



How can Digital Transformation help me handle information overload?

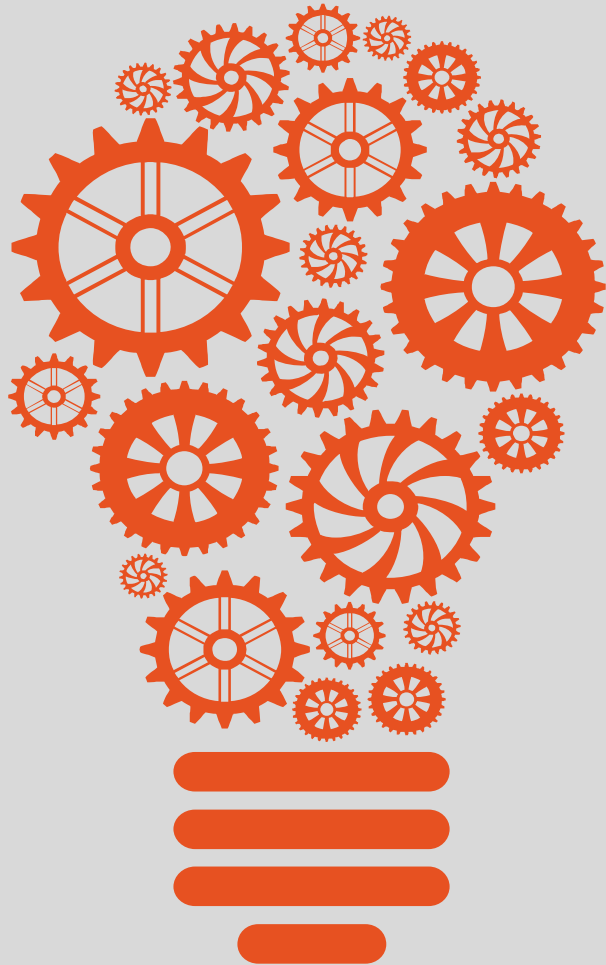
(4 Petabytes incoming!)

Dave Camden, Director
2nd October, 2019

Final V1.0



Case study – Supermajor acquisition



Digital Transformation

Initiatives must consider how vast stores of legacy information can be integrated into the future digital environment.



Case Study

This case study will look at a company acquisition and shows how a large incoming data set was handled.



Applicability

The methodologies and approaches could be applied to any digital transformation project.



The challenge

Vast amounts of structured and unstructured data.
Desire to do clever stuff with it.
How can it be done?

Holistic Information Management

Consulting

Software

Services

Oil & Gas Experience



Evergreen solutions

- Operate more efficiently
- Make better decisions
- Reduce risk
- Augment existing solutions



Flare was formed in 1998 and has been evolving and applying its holistic information management approach since then

Clients include supermajors, independents, governments, standards bodies and service companies



The project: 4 Petabytes incoming!



Business Context

- Company acquisition: A multinational “**Supermajor**” oil and gas company acquired a large oil and gas company “**OilCo**” to strengthen its portfolio



Objectives

- Merge, reconcile and re-distribute information assets
- Deletion of non-entitled datasets
- Deriving value from data



Challenges

- Limited timeframe
- OilCo had over 20 years of legacy data
 - 300m+ unstructured files (4 Pb) in two locations (3Pb and 1Pb split)
- Non-connected infrastructures

Goals and Requirements



Goals

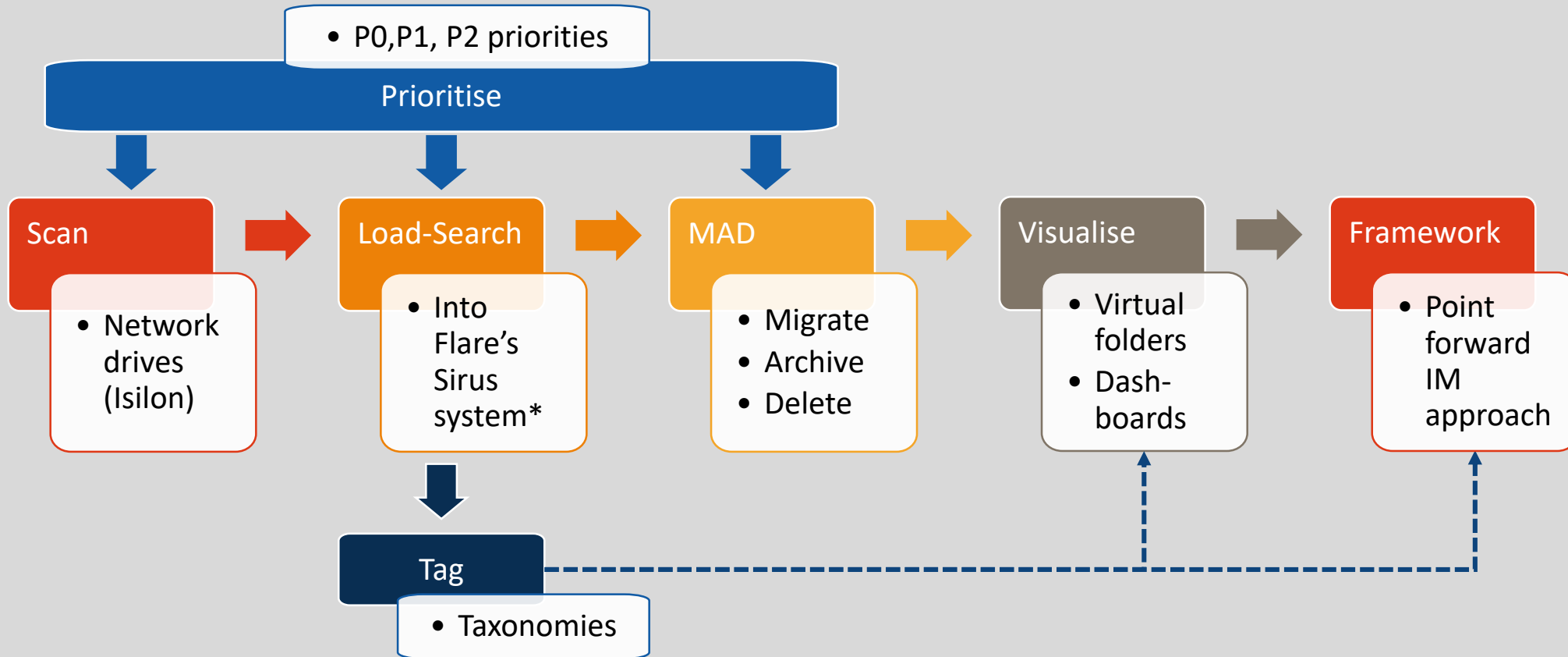
- Organise and make available high value information from OilCo
- Maintain business continuity



Requirements

- Manage project
- Communicate with all stakeholders
- Create searchable, integrated index of 300 million + items
- Find, prioritise and label 'information packages' for Migration, Archive or Deletion (MAD approach)
- Transport Migration and Archive packages to destinations
- Build new environment for incoming information

Approach



* *Sirus is Flare's graph based Information Management solution for Oil and Gas*

Disk Scanning



Infrastructure, accessibility and sheer scale meant a pragmatic approach to disk scanning was required

- Disk areas prioritised in 3 groups: P0, P1 and P2
- Disks scanned for main file/folder attributes from the Isilon system
- No checksums were calculated
- No content reading / scraping
- Scan files run through Flare's QC process to detect and fix issues

Scan Loading



The quality-controlled scans were loaded into a graph database system that allowed the original disk structures to be visualised and analysed. URLs link to the disk files.



The 300m+ items were organised based on client requirements into different disk areas on a tabbed interface. Each community can gain access to what they need. Tabs visible based on user.

Tabs

The screenshot shows a software interface with a top navigation bar containing tabs for various data management functions: Data Management, Field, Gathers, Geophys, Geophysics, GIS, N: PDD, N Drive, O Drive, Petrophysics, Petrophysics NEW, Seis, and Survey. Below the navigation bar, there are sections for 'Document folders', 'Scheduled publishing', and 'Publishing rules'. A search bar is present with the text 'Search content' and a search icon. The main interface is divided into two panels. The left panel, titled 'Disk Folder Structure', displays a table with columns for '#', 'File size', and a tree view of folder paths. The right panel, titled 'Contents and Search Results', shows a search results table with columns for 'Title', 'Author', 'Date', 'P..', 'A..', and 'File size'. The right panel also includes a search icon and a 'Tagging' dropdown menu.

Folder (1 items selected)	#	File size
0 / 64,509,724		898.56 TB
0 / 64,509,724		898.56 TB
0 / 64,509,724		898.56 TB
0 / 60,128,440		808.69 TB
0 / 44,306		1.33 TB
0 / 220,233		2.04 TB
0 / 71,727		6.93 TB
4 / 592,500		4.80 TB
6 / 1,077		9.35 GB
12 / 36,561		133.02 GB
0 / 35,737,505		49.13 TB
		95.03 GB
		136.74 GB
		18.36 GB
		18.07 TB
		7.03 GB
4 / 1,882		1.34 TB
0 / 130,971		9.78 TB
4 / 393,522		291.41 GB
4 / 18,596		79.51 GB
0 / 23		70.29 TB
5 / 2,140,806		16.84 GB
1,036 / 6,254		12.76 TB
29 / 320,379		11.14 TB
5 / 584,528		243.56 GB
1 / 113,077		7.54 TB
22 / 670,152		81.90 GB
2 / 20,072		8.98 TB
13 / 1,297,974		



Searching

- ◆ The 300m+ items were organised into different areas on a tabbed interface
- ◆ Initially searching used string matches on path/file names
 - The key benefits are speed, document counts, roll-ups and folder search/distribution
 - ***The results surprised the client – able to search the entire data set in seconds***
- ◆ Search display:
 - LISTING: the familiar ‘results listing’ with links to which folders files are in. Fine for small numbers of results.
 - DISTRIBUTION: where results number 1,000’s – 100,000’s, an innovative and highly performant ‘display search results in folders’ approach was designed by Flare. This visually highlighted which disk areas contained the bulk of the search results, facilitating the partitioning of the disk content.

Flare Folder Search Results

Number of search results File count: in folder/all files Cumulative file size

View columns Contents (13039) Treemap Pack Duplicate check Allocation Pruning Tagging

Folder (1 items selected) (search : porosity)

#	File size	Title	Author	Date	P..	A..	File size
(13039) 0 / 64,509,724	898.56 TB			02 Oct 2018			27.00 kB
(13039) 0 / 64,509,724	898.56 TB			01 Oct 2018			298.00 kB
(13039) 0 / 64,509,724	898.56 TB			2016			1.77 MB
(7486) 0 / 60,128,440	808.69 TB			2017			2.29 kB
(683) 0 / 44,306	1.33 TB			2018			2.88 kB
(49) 0 / 220,233	2.04 TB			2018			36.57 kB
(39) 0 / 71,727	6.93 TB			2016			1.54 MB
(7) 4 / 592,500	4.80 TB			2018			31.27 kB
6 / 1,077	9.35 GB			2018			32.83 kB
12 / 36,561	133.02 GB			2018			2.05 MB
(47) 0 / 35,737,505	49.13 TB			2018			1.25 kB
1 / 7,491	95.03 GB			2018			31.27 kB
32,113	136.74 GB			2018			101.43 kB
2,044	18.36 GB			2016			1.98 MB
16,511	18.07 TB			2018			3.20 MB
1,882	7.03 GB			2016			0 bytes
(21) 0 / 130,971	1.34 TB			2016			1.54 MB
(611) 4 / 393,522	9.78 TB			2018			31.27 kB
(1) 4 / 18,596	291.41 GB			2018			35.17 kB
0 / 23	79.51 GB			2018			39.08 kB
(210) 5 / 2,140,806	70.29 TB			01 Oct 2018			39.08 kB
1,036 / 6,254	16.84 GB			14 Apr 2016			1.98 MB
(19) 29 / 320,379	12.76 TB			01 Oct 2018			35.17 kB
(196) 5 / 584,528	11.14 TB			27 Sep 2018			1021.00 bytes
(5) 1 / 113,077	243.56 GB			01 Oct 2018			2.19 MB
(264) 22 / 670,152	7.54 TB			01 Oct 2018			2.19 MB
(6) 2 / 20,072	81.90 GB			01 Oct 2018			2.19 MB
(173) 11 / 1,297,974	8.98 TB			01 Oct 2018			36.73 kB

Disk Folder Structure

(13039) 0 / 64,509,724 898.56 TB

(7486) 0 / 60,128,440 808.69 TB

(683) 0 / 44,306 1.33 TB

(49) 0 / 220,233 2.04 TB

(39) 0 / 71,727 6.93 TB

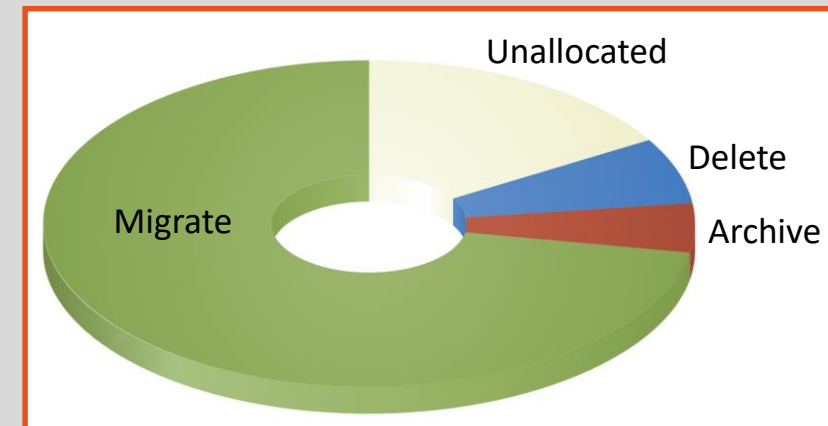
(7) 4 / 592,500 4.80 TB

6 / 1,077 9.35 GB

MAD Partitioning

- ◆ Used a drag/drop into 'basket' system to collate 'folder hierarchy lists' (areas of the disk structure)
- ◆ Basket types are 'Migrate' (🛒), 'Archive' (🛒) and 'Delete' (🛒)
- ◆ The baskets are permanent, and record
 - what information has been selected
 - what process is to be carried out
- ◆ The overall progress of each disk area is monitored to measure project progress
- ◆ The resultant lists are passed on to the IT service group to carry out the required action

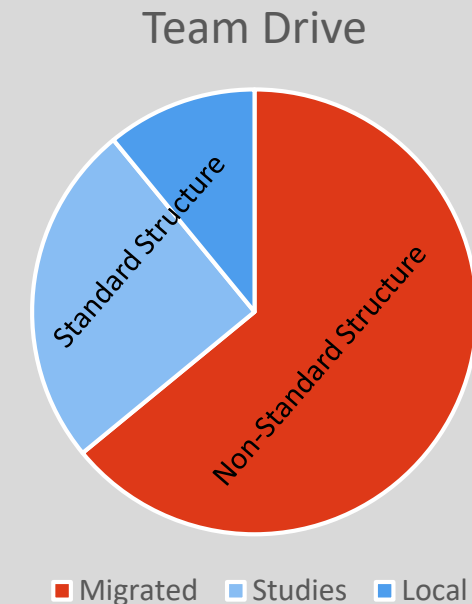
(13039) 0 / 64,509,724	898.56 TB
(7486) 0 / 60,128,440	808.69 TB 🛒
(683) 0 / 44,306	1.33 TB
(49) 0 / 220,233	2.04 TB
(39) 0 / 71,727	6.93 TB 🛒
(7) 4 / 592,500	4.80 TB 🛒
6 / 1,077	9.35 GB





Reception of Incoming Data

- ◆ Supermajor regional office received their Migrated data and loaded to new Team Drives. Windows for working files, Unix for application data
- ◆ Data for Migration to other Supermajor business regions was channelled through a centralised data handling group
- ◆ Regional office populated new 'team-drive' from various sources
 - Migrated data from the 4Pb drives – in original structure
 - Study groups – in 'company standard' structure
 - In-situ, on-going work – also in 'company standard' structure
 - Approx. 40m files



Re-structuring with Virtual Folders

- ◆ Resulting folder structure a mix of 'company standard' plus what came from Migration
- ◆ Decision made to keep 'original' folder structures on disk ...
- ◆ .. and create high-level additions in Flare Sirius using virtual folders
- ◆ A virtual folder shows content which is tagged with specific, standardised metadata
- ◆ 'Target' folders were identified and tagged using a mix of methods :
 - Manually
 - Automatically

Sirus folders representing Isilon structure

Folder (1 items selected)	#	Title	Date
AutoTagTestFiles	1 / 14,778	MapS1	15 Mar 2018
AutoTagTest01	0 / 14,778	MapS2	15 Mar 2018
SW	0 / 1		
WBS11 Extract 8	14,505 / 14,575		
SWells	2 / 200		
2D seismic	0 / 2		
Maps	2 / 2		
3D seismic			

MapS1 = Denmark '2D seismic' map EXP-G
 MapS2 = Norway 'well location' map OPS

Sirus Virtual folders

Folder	Content
Exploration	
DK	Documents, Wells, 2D Seismic, 3D Seismic
Approved	DK, Documents, Wells, 2D Seismic, 3D Seismic
OPS	Norway, Documents, Wells, 2D Seismic, 3D Seismic

Appears in

Results



Outcome so far

- Users have access to multiple data sets with one search
- Successful integration with minimal disruption
- P0 and P1 migration completed over a 6-month period
- P2 is part-complete and ongoing
- Audit record of what happened to 300m files
- New Regional Team Drive
 - has stayed operational throughout
 - has been 'live scanned' and is continuously monitored to maintain an up-to-date index
- Have built a foundation for moving forward



In Progress/Next Steps

- ◆ Index and tag team drives using standardised terms from a set of taxonomies (actually an ontology)
 - ‘Asset’ taxonomies of wells, fields, licences, countries etc.
 - ‘Context’ taxonomies based on deliverables (work products) plus relationships to Topics, Disciplines etc.
 - All taxonomies hierarchical, inter-related and include synonyms
- ◆ Tagging
 - Manual
 - Business Rules
 - Inheritance – see below
 - Automatic – joint analytics / ML program to improve accuracy and capabilities
- ◆ Continue drive monitoring to capture any changes on the file system
- ◆ Hide ‘uninteresting’ areas
- ◆ Modelling key processes
 - Identify deliverables by process step
 - Simplify tagging using drag/drop and inheritance (leveraging Flare’s comprehensive taxonomy)
 - Simultaneously track progress of project vs deliverables



Learnings (i)

- ◆ Project controls
 - Develop a bespoke way to keep a track of progress accounting for large data volumes
 - Partition the data and prioritise based on business needs – involve end-users
- ◆ Be selective to avoid information overload
 - Don't handle everything
 - Work at folder-, not file-level where appropriate
 - 'Hide' low-level application files
- ◆ Big network-based data sets (100m files / Petabytes)
 - Require a fast search system - native file systems often too slow/inflexible to be practical
 - Indexing files and moving to alternative storage takes a LOT of time

Learnings (ii)

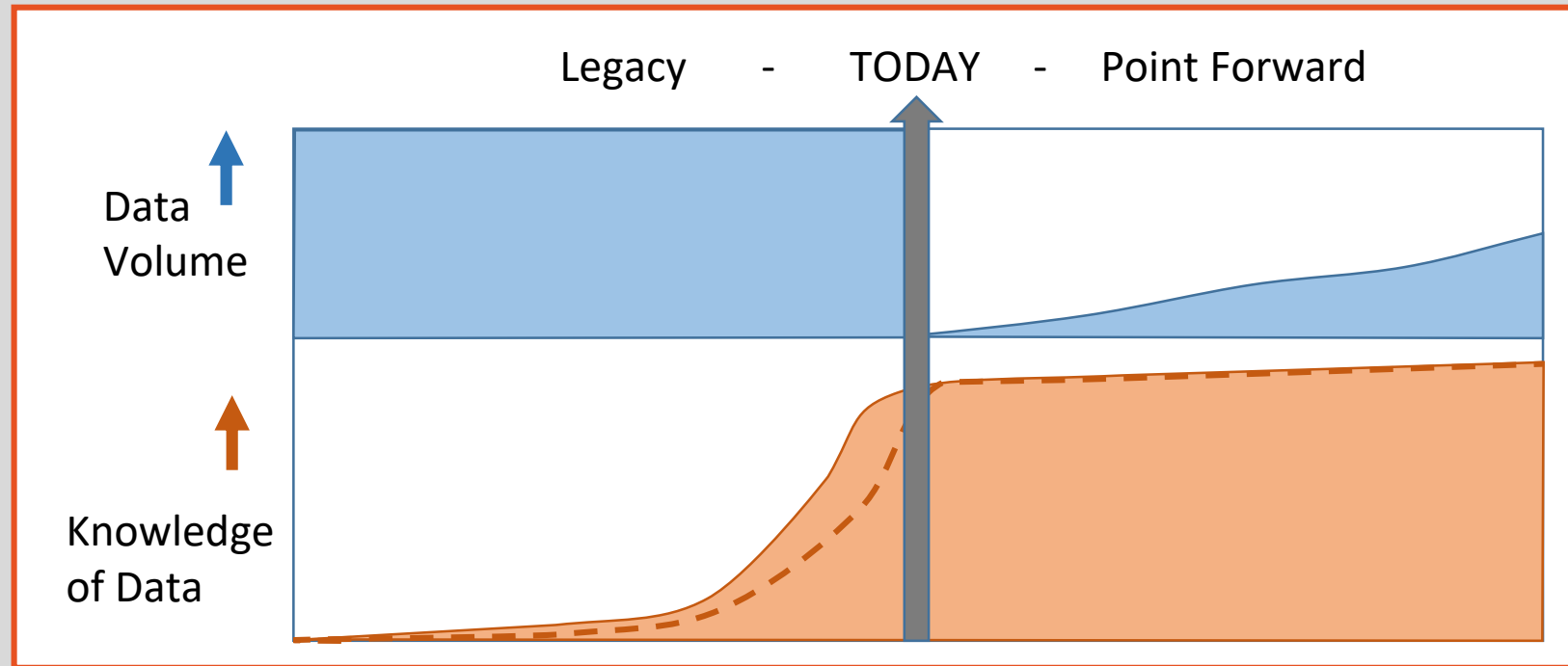
- ◆ Keep it simple
 - Milliseconds count! Use the fastest, simplest methods initially
 - Use more complex methods progressively
 - Don't aim for 100% - reasonable endeavours then improve incrementally
- ◆ Folders
 - Useful for project-based working, but need augmenting with metadata and 'smart' searching. Folders are less useful for enterprise-wide storage.
- ◆ Deletion
 - Data past their retention period or no longer useful should be deleted, but must know what you have in order to do this – good metadata is essential, especially point-forward
 - Keep a record of what has been deleted



Legacy and Point Forward Information



Legacy



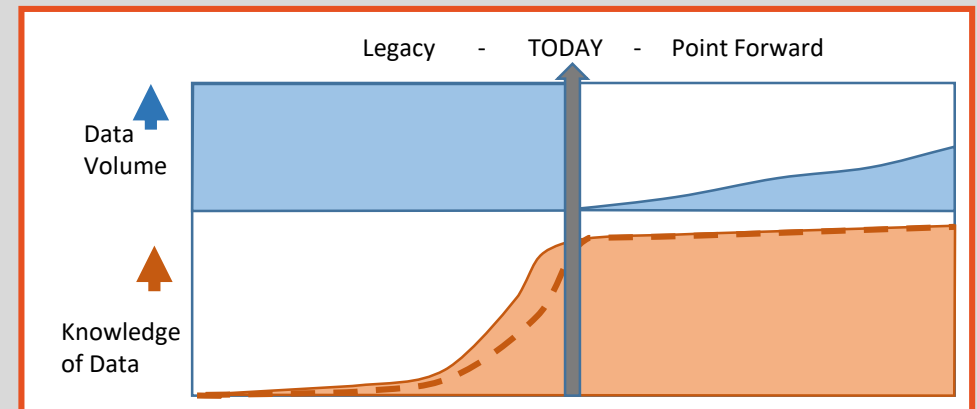
- Work so far has been on the legacy information: this system index will stay in place to provide an historical searching capability and audit record of what was done
- Other legacy/current information sources are being tagged to provide an integrated search across multiple information systems
- Large, often poorly structured data sets with little metadata
- Original creators are gone! Automate as much as possible

Legacy and Point Forward Information



Point Forward

- Creators are in place - opportunities for knowledge capture
- Opportunity to create better-structured, higher-value data sets
- Don't rely on folder structures alone – also use metadata tags and virtual folders to create alternative views, dashboards and support searching
- Understand concept of 'information objects' that are useful to the business
- Add standard metadata to support management and searching



Benefits and Enablers



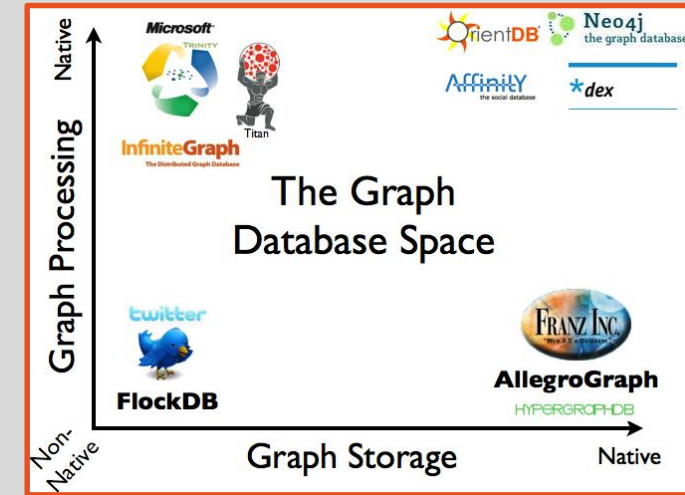
Benefits

- Unified corporate information assets
 - Teams' information needs effectively addressed
 - Rapid search capabilities
 - Custom views
 - Virtual folders
 - Information objects
 - Search 'in folders' (distribution)
 - Hiding uninteresting
 - Foundation for the future!
 - Integration, support analytics
- ◆ This project was dealing with traditional company-based infrastructure, but the approach is relevant to Cloud-based environments
- Can you easily move 4Pb of content and its metadata from AWS to Google Cloud?



Enablers

- Graph Database



- Metadata tags
- Oil and Gas taxonomies
 - Asset
 - Context
 - Other taxonomies
- Holistic IM approach



Thank you for listening. Any questions?



Flare Solutions Limited
3, Acorn Business Centre, Northarbour Road,
Cosham, Portsmouth, PO6 3TH, UK

Dave Camden
Director

d.camden@flare-solutions.com

+44 7703 234 891

+44 1892 785 007

Tel: +44 203 397 7766

Fax: +44 870 460 2543

enquiries@flare-solutions.com

www.flare-solutions.com