# Overview of HPCC Systems and Case Studies

Presenter: Brian Bounds, Director of Software Engineering
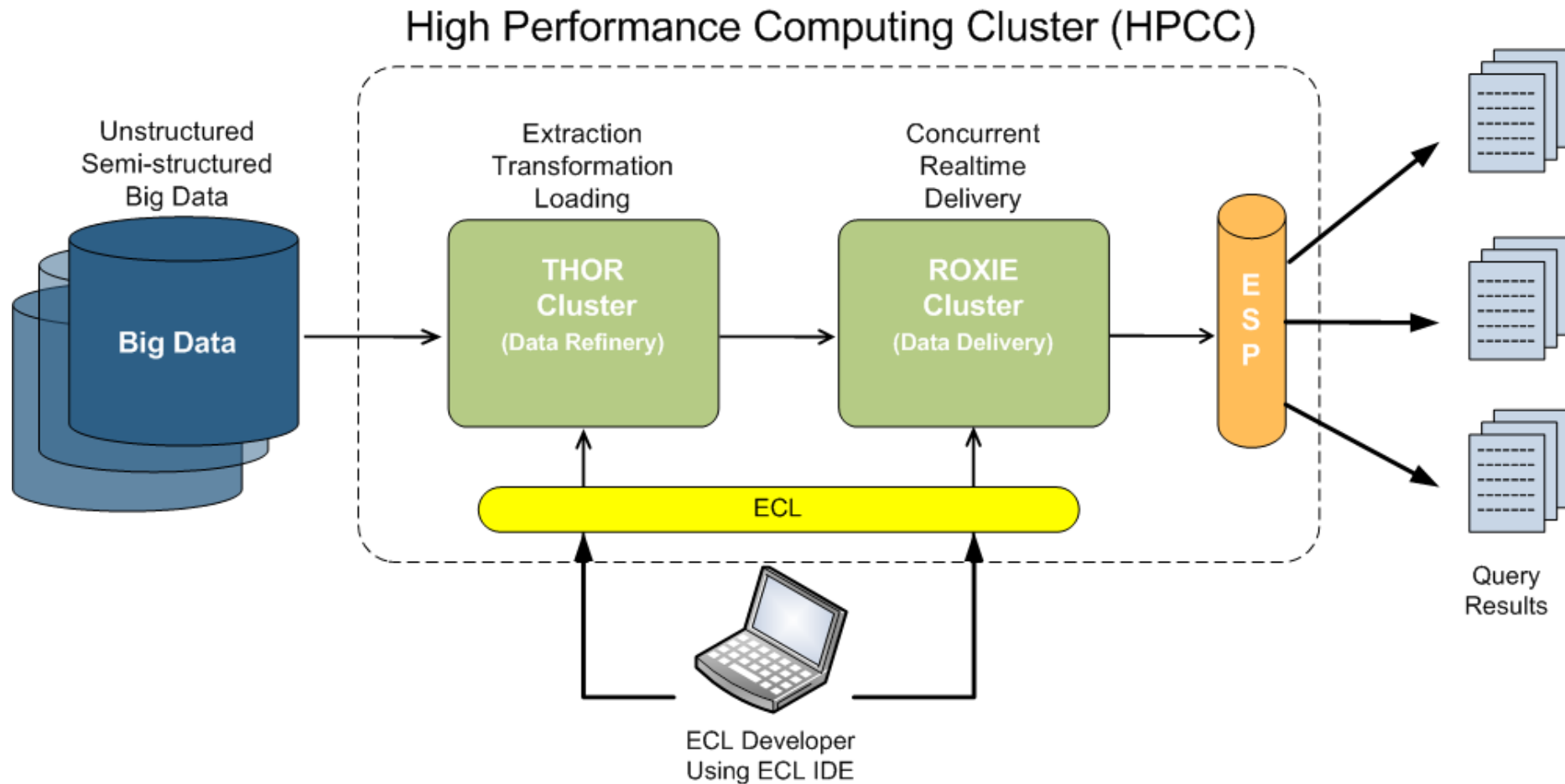November 2014

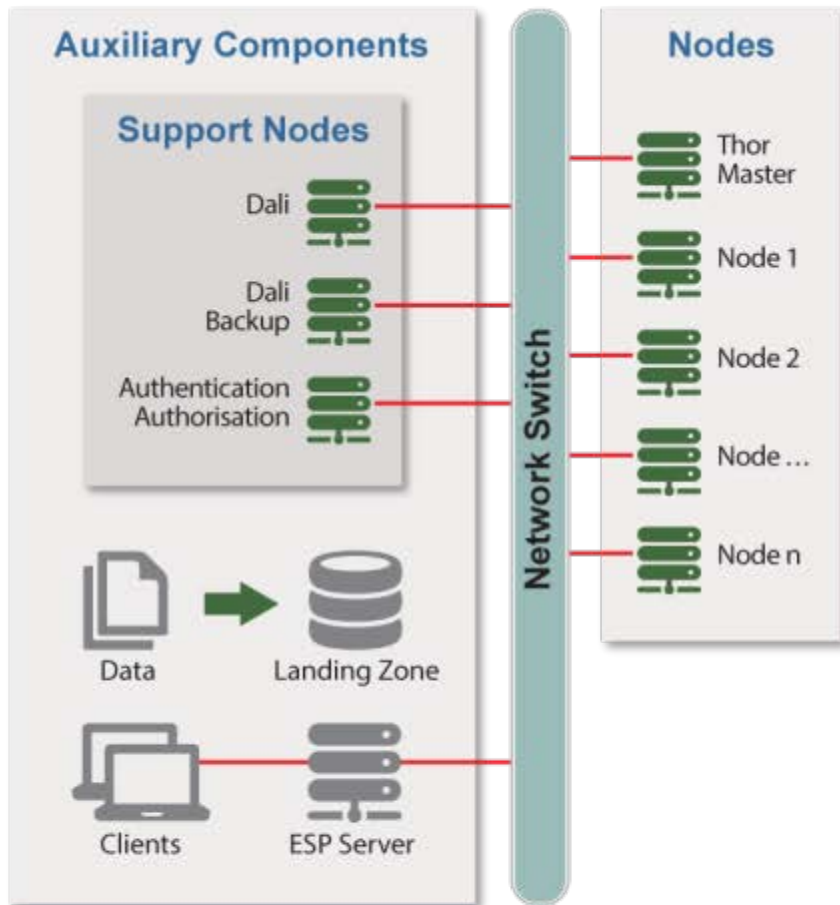**LexisNexis**

Risk Solutions

# Context

- About Brian
- HPCC Systems in the Context of LexisNexis
- Isn't Big Data just Data?

- The HPCC Systems platform includes:
  - Thor: batch oriented data manipulation, linking and analytics engine
  - Roxie: real-time data delivery and analytics engine

- A high level declarative dataflow language: ECL
  - Implicitly parallel
  - No side effects
  - Code/data encapsulation
  - Extensible
  - Highly optimized
  - Builds graphical execution plans
  - Compiles into C++ and native machine code
  - Common to Thor and Roxie

- An extensive library of ECL modules, including data profiling, linking, graph analytics, and Machine Learning
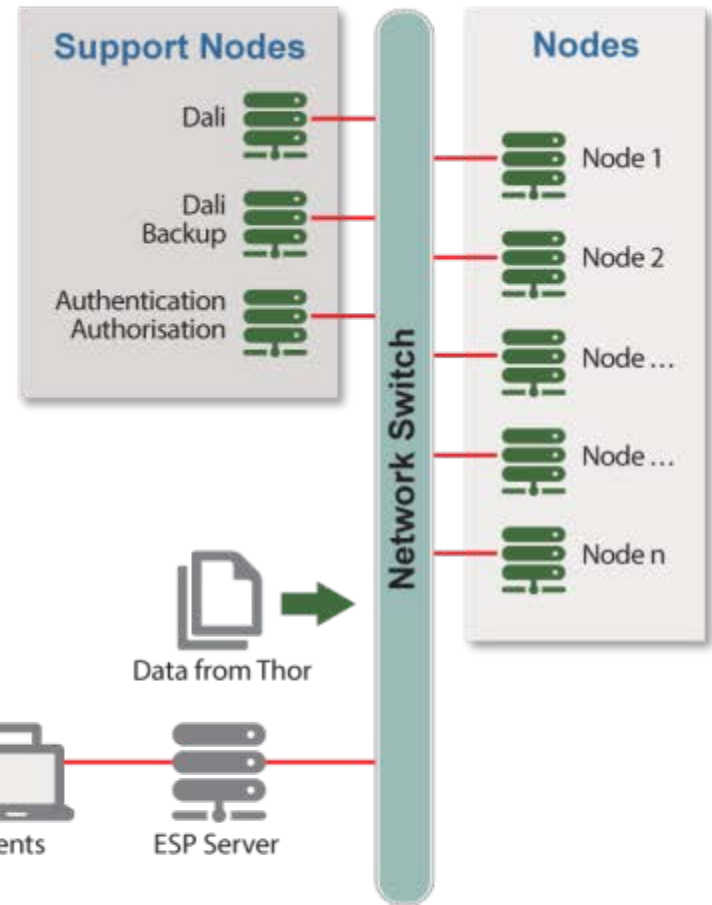
High Performance Computing Cluster (HPCC)

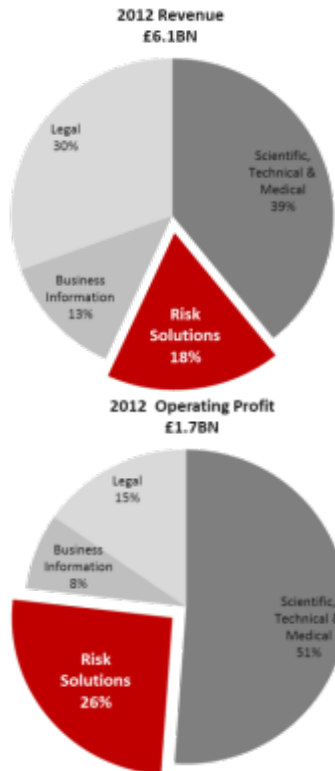# Detailed HPCC Systems Platform Architecture

# Drea's HPCC Overview

- The Programming Language (ECL)

- The Delivery Engine (ROXIE)

- Enterprise Readiness

- Big Data ... becomes Data (that might be Big)

# Case Study #1 (Enterprise) – LexisNexis Risk Solutions

# Case Study #1 (Enterprise) – LexisNexis Risk Solutions

We are among the largest providers of risk solutions in the market today

**Reed Elsevier is a world leading provider of information solutions.**
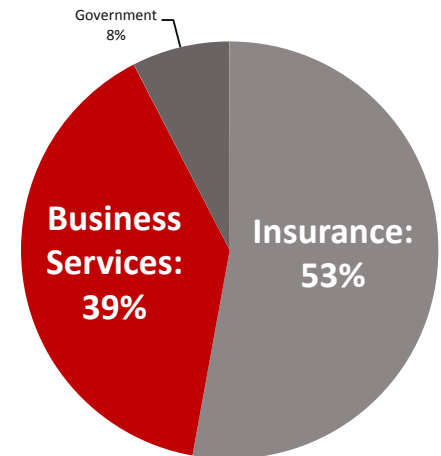
**LexisNexis Risk Solutions has seen sustained revenue and profit growth**

**LexisNexis Risk Solutions is a leading provider in the U.S. across Business Services, Insurance and Government segments.**



2012 Revenue
£6.1BN

Legal 30%

Scientific, Technical & Medical 39%

Business Information 13%

Risk Solutions 18%

2012 Operating Profit
£1.7BN

Legal 15%

Business Information 8%

Scientific, Technical & Medical 51%

Risk Solutions 26%



LexisNexis Risk Solutions Revenue: $M



LexisNexis Risk Solutions Revenue by Segment

Government 8%

Business Services: 39%

Insurance: 53%

LexisNexis

## We are among the largest providers of risk solutions in the market today

LexisNexis Risk Solutions is a leading provider in the U.S. across Business Services, Insurance and Government segments.

*Our customers include:*
- 99 of the top 100 US banks
- 90% of the Fortune 500
- 100% of US P&C insurance carriers
- All 50 US states, 70% of local governments and 80% of US federal agencies
- 97 of Am Law 100 firms

LexisNexis Risk Solutions Revenue by Segment

Government 8%

Business Services: 39%

Insurance: 53%

LexisNexis

# Case Study #1 (Enterprise) – LexisNexis Risk Solutions

**We have a unique set of capabilities: Data, Linking, Analytics, and Product Development**

| Data Technology | | Vast Data Resources | | Linking & Analytics | | Industry-Specific Expertise & Delivery | | Customer-Focused Solutions |
|---|---|---|---|---|---|---|---|---|
| | **+** | | **+** | | **+** | | **=** | |
| • Speed<br>• Capacity<br>• Cost savings | | • Process<br>• Sources<br>• Coverage | | • Advanced linking & analytics<br>• Accuracy & efficiency<br>• Protect private information | | • Aligned with our customers' industries<br>• Deep industry expertise | | • Predict, manage and assess risk across many industries. |

LexisNexis®
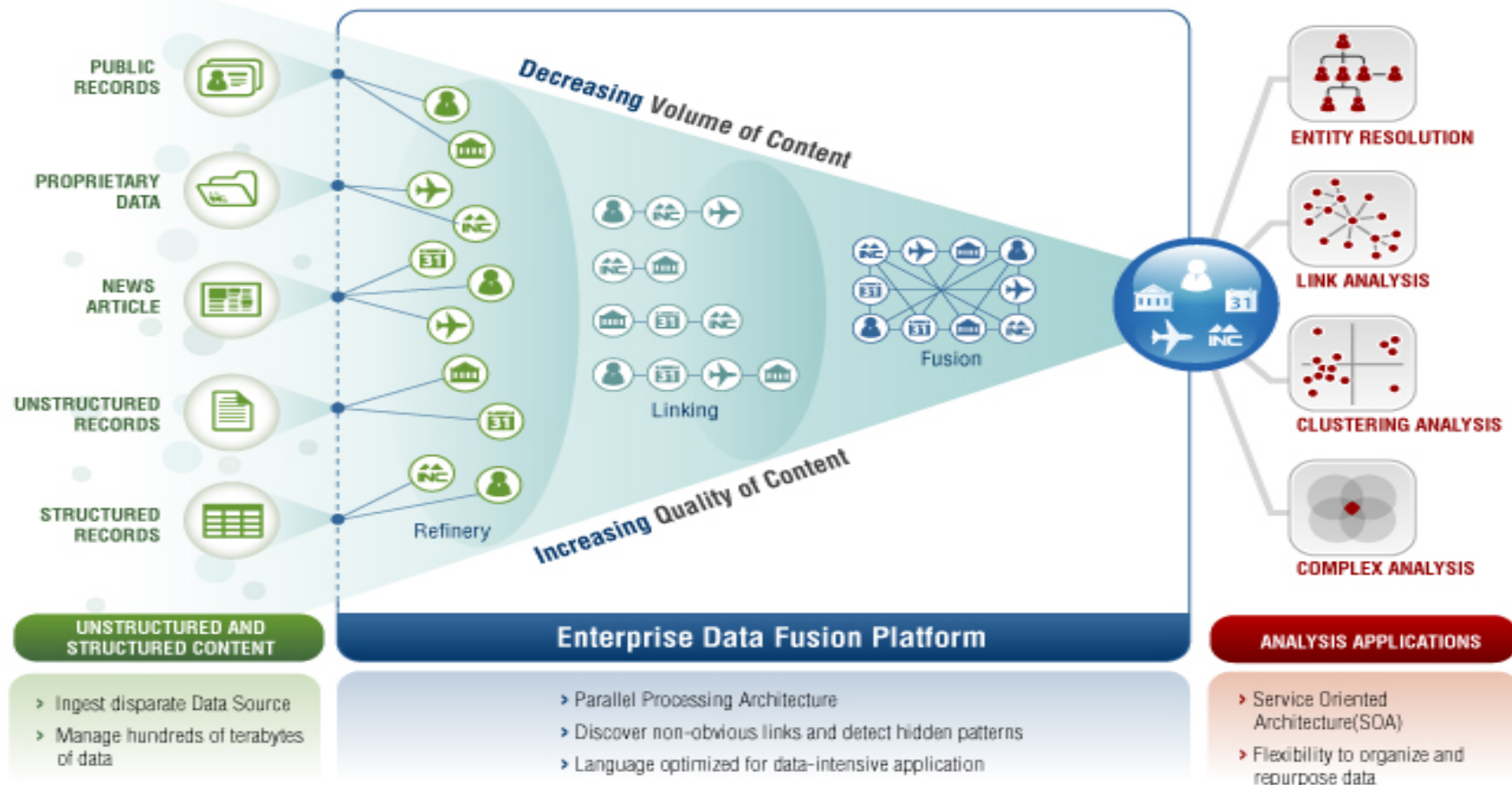
# Case Study #1 (Enterprise) – LexisNexis Risk Solutions

## Access to more than 25 billion public record filings
## Break-down of record counts for the more popular data sets:

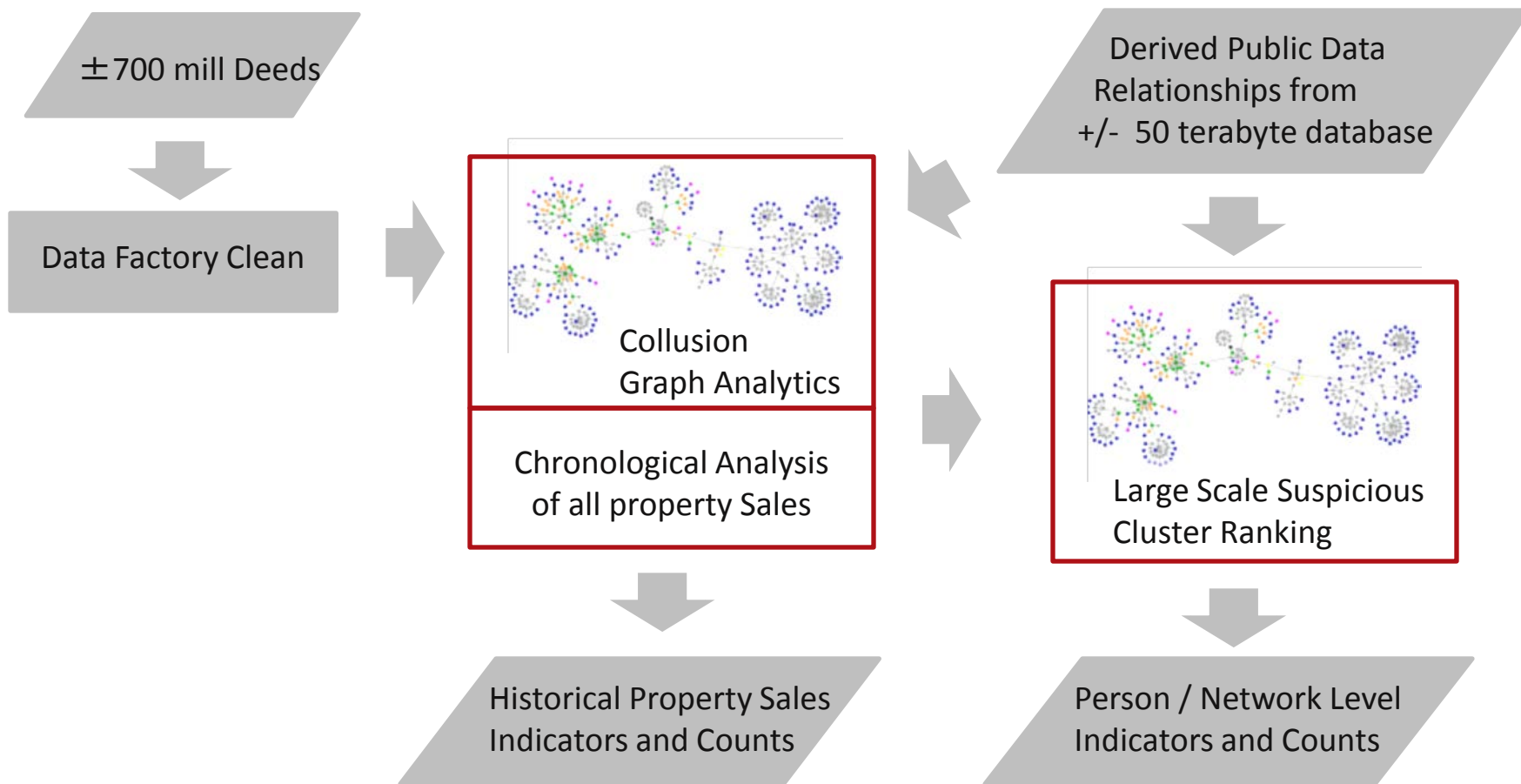| Data Source | # of records | | Data Source | # of records |
|---|---|---|---|---|
| Associates/Relatives | 1.8 Billion | | People at Work | 1.5 Billion |
| Bankruptcy | 23 Million | | Private Phones | 172 Million |
| Business BDID's | 283 Million | | Professional Licenses | 94 Million |
| Business People Links | 959 Million | | Property | 2.5 Billion |
| Canadian Phones | 62 Million | | Sex Offenders | 550,000 |
| Consumer Header | 10.8 Billion | | SSN's | 7.2 Billion |
| Criminal | 216 Million | | Student Records | 38 Million |
| Date of Birth | 5.2 Billion | | TIN | 2.9 Million |
| Death | 98 Million | | Unique ADLs - active | 257 Million |
| Drivers Licenses | 397 Million | | Utility | 645 Million |
| EDA Phones | 124 Million | | Vehicle Titles | 635 Million |
| FEINs | 10.4 Million | | Vehicle Registrations | 2.5 Billion |
| Historical Phones | 800 Million | | White Pages | 116 Million |
| Hunting and Fishing Licenses | 67 Million | | Wireless Phones | 101 Million |
| Liens and Judgments | 244 Million | | Yellow Pages | 14 Million |

Case Study #2 (Boil the Ocean) – Property  Transaction Risk

**Three core transaction variables measured**

- Velocity

- Profit (or not)

- Buyer to Seller Relationship Distance
(Potential of Collusion)

Flipping

Collusion

Profit

±700 mill Deeds

Data Factory Clean

Collusion
Graph Analytics

Chronological Analysis
of all property Sales

Derived Public Data
Relationships from
+/- 50 terabyte database

Large Scale Suspicious
Cluster Ranking

Historical Property Sales
Indicators and Counts

Person / Network Level
Indicators and Counts

**Large scale measurement of influencers strategically placed to potentially direct suspicious transactions.**

- All data on one supercomputer measuring over a decade of property transfers nationwide.

- Data Products to turn other Data into compelling intelligence.

- Large Scale Graph Analytics allow for identifying known unknowns.

- Florida Proof of Concept
  - Highest ranked influencers
    - Identified known ringleaders in flipping and equity stripping schemes.
    - Typically not connected directly to suspicious transactions.
  - Known ringleaders not the Highest Ranking.

- Clusters with high levels of potential collusion.
- Clusters offloading property, generating defaults.
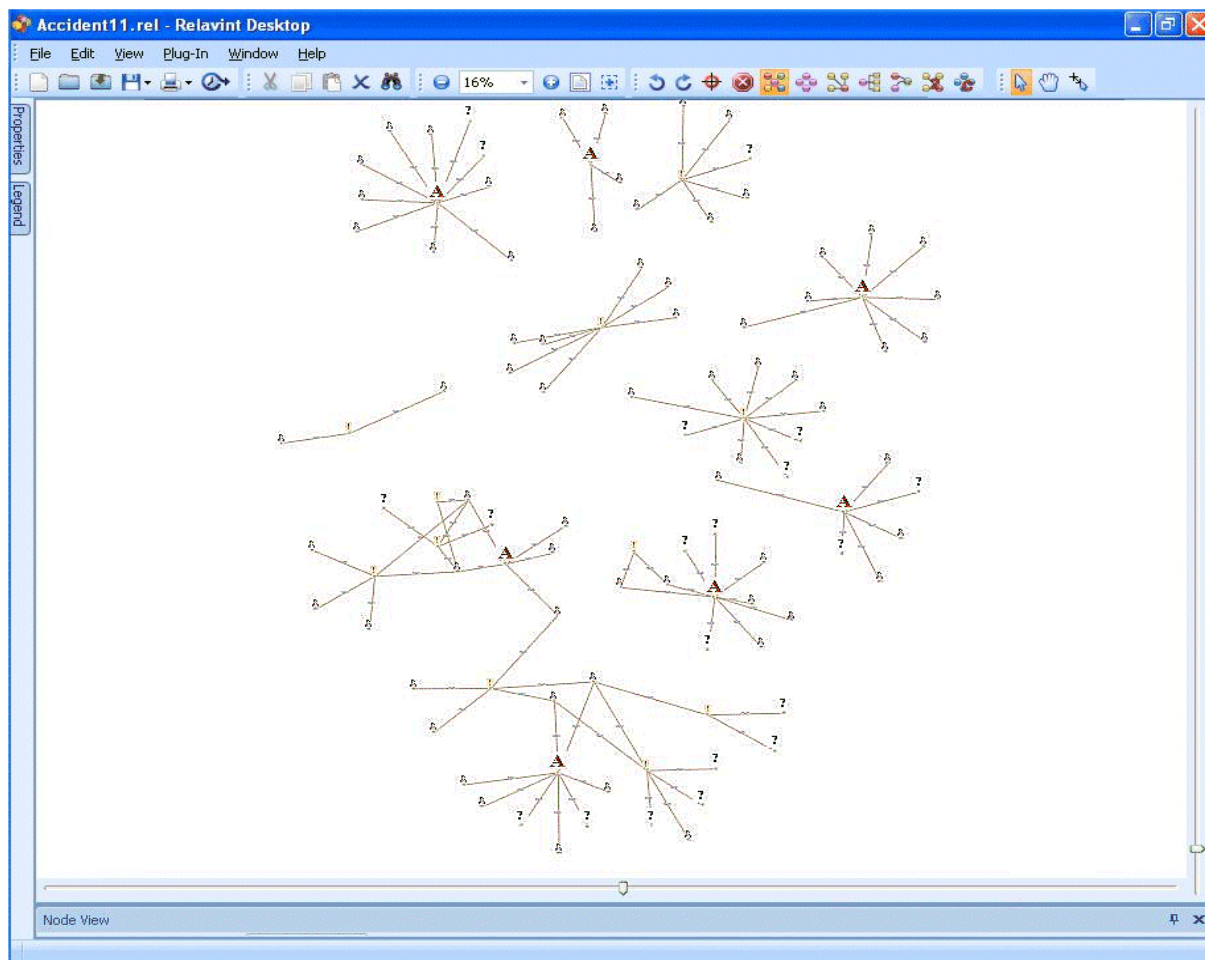- Agile Framework able to keep step with emerging schemes in real estate.

Case Study #3 (Serendipity) – Family Ties

## Scenario

This view of carrier data shows seven known fraud claims and an additional linked claim.

The Insurance company data **only finds a connection between two of the seven claims**, and only identified one other claim as being weakly connected.
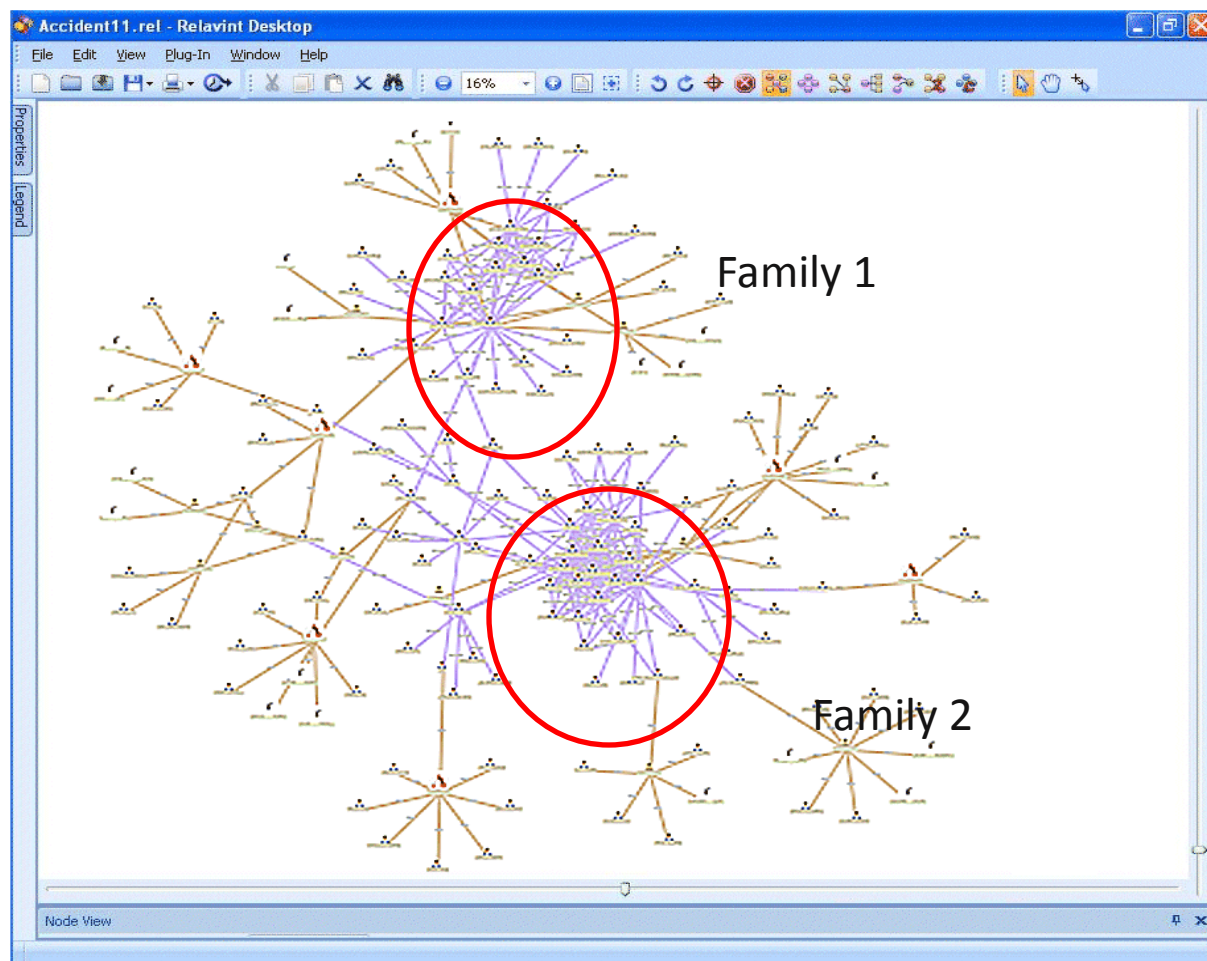
## Task

After adding the LexID to the carrier Data, LexisNexis HPCC technology then explored 2 additional degrees of relative separation

## Result

The results showed **two family groups interconnected on all of these seven claims**.

The links were much stronger than the carrier data previously supported.

Case Study #4 (Tactical) – 30 Hour Job

Objective:

- Re-engineer long-running legacy process (proof-of-concept)
- 3m+ rows in … 500m rows out … 30+ hours
- Use similar hardware to maximize comparability

Pre-Coding Set Up:

- Legacy developers … completed on-line ECL training
- 1 ECL developer … exposed to legacy process and data
- Dump of input data and known result files
- Create HPCC hardware environment comparable to legacy environment
    - Amazon AWS 4 x m1.xlarge total (3 x m1.xlarge Thor Slaves)
    - 12-slave CPUs

Coding:

- Meet-up for 1 week coding session

# Case Study #4 (Tactical) – 30 Hour Job

Legacy Environment:
- Oracle on Intel SMP
- 16 cores

HPCC Environment:
- Amazon AWS
- 1 x m1.xlarge (support node)
- 3 x m1.xlarge (Thor Slaves)
- 12-slave Cores

Results:
- 3 Days of Coding
- 450-ish Lines of ECL
- Legacy Run-Time: 30+ hours
- HPCC Run-Time: 1.5 hours

# Additional Infomation

- LexisNexis Open Source HPCC Systems Platform: http://hpccsystems.com

- Free Online Training: http://learn.lexisnexis.com/hpcc

- SALT: http://hpccsystems.com/products-and-services/products/modules/SALT

- Machine Learning portal: http://hpccsystems.com/ml

- The HPCC Systems blog: http://hpccsystems.com/blog

- Community Forums: http://hpccsystems.com/bb

- Our GitHub portal: https://github.com/hpcc-systems

- JIRA: https://track.hpccsystems.com

Thank you!