

The Essentials of Data Science and Machine Learning: How Machine Learning Extracts Knowledge From Data

Peter Krensky

AI Renaissance or Apocalypse?



Key Issues

1. How do you navigate the hype and semantics and deal with the shadow of AI?
2. How does it all work?
3. What can you do right now with data science and machine learning?

Key Issues

1. How do you navigate the semantics and deal with the shadow of AI?
2. How does it all work?
3. What can you do with data science and machine learning right now?

se-man-tics (sə'man(t)iks) *n.*
The meanings of words and
phrases in a particular context

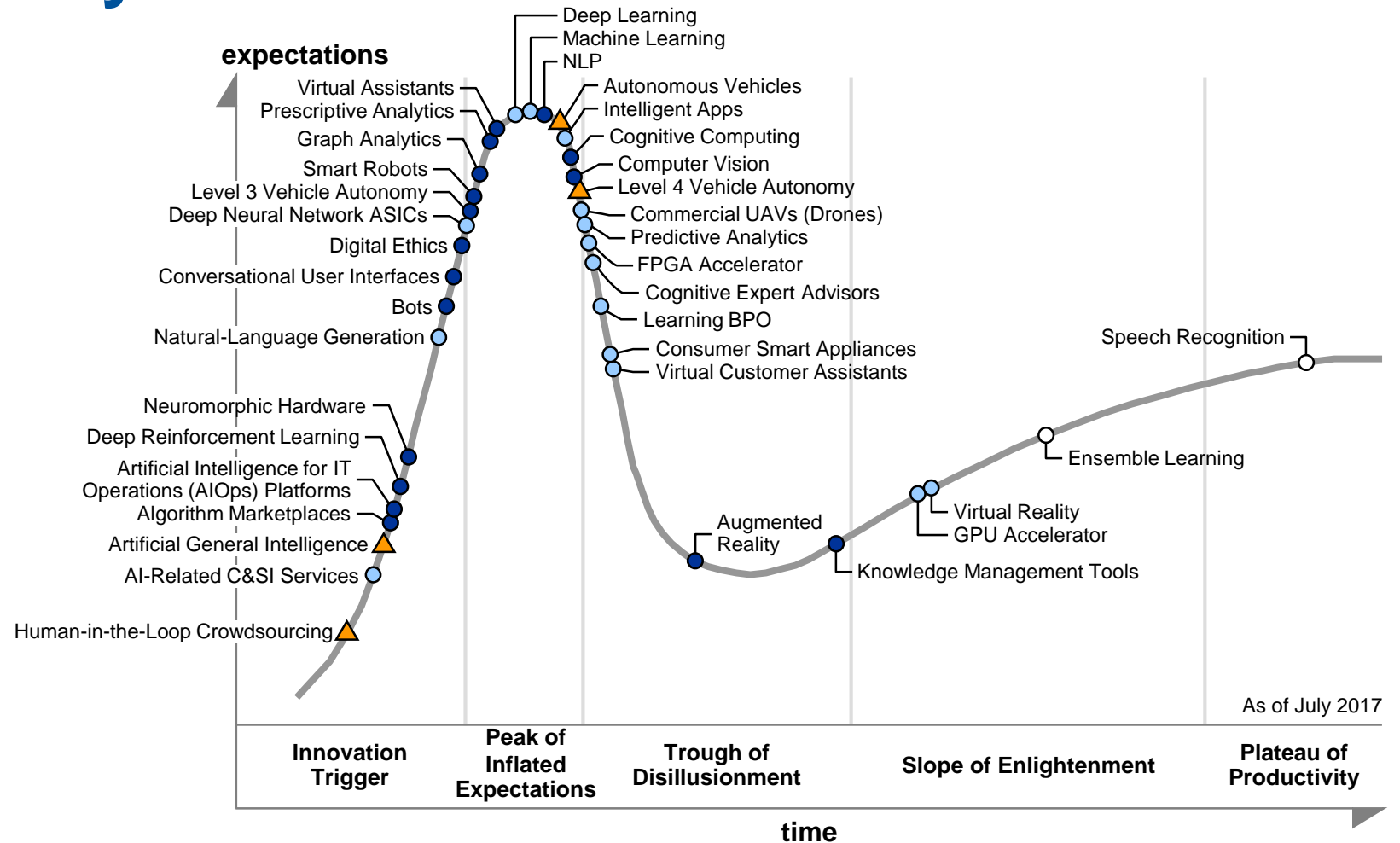
The Gartner AI Hype Cycle

Very early maturity levels

- **86%** of tech profiles (dots) headed to the bottom of Trough of Disillusionment
- **54%** not expected to plateau and deliver reliable productivity for mainstream buyers until 2022 or later

Huge potential

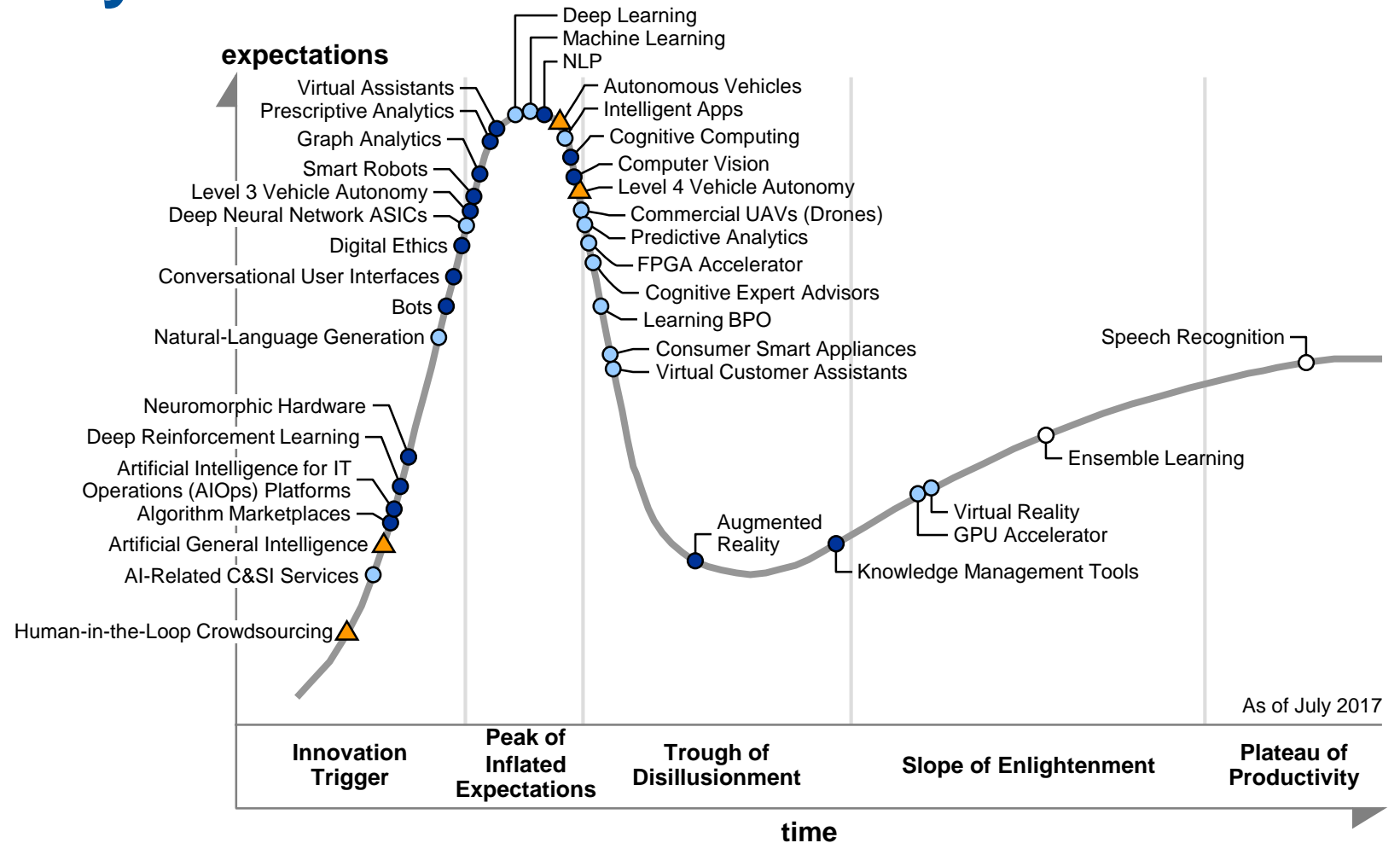
- **41%** offer transformational benefits
- **44%** offer high benefits



The Gartner AI Hype Cycle

Though the lens of D&A

- **>50%** of hype has nothing do with data and analytics
- Marketers stick the AI moniker on anything remotely algorithmic
- Key data science concepts are crowded around the peak

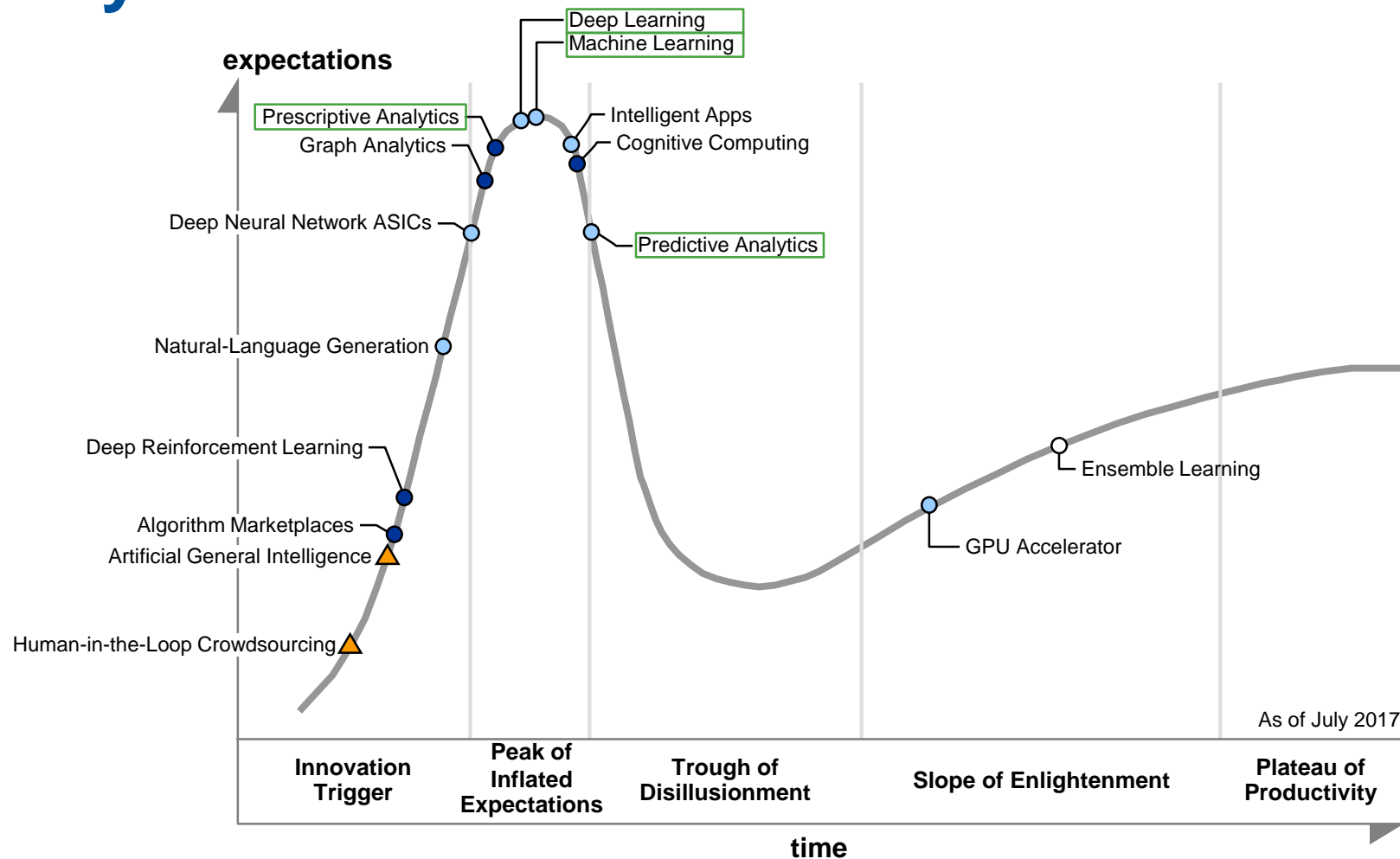


Source: ["Hype Cycle for Artificial Intelligence, 2017,"](#) 24 July 2017 (G00314732)

The Gartner AI Hype Cycle

Through the lens of D&A

- >50% of hype has nothing to do with data and analytics
- Marketers stick the AI moniker on anything remotely algorithmic
- Key data science concepts are crowded around the peak



Years to mainstream adoption:

- less than 2 years ● 2 to 5 years ● 5 to 10 years ▲ more than 10 years ✕ obsolete before plateau

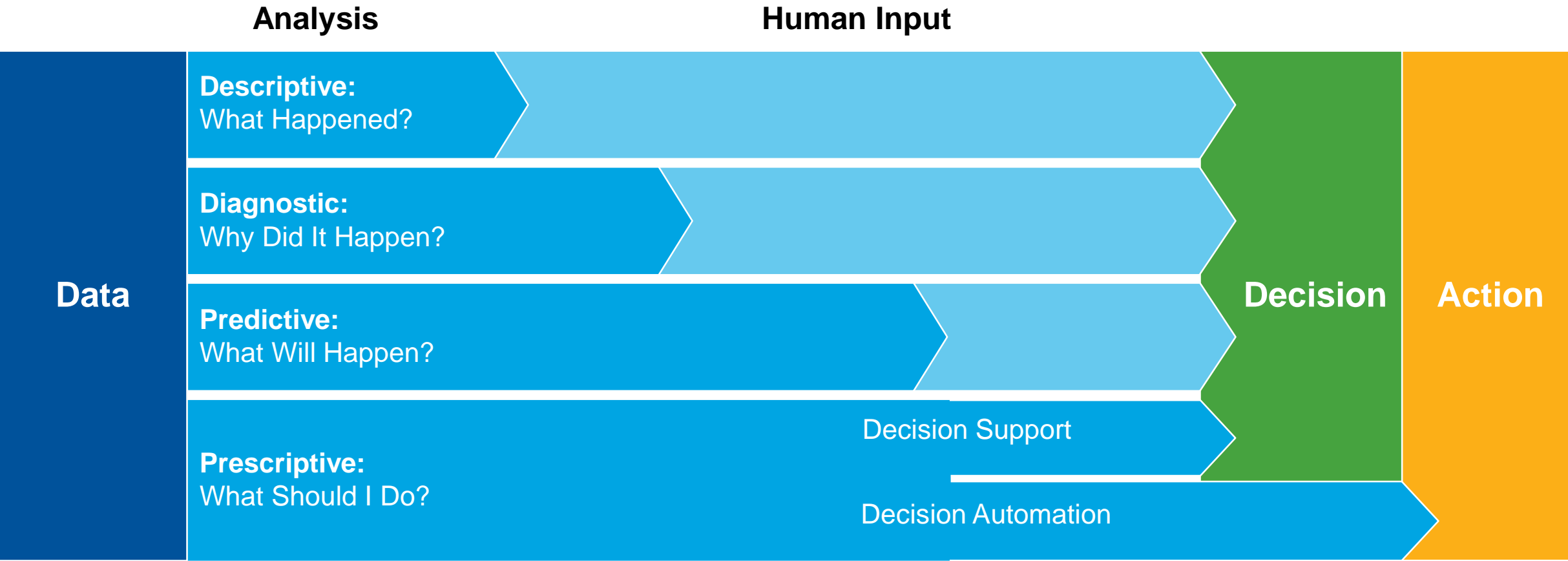
Hierarchy of Terms From a D&A Perspective

Artificial intelligence can be defined as a general approach to the simulation of cognitive processes by means of computer programs.

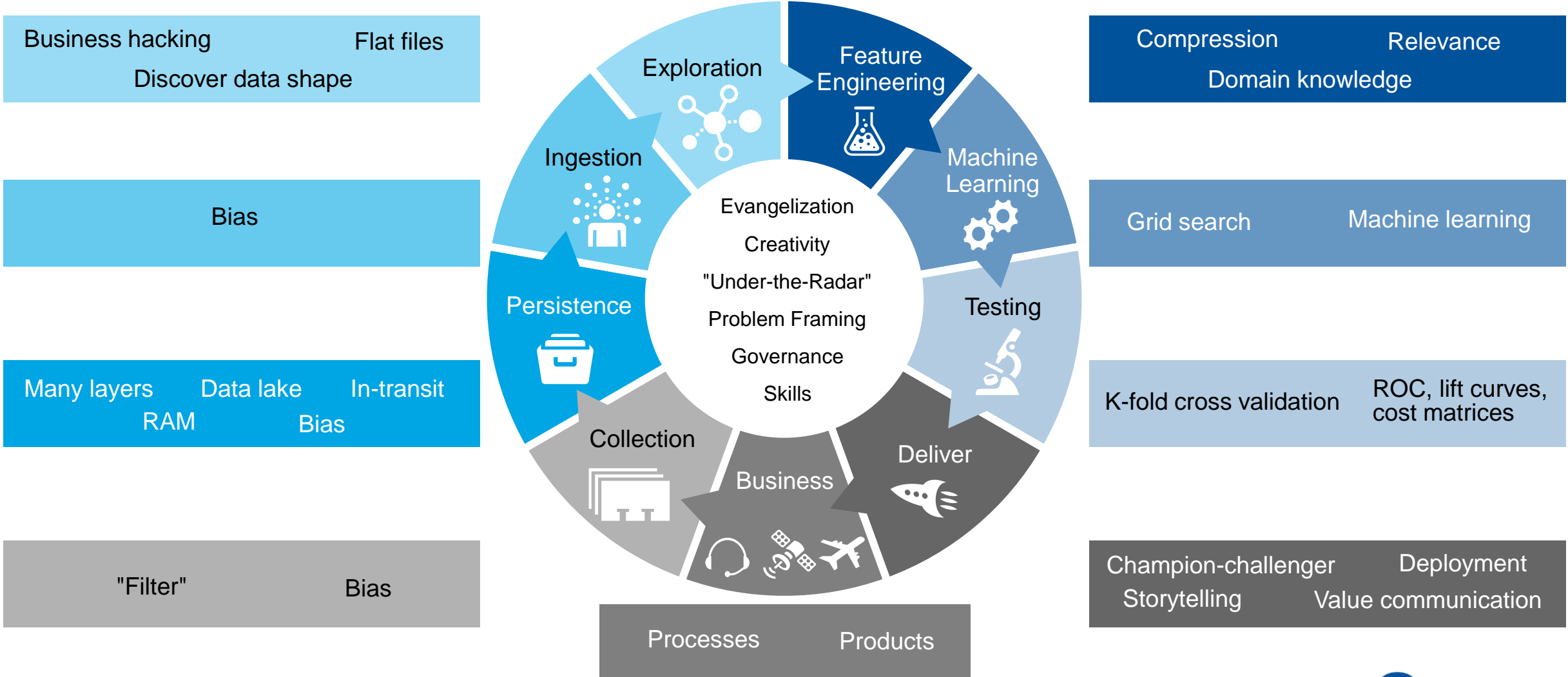
Data science and machine learning at its most basic is the practice of using algorithms to parse data, capture knowledge, learn from it, make a deterministic or predictive model and deploy that output into a business decision.

Deep learning is a subset of ML algorithms that creates knowledge from multiple layers of information processing.

The Analytics Spectrum



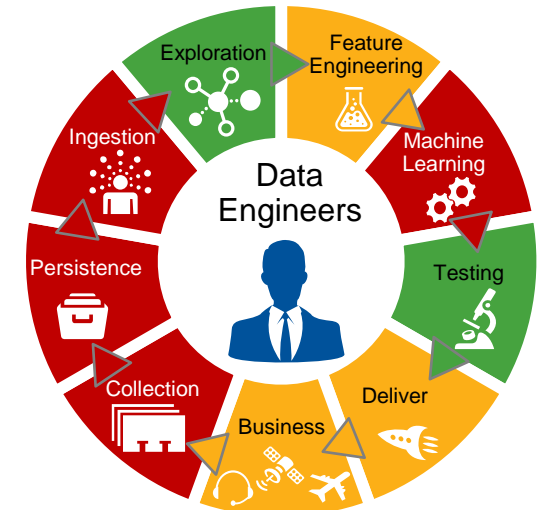
Data Science Is a Model Creation Cycle



Skill Map



Knowledge Spectrum



... at the Core, ML Is About Creating Mappings Informed by Input/Output Pairs

Type of Problem	Inputs	Outputs
Loan Application	Application data	Will the applicant repay the loan? (0 or 1)
Demand Prediction	Market situation	How many products will be bought? (n)
Self-Driving Cars	Car sensory data	Break, accelerate, tilt the wheel? (x, y)
Propensity to Buy	Profile and transactions	Will the customer buy or not? (0 or 1)
Failure Prediction	Sensor readings	Will a failure happen with 4 weeks (0 or 1)
Customer Churn	Profile and activities	Will customer cancel the contract? (0 or 1)
Medical Diagnosis	Pixel data from a retinal scan	Will the disease break out? (0 or 1)
Advertisement	Ad + context + user profile	Will the user click on ad? (0 or 1)

Landscape of ML Solutions

DYI

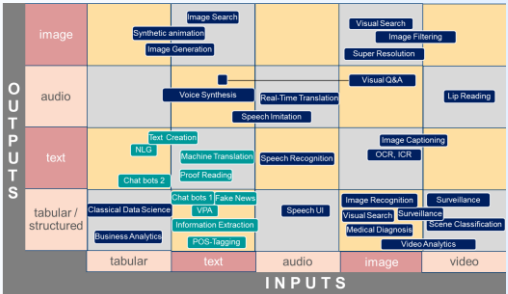
Salesforce Einstein

SAP Clea



Business Users

Embedded Machine Learning



Application Engineers

Machine-Learning APIs

Magic Quadrant Placeholder

ABILITY TO EXECUTE

COMPLETENESS OF VISION

As of June 2012

From "Magic Quadrant for Data Science Platforms." xx February xxxx (G00xxxxxx)

Data Scientists



Augmented Analytics

Data Science and Machine-Learning Platforms



ML Engineers

R, Python, Scala, Matlab

Deep-Learning Frameworks

Deep-Learning Cloud Platforms

Deep-Learning Hardware

Nvidia, AMD, IBM, Intel

Intel Nervana
Microsoft Azure
Rescale AWS

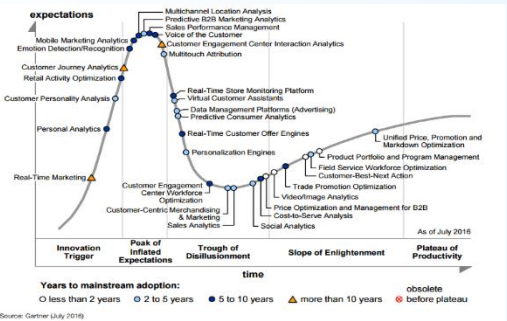
Google Cloud Platform

Data Analysis Software

Data Analysts



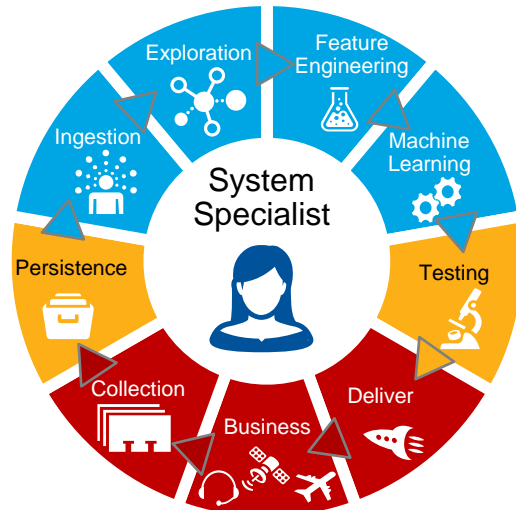
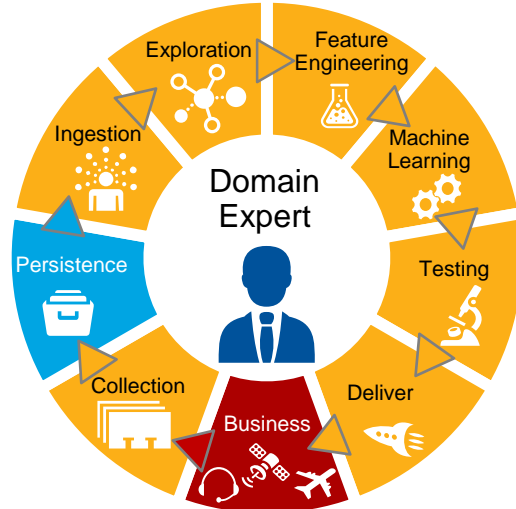
Buy



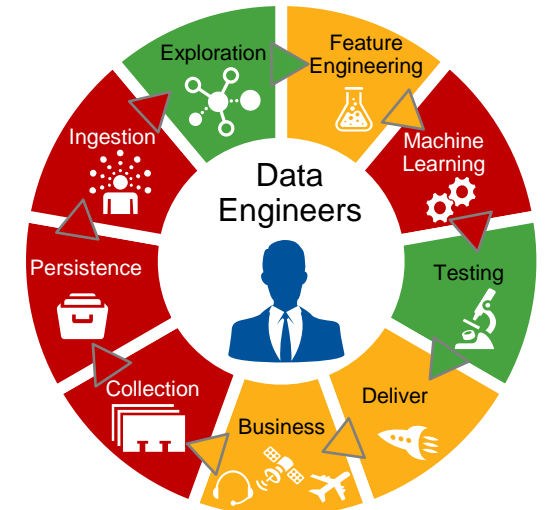
Key Issues

1. How do you navigate the semantics and deal with the shadow of AI?
2. How does it all work?
3. What can you do right now with data science and machine learning?

Skill Map



Knowledge Spectrum



Citizen Data Scientist

Finding, Keeping, Nurturing Skills

By 2020, >40% of data science tasks will be automated

Core Data Scientists



Citizen Data Scientists



Academia

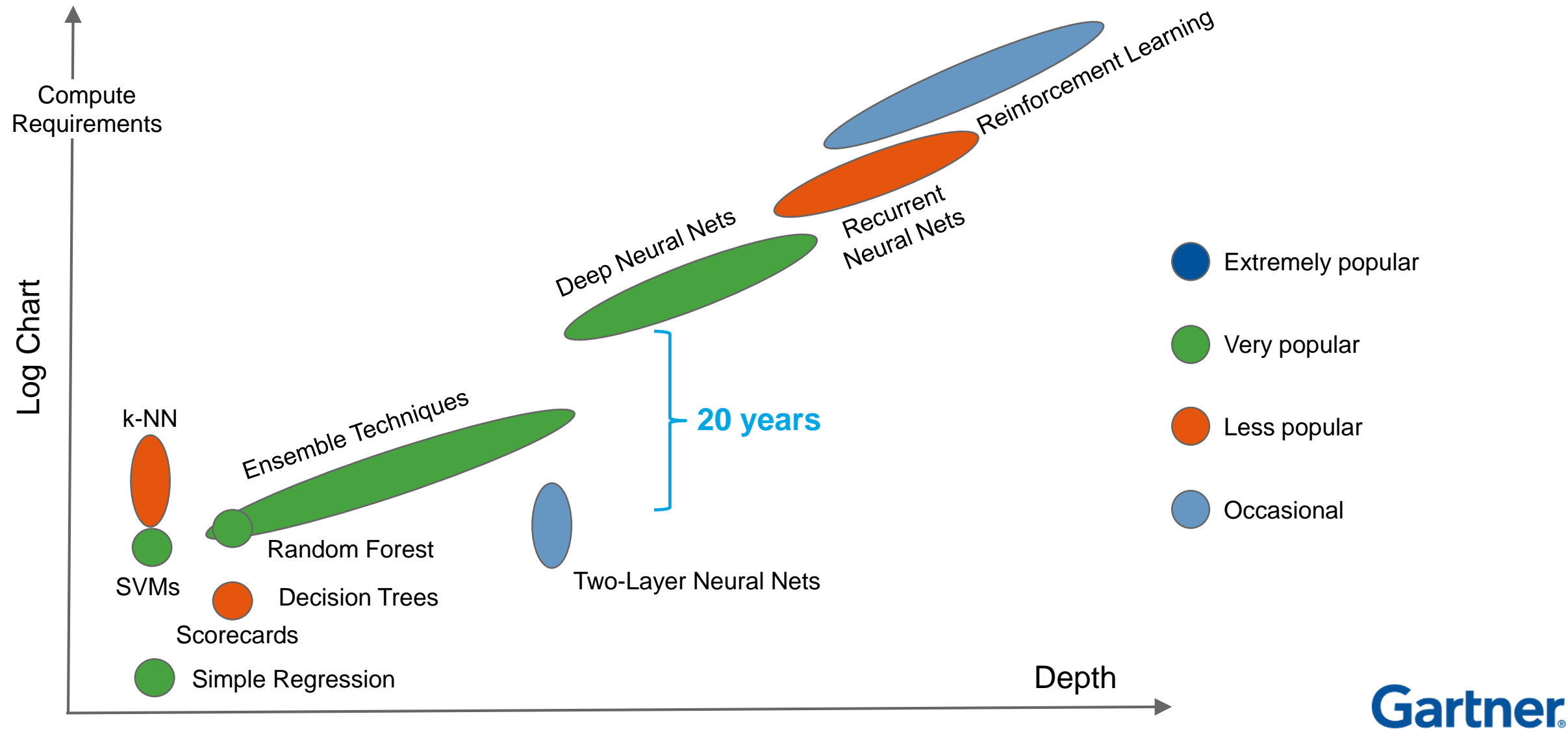


Consultants/
Freelancers

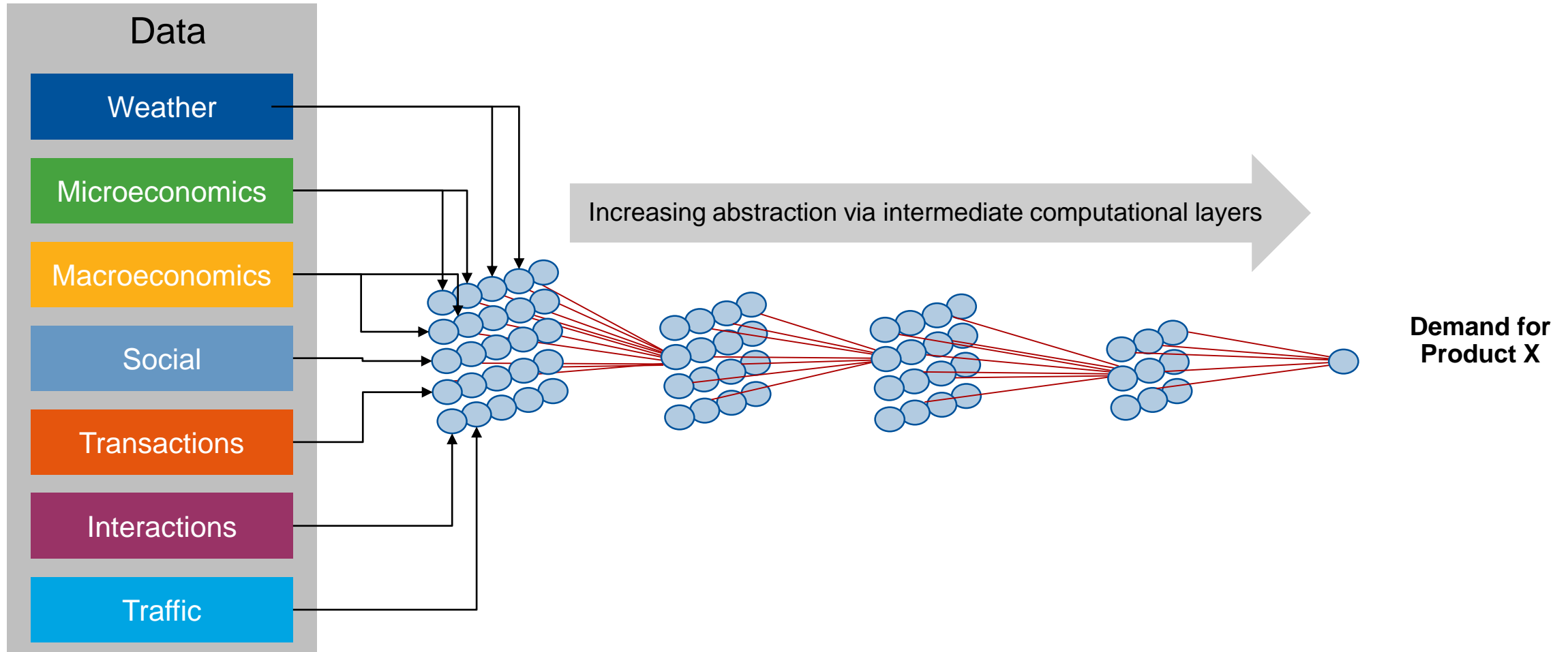


79% of DS teams currently use R
75% currently use Python

There Is a Zoo of Machine-Learning Approaches Out There ...

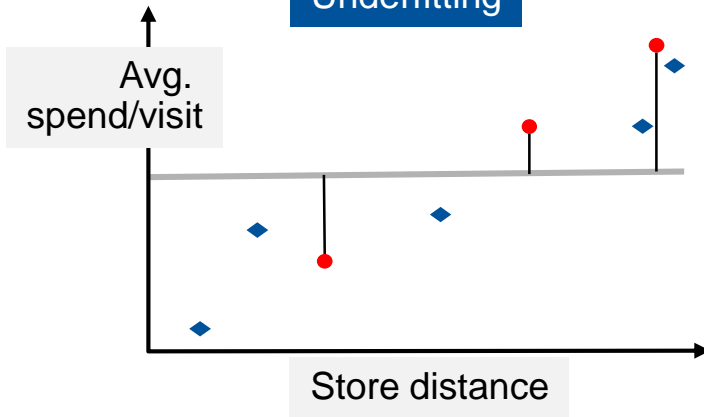


Deep Learning Addresses One of the Biggest "Big Data" Challenges: Data Fusion



The Challenge of Machine Learning: Under and Overfitting

Underfitting

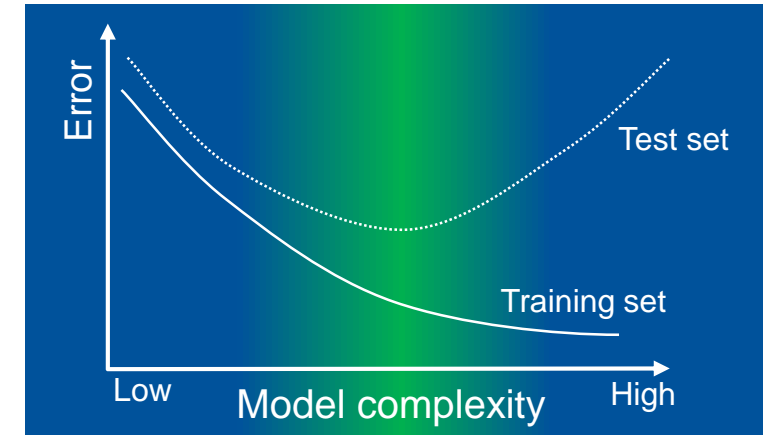
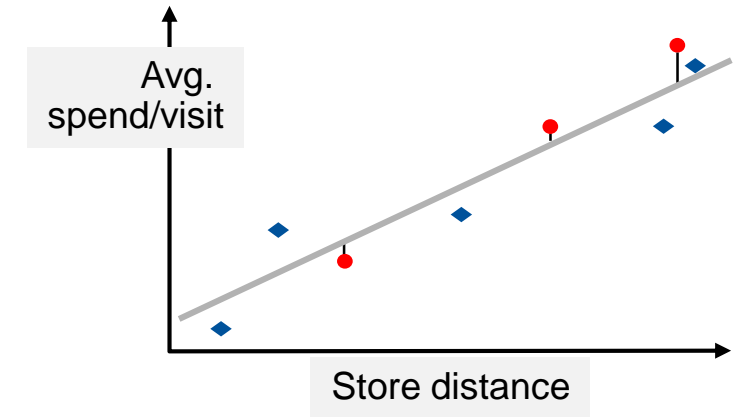


- Predictor is too "simplistic"
- Cannot capture the pattern

Overfitting

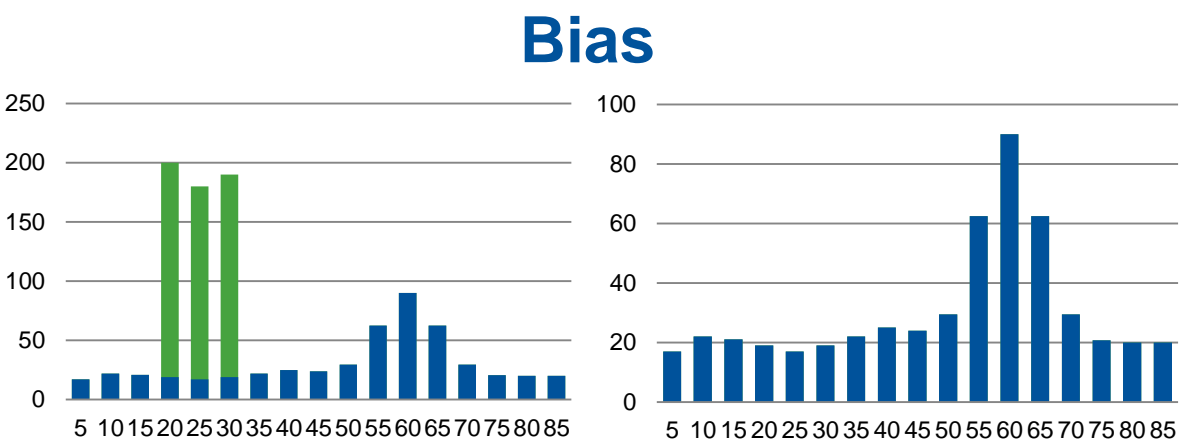
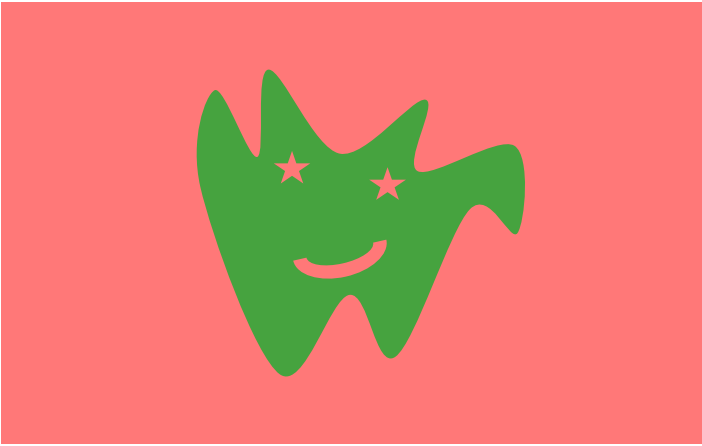


- Easy to be good on the training data
- Predictor is too "powerful"
- Rote learning

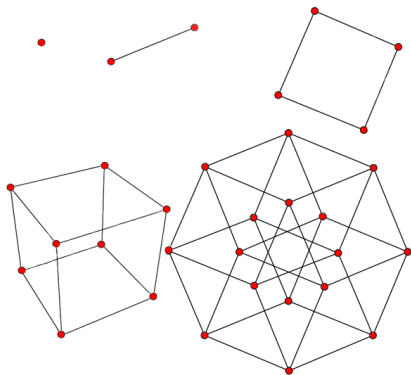


Machine Learning — Further Issues

Boundary
problem

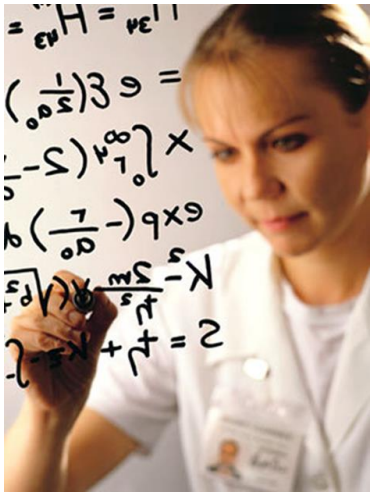


Poor data



No. of required data points
to just occupy the corners

	Dimensions
2	1
4	2
8	3
16	4
256	8
65.536	16
~ 4 billion	32
~ 4 billion x 4 billion	64



Talent shortage
and the danger
of automation

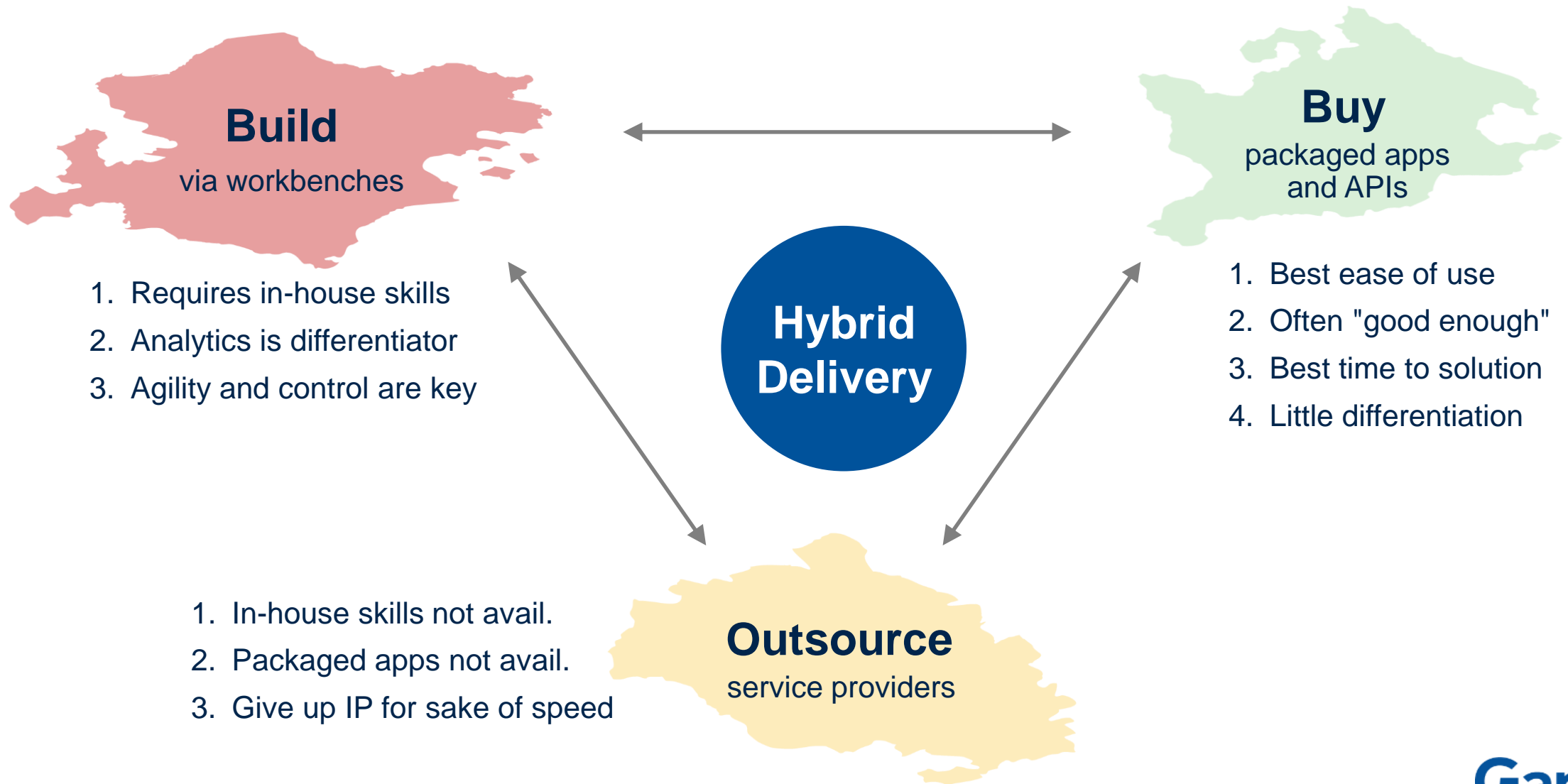
Key Issues

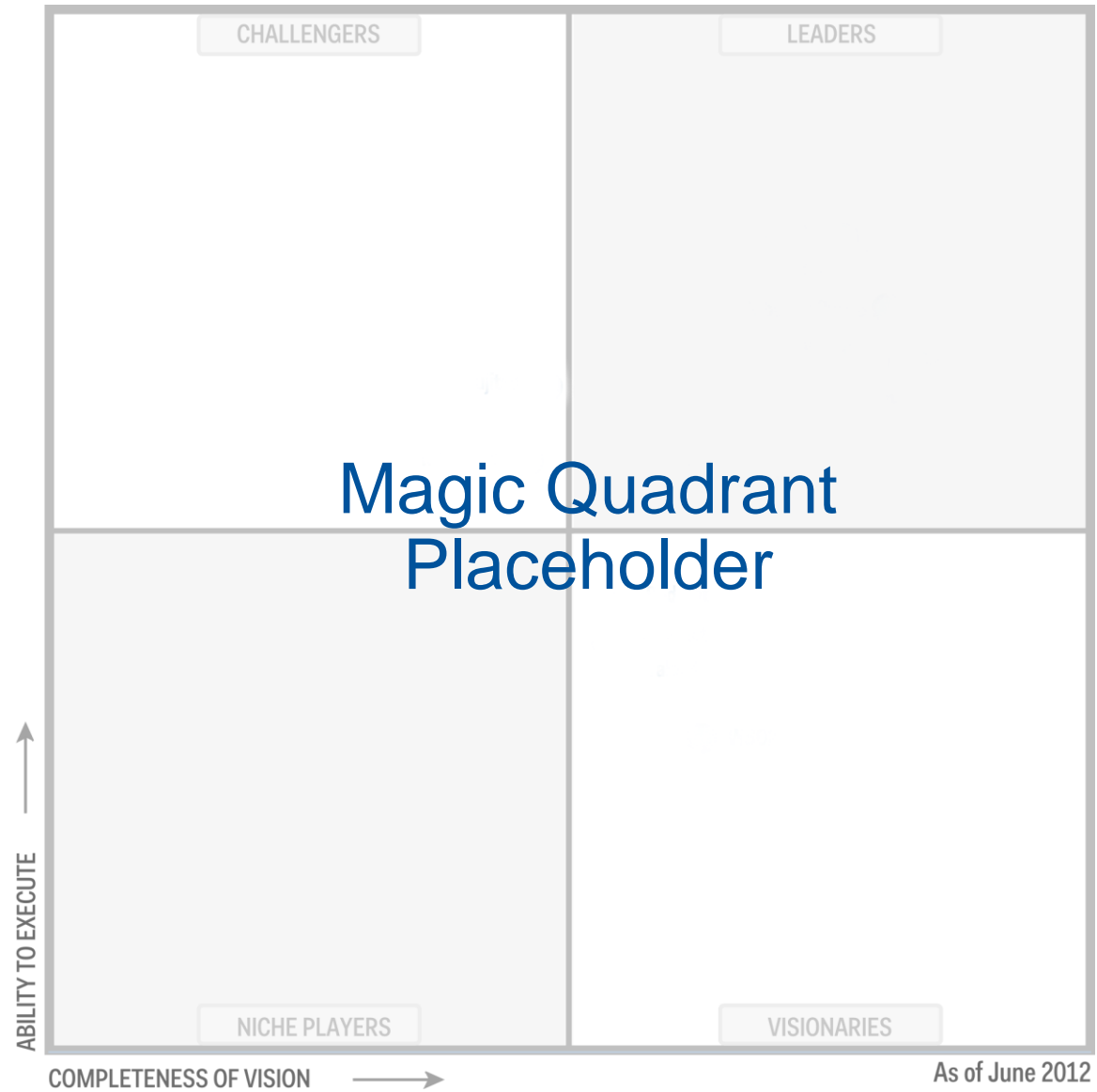
1. How do you navigate all semantics and deal with the shadow of AI?
2. How does it all work?
3. What can you do with data science and machine learning right now?

Build, Buy or Outsource?

	A) Generic Data Science Platforms	B) Packaged Analytics Applications	C1) Global IT Service Provider	C2) Global DS Service Provider	C3) Special DS Service Provider	C4) Crowdsourcing Platforms
Ease of Use						
Time to Solution						
Solution Quality						
Cost-Effectiveness						
Learning Experience						
Agility/Granularity of Control						
	Build	Buy		Outsource		

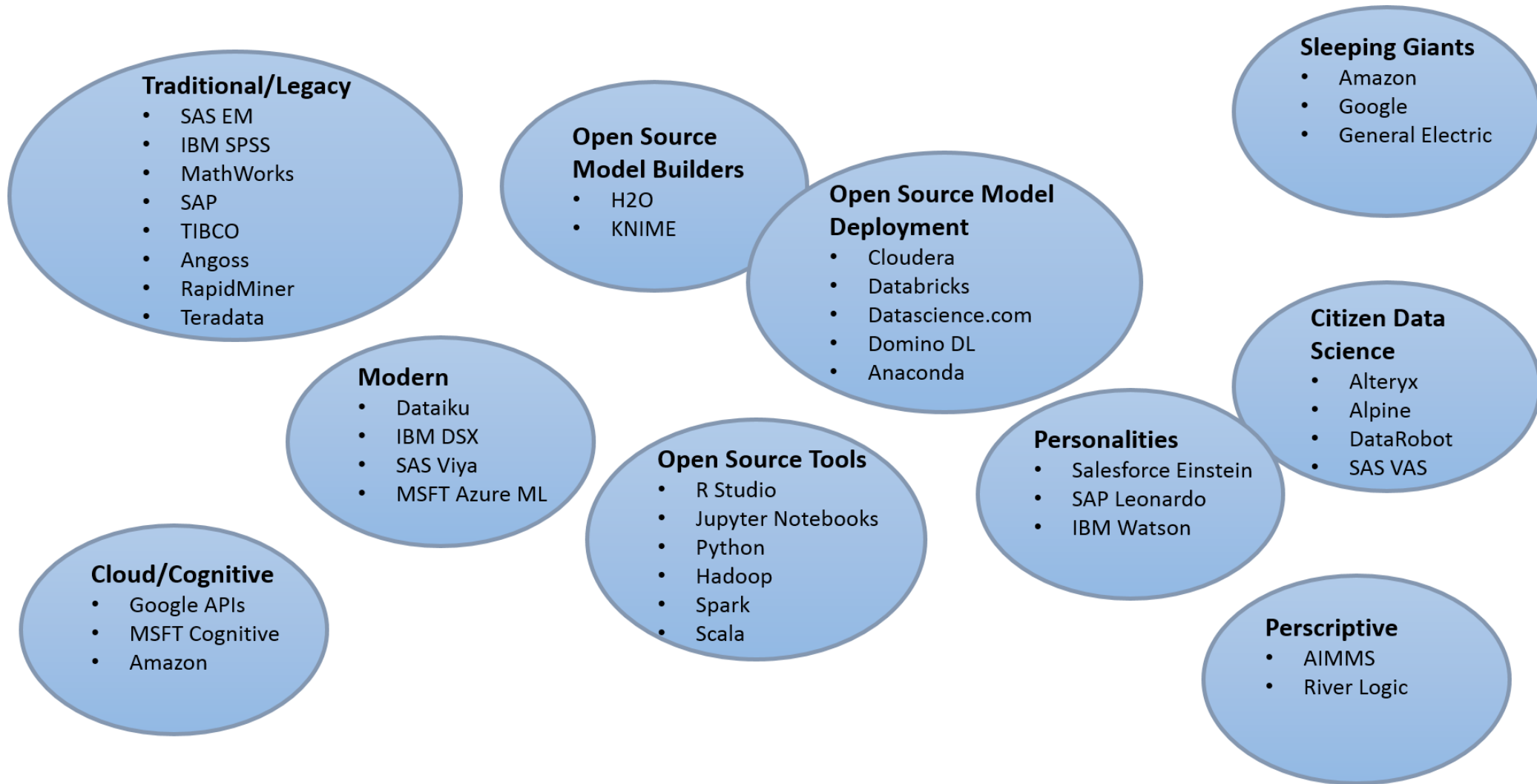
Delivery Models Blurring





(From "Magic Quadrant for Data Science and Machine Learning Platforms," XX month 20XX)

DS/ML Clusters — Apples to Oranges



Challenges and Pitfalls

Talent Gap

Scope Creep

Poor Data Quality

Fear and Misunderstanding

Culture and Territorialism

User Inertia and
the Analytical Status Quo



Upskilling

- Physicists
- Chemistrists
- Biologists
- Engineering Disciplines
- Social Scientists
- Computer Scientists
- Statistician
- Operations Researcher
- Mathematicians
- Industrial Engineers
- MBAs
- Astronomers
- Data Analysts
- Actuaries
- Risk Managers
- Control Engineers
- Financial Accountants
- Quality Specialists (Six Sigma)

- Cast a wide net
- Identify candidates and take an inventory of skills
- Allow time for training promising candidates
- Simultaneously educate analytics consumers

Math-Literate
Professionals



Quantitative
Professionals



Citizen
Data Scientists

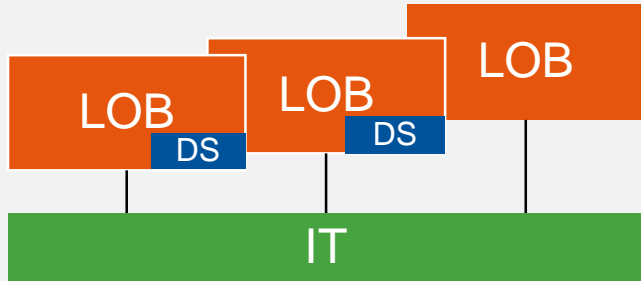


Data Scientists

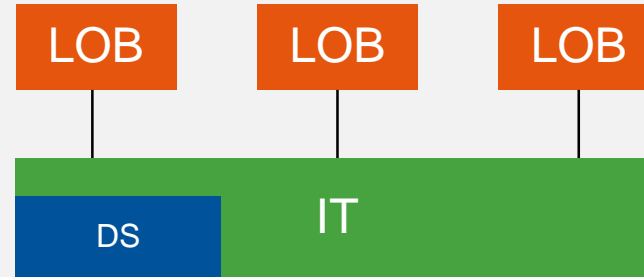


Where to Put the Data Scientists?

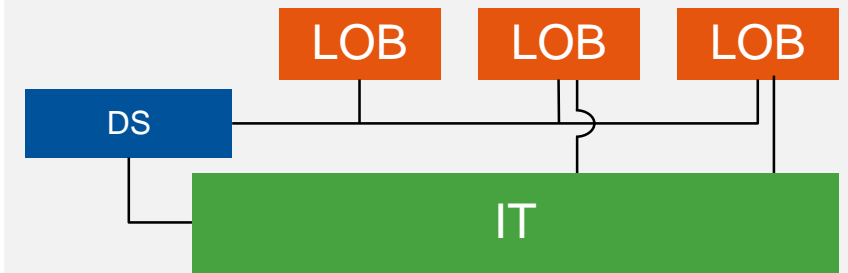
Data Scientists @ LOBs



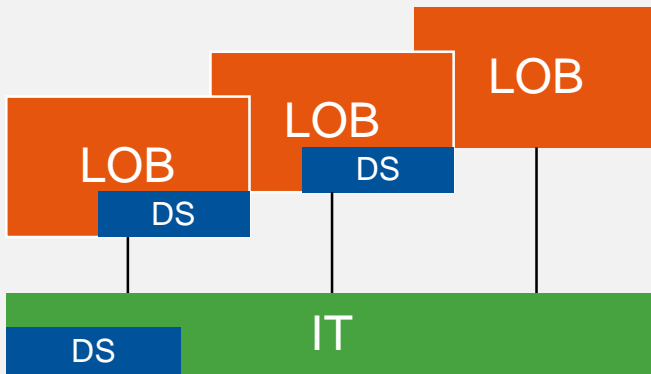
Data Scientists @ IT



Data Scientists as Separate BU

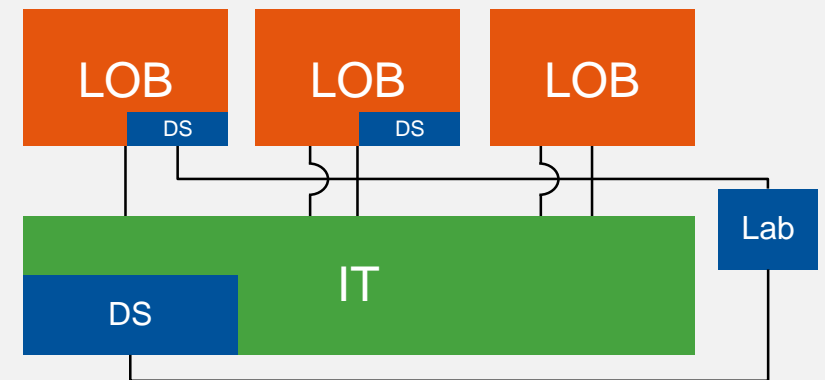


Scattered Experts



- Business Intimacy
- Knowledge Sharing
- Agility
- Cross-Functional View
- Proximity to Process and Data

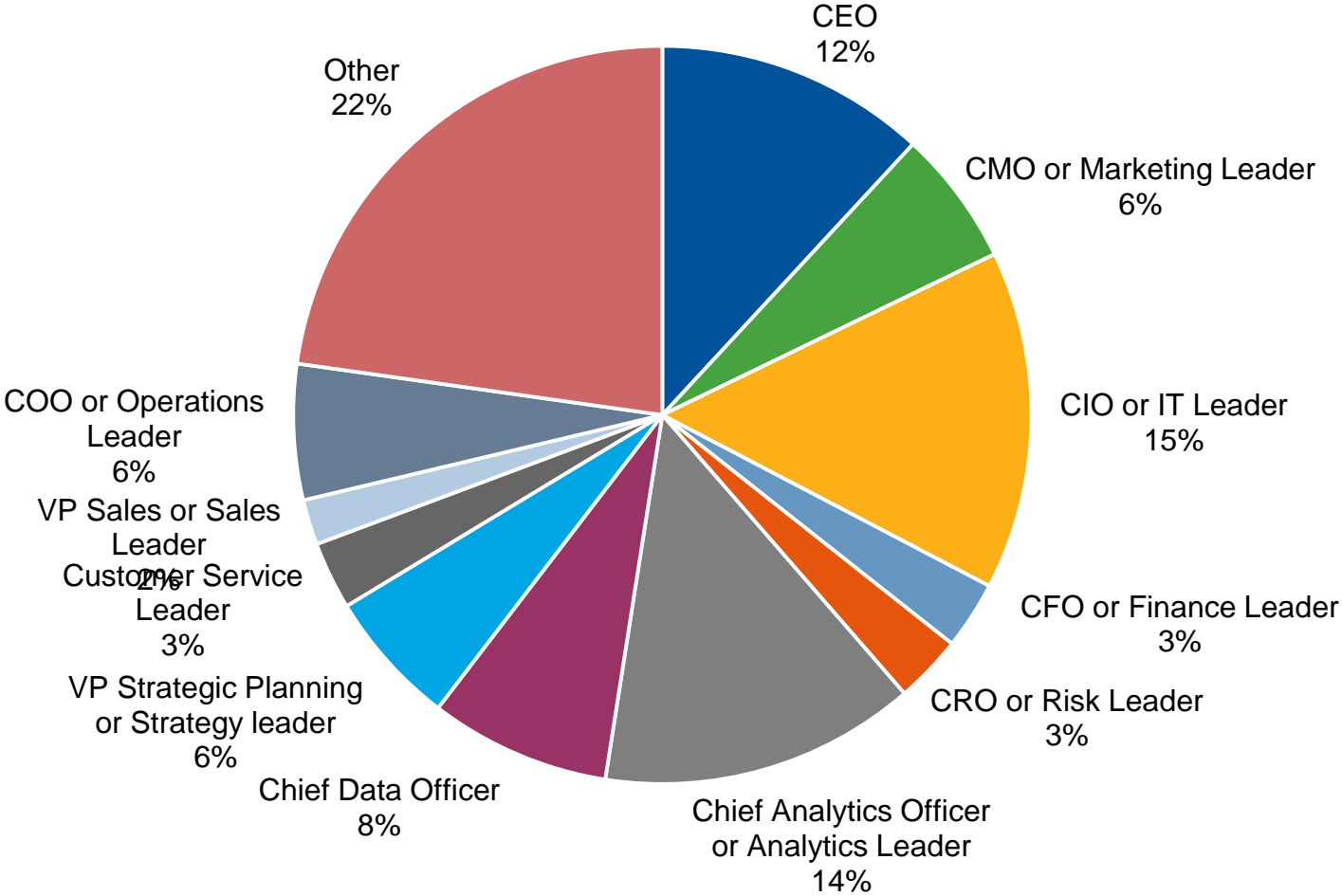
Data Science Lab/CoE



Source: ["Organizational Principles for Placing Data Science and Machine Learning Teams."](#) (G00325989)

DS = Data Science; LOB = Line of Business

Data Science Teams are Found Everywhere — Only 15% Report to IT



Recommendations

- ✓ Don't fear the shadow of AI.
- ✓ Cut through the hype, start focused and stay focused.
- ✓ Always look first at packaged applications.
- ✓ Be prepared for significant staffing and communication challenges.
- ✓ Take an inventory of data science skills and support upskilling initiatives.

AI Renaissance or Apocalypse?



Recommended Gartner Research

- ▶ [How to Start a Machine-Learning Initiative With Less Anxiety](#)
Svetlana Sicular (G00331893)
- ▶ [Leading Upskilling Initiatives in Data Science and Machine Learning](#)
Peter Krensky, Shubhangi Vashisth and Douglas Laney (G00334219)
- ▶ [Innovation Insight for Deep Learning](#)
Alexander Linden, Tom Austin and Svetlana Sicular (G00319191)
- ▶ [Machine Learning: FAQ From Clients](#)
Shubhangi Vashisth, Alexander Linden and Others (G00327948)
- ▶ [Hype Hurts: Steering Clear of Dangerous AI Myths](#)
Tom Austin, Alexander Linden and Mike Rollings (G00324274)