

---

---

# اکنون کلان داده

---

---

**مؤلف:**

اوریلی

**مترجمین:**

دکتر حامد رجبی قمی

مهندس مریم کاردگر



فن آوری نوین

---

---

عنوان و نام پدیدآور	:	اکتون کلان داده/ مولف اوریلی؛ مترجمین حامد رجبی قمی، مریم کاردگر.
مشخصات نشر	:	بابل: فناوری نوین، ۱۴۰۰.
مشخصات ظاهری	:	۱۲۸ ص.
شابک	:	۹۷۸-۶۲۲-۷۳۹۳-۳۳-۰۰: ریال ۵۰۰۰۰۰
وضعیت فهرست نویسی	:	فیبا
یادداشت	:	عنوان اصلی: ۲۰۱۴، <b>Big data now</b>
یادداشت	:	کتاب حاضر در سالهای مختلف توسط ناشران و مترجمان متفاوت منتشر شده است.
موضوع	:	داده‌های کلان
موضوع	:	<b>Big data</b>
موضوع	:	داده کاوی
موضوع	:	<b>Data mining</b>
موضوع	:	هوش مصنوعی
موضوع	:	<b>Artificial intelligence</b>
شناسه افزوده	:	رجبی قمی، حامد، ۱۳۶۰-، مترجم
شناسه افزوده	:	کاردگر، مریم، ۱۳۶۶-، مترجم
شناسه افزوده	:	شرکت اوراییلی مدیا
شناسه افزوده	:	<b>O'Reilly Media, Inc</b>
رده بندی کنگره	:	۷۶/۹QA
رده بندی دیویی	:	۰۰۵/۷۴
شماره کتابشناسی ملی	:	۷۵۷۱۰۲۵
وضعیت رکورد	:	فیبا

## @fanavarienovinpub

تلفن: ۰۱۱-۳۲۲۵۶۶۸۷

بابل، کد پستی ۷۳۴۴۸-۷۱۶۷

فن آوری نوین

اکتون کلان داده

ترجمه: حامد رجبی قمی، مریم کاردگر.

نوبت چاپ: چاپ اول

سال چاپ: بهار ۱۴۰۰

شمارگان: ۲۰۰

قیمت: ۵۰۰۰۰ تومان

نام چاپخانه و صحافی: دفتر فنی سورنا

شابک: ۹۷۸-۶۲۲-۷۳۹۳-۳۳-۰۰

نشانی ناشر: بابل، چهارراه نواب، کاظم بیگی، جنب مسجد منصور کاظم بیگی، طبقه اول

طراح جلد: کانون آگهی و تبلیغات آبان (احمد فرجی)

**فروشگاه و پخش کتاب چاپی: تهران، تلفن ۶۶۴۰۰۱۴۴-۶۶۴۰۰۲۲۰**

تهران، خ اردیبهشت، نبش وحید نظری، پلاک ۱۴۲ تلفکس: ۶۶۴۰۰۱۴۴-۶۶۴۰۰۲۲۰

# فهرست مطالب

فصل اول: مقدمه.....	۵
فصل دوم: افزایش سرعت با کلان داده.....	۶
۱-۲. کلان داده چیست؟.....	۶
۱-۱. کلان داده چه شکلی هستند؟.....	۶
۲-۱-۲. در عمل.....	۱۱
۲-۲. Apache Hadoop چیست؟.....	۱۳
۲-۲-۱. هسته هادوپ: نکاشت کاهش.....	۱۳
۲-۲-۲. لایه‌های پایین‌تر هادوپ: HDFS و نکاشت کاهش.....	۱۴
۲-۲-۳. بهبود قابلیت برنامه‌نویسی: Pig و Hive.....	۱۴
۲-۲-۴. بهبود دسترسی به داده: Flume و Sqoop، HBase.....	۱۵
۲-۲-۵. هماهنگی و گردش کار: Oozie و Zookeeper.....	۱۶
۲-۲-۶. مدیریت و استقرار: Whirr و Ambari.....	۱۷
۲-۲-۷. یادگیری ماشین: Mahout.....	۱۷
۲-۲-۸. کاربرد هادوپ.....	۱۷
فصل سوم: ابزارها، تکنیک‌ها و استراتژی‌های کلان داده.....	۲۱
۳-۱. طراحی محصولات کلان داده.....	۲۱
۳-۱-۱. محصولات داده‌ای مبتنی بر هدف.....	۲۱
۳-۱-۲. خط مونتاژ مدل: یک مطالعه موردی از گروه تصمیمات بهینه.....	۲۲
۳-۱-۳. رویکرد پیش‌رانه برای سیستم‌های توصیه.....	۲۷
۳-۱-۴. بهینه‌سازی ارزش طول عمر مشتری.....	۳۰
۳-۱-۵. بهترین روش‌های محصولات داده‌ای فیزیکی.....	۳۲
۳-۱-۶. آینده محصولات داده‌ای.....	۳۶
۳-۲. ساختن محصولات یادگیری ماشین بزرگ چه بهایی دارد.....	۳۷
۳-۲-۱. پیشرفت در یادگیری ماشین.....	۳۸
۳-۲-۲. مسائل جالب هرگز از قفسه خارج نمی‌شوند.....	۳۹
۳-۲-۳. تعریف مسئله.....	۴۱
فصل چهارم: کاربرد کلان داده‌ها.....	۴۳
۴-۱. ماجرای صفحات گسترده.....	۴۳
۴-۱-۱. بینشی بر داشبورد.....	۴۴
۴-۱-۲. مصاحبه کامل.....	۴۵
۴-۲. کاوش ادبیات نجومی.....	۴۵
۴-۲-۱. مصاحبه با رابرت سیمپسون: پشت پرده و آنچه پیش می‌آید.....	۵۰
۴-۲-۲. علم بین شکافها.....	۵۳
۴-۳. قسمت تاریک داده.....	۵۴
۴-۳-۱. چشم‌انداز انتشار دیجیتال.....	۵۵
۴-۳-۲. حریم شخصی در طراحی.....	۵۶
فصل پنجم: در کلان داده به دنبال چه چیزی هستیم.....	۵۸
۵-۱. کلان داده مسئله حقوق شهروندی نسل ماست، و ما آن را نمی‌شناسیم.....	۵۸
۵-۲. سه نوع کلان داده.....	۶۳
۵-۲-۱. Enterprise BI ۲,۰.....	۶۴
۵-۲-۲. مهندسی عمران.....	۶۵
۵-۲-۳. بهینه‌سازی ارتباط با مشتری.....	۶۷
۵-۲-۴. Headlong into the Trough.....	۶۹

۶۹.....	۵-۳. علم خودکار، داده‌های عمیق، و پارادوکس اطلاعات
۷۰.....	۵-۳-۱. علم (نیمه) خودکار
۷۱.....	۵-۳-۲. داده‌های عمیق
۷۲.....	۵-۳-۳. پارادوکس اطلاعات
۷۵.....	۵-۴. مرغ و تخم‌مرغ راه‌کارهای کلان داده‌ها
۷۷.....	۵-۵. نگاهی بر نقد بصری سازی
۷۸.....	۵-۵-۱. اکوسیستم بصری سازی
۷۹.....	۵-۵-۲. غیرمنطقی بودن نیازها: غذای سریع تا غذای خوب
۸۱.....	۵-۵-۳. افزایش نقد
۸۴.....	۵-۵-۴. حرف‌های پایانی
<b>۸۵.....</b>	<b>فصل ششم: کلان داده، بهداشت و درمان</b>
۸۵.....	۱-۶. حل مسئله Wanamaker برای بهداشت و درمان
۸۷.....	۱-۱-۶. مؤثرتر ساختن مراقبت‌های بهداشتی
۹۱.....	۲-۱-۶. داده‌های بیش‌تر، منابع بیش‌تر
۹۲.....	۳-۱-۶. پرداختن به نتایج
۹۴.....	۴-۱-۶. فعال‌سازی داده
۹۷.....	۵-۱-۶. ساختن سیستم مراقبت سلامتی که می‌خواهیم
۹۸.....	۶-۱-۶. توصیه‌هایی برای خواندن
۹۸.....	۲-۶. دکتر فرزاد مستشری در زمینه زیرساخت اطلاعات سلامتی برای بیماران الکترونیکی مدرن
۱۰۲.....	۳-۶. جان ویلینگس در مورد ریسک‌ها و پاداش‌های مربوط به داده‌های سلامتی صحبت می‌کند
۱۰۸.....	۴-۶. بحث با استر دیسون در خصوص داده‌های سلامتی "مراقبت‌های سلامتی پیش‌گیرانه" و چیز عظیم بعدی
۱۲۲.....	۵-۶. پنج عنصر اصلاحی که ارائه‌دهندگان سلامتی احتمالاً چیزی در مورد آن نشنیده‌اند

در اولین نسخه "Big Data Now"، تیم O'Reilly تولد و توسعه زود هنگام ابزارهای داده و علم داده را دنبال کردند. حال در این نسخه دوم، می‌خواهیم ببینیم زمانی که کلان داده رشد می‌کنند چه اتفاقی رخ می‌دهد: چه چیزی در حال اعمال است، کجا نقشی ایفا می‌کند و پیامدهای - به‌طور یکسان خوب و بد - رشد داده.

نسخه ۲۰۱۰۲، Big Data Now را در پنج بخش سازمان‌دهی کردیم:

**افزایش سرعت با کلان داده** - اطلاعات ضروری در مورد ساختارها و تعاریف کلان داده

**ابزارهای کلان داده، تکنیک‌ها و استراتژی‌ها** - راهنمایی کارشناسانه برای تبدیل نظریه‌های کلان داده به محصولات کلان داده.

**کاربرد کلان داده** - نمونه‌هایی از کلان داده‌ها در عمل، شامل نگاهی به ضعف داده‌ها

**آنچه در کلان داده دنبال می‌شود** - افکار در خصوص نحوه تکامل کلان داده و نقش آن در صنایع و حوزه‌ها.

**کلان داده و مراقبت سلامتی** - یک بخش ویژه برای بررسی احتمالات در زمانی که داده و مراقبت سلامتی با یکدیگر همراه می‌شوند.

علاوه بر Big Data Now، با تحلیلاتی که بر روی O'Reilly Radar انجام می‌دهیم و از طریق Strata coverage و events series، می‌توانید به توسعه‌های اخیر داده دسترسی داشته باشید.

## افزایش سرعت با کلان داده

### ۱-۲. کلان داده چیست؟

کلان داده، داده‌ای با ظرفیت پردازشی بیش‌تر از ظرفیت پردازشی سیستم‌های پایگاه داده متعارف است. داده بیش‌ازحد بزرگ است، خیلی سریع رشد می‌کند، یا در محدودیت‌های معماری پایگاه داده شما نمی‌گنجد. برای دستیابی به منفعتی از این داده، باید روش دیگری برای پردازش آن انتخاب کنید.

شعار فناوری اطلاعات داغ ۲۰۱۲، کلان داده به‌عنوان روش‌های مقرون‌به‌صرفه برای کاهش حجم، سرعت و تنوع داده‌های انبوه در حال ظهور هستند.

در این میان، داده حاوی الگوها و اطلاعات ارزشمندی است که قبلاً به دلیل حجم کار موردنیاز برای استخراج‌شان پنهان بودند. با مشارکت شرکت‌هایی مثل Walmart و Google این توانایی تا حدودی میسر گشته است اما با هزینه بسیار بالا. تجهیزات سخت‌افزاری امروزی، معماری ابر و نرم‌افزارهای منبع باز پردازش داده کلان داده را با حداقل منابع امکان‌پذیر ساخته است. پردازش کلان داده حتی برای گارازهای راه‌انداز کوچکی که می‌توانند زمان سرویس‌دهنده در ابر را با هزینه پایین اجاره کنند نیز میسر شده است.

ارزش کلان داده برای یک سازمان دو نوع است: استفاده تحلیلی و ارائه محصولات جدید. تحلیل داده‌های بزرگ، بیش‌های پنهان موجود در داده (از جمله تأثیر نظیر به نظیر بر مشتریان، آشکار شده با تحلیل تراکنش‌های خریداران، و داده‌های اجتماعی و جغرافیایی) که پردازش آن‌ها بسیار هزینه‌بر بود را آشکار می‌کند. علی‌رغم ماهیت نسبتاً ایستای گزارش‌های از پیش تعیین‌شده، این که بتوانیم هر قلم داده را در یک زمان معقول پردازش کنیم، نیاز مبرم به نمونه‌برداری را حذف و یک رویکرد تحقیقاتی برای داده‌ها باز می‌کند.

راه‌اندازی‌های موفق وب در ده گذشته، نمونه‌های بزرگی از کلان داده به‌کاررفته به‌عنوان یک فعال‌ساز برای محصولات و دستگاه‌های جدید است. به‌عنوان مثال، فیس‌بوک با ترکیب سیگنال‌های بسیاری از واکنش‌های کاربران و دوستان‌شان، توانست تجربیات بسیار شخصی کاربر را کشف و نوع جدیدی از تبلیغات را ایجاد کند. این که بخش زیادی از ایده‌ها و ابزارهای پایه کلان داده از گوگل، یاهو، آمازون و فیس‌بوک پدید آمده است، تصادفی نیست.

ظهور کلان داده در شرکت‌ها یک همتای ضروری را برای آن ارمغان آورده است: چابکی. بهره‌برداری موفق از مقادیر در داده‌های بزرگ نیازمند آزمایش و اکتشاف است. چه محصولات جدیدی ایجاد کنیم یا به دنبال روش‌هایی برای دستیابی به مزیت رقابتی باشیم، شغل نیازمند کنجکاوی و کارآفرینی است.

### ۱-۱-۲. کلان داده چه شکلی هستند؟

به‌عنوان یک اصطلاح کلی، همانند اصطلاح "ابر" که فناوری‌های متنوعی را پوشش می‌دهد، اصطلاح "کلان داده" نیز می‌تواند بسیار مبهم باشد. داده‌های ورودی در سیستم‌های کلان داده می‌توانند برگرفته از شبکه‌های

## ۷ افزایش سرعت با کلان داده

اجتماعی، لاگ‌های سروری وب، حسگرهای جریان ترافیکی، تصاویر ماهواره، جریان‌های پخش صوتی، تعاملات بانکی، MP۳های موسیقی راک، محتوای صفحات وب، اسکن اسناد دولتی، مسیرهای GPS، سنجش از راه دور خودروها، اطلاعات بازار مالی، و غیره باشد. آیا همه این‌ها یک چیز یکسان هستند؟

برای روشن کردن مباحث، معمولاً از سه ۷ (برای حجم<sup>۱</sup>، سرعت<sup>۲</sup>، و تنوع<sup>۳</sup>) برای مشخص کردن جنبه‌های مختلف کلان داده استفاده می‌شود که لنز مفیدی برای مشاهده و درک ماهیت داده‌ها و پلت‌فرم‌های نرم‌افزاری موجود برای بهره‌برداری از آن‌ها هستند. احتمالاً تا حدودی با یکی از این ۷ها سروکار خواهید داشت.

### حجم

مزیت به دست آمده از توانایی پردازش حجم بزرگی از اطلاعات، جذابیت اصلی تجزیه و تحلیل کلان داده است. داشتن داده‌های بیش تر منجر به داشتن مدل‌های بهتر می‌شود: بیت‌های ساده ریاضی می‌توانند بدون دلیل برای حجم-های زیاد داده مؤثر باشند. آیا اگر شما بتوانید با در نظر گرفتن حدوداً ۳۰۰ فاکتور به جای ۶ فاکتور، پیش‌بینی انجام دهید، در این صورت می‌توانید تقاضا را بهتر پیش‌بینی کنید؟ این مسئله، مهم‌ترین چالش در ساختارهای IT متعارف است که نیازمند فضای ذخیره‌سازی مقیاس‌پذیر و یک روش توزیعی برای پرس و جوها است. بسیاری از شرکت‌ها دارای مقدار زیادی داده بایگانی‌شده، احتمالاً در قالب لاگ‌ها هستند اما توانایی پردازش آن‌ها را ندارند.

با فرض این که حجم داده‌ها بیش تر از حدی است که زیرساخت‌های پایگاه داده رابطه‌ای بتواند با آن مقابله کند، گزینه‌های پردازشی به‌طور گسترده به انتخاب بین معماری‌های پردازش موازی انبوه-انبار داده‌ها یا پایگاه داده‌هایی مثل Greenplum و روش‌های مبتنی بر Apache Hadoop تقسیم می‌شوند. این انتخاب اغلب با درجه‌ای که یکی از ۷های دیگر-تنوع- در آن نقش خواهد داشت، اعلان می‌شود. معمولاً، روش‌های انبار داده‌ها شامل طرح‌های از پیش تعیین‌شده، مجموعه داده‌های منظم و به‌آرامی تکامل یافته هستند. از طرفی دیگر، Apache Hadoop هیچ جایی در ساختار داده‌هایی که پردازش می‌کند، ندارد.

در حقیقت، هادوپ پلت‌فرمی برای توزیع مسائل محاسباتی بر روی تعدادی سرویس‌دهنده است و به‌عنوان اولین نرم‌افزار منبع باز منتشرشده، از روش نگاهش کاهش<sup>۴</sup> که گوگل در استفاده از آن پیش‌گام است، برای جمع‌آوری شاخص‌های جست‌وجویش استفاده می‌کند. MapReduce در هادوپ شامل توزیع یک مجموعه داده بر روی

<sup>۱</sup> volume

<sup>۲</sup> velocity

<sup>۳</sup> variety

<sup>۴</sup> MapReduce

چندین سرویس دهنده و انجام عملیاتی بر روی داده است: گام "نگاشت"<sup>۵</sup>. سپس نتایج جزئی مجدداً باهم ترکیب می‌شوند: گام "کاهش"<sup>۶</sup>.

هادوپ برای ذخیره داده از سیستم فایل توزیعی خودش، HDFS که دسترسی به داده را از چندین گره محاسباتی میسر می‌سازد، استفاده می‌کند. یک الگوی استفاده معمول از هادوپ شامل سه گام است:

❖ بارگذاری داده در HDFS

❖ عملیات نگاشت کاهش و

❖ بازیابی نتایج از HDFS

این فرآیند که ماهیتاً یک عمل دسته‌ای است، برای وظایف محاسباتی تحلیلی یا غیرتعاملی مناسب است. به همین دلیل هادوپ خودش یک روش پایگاه داده‌ای یا انبار داده نیست بلکه می‌توان به‌عنوان یک یک مکمل تحلیلی<sup>۷</sup> عمل کند.

یکی از معروف‌ترین کاربران هادوپ، فیس‌بوک است که مدلش از الگوی پیروی می‌کند. پایگاه داده MySQL داده‌های اصلی را ذخیره می‌کند و سپس این داده‌ها در هادوپ منعکس می‌شود که در آن محاسباتی از جمله ایجاد توصیه‌هایی بر اساس علائق دوستان شما انجام می‌شود. سپس فیس بوک نتایج را به‌منظور استفاده در صفحات کاربران به MySQL برمی‌گرداند.

## سرعت

اهمیت سرعت داده‌ها- نرخ افزایشی جاری شدن داده در یک سازمان- مطابق الگویی مشابه الگوی حجم است. مسائلی که قبلاً محدود به بخش‌های صنعتی بودند اکنون در محیط‌های وسیع‌تری وجود دارند. شرکت‌های تخصصی مثل معامله‌گران مالی از سیستم‌هایی که با انتقال سریع داده مقابله می‌کردند، به نفع خود استفاده کردند. حال نوبت ماست.

چرا این طور است؟ عصر اینترنت و تلفن همراه بدین معنی است که شیوه تحویل محصول و سرویس به مشتری به‌طور فزاینده‌ای ابزاری شده است و جریان داده‌ای را به ارائه‌دهنده برمی‌گرداند. خرده‌فروشان آنلاین قادرند با هر کلیک و تعامل، تاریخچه بزرگی از مشتریان را کامپایل کنند: نه فقط فروش‌های نهایی. کسانی که می‌توانند با پیشنهاد خریدهای اضافی، برای رسیدن به مزیت رقابتی، به‌سرعت از اطلاعات استفاده کنند. در عصر گوشی‌های

<sup>۵</sup> map

<sup>۶</sup> reduce

<sup>۷</sup> analytical adjunct



## ۹ افزایش سرعت با کلان داده

هوشمند، از آنجایی که مشتریان به همراه خود منبعی از تصاویر جغرافیایی و داده‌های صوتی را حمل می‌کنند، مجدداً نرخ جریان داده افزایش کرده است.

مسئله فقط سرعت داده‌های ورودی نیست: مثلاً انتقال سریع داده‌ها به فضای ذخیره‌سازی عمده برای پردازش دسته‌ای آتی نیز امکان‌پذیر است. اهمیت در سرعت چرخه بازخورد، گرفتن داده از داده به واسطه یک تصمیم نهفته است. یک IBM تجاری این نکته را بیان می‌کند که اگر تصویر لحظه‌ای مکان ترافیک را نداشتید، نمی‌توانستید از جاده عبور کنید. زمان‌هایی وجود دارد که نمی‌توانید منتظر گزارشی برای اجرا یا تکمیل کرد هادوپ باشید. اصطلاح صنعتی این‌گونه انتقال سریع داده "جریان داده‌ها" یا "پردازش رویداد پیچیده" نام دارد. اصطلاح دوم که در دسته محصولات قبل از پردازش جریانی داده‌ها قرار دارد به ارتباطات گسترده‌تری دست یافته است و به نظر می‌رسد به نفع جریان داده‌ها، کاهش یافته است. دو دلیل اصلی برای در نظر گرفتن پردازش جریانی وجود دارد: اولی زمانی است که داده‌های ورودی خیلی سریع و کامل ذخیره می‌شوند: برای حفظ الزامات ذخیره‌سازی عملی، همانند جریان‌های داده، باید چندین سطح تحلیلی رخ دهد. در پایان مقیاس، Large Hadron Collider در CERN داده‌های خیلی زیادی تولید کرده است که دانشمندان باید اکثریت آن را حذف کنند- با امید این که هیچ چیز مفیدی را از بین نبرند. دلیل دوم در نظر گرفتن پردازش جریانی، جایی است که برنامه به داده‌ها پاسخ سریع می‌دهد. به دلیل ظهور برنامه‌های موبایل و بازی‌های آنلاین، این وضعیت به میزان زیادی مرسوم شده است.

دسته‌بندی محصولات برای مدیریت جریان داده‌ها به محصولات انحصاری از جمله جریان‌های InfoSphere در IBM و lesspolished تقسیم می‌شود و همچنان چارچوب‌های منبع باز نشأت گرفته از صنعت وب نیز وجود دارند: طوفان تویتر و یاهو SF.

همان‌طور که در بالا ذکر گردید، این مسئله فقط در مورد داده‌های ورودی نیست. سرعت خروجی سیستم نیز ممکن است بسیار مهم باشد. هرچه حلقه بازخورد محدودتر باشد، مزیت رقابتی بیش‌تر خواهد بود. نتایج ممکن است به‌طور مستقیم به یک محصول مثل توصیه‌های فیس‌بوک یا داشبوردهای به‌کاررفته برای تصمیم‌گیری بروند که برای افزایش سرعت به‌ویژه در وب نیاز است و توسعه ذخیره‌سازی مقادیر کلیدی و پایگاه داده‌های ستونی را برای بازیابی سریع اطلاعات از قبل محاسبه‌شده، بهینه‌سازی کرده است. این پایگاه داده‌ها بخشی از یک **دسته چتری**<sup>۸</sup> شناخته‌شده به نام NoSQL را تشکیل می‌دهند و زمانی به‌کاررفته می‌روند که مدل‌های رابطه‌ای مناسب نباشند.

---

<sup>۸</sup> umbrella category

## نوع

داده‌ها به ندرت در فرمی بی‌نقص و آماده برای پردازش هستند. یک موضوع رایج در سیستم کلان داده این است که داده‌های منبع متنوع هستند و دارای ساختارهای رابطه‌ای منظم نیستند. این داده‌ها می‌توانند متنی از شبکه‌های اجتماعی، تصویر یا داده خامی از یک منبع حسگر باشند. هیچ کدام برای ادغام شدن در یک برنامه آماده نیستند.

حتی در وب، که ارتباطات کامپیوتر به کامپیوتر باید تضمین‌هایی به همراه داشته باشد، ماهیت داده‌ها کثیف است. مرورگرهای مختلف داده‌های مختلفی ارسال می‌کنند، کاربران داده‌ها را دریغ می‌کنند، آن‌ها ممکن است از نسخه-های نرم‌افزاری مختلف یا از فروشندگان مختلف برای برقراری ارتباط با شما استفاده کنند و می‌توانید شرط ببندید که اگر بخشی از فرآیند توسط انسان انجام شود، خطا و ناسازگاری وجود خواهد داشت.

یک استفاده رایج پردازش داده‌های بزرگ، گرفتن داده‌های غیر ساخت‌یافته و استخراج معانی مرتب، برای مصرف توسط انسان یا به‌عنوان داده‌های ورودی برای یک برنامه است. این چنین مثالی یک قطعه‌نامه موجودیت است، فرآیند تعیین دقیق چیزی که نام به آن اشاره دارد. آیا این‌جا شهر لندن، انگلستان یا لندن تگزاس است؟ تا زمانی که منطق تجاری شما دخیل است، نمی‌خواهید حدس زده شوید.

فرآیند انتقال از داده‌های منبع به داده‌های کاربردی پردازش شده شامل از دست رفتن اطلاعات است. هنگامی که مرتب می‌شوید، پرتاب وسایل را ترک خواهید کرد. این اصل کلان داده است: زمانی که می‌توانید همه چیز را حفظ کنید. ممکن است در بیت‌هایی که دور می‌ریزید، سیگنال‌های مفیدی وجود داشته باشد. اگر داده‌های منبع را از دست بدهید، هیچ بازگشتی وجود ندارد.

علی‌رغم محبوبیت و شناخت ماهیت پایگاه داده‌های رابطه‌ای، این‌طور نیست که داده‌ها همواره باید در مقصد باشند، حتی اگر مرتب‌شده باشند. انواع خاصی از داده‌ها برای کلاس‌های خاصی از پایگاه داده بهتر هستند. به‌عنوان مثال، اسناد کدگذاری شده در XML، زمانی که در یک XML اختصاصی مثل MarkLogic ذخیره شوند، متنوع‌ترین اسناد هستند. روابط شبکه‌های اجتماعی دارای ماهیت گروهی هستند، و پایگاه داده‌های گرافی مثل Neo4J عملیات را ساده‌تر و مؤثرتر بر روی آن‌ها انجام می‌دهند.

حتی در مواردی که هیچ عدم سازگاری در نوع داده رادیکالی وجود ندارد، عیب پایگاه داده رابطه‌ای مربوط به ماهیت ایستای طرح‌هایش است. در یک محیط اکتشافی چابک، نتایج محاسبات با کشف و استخراج سیگنال‌های بیش‌تر تکامل خواهند یافت. پایگاه داده‌های NoSQL نیمه ساخت‌یافته به‌منظور ایجاد انعطاف‌پذیری، این نیازمندی

## افزایش سرعت با کلان داده ۱۱

را بر آورده می‌سازند: آن‌ها ساختار کافی برای سازمان‌دهی داده فراهم می‌کنند، اما نیازی به استخراج طرح داده قبل از ذخیره‌سازی آن ندارند.

### ۲-۱-۲. در عمل

ماهیت کلان داده را کشف و چشم‌انداز کلان داده را از یک سطح بالا بررسی کردیم. به‌طور معمول، هنگام استقرار، ابعادی هستند که باید فرای انتخاب ابزار در نظر گرفته شوند.

### ابری یا در خانه؟

در حال حاضر اکثریت راه‌حل‌های کلان داده در سه فرم ارائه می‌شوند: تنها نرم‌افزار، به‌عنوان یک دستگاه<sup>۹</sup> یا مبتنی بر ابر. تصمیمات در خصوص این که کدام مسیر باید انتخاب شود، به فاکتورهای دیگری مثل محلیت داده، حریم خصوصی و تنظیمات، منابع انسانی و الزامات پروژه بستگی دارد. بسیاری از سازمان‌ها یک راه‌حل ترکیبی را انتخاب می‌کنند: استفاده از منابع ابری مبتنی بر تقاضا برای تکمیل کارها در خانه.

### کلان داده بزرگ است

این یک حقیقت محض است که داده‌ای که برای پردازش عادی بسیار بزرگ باشد، برای انتقال به‌جای دیگر نیز بزرگ است. اولویت‌های IT در حال تغییر هستند: این برنامه است که باید انتقال داده شود، نه داده. اگر بخواهید داده‌های سرشماری ایالت متحده را تحلیل کنید، اجرای کد آن در پلت‌فرم سرویس‌های تحت وب آمازون که میزبان این گونه داده‌های محلی است بسیار ساده‌تر است و دیگر نیازی به صرف زمان یا هزینه برای انتقال آن ندارید. حتی اگر داده برای انتقال خیلی بزرگ نباشد، محلیت می‌تواند مسئله‌ساز باشد، به‌خصوص با این به‌روزرسانی سریع داده‌ها. سیستم‌های معاملات مالی برای به دست آوردن سریع‌ترین اتصال به داده‌های منبع، در مراکز داده تجمع کردند، بنابراین دلیل، تفاوت میلی‌ثانیه‌ای در زمان پردازشی نماد یک مزیت رقابتی است.

### کلان داده کثیف است

همه‌چیز به زیرساخت مربوط نیست. کاربران کلان داده به‌طور پیوسته گزارش می‌دهند که ۸۰ درصد تلاش‌های انجام‌شده برای مقابله با داده‌ها به پاک کردن آن در ابتدا مربوط می‌شود، همان‌طور که Pete Warden در واژه‌نامه کلان داده خود مشاهده نمود: "من احتمالاً زمان بیش‌تری را برای تبدیل داده‌های منبع کثیف به چیزی قابل‌استفاده صرف می‌کنم، نسبت به این که فرآیند تحلیل داده را انجام دهم."

<sup>۹</sup> appliance

به دلیل هزینه بالای جمع‌آوری و پاک‌سازی داده، این که آنچه را که واقعاً نیاز دارید، به خود بسپارید. بازار داده-ها ابزارهایی برای دستیابی به داده‌های مشترک هستند و اغلب می‌توانید به بهبود آن‌ها کمک کنید. کیفیت نیز ممکن است متغیر باشد اما به‌طور فزاینده‌ای یک معیار رقابتی در بازار داده‌ها خواهد بود.

### فرهنگ

پدیده کلان داده به ظهور علم داده‌ها بسیار وابسته است، رشته‌ای که ریاضیات، برنامه‌نویسی، و غریزه علمی را باهم ترکیب می‌کند. بهره‌مندی از ابزارهای داده به معنی سرمایه‌گذاری در تیم‌هایی با این مهارت و اتصال آن‌ها به تمایلات سازمانی به‌منظور درک و استفاده از داده‌ها است.

در این گزارش، "ساختن تیم‌های علوم داده"، D.J. Patil دانشمندان داده‌ای را با معیارهای کیفی زیر مشخص کرده است:

- ❖ **تخصص فنی:** بهترین دانشمندان داده‌ای معمولاً دارای تخصص‌های عمیق در برخی رشته‌های علمی هستند.
- ❖ **کنجکاوی:** تمایل به بررسی عمیق و کشف و تبدیل یک مسئله به یک مجموعه فرضیات واضحی که قابل آزمایش باشد.
- ❖ **داستان‌سرایی:** توانایی استفاده از داده‌ها برای گفتن یک داستان و قادر به برقراری ارتباط مؤثر با آن.
- ❖ **باهوشی:** توانایی برخورد با مسئله به روش‌های خلاق مختلف.

ماهیت گسترده پروژه‌های تجزیه و تحلیل کلان داده می‌تواند دارای جنبه‌های ناراحت‌کننده‌ای باشد: داده باید به چندین سیلو شکسته شود تا قابل استخراج باشد و سازمان باید نحوه برقراری ارتباط و تفسیر نتایج تحلیل‌ها را یاد بگیرد.

مهارت‌های داستان‌سرایی و باهوشی فاکتورهای بسیار مهمی هستند که در نهایت تعیین می‌کنند که آیا مزایای کارهای تحلیلی در یک سازمان جذب می‌شوند یا خیر. هنر و تجربه بصری سازی داده در ایجاد پل بر روی شکاف بین انسان و کامپیوتر به‌منظور ایجاد هدفمند بینش تحلیلی بیش‌ازپیش مهم تر شده است.

### بدانید کجا می‌خواهید بروید

در نهایت، به خاطر داشته باشید که کلان داده، نوشدارو نیست. می‌توانید الگوها و سرنخ‌هایی در داده‌های خود پیدا کنید، اما بعد چی؟ Christer Johnson، سرگروه IBM در تحلیل پیشرفته در آمریکای شمالی توصیه زیر را برای شروع کار تجارت با کلان داده گفته است: اول، تصمیم بگیرید که چه مسئله‌ای را می‌خواهید حل کنید. انتخاب یک مسئله تجاری واقعی، مثلاً چگونه می‌توانید برای افزایش هزینه‌های مشتری، استراتژی تبلیغ خود را تغییر دهید،

## افزایش سرعت با کلان داده ۱۳

در پیاده‌سازی به ما کمک می‌کند. علاوه بر روحیه شرکتی، داشتن یک هدف مشخص نیز به شدت برای کلان داده مفید است.

## ۲-۲. Apache Hadoop چیست؟

Apache Hadoop نیروی محرک رشد صنعت کلان داده بوده است که به همراه فناوری‌های مربوطه مثل Hive و Pig اغلب خواهید شنید. اما چه کاری انجام می‌دهد و چرا به همه استراتژی‌هایش مثل Oozie و Zookeeper و Flume نیاز دارید؟

Hadoop توانایی پردازش ارزان حجم زیادی از داده را بدون در نظر گرفتن ساختارش به ارمغان می‌آورد. منظور ما از حجم زیاد، ۱۰ تا ۱۰۰ گیگابایت و بالاتر است. تفاوت این توانایی پردازش نسبت به گذشته چیست؟ انبار داده‌های سازمانی موجود و پایگاه داده‌های رابطه‌ای در پردازش داده‌های ساخت یافته پیشرفته هستند و می‌توانند حجم عظیمی از اطلاعات را ذخیره کنند، اما هزینه‌ای نیز دارند: این الزام برای ساختار، انواع داده‌هایی که می‌توانند پردازش شوند را محدود می‌کند و منجر به یک حالت سکون می‌شود که باعث می‌شود انبار داده‌ها برای اکتشاف چابک داده‌های ناهمگن کلان مناسب نباشد. میزان تلاش موردنیاز برای انبار داده اغلب بدین معنی است که منابع داده‌ای با ارزش درون سازمان‌ها هرگز استخراج نمی‌شوند. این‌جا است که Hadoop می‌تواند تفاوت بزرگی ایجاد کند.

این تحقیق مؤلفه‌های اکوسیستم هادوپ را بررسی و توابع هر کدام را شرح می‌دهد.

### ۱-۲-۲. هسته هادوپ: نگاهت کاهش

در پاسخ به مسئله ایجاد شاخص‌های جست‌وجوی تحت وب، چارچوب نگاهت کاهش ایجاد شده در گوگل یک نیروگاه قوی برای اکثر سیستم‌های پردازش کلان داده امروزی است. علاوه بر هادوپ، نگاهت کاهش نیز در پایگاه داده‌های MPP و NoSQL از جمله Vertica یا MongoDB مشاهده می‌شود.

نوآوری مهم نگاهت کاهش، قابلیت گرفتن یک پرس‌وجو در یک مجموعه داده، تقسیم آن، و سپس اجرای موازی آن در چندین گره است. توزیع محاسبات مسئله بزرگ بودن داده به‌منظور قرار گرفتن بر روی یک ماشین را حل می‌کند. با ترکیب این تکنیک با سرویس‌دهنده‌های لینوکس یک روش مقرون‌به‌صرفه برای آرایه‌های محاسباتی عظیم به دست می‌آید.

در هسته، هادوپ یک پیاده‌سازی نگاشت کاهش منبع باز است. مطابق گزارش یاهو، هادوپ در سال ۲۰۰۶ پدیدار گشت، و به گفته سازنده آن، Doug Cutting در سال ۲۰۰۸ به قابلیت "مقیاس وی" دست یافت. با تکامل پروژه هادوپ، مؤلفه‌های بیش‌تری برای بهبود قابلیت استفاده و کارکرد به آن اضافه گردید. نام "هادوپ" برای نمایش کل این اکوسیستم به وجود آمده است که به موازات ظهور لینوکس بوده است: این نام شدیداً به هسته لینوکس اشاره دارد، اما به‌عنوان نام سیستم‌عامل کامل نیز پذیرفته شده است.

### ۲-۲-۲. لایه‌های پایین‌تر هادوپ: HDFS و نگاشت کاهش

در بالا در مورد توانایی نگاشت کاهش در توزیع محاسبات بر روی چند سرویس‌دهنده صحبت کردیم. برای انجام محاسبات، هر سرویس‌دهنده باید به داده‌ها دسترسی داشته باشد. این نقش HDFS، به‌عنوان سیستم فایل توزیع‌شده هادوپ است.

HDFS و نگاشت کاهش قوی هستند. سرویس‌دهنده‌ها در خوشه هادوپ ممکن است شکست بخورند و فرآیند محاسبات را قطع نکنند. HDFS تضمین می‌کند که داده‌ها به‌واسطه افزونگی در خوشه، تکرار می‌شوند. با تکمیل محاسبات، گره نتایج خود را در HDFS می‌نویسد.

هیچ محدودیتی برای داده‌هایی که توسط HDFS ذخیره می‌شوند، وجود ندارد. داده‌ها می‌توانند بدون ساختار و بدون طرح باشند. در عوض، پایگاه داده‌های رابطه‌ای نیاز دارند تا داده‌ها ساخت یافته باشد و طرح‌ها قبل از ذخیره‌سازی داده تعریف شده باشند. در HDFS، ارائه‌دهنده کد مسئول معنی‌دار بودن داده‌ها است. برنامه‌نویسی هادوپ در سطح نگاشت کاهش یک نمونه کار با API‌های جاوا است که در آن فایل داده‌ها به‌صورت دستی در HDFS بارگذاری می‌شود.

### ۲-۲-۳. بهبود قابلیت برنامه‌نویسی: Pig و Hive

کار کردن مستقیم با API‌های جاوا می‌تواند خسته‌کننده و مستعد خطا باشد و همچنین استفاده از هادوپ را به برنامه‌نویسان جاوا محدود می‌کند. هادوپ دو راهکار برای برنامه‌نویسی راحت‌تر هادوپ ارائه کرده است.

❖ Pig یک زبان برنامه‌نویسی است که وظایف رایج در کار کردن با هادوپ را ساده می‌کند: بارگذاری داده، نمایش تبدیلات انجام‌شده بر روی داده‌ها و ذخیره‌سازی نتایج نهایی. عملیات داخلی Pig می‌توانند با داده‌های نیمه ساخت یافته مثل فایل‌های لاگ سروکار داشته باشند و این زبان با استفاده از جاوا قابل بسط است تا بتواند انواع داده‌های سفارشی و تبدیلات را پشتیبانی کند.

## افزایش سرعت با کلان داده ۱۵

❖ Hive به هادوپ قابلیت کار کردن به عنوان یک انبار داده را می‌دهد و ساختار را به سوی داده‌ها در HDFS سوق می‌دهد و سپس اجرای پرس‌وجو بر روی داده را با استفاده از نحوی همانند نحو SQL میسر می‌سازد. همانند Pig، قابلیت‌های هسته Hive نیز قابل گسترش هستند.

انتخاب بین Hive و Pig ممکن است گیج‌کننده باشد. Hive برای امور ذخیره‌سازی داده‌های با ساختار ایستا و نیازمند تجزیه و تحلیل مکرر مناسب‌تر است. نزدیک بودن Hive به SQL باعث شده است تا گزینه ایده آلی برای ادغام هادوپ با دیگر ابزارهای هوشمند تجاری باشد.

Pig به توسعه‌دهنده چابکی بیش‌تری برای کشف مجموعه داده‌های بزرگ، توسعه اسکریپت‌های موجز برای انتقال جریان داده به منظور پیوستن به برنامه‌های بزرگ را می‌دهد. Pig یک لایه نازک‌تر نسبت به Hive در هادوپ است و مزیت اصلی آن کاهش چشم‌گیر حجم کد موردنیاز در مقایسه با استفاده مستقیم از API‌های جاوا در هادوپ است. بدین ترتیب، مخاطب موردنظر Pig در درجه اول توسعه‌دهنده نرم‌افزار است.

### ۴-۲-۲. بهبود دسترسی به داده: HBase، Sqoop و Flume

هادوپ اساساً یک سیستم دسته‌گرا است. داده‌ها در HDFS بارگذاری، پردازش و سپس بازیابی می‌شوند. این نوعی محاسبات بازگشتی و اغلب به دسترسی تعاملی و تصادفی به داده‌ها نیاز است.

HBase یک پایگاه داده ستون‌گرایی است که در بالای HDFS اجرا می‌شود. هدف این پروژه که پس از BigTable گوگل مدل‌سازی شده است، میزبانی میلیاردها سطر داده برای دسترسی سریع است. نداشتن کاهش می‌تواند از HBase هم به عنوان یک منبع و هم به عنوان یک مقصد برای محاسباتش استفاده کند، و از Hive و Pig نیز می‌توان به صورت ترکیبی با HBase استفاده نمود.

HBase برای دسترسی تصادفی به داده محدودیت‌هایی را در نظر گرفته است: کارایی Hive با HBase، ۴ تا ۵ مرتبه کندتر از HDFS است و حداکثر حجم داده‌ای که می‌توان در HBase ذخیره نمود، تقریباً در حد پتابایت است در حالی که حداکثر حجم ذخیره‌سازی HDFS بیش از ۳۰PB است.

HBase برای تجزیه و تحلیل داده‌ها مناسب نیست و بیش‌تر برای ادغام کلان داده به عنوان بخشی از یک برنامه بزرگ مناسب هستند. موارد استفاده شامل ورود، شمارش، و ذخیره‌سازی داده‌های سری زمانی هستند.

## تاریخچه هادوپ

توسعه، پیکربندی و نظارت	Ambari
جمع آوری و ورود لاگ و داده‌های رویدادی	Flume
تغییر مقیاس پایگاه داده ستون‌گرا به میلیاردها سطر	HBase
اشتراک طرح و نوع داده در Pig، Hive و نگاشت کاهش	HCatalog
فایل سیستم افزونه توزیع شده برای هادوپ	HDFS
دسترسی به انبار داده‌ها همانند SQL	Hive
کتابخانه‌ای از الگوریتم‌های یادگیری ماشین و داده کاوی	Mahout
محاسبات موازی در خوشه‌های سروری	MapReduce
زبان برنامه‌نویسی سطح بالا برای محاسبات هادوپ	Pig
تنظیم (ارکستراسیون) و مدیریت گردش کار	Oozie
ورود داده از پایگاه داده‌های رابطه‌ای	Sqoop
استقرار ابری-آگنوستیک خوشه‌ها	Whirr
مدیریت پیکربندی و هماهنگ‌سازی	Zookeeper

## ورود و خروج داده

بهبود قابلیت همکاری با مابقی جهان داده‌ها با Sqoop و Flume میسر گردید. Sqoop یک ابزار طراحی شده برای ورود داده از پایگاه داده‌های رابطه‌ای به هادوپ با ورود مستقیم به HDFS یا Hive است. Flume برای ورود مستقیم جریانی از لاگ‌ها به HDFS طراحی شده است.

دوستانه بودن SQL در Hive بدین معنی است که می‌توان از آن به‌عنوان یک نقطه ادغام با جهان وسیعی از ابزارهای پایگاه داده که قابلیت برقراری ارتباط از طریق درایورهای پایگاه داده JDBC یا ODBC را دارند، استفاده نمود.

## ۵-۲-۲. هماهنگی و گردش کار: Oozie و Zookeeper

با رشد مجموعه سرویس‌های در حال اجرا به‌عنوان بخشی از یک خوشه هادوپ، به سرویس‌های هماهنگ‌سازی و نام‌گذاری نیاز است. با توجه به این که آمادورفت گره‌های محاسباتی امکان‌پذیر است، اعضای خوشه باید با یکدیگر همگام باشند، باید بدانند کجا به سرویس‌ها دسترسی دارند و نحوه پیکربندی خود را نیز باید بدانند که اهداف Zookeeper است.



## افزایش سرعت با کلان داده ۱۷

سیستم‌های تولیدی که از هادوپ استفاده می‌کنند اغلب دارای خط لوله‌های پیچیده‌ای از تبدیلات وابسته به یکدیگر هستند. به‌عنوان مثال، ورود یک دسته جدید از داده‌ها منجر به ورودی‌هایی می‌شوند که سپس منجر به محاسبات جدیدی در مجموعه داده‌های وابسته می‌شوند. مؤلفه Oozie ویژگی‌هایی برای مدیریت جریان کاری و وابستگی‌ها فراهم و نیاز به توسعه‌دهندگان برای کد کردن راه‌حل‌های سفارشی را از بین می‌برد.

### ۶-۲-۲. مدیریت و استقرار: Whirr و Ambari

یکی از ویژگی‌های رایج افزوده‌شده به هادوپ از طریق توزیع‌کنندگانی مثل IBM و مایکروسافت، نظارت و مدیریت است. در مراحل اولیه، هدف Ambari افزودن این ویژگی‌ها به پروژه اصلی هادوپ است. Ambari برای کمک به مدیران سیستم به‌منظور راه‌اندازی و پیکربندی هادوپ، ارتقاء خوشه‌ها و نظارت بر خدمات در نظر گرفته شده است و از طریق API می‌تواند با دیگر ابزارهای مدیریت سیستم ادغام شود.

Whirr به‌عنوان یک بخش نه‌چندان محدود، یک مؤلفه شدیداً مکمل است که روشی برای اجرای سرویس‌ها شامل هادوپ بر روی پلت‌فرم‌های ابری را ارائه می‌نماید. Whirr یک ابر خنثی است و در حال حاضر از سرویس‌های Amazon EC2 و Rackspace پشتیبانی می‌کند.

### ۷-۲-۲. یادگیری ماشین: Mahout

داده‌های هر سازمانی مختلف و متناسب با نیازهایشان است. اما تنوع تحلیل‌های انجام‌شده بر روی این داده‌ها بسیار کم است. پروژه Mahout کتابخانه‌ای از پیاده‌سازی‌های محاسبات تحلیلی رایج توسط هادوپ است و موارد استفاده آن شامل فیلتر مشارکتی کاربر، توصیه‌های کاربری، خوشه‌بندی و دسته‌بندی است.

### ۸-۲-۲. کاربرد هادوپ

به‌طورمعمول، از هادوپ به شکل توزیع‌شده استفاده می‌شود. همانند لینوکس که قبل از هادوپ بود، فروشندگان مؤلفه‌های اکوسیستم Apache Hadoop را باهم ادغام و ابزارها و ویژگی‌های مدیریتی را در آن گنجانده‌اند. اگرچه نصب ابری و مدیریت‌شده نگاهش کاهش هادوپ به‌خودی‌خود توزیع‌شده نیست، اما در سرویس نگاهش-کاهش کشسان آمازون در دسترس است.

### چرا کلان داده، بزرگ است: سیستم عصبی دیجیتال

کل داده‌های موجود در "کلان داده" از کجا آمده‌اند؟ و چرا کلان داده فقط نگرانی شرکت‌هایی مثل فیس‌بوک و گوگل نیست؟ پاسخ این است که شرکت‌های تحت وب پیش‌گام این عرصه هستند. به‌وسیله شبکه‌های اجتماعی،

تلفن همراه و فناوری ابری، یک انتقال مهم رخ داده است که همه ما را به سمت جهان داده‌ها سوق داده و این شرکت‌ها امروزه ساکن هستند.

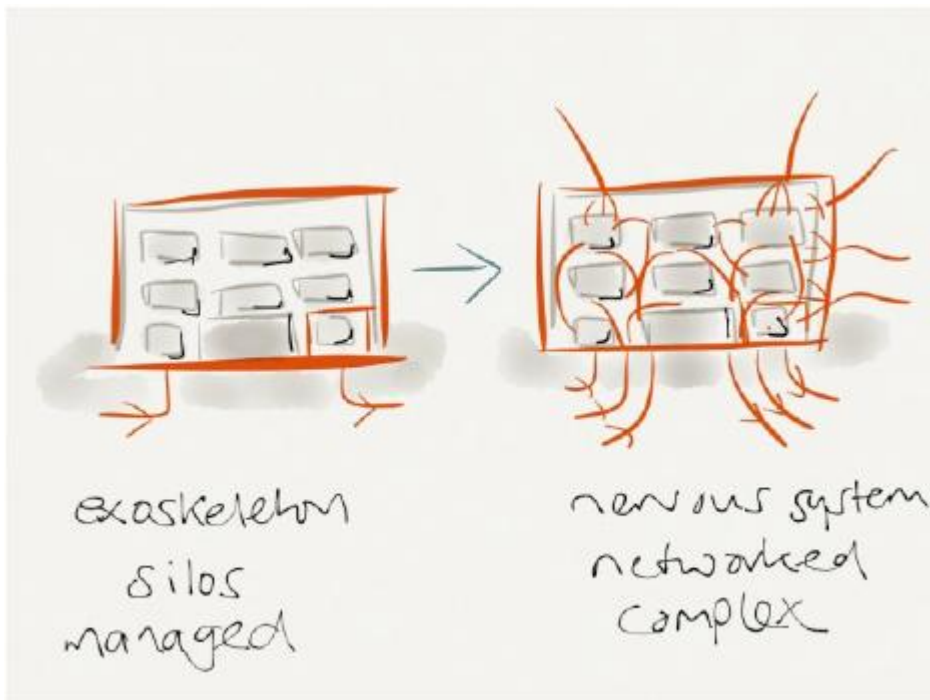
### از اسکلت خارجی به یک سیستم عصبی

تا چند سال پیش، کاربرد اصلی سیستم‌های کامپیوتری در جامعه و به‌خصوص تجارت، به‌عنوان یک سیستم پشتیبانی دیجیتال بوده است. برنامه‌های دیجیتالی شده که در دنیای واقعی وجود دارند، شامل برنامه واژه پردازش، برنامه حقوق و دستمزد و برنامه موجودی هستند. این سیستم‌ها از طریق، فروشگاه‌ها، مردم، تلفن، حمل و نقل و غیره رابط‌هایی با دنیای واقعی برقرار کردند. عبارت جدید "اداره بدون کاغذ" به انتقال پردازش‌های کاغذی پیشین به کامپیوتر اشاره دارد. این سیستم‌های کامپیوتری تشکیل یک اسکلت دیجیتال دادند و از تجارت در جهان واقعی پشتیبانی می‌کنند.

ورود اینترنت و وب، منجر به افزودن ابعاد جدیدی شده و عصر تجارت کاملاً دیجیتال را به ارمغان آورده است. تعامل با کاربران، پرداخت‌ها و اغلب تحویل محصول توسط سیستم‌های کامپیوتری به‌طور کامل قابل انجام است. داده فقط درون اسکلت باقی نمی‌ماند، بلکه یکی از عناصر کلیدی در عملیات است. در دورانی هستیم که تجارت و جامعه به یک سیستم عصبی دیجیتال دستیابی دارند.

همان‌طور که طرح زیر نشان می‌دهد، یک سازمان دارای سیستم عصبی دیجیتال با تعداد زیادی از جریان‌های ورودی و خروجی داده و یک شبکه سطح بالا مشخص می‌شود که هر دو به‌صورت داخلی و خارجی منجر به افزایش جریان داده و در نتیجه افزایش پیچیدگی می‌شوند.

این انتقال نماد مهم بودن کلان داده است. تکنیک‌های ارائه‌شده برای مقابله با اطلاعات به‌هم‌پیوسته و ناهمگن به‌دست‌آمده از شرکت‌های تحت وب گسترده، ابزارهای اصلی ما برای انتقال به عملیات دیجیتال بومی خواهند بود. نمونه‌های اولیه این فناوری شامل کشف تقلب در معاملات مالی و اشکال‌زدایی و بهبود روند استخدام در HR است: و در حال حاضر تقریباً هر کسی به جریان عظیم اطلاعات در شبکه‌های اجتماعی توجه دارد.



### نمودار انتقال

با گسترش فناوری در تجارت، هر قدم برداشته شده منجر به جهشی در حجم داده شده است. سؤال منطقی افراد درگیر با کلان داده این است که بپرسند چرا زمانی که تجارت آن‌ها گوگل یا فیس بوک نیست، کلان داده به آن‌ها اعمال می‌شود؟

پاسخ این سؤال به توانایی تجارت‌های اینترنتی در انجام ۱۰۰ درصد آنلاین فعالیت‌های‌شان بستگی دارد. سیستم عصبی دیجیتال آن‌ها به راحتی از ابتدا تا انتهای عملیات گسترش می‌یابد. اگر کارخانه، فروشگاه یا بخش‌های دیگری از جهان واقعی را در تجارت خود دارید، بیش‌تر درگیر ترکیب آن‌ها با سیستم عصبی دیجیتال هستید. اما "بیش‌تر" بدین معنی نیست که اتفاق نخواهد افتاد. موتور وب، رسانه‌های اجتماعی، تلفن همراه و ابر تجارت‌های بیش‌تری را به جهان مبتنی بر داده هدایت می‌کند. در انگلستان، سرویس دیجیتال دولتی تحویل سرویس به شهروندان را یکپارچه‌سازی می‌کند. نتایج شامل بهبود تجربیات شهروندان است و برای اولین بار بسیاری از ادارات توانستند تصویر واقعی نحوه انجام آن‌ها را مشاهده کنند. در هر خرده‌فروشی، شرکت‌هایی نظیر American، Square،