



High Performance Computing Cluster



Collection – Structured, unstructured and semi-structured data from multiple sources

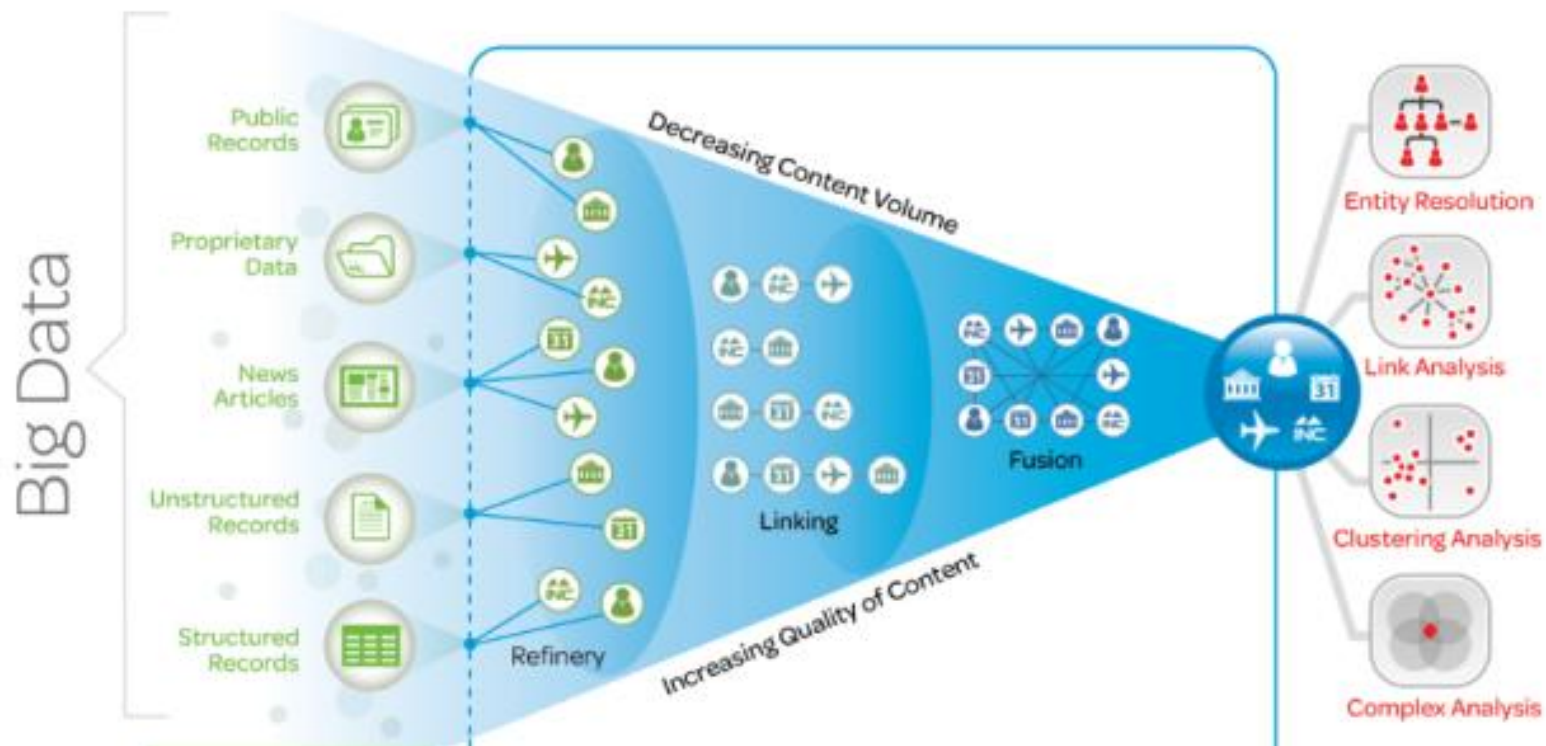
Ingestion – loading vast amounts of data onto a single data store

Discovery & Cleansing – understanding format and content; clean up and formatting

Integration – linking, entity extraction, entity resolution, indexing and data fusion

Analysis – Intelligence, statistics, predictive and text analytics, machine learning

Delivery – querying, visualization, real time delivery on enterprise-class availability





High Performance Computing Cluster (HPCC)

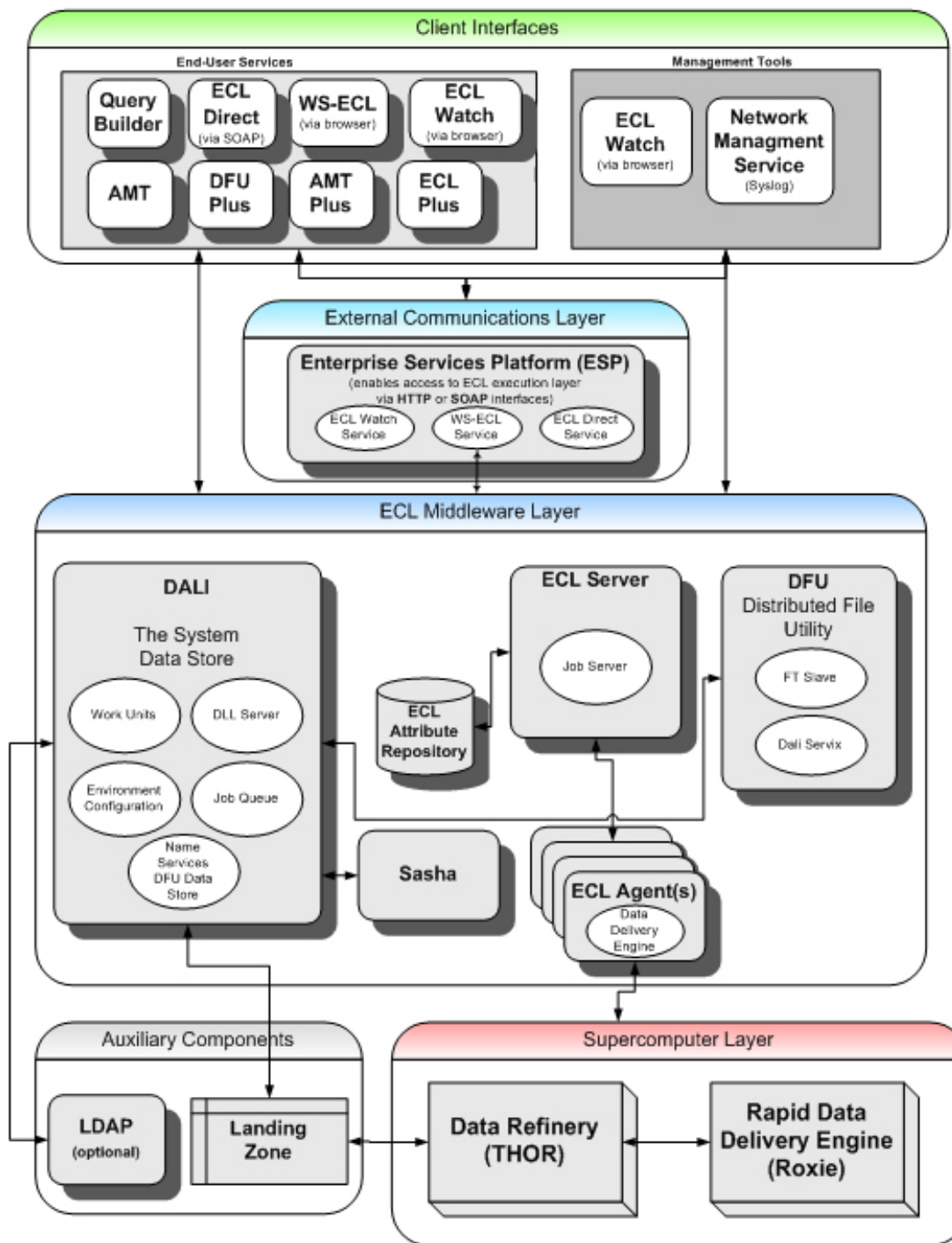
- Large-scale data storage and analytics platform
- Developed by LexisNexis Risk Solutions in the early 2000s
- Released as open-source in 2011
- Configurations to support both parallel batch data processing and online query applications using indexed data files
- Enterprise Control Language (ECL)



High Performance Computing Cluster

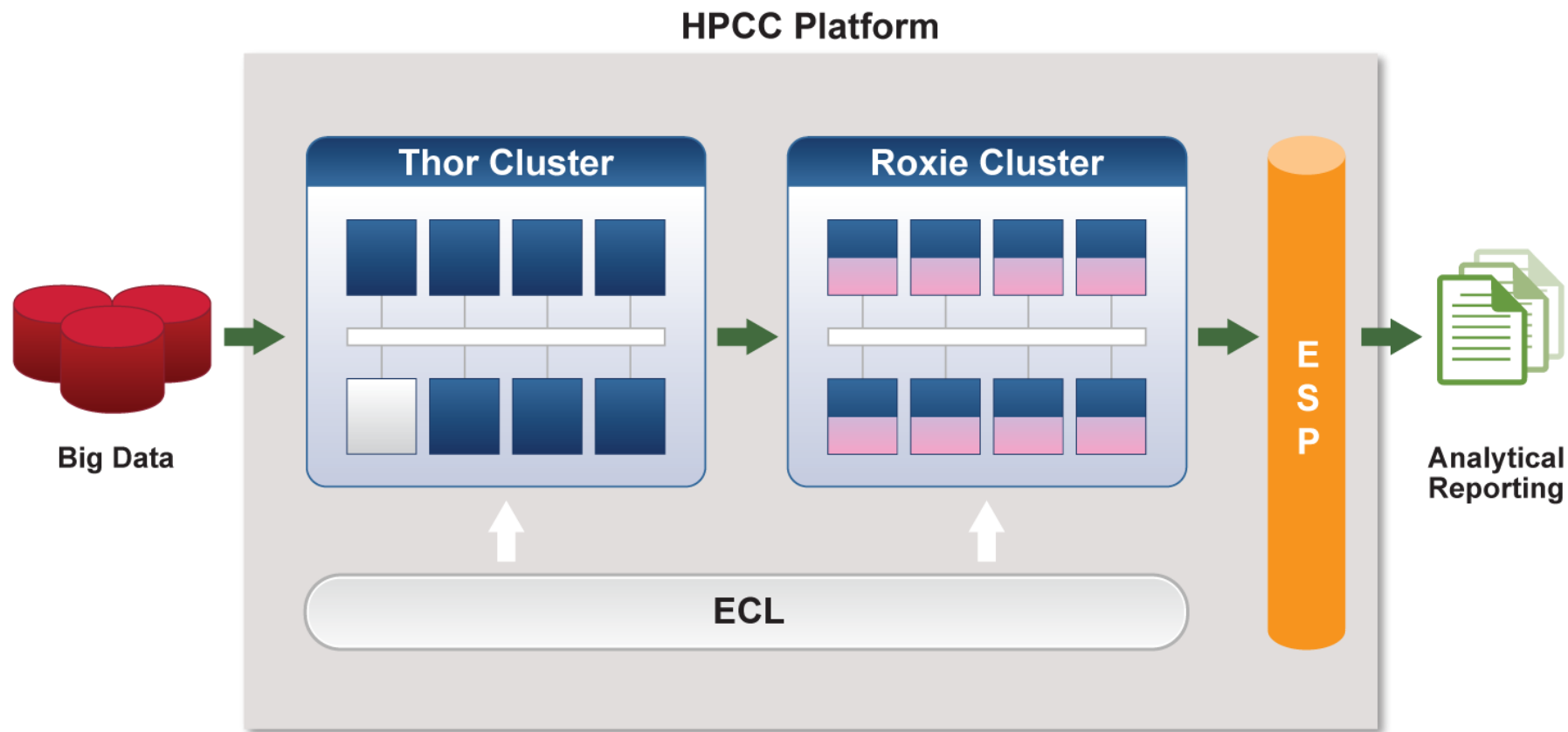
HPCC

- Consists of two processing environments
 - Thor (Data Refinery)
 - Roxie (Data Delivery Engine)
- Software and Middleware
- Enterprise Control Language (ECL)
- Client Interfaces
 - ECL Watch
 - ECL IDE
 - ECL Plugin for Eclipse



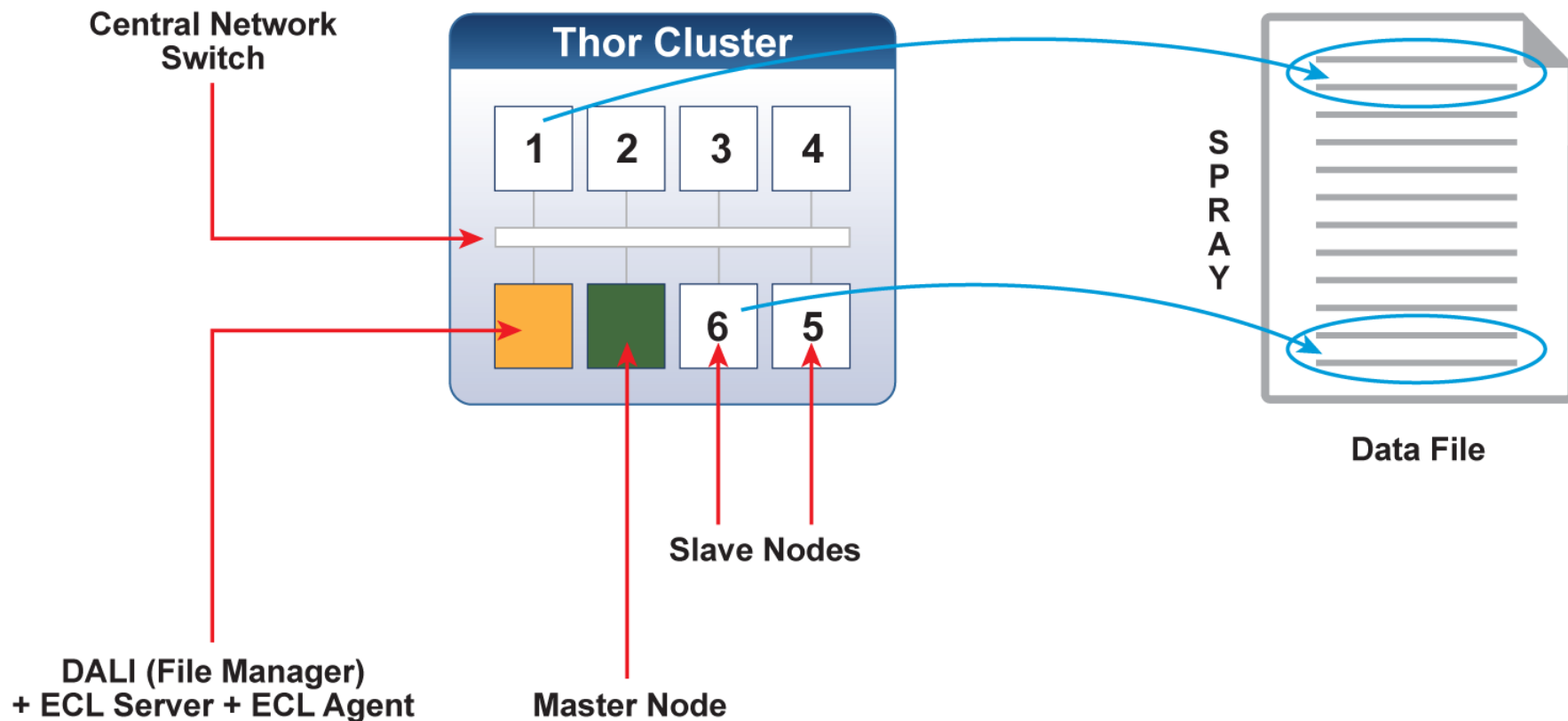


HPCC Systems Architecture



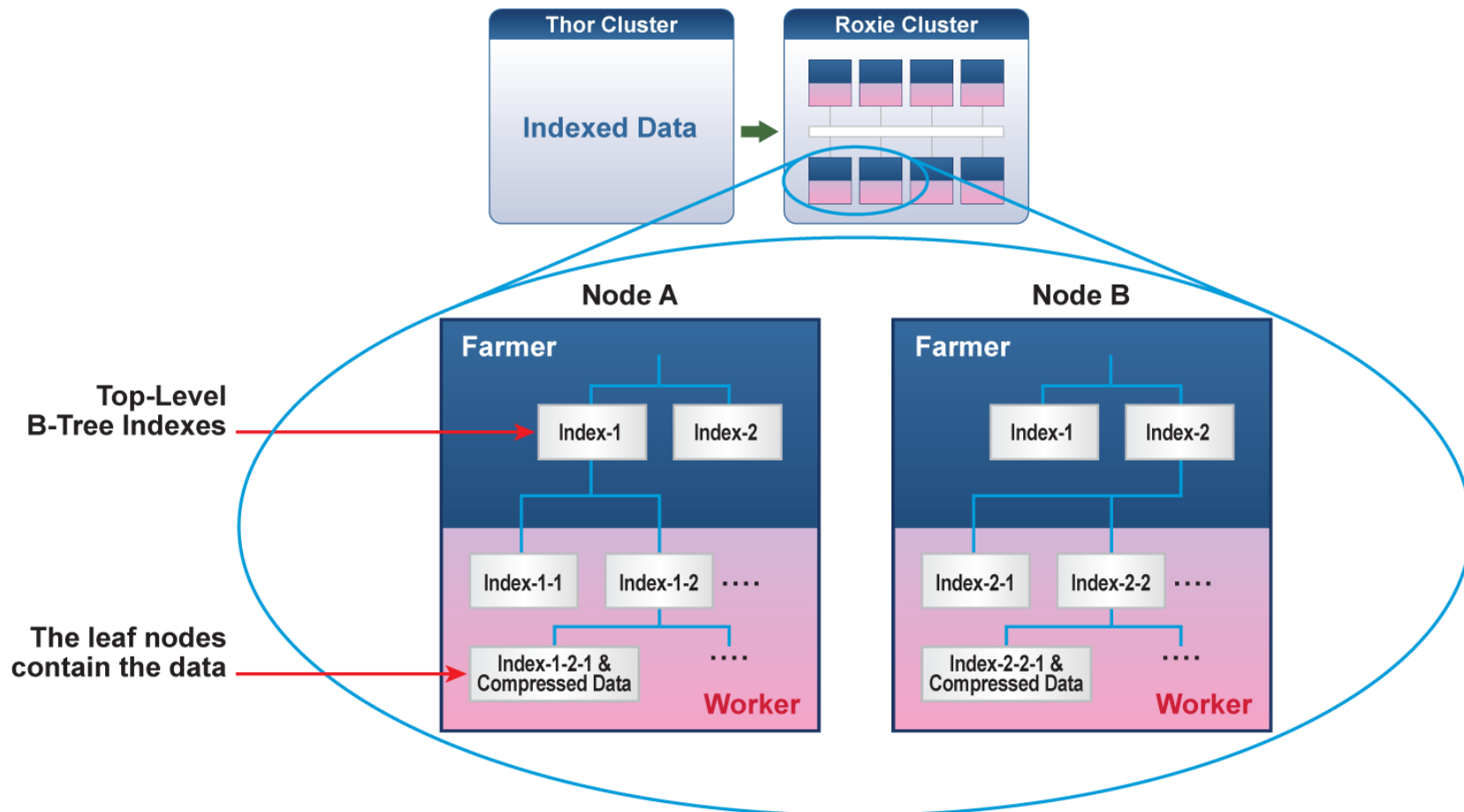


Thor Logical Architecture



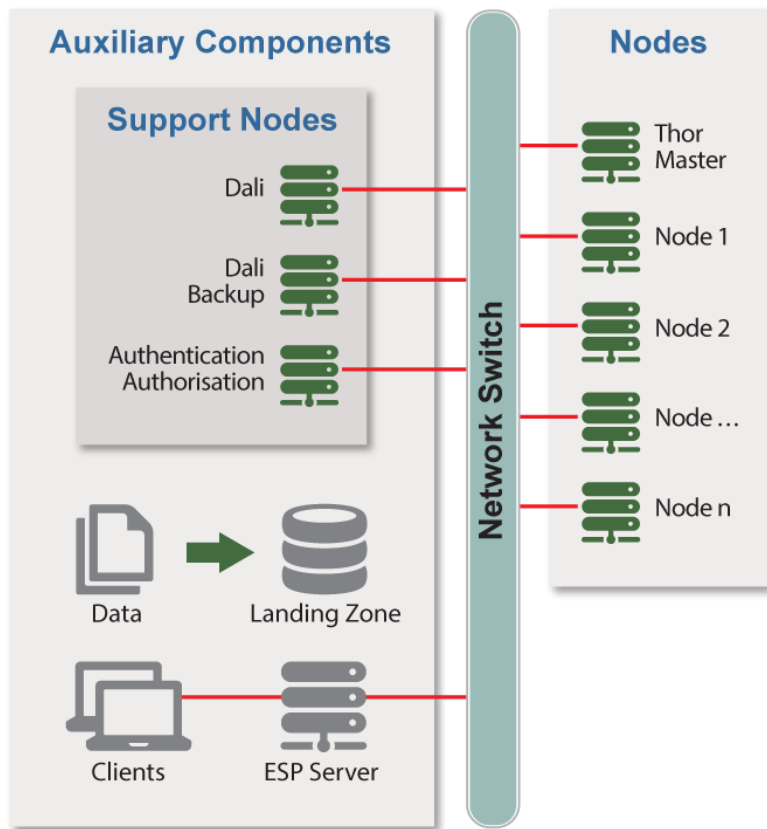


Roxie Logical Architecture



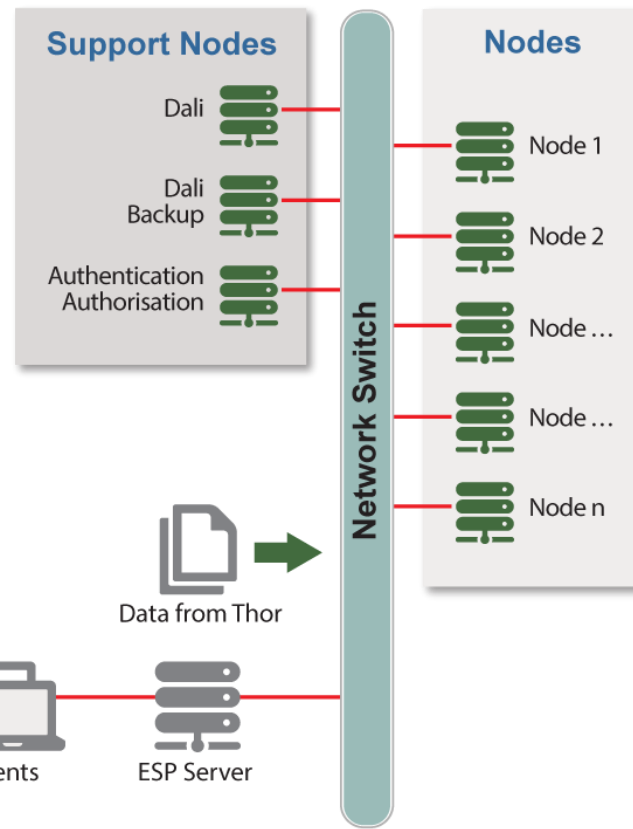


Thor/Roxie Physical Architecture



Thor

(Batch Job Execution Engine + DFS)
Physical Layout Schematic Diagram

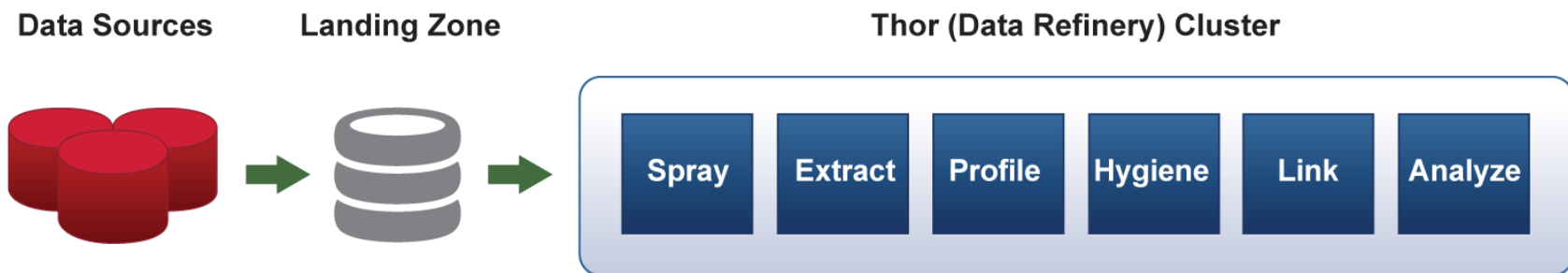


Roxie

(Rapid Data Delivery Engine)
Physical Layout Schematic Diagram



Data Refinery Process





- ```
// Initialize output log
log_out_init := project(log_init,
 transform(layout_logout,
 self := left,
 self := []));

// Create error log
outerrorfile := join(log_seq,
 log_out_init,
 left.linenum = right.linenum,
 transform(recordof(log_seq),
 self := left),
 left only,
 hash);

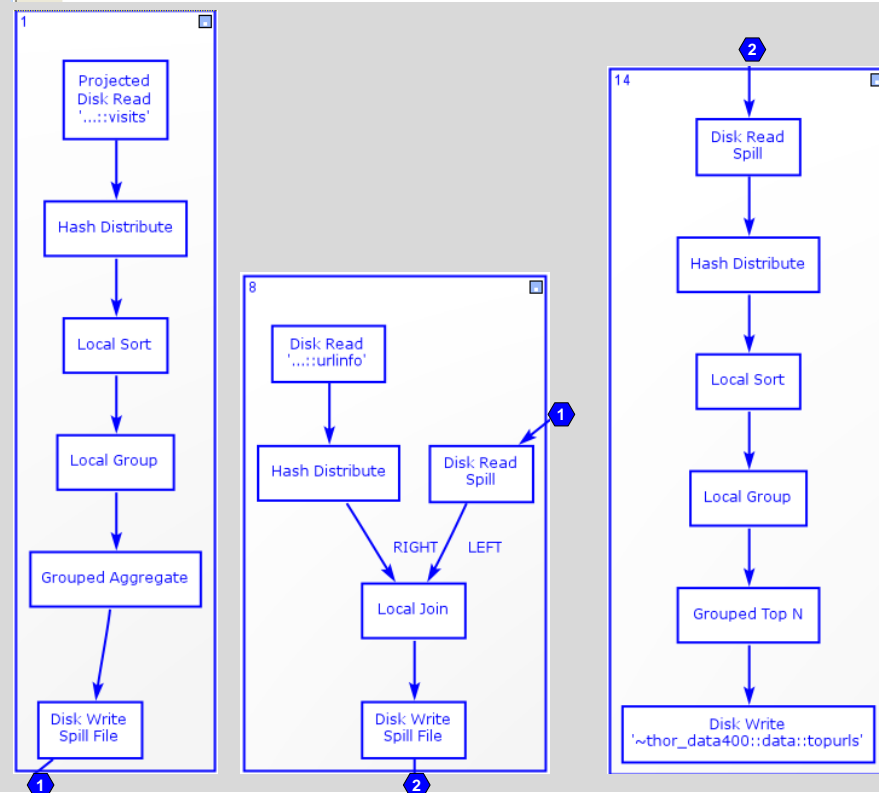
// Denormalize key value pairs
outlogfile := sort(denormalize(distribute(log_seq,
 sort(distribute(key_val,
 left.linenum = right.linenum,
 transform(layout_logout,
 self.keyvals := left,
 row({right.linenum,
 self := left}),
 left.linenum = right.linenum,
 transform(recordof(log_seq),
 self := left),
 left only,
 hash)
 left.linenum = right.linenum,
 transform(recordof(log_seq),
 self := left),
 left only,
 hash)
 left.linenum = right.linenum,
 transform(recordof(log_seq),
 self := left),
 left only,
 hash)
 left.linenum = right.linenum,
 transform(recordof(log_seq),
 self := left),
 left only,
 hash);
```



# Enterprise Control Language (ECL)

- 1 line of ECL is roughly equivalent to 120 lines of C++
- ECL primitives that act upon datasets include:
  - SORT,
  - ROLLUP,
  - DEDUP,
  - ITERATE,
  - PROJECT,
  - JOIN,
  - NORMALIZE,
  - DENORMALIZE,
  - PARSE,
  - DISTRIBUTE
  - Etc.

```
Go Queue: dev_edserver_ Cluster: thor400_88_de More
1 // Sample ECL Code
2 layout_visits := RECORD string user; string url; string time; END;
3 visits := DATASET('~thor_data400::data::visits', layout_visits, FLAT);
4
5 layout_urlInfo := RECORD string url; string category; string pRank; END;
6 urlInfo := DATASET('~thor_data400::data::urlInfo', layout_urlInfo, FLAT);
7
8 // Distribute Visits by URL, Count visits by URL
9 layout_visitCounts := RECORD visits.url; visits_cnt := COUNT(GROUP); END;
10 visitCounts := TABLE(DISTRIBUTE(visits, HASH(url)), layout_visitCounts, url, LOCAL);
11
12 // Distribute Category by URL, Join category to URLs
13 visitCountsCat := JOIN(visitCounts, DISTRIBUTE(urlInfo, HASH(url)), LEFT.URL=RIGHT.URL, LOCAL);
14
15 // Distribute and Group by Category, Output top 10 URLs for each category
16 topUrls := TOPN(GROUP(DISTRIBUTE(visitCountsCat, HASH(category)), category, ALL, LOCAL), 10, -visits_cnt);
17 OUTPUT(topUrls, '~thor_data400::data::topurls', OVERWRITE);
```





# Enterprise Control Language (ECL)

- 1 line of ECL is roughly equivalent to 120 lines of C++
- ECL primitives that act upon datasets include: SORT, ROLLUP, DEDUP, ITERATE, PROJECT, JOIN, NORMALIZE, DENORMALIZE, PARSE, CHOSEN, ENTH, TOPN, DISTRIBUTE



# Enterprise Control Language (ECL)

```
// First declare a dataset with one column containing a list of strings
// Datasets can also be binary, CSV, XML or externally defined structures

D := DATASET([{'ECL'},{'Declarative'},{'Data'},{'Centric'},{'Programming'},{'Language'}],{STRING Value;});
SD := SORT(D,Value);
output(SD)
```

- OUTPUT(SD)
  - What is an SD?
- SD := SORT(D,Value);
  - SD is a D that has been sorted by 'Value'
  - What is a D?
- D :=  
DATASET([{'ECL'},{'Declarative'},{'Data'},{'Centric'},{'Programming'},{'Language'}],{STRING Value;});
  - D is a dataset with one column labeled 'Value' and containing the following list of data.



# HPCC vs HADOOP

| Hadoop Name/Term         | ECL equivalent              | Comments                                                                                                      |
|--------------------------|-----------------------------|---------------------------------------------------------------------------------------------------------------|
| MAPing within the MAPper | PROJECT/TRANSFORM           | Takes a record and converts to a different format; in the Hadoop case the conversion is into a key-value pair |
| SHUFFLE (Phase 1)        | DISTRIBUTE(,HASH(KeyValue)) | The records from the mapper are distributed dependent upon the KEY value                                      |
| SHUFFLE (Phase 2)        | SORT(,LOCAL)                | The records arriving at a particular reducer are sorted into KEY order                                        |
| REDUCE                   | ROLLUP(,Key,LOCAL)          | The records for a particular KEY value are now combined together                                              |





# HPCC vs HADOOP

| Feature                 | HPCC               | Hadoop     |
|-------------------------|--------------------|------------|
| Distributed File System | Thor DFS/Roxie DFS | Hadoop DFS |
| Database Capability     | native to DFS      | HBase      |
| Data Warehouse          | Roxie              | Hive       |
| Programming Languages   | ECL                | Java/Pig   |



## **Using HPCC Systems to Manage Academic Data**



# Managing the Academic Data Lifecycle

- Single platform for everything
- Supports Parallelization
- Supports Scalability
- Embedded support for external languages and libraries for specialized capability



# Managing the Academic Data Lifecycle

- Data required for scholarly data research
  - Many different sources
  - Various formats
- Aggregating these sources into a cohesive structure requires a tool that supports a data-intensive approach for:
  - Preprocessing
  - Integration
  - Analysis

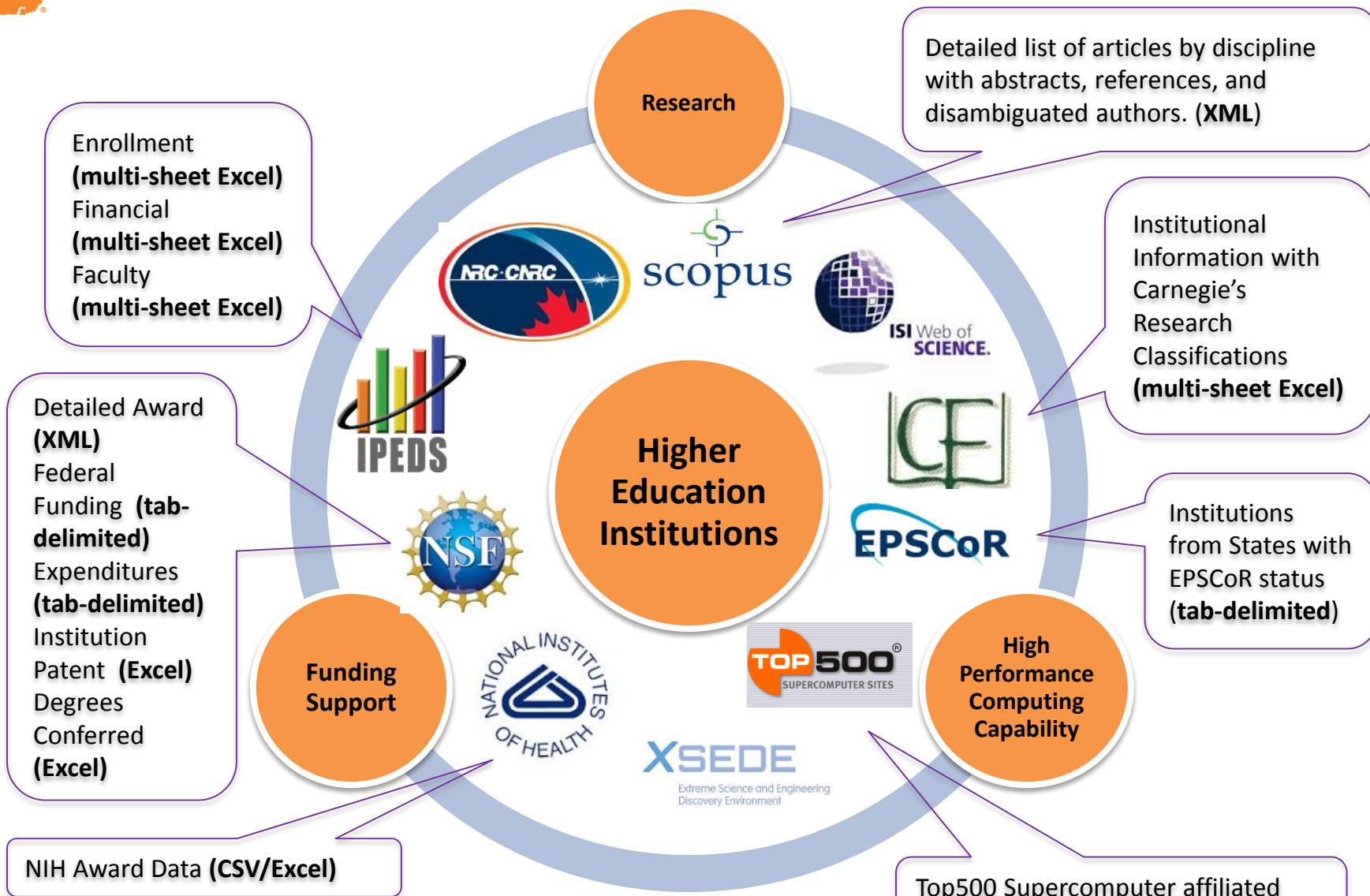


# Managing the Academic Data Lifecycle





# Managing the Academic Data Lifecycle





# Managing the Academic Data Lifecycle

## Ingesting Tabular Data

```
1 Ipeds_Layout := RECORD
2
3 STRING UNITID;
4 STRING INSTNM;
5 STRING ADDR;
6 STRING CITY;
7 STRING STABBR;
8 STRING ZIP;
9 STRING CONTROL;
10 STRING HBCU;
11 STRING MEDICAL;
12 STRING TRIBAL;
13 STRING CCBASIC;
14
15 END;
16
17 EXPORT Ipeds := DATASET('ipeds::institutional_
18 characteristics', Ipeds_Layout, CSV(HEADING(1)));
```

```
1 IMPORT $.IPEDS;
2
3 COUNT IPEDS;
4 OUTPUT(IPEDS, {unitid, instnm}, NAMED('ID_NAME'));
```



# Managing the Academic Data Lifecycle

## Ingesting XML Data

```
1 ProgramsXML:= RECORD
2 UNICODE program {XPATH('')};
3 END;
4
5
6 AwardXML:= RECORD
7
8 UNICODE awardNumber {XPATH('AwardNumber')};
9 UNICODE amountToDate {XPATH('AwardedAmountToDate')};
10 UNICODE title {XPATH('Title')};
11 UNICODE awardInstrument {XPATH('AwardInstrument')};
12 UNICODE startDate {XPATH('StartDate')};
13 UNICODE expirationDate {XPATH('ExpirationDate')};
14 UNICODE directorate {XPATH('NSFDirectorate')};
15 UNICODE nsfOrganization {XPATH('NsfOrganization')};
16 DATASET (ProgramsXML) programs {XPATH('Program')};
17 UNICODE abstract {XPATH('Abstract')};
18
19 END;
20
21 EXPORT NsfAwards := DATASET('nsfdata::nsfxml',AwardXML,XML('AwardsList/Award',NOROOT));
```





# Managing the Academic Data Lifecycle

## Cleaning Data

```
1 IMPORT Std;
2
3 EXPORT CleanForTokens(STRING s) := FUNCTION
4 sRestrictChars := ' ' + REGEXREPLACE(' [^ - A-Z 0-9 \'] ', Std.Str.ToUpperCase(s), ' ') + ' ';
5 sStripPunctEnds := REGEXREPLACE(' (-) | (-) | (\') | (\') ', sRestrictChars, ' ');
6 sRemoveNumberOnly := REGEXREPLACE(' [-0-9 \'] + (?= [])', sStripPunctEnds, ' ');
7 sCompressSpaces := REGEXREPLACE(' [] +', sRemoveNumberOnly, ' ');
8 sNormalizePossessives := REGEXREPLACE(' \'S ', sCompressSpaces, ' ');
9 sSplitContraction01 := REGEXREPLACE(' \'RE ', sNormalizePossessives, ' ARE ');
10 sSplitContraction02 := REGEXREPLACE(' \'LL ', sSplitContraction01, ' WILL ');
11 sSplitContraction03 := REGEXREPLACE(' I \'M ', sSplitContraction02, ' I AM ');
12 RETURN sSplitContraction03;
13 END;
```

```
1 IMPORT $.Ipeds;
2 IMPORT $.CleanForTokens;
3
4 $.Ipeds Standardize($.Ipeds Le) := TRANSFORM
5
6 SELF.INSTNM := $.CleanForTokens(Le.INSTNM);
7 SELF.ADDR := $.CleanForTokens(Le.ADDR);
8 SELF.CITY := $.CleanForTokens(Le.CITY);
9 SELF.STABBR := $.CleanForTokens(Le.STABBR);
10 SELF.ZIP := Le.ZIP[..5];
11 SELF := Le;
12
13 END;
14
15 EXPORT Standard_IPEDS := PROJECT($.Ipeds, Standardize(LEFT));
```



# Managing the Academic Data Lifecycle

## Linking Data

```
1 IMPORT $;
2
3 NsfIpeddsRec := RECORD
4
5 STRING unitid;
6 UNSIGNED awardnumber;
7
8 END;
9
10 NsfIpeddsRec JoinThem($.Awards Le, $.IPEDS Ri) := TRANSFORM
11
12 SELF.unitid := Le.unitid;
13 SELF.awardnumber := Ri.awardnumber;
14
15 END;
16 EXPORT NsfIpedds := JOIN($.IPEDS,$.Awards, LEFT.name = RIGHT.institution, JoinThem(LEFT,RIGHT));
```



# Managing the Academic Data Lifecycle

## Examples of Scholarly Data Links

