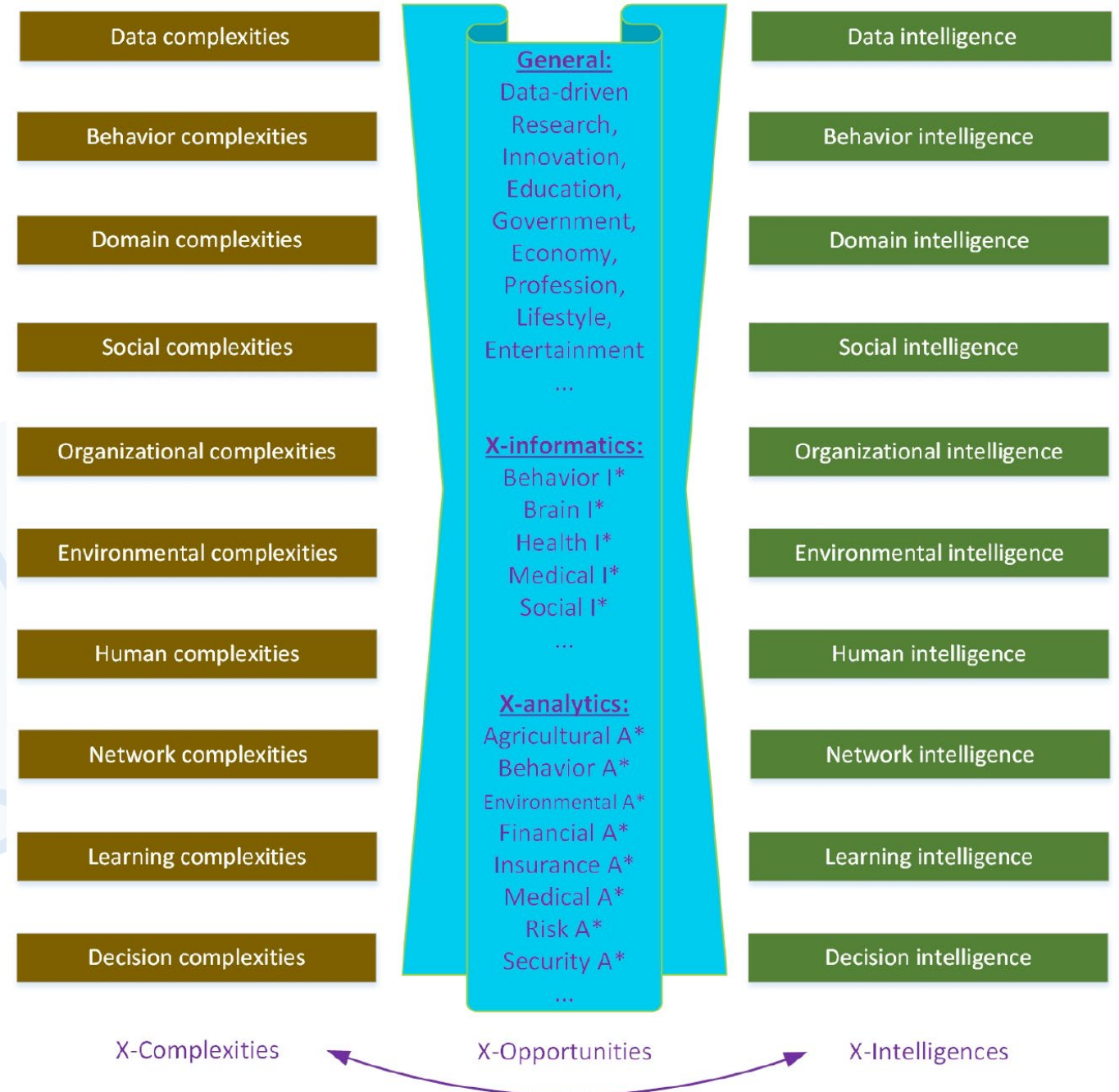
# X intelligences

# X analytics

- E.g., deep learning on small, noisy, inconsistent, evolving case data
- Contextual factors: ethnic, social, cultural, persona, ...
- Epidemiological knowledge on its transmission, ...

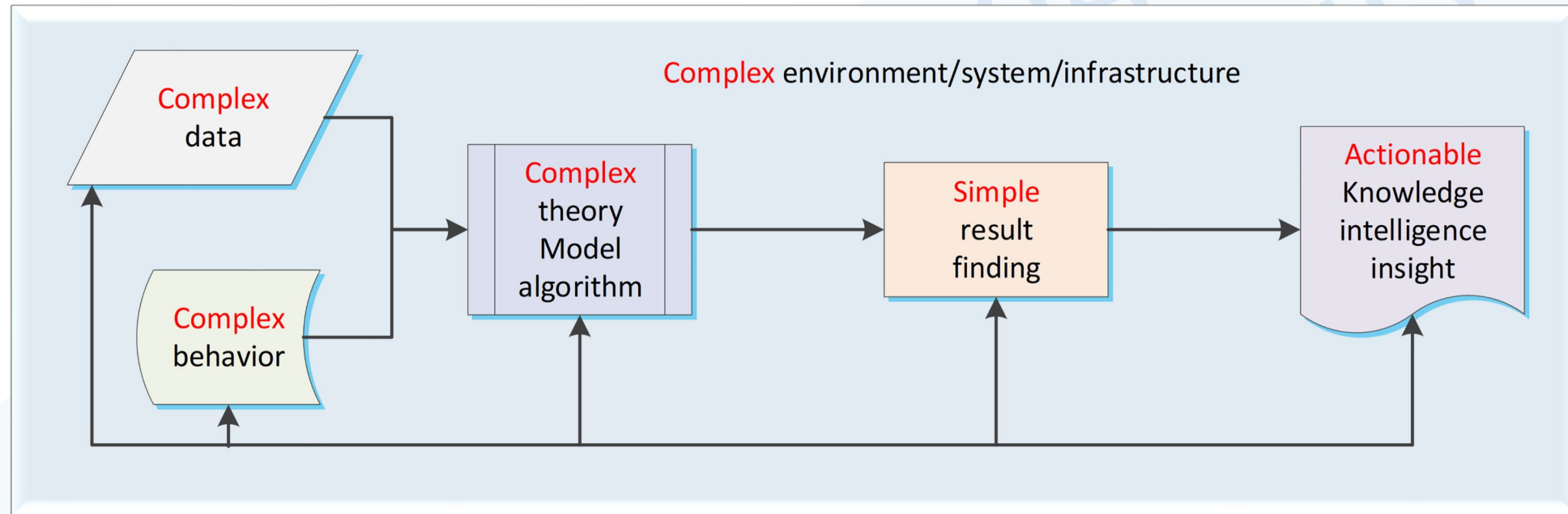
Data science: Challenges and directions, Communications of the ACM, 2017

Data science thinking, Springer, 2018



**Fig. 1.3** The new X-generations: X-complexities, X-intelligence, and X-opportunities.

# Complex real world vs. often simple, specific solutions and results

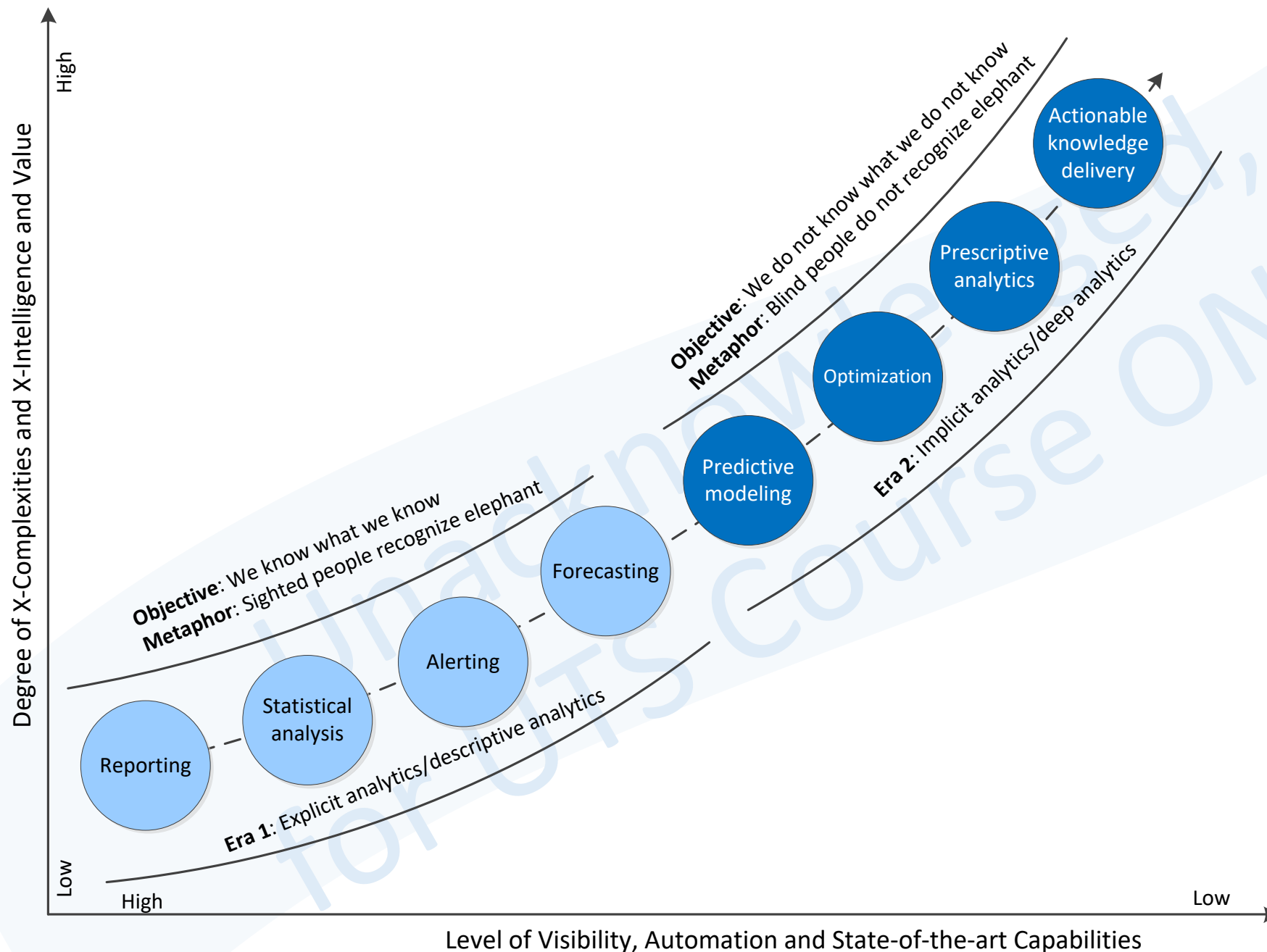


- Enterprise data science is a complex process and system
- Understanding and quantifying complexities is a major task for new strategic and value proposition



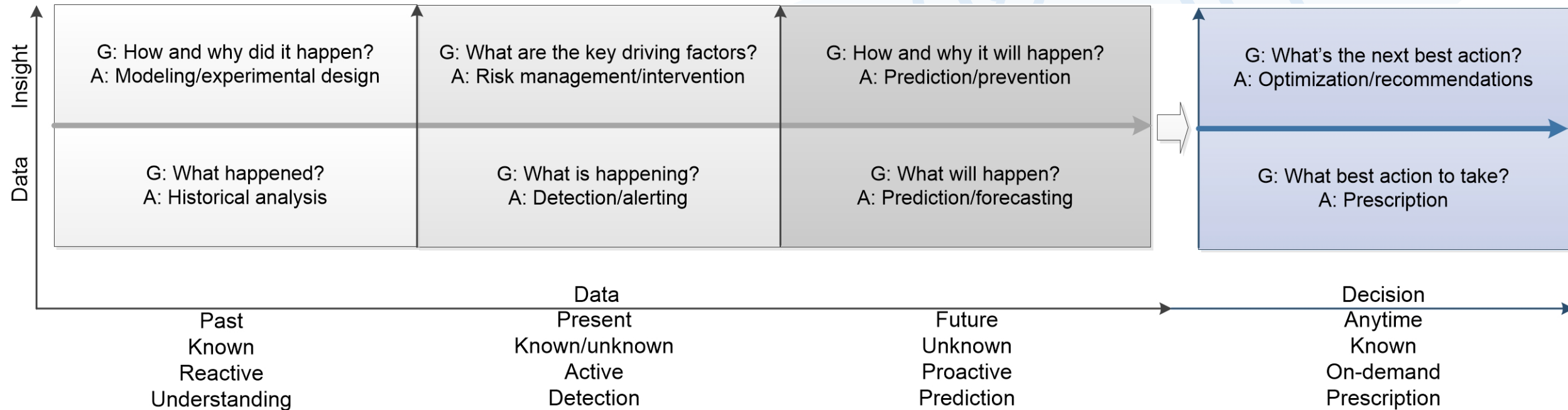
# Enterprise Data Science

# Analytics spectrum



- Present enterprise data science focus on explicit and descriptive analytics
- Enterprise data science shift to implicit and deep enterprise-wide analytics

# Data → Insight/Decision



- What → why → how
- Lifelong: past → present → future
- Known → unknown → known →
- Reactive → active → proactive
- Understanding → detection → prediction → intervention

Data science  
thinking,  
Springer, 2018

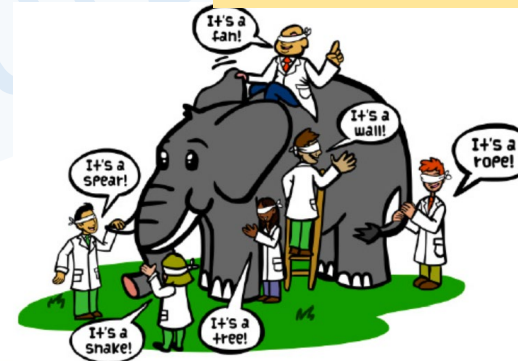
# Big data analytics era

## Major effort:

Shallow learning  
Descriptive analytics  
Explicit analytics  
Off-the-shelf analytics  
...

## Minor effort:

Deep analytics  
Advanced analytics  
Implicit analytics  
Specialised analytics  
...



Analytics Spectrum

Sighted people  
recognize elephant

Blind people  
recognize elephant



We do not know what we do not know: challenges, solutions, gaps, opportunities



Fig. 3.7 Data science: The unknown world.

Data science: Challenges and directions,  
Communications of the ACM, 2017  
Data science thinking, Springer, 2018

Explore the unknowns

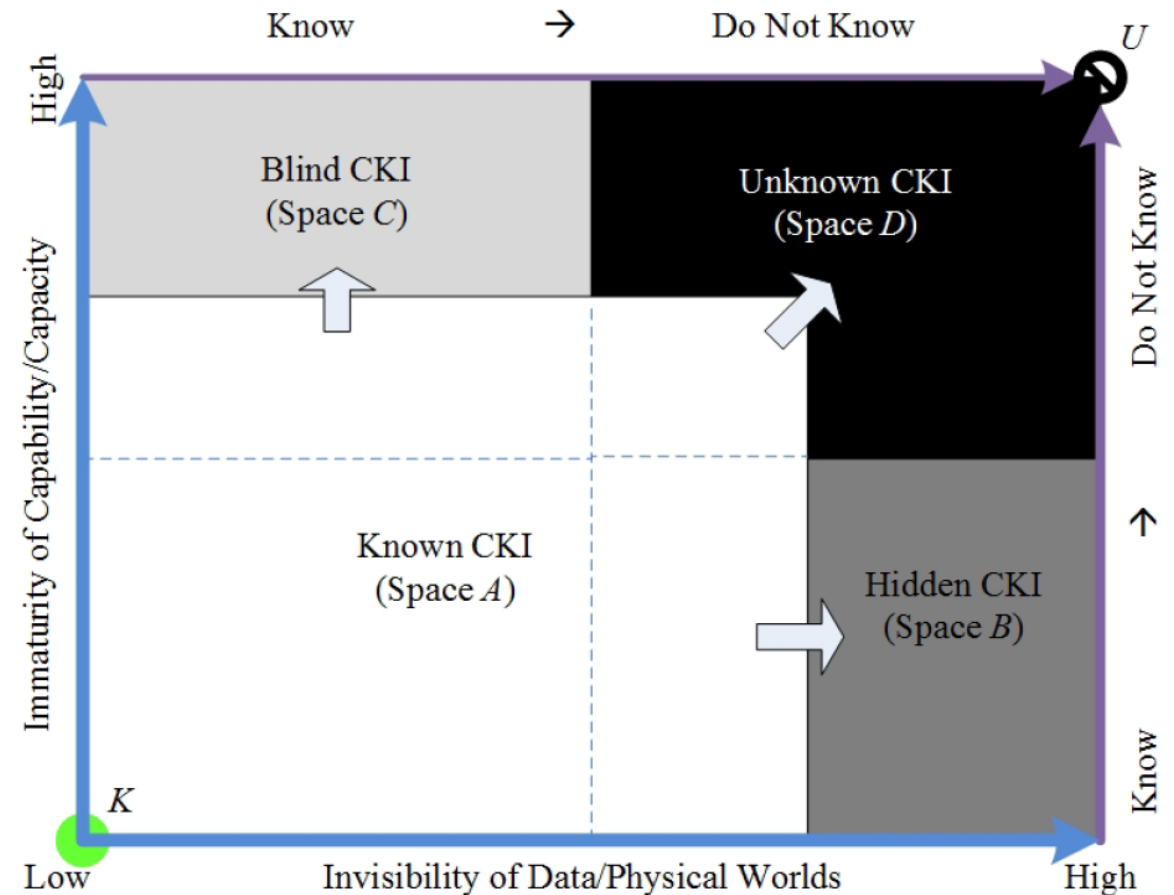
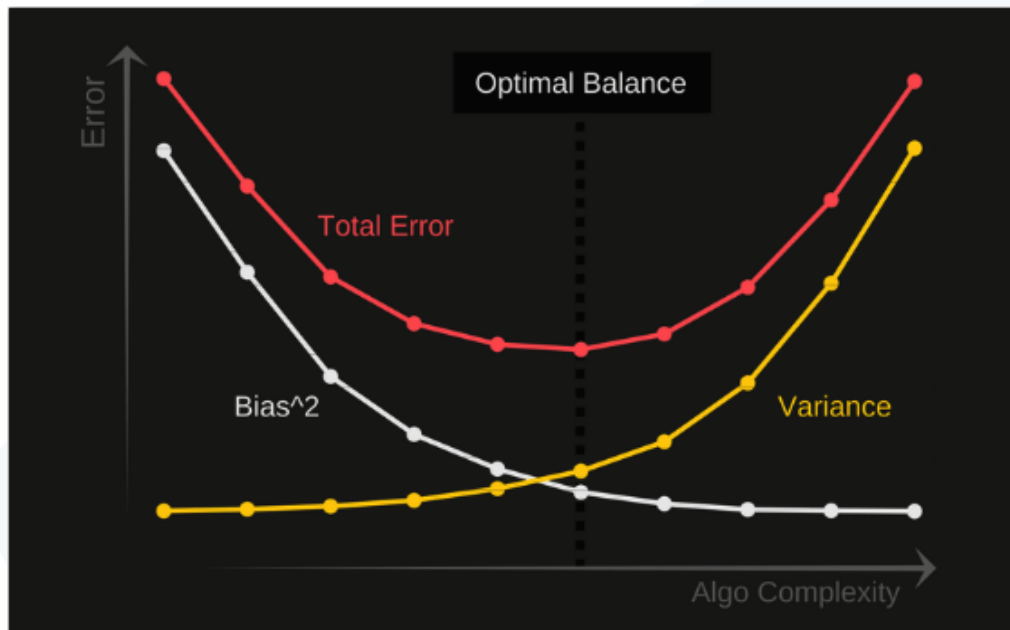


Fig. 3.6 Data science: Known-to-unknown research evolution.

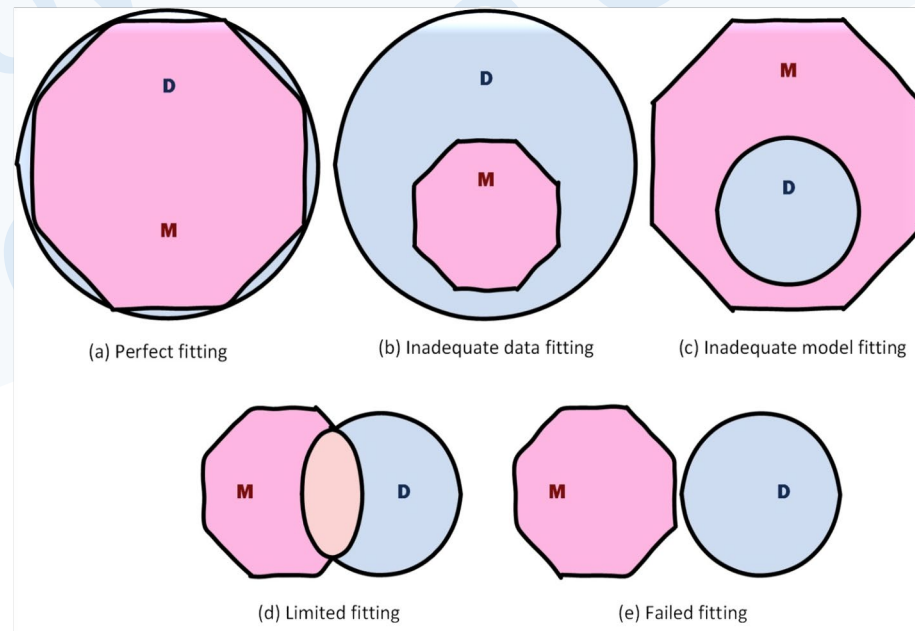
# Fitting gaps between model capacity vs. data potential

- Beyond fitting: data, model (parameter)
- Appropriate processes: sampling, validation etc.
- Appropriate thinking, theory, design, fitness, etc.

- Proper processes: sampling, validation etc.
- Proper thinking, design, fitness, etc.



$$\text{Total Error} = \text{Bias}^2 + \text{Variance} + \text{Irreducible Error}$$

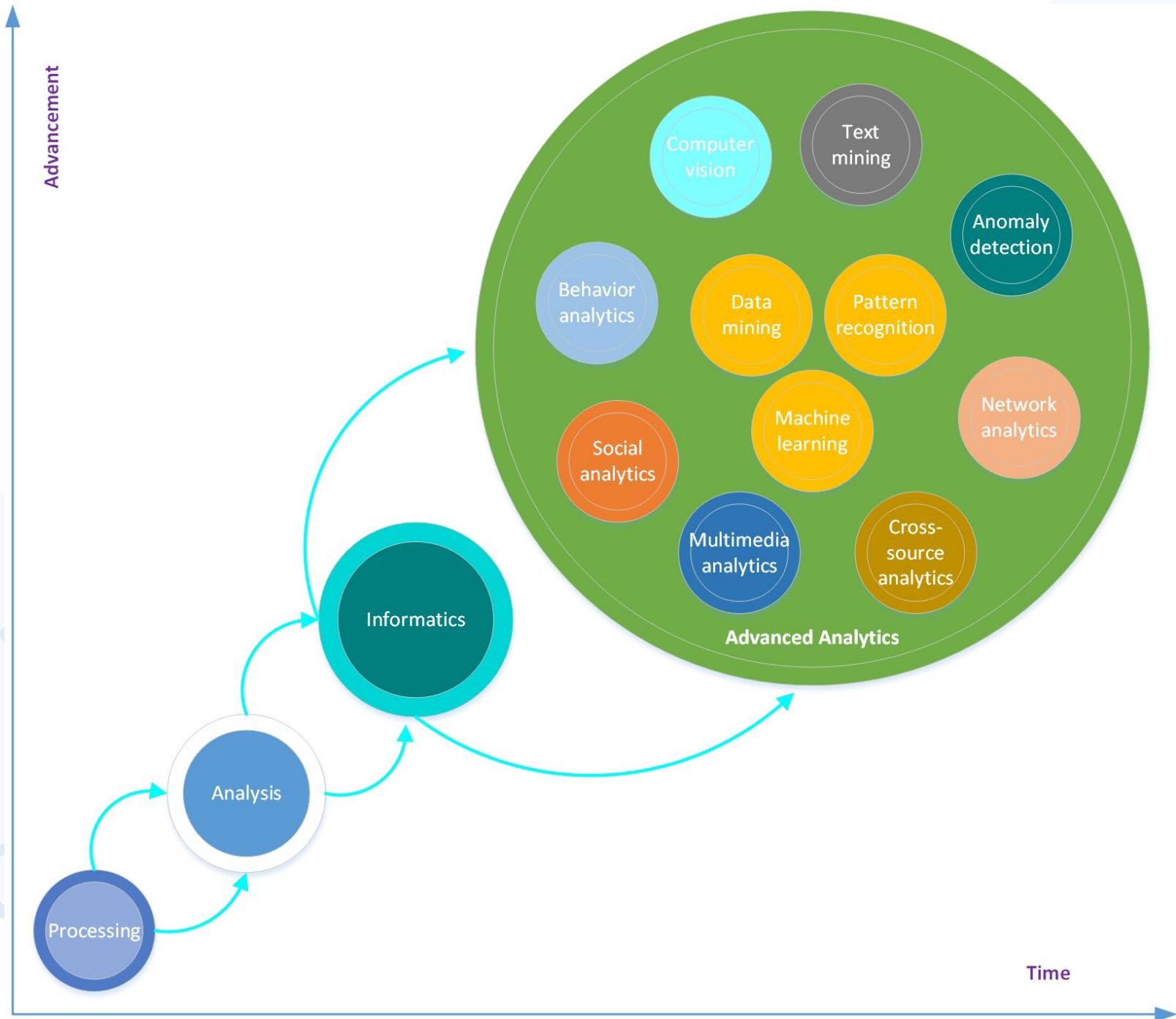


L. Cao, Data science thinking, Springer, 2018

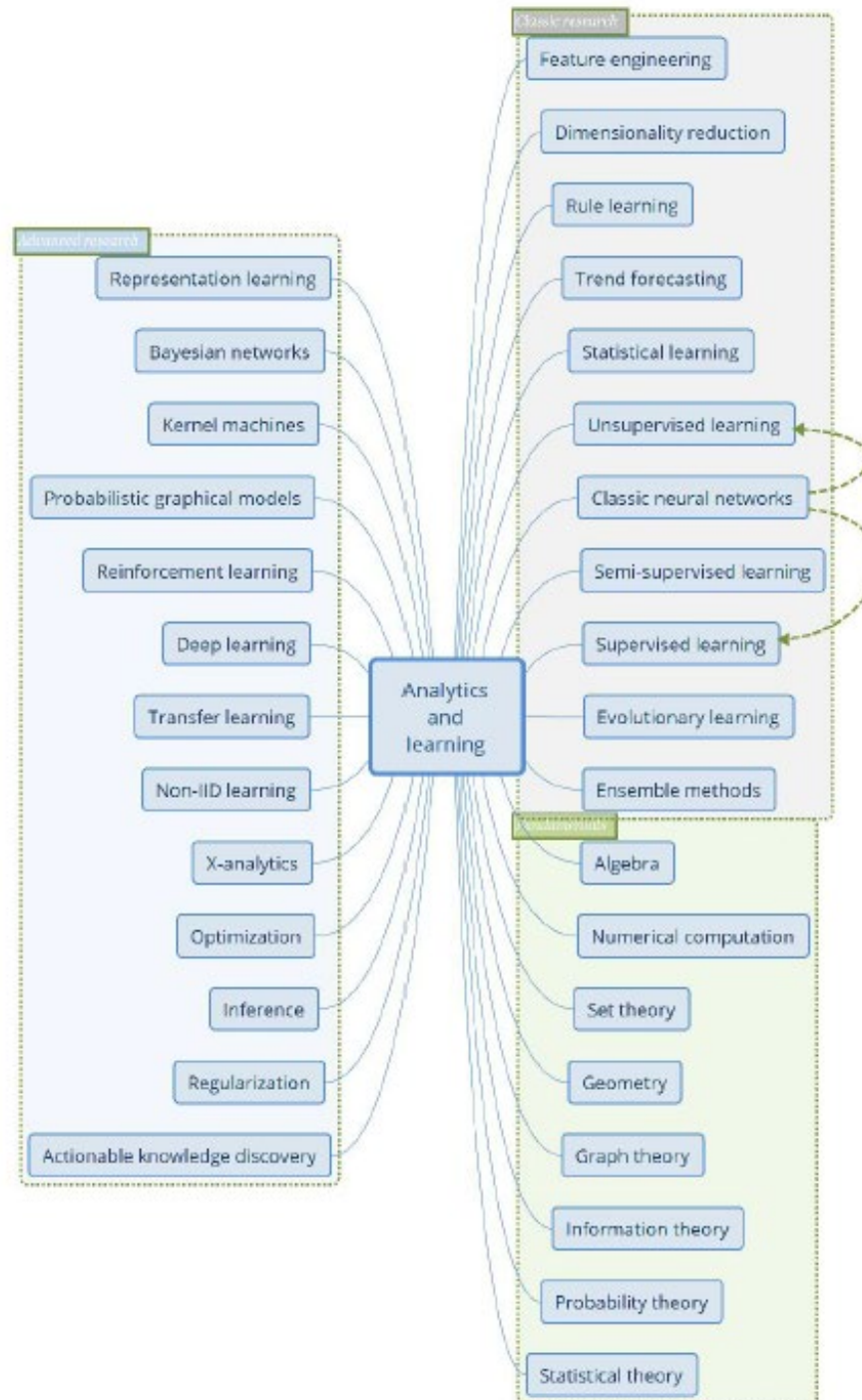
Data science thinking, Springer, 2018

# Advanced analytics

- Transactions
- Demographics
- Behaviors
- Communications
- Interactions
- Networks
- Intent
- Emotion



# Analytics and learning techniques

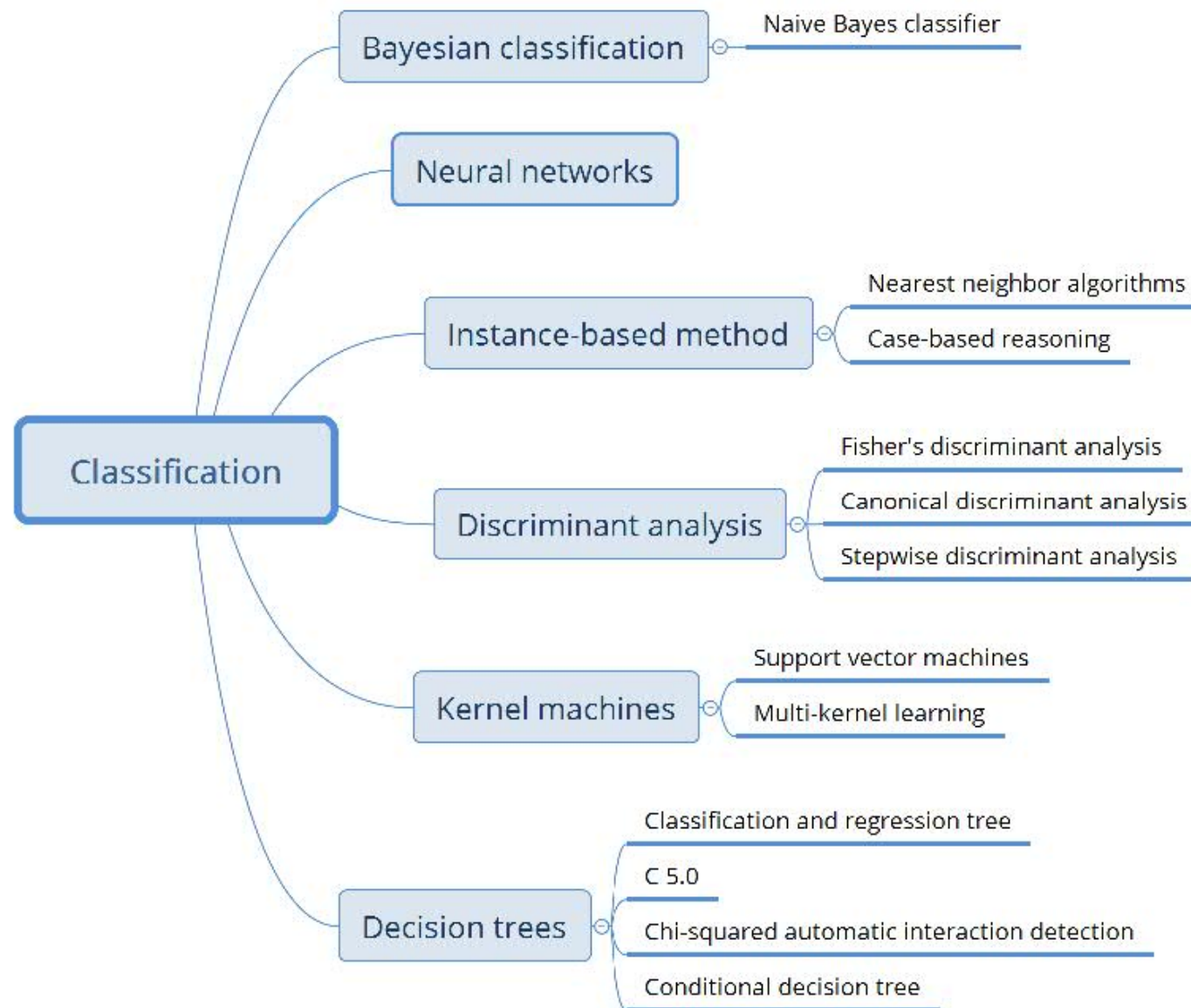


- Data science foundations
- Classic theories and methods
- Advanced theories and methods
- Data science thinking and design

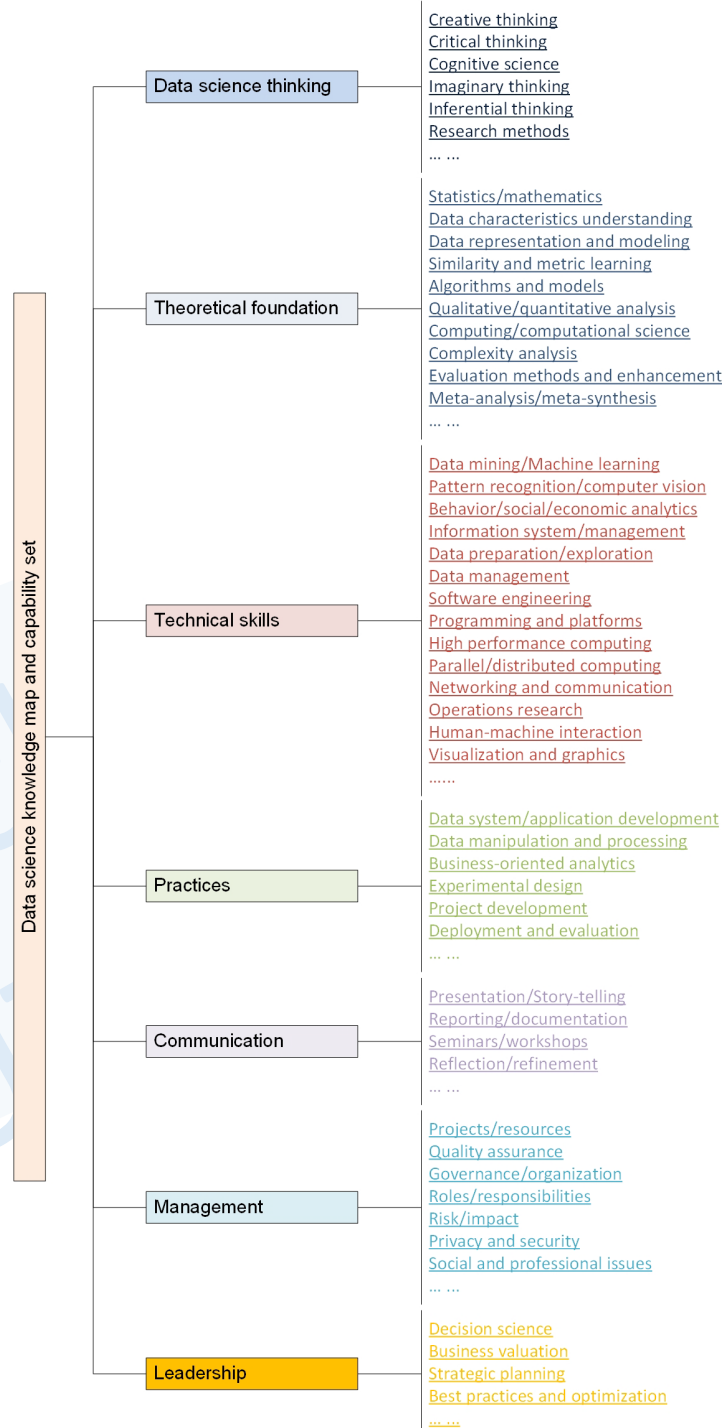


# Classifiers

- Popular shallow classifiers may not work well
- Pros and cons
- Fundamental assumptions for each method
- Different design methods
- Semi-supervised classification
- Unsupervised classification



# Knowledge map & capability set of a qualified data scientist



While it is impossible to achieve everything as we want,

it is hoped that the MBA Innovation Studio will empower you with some of enterprise data science thinking, knowledge, experience, and readiness

# Data scientists vs. data engineers

- Domain understanding
- Constraint understanding
- Data characteristics/complexities understanding
- Business-analytics problem transformation
- Data usage planning
- Analytics project building
- Data preparation
- Feature engineering
- Knowledge discovery
- Insight extraction
- Communications of results
- Analytics operationalization
- Solution marketing
- Project management

Data Scientists

- Data/system requirement engineering
- Data system selection, installation, maintenance, management
- Data acquisition, extraction, integration
- Data quality enhancement
- Data transformation
- Data matching, loading, and sharing
- Data exploration
- Enhance data competency
- Analytics programming
- Data discovery system development
- Data discovery performance enhancement
- Computational performance enhancement
- Data-related social issue management
- Data system risk management

Data Engineers

**Responsibilities of Data Scientists  
vs. Data Engineers**

- Everyone claims to be a data scientist or is doing data science
- So what about you?

# Data Economy & Case Studies

Unacknowledged,  
for UTS Course ONLY



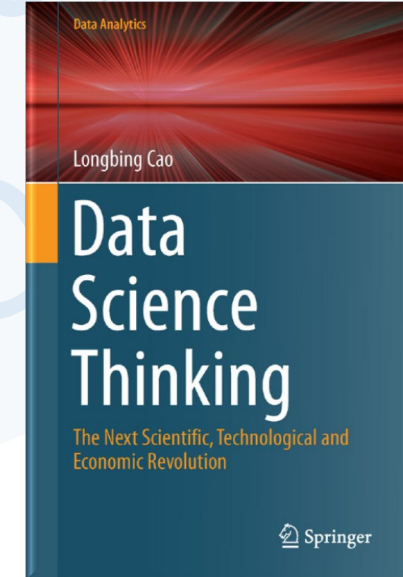
# Data economy family



# Case studies - discussion

## Chapter 9: 27 domains

- Advertising
- Aerospace and astronomy
- Arts, creative design and humanities
- Bioinformatics
- Consulting services
- Ecology and environment
- E-commerce and retail
- Education
- Engineering
- Finance and economy

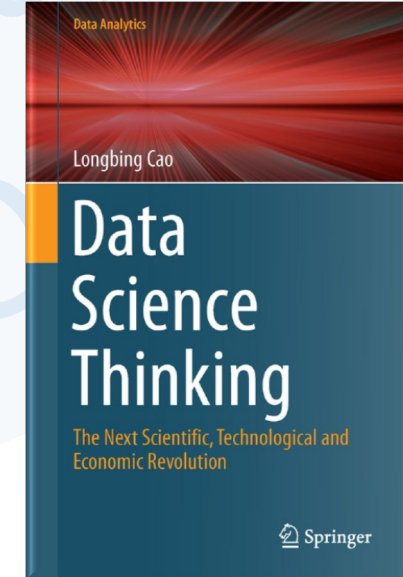


<https://link.springer.com/content/pdf/10.1007%2F978-3-319-95092-1.pdf>

# Case studies - discussion

## Chapter 9: 27 domains

- Gaming industry
- Government
- Healthcare and clinics
- Living, sports and entertainment
- Management, operations and planning
- Marketing and sales
- Medicine
- Physical-cyber-social society, networks
- Publishing and media

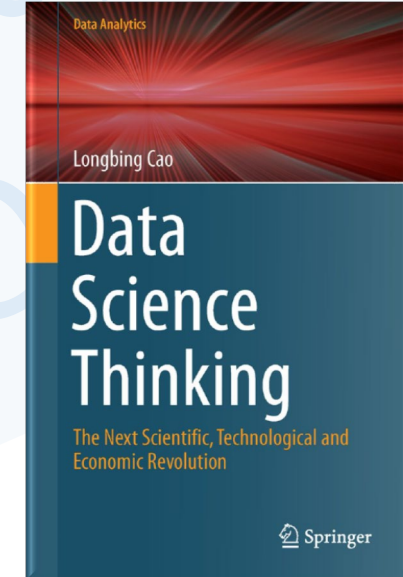


<https://link.springer.com/content/pdf/10.1007%2F978-3-319-95092-1.pdf>

# Case studies - discussion

## Chapter 9: 27 domains

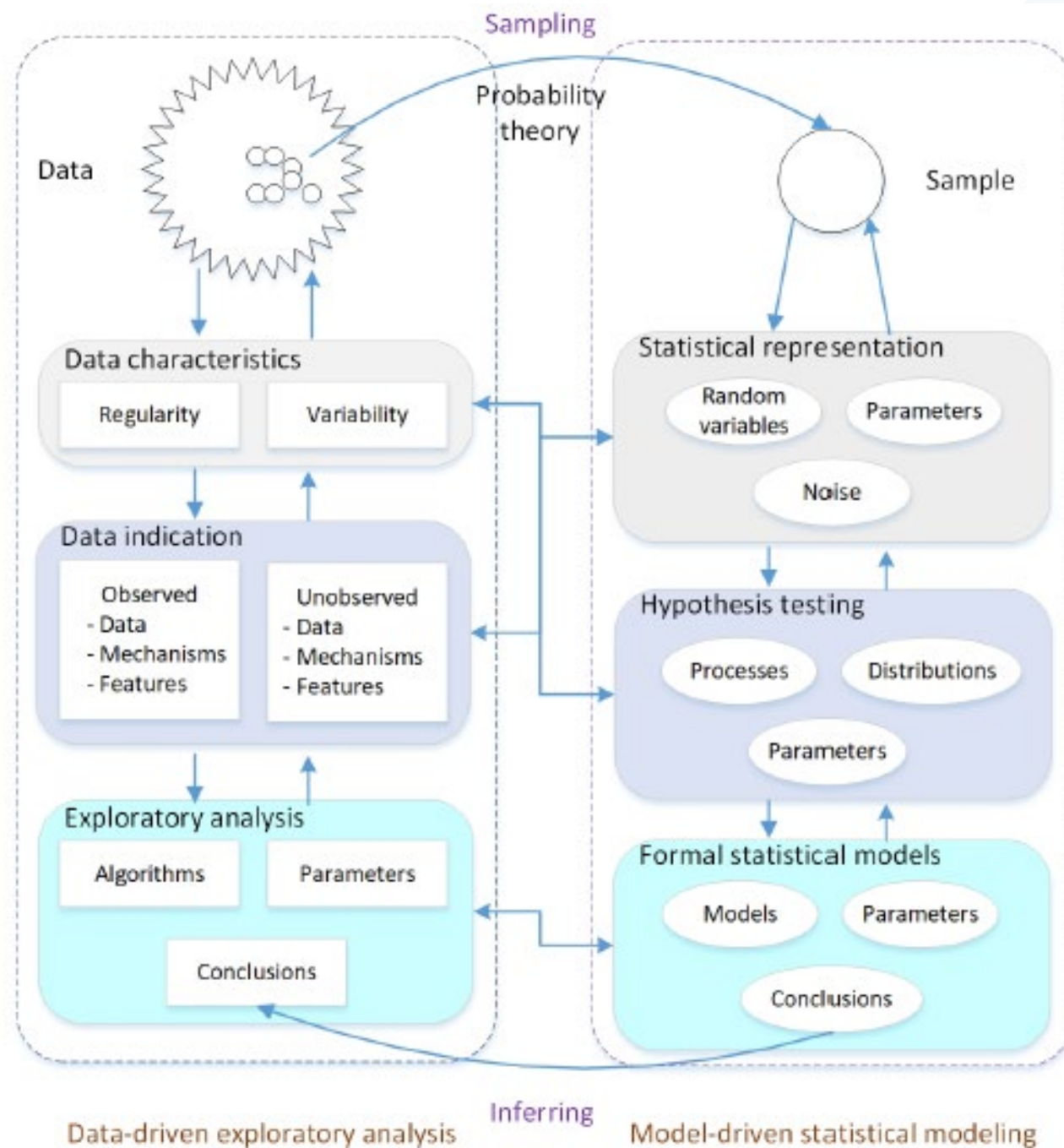
- Recommendation services
- Science
- Security and safety
- Social sciences and problems
- Sustainability
- Telecom and mobile services
- Tourism and travel
- Transportation



<https://link.springer.com/content/pdf/10.1007%2F978-3-319-95092-1.pdf>



# Domain, data + model- driven statistical inference



# Evaluation - Knowledge actionability

- Interestingness & actionability

$$Int(p) = I(t_i(p), b_i(p))$$

$$Int(p) = t_o(\mathbf{x}, \tilde{p}) \wedge t_s(\mathbf{x}, \tilde{p}) \wedge b_o(\mathbf{x}, \tilde{p}) \wedge b_s(\mathbf{x}, \tilde{p})$$

$$Int(p) \rightarrow \hat{I}(\hat{t}_o(), \hat{t}_s(), \hat{b}_o(), \hat{b}_s()) \\ = \alpha \hat{t}_o() + \beta \hat{t}_s() + \gamma \hat{b}_o() + \delta \hat{b}_s()$$

$$AKD^{e,\tau,m \in M} \longrightarrow O_{p \in P}(Int(p)) \\ \rightarrow O(\alpha \hat{t}_o()) + O(\beta \hat{t}_s()) + \\ O(\gamma \hat{b}_o()) + O(\delta \hat{b}_s())$$

$$act(p) = O_{p \in P}(Int(p)) \\ \rightarrow O(\alpha \hat{t}_o(p)) + O(\beta \hat{t}_s(p)) + \\ O(\gamma \hat{b}_o(p)) + O(\delta \hat{b}_s(p)) \\ \rightarrow t_o^{act} + t_s^{act} + b_o^{act} + b_s^{act} \\ \rightarrow t_i^{act} + b_i^{act}$$

Table 3. Measurement of interest of data-driven versus domain-driven data mining.

Interest		Traditional Data-Driven	Domain-Driven
Technical	Objective	Technical objective <i>tech_obj()</i>	Technical objective <i>tech_obj()</i>
	Subjective	Technical subjective <i>tech_subj()</i>	Technical subjective <i>tech_subj()</i>
Business	Objective	—	Business objective <i>biz_obj()</i>
	Subjective	—	Business subjective <i>biz_subj()</i>
Integrative		—	Actionability <i>act()</i>

- Technical and business evaluation
- Objective and subjective evaluation
- Holistic/systematic evaluation

Knowledge Actionability: Satisfying Technical and Business Interestingness, International Journal of Business Intelligence and Data Mining (IJBIDM)

Many classic methods may not be actionable: Frequent pattern

mining and association rule

Combined Mining: Analyzing Object and Pattern Relations for Discovering and Constructing Complex but Actionable Patterns, WIREs Data Mining and Knowledge Discovery, 3(2): 140-155, 2013

# Combined mining

- Combined sources/patterns/methods, etc.

- Pair patterns

$$\mathcal{P} : \begin{cases} X_1 \rightarrow T_1 \\ X_2 \rightarrow T_2 \end{cases} \quad \mathcal{P} : \begin{cases} X_p \rightarrow T_1 \\ X_p \wedge X_e \rightarrow T_2 \end{cases}$$

- Cluster patterns

$$\mathcal{P} : \begin{cases} X_1 \rightarrow T_1 \\ \dots \\ X_k \rightarrow T_k \end{cases}$$

- Derivative patterns

$$\mathcal{P} : \begin{cases} X_p \rightarrow T_1 \\ X_p \wedge X_{e,1} \rightarrow T_2 \\ X_p \wedge X_{e,1} \wedge X_{e,2} \rightarrow T_3 \\ \dots \\ X_p \wedge X_{e,1} \wedge X_{e,2} \wedge \dots \wedge X_{e,k-1} \rightarrow T_k \end{cases}$$

An Example of Combined Pattern Clusters

Clusters	Rules	$X_p$	$X_e$		$T$	$Cnt$	$Conf$ (%)	$I_r$	$I_c$	$Lift$	$Cont_p$	$Cont_e$	$Lift$ of $X_p \rightarrow T$	$Lift$ of $X_e \rightarrow T$
		demographics	arrangements	repayments										
$\mathcal{P}_1$	$P_5$	marital:sin &gender:F &benefit:N	irregular	cash or post	A	400	83.0	1.12	0.67	1.80	1.01	2.00	0.90	1.79
	$P_6$		withhold	cash or post	A	520	78.4	1.00		1.70	0.89	1.89	0.90	1.90
	$P_7$		withhold & irregular	cash or post & withhold	B	119	80.4	1.21		2.28	1.33	2.06	1.10	1.71
	$P_8$		withhold	cash or post & withhold	B	643	61.2	1.07		1.73	1.19	1.57	1.10	1.46
	$P_9$		withhold & vol. deduct	withhold & direct debit	B	237	60.6	0.97		1.72	1.07	1.55	1.10	1.60
	$P_{10}$		cash	agent	C	33	60.0	1.12		3.23	1.18	3.07	1.05	2.74
$\mathcal{P}_2$	$P_{11}$	age:65+	withhold	cash or post	A	1980	93.3	0.86	0.59	2.02	1.06	1.63	1.24	1.90
	$P_{12}$		irregular	cash or post	A	462	88.7	0.87		1.92	1.08	1.55	1.24	1.79
	$P_{13}$		withhold & irregular	cash or post	A	137	85.7	0.96		1.86	1.18	1.50	1.24	1.57
	$P_{14}$		withhold & irregular	withhold	C	50	63.3	2.91		3.40	2.47	4.01	0.85	1.38

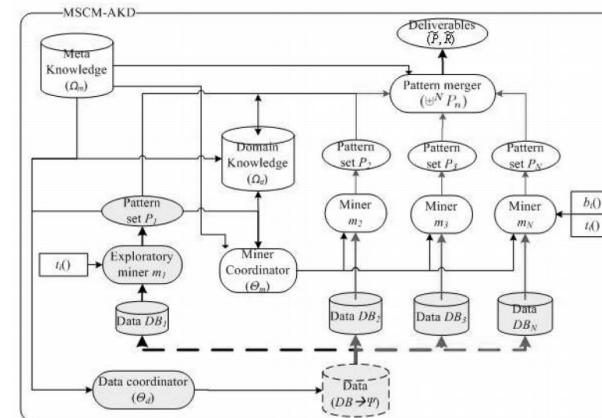


Fig. 5. Multisource + combined-mining-based AKD.

- Flexible Frameworks for Actionable Knowledge Discovery, IEEE Trans. on Knowledge and Data Engineering
- Combined Mining: Analyzing Object and Pattern Relations for Discovering and Constructing Complex but Actionable Patterns

# Impactful & impact-target collective behaviors

- Conditional impact ratio (*Cir*)

$$\begin{aligned}
 Cir(Q\bar{T}|P) &= \frac{Prob(Q\bar{T}|P)}{Prob(Q|P) \times Prob(\bar{T}|P)} \\
 &= \frac{Prob(PQ \rightarrow \bar{T}) / Prob(P)}{(Prob(PQ) / Prob(P)) \times (Prob(P \rightarrow \bar{T}) / Prob(P))} \\
 &= \frac{Prob(PQ \rightarrow \bar{T}) / Prob(PQ)}{Prob(P \rightarrow \bar{T}) / Prob(P)}.
 \end{aligned}$$

- Conditional Piatetsky-Shapiro's (P-S) ratio (*Cps*)

$$\begin{aligned}
 Cps(Q\bar{T}|P) &= Prob(Q\bar{T}|P) - Prob(Q|P) \times Prob(\bar{T}|P), \\
 &= \frac{Prob(PQ \rightarrow \bar{T})}{Prob(P)} - \frac{Prob(PQ)}{Prob(P)} \times \frac{Prob(P \rightarrow \bar{T})}{Prob(P)}
 \end{aligned}$$

- Business impact and utility

$$\alpha_s = \sum b_i \times p_i \times v_i$$

$$\beta_s = \sum |b_i| \times \beta_i \times p_i \times v_i$$

$$SR = (R_p - R_f) / \sigma_p$$

$$TR = \frac{\sum_{i=1}^u AskPrice_i * AskVolume_i - \sum_{j=1}^v BidPrice_j * BidVolume_j}{TotalInvestment}$$

$$IR = (\sum_{i=1}^n (Index_{i+1} - Index_i) / Index_i) / n$$

$$\begin{cases}
 PLN \rightarrow T \\
 PLN, DOC \rightarrow T \\
 PLN, DOC, DOC \rightarrow T \\
 PLN, DOC, DOC, DOC \rightarrow T \\
 PLN, DOC, DOC, DOC, REA \rightarrow T \\
 PLN, DOC, DOC, DOC, REA, IES \rightarrow T.
 \end{cases}$$

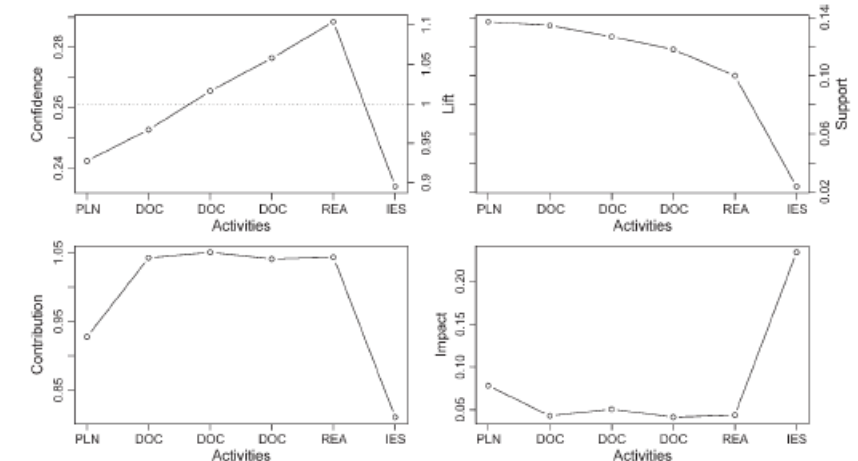
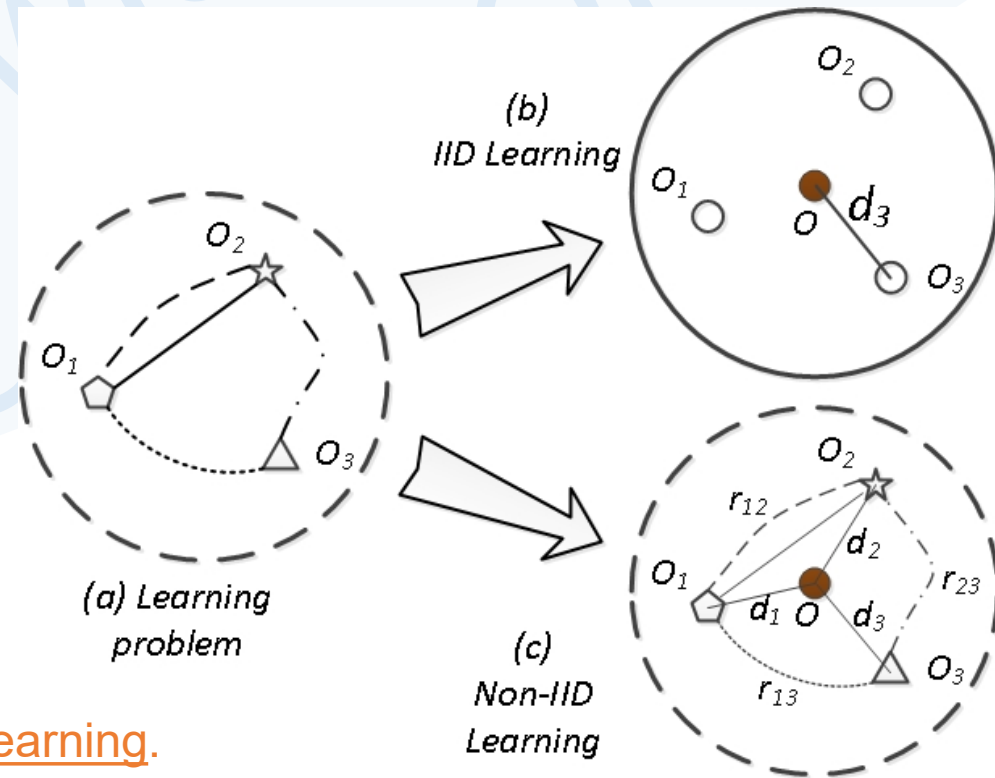


Fig. 3. Dynamic charts showing the dynamics of incremental cluster patterns.

Mining Impact-Targeted Activity Patterns in Imbalanced Data, IEEE Trans. on Knowledge and Data Engineering

# Beyond i.i.d. – Non-IID learning

- Outcomes to be delivered by IID analytical/learning methods/algorithms on non-IID data could be:
  - incomplete
  - biased, or even
  - misleading



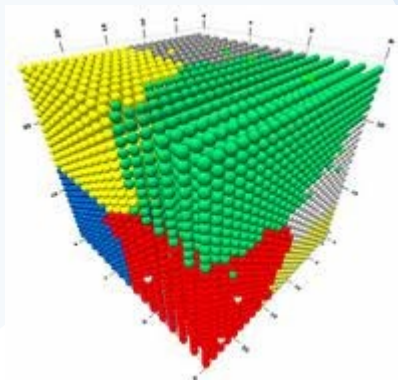


# IID K-means

acct_id	trade_date	trade_time	sec_code	trade_price	trade_vol	trade_dir	seat_code	trade_bal
210266501	20090106	112138	600331	5.63	200	B	51721	200
315726605	20090106	92500	600477	7.4	400	B	73061	2000
315726605	20090106	92500	600477	7.4	1200	B	73061	3200
315726605	20090106	145838	600477	7.64	1600	S	73061	1600
315726605	20090107	93952	600477	7.67	1600	B	73061	3200
315726605	20090106	92500	600547	48	400	B	73061	1200
315726605	20090106	95552	600547	49.14	200	S	73061	1000
315726605	20090106	95756	600547	49.1	200	S	73061	800
783486703	20090106	92500	600001	3.32	1000	B	46451	6000
783486703	20090106	92500	600001	3.32	1000	B	46451	7000



**Clustering**



Objective functions:

-K-means

$$\arg \min_{\mathbf{S}} \sum_{i=1}^k \sum_{\mathbf{x}_j \in S_i} \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2$$

-FCM

$$J_{\text{FCM}}(\boldsymbol{\mu}, \mathbf{A}) = \sum_{i=1}^c \sum_{j=1}^n (\mu_{ij})^m \|\mathbf{x}_j - \mathbf{a}_i\|^2$$

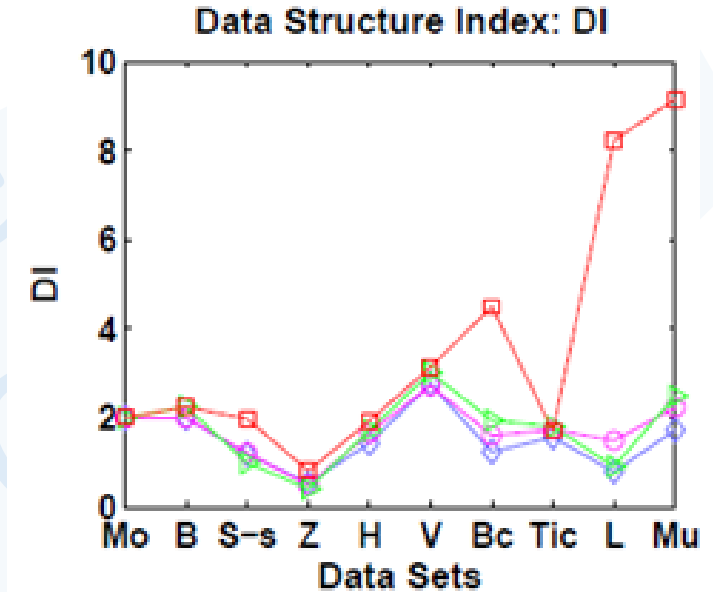
$$\sum_{i=1}^c \mu_{ij} = 1 \quad \text{for all } j \in J.$$

Note:

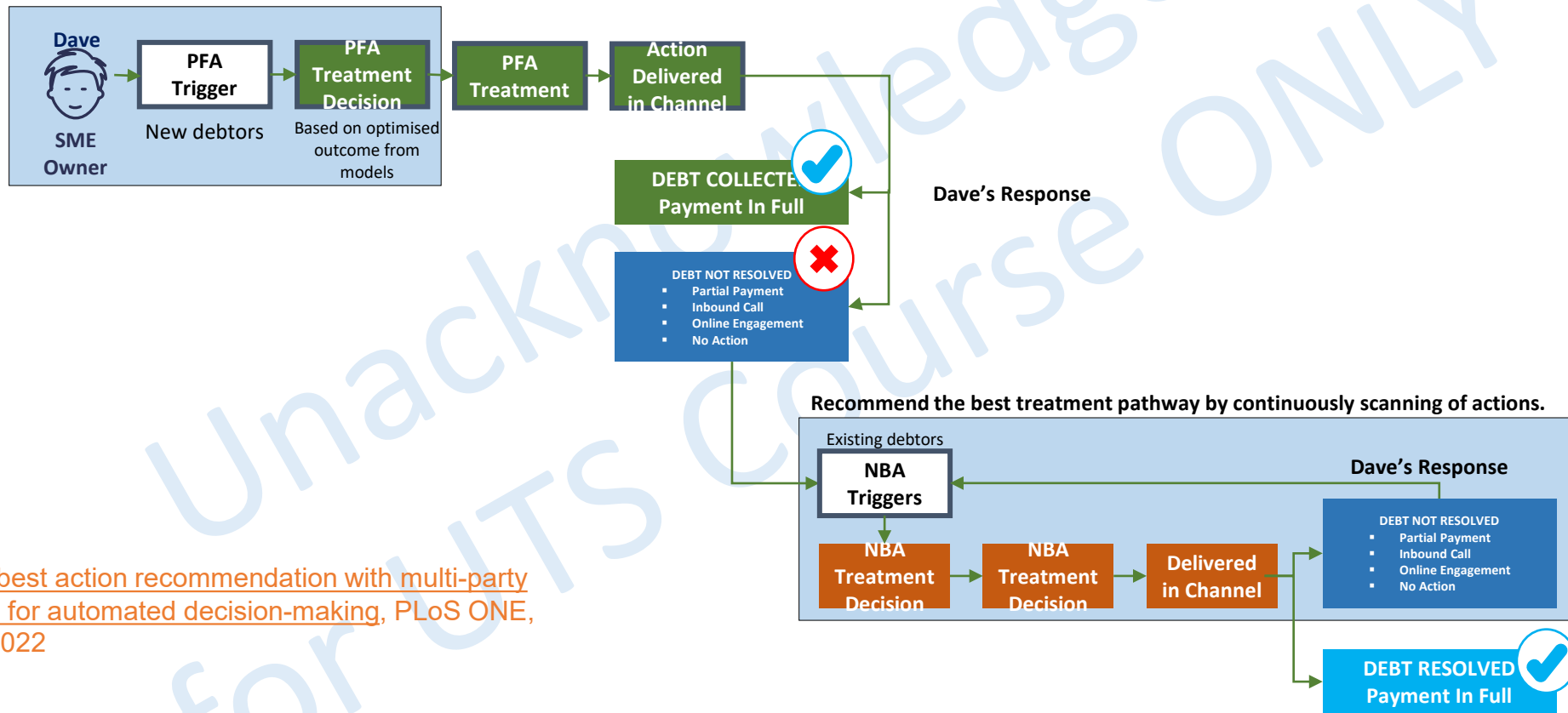
-  $X_j$  Individual objects only!

Question:

- How about  $X_{j1}$  and  $X_{j2}$  dependent?



# Non-linear, Tailored Client Engagement



[Personalized next-best action recommendation with multi-party interaction learning for automated decision-making, PLoS ONE, 17\(1\): e0263010, 2022](#)

# PNBA learning framework

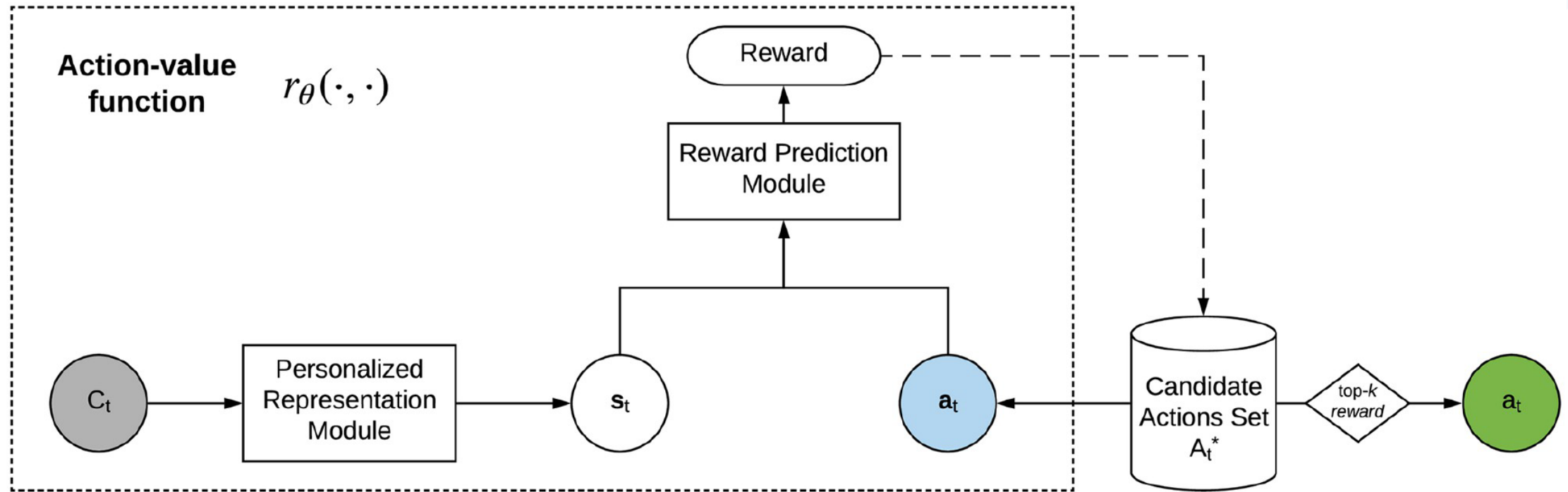
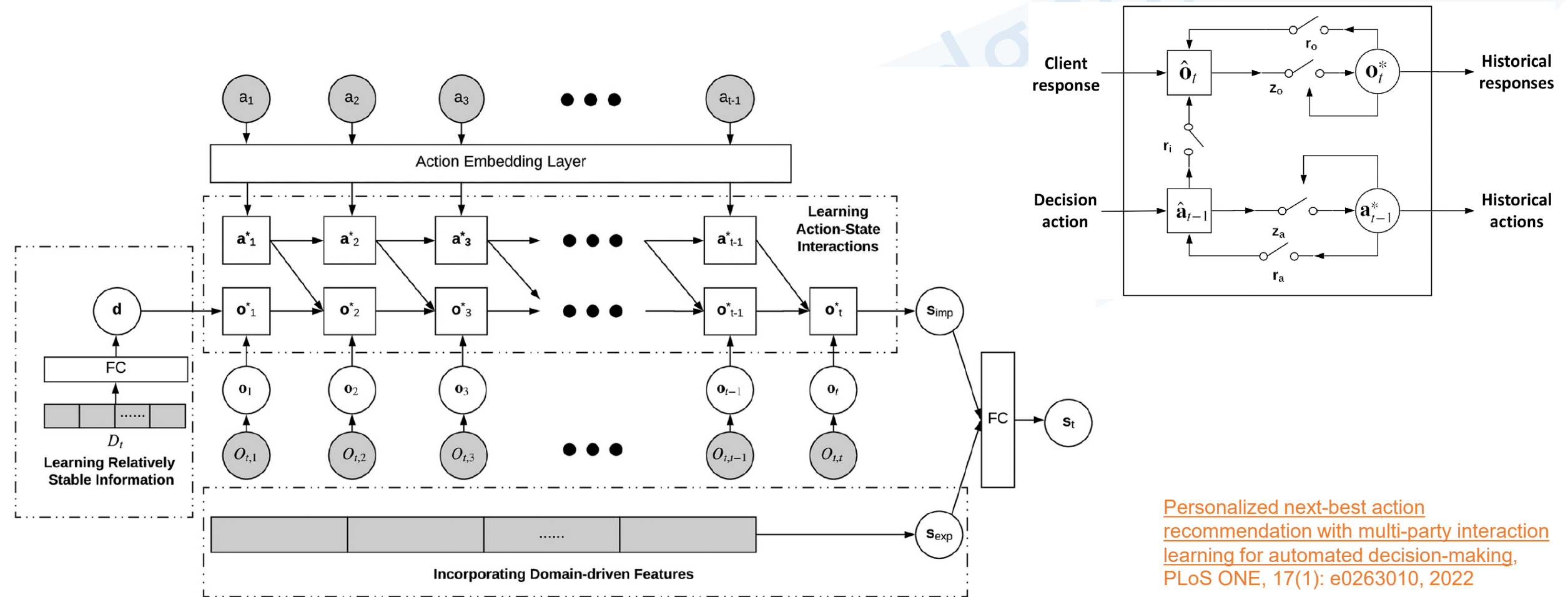


Fig 2. The framework for modeling the next-best action-oriented personalized decision-making.

[Personalized next-best action recommendation with multi-party interaction learning for automated decision-making, PLoS ONE, 17\(1\): e0263010, 2022](#)

# Learn personalized interaction representation



Personalized next-best action recommendation with multi-party interaction learning for automated decision-making, PLoS ONE, 17(1): e0263010, 2022

Fig 3. A reinforced coupled recurrent network to learn personalized client representation.

# PNBA: Case studies

- Non-Markovian NBA recommendation

Table 2. Average reward lift for 10 actions recommended by 11 deep models over the review measured by domain-driven debt collection rules.

Model	A1	A2	A3	A4	A5	A6	A7	A8	A9	A10	Total_Avg	Action_Avg
CRN_IMB	5	4	3.0534	2.8752	6.8	2.1415	2.6984	3.3567	1.6772	2.9969	2.5569	3.4599
CRN	2.1957	3.5383	2.2068	2.6616	3.216	2.074	2.326	2.6277	1.7654	2.3425	2.1942	2.4954
WD	2.604	1.5992	2.0979	2.2798	3.2239	1.9824	2.2629	2.6967	0.9899	2.312	2.1089	2.2049
LSTM	0.9722	1.0987	0.9391	0.974	1.1272	1.0159	0.897	1.1097	1.1024	1.0847	1.0013	1.0321
WD_LSTM	2.0471	1.2731	1.9709	2.4755	2.2217	1.8129	2.0816	2.1909	1.1405	2.105	1.9198	1.9319
WD_Res_LSTM	1.7247	0.8219	1.7007	1.9816	2.4985	1.8164	1.9851	2.0921	0.8285	1.967	1.8488	1.7416
WD_Multi_LSTM	1.684	1.0468	1.6591	1.774	1.6924	1.7083	1.671	2.1678	1.2222	1.8098	1.7161	1.6435
GRU	0.5783	0.0865	0.9852	1.1201	1.5022	0.9154	0.861	0.9463	1.0347	1.0416	0.9345	0.9071
WD_GRU	1.0049	0.6397	1.3454	1.7369	2.1271	1.6489	1.6049	2.1562	0.665	1.6602	1.611	1.4589
WD_Res_GRU	1.4488	1.1333	1.7364	1.3479	2.2259	1.6932	1.7091	1.9582	1.2507	1.8869	1.7248	1.6391
WD_Multi_GRU	1.6329	1.8399	1.9114	1.7949	1.8781	1.8206	2.0276	1.7613	1.0508	2.2347	1.8959	1.7952
$\Delta$ _IMB	92.01%	117.40%	45.55%	16.15%	110.92%	8.03%	19.25%	24.47%	34.10%	29.62%	21.24%	56.92%
$\Delta$	-15.68%	92.31%	5.19%	7.52%	-0.25%	4.62%	2.79%	-2.56%	41.15%	1.32%	4.04%	13.18%



# Sampling non-IIDness in deep networks

- Sampling non-IIDnesses
- Distributional vulnerability
- High confidence prediction on out-of-distribution samples

Label vs  
distributional  
fitting

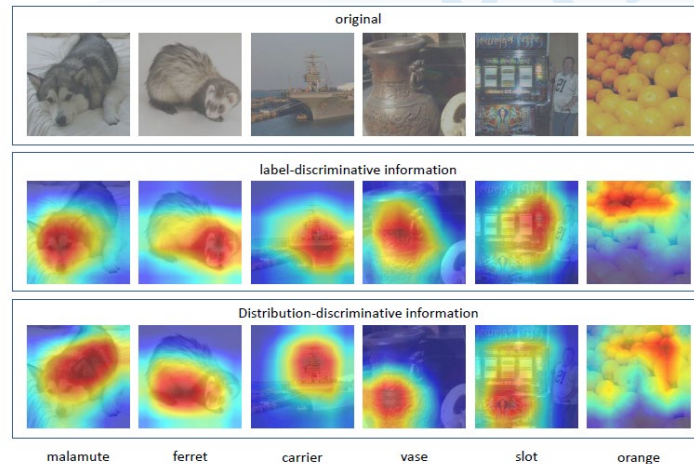


Fig. 6.8 DRL: Heat maps of Grad-CAM for label- and distribution-discriminative representations.

Red regions correspond to high scores for class, while blue regions correspond to low scores. The figure is best viewed in color.

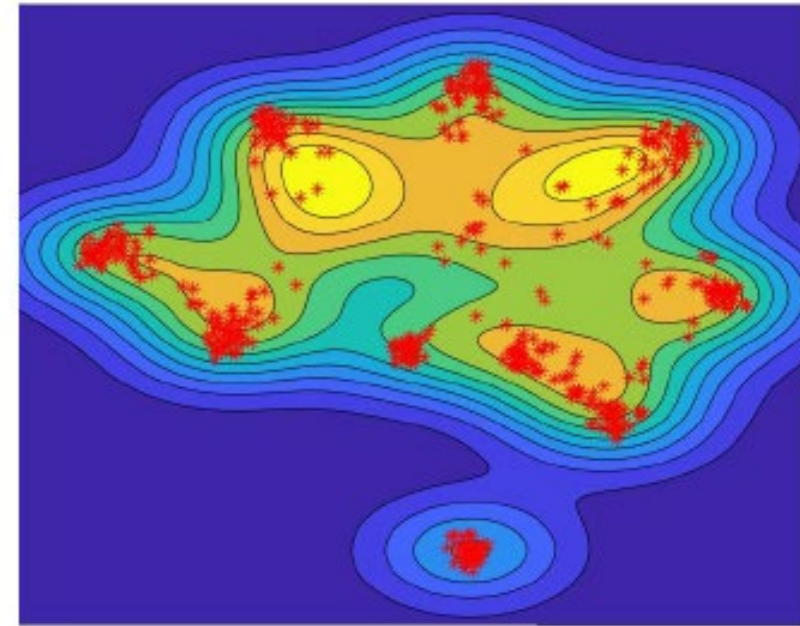
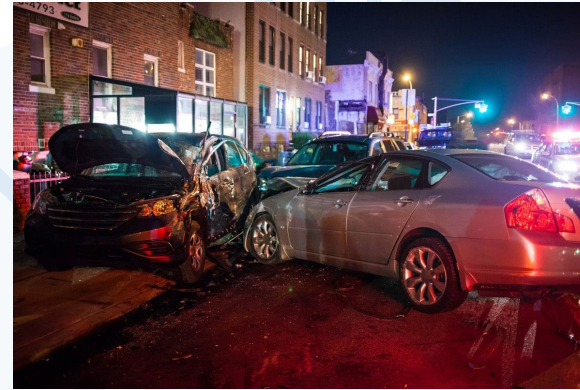


Fig. 1. The heat map of prediction confidence by ResNet18 on CIFAR10. The embedding results are constructed by t-SNE [7]. Red points correspond to training in-distribution samples. Yellow regions correspond to high confidence for predictions, while blue regions correspond to low confidence. ResNet18 assigns high-confidence predictions on samples located in the regions outside the training in-distribution samples, i.e., out-of-distribution samples. It shows ResNet18 does not discriminate between in- and out-of-distribution samples. The figure is best viewed in color.

# Various vulnerabilities of deep networks

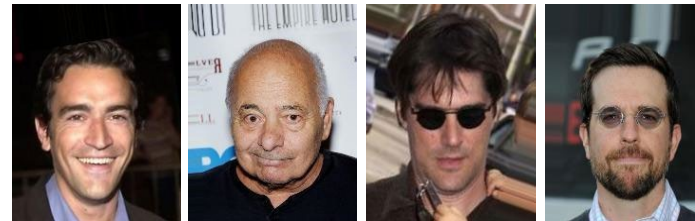
- Network, distribution, domain... vulnerability
  - **Face recognition**: incorrectly recognize a stranger as a person authenticated by the system
  - **Driverless car**: a high-confidence action at an unknown situation, which should be passed to the human driver for handling, may cause a serious accident
  - **Deep fake**



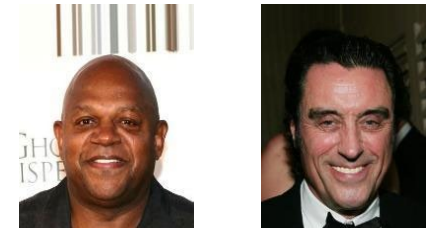
[https://fee.org/media/27691/car\\_crash.jpg?anchor=center&mode=crop&width=1200&rnd=131672371490000000](https://fee.org/media/27691/car_crash.jpg?anchor=center&mode=crop&width=1200&rnd=131672371490000000)



ID: Authenticated



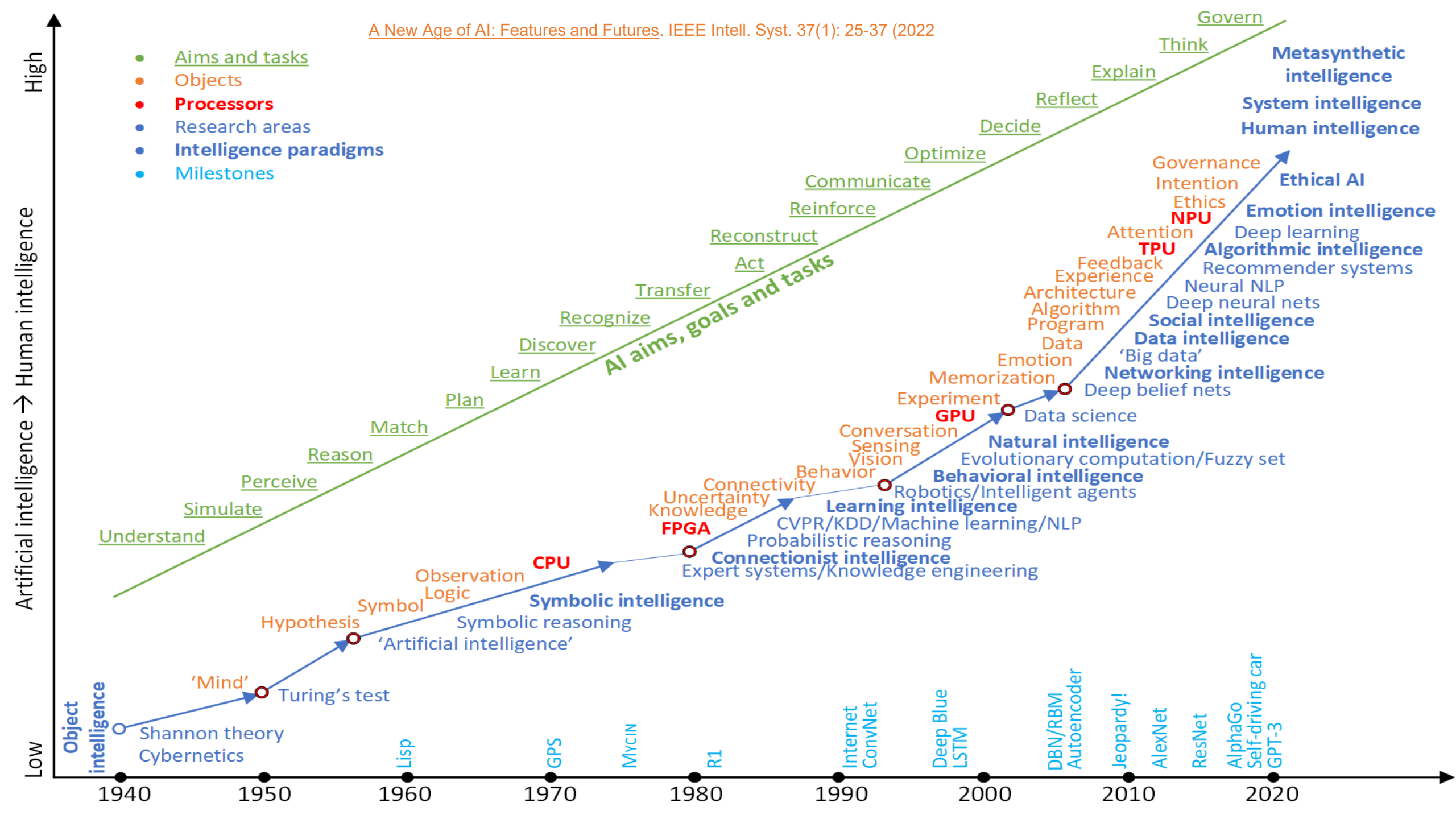
OOD: Unauthenticated



Mistakenly recognized as

# Concluding Remarks

Unacknowledged,  
for UTS Course ONLY

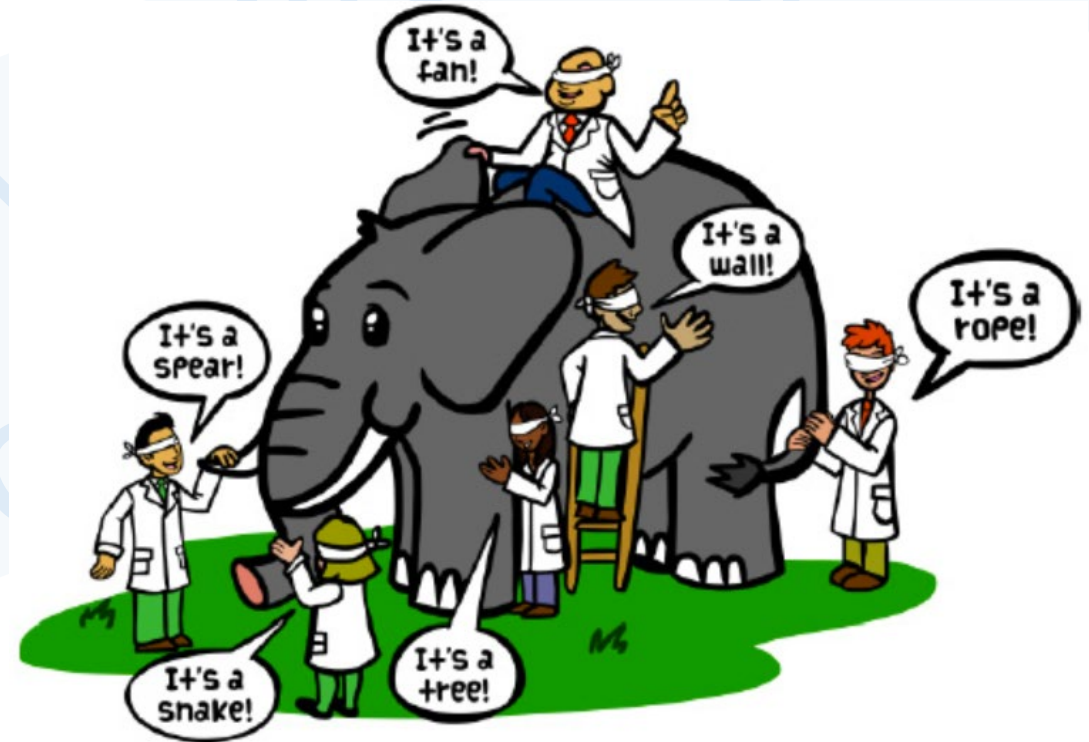




# We do not know what we do not know

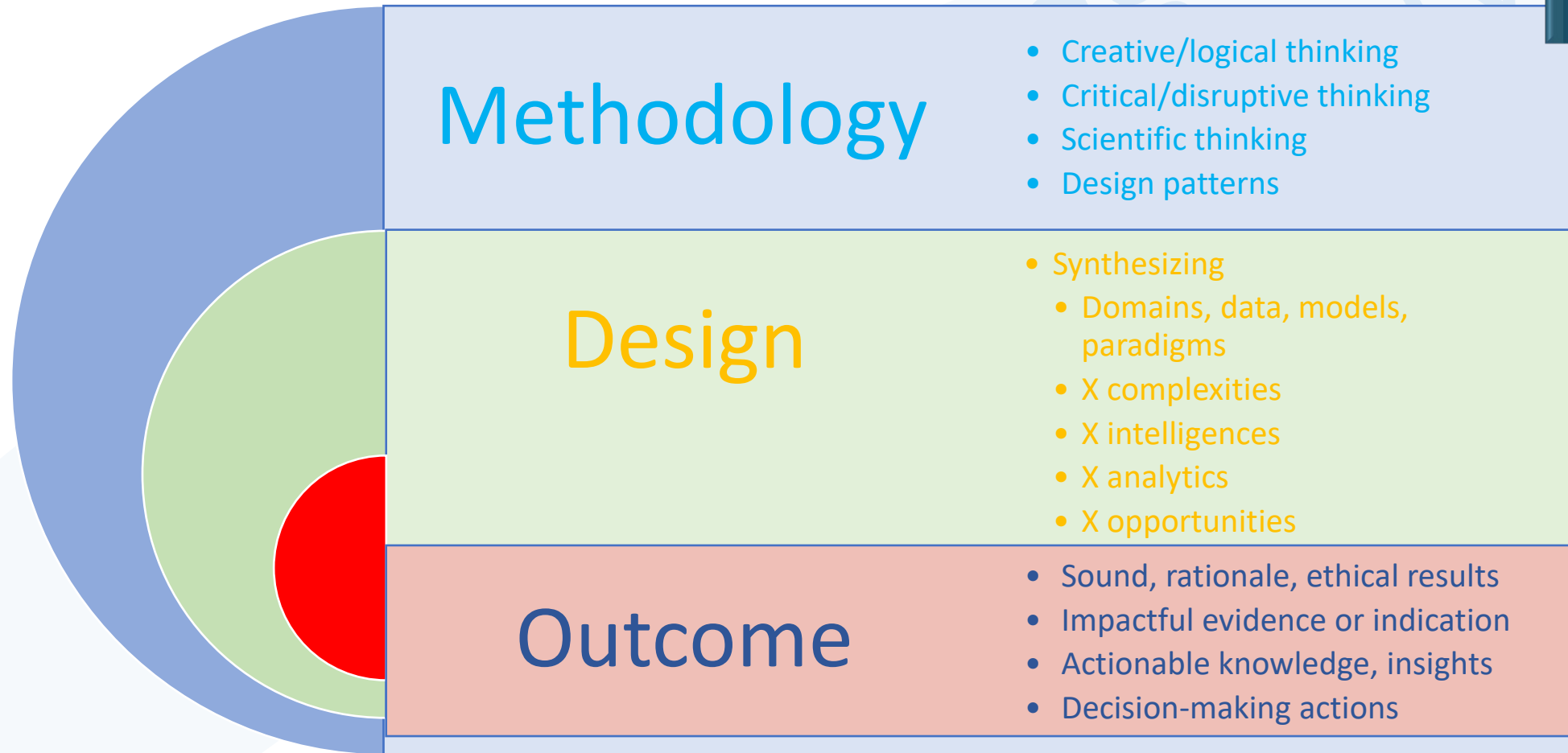
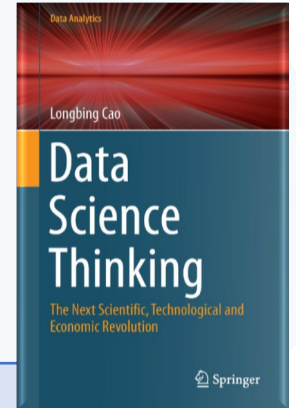
- How can AI enable blind people to tell a genuine story about elephant?

- Beyond IID thinking
- Couplings between parts
- Heterogeneities between parts
- From touching/representation → analysis → reasoning/inference → summarization/integration
- Local – global picture (known → unknown)/optimization





# Data science thinking matters



# References

[www.datasciences.org](http://www.datasciences.org)

- David Donoho. [50 years of data science](#).
- Longbing Cao: Data Science Thinking, Springer, 2018. <https://www.springer.com/us/book/9783319950914>
- [Data Science: A Comprehensive Overview](#). ACM Computing Surveys, 50(3), 43:1-42, 2017. <https://arxiv.org/abs/2007.03606>
- [Data Science: Nature and Pitfalls](#). IEEE Intelligent Systems, Volume: 31, Issue: 5, 66-75, 2016. <https://arxiv.org/abs/2006.16964>
- [Data Science: Challenges and Directions](#). Communications of the ACM, Vol. 60 No. 8, Pages 59-68, 2017. <https://arxiv.org/abs/2006.16966>
- [Data Science: Profession and Education](#). IEEE Intelligent Systems, 34(5): 35-44, 2019. [https://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=3755747](https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3755747)

<https://link.springer.com/content/pdf/10.1007%2F978-3-319-95092-1.pdf>

# References

[www.datasciences.org](http://www.datasciences.org)

- [Data Mining for Business Applications](#), ISBN: 978-0-387-79420-4, Springer, 2008
- [AI in Finance: Challenges, Techniques and Opportunities](#), SSRN, 1-40, 2021
- [Data science and AI in FinTech: An overview](#). Int. J. Data Sci. Anal. 12(2): 81-99, 2021
- [Effective detection of sophisticated online banking fraud on extremely imbalanced data](#), World Wide Web, 16:449–475, 2013
- [Copula-Based High Dimensional Cross-Market Dependence Modeling](#), DSAA2017 Research Track, 734-743, 2017
- An effective contrast sequential pattern mining approach to taxpayer behavior analysis, World Wide Web 19(4): 633-651, 2016
- Personalized next-best action recommendation with multi-party interaction learning for automated decision-making, PLOS One, 17(1): e0263010, 2022
- Table2Vec-automated universal representation learning of enterprise data DNA for benchmarkable and explainable enterprise data science. Sci Rep 11, 23957, 2021
- Combined mining: Analyzing object and pattern relations for discovering and constructing complex yet actionable patterns. WIREs Data Mining Knowl. Discov. 3(2): 140-155, 2013
- Domain Driven Data Mining, ISBN: 978-1-4419-5737-5, Springer, 2010