

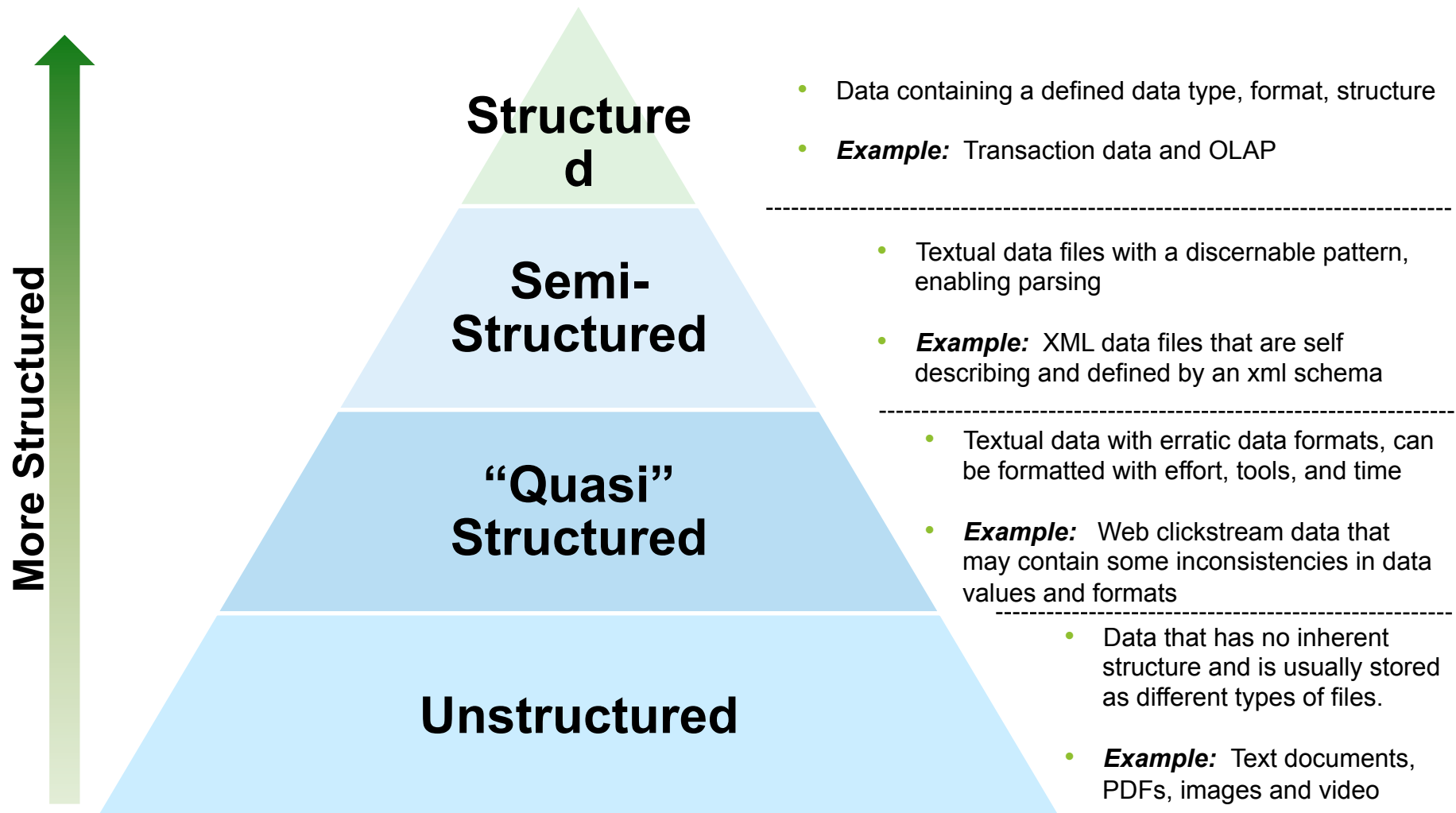
# Big Data Defined

- ***“Big Data” is data whose scale, distribution, diversity, and/or timeliness require the use of new technical architectures and analytics to enable insights that unlock new sources of business value.***
  - ▶ Requires new data architectures, analytic sandboxes
  - ▶ New tools
  - ▶ New analytical methods
  - ▶ Integrating multiple skills into new role of data scientist
- Organizations are deriving business benefit from analyzing ever larger and more complex data sets that increasingly require real-time or near-real time capabilities

Source: McKinsey May 2011 article *Big Data: The next frontier for innovation, competition, and productivity*

# Big Data Characteristics: Data Structures

## Data Growth is Increasingly Unstructured



# Four Main Types of Data Structures

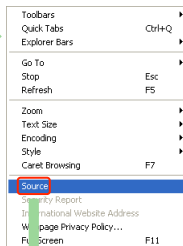
## Structured Data

SUMMER FOOD SERVICE PROGRAM 1)				
(Data as of August 01, 2011)				
Fiscal Year	Number of Sites	Peak (July) Participation	Meals Served	Total Federal Expenditures 2)
	-----Thousands-----		--Mil--	--Million \$--
1969	1.2	99	2.2	0.3
1970	1.9	227	8.2	1.8
1971	3.2	569	29.0	8.2
1972	6.5	1,080	73.5	21.9
1973	11.2	1,437	65.4	26.6
1974	10.6	1,403	63.6	33.6
1975	12.0	1,785	84.3	50.3
1976	16.0	2,453	104.8	73.4
TQ 3)	22.4	3,455	198.0	88.9
1977	23.7	2,791	170.4	114.4
1978	22.4	2,333	120.3	100.3
1979	23.0	2,126	121.8	108.6
1980	21.6	1,922	108.2	110.1

## Semi-Structured Data



View -> Source

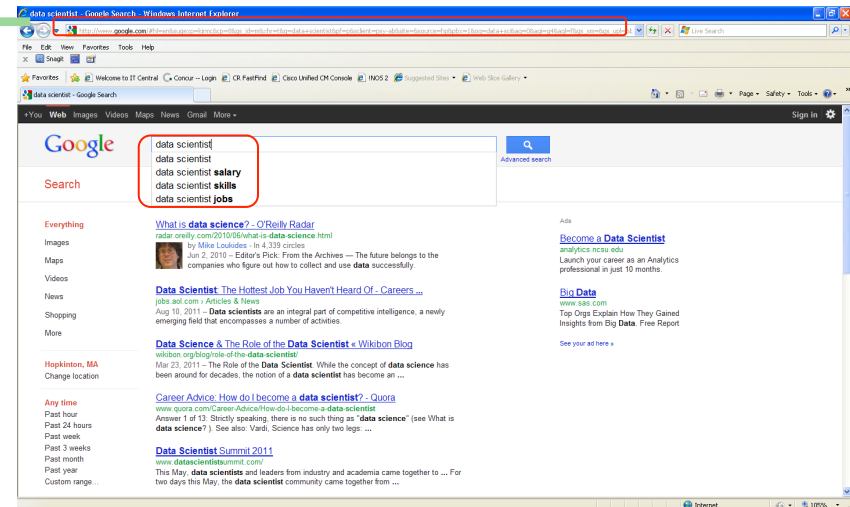


```

1  <!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.0 Transitional//EN" "http://www.w3.org/TR/xhtml1/DTD/xhtml1-trans
2  <html xmlns="http://www.w3.org/1999/xhtml">
3
4      <head>
5          <meta http-equiv="Content-Type" content="text/html; charset=UTF-8" />
6          <META name="y_key" content="859b402e1c9acc">
7          <link rel="canonical" href="http://www.emc.com/index.htm" />
8          <META NAME="verify-v1" CONTENT="yItz9VOP4eV0jFdiPeVfRP2g4qtWFE0IUvVMfSU
9          <title>EMC - Data Recovery, Cloud Computing, and Storage Hardware</title>
10         <META NAME="description" CONTENT="EMC is a leading provider of storage hardware solutions th
11         data recovery and improve cloud computing." />
12         <META NAME="keywords" CONTENT="emc,network storage,data recovery,information management
13         software,nas storage,information protection,information management" />
14         <!-- Start :stylesheet includes -->
15         <link rel="stylesheet" href="/_admin/css/styles.css" />
16         <link rel="stylesheet" href="/_admin/css/styles_nav.css" />
17         <!--(if IE)>

```

## Quasi-Structured Data

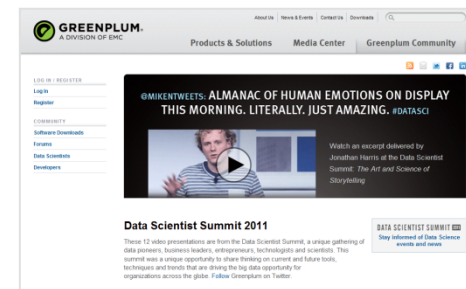


http://www.google.com/#hl=en&sugexp=kjrmc&cp=8&gs\_id=2m&xhr=t&q=data+scientist&pq=big+data&pf=p&scient=psyb&source=hp&pbx=1&oq=data+sci&aq=0&aqi=g4&aql=f&gs\_sm=&gs\_upl=&bav=on.2,or\_r\_gc\_r\_pw,.cf.osb&fp=d566e0fbd09c8604&biw=1382&bih=651

## Unstructured Data

The Red Wheelbarrow, by William Carlos Williams

so much depends  
upon  
  
a red wheel  
barrow  
  
glazed with rain  
water  
  
beside the white  
chickens.



# Data Repositories, An Analyst Perspective

## Data Islands “Spreadmarts”

*Isolated data marts*



- Spreadsheets and low-volume DB's for recordkeeping
- Analyst dependent on data extracts

## Data Warehouses

*Centralized data containers  
in a purpose-built space*



- Supports BI and reporting, but restricts robust analyses
- Analyst dependent on IT & DBAs for data access and schema changes
- Analysts must spend significant time to get extracts from multiple sources

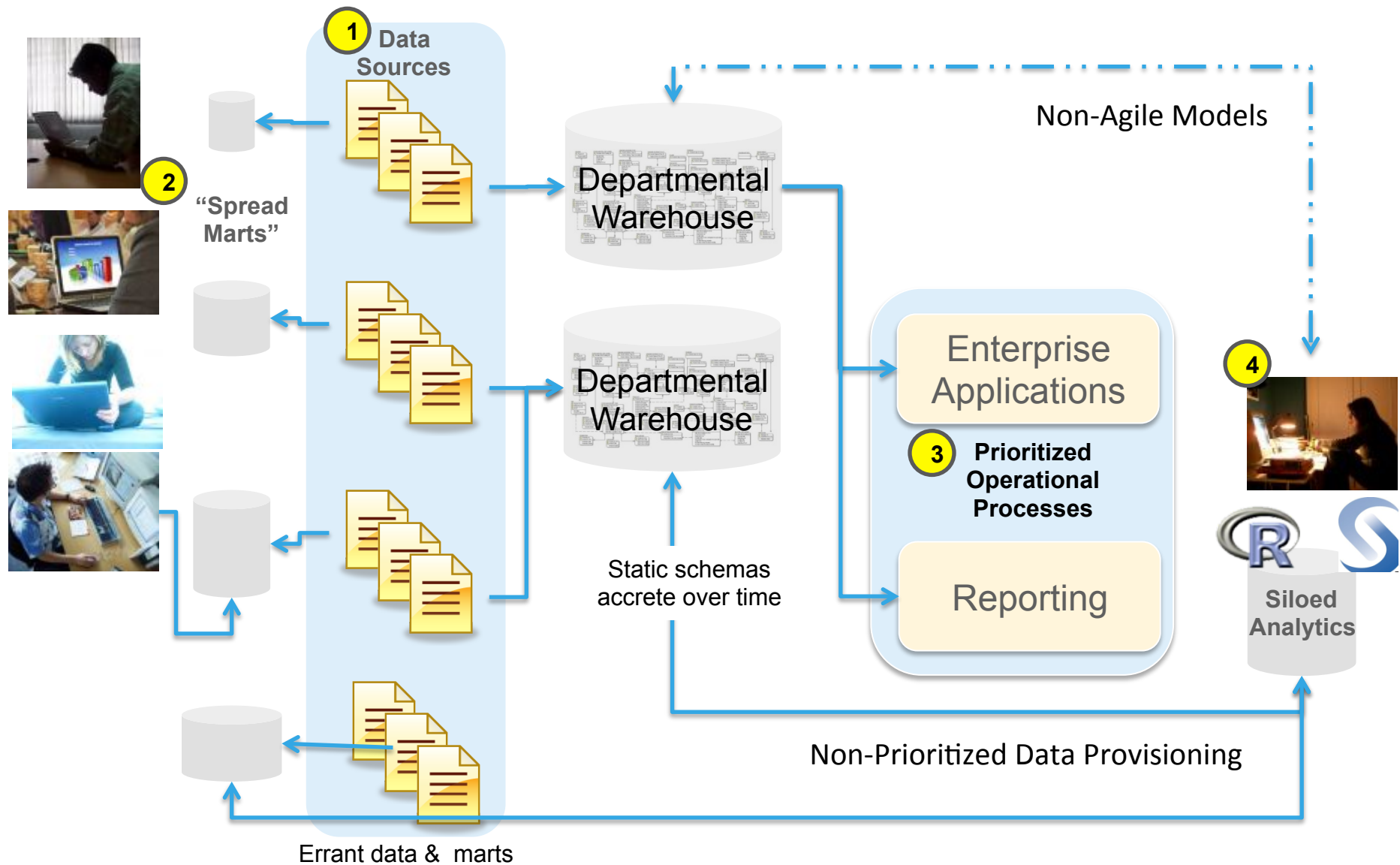
## Analytic Sandbox

*Data assets gathered from multiple  
sources and technologies for analysis*



- Enables high performance analytics using in-db processing
- Reduces costs associated with data replication into "shadow" file systems
- “Analyst-owned” rather than “DBA owned”

# A Typical Analytical Architecture



# Implications of Typical Architecture for Data Science

- High-value data is hard to reach and leverage
- Predictive analytics & data mining activities are last in line for data
  - ▶ Queued after prioritized operational processes
- Data is moving in batches from EDW to local analytical tools
  - ▶ In-memory analytics (such as R, SAS, SPSS, Excel)
  - ▶ Sampling can skew model accuracy
- Isolated, *ad hoc* analytic projects, rather than centrally-managed harnessing of analytics
  - ▶ Non-standardized initiatives
  - ▶ Frequently, not aligned with corporate business goals

Slow  
“time-to-insight”  
&  
reduced  
business impact

# Considerations for Big Data Analytics

## Criteria for Big Data Projects

1. Speed of decision making
2. Throughput
3. Analysis flexibility

## New Analytic Architecture

### Analytic Sandbox

*Data assets gathered from multiple sources and technologies for analysis*



- Enables high performance analytics using in-db processing
- Reduces costs associated with data replication into "shadow" file systems
- "Analyst-owned" rather than "DBA owned"