**Session 11413**
# Issues in Big Data: Analytics

**Tom Deutsch, tdeutsch@us.ibm.com**
Program Director, Big Data

**Bob Foyle, bfoyle@us.ibm.com**
Sr. Product Manager IBM Content Analytics with Enterprise Search

# Abstract

- OK, so you've got all the data lying around and you've decided to make it work for you. Now, what you want to know is how to make sense of it. How do you do that? What do you get? Then, what do you do with it once you've gotten it? How do you know if the analysis is to be trusted?

- The speakers will address those issues, present IBM best practices, and discuss some real-world use cases.

# Agenda

**1** Big Data – Definitions & A Quick Story

**2** Technology Primer
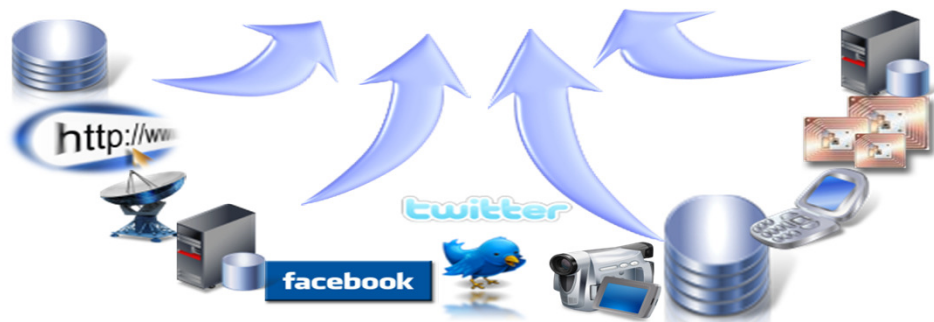
**3** Dive on Content Analytics

**4** Thoughts On Getting Started

# Conventional Definition of "Big Data"

- Never before possible
- Solely defined by large volumes
- Only Unstructured Data
- Valuable insight, but difficult to extract
- Basically an ETL environment

## This definition is <u>wrong</u>

# Functional Big Data Definition

**Dealing with** information management **challenges** that don't natively fit with traditional approaches to handling the problem

# Functional Big Data Definition, Part 2

The technologies that deal with these problems are **broad and diverse, it is <u>not</u> just Hadoop**

## Giving Rise to "Fit For Purpose Architectures"

The introduction of **purpose-built technology** focused on a **specific computing problem** that is compelling better than existing technologies

# Giving Rise to "Fit For Purpose Architectures"

## Match the compute problem to the best way to handle it **rather than assuming SQL by default**

# …Which In-turn Enables Paradigm Shifts*

- **The Idea Of the Super-set**
- **Move from Sampling to "Absolute Knowledge"**
- **Combining Structured and Unstructured Analytics**
- **Importance of Augmented Decision Making**
- Streaming Computing Paradigm
- Lowering The Cost of Experimentation
- Changing Role Of Archiving

(* but the Laws of Gravity Still Apply!)

# Agenda

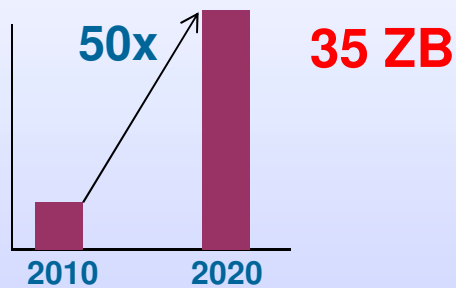| 1 | **Big Data – Definitions & A Quick Story** |
|---|---|
| 2 | **Technology Primer** |
| 3 | **Dive on Content Analytics** |
| 4 | **Thoughts On Getting Started** |

# The characteristics of big data

**Cost efficiently processing the growing Volume**

50x

35 ZB

2010  2020

**Responding to the increasing Velocity**

**30 Billion**
RFID sensors and counting

**Collectively Analyzing the broadening Variety**

**80%** of the worlds data is unstructured

**Establishing the Veracity of big data sources**

**1 in 3** business leaders don't trust the information they use to make decisions

# Hadoop is Well Suited for Handling Certain Types of Big Data Challenges

Analyzing larger volumes may provide better results

Deriving new insights

from combinations of data types

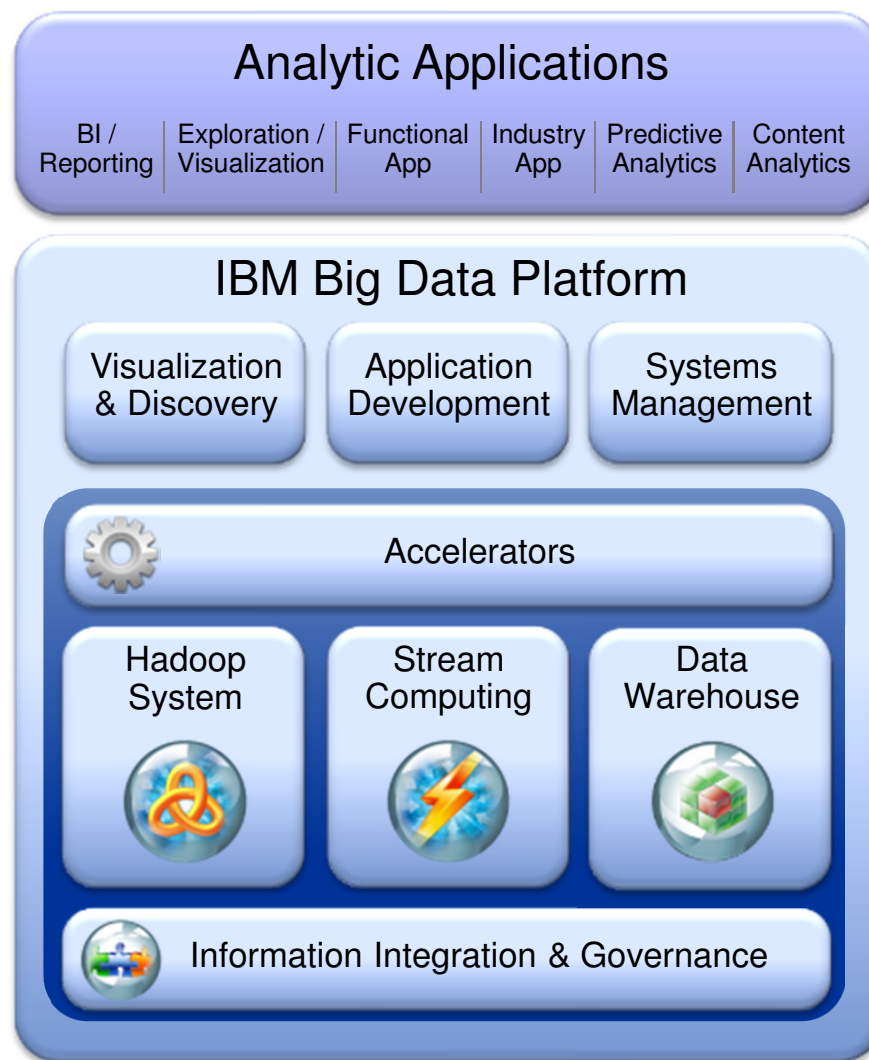Larger data volumes are cost prohibitive with existing technology

Exploring data –

a sandbox for ad-hoc analytics

# IBM Big Data Strategy: Move the Analytics Closer to the Data

**New analytic applications drive the requirements for a big data platform**

- **Integrate and manage the full variety, velocity and volume of data**
- **Apply advanced analytics to information in its native form**
- **Visualize all available data for ad-hoc analysis**
- **Development environment for building new analytic applications**
- **Workload optimization and scheduling**
- **Security and Governance**

## Analytic Applications

| BI / Reporting | Exploration / Visualization | Functional App | Industry App | Predictive Analytics | Content Analytics |
|---|---|---|---|---|---|

## IBM Big Data Platform

| Visualization & Discovery | Application Development | Systems Management |
|---|---|---|

Accelerators

| Hadoop System | Stream Computing | Data Warehouse |
|---|---|---|

Information Integration & Governance

# InfoSphere BigInsights Brings Hadoop to the Enterprise

- **Manages a wide variety and huge volume of data**

- **Augments open source Hadoop with enterprise capabilities**

  - Performance Optimization

  - Development tooling

  - Enterprise integration

  - Analytic Accelerators

  - Application and industry accelerators

  - Visualization

  - Security

- **Provides Enterprise Grade Hadoop analytics**
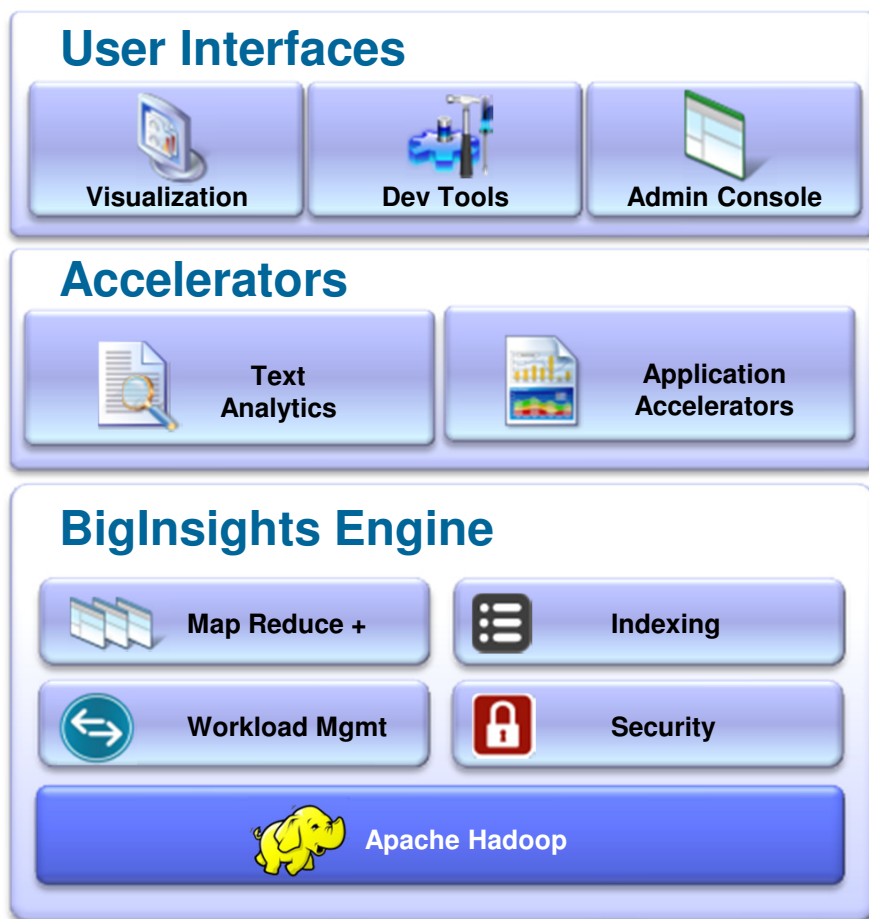
# IBM Significantly Enhances Hadoop



**IBM Innovation**

- **Scalable**
  - New nodes can be added on the fly.

- **Affordable**
  - Massively parallel computing on commodity servers

- **Flexible**
  - Hadoop is schema-less, and can absorb any type of data.

- **Fault Tolerant**
  - Through MapReduce software framework

+

- **Performance & reliability**
  - Adaptive MapReduce, Compression, Indexing, Flexible Scheduler

- **Analytic Accelerators**

- **Productivity Accelerators**
  - Web-based UIs
  - Tools to leverage existing skills
  - End-user visualization

- **Enterprise Integration**
  - To extend & enrich your information supply chain.

IBM

# InfoSphere BigInsights – A Closer Look

## User Interfaces

| Visualization | Dev Tools | Admin Console |

## Accelerators

| Text Analytics | Application Accelerators |

## BigInsights Engine

| Map Reduce + | Indexing |
| Workload Mgmt | Security |

**Apache Hadoop**

## Integration

**Databases**

**Content Management**

**Information Governance**

## *More Than Hadoop*

- **Performance & workload optimizations**

- **Unique text analytic engines**

- **Spreadsheet-style visualization for data discovery & exploration**

- **Built-in IDE & admin consoles**

- **Enterprise-class security**

- **High-speed connectors to integration with other systems**

- **Analytical accelerators**

# The Only Platform to Support Multiple Hadoop Distributions

**User Interfaces**

**Accelerators**

**BigInsights Engine**

Map Reduce +

Indexing

Workload Mgmt

Security

**IBM** tested & supported open source components

**cloudera**
Distribution of Hadoop open source components

*future*

Integration

- **Provides a rich set of big data analytics and accelerators on top of open source**

- **Delivers a comprehensive big data platform on top of open source, that addresses all big data requirements.**

# Agenda

**1** Big Data – Definitions & A Quick Story

**2** Technology Primer

**3** Dive on Content Analytics

**4** Thoughts On Getting Started

# Text Analytics is the basis for Content Analytics

**What is Text Analytics?**

*Text Analytics* (NLP*) describes a set of linguistic, statistical, and machine learning techniques that allow text to be analyzed and key information extraction for business integration.



PC 143 (Hunter)
15 June 2006 23:47
Suspect identified himself as John Setsuko. Matched description given by night club doorman (IC1, Male, Ag 22-24 yrs, blue Everton shirt). Stopped whilst driving White Ford Mondeo, W563 WDL. Address given as 22 East Dene Ridge, Copdock, Ipswich. Searched at scene and found in possession of 1oz Cannabis Resin and lockable pocket knife.

| | |
|---|---|
| Arresting_Officer | PC 143 |
| Arrest_Date_Time | 15/06/2006 : 23:47 |
| Suspect_Forename | John |
| Suspect_Surname | Setsuko |
| Suspect_VRN | W563WDL |
| Suspect_Vehicle_Color | White |
| Suspect_Vehicle_Make | Ford Mondeo |
| Suspect_Addr_Street | 22 East Dene Ridge |
| Suspect_Addr_Town | Ipswich |
| Evidence_1_Description | 1 oz Cannabis Resin |
| Classification | Drug possession |



**What is Content Analytics?**

*Content Analytics* (Text Analytics + Mining) refers to the text analytics process plus the ability to visually identify and explore trends, patterns, and statistically relevant facts found in various types of content spread across internal and external content sources.
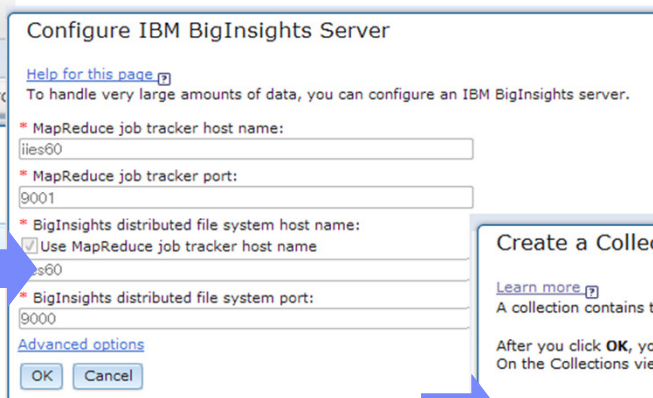
\* Natural Language Processing

IBM

# Search and Analytics quick look

- 8 views for analysis, exploration and investigation

- Learn about different ways to discover Rapid Insights from content

- Easy to use to search and analyze

**Facets**

**Dashboard**

**Time Series**

**Deviations / Trends**

**Document Analysis**

**Connections**

**Facet Pairs**

**Enterprise Search**

**Sentiment**

20

2 IBM Corporation

# New in ICAwES v3.0 - Seamless Scale-out with BigInsights / Hadoop

- **Select "Configure BigInsights Server"**

**IBM** IBM Content Analytics

Collections | System | Security

Data Listener | Server Configuration | Parse | Searc

Help for this page

Add Server | Configure IBM BigInsights Server

Server Information

- **Enter the BigInsights Server Information**

Configure IBM BigInsights Server

Help for this page
To handle very large amounts of data, you can configure an IBM BigInsights server.

* MapReduce job tracker host name:
iies60

* MapReduce job tracker port:
9001

* BigInsights distributed file system host name:
☑ Use MapReduce job tracker host name
es60

* BigInsights distributed file system port:
9000

Advanced options

OK | Cancel

- **Specify "Use IBM BigInsights" as a Collection setting**

Create a Collection

Learn more
A collection contains the various sources that users can search with a single query.

After you click **OK**, you return to the Collections view.
On the Collections view, click your new collection and add content by adding a crawler or import

General options

* Collection name:
NHTSA

* Collection type:
Text analytics collection

☐ Use IBM BigInsights

…configuration files and ICA libraries, UIMA PEARs (including custom PEAR) and other required modules will be distributed to BIgInsights servers automatically for that collection

Admin user can confirm the setting on Topology View

Learn more
Add Server | Configure IBM BigInsights Server | Start Selected Servers | Stop Selected Servers

Server Information

| | Host name | Crawler | Controller | Document Pro | Search | Backup | Status | Actions |
|---|---|---|---|---|---|---|---|---|
| | node0.ibm.com | | | | | | ● | × |
| | node3.ibm.com | ✓ | ✓ | | ✓ | | ● | ■× |
| | node1.ibm.com | | | ✓ | | | ● | ■× 📄 |
| | node2.ibm.com | | | ✓ | | | ● | ■× 📄 |
| | hadoop.ibm.com | | | ✓ | | | ● | |

Topology View

# New in ICAwES v3.0 - Seamless Scale-out with BigInsights / Hadoop



**IBM InfoSphere BigInsights**

**MetaTracker**

**IBM Content Analytics**

Analytics Flow

Pre-Processing

RDS

Cache

System-T Analysis

UIMA Analysis

Indexing

ICA GA

Main Index

Slave Index

Big Index

**Data Flow**

**Job Flow controlled by MetaTracker**

**Regular OS**

Various Data source

Crawler

Importer

RDS

**Analytics Flow**

Indexing Service Process

Cache

Local Analysis (UIMA base)

Global Analysis

Index

Text Mining/ Search Runtime

Exporter

Text Miner UI /Custom UI

Another App.

# Documents View



**View the Documents based on current analysis criteria**

- Flagging of documents
- Auto highlighting of key entities based on current query
- Export documents to be used in another system.



- **Document Analysis Dialog**
  - Preview the whole document
  - Display document metadata and annotated facets with highlighted texts

# Facets View



**Show the frequency and correlation indices of keywords belong to the selected facet**

- *Add selected keyword to current search condition*

- *Quick filter*
- *Multi-column sort*

# Deviations View and Trends View



- *Deviations View*
  - *show how keyword deviates from average occurrence between other keywords.*
  - *Identify patterns that are cyclic.*



- *Trends View*
  - *show how keyword in each time period deviates from average occurrence (against the whole current matched documents)..*
  - *Predict pattern behaviors in future.*

© 2012 IBM Corporation

# Facet Pairs View

- **Show the correlation between keywords belong to two different facets**
- **3 view modes**

**Bird's eye view**

**Look at whole area**

**Table view**

**Grid view**

**Quick filter and sort**

**See in detail**

# Connections View links highly correlated terms

- Show relationship between multiple facet values
- Connections between nodes represents correlation between two facet values
- Color of line represents the importance of correlation index (red is the highest)



Identify relations between "FORD", "blow" and "FIRESTONE"

# New in ICAwES v3.0 – Sentiment Analytics



**Sentiment Expressions**

- Lists positive or negative expressions for selected facet value with colors

- Color represents the rank of correlation

# Create Dashboard Views for Executive Summaries



**Click here**

# ICA integrates with Cognos BI reports

From ICA Text Miner, a user can:
- Issue a request to create a report
- List the created reports
- Open the created report
- Delete the created report
- Cognos reports can link to and from Text Miner

IBM

# IBM Content Analytics adds value to Big Data...

### Healthcare Analytics
- **Analyzing:** E-Medical records, hospital reports
- **For:** Clinical analysis; treatment protocol optimization
- **Benefits:** Better management of chronic diseases; optimized drug formularies; improved patient outcomes

### Customer Care
- **Analyzing:** Call center logs, emails, online media
- **For:** Buyer Behavior, Churn prediction
- **Benefits:** Improve Customer satisfaction and retention, marketing campaigns, find new revenue opportunities

### Crime Analytics
- **Analyzing:** Case files, police records, 911 calls...
- **For:** Rapid crime solving & crime trend analysis
- **Benefits:** Safer communities & optimized force deployment

### Insurance Fraud
- **Analyzing:** Insurance claims
- **For:** Detecting Fraudulent activity & patterns
- **Benefits:** Reduced losses, faster detection, more efficient claims processes

### Automotive Quality Insight
- **Analyzing:** Tech notes, call logs, online media
- **For:** Warranty Analysis, Quality Assurance
- **Benefits:** Reduce warranty costs, improve customer satisfaction, marketing campaigns

### Social Media for Marketing
- **Analyzing:** Call center notes, SharePoint, multiple content repositories
- **For:** churn prediction, product/brand quality
- **Benefits:** Improve consumer satisfaction, marketing campaigns, find new revenue opportunities or product/brand quality issues

Energy  Traffic  Food  Infrastructure  Retail  Intelligence  Stimulus  Banking  Telecom  Retail  Intelligence  Stimulus  Banking  Telecom  Oil  Healthcare  Cities  Water  Public safety

# Interaction between Big Data and Content Analytics

- **IBM Content Analytics running on top of BigInsights for full analytics and rapid insights**

- **BigSheets and IBM Content Analytics – Filter with BigSheets then analyze in ICA for graphical view of analytic results**

- **IBM Content Analytics as a NLP engine to add define concepts and add depth to Cognos Consumer Insights**

## Agenda

| 1 | **Big Data – Definitions & A Quick Story** |
|---|---|
| 2 | **Technology Primer** |
| 3 | **Dive on Content Analytics** |
| 4 | **Thoughts On Getting Started** |

# Thoughts On Getting Started

# Use Case Selection Criteria

- **Well proven path**
- **Can be done off-line/non-disruptive to existing systems**
- **Intuitive that there is low hanging fruit for additional insights**
- **Data set is already stored, but under-instrumented or overly summarized**
- **Initial findings can be arrived at 3 weeks or less**
- **Initial use cases have to be accretive to next set**
- **Initial use cases have to leverage common technology for next set of use cases**

# One Example Of How To Think About It

**Start here**

**Grow to**

Value

High

Predictive
Cross Channel
Modeling

Cross Channel
Communications

Customer
Email

Active
Listening To
Social Media @
Individual Level

Agent Notes

Passive
Listening To
The Social
Media Crowd

Web Log

Complexity

Low

High

# Use Case Starting Points

- **Mobile/Web Log**
  - Why: Reduce latency to under a day instead of weeks from 3rd parties
  - Why: Have the full data set for cross-correlation for other use cases
- **Agent Notes**
  - Why: Has proven to be an importance source of wisdom from the branches and Contact Centres
  - Why: An element of crowd sourcing
  - Why: Pulse of our internal teams whom we need to market to
- **Customer Email**
  - Why: They are often your best direct source of "pulse" of the business
  - Why: Lynchpin of crowd sourcing
  - Why: Six Sigma feedback component to cost take out
- **Passive Listening To The Social Media Crowd**
  - Why: Important early warning source of information
  - Why: You need the data – not summarization – for more advanced analytics
- **Cross Channel Communications**
- **Active Listening To Social Media @ Individual Level**
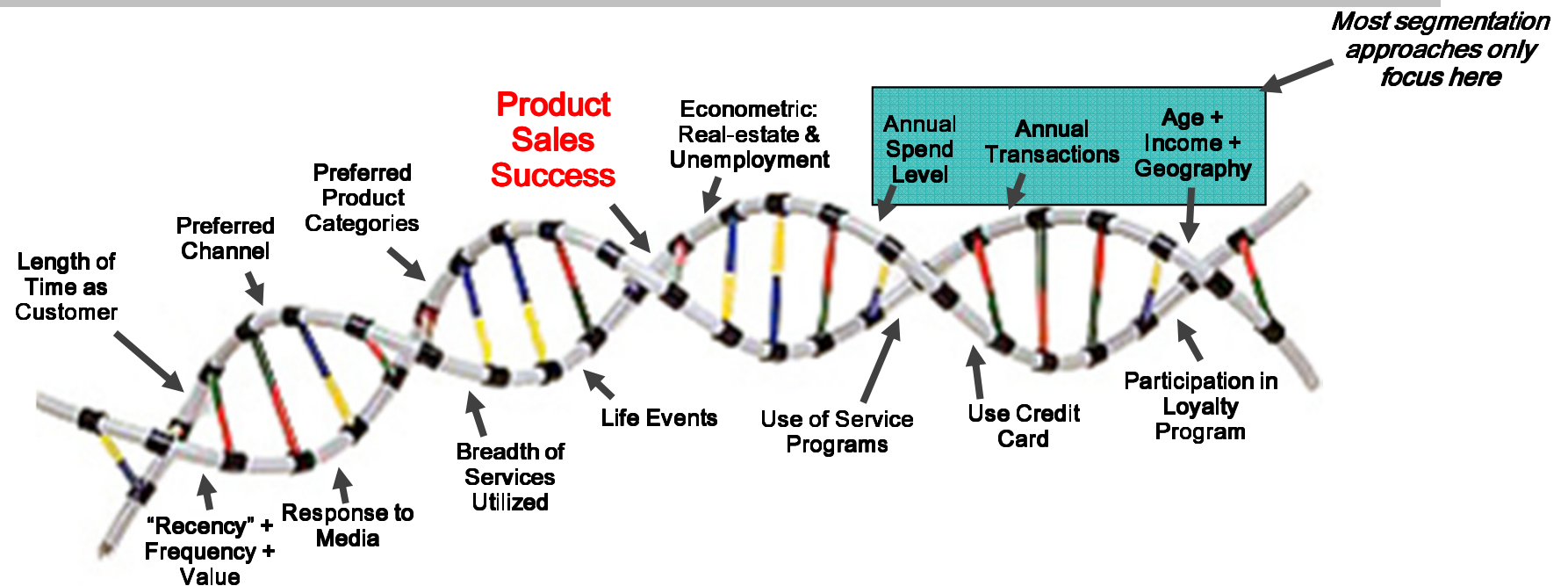- **Predictive Cross Channel Modeling**

# Examples Of Insight / Value Delivered

- **Web Log** *
  - Same day feedback on A|B testing new design (if sample is big enough)
  - Same day feedback on change in font/images/icons wording on behavior (if sample if big enough)
  - Input source into cross-channel behavior modeling
- **Agent Notes** *
  - Additional insight into customer "state" based on call handling
  - Feedback on features / policies not captured in formal system – nuance of customer opinion
  - Detecting of otherwise un-surfaced issues both at segment and individual customer level
  - Detection of staff dissatisfaction
  - Detection of otherwise hidden life events
- **Customer Email** *
  - Direct feedback on policy/fees
  - Notice of customers planning to attrition
  - Detecting of otherwise un-surfaced issues both at segment and individual customer level
  - (in private / brokerage situations early warning of law suits / advisor attrition)
  - Detection of otherwise hidden life events
  - Basis for cross-interaction behavioral / acceptance modeling
- **Passive Listening To The Social Media Crowd** *
  - Hour-by-hour feedback on announcements, offerings, policy changes
  - Taking the base reporting the rest of the org expects off the table so more important things can be focused on
  - Crowd feedback on competitor offerings
  - Path to SMARC

**\* All are expanded "feature vectors"**

# Case Study: Improvement Of Action Cluster Segmentation Approach

- Use of under-utilized data to expand modeled variables – "Feature Vectors" – The customers response to the firm's value proposition
- Each feature vector is like a gene strand, which describes a facet, or set of customer behavior traits



*Most segmentation approaches only focus here*

Product Sales Success

Econometric: Real-estate & Unemployment

Annual Spend Level

Annual Transactions

Age + Income + Geography

Preferred Product Categories

Preferred Channel

Length of Time as Customer

Life Events

Use of Service Programs

Use Credit Card

Participation in Loyalty Program

Breadth of Services Utilized

"Recency" + Frequency + Value

Response to Media

# Wait – How Do You Know If This Is Working?

- **Quick – how many of you were part of an experiment this week?**

- **Back testing**
- **A|B real world testing**
- **"Test Fast, Fail Fast, Correct Fast"**
- **Concept of MVP (minimum viable product)**

# THINK
# BIG