

O'REILLY®



# Getting Started with Artificial Intelligence

Second Edition

A Practical Guide to Building  
Enterprise Applications

Tom Markiewicz & Josh Zheng

REPORT



# Build Smart Build Secure



Take your journey to AI further and  
integrate AI into your next app today.

[ibm.biz/buildwithAI](https://ibm.biz/buildwithAI)



SECOND EDITION

---

# Getting Started with Artificial Intelligence

*A Practical Guide to Building  
Enterprise Applications*

*Tom Markiewicz and Josh Zheng*

Beijing • Boston • Farnham • Sebastopol • Tokyo

**O'REILLY®**

## Getting Started with Artificial Intelligence

by Tom Markiewicz and Josh Zheng

Copyright © 2020 O'Reilly Media, Inc. All rights reserved.

Printed in the United States of America.

Published by O'Reilly Media, Inc., 1005 Gravenstein Highway North, Sebastopol, CA 95472.

O'Reilly books may be purchased for educational, business, or sales promotional use. Online editions are also available for most titles (<http://oreilly.com>). For more information, contact our corporate/institutional sales department: 800-998-9938 or [corporate@oreilly.com](mailto:corporate@oreilly.com).

**Acquisitions Editor:** Rebecca Novack

**Development Editor:** Melissa Potter

**Production Editor:** Kristen Brown

**Copyeditor:** Penelope Perkins

**Interior Designer:** David Futato

**Cover Designer:** Randy Comer

**Illustrator:** Kate Dullea

December 2017: First Edition

September 2020: Second Edition

### Revision History for the Second Edition

2020-09-30: First Release

The O'Reilly logo is a registered trademark of O'Reilly Media, Inc. *Getting Started with Artificial Intelligence*, the cover image, and related trade dress are trademarks of O'Reilly Media, Inc.

The views expressed in this work are those of the authors, and do not represent the publisher's views. While the publisher and the authors have used good faith efforts to ensure that the information and instructions contained in this work are accurate, the publisher and the authors disclaim all responsibility for errors or omissions, including without limitation responsibility for damages resulting from the use of or reliance on this work. Use of the information and instructions contained in this work is at your own risk. If any code samples or other technology this work contains or describes is subject to open source licenses or the intellectual property rights of others, it is your responsibility to ensure that your use thereof complies with such licenses and/or rights.

This work is part of a collaboration between O'Reilly and IBM. See our [statement of editorial independence](#).

978-1-492-08341-2

[LSI]

---

# Table of Contents

<b>1. Introduction to Artificial Intelligence.....</b>	<b>1</b>
The Market for Artificial Intelligence	2
Avoiding an AI Winter	3
Artificial Intelligence, Defined?	4
Applications in the Enterprise	7
Next Steps	8
<b>2. Natural Language Processing.....</b>	<b>11</b>
Overview of NLP	12
The Components of NLP	13
Enterprise Applications of NLP	15
How to Use NLP	18
Challenges of NLP	20
Summary	21
<b>3. Chatbots.....</b>	<b>23</b>
What Is a Chatbot?	24
The Rise of Chatbots	24
How to Build a Chatbot	26
Challenges of Building a Successful Chatbot	29
Summary	30
<b>4. Computer Vision.....</b>	<b>31</b>
Capabilities of Computer Vision for the Enterprise	32
How to Use Computer Vision	36
Computer Vision on Mobile Devices	38
Best Practices	39

Use Cases	41
Existing Challenges in Computer Vision	43
Implementing a Computer Vision Solution	44
Summary	45
<b>5. AI Data Pipeline.....</b>	<b>47</b>
Preparing for a Data Pipeline	49
Sourcing Big Data	50
Storage: Apache Hadoop	50
Discovery: Apache Spark	52
Automation	53
Summary	57
<b>6. Hybrid Clouds.....</b>	<b>59</b>
Overview of Clouds	60
Using AI on Hybrid Clouds	63
Practical Solutions	65
The Future of AI on Hybrid Clouds	67
Summary	67
<b>7. Looking Forward.....</b>	<b>69</b>
What Makes Enterprises Unique?	70
Current Challenges, Trends, and Opportunities	71
Scalability	76
Social Implications	77
Summary	78

# Introduction to Artificial Intelligence

*In the future, AI will be diffused into every aspect of the economy.*

—Nils J. Nilsson, founding researcher, Artificial Intelligence & Computer Science, Stanford University

Beyond the buzzwords, media coverage, and hype, artificial intelligence techniques are becoming a fundamental business growth component across a wide range of industries. And while the various terms (algorithms, transfer learning, deep learning, neural networks, NLP, etc.) associated with AI are thrown around in meetings and product planning sessions, it's easy to be skeptical of the potential impact of these technologies.

Today's media represents AI in many ways, both good and bad—from the fear of machines taking over all human jobs and portrayals of evil AIs via Hollywood to the much-lauded potential of curing cancer and making our lives easier. Of course, the truth is somewhere in between.

While there are valid concerns about how the future of artificial intelligence will play out (and the social implications), the reality is that the technology is currently used in companies across all industries.

AI is used everywhere—IoT (Internet of Things) and home devices, commercial and industrial robots, autonomous vehicles, drones, digital assistants, and even wearables. And that's just the start. AI

will drive future user experiences that are immersive, continuous, ambient, and conversational. These conversational services (e.g., chatbots and virtual agents) are currently exploding, while AI will continue to improve these contextual experiences.

Despite several stumbles over the past 60–70 years of effort on developing artificial intelligence, the future is here. If your business is not incorporating some AI, you'll quickly find you're at a competitive disadvantage in today's rapidly evolving market.

Just what are these enterprises doing with AI? How are they ensuring an AI implementation is successful (and provides a positive return on investment)? These are only a few of the questions we'll address in this book. From natural language understanding to computer vision, this book will provide you a high-level introduction to the tools and techniques to better understand the AI landscape in the enterprise and initial steps on how to get started with AI in your own company.

It used to be that AI was quite expensive and a luxury reserved for specific industries. Now it's accessible for everyone in every industry, including you. We'll cover the modern enterprise-level AI techniques available that allow you to create models efficiently and implement them into your environment. While not meant as an in-depth technical guide, the book is intended as a starting point for your journey into learning more and building applications that implement AI.

## The Market for Artificial Intelligence

The market for artificial intelligence is already large and growing rapidly, with numerous research reports indicating a growing demand for tools that automate, predict, and quickly analyze. **Estimates from IDC** predict revenue from artificial intelligence will reach \$98.4 billion by 2023 with a compound annual growth rate (CAGR) of 28.5% over the forecast period, with nearly half of that going to software. Additionally, investment in AI and machine learning companies has increased dramatically—**AI startups raised \$26.6 billion in funding** in 2019. Clearly, the future of artificial intelligence appears healthy.



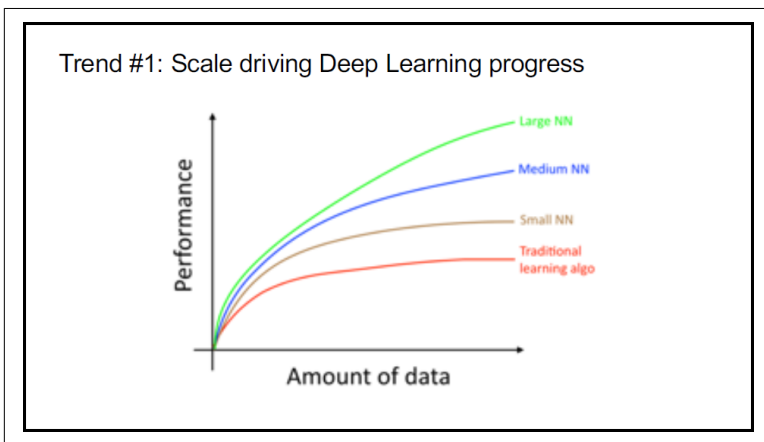
# Avoiding an AI Winter

Modern AI, as we know it, started in earnest in the 1950s. While it's not necessary to understand the detailed history of AI, it is helpful to understand one particular concept—the *AI winter*—as it shapes the current environment.

There were two primary eras of artificial intelligence research where high levels of excitement and enthusiasm for the technology never lived up to expectations, causing funding, interest, and continued development to dry up. The buildup of hype followed by disappointment is the **definition of an AI winter**.

So why are we now seeing a resurgence in AI interest? What's the difference today that's making AI so popular in the enterprise, and should we fear another AI winter? The short answer is likely no—we expect to avoid an AI winter this time around due primarily to much more (or big) data and the advent of better processing power and graphics processing units (GPUs). From the tiny supercomputers we all carry in our pockets to the ever-expanding role of IoT, we're generating more data now, at an ever-increasing rate. For example, **IDC estimates** that 180 zettabytes of data will be created globally in 2025, up from less than 10 zettabytes in 2015.

**Andrew Ng**, cofounder of Coursera and Stanford adjunct professor, often presents **Figure 1-1** in his courses on **machine learning** and **deep learning**.



*Figure 1-1. Deep learning performance (image courtesy of Andrew Ng, DeepLearning.ai course on Coursera)*

Conceptually, this chart illustrates how the performance of deep learning algorithms improves with an increasing amount of data—data that we now have in abundance, and that is growing at an exponential rate.

So why does it matter whether we understand the AI winter? Well, if companies are going to invest in AI, it's essential to avoid the hyperbole of the past and keep expectations based on reality. While the current situation is much more likely to justify the enthusiasm, that excitement still needs to be tempered with real-world results to avoid another AI winter. As (future) AI practitioners, that's something we can all agree would challenge your businesses.

## Artificial Intelligence, Defined?

The market for artificial intelligence is immense, but what are we genuinely discussing? While it sounds great to say you're going to implement AI in your business, what does that mean in practical terms? Artificial intelligence is quite a broad term and is, in reality, an umbrella over a few different concepts.

So, to get started and keep everyone on the same page, let's briefly discuss some of the terms associated with AI that are often confused or interchanged: artificial intelligence, machine learning, and deep learning.

### Artificial Intelligence

Over the years, there have been numerous attempts at precisely defining AI. While there's never really been an official, accepted definition, a few high-level concepts shape and define what we mean. The descriptions range from artificial intelligence merely aiming to augment human intelligence to thinking machines and true computer intelligence. For our purposes, we can think of AI as being any technique that allows computers to bring meaning to data in similar ways to a human.

While most AI is focused on specific problem domains (natural language processing or computer vision, for example), the idea of **artificial general intelligence**, or having a machine perform any task a human could (also commonly referred to as “strong AI”), is still more than 10 years out **according to Gartner**.

## Machine Learning

In 1959 Arthur L. Samuel, an IBM researcher and Stanford professor, is said to have stated that machine learning is the “field of study that gives computers the ability to learn without being explicitly programmed,” thereby **becoming the originator of the term**.

Essentially, machine learning is a subset of AI focused on having computers provide insights into problems without explicitly programming them. Most of the tools and techniques that today refer to AI are representative of machine learning. There are three main types of machine learning—supervised learning, unsupervised learning, and reinforcement learning.

By looking at many labeled data points and examples of historical problems, supervised learning algorithms can help solve similar problems under new circumstances. Supervised learning trains on large volumes of historical data and builds general rules to be applied to future problems. The better the training set data, the better the output. In supervised learning, the system learns from these human-labeled examples.

While supervised learning relies on labeled or structured data (think rows in a database), unsupervised learning trains on unlabeled or unstructured data (the text of a book). These algorithms explore the data and try to find structure. Here, widely used unsupervised learning algorithms are cluster analysis and market basket analysis. Naturally, this tends to be more difficult, as the data has no preexisting labels to assist the algorithms in understanding the data. While more challenging to process, as we’ll discuss later, unstructured data makes up the vast majority of data that enterprises need to process today.

Finally, there’s reinforcement learning as a machine learning technique. Reinforcement learning takes an approach similar to behavioral psychology. Instead of training a model with predefined training sets (i.e., where you know the prescribed answers in advance) as in supervised learning, **reinforcement learning rewards the algorithm** when it performs the correct action (behavior). Reinforcement learning resists providing too much training and allows the algorithm to optimize itself for performance-based rewards.

Reinforcement learning was initially conceived in 1951 by Marvin Minsky. Still, as was the case with many AI implementations, the

algorithms were held back by both the scale of data and computer processing power needed for effectiveness. Today, we see many more successful examples of reinforcement learning in the field, with Alphabet subsidiary DeepMind's AlphaGo one of the more prominent. Another notable application of reinforcement learning is the development of self-driving cars.

## Deep Learning

No book on developing AI applications in the enterprise would be complete without a discussion of deep learning. One way to understand deep learning is to think of it as a subset of AI that attempts to develop computer systems that learn using neural networks like those in the human brain. While machine learning is mostly about optimization, deep learning focuses on creating algorithms to simulate how the human brain's neurons work.

Deep learning algorithms are composed of an interconnected web of nodes called neurons and the edges that join them together. **Neural nets receive inputs, perform calculations, and then use this output to solve given problems.**

One of the ultimate goals of AI is to have computers think and learn like human beings. Using neural networks based on the human brain's decision-making process, deep learning is a set of computational techniques that moves us closest to that goal.

According to Stanford professor Chris Manning, around the year 2010, **deep learning techniques started to outperform other machine learning techniques.** Shortly after that, in 2012, Geoffrey Hinton at the University of Toronto led a team creating a neural network that learned to recognize images. Additionally, that year Andrew Ng led a team at Google that created a system to **recognize cats in YouTube videos** without explicitly training it on cats.

By building multiple abstraction layers, deep learning can also **solve complex language problems** where it's essential to discern meaning. Traditional AI follows a linear programming structure, limiting the ability for a computer to learn without human intervention. However, deep learning replaces this existing approach with self-learning and self-teaching algorithms enabling more advanced data interpretation. For example, traditional AI could read a zip code and mailing address off an envelope, but deep learning could also infer

from the envelope's color and date sent that the letter contains a holiday card.

Deep learning's closer approximation to how the human brain processes information via neural networks offers numerous benefits. First, the approach delivers a significantly improved ability to process large volumes of unstructured data, finding meaningful patterns and insights. Next, deep learning offers the ability to develop models across a wide variety of data, highlighted by advanced approaches toward text, images, audio, and video content. Finally, deep learning's massive parallelism allows for the use of multiple processing cores/computers to maximize overall performance as well as faster training with the advent of GPUs.

As with all advances in AI, deep learning benefits from much more training data than ever before, faster machines and GPUs, and improved algorithms. These neural networks are driving rapid advances in speech recognition, image recognition, and machine translation, with some systems performing as well or better than humans.

## Applications in the Enterprise

From finance to cybersecurity to manufacturing, there isn't an industry that will not be affected by AI. But before we can discuss and examine building applications in the enterprise with AI, we first need to define what we mean by enterprise applications. There are two common definitions of "the enterprise":

- A company of significant size and budget
- Any business-to-business commerce (i.e., not a consumer)

So Geico, Procter & Gamble, IBM, and Verizon would all be enterprise companies. And by this definition, any software designed and built internally would be considered enterprise software. On the other hand, a small company or startup could develop applications to be used by businesses (whether large or small), and this would still be considered enterprise software. However, a photo-sharing app for the average consumer would not be considered enterprise software.

This is probably obvious, but since we're discussing enterprise applications with AI in this book, it's important to be explicit about what

we mean when talking about the enterprise. For the rest of the book, the context of an enterprise will be that the end user is a business or business employee.

Does that mean if you are building the next great consumer photo-sharing application that this book will be of no use? Absolutely not! Many (if not most) of the concepts discussed throughout the book can be directly applied to consumer-facing applications as well. However, the use cases and discussion will be centered on the enterprise.

## Next Steps

This book is intended to provide a broad overview of artificial intelligence, the various technologies involved, relevant case studies and examples of implementation, and the current landscape. We'll also discuss the future of AI and where we see it moving in the enterprise from a current practitioner's viewpoint.

While the book focuses on a more technical audience, specifically application developers looking to apply AI in their business, the general topics and discussion will help anyone interested in how artificial intelligence will impact the enterprise.

Note that this book is not a substitute for a course or deep study on these topics, but we hope the material covered here jump-starts the process of your becoming a proficient practitioner of AI technologies.

While there are many options for implementing the various AI techniques we'll discuss in this book—coding in-house from scratch, outsourcing, the use of open source libraries, software-as-a-service APIs from leading vendors—our familiarity and deep experience revolves around IBM Watson, so we'll be using those code examples to illustrate the individual topics as necessary.

It's important to note that, while we believe IBM provides a great end-to-end solution for applying AI in the enterprise, it's not the only option. All the major cloud computing providers offer similar solutions. Also, there are numerous open source tools that we'll discuss later.

In the next few chapters, we'll discuss some of the more common uses for AI in the enterprise (natural language processing, chatbots,

and computer vision). Then we'll discuss the importance of a solid data pipeline, followed by a look forward at the challenges and trends for AI in the enterprise.





# Natural Language Processing

*Our intelligence is what makes us human, and AI is an extension of that quality.*

—Yann LeCun, professor, New York University

Humans have been creating the written word for thousands of years, and we've become pretty good at reading and interpreting the content quickly. Intention, tone, slang, and abbreviations—most native speakers of a language can process this context in both written and spoken word quite well. But machines are another story. As early as the 1950s, computer scientists began attempts at using software to process and analyze textual components, sentiment, parts of speech, and the various entities that make up a body of text. Until relatively recently, processing and analyzing language has been quite a challenge.

Ever since IBM's Watson won on the game show *Jeopardy!*, the promise of machines being able to understand language has slowly edged closer to reality. In today's world, where people live out their lives through social media, the opportunity to gain insights from the millions of words of text being produced every day has led to an arms race. New tools allow developers to easily create models that understand words used in the context of their industry. This leads to better business decisions and has resulted in a high-stakes competition in many industries to be the first to deliver.

A 2017 study from IBM reported that 90% of the world's data had been created in the past two years, and that 80% of that data was unstructured. Insights valuable to the enterprise are hidden in

this data, ranging from emails to customer support discussions to research reports. This information is incredibly useful if it can be found, interpreted, and utilized. When an enterprise can harness this massive amount of unstructured data and transform it into something meaningful, there are endless possibilities for improving business processes, reducing costs, and enhancing products and services.

Alternatively, those companies without the ability to handle their unstructured data suffer lost revenue, missed business opportunities, and increased costs, all likely without their knowledge.

Interpreting this unstructured data is quite difficult. In fact, processing human-generated (not machine) words (or natural language) is considered an *AI-hard* or *AI-complete* problem. In other words, it's a challenge that brings the **full effort of AI** to bear on the problem and isn't easily solved by a single algorithm designed for a particular purpose.

In this chapter, we'll give an overview of NLP, discuss some industry examples and use cases, and look at some strategies for implementing NLP in enterprise applications.

## Overview of NLP

Natural language processing is essentially the ability to take a body of text and extract meaning from it using a computer. While computational language is very structured (think XML or JSON) and easily understood by a machine, written words by humans are quite messy and unstructured—meaning when you write about a house, friend, pet, or a phone in a paragraph, there's no explicit reference that labels each of them as such.

For example, take this simple sentence:

I drove my friend Mary to the park in my Tesla while listening to music on my iPhone.

This is an easily understandable sentence for a human reader and paints a clear picture of what's happening. But for a computer, not so much. For a machine, the sentence would need to be broken down into its structured parts. Instead of an entire sentence, the computer would need to see both the individual parts or entities along with the relations between these entities.

Humans understand that Mary is a friend and that a Tesla is likely a car. Since we have the context of bringing our friend along with us, we intuitively rule out that we're driving something else, like a bicycle. Additionally, after many years of popularity and cultural references, we all know that an iPhone is a smartphone.

A computer understands none of the above without assistance. Now let's look at how that sentence could be written as structured data from the outset. If developers had made time in advance to structure the data in our sentence, in XML you'd see the following entities:

```
<friend>Mary</friend>  
<car>Tesla</car>  
<phone>iPhone</phone>
```

But obviously, this can't happen on the fly without assistance. As mentioned previously, we have significantly more unstructured data than structured. And unless time is taken to apply the correct structure to the text in advance, we have a massive problem that needs solving. This is where NLP enters the picture.

Natural language processing is needed when you wish to mine unstructured data and extract meaningful insight from text. General applications of NLP attempt to identify common entities from a body of text; but when you start working with domain-specific content, a custom model needs training.

## The Components of NLP

To understand NLP, we first need to understand the components of its model. Specifically, natural language processing lets you analyze and extract key metadata from text, including entities, relations, concepts, sentiment, and emotion.

Let's briefly discuss each of these aspects that can be extracted from a body of text.

### Entities

Likely the most common use case for natural language processing, entities are the people, places, organizations, and things in your text. In our initial example sentence, we identified several entities in the text—friend, car, and phone.

## Relations

How are entities related? Natural language processing can identify whether there is a relationship between multiple entities and tell the type of relation between them. For example, a “createdBy” relation might connect the entities “iPhone” and “Apple.”

## Concepts

One of NLP’s more magical aspects is extracting general concepts from the body of text that may not explicitly appear in the corpus. This is a potent tool. For example, an analysis of an article about Tesla may return the concepts “electric cars” or “Elon Musk,” even if those terms are not explicitly mentioned in the text.

## Keywords

NLP can identify the important and relevant keywords in your content. This allows you to create a base of words from the corpus that are important to the business value you’re trying to drive.

## Semantic Roles

Semantic roles are the subjects, actions, and the objects they act upon in the text. Take the sentence, “Intel bought a company.” In this sentence, the subject is “Intel,” the action is “bought,” and the object is “company.” NLP can parse sentences into these semantic roles for various business uses—for example, determining which companies were acquired last week or receiving notifications any time a particular company launches a product.

## Categories

Categories describe what a piece of content is about at a high level. NLP can analyze text and then place it into a hierarchical taxonomy, providing categories to use in applications. Depending on the content, categories could be sports, finance, travel, computing, and so on. Possible applications include placing relevant ads alongside user-generated content on a website or displaying all the articles talking about a particular subject.

## Emotion

Whether you're trying to understand the emotion conveyed by a post on social media or analyze incoming customer support tickets, detecting emotions in text is extremely valuable. Is the content conveying anger, disgust, fear, joy, or sadness? Emotion detection in NLP will assist in solving this problem.

## Sentiment

Similarly, what is the general sentiment in the content? Is it positive, neutral, or negative? NLP can provide a score as to the level of positive or negative sentiment of the text. Again, this proves to be extremely valuable in the context of customer support. This enables continued, automatic understanding of sentiment related to your product.

Now that we've covered what constitutes natural language processing, let's look at some examples to illustrate how NLP is currently used across various industries.

# Enterprise Applications of NLP

While there are numerous examples of natural language processing being used in enterprise applications, the following are some of the best representations of NLP's power.

## Social Media Analysis

One of the most common enterprise applications of natural language processing is social media monitoring, analytics, and analysis. **Over 500 million tweets are sent per day.** How can we extract valuable insights from them? What are the relevant trending topics and hashtags for a business? Natural language processing can deliver this information and more by analyzing social media. Not only can sentiment and mentions be mined across all this user-generated social content, but specific conversations can also be found to help companies better interact with customers.

Additionally, when an incident occurs in real time, applying NLP to monitor social media provides a distinct advantage to help businesses react immediately with the appropriate understanding of the issue at hand.

## Customer Support

A recent study has shown that **companies lose more than \$62 billion annually due to poor customer service, a 51% increase since 2013**. Therefore, there's a need for methods to improve customer support.

Companies are using natural language processing in a variety of ways in customer support. For each incoming support ticket, the content can be analyzed to obtain its sentiment, relevant keywords, and categorization. This process can be used to route the support ticket faster to the correct representative and, in some cases, to automatically respond to the request (this can then be extended with chatbots, as we'll see in the next chapter).

Natural language processing can also assist in making sure support representatives are both consistent and nonaggressive (or any other trait the company is looking to minimize) in their language. When preparing a reply to a support question, an application that incorporates NLP can provide a suggested vocabulary to assist this process.

These customer support approaches can make the overall system much faster, more efficient, and easier to maintain, and subsequently reduce costs over a traditional ticketing system.

## Business Intelligence

According to Gartner, **the business intelligence (BI) software market reached \$24.8 billion in 2019**. Unfortunately, one of the common problems associated with BI is the reliance on running complex queries to access the mostly structured data. This presents two significant problems. First, how does a company access the bigger set of unstructured data, and second, how can this data be queried on a more ad hoc basis without the need for developers to write complex queries?

The inability to use unstructured data, both internal and external, for business decision making is a critical problem. Natural language processing allows all users, especially non-technical experts, to ask questions of the data instead of needing to write a complex query of the database. This allows the business users to ask questions of the data without having to request developer resources to make it happen. This democratizes BI within the enterprise and frees up crucial development time for developers in other areas. Additionally, this

significantly improves overall productivity in the organization and allows for a potential reduction in staff for a particular project or application implementation.

## **Content Marketing and Recommendation**

As it becomes harder to reach customers with advertising, companies now look to content marketing to produce unique stories that drive traffic and increase brand awareness. Not only do they look for new content to create, but companies also want better ways to recommend more relevant content to their readers. Everyone is familiar with recommended articles that are merely clickbait with little value or applicability to your interests.

Also, as more people use ad blockers, the traditional method of monetizing content is rapidly waning. In response, this leads businesses to engage in more compelling ways, primarily through better content and unique storytelling.

Natural language processing enables companies publishing content to analyze all the articles, blog posts, and customer comments and reviews to both understand what to write about as well as produce more interesting and relevant topics to readers. Additionally, massive amounts of trend data can also be gleaned from this newly processed content, providing additional insights.

## **Additional Topics**

We have discussed just a few industry examples, but there are many more. For example, natural language processing is used in brand management. Customers are talking about brands every day across multiple channels. How does a company monitor what's said about the brand and understand the content and sentiment? Relatedly, market intelligence is another area often improved through natural language processing.

While more specific to a particular domain or industry, there are also other examples that illustrate the power of natural language processing for improving business results. An example of this is the legal industry. NLP is being used by numerous companies to review cases and other legal documents to alleviate the need for expensive lawyers and paralegals. Not only do legal professionals save time by not having to read every word personally, but the firms also reduce error rates by having a machine process many thousands of words

quickly as opposed to a human reader who can soon tire. Interestingly, while one may think this leads to a reduction in jobs (particularly for the relatively lower-cost paralegals and legal assistants), it has improved their efficiency instead, allowing them to spend their time doing more and higher-rate billable work.

### Practical Tip

Now that you've read some examples of natural language processing used in the enterprise, take a minute to think about your industry. First, in what areas have you seen the approach applied in your field? Second, brainstorm some examples of how NLP can be used in your company. Finally, start to think of what you may need to implement these as solutions. We'll discuss options as the book proceeds, but challenge yourself to think of what's required to improve your applications with NLP.

## How to Use NLP

Now that we've provided an overview of natural language processing and given some industry examples, let's look at some of the strategies for implementing NLP in an application.

There are many solutions for natural language processing. Starting with open source software projects, a few of the more popular include Apache NLP, Stanford CoreNLP, NLTK for Python, and SyntaxNet.

While these are some of the more popular options, there's a collection of open source libraries for natural language processing in almost every programming language. For example, if you use Ruby, you can find a collection of small libraries at <http://rubynlp.org>. The same goes for PHP: <http://php-nlp-tools.com>. At this point, there's typically no need to reinvent the wheel, or in this case, the algorithm!

Nevertheless, while there are many options to implement natural language processing using open source as a starting point, from a cost-benefit perspective, it can often make sense for enterprise applications to utilize one of the numerous third-party services.

Currently, several companies provide APIs offered as software as a service (SaaS). From IBM Watson Natural Language Understanding



(NLU) to Azure Text Analytics to Amazon Comprehend, utilizing a hosted service API can reduce developer time and save these vital resources for other application development aspects.

When evaluating whether to build in-house, outsource, or use hosted APIs, ask yourself the following important question: how much of a core component to your business is artificial intelligence? Your answer can then drive the technical level of expertise required for your enterprise application. For example, if you're an ecommerce company attempting to add intelligence to your customer support system, it would be more appropriate to start with hosted APIs, as better customer support improves the business but isn't your core functionality.

Alternatively, companies like Amazon and Netflix rely on recommendation engines, which assist in creating a personalized experience for users, as core functions of their business. According to McKinsey, **these recommendation algorithms produce 35% of Amazon purchases and 75% of Netflix viewings**. In this case, they would employ machine learning engineers and data scientists to improve this part of the application continually.

### Practical Tip

When comparing NLP tools, take care to examine what the service includes. Most third-party providers bundle together several algorithms to create their product. Either plan to mix and match to meet your needs or carefully examine what the specific natural language processing API offers and whether it meets your application's needs.

## Training Models

If you develop natural language processing from scratch in your enterprise, you'll create custom models by default. But when you're using third-party solutions or open source options, the out-of-the-box solution will cover only the majority of cases, and it will be decidedly non-domain-specific. If you want to improve the accuracy and reliability of your output, you'll want to create and train a custom model. This is especially true if you're using a third-party service.

While there are various ways to accomplish training a model, the details are beyond the scope of this book, as they vary depending on the particular solution.

Using IBM Watson NLU as an example, you can train a custom model using the Watson Knowledge Studio (WKS). WKS is a web-based tool that enables domain experts to train a custom natural language processing model without programming. Both developers and non-technical end users can upload relevant documents and then annotate them for their domain-specific entities and relations. They can then use this data to train a custom model via machine learning and publish it to the Watson NLU APIs for use in their applications.

## Challenges of NLP

Despite being a robust technology at this point, natural language processing isn't always a perfect solution. While we've previously discussed NLP's numerous benefits, two major areas still prove to be a challenge for its implementation in enterprise applications.

First, natural language processing works best with massive data sets: the more data, the better the accuracy. While the data set's necessary size depends on the actual application, more data is better in general.

Second, natural language processing isn't a magic bullet. After you've been exploring and working on NLP some, it's easy to think you'll obtain easy answers to questions. When you're testing out NLP, the results tend to be very accurate, as the tendency is to input relatively straightforward bodies of text for testing. Unfortunately, human languages have many nuances. Think of all the phrases and words that are open to interpretation. Concepts like sarcasm are still quite hard to understand via natural language processing. Slang, jargon, and humor are hard to process as well. There's a tremendous amount of ambiguity in language that is only understood from the context. Additionally, handling spelling mistakes and errors in grammar is especially tricky.

What's the best way to handle these challenges then? Until the technology catches up and increases accuracy in these cases, the best approach is to know they exist and filter/review the content going through natural language processing as much as possible. While this

isn't an optimal solution in and of itself, paying attention to your preprocessed content beforehand and filtering any questionable content in advance is the best option.

### **Practical Tip**

Take a minute and visit your Twitter, Facebook, or LinkedIn feeds. Read through the posts and imagine being able to programmatically read and understand every piece of text almost instantaneously. What would you do with that insight? How could you incorporate this new knowledge into your enterprise application?

## **Summary**

Natural language processing is a powerful tool used in a wide range of enterprise applications. Since text appears almost everywhere, NLP provides an essential building block for all enterprise applications utilizing artificial intelligence.

In this vein, natural language processing also forms the backbone for creating conversational applications, more commonly known as chatbots. In the next chapter, we'll discuss these in more detail.



# Chatbots

*AI is a tool. The choice about how it gets deployed is ours.*

—Oren Etzioni, professor, University of Washington

Dr. Ashok Goel teaches a class called Knowledge-Based Artificial Intelligence every semester at Georgia Tech. It's a massive class of around 300 registered students. Understandably, being a teaching assistant (TA) for this course can be quite the challenge as 300 students post roughly 10,000 messages in the online forums. But during the spring semester of 2017, one particular TA was especially good at her job. She answered the students' questions with excellent efficiency. She would even make herself available during late hours the night before a deadline. Naturally, the students loved her and gave her rave reviews.

You can probably guess how the story ends. At the end of the semester, most students were surprised to find out that the TA, named Jill Watson, was actually a chatbot. Dr. Goel and his team built Jill by tracking down all the questions that had ever been asked on the course's online forum (about 40,000 postings in all). They then trained Jill to answer these questions using IBM Watson. Jill wasn't very good at first, but she learned from her mistakes and was eventually giving answers with 97% certainty.

It's stories like this that sparked the interest of developers and entrepreneurs and fueled interest in chatbots. In this chapter, we will explore a few aspects of this topic. We'll talk about what a chatbot is, why now is an excellent time to build one, and a few important considerations for creating a successful chatbot.

# What Is a Chatbot?

So what exactly is a chatbot? A chatbot is a way to expose a business's service or data via a natural language interface. It's important to understand that as an interface, the chatbot is only as good as the underlying service or data. So if you're thinking about building a chatbot, first make sure your services and data are solid. How do you know if they're in good shape? Imagine the service you're providing has a traditional web or mobile interface. Would that still be useful? If the answer is "no," your service is not ready for a chatbot interface either.

If built well, chatbots can help your business cut costs and establish additional revenue streams. A virtual customer support agent can reduce head count and exponentially scale your real-time customer support capabilities. A conversational commerce chatbot gives your business a whole new channel on which you can engage with your customers via messaging platforms.

## The Rise of Chatbots

Chatbots have been around for a long time. Twenty years ago, they would only look for a couple of words from a very highly restricted set of commands: "Is this correct? Type yes or no." But as you can see from the Georgia Tech story, the chatbots of today look entirely different.

There are a few reasons why this is happening now, especially in the enterprise:

- The availability of NLP capabilities discussed in the prior chapter, and particularly those in the cloud
- The proliferation of popular messaging platforms such as Slack and Facebook Messenger
- The push for natural language interfaces

The next several sections will elaborate on each.

## Natural Language Processing in the Cloud

The availability of natural language processing capabilities in the cloud has been the most potent force behind the rise of chatbots. NLP, specifically text classifiers and entity extractors, powers a few of the core functionalities inside a chatbot. If you've read the previous chapter or have experience in machine learning, you'll know that these NLP techniques are not new. The problem has been the difficulty in utilizing them for people outside the research community. Open source solutions have made these methods more accessible, but the arrival of cloud APIs, with a superior user experience, has enabled many more enterprises to use this technology.

## Proliferation of Messaging Platforms

Messaging apps have come to dominate our mobile app usage. **Recent data** shows they have surpassed social networks in monthly active users. As more users spend time on messaging apps, companies are looking at ways to reach users through these channels. It turns out there is a large amount of contextual data buried in these messages. We make dinner plans, inquire about stores, and look to purchase goods. Companies are now looking to help users by embedding chatbots inside message channels to answer questions or assist with various tasks.

## Natural Language Interface

In most human-machine interactions, users translate their intentions into a series of keystrokes and button clicks. The machine then responds via pixels on a screen. Wouldn't it be nice to talk to computers the way we talk to each other? This desire for natural language interfaces has always been around. In the early days of search engines, people liked Ask Jeeves because it allowed its users to search the web using natural language. Now, the proliferation of devices such as the Amazon Echo has drawn developers toward the idea of a voice-controlled home. After all, home appliances are notorious for their clunky user interfaces, and to replace them with smart agents that we could talk to seems like a much better user experience.

# How to Build a Chatbot

A chatbot has a frontend and a backend. The frontend is the messaging channel where the chatbot interacts with the user. The backend is the application logic, the persistence stores, and the supporting services.

## The Messaging Channel

There are many messaging channels available. You can leverage an existing one, such as Slack or Facebook Messenger. You can also build your own messaging layer such as a custom website or mobile app. Choosing the right channel depends on how you plan to engage your users. If you're a bank with a popular mobile app, you should expose your chatbot there. If you're a small business with an active Facebook page, integrating your chatbot with Facebook Messenger is a good idea.

## The Backend

Let's first discuss and expand upon one NLP technology in particular—text classifiers. It was not until recently that text classifiers became easy to use and available on the cloud. They're an important part of chatbots. If you were to approach building a customer support chatbot from scratch, the task would probably seem overwhelming. After all, even in a narrow domain such as customer support, there can still be an intractable number of different requests. How can you have an answer to every single one of them? The key insight is that though customer requests can be expressed in an almost infinite number of ways, the solution space is magnitudes smaller. Here's a more specific example. When talking to a customer service agent, a customer might say any of the following:

- “How come I can't log into my account anymore?”
- “I forgot my password.”
- “It says my password is incorrect.”
- “I'm locked out of my account.”

Luckily, all these requests have the same solution, which is to reset the customer password. So the core functionality of a chatbot is to



map all possible user inputs into a much smaller set of responses. This type of pattern recognition is what text classifiers do best.

Dr. Goel reached the same conclusion when building Jill Watson: “One of the secrets of online classes is that the number of questions increases if you have more students, but the number of different questions doesn’t really go up. Students tend to ask the same questions over and over again.”

### **Using a framework versus building your own**

The first significant decision in building a chatbot backend is to decide whether to leverage an existing framework. There are pros and cons to using one, and if chosen correctly, a framework can provide a large part of the solution with little effort. But this also means giving up control to the framework itself. It’s straightforward to build a chatbot within it, but impossible to customize or integrate the chatbot in a way outside the design of the framework. There are many chatbot frameworks available, including Google Cloud Dialogflow, Wit.ai, Microsoft Bot Framework, Amazon Lex, and IBM Watson Assistant.

### **Anatomy of a chatbot backend**

If you decide to use a framework, your chatbot’s backend will most likely consist of three main parts: intents, entities, and dialog. These can then be integrated with one or more messaging channels. Extra features such as sentiment analysis, human intervention, or personality can be added as well.

**Intent.** You usually start building a chatbot with intents. An intent is the purpose of a user’s input. This can be a question about your business hours or a complaint about the registration process. The response of your chatbot to this intent is entirely up to you. It might be a paragraph that answers the question or an action such as starting the password reset process. Another way to think about intents is that they’re the verbs for your chatbots to act on. They dictate what your chatbots will do next.

To train your chatbot to recognize an intent, first determine the action you’d like to map to this intent—for example, provide information on business hours. Then supply the framework with examples of user inputs that would require this action. The business hours scenario would include example inputs such as the following:

“What time are you open?” “Are you open on weekends?” “I can’t find your hours of operation.” Usually, a minimum of five inputs is needed, but more is better. Remember, it might be tempting to make up these examples, but it is always better to draw from past data. The more similar these examples are to real-world user requests, the better your chatbot will perform. Typically, an intent (really a text classifier) is built from these examples, so the service will recognize similar inputs in the future, even if they’re not exact matches.

**Entities.** If intents are the verbs for a chatbot to act on, then entities are the nouns. They’re the keywords in a user’s input. If a particular word differentiates one input from another, it probably should be an entity.

For example, if a user wants to find the business hours of a bank’s branch, the bank’s location would be one of the entities. Specifically, if a user asks, “What are the hours for the Austin branch?” *provide business hours* would be the intent, and *Austin* would be the entity.

**Dialog.** Dialog is the conversation flow. It’s usually represented as a directed graph where each node represents one exchange in the conversation. Together, it is the combined structure of all your possible conversations. Since this can become quite complex, most chatbot platforms provide a UI to help you visualize the process.

**Context variable.** A context variable contains the information shared between the framework and your application. It’s the way to exchange information between your business logic and the framework. For example, Watson Assistant provides a context object that lets you store any key/value pair as context variables.

**Human in the loop.** Embedding a human into your chatbot has two significant benefits.

First, humans are an excellent way to bootstrap your chatbot before it gathers enough data to operate autonomously. If the chatbot is not confident about its response, it can ask the human for approval or edits before sending it to the end user. Another setup is to have the chatbot provide multiple responses, and have the human choose the most appropriate one.

The second benefit is that humans are the best way to prevent your chatbot from utterly failing your users. There are a few ways to detect if the conversation needs to be routed to a human:

- The chatbot doesn't understand the user's input—this usually means the user input doesn't match any of your established intents.
- The conversation is taking too long, or a circular pattern is detected.
- Negative sentiment is caught in the user's input.
- The user directly asks to talk to a real person.

While seemingly counterintuitive when discussing AI, adding a human to your chatbot process ensures more seamless interactions with your users.

## Challenges of Building a Successful Chatbot

There are many things to get right for a chatbot to be useful. One of the most important is to define the project scope correctly. It needs to be broad enough for the chatbot to be helpful, yet narrow enough so that you're not wasting time building artificial general intelligence. Specifically, this means capturing as many user requests as possible, yet still reconciling the nuanced differences between each one.

This is not an easy problem. For example, one travel agency tried to deploy a vacation planning chatbot. A critical component was a vocabulary base large enough to recognize all the destinations and their everyday variations. It turns out there were over ten ways people could refer to the Cayman Islands, even assuming all spellings were correct. It took the company months to build a list that could confidently capture all the variations for this one destination.

### Practical Tip

It takes at least six months to create a chatbot that's useful, so make sure to give yourself or your development team enough of a runway.

Unfortunately, there's not a one-size-fits-all solution. Your best option is to put the chatbot in front of real customers and iterate through user feedback. If you're a product manager or developer, this should sound familiar. Keep in mind that chatbots are still a nascent field in the enterprise, so you'll likely be a trailblazer. Nothing but trial and error will help you discover the problems and edge cases specific to your domain.

## Summary

Remember, it's easy to get started with chatbots, but it takes patience and hard work to create a truly successful one. Don't try to develop general intelligence! Define the correct scope for your chatbot and keep in mind that context is vital, as the medium has much less bandwidth for communication. We'll leave you with this analogy from one of our colleagues. Building a chatbot is much like training a new hire—they probably only know a little on their first day. Still, through coaching and supervision, they'll eventually become a productive employee. It's essential for your users to see your chatbot as a human, but it's equally crucial for you to do the same. In the next chapter, we'll switch gears and talk about computer vision.

# Computer Vision

*Some people call this artificial intelligence, but the reality is this technology will enhance us. So instead of artificial intelligence, I think we'll augment our intelligence.*

—Ginni Rometty, Executive Chairman of IBM

Take a moment and look up from this book. Examine the room around you and take a quick inventory of what you see. Perhaps a desk, some chairs, bookshelves, and maybe even your laptop. Identifying these items is an effortless process for a human, even a young child.

Speaking of children, it's quite easy to teach them the difference between multiple objects. Over time, parents show them items or pictures and then repeat the name or description. Show them a picture of an apple, and then repeat the word *apple*. In the kitchen, hand them an apple, and then repeat the word *apple*. Eventually, through much repetition, the child learns what an apple is along with its many color and shape variations—red, green, yellow. Over time, we provide information as to what is a correct example and what isn't. But how does this translate to machines? How can we train computers to recognize patterns visually, as the human brain does?

Training computer vision models is done in much the same way as teaching children about objects. Instead of a person being shown physical items and having them identified, however, the computer vision algorithms are provided many examples of images that have been tagged with their contents. In addition to these positive

examples, negative examples are also added to the training. For example, if we're training for images of cars, we may also include negative examples of airplanes, trucks, and even boats.

Giving computers the ability to see opens up a world of possibilities for many industries—from recognizing human faces for security to automating processes that would take a human days, if not weeks. Let's take one industry as a quick example—insurance. Using computer vision to quickly, accurately, and objectively receive an automated analysis of images—for instance, of a fender bender or a weather-damaged roof—allows insurance companies to deliver better service to their clients.

Machine learning is pattern recognition through learned examples. Nothing exemplifies this more than computer vision. Just as NLP provides a core functionality of AI in the enterprise, computer vision extends a machine's ability to recognize and process images.

## Capabilities of Computer Vision for the Enterprise

So just what can computer vision do for our applications? Before diving into the details of computer vision and how to use it in the enterprise, let's examine some of the specific capabilities computer vision can bring to your applications.

### Image Classification and Tagging

Computer vision's most core functionality, general image tagging and classification, allows users to understand an image's content. While you'll often see image *classification* and *tagging* used interchangeably, it's best to consider classification as assigning an image to one or more categories. In contrast, tagging is an assignment of a single word (or multiple words) describing the image. When an image is processed, various keyword tags or classes are returned, describing the image with varying confidence levels. Based on an application's needs, these can be used to identify the image contents. For example, you may need to find images with the contents of "male playing soccer outside" or organize images into visual themes such as cars, sports, or fruits.

Words, phrases, and other text are frequently part of images. In our day-to-day lives, we come across street signs, documents, and advertisements. Humans see the text, read it, and comprehend it rather easily. For machines, this is an entirely different challenge. Via optical character recognition (OCR), computers can extract text from an image, enabling a wide range of potential applications. From language translation to mobile apps that assist the visually impaired, computer vision algorithms equip users to pull words out of a picture into readily usable text in applications.

In addition to general image tagging, computer vision can be used for more specific tasks. Some of the more common are the ability to detect logos and food items in images. Another frequent application of computer vision is in facial detection. With training, computer vision algorithms allow developers to recognize faces, sometimes getting even more specialized to detect celebrities.

## Object Localization

Another capability of computer vision is **object localization**. Sometimes your application's requirements will include not just classifying what is in the image, but also understanding the position of the particular object in relation to everything else. This is where object localization comes into play. Localization finds a specific object's location within an image, displaying the results as a bounding box around the specified object. Similarly, object detection then identifies a varying number of objects within the image. An example is the ability to recognize faces or vehicles in an image. **Figure 4-1** shows an example of object detection with dogs.

There are some challenges associated with object localization, however. Often, objects in an image overlap, making it difficult to ascertain their specific boundaries. Another challenge is visually similar items. When the colors or patterns match their background in an image, it can again be difficult to determine the objects.

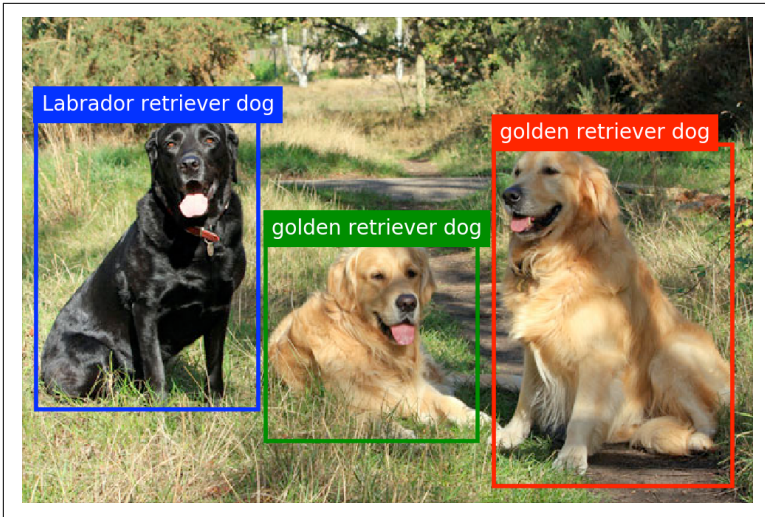


Figure 4-1. Object detection

## Custom Classifiers

Most of the time, you don't need to recognize everything. If you're looking to identify or classify only a small set of objects, custom classifiers could be the right tool. Most of the large third-party platforms provide some mechanism for building custom visual classifiers, allowing you to train the computer vision algorithms to recognize specific content within your images. Custom classifiers extend general tagging to meet your application's particular needs to identify your visual content. They primarily exist to gain higher accuracy by reducing the search space of your visual classifier.

At a high level, when creating a custom classifier, you'll need to have a collection of images that are identified as positive and negative examples. For example, if you were training a custom classifier on fruits, you'd want to have positive training images of apples, bananas, and pears. For negative examples, you could have pictures of vegetables, cheeses, or meats (see [Figure 4-2](#)).



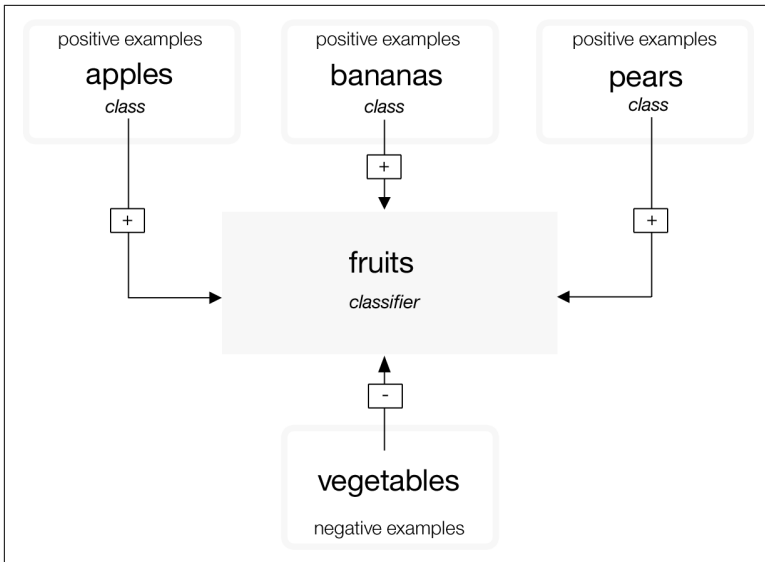


Figure 4-2. Creating a custom classifier

Most computer vision uses deep learning algorithms (discussed in “[Deep Learning](#)” on page 6), specifically **convolutional neural networks (CNNs)**. If you’re interested in building CNNs from scratch and going more in-depth with deep learning for computer vision, there are a wide variety of resources available. We recommend Andrew Ng’s **Deep Learning** specialization on Coursera as well as [fast.ai](#). Additionally, the following resources dive more in-depth into relevant computer vision topics:

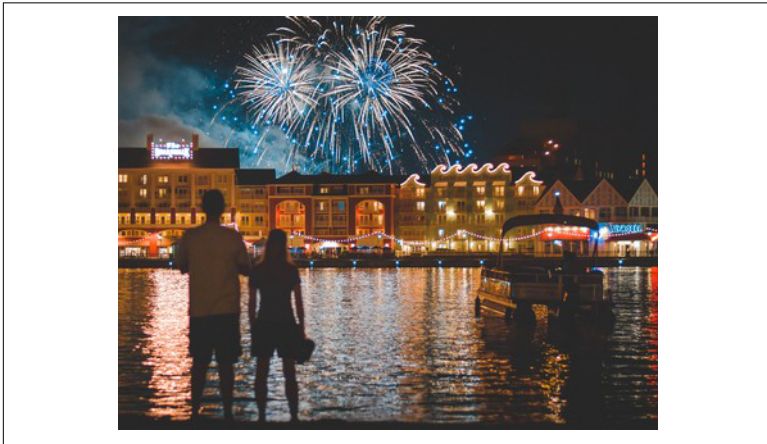
- *Programming Computer Vision with Python* (Jan Erik Solem, O’Reilly)
- *Mastering Computer Vision with TensorFlow 2.x* (Krishnendu Kar, Packt Publishing)
- *Learning OpenCV 3* (Gary Bradski and Adrian Kaehler, O’Reilly)

In the next section, we’ll start to look at how to use computer vision in enterprise applications.

# How to Use Computer Vision

When discussing how to use computer vision, we'll look at an example of the most common output—general tagging. As covered earlier, general tagging in computer vision returns to the user the overall items contained in the subject image. Frequently, depending on the algorithm or service used, the confidence levels are also returned. Based on your prototypes and testing, you can then use these scores to set your own thresholds for applying the returned tags based on your application's needs.

To demonstrate computer vision's tagging, we'll use the IBM Watson Visual Recognition service. Let's start by sending the image in [Figure 4-3](#) to the API.



*Figure 4-3. Fireworks over a harbor*

We'll make a simple request using `curl` to the API, referencing our image as the image URL parameter:

```
curl "https://gateway-a.watsonplatform.net/visual-recognition/  
api/v3/classify?api_key=YOURAPIKEY&url=  
https://visual-recognition-demo.ng.bluemix.net/images/samples/  
7.jpg&version=2016-05-20"
```

Submitting this API request returns the following JSON:

```
{
  "classes": [
    {
      "class": "harbor",
      "score": 0.903,
      "type_hierarchy": "/shelter/harbor"
    },
    {
      "class": "shelter",
      "score": 0.903
    },
    {
      "class": "Fireworks",
      "score": 0.558
    },
    {
      "class": "waterfront",
      "score": 0.5
    },
    {
      "class": "blue color",
      "score": 0.9
    },
    {
      "class": "steel blue color",
      "score": 0.879
    }
  ],
  "classifier_id": "default",
  "name": "default"
}
```

As you can see from these JSON results, the API returns varying classes with confidence scores on what it thinks the image contains. If we look at the picture again and compare it with the returned keywords, the essential elements of the image are returned with high confidence. To improve these scores for groups of images, you'd need to train the computer vision model using custom classifiers. We'll cover this more in an upcoming section, but mainly it provides the algorithm with many images of both positive and negative examples for training. A useful demo of this service can be found [online](#).

# Computer Vision on Mobile Devices

Another fascinating application of computer vision is that **algorithms may run locally on mobile devices** and IoT devices. Machine learning is now able to run on devices without necessarily having to connect to the cloud. Without sending and receiving data from the cloud, processing speed is increased, and security is improved.

There are some pros and cons associated with this approach, however. Let's start with the positive aspects of running computer vision locally. First, the speed and refresh rate of object detection are significantly faster than sending and receiving data from the cloud. Specifically, not needing to wait for the cloud to respond to a request speeds the process considerably. This is quite important for consumer products, such as virtual reality (VR) and augmented reality (AR), where our eyes are very much used to an immediate reaction. Quick response times for certain queries and questions make mobile processing more natural for responding in specific applications. Privacy is another benefit, as all computation is done locally, and no personal or private data is sent to remote servers. If the application needs data analysis, it only has to send the inference result rather than the actual data, thus preserving privacy.

Some examples of areas where **running computer vision locally** are advantageous include language translation, drones, autonomous vehicles, the self-monitoring of industrial IoT devices, and the ability to diagnose health issues without sending private data to the cloud.

While these are some compelling reasons to use computer vision locally on mobile devices, there are some drawbacks to discuss as well. For mobile applications, battery life and resource allocation are critical issues. Running computer vision locally uses much more computing power and requires optimization, so this is a factor to take into consideration when building your application. Additionally, it's more challenging to deploy computer vision applications locally, but is quickly becoming much easier. These days, you can even run a detection and classification system on a Raspberry Pi.

However, there is a hybrid approach that solves some of the issues—using a local model for simple object detection such as “Is that a dog?” and then sending the object to the cloud to determine the breed of dog. Another example of this hybrid approach is detecting

a product while a user is in a grocery store, but then sending the data to the cloud to retrieve the actual pricing information.

The key to running models locally is **model compression**. You might need a large model in the cloud if you're trying to recognize everything in the world, but the model can be much smaller if you're merely interested in recognizing your family members' faces. There are now **ways to reduce the parameters of visual recognition models** by magnitudes and thus effectively **reduce their size**.

## Best Practices

Next, we'll look at some best practices in developing applications with computer vision. We'll discuss what makes good training images shortly, but a general rule of thumb is that the more high-quality images you can provide for training, the better the results.

The accuracy of custom classifiers depends on the quality of your training data as well as the process. **Data from IBM Watson** has shown that clients "who closely controlled their training processes have observed greater than 98% accuracy for their use cases." This accuracy was based on the **ground truth** for their classification problem and data set.

## Quality Training Images

Let's now take a quick look at some guidelines for good training images. There are several characteristics of good training images. The images in your training and testing sets should ideally resemble each other in as many ways as possible. For example, be sure they're similar yet varying regarding angle, lighting conditions, distance to the subject, and the subject's size.

Also, make sure the training images represent what the test image will show. For example, if your test images show baskets of apples, then a close-up of a single apple would not be a good training image (Figures 4-4 and 4-5). Just because there's an apple in the picture doesn't mean it meets the criteria. Instead, you'd want many varying images of baskets of apples.



Figure 4-4. A single apple on a table (photo courtesy of *adrianbartel*, Creative Commons License 2.0)



Figure 4-5. Multiple apples in a basket (photo courtesy of *Mike Mozart*, Creative Commons License 2.0)

## Use Cases

Now that we've discussed computer vision in some detail, let's look at some industry examples and use cases.

### Satellite Imaging

When a drought in California reached a crisis level in April 2015, Governor Jerry Brown issued the state's first mandatory water restrictions. All cities and towns were instructed to reduce water usage by 25% in 10 months. Achieving this required measures more effective than just asking residents to use less water. Specifically, the state needed to run ad campaigns targeted at property owners using more water than necessary. Unfortunately, the government didn't even have water consumption data on such a granular level.

Scientists at OmniEarth came up with the idea of analyzing aerial images to identify these property owners. They first trained IBM Watson's Visual Recognition service on a set of aerial images containing individual homes with different topographical features, including pools, grass, turf, shrubs, and gravel. They then fed a massive amount of similar aerial imagery to Watson for classification. Partnering with water districts, cities, and counties, the scientists at OmniEarth could then quickly identify precisely which land parcels needed to reduce water consumption, and by how much. For example, they identified swimming pools in 150,000 parcels in just 12 minutes.

Armed with this knowledge, **OmniEarth helped water districts** make specific recommendations to property owners and governments. Such proposals included replacing a patch or percentage of turf with mulch, rocks, or a more drought-tolerant species, or draining and filling a swimming pool less frequently.

### Video Search in Surveillance and Entertainment

The proliferation of cameras in recent years has led to an explosion in video data. Though videos contain numerous insights, these are hard to extract using computers. In many cases, such as home surveillance videos, the only viable solution is still human monitoring. That is, a human sits in an operations center 24/7, watching screens and raising an alert when something happens, or reviewing dozens or even hundreds of hours of past footage to identify key events.

**BlueChasm** is a company looking to tackle this problem using computer vision. The founder, Ryan VanAlstine, believes that if successful, video can be a new form of sensor where traditional detection and human inspection fail. BlueChasm's product, VideoRecon, can watch and listen to videos, identifying key objects, themes, or events within the footage. It will then tag and timestamp those events and then return the metadata to the end user.

The industry that the company plans to focus on first is law enforcement. Ryan explains: "Imagine you work in law enforcement, and you knew there was a traffic camera on a street where a burglary occurred last night. You could use VideoRecon to review the footage from that camera and create a tag whenever it detected a person going into the burgled house. Within a few minutes, you would be able to review all the sections of the video that had been tagged and find the footage of the break-in, instead of having to watch hours of footage yourself." Once a video is uploaded, IBM Watson's Visual Recognition is used to analyze the video footage and identify vehicles, weapons, and other objects of interest.

BlueChasm is also looking to help media companies analyze movies and TV shows. It wants to track the episodes, scenes, and moments where a particular actor appears, a setting is shown, or a line of dialog is spoken. This metadata could help viewers find a classic scene from a show they'd like to rewatch by simply typing in some search terms and navigating straight to the relevant episode. More generally, for any organization that manages an extensive video archive, the ability to search and filter by content is an enormous time-saver.

## **Additional Examples: Social Media and Insurance**

Just as NLP was found to be incredibly useful in social media and content discovery, computer vision also proves to be similarly beneficial. For example, a company called **Ampsy** uses computer vision to grab all the historical images shared by an audience on social media to retrieve a list of activities, interests, people, and places for an influencer. Additionally, Ampsy uses custom visual classifiers to train on specific corporate logos. It can then detect all the images in a particular event where the advertiser's logos are captured. The company is then offering its users the capability to search for their logos in the collected images.



Similarly, social media analytics company **iTrend** uses general tagging to understand the content within images gathered from social media, blogs, and live streaming feeds. The company claims that it can analyze 20–50 times more data than its competition and then provide up to 80% of its findings as actionable insights.

Finally, another industry using computer vision in innovative ways is insurance. Primarily used in processing claims, visual recognition techniques and augmented reality are used by **insurance companies** to make the claims process much faster, decreasing costs and increasing customer satisfaction. They're accomplishing this through applications that allow customers to take multiple photos of accident and vehicle damage and send these to the insurer's servers for processing, eliminating much of the manual human processing that was previously required. According to **Accenture**, 82% of insurers believe AI-driven automation will be embedded in every aspect of their business over the next five years. Additionally, 35% of insurers report more than 15% in cost savings from this automation over the past two years. This is an industry using AI and computer vision techniques to their fullest potential to improve their business.

## Existing Challenges in Computer Vision

In **Chapter 2** on NLP, we discussed the challenges of extracting meaningful information from text, and doing the same for images is just as challenging. Let's forget about recognizing tens of thousands of objects, a complex classification problem, and instead examine a much simpler case of just identifying cats versus dogs in an image.

If you didn't have experience in machine learning before approaching this problem, you might start by writing an algorithm that separates the animals by color. But what if the pictures are black and white? Maybe you'd then try to separate them by color *and* texture. But what if it's a Golden Retriever instead of a Labrador? You get the point—building a successful visual classification system requires more than a rote list of rules, because you can always find exceptions that break them. The correct approach is to develop an algorithm that will automatically generalize to a set of rules based on a large enough data set. Over time, as you train the vision model by providing more positive and negative examples to learn from, the algorithm improves, ultimately providing an incredibly useful method for identifying objects within images.

Similar to our discussion of NLP challenges, several applications of computer vision still pose considerable challenges to implementation. First is facial recognition. While computer vision is very capable of face detection (detecting the presence of faces), face recognition (identifying individuals) does not come easily out of the box. Large data sets for training on individual faces are needed for this to be an effective process.

Another area of concern is detecting details. If you want to classify an image based on a small section of it or the details scattered within it, this tends to be challenging. Some services analyze an entire image when training, so they may struggle on classifications that depend on small details. A solution to this problem is **breaking down the image into smaller pieces** or zooming into the image's relevant parts.

## Implementing a Computer Vision Solution

As with the other AI solutions discussed in this book, you'll always face the decision of build versus buy when determining an implementation. As in previous chapters, we'll look briefly at some of the popular open source options and the SaaS services available.

As always, you need to try these yourself via testing and prototyping to see which meets your needs, including but not limited to price, accuracy, and the available code libraries.

Regarding open source, two of the primary options for building applications using computer vision are OpenCV and SimpleCV. They both provide **access to computer vision**, with OpenCV being a library available for many programming languages and SimpleCV a framework incorporating several libraries written in Python.

In addition to IBM Watson Visual Recognition, other companies are providing similar computer vision services as APIs. These include Clarifai, Google Cloud Vision, Azure Cognitive Services, and Amazon Rekognition. Each service has its pros and cons, so we again recommend testing and building prototypes for those that appear to meet your needs to find the proper fit.

## Summary

Computers “see” very differently than humans, and almost always within the context of how they are trained. The computer vision topics covered in this chapter have hopefully inspired you to use vision as a powerful tool to increase the utility of your enterprise applications. As we’ve discussed, numerous industries are already taking advantage of computer vision with great success. How might you incorporate computer vision in your applications?

Now that we’ve covered some of the significant use cases of AI in the enterprise (NLP, chatbots, and computer vision), let’s take a look at our data and how it gets processed in AI data pipelines.



# AI Data Pipeline

*In God we trust; all others bring data.*

—W. Edwards Deming

There are now **more mobile devices than people** on the planet, and each is collecting data every second on our habits, physical activity, locations traveled, and daily preferences. Daily, we create **2.5 quintillion bytes of data** from a wide variety of sources. And it's coming from everywhere. Just think of all the sources collecting data—IoT sensors in the home, social media posts, pictures, videos, all our purchase transactions, as well as GPS location data monitoring our every move.

Data is even touted as being **more important and valuable than oil**. For that reason, companies are creating vast repositories of raw data (typically called *data lakes*)—both historical and real-time. Being able to apply AI to this enormous quantity of data is a dream of many companies across industries. To do so, you have to pick the right set of tools not only to store the data but also to access it as efficiently as possible. Current tools are evolving, and how you store and present your data must change accordingly. Failure to do so will leave you and your data behind. To illustrate this point, MIT professor Erik Brynjolfsson performed a **study** that found firms using data-driven decision making are 5% more productive and profitable than competitors. Additional **research** shows that organizations using analytics see a payback of \$9.01 for every dollar spent.

As we've seen so far, if large amounts of high-quality data are a prerequisite for successfully implementing AI in the enterprise, then a process for obtaining and preparing the data is equally critical.

In previous chapters, we've covered some of the major applications of AI in the enterprise, from NLP to chatbots to computer vision. We've also discussed the importance of data to all of these implementations. What we've implied, yet haven't actually touched on to this point, is the concept of a data pipeline that forms the backbone for all these AI implementations.

Whether you use off-the-shelf technology or build your own, these AI solutions can't be effective without a data pipeline. With the numerous third-party solutions in the market, like IBM Watson, it's easy to forget that you need to have a data pipeline with even the simplest implementation. For example, in a computer vision solution, you still need to find representative images, use them to train the algorithm, and then provide a mechanism for repeating this loop with new and better data as the algorithm improves. With NLP, you still need to feed the APIs source text to process and then often train custom models with data from your domain. With chatbots, your initial data pipeline would focus on the known questions and answers from your existing customer support logs, and then build a process to capture new data to feed back into the chatbot. From SaaS to developing your own AI solutions, data pipeline needs grow even larger—and more critical to the overall implementation of AI in enterprise applications.

Not only is the data pipeline a crucial component of performing AI, but it also applies elsewhere in the enterprise—specifically in analytics and business intelligence. While the intricacies of creating a full AI data pipeline are outside the scope of this book, next we'll provide a high-level guide for getting started.

So what exactly is a data pipeline for AI? **Dataconomy** defines a data pipeline as “an ideal mix of software technologies that automate the management, analysis and visualization of data from multiple sources, making it available for strategic use.” Data preparation, a data platform, and discovery are all significant pieces of an effective pipeline. Unfortunately, a **data pipeline** can be one of the most expensive parts of the enterprise AI solution.

The key to this process is making sure data can be accessed in an integrated manner instead of sitting in different silos, both internal and external to the enterprise. This ability to access and analyze real-time, or at least recent, data is key to an AI data pipeline (Figure 5-1).

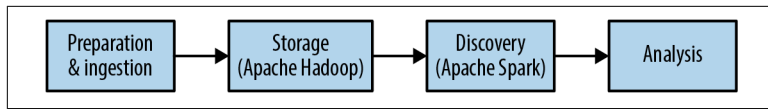


Figure 5-1. AI data pipeline

As we'll discuss in the next few sections, outside of the actual data itself, the most popular components of an AI data pipeline are Apache Hadoop, Apache Spark, and automation. In fact, an **IBM study** of organizations in 90 countries, each with more than 1,000 employees, showed that 57% of surveyed firms either already had or planned to implement a pipeline based on Hadoop/Spark.

## Preparing for a Data Pipeline

At a high level, there are several areas enterprises must focus on to effectively build and maintain a data pipeline. Before even beginning the process, the critical consideration is to make sure all stakeholders have bought into the idea of having a data pipeline. There are no shortcuts here, and enterprises must be sure of their preparedness in several areas. First, are your data storage and cloud practices solid? Are employees trained and aware of them? Next, where is the data stored? Frequently in the enterprise, the data is located in numerous silos, controlled and maintained by different groups in the company. Has everyone bought into the project and are they prepared to break down these data silos to build the AI pipeline? Finally, is there a process for cleaning data and repairing any issues with metadata? It's crucial that the enterprise embrace modern data science practices and not just rehash its existing business intelligence systems.

# Sourcing Big Data

Like AI itself, the term *big data* has various definitions depending on who you talk to. But generally speaking, once the data you're collecting is too large to store in memory or your existing storage systems, you're entering big data territory. Depending on the infrastructure in your enterprise, this will vary considerably. As the adage goes, you'll know it when you see it! For our purposes, we'll define big data as bringing structured and unstructured data together in one place to do some analysis with AI.

IBM describes big data as having **four major dimensions**: volume, velocity, variety, and veracity. *Volume* refers to the amount of data needed. As discussed in previous chapters, we're inundated with data, and it isn't slowing down. From the over 500 million tweets per day to the **350 billion annual meter readings** to better predict power consumption, enterprises generate large amounts of information. How fast are you storing and processing your data? Many applications are extremely time-sensitive, so the *velocity* of the data with respect to your storage and application processing is critical. For example, in areas like fraud detection and customer support, accessing the data quickly is essential. *Variety* refers to the diversity of the collected data, from structured to unstructured. This includes text, video, audio, clicks, IoT sensor data, payment records, and more. Finally, big data must have *veracity*. How accurate or trustworthy is the data? When **1 in 3 business leaders don't trust the data** they need to make decisions, it's clear that the vast majority of data has accuracy issues.

Returning to our previous discussion of AI winters and how the convergence of various trends has allowed AI to flourish again, a subtheme of this is how big data technology has come to the forefront. Commodity hardware, inexpensive storage, open source software and databases, and the adoption of APIs all have enabled **new strategies for creating data pipelines**.

## Storage: Apache Hadoop

Originally written in Java by Doug Cutting, Hadoop is an open source framework for distributed processing and computing of large data sets using MapReduce for parallel processing. Incredibly popular and rapidly growing, it's been estimated that the **global mar-**



ket for Hadoop will reach \$340 billion by 2027. So just what is Hadoop? IBM Analytics defines it as “an open source platform providing highly reliable, scalable, distributed processing of large data sets using simple programming models. Hadoop is built on clusters of commodity computers, providing a cost-effective solution for storing and processing massive amounts of structured, semi- and unstructured data with no format requirements.”

Hadoop can store data from many sources, serving as a centralized location for storing data needed for machine learning. Apache Hadoop is itself an ecosystem, made popular by its ability to run on commodity hardware. Two major Hadoop concepts we'll discuss that are relevant to an AI data pipeline are HDFS and MapReduce. Built to support MapReduce, the Hadoop Distributed File System (HDFS) can process both structured and unstructured data, resiliently enabling scalable storage across multiple computers. HDFS is a purpose-built filesystem for storing big data, while MapReduce is a programming paradigm that refers to two distinct tasks that are performed: **map and reduce**. The map job takes a set of data and converts it to key/value pairs. The reduce job then takes this output from the map job and combines it into a smaller set of key/value pairs for summary operations.

Like other powerful programming tools, MapReduce allows developers to write code without needing to understand the underlying complexion of distributed systems.

## Hadoop as a Data Lake

Hadoop is often used as a data lake. Again, while definitions vary, data lakes are typically considered shared storage environments for large amounts of varying data types, both structured and unstructured. This data can then be used for a variety of applications, including analytics and machine learning.

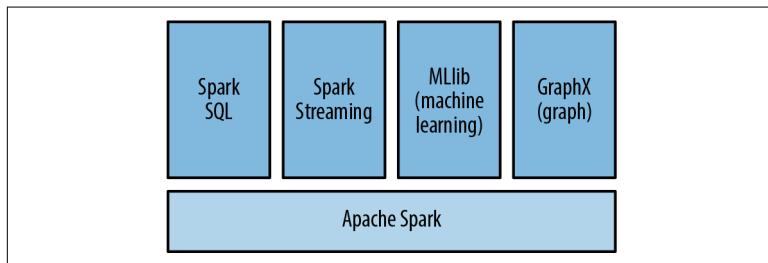
The main feature of a data lake is the **ability to centrally store and process raw data** where before it would be too expensive to do so. In contrast to data warehouses, which store structured, processed data, **data lakes** store large amounts of raw data in its native format, including structured, semistructured, and unstructured data. Hadoop shines in storing both this structured as well as unstructured data, thus making it an excellent tool for data lakes.

While data lakes have numerous benefits, from supporting data discovery to analytics and reporting, they do come with a caveat. As an **IBM report** stated: “Without proper management and governance, a data lake can quickly become a data swamp.”

## Discovery: Apache Spark

Created in 2009 at the University of California, Berkeley AMPLab, Apache Spark is an open source distributed computing framework that uses **in-memory processing to speed up analytic applications**. Written in Scala (though also supporting Java, Python, Clojure, and R), Spark is not necessarily a replacement for Hadoop, but is instead complementary and can work on top of Hadoop, taking advantage of Hadoop’s previously discussed benefits.

Contributing to its popularity among data scientists, the technology is extremely fast. According to **Databricks**, Apache Spark can be up to 100× faster than Hadoop MapReduce for large-scale data processing. This increased speed enables it to **solve machine learning problems** at a much greater scale than other solutions. Additionally, Spark comes with built-in libraries for working with structured data, streaming/stream processing, graphs, and machine learning (**Figure 5-2**). Numerous **third-party projects** have also created a thriving ecosystem around Spark.



*Figure 5-2. Apache Spark stack*

Spark itself doesn’t have a persistent data store and instead keeps data in memory for processing. It’s important to reiterate that Spark isn’t a database. Instead, it connects to external data sources, usually Hadoop’s HDFS, but also anything commercial or open source that developers are already using or familiar with, such as HBase, Cassandra, MapR, MongoDB, Hive, Google Cloud, and Amazon S3. Selecting a database for your application is outside the scope of this

book, but it's helpful to know that Spark supports a wide variety of popular database options.

## Spark Versus MapReduce

While using Spark doesn't necessarily preclude the use of MapReduce, it does in many ways compete with MapReduce. For our purposes, there are two main differences to consider.

The most significant difference is where the data is stored while processing. MapReduce stores the data to disk, continually writing I/O, while **Spark keeps the data in memory**. Writing to disk is much slower, so Spark often sees **performance gains of 100× over MapReduce**. Additionally, development is considered **easier and more expressive with Spark**, because in addition to map and reduce, Spark also has filter, join, and group-by functions.

## Machine Learning with Spark

As previously mentioned, Spark has a built-in module for machine learning called **MLlib**. This is an integrated, scalable machine learning library consisting of common learning algorithms and utilities, including classification, regression, clustering, and collaborative filtering.

Having this library native to Spark makes machine learning much more accessible, easy, and scalable for developers. It's easy to get started locally from a command line and then move on to full cluster deployments. With built-in machine learning, Spark can become a foundation to build data-centric applications across the organization. Since much of machine learning centers on repeated iterations for training, **Spark's ability** to store data in memory makes this process much faster and more efficient.

## Automation

Now that we understand the basics of a data pipeline, let's look at another emerging aspect—automation. As we've previously discussed, data is everywhere, and the raw amount is not an issue. However, data that is ready and useful for AI is often lacking.

We moved through the steps of preparation, storage, discovery, and now we're at the analysis stage in our pipeline discussion. As part of this step, we need to look at two areas—data cleaning and model

preparation—and how new advances in the field impact them. While building machine learning models from prepared data, data science professionals will often have to go back to clean, modify, and update the data. While we talk about these as two separate steps, [cleaning a data set and training a model](#) are often interwoven.

Unfortunately, building machine learning models can be an expensive proposition. Data from [Alteryx](#) shows that deployments can take six to nine months, with only 13% of the models making it to production, and it can cost in excess of \$250,000 to deploy a single model. Why is creating and deploying machine learning models so expensive? Generally speaking, it comes down to two high-level factors. First, it takes time to build models, with much of this being tedious and time-consuming work. Second, the talent needed to do this work is limited and expensive. While many companies are now trying to rapidly expand their talent base in this area, the stark reality is that there is a lack of qualified employees. Combine these two factors, and you can see why building and deploying models can take a considerable amount of time. So how is this being addressed in the enterprise? One of the ways is through the use of automated machine learning or AutoML.

Most AutoML systems consist of two main areas: data preparation and model prediction. Data scientists, developers, or other team members can import their data, add and identify labels, and then push a button to [generate a trained and optimized model](#) ready to make predictions using either a web GUI or a REST API for use in applications.

Automating machine learning sounds impressive. But what can truly be automated, and what are the benefits of implementing AutoML? [Wikipedia defines automated machine learning](#) as “automating the process of applying machine learning to real-world problems.” And [KDnuggets](#) describes it as “the automated process of algorithm selection, [hyperparameter](#) tuning, iterative modeling, and model assessment.”

For AutoML, it’s better to consider the technology as a toolset to help the overall pipeline, democratizing many of the functions previously reserved for data scientists, rather than as full automation. These systems are not ready to completely replace humans in the process, but can make their jobs easier, leading to higher

productivity. AutoML aides the process, freeing humans from repetitive, less challenging, or demanding tasks better done by machines.

## Why Automated AI/ML?

Now that we have a general understanding of automated machine learning, the next question is, why do we need it? There are three significant factors in favor of automation, and they're simplified to talent, time, and trust.

### Talent

There's a high demand for data scientists and machine learning practitioners, but not enough skills to meet the need. Also, the rare talent that is available is costly. For many projects, there are not enough people with the requisite knowledge.

Data scientists used to be worried about losing jobs to automation, but they're now finding those fears were mostly unfounded. Their specific expertise and domain knowledge are crucial to machine learning projects, and AutoML frees up time to focus on the challenging and domain-specific aspects, leaving the mundane work to automated systems.

In addition, AutoML solutions can allow non-data scientists to participate in the process. Where before, most, if not all, tasks in the AI lifecycle were performed by data scientists and specialists, AutoML now allows non-specialized employees to assist. From importing and cleaning data to selecting the optimized models generated by the AutoML, non-data scientists can contribute and free up the specialists to work on more challenging problems that suit their skill sets.

### Time

Creating and optimizing models is hugely time-consuming. Additionally, some tasks are quite mundane—precisely the types of functions machines are better at than humans. AutoML can save time and help people be more efficient in the overall process. It can remove the manual, time-consuming, and tedious processes that exist as part of the AI lifecycle.

Even for experienced data scientists, it takes a lot of time to get models tuned and trained. And as **data becomes increasingly**

**complex**, progress slows, and the AI data workflows become larger and more complex. Using AutoML can reduce errors associated with doing much of this work, saving a considerable amount of overall time.

## Trust

As we'll discuss more in **Chapter 7**, the issue of trust in our AI predictions and machine learning models is becoming paramount. When these systems are a black box and blindly spit out answers, it's reasonable to ask how the AI arrived at a decision. The questions include: what factors went into the system, and can we ascertain any bias in the data, process, or model? These are fundamental questions to ask as AI is increasingly integrating into the most critical areas of our lives. There is growing **skepticism of AI systems and processes**. AutoML software can help open up this black box with monitoring, thus providing the necessary risk management for many regulated industries.

## Limitations of AutoML

As discussed above, automation technologies are excellent for assisting the process, but not as a replacement for people. So if you're using an AutoML solution with the expectation it will replace a human operator, you're setting up your organization for disappointment and possibly failure. If, on the other hand, the goal is to augment your data scientists, make their lives easier, and improve effectiveness and productivity, you'll be more successful in implementing an AutoML project.

Two additional areas prove to be limitations. The first is **domain knowledge**, which is difficult to automate as it's one of the primary aspects data scientists bring to the solution. The deep understanding of a set of business problems is invaluable and uniquely human at this point. It's challenging to automate domain knowledge, so practitioners can focus on this while automation starts a level above. Finally, once the system has produced results, their interpretation is still reliant on a person. For example, the AutoML system may return several models with associated accuracy values to evaluate. The results may be ranked, but it's still up to a person to make the final decision on which to choose based on their expertise. These limitations prove that AutoML is excellent for assisting data scientists, but cannot replace their insights and expertise.

## Existing Automated Machine Learning Solutions

Developers, both in open source and commercial products, have been working on solutions to create and improve on AutoML. The solutions from enterprise vendors include IBM AutoAI, Microsoft Azure AutoML, H2O.ai, Google Cloud AutoML, and Amazon SageMaker. The numerous open source software solutions include [TPOT](#), [AutoKeras](#), [NNI](#), [Auto-sklearn](#), and [AdaNet](#). By simplifying the AI lifecycle management cycle, these systems can automate data preparation, model development and selection, feature engineering, and hyperparameter optimization.

## Summary

Big data is being captured everywhere and rapidly increasing. Where is your data stored, and what can you do now to future-proof it for machine learning? While the enterprise typically embraces the tried and true—in this case, Hadoop and Spark—other technologies show great promise and are being embraced in research as well as startup companies. Some of these open source projects include TensorFlow, Caffe, Torch, Chainer, and Theano.

And automated machine learning is a growing trend—Gartner predicts that by 2021, [AI augmentation will recover 6.2 billion hours of worker productivity](#). We're only going to see more applications of automation and assistance in the future and anticipate that these tools will continue to improve, offering more capabilities to companies and their data scientists.

Now that we've learned about the data pipeline needed for AI/ML work, let's look at an increasingly important consideration when working with enterprise AI—hybrid and multiclouds.





# Hybrid Clouds

*Cloud is about how you do computing, not where you do computing.*

—Paul Maritz, CEO of VMware

By this point, you probably have a good understanding of the cloud (and cloud computing.) We tend to equate all the cloud-related technology we use in our daily lives to the public version—from Netflix to Gmail to Amazon, almost everything we do or deal with daily has public cloud roots. But we actually use more than just the public cloud. In previous chapters, most of the topics centered around using public APIs and public cloud computing to perform AI tasks. But what happens when that doesn't fit an organization's requirements? We'll discuss numerous situations where utilizing only a single cloud type doesn't make sense, and we need to look to hybrid and multicloud solutions, especially for AI.

Many readers will be surprised to learn how prevalent hybrid and multicloud are in the enterprise. According to a study by [Flexera](#), 93% of enterprises have a multicloud strategy, and 87% a hybrid cloud strategy. Of these, they average using 2.2 public and 2.2 private clouds each. And [Gartner data](#) shows that 81% of public cloud users choose two or more providers. But while most enterprises are using the cloud, they've yet to embrace it fully. [McKinsey estimates](#) that less than 20% of enterprise workloads have moved to the cloud, presenting a tremendous opportunity. The numbers may vary slightly, but the direction is an increasing use of hybrid and multicloud environments for most enterprises.

In the next several sections, we'll discuss these trends along with the different cloud types and their use in AI.

## Overview of Clouds

Why is the percentage of workloads using the cloud still low? What holds these companies back? Historically, there have been three main considerations, which we'll discuss in the next section: security, compliance, and cost.

### Concerns about the Cloud

Security is critical for almost every enterprise. Many companies are concerned (rightly) about storing their data in the public cloud. What happens if there's a breach, and how will they handle the subsequent liability concerns?

Next is the related issue of compliance. Certain industries have to adhere to specific regulations on data usage and storage. Most often, this applies to customer data. For example, the organization may have personally identifiable data that cannot be stored in public cloud databases. Requirements around data privacy or industry regulations frequently prevent storing, processing, and manipulating data on the public cloud. Specifically, there are considerable regulations around customer data in the financial and healthcare sectors.

Finally, there are cost issues. We often assume using the public cloud is less expensive (especially the on-demand computing aspect), and it can be, but that's not always the case. Suppose you include the total costs involved with moving workloads to the public cloud from in-house and include data transfer costs. In that case, it can often be considerably more expensive, especially for AI workloads. Depending on an organization's requirements, the public cloud may not be a complete option. And as we'll see later in this chapter, while these factors apply to most common applications, they prove to be acute issues in AI-related ones.

And the challenge is not just managing multiple clouds with a single vendor. Today's computing needs also dictate the ability to compute across different vendors and their cloud resources. This creates numerous challenges. How can an enterprise manage these workloads across multiple vendors, data stores, and clouds? As we'll discuss soon, the answer is a hybrid cloud (usually managed with

Kubernetes.) Before we discuss them, though, let's step back and look at the major cloud types and the building blocks needed for their use, starting with containers.

There are several different types of computing clouds: public, private, hybrid, and multicloud. Sometimes these are also referred to as deployment models versus generic cloud types. We're going to briefly cover the four significant clouds you'll encounter, especially when considering AI. So, just what are these different cloud options? How are they similar, as well as different? And what do they mean for artificial intelligence?

## Public and Private

At the outset of this chapter, we mentioned **public clouds**. Their primary advantage is that organizations do not pay for the hardware, software, and related infrastructure, which are all maintained by the cloud provider. Public cloud computing can provide cost savings, but more importantly, it provides flexibility to run workloads as needed, paying only for the resources required and used. Unfortunately, you're also paying for all the data transfer in the public cloud, which is particularly crucial in AI and data science. The public cloud can get expensive, offsetting any potential cost savings.

On the other hand, a private cloud is a set of computing infrastructure owned and run by a single organization. Location is not a requirement for the private cloud definition, and it can be operated either inside or outside the company, run internally or externally by a vendor. The key is that everything is for a single company with nothing shared between organizations. Private clouds provide the advantages of more control, improved compliance and governance, and increased security.

## Hybrid

Attempting to deliver the best features of both, **hybrid clouds** use a mix of private and public plus orchestration to coordinate everything. Depending on requirements, workloads can be distributed across a combination of both to take advantage of the individual benefits. For example, sensitive data can be kept on the private cloud, while applications and workloads that need elastic scaling can run on the public cloud. Alternatively, an organization can consider its critical versus non-critical workloads. Where should each run?

For example, critical ones could be moved to a hybrid model to ensure availability and **continuity** in emerging crises. Again, hybrid clouds provide the option to make the best decision for each.

Additionally, **hybrid clouds** offer a host of other benefits. Typically they provide increased security, built-in at the container level, and offer better flexibility for enterprise workloads. When using hybrid, the organization can deploy and run its applications where they perform best instead of relying only on the infrastructure at hand. Resources can be used where required to manage workloads, placing those that need to scale in the correct environments. This provides a higher level of agility as the enterprise infrastructure grows.

## Multicloud

Finally, we have multiclouds. They seem somewhat the same as hybrid clouds at a high level, and it sounds like merely using many different clouds—and that’s partially true. While these two terms are sometimes used **interchangeably**, in the simplest terms, hybrid clouds are a combination of using public and private clouds, while multicloud is the use of clouds from different vendors to run a mix of varying workloads. Multiclouds enable organizations to take advantage of each cloud’s various value propositions and customize the computing environment to their specific needs. Also, vendor lock-in is a genuine concern for many enterprises. Using a multicloud environment de-risks the fear of having too much reliance on any single vendor for all their computing needs. As requirements become more sophisticated, we can expect multicloud deployments to become more prevalent.

Despite the benefits, multicloud management is still an issue. Once an organization starts adding in clouds from different vendors, in addition to those already existing in a hybrid cloud environment, how does it manage all the moving parts? As we’ll see in an upcoming section, the major cloud vendors provide software services that effectively solve this challenge for the enterprise.

## Cloud Construction

Now that we’ve discussed the different types of clouds, let’s examine their construction. Once we understand that aspect, we can then look to using artificial intelligence on these different clouds, and expressly, on hybrid clouds.

Containers have quickly become the basic cloud computing unit, with **Docker** being the most familiar to developers and therefore synonymous with containers. With virtualization and virtual machines, the hardware is replicated, including its operating system, requiring considerable overhead. However, with containerization, you're **virtualizing the processes** instead. The container wraps everything needed for an application, including code and any dependencies. The footprint is smaller, and the containers can be distributed across the overall infrastructure much easier, running consistently in any environment. These factors have popularized containers over virtual machines, and according to **451 Research**, the market for container technologies is expected to reach \$4.3 billion by 2022.

But as you can imagine, as container usage increases, coordinating them on each deployment or update becomes increasingly tricky. Many of the benefits are lost due to the increased complexity of managing so many containers, frequently thousands of them. Kubernetes, **an open source platform released by Google in 2014**, provides for easy orchestration of many containers in an environment and has become a popular option. And while Docker is the most popular container used with Kubernetes, you're not restricted to using it (others include rkt, containerd, and LXD). Kubernetes makes your applications much more portable across multiple clouds with improved scaling, making it the backbone for hybrid and multicloud implementations. Applications can be deployed and updated much faster with no downtime.

## Using AI on Hybrid Clouds

Particularly useful for AI, hybrid clouds provide the opportunity to utilize the best aspects of public and private clouds for applications and workloads to deliver AI at scale. Using a hybrid cloud for these projects can provide some key benefits, which we'll cover in the next sections.

## Privacy and Security

One of the primary advantages of a hybrid environment for AI is data privacy. As mentioned at the beginning of this chapter, there are numerous industry regulations regarding working with data on public clouds. Before hybrid clouds, that data may have been trapped behind a firewall and unavailable for any AI implementa-

tions. Running AI on hybrid clouds provides the opportunity to keep private data, or data with restrictions, separate from other data in the public cloud. By utilizing the best features of both, with data in public and private clouds, enterprises can now expand their AI initiatives in ways previously inaccessible, with greater control over data and privacy.

## **Portability and Scalability**

In a hybrid, multicloud environment, organizations have full control over where they want their applications deployed and data stored. Utilizing the benefits of containers with Kubernetes, the hybrid cloud enables them to easily deploy an application when and where needed. This is especially important when performing AI work, as workloads may not be consistent. There may be times when models need to be run, and the scaling of resources is necessary. When they're not being used, resources can be decreased. A hybrid cloud also makes data more accessible to team members. Data can be stored in different places, but it becomes much simpler to collect and organize it, making analytics more collaborative.

## **Processing and Computing Power**

Why make a considerable investment in powerful GPU processors for all your applications when you only occasionally need to run computing-intensive workloads like AI/ML? Hybrid clouds allow you to use the necessary computing power when and where it is required, saving costs and resources. They can scale in response to a load and then be deprovisioned as needed, such as after a seasonal event like a holiday sale. The appropriate use of computing resources can now happen across the enterprise infrastructure.

## **Example of AI on Hybrid Clouds**

Let's now look at an example of AI effectively implemented using hybrid clouds. Here, we'll take a look at a fictional mortgage company that needs a customer support application.

During the mortgage process, customers go back and forth with different employees, including the loan originator and loan processor. Especially on a home purchase, the customer is usually in a time crunch, with their loan file needing to be processed promptly. Many records are collected—especially PDFs of mortgage documents and

the customer's financial statements. The customer frequently contacts the company with questions, which can be reasonably answered, but take time away from employees and could be more efficiently handled with automation. As this customer data is private financial information, a secure, private, and integrated solution is necessary. Due to the data's sensitive nature, in this scenario it wouldn't be recommended to use a solely public cloud solution.

The company could use artificial intelligence in a few areas: NLP to process the questions, chatbots to help with the back and forth communication and interface to answer customer questions, and intelligent search on the documents, to name a few of the more obvious ones.

This is an ideal use case for a hybrid cloud. Sensitive customer data can be stored behind the firewall on a secure, private cloud. The customer support application, leveraging a chatbot, could be deployed to the public cloud and scaled up or down rapidly as needed. Then, it can still interface with data on the private cloud. Thus, a hybrid cloud could genuinely enable a new solution for this mortgage company, providing significant cost and time savings while making its customer support agents more efficient.

## Practical Solutions

As you can see from this discussion, while solutions like Kubernetes make the management more straightforward, there's still a considerable amount of work needed to run and maintain hybrid solutions. Often an enterprise will need a dedicated team to manage everything. This is where the various vendor solutions come into play. To do all the work yourself takes talent and resources your organization may not have.

Fortunately, many companies offer hybrid solutions designed for data and AI work, including Google Anthos, Microsoft Azure Stack, IBM Cloud Pak for Data, and AWS Outposts. These all have different features for successfully running applications on hybrid clouds. Each provides hybrid capabilities and the functionality of keeping your data secure. They mainly vary on their foundational technology and what specific services they include. Still, each offers the main features we discussed regarding data portability, scalability, and privacy.

As we've done throughout this book, we'll look at an example from IBM to illustrate the concepts. IBM Cloud Pak for Data provides an integrated platform consisting of a set of unified data and AI services running on **Red Hat OpenShift**. Built for developing and managing containerized applications, Cloud Pak for Data also enables interactions to popular open source tools.

Key features include the ability to deploy to any cloud and providing a full solution to collect, organize, and analyze data, particularly incorporating AI/ML. For example, an organization may have most of its infrastructure on AWS but wants to use IBM Watson for its AI/ML capabilities. Cloud Pak for Data would enable this as an option. Products like this assist in avoiding vendor lock-in and can run anywhere. As we discussed previously, one of the main benefits of hybrid and multicloud environments is the option to pick services that best meet the organization's needs and then deploy them on systems from any vendor. An organization may prefer a software solution from one vendor, but run the majority of its infrastructure on the software of another. These types of hybrid, multicloud solutions provide the ability to do this. Additional features include data virtualization and warehousing, as well as data discovery, search, and analytics capabilities. An integrated platform allows developers, data scientists, and other business users to collaborate and work with data more easily, thus saving costs and increasing efficiency. It helps streamline work by creating a pipeline for collecting, organizing, analyzing, and consuming data.

Further, there are distinct benefits for enterprise developers when working with hybrid software solutions. Instead of having to jump between different developer tools to get their jobs done, solutions like Cloud Pak for Data offer the necessary utility integration. A single platform for data management and analysis allows them to easily manage data connections and access analysis tools. Finally, developers often don't like to spend their time worrying about DevOps and monitoring. These platforms offer built-in capabilities for core operational services, including logging, monitoring, and security, so developers can spend more of their precious time focusing on their applications and code.

We can expect more focus on these solutions from major vendors. As we've seen, the benefits are too great to ignore in the current environment, and the ability to leverage each solution's strength will be paramount in the future.



# The Future of AI on Hybrid Clouds

What should we expect from hybrid clouds moving forward? A trend we're currently seeing is the increasing use of open source on hybrid clouds. Projects like [Open Data Hub](#) prove to be an excellent alternative to commercial solutions when data scientists desire an integrated platform to perform their work. Based on Red Hat OpenShift, Open Data Hub is a complete data science and AI platform. Integration with Jupyter Notebook, TensorFlow, Apache Spark, and scikit-learn all allow for model building and development. One of the best features of Open Data Hub is that it only needs a running Red Hat OpenShift cluster and an S3 object store like Ceph. With just a few clicks, a basic installation can be up and running on top of OpenShift for users to start working on their projects. Moving forward, we're likely to see more attempts to make it easier to do data science and AI work on open source platforms.

Finally, many enterprises will increasingly use hybrid cloud and multicloud strategies together. As they continue to modernize applications that are still mostly on-premise, the hybrid and multicloud use cases will dramatically increase.

## Summary

As we've seen throughout this chapter, hybrid clouds are powering the future of enterprise AI. Solving many data issues facing enterprises, they'll pave the way for use cases previously unattainable. Are there potential opportunities in your organization that using a hybrid cloud for AI will open up?

In the next chapter, we'll wrap up the discussion of AI in the enterprise with a look at how to move forward and begin your journey.



---

# Looking Forward

*Over the next decade, AI won't replace managers, but managers who use AI will replace those who don't.*

—Erik Brynjolfsson and Andrew McAfee, [Harvard Business Review](#)

Throughout this book, we've discussed practical techniques for implementing AI in the enterprise. Ranging from NLP to chatbots to computer vision, these technologies provide businesses not only significant cost savings over the long term but also the ability to solve problems they previously could not.

While we've discussed the benefits and methods of implementing AI in the enterprise, there are a few additional aspects any AI practitioner should keep in mind for the future. The next few sections discuss some of these forward-looking areas in the enterprise that artificial intelligence will impact.

First, let's recap our initial discussion of AI in the first chapter. Despite what we see in the media, it's important to remember that we're not talking about building general artificial intelligence in the enterprise (at least in the relatively short term anyway!). Instead, we're talking about creating *augmented* intelligence for the enterprise. Augmenting human intelligence is more about scaling our human capabilities and helping employees make better decisions, versus creating a newly intelligent life-form.

Enterprise technologies that automate and detect patterns can now advise and enhance human expertise, empowering both employees

and applications to make richer, more data-driven decisions. This is just an extension of what we've already been doing with computers, but helping humans to make these better decisions is now the real problem we're trying to solve in the enterprise.

Most machines can't read the majority of data that's created, as it's not structured. As discussed in [Chapter 2](#), 90% of the world's data was created in the last two years, and 80% of that data is unstructured. So again, dealing with messy, unstructured data becomes a primary focus for developers of enterprise applications using AI.

Next, we'll discuss some aspects of enterprises that make them particularly unique compared to other types of companies. Then we'll examine some current challenges, trends, and opportunities for bringing AI to the enterprise. We'll wrap up the chapter by exploring scalability and the societal impact of AI on an organization.

## What Makes Enterprises Unique?

In [Chapter 1](#), we laid out a basic definition of an enterprise to make sure we were all on the same page. As a refresher, we defined an enterprise as a company whose end users were other businesses or business employees. In this section, we'll expand on that by looking at some of the unique features of enterprises.

First, let's come back to data. For enterprise companies, their data is of the highest importance. Enterprise data is not a commodity to be monetized or traded. Not only are there privacy and personal data issues to worry about, but there's also a primary concern for protecting data that becomes valuable intellectual property (IP). Additionally, enterprise companies must have control and choice over how this data, now critical IP, is handled. Even if they let the data move somewhere else, there are not enough people (likely in the world) with the skill required to label *their* specific data. All these reasons make managing data in the enterprise a much more complicated endeavor than it is for their consumer-focused counterparts.

Another area specific to the enterprise is its overall technology stack and how AI integrates into it. Technology in the enterprise is typically focused on solving complex but mundane problems. So while adding AI to applications sounds sexy, there are many stages that developers must go through to integrate the technology. Any new implementation of AI in applications must adapt to the usually

numerous existing enterprise processes. The positive side of this more extended undertaking is that the AI implementation can often be customized to business-specific problems, thus providing higher orders of value to the firm.

Finally, most enterprise companies need domain and industry expertise. For most machine learning problems, they require unique and pretrained data sets. Also, there are often industry-specific concerns to consider. For example, as mentioned in the hybrid cloud chapter, the **financial**, **healthcare**, and **human resources** industries all have distinct sets of regulations that must be addressed and followed. Other industries, such as cybersecurity, IoT, and customer care/support, all have special requirements for data management. While each enterprise is unique in many ways, the key is building a set of customization tools that lets each company define the solution to its needs.

## Current Challenges, Trends, and Opportunities

In the future, developers face several challenges to implementing enterprise AI, all of which fortunately reflect some of the industry's future trends. We'll discuss these in the next few sections.

### Little Data

Enterprises have smaller or less data compared to consumer applications. We've talked a lot about big data in previous chapters, but the truth is that while big data exists in the enterprise, it's typically on a much smaller scale than for consumer-facing applications. In addition to less data, the information that does exist is frequently locked up, inaccessible, and unable to be shared for various reasons (including IP and regulations). So the crucial issue here is how to obtain high accuracy in the algorithms using much less data. The typical solution to working efficiently with a small amount of data is a technique called *transfer learning*.

Transfer learning is an approach that can use much smaller amounts of data, as opposed to deep learning, which needs lots of data. **Transfer learning** allows the leveraging of existing labeled data from a different (though closely related) domain than the one in which we're trying to solve the business problem. What's interesting about this approach is that even as humans, we learn from relatively small amounts of data—we don't learn from scratch, nor do we learn from

large amounts of data. We're able to learn from the collection of experiences we've gained somewhere else. Transfer learning is an approach that mimics this human process, which parallels deep learning's relationship to how the brain works.

In the enterprise, there are different workspaces or domain areas, so how do we use what we've learned in one workspace to inform another in an efficient manner? This is where transfer learning comes into play, and it will be a driving force in the future.

Additionally, **one-shot** learning is the ability to solve problems using an even smaller amount of training data. Most frequently used in computer vision problems, one-shot learning attempts to classify images using a single or very few training images.

## Inaccessible Data Formats

While the collective AI industry has made prodigious progress, there are still areas, especially in NLP, that prove a particular challenge for the enterprise. Specifically, much knowledge in companies is contained in PDFs, Word documents, and other proprietary formats. Within these documents, tables and graphs are extremely tough to understand for NLP algorithms.

Being able to understand these tables and graphs is a challenging problem. It's much easier to recognize images given enough high-quality data. Given a structured table in a PDF, there are infinite combinations. Companies like IBM are actively working on solutions to this problem and think they are close to a solution. Again, in the enterprise, these are some of the unglamorous issues you have to deal with that are quite essential to the business. You can't just throw these documents away, as they're often half of your data.

## Ensuring Accountability and Interpretability

Even if we use machine learning algorithms (either developed in-house or through third parties) that we deem effective, there's often the issue of explaining the results and defending the algorithm to stakeholders. A question that frequently gets asked of developers is, how did you reach these conclusions? Deep learning, in particular, doesn't offer much transparency into the model and this can pose problems in some situations.

Many of these algorithms are black boxes. We can understand how the algorithm works technically. Still, the steps it takes from iterating over training data to the final predictive results are not usually laid out for anyone to review and interpret. Adding to the issue is that interpretability, by definition, is subjective for humans.

Trying to balance the quantitative algorithms with the qualitative interpretive processes is exacerbated by more algorithmic improvements. It's great to be innovative with our machine learning algorithms, but at what expense? Regardless of the precision, anything less than 100% can cause stakeholders to focus on the error rate instead of the success of high accuracy. Here's a familiar expression: "80% accuracy, that sounds great, but why was there a 20% error rate?" There needs to be a trade-off between accuracy and interpretability. But what we've seen is that often **more accuracy means less interpretability**. This lack of data interpretability then poses issues for companies in many industries, especially those with regulatory requirements, and any errors in these highly regulated industries are much more costly.

So the issue for an enterprise developer now becomes, how can I explain the data to both internal teams and outside regulators? Again, the predictions from machine learning models are opaque and hard to understand. How did our algorithms get to this conclusion? And how reproducible are the results? When they're a black box, this becomes quite a challenging problem. Audit trails are vital to tackling this problem. Tracking a prediction to each model's unique heritage is critical to regulatory compliance. Additionally, enforcing access controls for model sharing and deployment ensures data security and application stability.

Outside the enterprise's needs (executives and other employees), external forces are coming into play. In May 2018, the EU's General Data Protection Regulation (GDPR) took effect, granting consumers a limited legal "right to explanation" from organizations that use algorithmic decision making. The **GDPR** applies to any company "if they collect or process personal data of EU residents." And while many experts believe this type of legislation can potentially **slow down the development and use of AI**, the political trend is clear, and we can expect to see more of these initiatives in the future. For example, there are nine other countries with **similar laws**, either existing or soon to be adopted. While the ramifications of these laws are still in question, the fact remains that governments have

demonstrated their willingness to legislate ways for companies to shed light on their algorithmic decision making. For more information on possible solutions, see Patrick Hall's [excellent overview](#) of striking a balance between accuracy and interpretability.

Directly related to the issue of regulation concerns is bias in the data and algorithms. It's critical to understand any partiality in the data or algorithms. There have been some high profile [examples of bias](#) from large companies. The concern is if developers are not responsible for making sure bias is removed from the outset—both from their machine learning algorithms as well as the data—then governments will regulate this for them. It's much better to self-regulate, not only from a legal perspective but also from a social responsibility perspective, as biased results have real-world impacts.

This is a thorny problem, as most developers and data scientists are likely not planning to create bias initially. Unfortunately, it is challenging to identify our own biases, let alone those that creep into our models over time. And as we discussed, this is even more difficult to detect when the algorithms are mostly black boxes. While humans can [self-reflect](#) and try to identify bias, our algorithms have no such built-in mechanism. Future work in this area will likely provide enormous overall gains for all stakeholders.

## Fairness, Trust, and Transparency

Because most industries increasingly rely on AI, we must ensure it is fair and produces results that are transparent and trustworthy. Unfortunately, this is easier said than done. We've now seen many high-profile instances in which AI has been shown to perpetuate bias and unfairness. While most of these unintended, biased results are a product of [unconscious bias](#), that often makes it more of a challenge to combat. When the bias is not explicit, it's harder to identify and thus eliminate.

When we talk about bias in AI, precisely what do we mean? Let's take a look at some specific examples that illuminate the issues. Research from the [Computation Policy Lab at Stanford](#) shows that speech recognition tools misunderstand black speakers twice as often as white, likely due to lack of diversity in training data, and they perform particularly poorly for black men. Facial recognition has also consistently shown [poor performance in identifying non-white faces](#). In contrast, research published in *Science* showed the



inadvertent **transfer of medical resources** to healthier white over sicker black patients. Other examples include the **Apple Card** assigning lower credit limits for women, **Google's social media** AI algorithm for monitoring hate speech being racially biased against African Americans, and the notorious issue of image and video **deepfakes**. This problem ends up being a **feedback loop of bias** and affects everything from the criminal justice system to credit scoring to job interviews. The list goes on the more you research.

How can enterprises then govern AI while scaling? How can we increase the explainability and transparency of AI to ensure fairness for all? Relatedly, what happens if you can't trust your AI models? What if they're not fair, ethical, or unbiased?

The key to any solution will include features to improve explainability, fairness, traceability, and bias detection to understand and explain how a model makes its predictions. Additionally, we can examine whether it's open and accepts community contributions to continually improve. For example, in the financial services or insurance industries, these systems would be crucial when making decisions on customer creditworthiness. These firms would need to articulate how and why a customer was declined a loan versus just presenting a simple decision. In this way, the company can show there was no bias in the results. Otherwise, it's open to interpretation and debate (and potential liability).

Services are working to solve this problem and create debiased models. One example is IBM Watson OpenScale, a platform that gives developers a clear view of their AI systems and the ability to add trust to their models. Across the AI lifecycle, a product like OpenScale tracks and measures outcomes, automating more of the process. Also, tools like the **AI Fairness 360 Toolkit** provide an open source library for developers and data scientists to identify, report, and then mitigate discrimination and bias in their models. To **mitigate discrimination** and bias in machine learning models, the toolkit contains fairness metrics and bias mitigation algorithms and is available in Python and R.

While this issue is disturbing, there's promise on the horizon. In addition to software vendors' solutions, there's an increasing awareness of the problem, which will hopefully lead to change. Best practices to avoid bias are being developed and shared. For example, the FTC offers **recommendations** on data and AI best practices to

prevent bias. So while this is a pressing issue, it's mostly one we've created as AI practitioners, which means it can also be solved by us, first through recognition and then better implementations. Moving forward, it's going to be critical that we improve our machine learning models and reduce risk. The key to this will be explainable and bias-free AI.

## Scalability

As enterprise AI demands grow, the ability to train large amounts of data is critical. But, in most cases, the training times for deep learning are incredibly long. Large models can take days or even weeks to train to desired levels of accuracy. The main culprit behind this performance lag is often a hardware limitation. Many popular open source deep learning libraries do not perform efficiently across multiple servers (nodes). So far, the solution has been to scale the number of GPUs on a single node, but performance gains have been limited to this point. Interestingly, the key to future improvement will be this ability to train large amounts of data and to scale the effort accordingly across many nodes.

**IBM Research** has recently shown promising results in the area of distributed deep learning ([Figure 7-1](#)). Researchers were able to linearly scale deep learning frameworks across up to 256 GPUs with up to 95% efficiency by building a library that connects to open source frameworks like TensorFlow, Caffe, Torch, and Chainer.

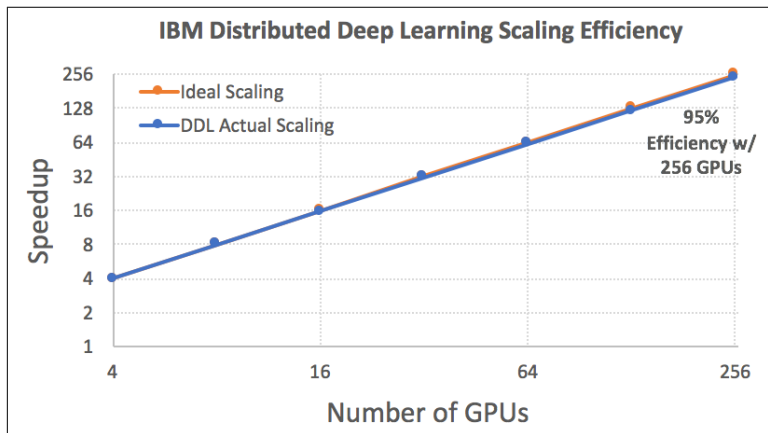


Figure 7-1. IBM Research's "Distributed Deep Learning" library

## Social Implications

While not necessarily a technical aspect of building AI applications in the enterprise, the societal implications of AI are also essential to discuss. Despite the stops and starts over the years, as well as concerns about the impact of automation on jobs and society, the potential of AI in the enterprise is only increasing—the benefits are too significant. For example, [McKinsey Global Institute](#) found that “45% of work activities could potentially be automated by today’s technologies, and 80% of that is enabled by machine learning.”

Additionally, 10% of all jobs in America are driving-related. What happens when AI-powered and automated, self-driving vehicles are ubiquitous, and 30 million jobs have been eliminated or vastly reduced in capacity? It begs the question: if any technology eliminates jobs, can it also create new ones? How do we effectively train workers, especially those displaced by modern technology, to work in these new industries?

The solutions to this problem are debatable, but one thing is sure—throughout history, advances in technology have created upheaval across industries and eliminated jobs. However, they’ve always [generated new ones](#) in their place. The key will likely be the retraining and continuing education of workers to take advantage of the new positions that will be created. According to the [World Economic Forum](#), while 75 million jobs will be lost, 133 million new jobs will be generated by AI by 2022.

Although we worry about the loss of jobs, the other side of the coin is that there’s a desperate, immediate need for more engineers and data scientists trained on applying AI, especially in the enterprise.

While not necessarily likely to slow down technological progress and innovations, this increasing shortage of engineering talent will affect how enterprises build applications. Without the ability to hire or train in-house expertise fast enough, companies will be forced to either outsource these capabilities to third parties (if they are not facing the same shortages) or rely more heavily on SaaS solutions for these skills.

## Summary

Building enterprise AI applications is more about augmenting than replacing human intelligence. In this chapter, we examined this topic as well as what makes enterprises unique, especially their data. We also discussed some of the challenges in enterprise AI including data confinement, little and inaccessible data, and social accountability. While each poses problems, they're also opportunities to make better applications and improve the enterprise development process.

We've covered a lot of ground in this book, but hopefully we've provided you a starting point to apply AI in your enterprise applications. The future is always cloudy, but one thing is sure—AI is here to stay and will be necessary for your enterprise applications to remain relevant today and into the future. Good luck with your development efforts!

## About the Authors

---

**Tom Markiewicz** is a developer advocate for IBM. He has a BS in aerospace engineering and an MBA. Before joining IBM, Tom was the founder of multiple startups. His preferred programming languages are Ruby and Swift. In his free time, Tom is an avid rock climber, trail runner, and skier.

**Josh Zheng** is a developer advocate for IBM. Before joining IBM, he was a full-stack developer at a political data mining company in Washington, DC. His favorite language is Python, followed by JavaScript. Since Josh has a background in robotics, he still likes to dabble in hardware as well. Outside of work, he likes to build robots and play soccer.