



Data Integration – Industry Practices and Opportunities for Improvement

PHMSA Risk Model Working Group Data Team (RMWG) Data & Risk Tolerance/Acceptability

> Wednesday, March 8, 2017 Peter Veenstra, TRC

Outline



- Background\Context
- Current Trends in Data Management
- RISK Models and modeling GIS Data Uncertainty
- Making this accessible to Small Operators
- Observations

Background



- Principal GIS Technologist
- Involved in GIS since 1992
- GIS and Remote Sensing Background
 - Tobler's first law of Geography
- Analyst, Programmer, Architect
- Involved in designing/implementing ISAT, PODS, APDM, UPDM, PODS Next Gen
- Current PODS BOD Member, Chairperson of PODS Next Gen





Context



- Are we looking at "Data Integration Practices and Opportunity" within the industry
 - For the industry as a whole ...?
 - For the regulators ...?
 - For the operators individually or as a whole ...?
- Difficulty in sharing of information and knowledge
 - Release of perceived competitive advantage and/or proprietary knowledge
 - Risk of exposure and damage to image and perception
 - Cost of doing business (is it worth it to spend this kind of money?)
 - Balance between under or over-proscribing regulatory requirements and compliance

Part 1 - Current Trends in Data Management



- Relational Databases and Pipeline Data Models
- Geography rears its ugly head
- Scalability for small operators
- Systems built in silos
- Operational Excellence

Part 1 – Relational Databases and Pipeline Data Models





2008			
2009	PODS 5.0 Released		
2010	PODS ESRI Spatial 5.0 Released	APDM 5.0 Released	
2011	PODS/ESRI Spatial 5.1 Released		
2012	PODS/ESRI Spatial/Open Spatial 5.1.1 Released		
2013	PODS 6.0 w/ Modularization Released	APDM 6.0 Released	
2014		APDM No Longer Supported	

So here we are ...













Part 1 – Relational Databases and Pipeline Data Models



- Relational Databases and Pipeline Data Models
 - -This seems to be our only option
 - RDBMS\GeoSpatial Data Models are excellent for storing the location of the centerline and managing the stationing LRS associated with that
 - RDBMS are well suited to managing ASSETS, but not necessarily CONDITIONS
- Flexibility for storing data easily
 - Geographic Coordinates vs Linear Referencing
- Cannot support unstructured data, un-validated data or the massive amounts of data
- Excellent for Fit-for-purpose data and applications Centerline and Assets
- Also excellent for tracking data edits via history and versioning (ESRI Geodatabase Versions/Long Transactions)

Part 1 – Geography rears its ugly head



 Geography in this context represents different parts of the United States and the World

- Different geo-political agencies look at standards and methods differently (regulations – US, Canada, UK, Europe, Russia, Australia)
- Geography allows us to integrate things in a common coordinate system but forces us to deal with different things
- Each geographic area has a different set of RISK (some are similar) but there is no one-size fits all



ource: Energy Information Administration, Office of Oil & Gas, Natural Gas Division, Gas Transportation Information System



Part 1 – Conversely, Geography is beautiful



- Coordinates unify different geographic information into a common plane of reference
- Creates inherent relationships between objects in space
 - Tobler's Law Of Geography
- Data can be overlaid with each other to form new data in a different manner than the standard SQL table join
- This is the basis for most 'classic' risk implementations





c) result of

subtractive overlay



Part 1 – Scalable for small operators



- Access to GIS systems, PODS or other GIS functionality
 - Complicated systems, detailed knowledge, specific requirements
- They are still bound by the same regulations and standards
- How do we support smaller operators
- Other industries have set up smaller data models and collectively built open source technologies for data management
- Agencies within this sphere could provide source data as services (Federal Geospatial Data Committee)
- Still not granular or specific enough for the pipeline industry but the model works

Part 1 – Systems built in Silos



- Risk requires data from many different systems
 - Work order management, document management, SCADA, Cathodic Protection
 - Pipeline GIS may not be the place to store this
 - Systems must be integrated after the fact
- Systems built on PMP
 - "My job, done now, on time, on budget, in isolation of the downstream effects".
 - Systems are compartmentalized and not built with the end-in-mind
 - There is little data exchange, data flow, data hand-off (this is HARD to do)
 - Varying levels of system complexity, completeness, and data quality/structure
 - If you are lucky there are coordinates being saved



Part 1 – Operational Excellence



- Lack of Operational Excellence for the "System of Record"
 - Fit for purpose data with quality sufficient for the task Data Excellence
 - Are end users challenged to get the right data into the system?
 - Is the system providing the end users with the right data?
 - Younger generation quality on demand (ADD)
- Organizations are relying on 'white' expertise
 - Subject Matter Experts know the data better than the systems how do we capture this?
- More is being done with less (continued separate of church and state)
 - Pipeline Integrity Engineers are being forced to learn GIS
 - GIS users need to understand how pipeline integrity data management works
- Digital Workflow is incomplete
 - Data is missing or is not collected or delivered
 - Field to enterprise work-flow is still kludge
 - Once data are gathered or collected, how are they integrated?
 - Many operators are gathering, massaging and formatting the same data (but separately)
- IT gets in the way
 - Need to enable solutions

Part 2 - Potential Trends for Pipeline Data Management



- The Cloud
- Unstructured Data (Big, IOT)
- Services Architecture (API)
- Data Integration\Lakes\Exposure
- Machine Learning

Adoption of Technology (Pipeline Industry)





https://setandbma.wordpress.com/2012/05/28/technology-adoption-shift/



- Scalability on demand
 - Development environments experimentation, learning, performance
 - Built on Open Source Technologies (organic infrastructure and Development)
- Collaboration
 - Shared computers and data stores (workspaces and servers)
 - migration, training, testing, data quality checking
 - Remove bureaucracy of IT
 - Rapid deployment and sizing without CAPEX process
 - Crowd (Company/Organizations) sourcing for data sharing
 - Build data storages and entry procedures for capturing shared data
 - Provides services (SMS)
 - Provide notifications
 - Provides map services



Data Services

- -High performance big table (tall and wide) performance (use SQL)
- Built in RDBMS technology and services (not just Oracle and MS SQL Server) but RDS (POSTGIS)
- Big Data (Hadoop, MongoDB, Google)
- Deep Storage (AWS S3)
- Notification Platform
 - Built in services architecture SMS (information on your phone)
 - Notifications something has changed, please respond (devices, location)
 - Push/Pull I have changed something, what is happening here?
 - Post to Default (ESRI Spatial Geodatabases)

Part 2 - Services Architecture



- Individual Cloud Based Components that talk to each other (with common technology) rather than monolithic software systems
- Constant and iterative data processing from local/disparate/siloed systems to shared data stores
- Processing on demand (as-needed) (FME Cloud)
- Once data is formulated serve it to systems that can easily consume it and provide access to it in a simple manner

Part 2 – Services Architecture





Compute



Storage



Networking & Content Delivery



Analytics



Messaging



Developer Tools



Artificial Intelligence



Business Productivity



Database

≣≡

Management Tools



Mobile Services



Desktop & App Streaming



Migration



Security, Identity & Compliance



Application Services



Ease of Use and Interaction (Value not complexity)





Connect the field to the office Grant your team the ability to access and update job specific information in real time on any mobile device.



Track all your progress to ensure your projects are traceable, verifiable, and complete.



Support for defined and custom workflows. *The nature of pipeline data and integrity management is highly flexible.*



Scalability to support any size project. Leveraging Cloud infrastructure and a managed secure cloud model.



Keep up to date in real time with changes and alerts that come down the pipeline. *Did my risk profile change?*



Affordable and quick to implement. No lengthy implementation cycle and offered in a SaaS infrastructure and with a simple pricing model.

Part 2 – Unstructured Data





Part 2 – Big Data



- Large amount of small computers holding JSON format data
- Query layer for retrieve and organize
- Reminiscent of main-frame computers (fast, but required black magic interface)
- Cloud providers support OOTB toolkits (Native SQL Interface)
- Requires the inclusion of geography or tags







- There are two parts to Hadoop:
 - HDFS split the data, put it on different nodes, replicate and manage it
 - MapReduce batch processing over an increasing amount of data
- Physically there are many parts
 - Scale out instead of scale up
 - Many computers/servers in a managed cluster
- Amazon, Mapr, Cloudera







- Hadoop uses key/value pairs as its basic data unit and is flexible enough to work with the less-structured data types
- Data can originate from any form, but it eventually transforms into key/value pairs for processing functions to work on
- Text Files
 - XML, txt, JSON
- Binary Files
 - Video dissected frame-by-frame



- Functional programming (MapReduce) instead of declaritive queries (SQL)
 - You query data by stating the result you want and let the database engine figure out how to derive it
 - MapReduce you specify the actual steps in processing the data, which is more analogous to an execution plan for a SQL engine
 - SQL has query statements, MapReduce has scripts and code
- Better get some Java (or Python)
- Pipeline Use Cases
 - Tracking location on a massive scale popcorn trail density,
 ILI density en masse (analytics vs other geography)

Part 2 – Big Table



- Massive number of rows and columns
- Dynamically segmented or merged tables
- Supported by providers or in private cloud
- Can be linked via standard SQL with other Big Tables
- Can store geometry via GeoJSON
- Allows for storage of complex data structures within an attribute
- Requires a mechanism for getting it out and displaying it
- LAYERLESS!!!
- Use cases Previous demonstration

Part 2 – Data Integration



Data Lake

- A data lake is a collection of storage instances of various data assets additional to the originating data sources. These assets are stored in a near-exact, or even exact, copy of the source format. (www.gartner.com/itglossary/data-lake/)
- Data Warehouse
 - A data warehouse is a storage architecture designed to hold data extracted from transaction systems, operational data stores and external sources. The warehouse then combines that data in an aggregate, summary form suitable for enterprise-wide data analysis and reporting for predefined business needs (http://www.gartner.com/it-glossary/datawarehouse/)
- Put the data into the hands of those who need to use it, in a format that they need to use it
- Publish it for consumption



Part 2 – Data Integration





Part 2 – Machine Learning and BI



- Machine Learning
 - Advanced machine learning algorithms are composed of many technologies (such as deep learning, neural networks and natural-language processing), used in unsupervised and supervised learning, that operate guided by lessons from existing information. (http://www.gartner.com/it-glossary/machine-learning/)
- Business Intelligence
 - Business intelligence (BI) is an umbrella term that includes the applications, infrastructure and tools, and best practices that enable access to and analysis of information to improve and optimize decisions and performance. (www.gartner.com/itglossary/business-intelligence-bi/)

Part 2 – Machine Learning



Machine Learning

- Advanced machine learning algorithms are composed of many technologies (such as deep learning, neural networks and natural-language processing), used in unsupervised and supervised learning, that operate guided by lessons from existing information. (http://www.gartner.com/it-glossary/machine-learning/)
- Requires data, systems and business\problem statement
- Similar to classification of remotely sensed imagery
 - Supervised and Unsupervised classification
- Classification of ILI Data and Alignment of ILI Data to GIS
- Identification of trends in incidents and pipeline defects compared to geographic conditions

Part 2 – Business Intelligence

Business Intelligence

- Business intelligence (BI) is an umbrella term that includes the applications, infrastructure and tools, and best practices that enable access to and analysis of information to improve and optimize decisions and performance. (www.gartner.com/it-glossary/business-intelligence-bi/)
- Requires data, systems and business\problem statement
- Allows for the visualization and analysis of data in ad-hoc fashion
- Tableau, Tibco Spotfire, Amazon QuickSite, ESRI Insight
- SME review and visualization for potential risk

Data-Driven Documents







Part 3 - Risk Models and Data Uncertainty



RISK Models

- Qualitative and Quantitative
- Machine Driven and SME Driven
- -Standard Methodology
- Require data
- Data Uncertainty
 - Uncertainty
 - Methods for modeling uncertainty
 - Modeling GIS data quality using uncertainty\







Network/Spatial-Based Risk Analysis



- GIS based Risk Application
- Geo-Aggregation
- Dynamic Segmentation
- Calculate Risk Score
- Risk Results Visualization













Part 3 – Risk Models



- Abstraction of reality (is a model)
 - Your tomato is my tomatoe
 - Geographic location changes the inputs
- Standard Methodology
 - Gather, convert, segment, calculate, organize, slice/dice
- Qualitative and Quantitative
 - Zero Incidents, Zero Cost, Zero loss of life
- Machine Driven and SME Driven
 - Somethings require a computer
 - All result require validation against experience and actual data (incidents, leaks)
- Require data, require complete data, require accurate data

- Or need a mechanism to deal with unknown, missing or uncertainty

Part 3 - Data Uncertainty







- Is data quality a measure of uncertainty?
- How do we measure quality and uncertainty?
- In PODS we model uncertainty:
 - Via attributes: SourceGCL, ValidityTolerance, DocumentCrossRef
 - Via linear events: CenterlineAccuracy
- Methods for modeling uncertainty
 - Bayesian (Probability), Dempster-Shafer (Evidence), Fuzzy Theory (Incomplete/Imprecise)
- Modeling GIS data quality using uncertainty
 - Location Quality/Precision/Accuracy
 - Attribute validity

Part 3 - Data Uncertainty



- Risk = Consequence x Likelihood
- Consequence is critical part of risk.
- Location is a critical part of consequence.
- Causality is an important part of likelihood.
- Accuracy of spatial dataset are suspect esp. public or purchased data sets.
- Incorporate Fuzzy Surfaces within data analysis.
- How you analyze or deal with this uncertainty might determine the priority of your actions



Part 3 - Data about the data



- How do we model what we don't have or do not know
 - Critical data is missing or unknown can we overcome this in our models?
- Data is not available
 - How do we get it?
 - Legacy conversion
 - Field Data Collection
- Data is stored in other systems
 - How do we integrate it?
 - Does it share a common coordinate reference? Does it share a tag or name reference?
- Everyone has their own way of doing it (different data owners)
 - Is there a way to make standardization easy, accessible and well understood?
- Data is stored in different frame of reference
 - All data will require some transformation to make it useful
 - Data has different structures how do we merge it?
- Data has different patterns that go beyond coordinates and PK-FK relationships (Semantics, Ontology)

Sidebar - Post to Default (Gateway to Information Management)

- ESRI Spatial Geodatabases have a technology called versioning
- Versions are long transactions that track the adds/deletes to a database
- At the root level is the DEFAULT version of the system/database
- When data is posted to DEFAULT it becomes accessible to standard queries against the base table
- Triggers and other technologies can be leverage to track 'post to default'
- This starts the data and information management process
 - Post to other systems
 - Perform QA/QC
 - Start processing and integration workbenches
 - Build other data stores





Part 3 - Small Operators



- How do we engage small operators in the GIS and RISK process?
 - Allow for flexibility on data management
 - Provide GIS solutions and RISK models for different levels of complexity
- PODS Next-Gen PODS LITE
 - Need to provide the minimum data set for a pipeline operator to understand the position of the pipeline centerline and to operate their pipeline safely
 - Provides the minimum required data set for submitting a complete NPMS submittal

Part 3 – PODS LITE (Part of PODS Next Gen)



Documentation

Logical Data Model Data Dictionary Governance (Model Usage, Editing Standards, Data Content Specification)

Reference Modes

Continuous (2d/3d), Interrupted (2d/3d) Referent Point and Offset (Milepost) XYZ Odometer

Software

Schema / Template Creation QA/QC Module Data Loader Module approval and documentation

Best Practices

Schema/Module Specification, Definition, Creation and Validation App Specification, Creation and Validation Managing History, Work Orders, Documents, Re-Routes and Activities Feature/Condition Provenance Managing Metadata Module Management Data Exchange





An implementation pattern is the sum of the following inputs for implementing PODS Next Gen database including the database type, the spatial storage type, and how edits are performed/managed.

Implementation Pattern Focuses on:

- Database Platform support for multiple RDBMS (Oracle, SQL Server, PostGres)
- Spatial Data Storage Type storing spatial data through native types (Geometry, SDO_Geometry, ST_Geometry) or the geodatabase (SDE_Geometry, ST_Geometry), or OGS WKT LineString
- Editing Paradigm Versioned, ArcObjects, Native SQL

Part 3 – PODS LITE (Implementation Patterns)





RDBMS Spatial Data Types: **Oracle Spatial** SQL Server Spatial



Hybrid



RDBMS Spatial Types ESRI Geodatabase





Relationa E Coordinates managed

in database tables

PostGIS native data format

- **RDBMS** Oracle/SQLServer/PostGIS, with APR, Shape column in ST_Geometry, Spatial SQL Library used to update data, Non-Versioned
- **Hybrid** Oracle/SQLServer/PostGIS, ArcSDE to store APR (centerline as feature class), Event/Feature tables stored in Relational Tables, Shape columns in either SDO or ST_Geometry, Editing is both ArcObjects (Centerline), and SQL (Event/Feature Tables), Partially Versioned
- Geodatabase Oracle/SQLServer/PostGIS, APR & Event/Feature data stored as feature classes, Shape column store as SDO or ST_Geometry, fully Versioned
- **Open Source** PostGIS, APR & Event/Feature data stored as relational tables, ST_Geometry (native spatial type), not Versioned.
- Relational Oracle/SQLServer/PostGIS, APR & Event/Feature data stored as relational tables, OGC WKT Linestring in Text/VarChar



Part 3 – PODS LITE (Modeling Approach)







- Conceptual Model
- Model Sections
- Reading the Drawing

Conceptual Model





A conceptual model is a representation of a system, made of the composition of concepts which are used to help people know, understand, or simulate a subject the model represents.

Next Gen Design Principle:

The system of record for pipeline centerlines for all pressurized containment assets, for the transport of product, for the safe operation of the pipeline, and to mitigate the potential consequence of failure.

Logical Model Design – Abstract Classes (I)

Abstract Classes





 Implementation of abstract classes to physical tables/ FeatureClasses





- HasLRSLF ?
- Coordinate based or LRS-based?



X**YQO**/I

Yes





Valve				
editResponseCL <d> (NN)</d>	positionSourceCL <d> (NN)</d>	locationCL <d> (NN)</d>	maintainLRSLF <y n=""> (NN)</y>	
Absolute	Coordinate Located	Offline	Yes	



- Logical Design:
 - -Hierarchy

PODS Next Gene



- Logical Design:
 - -MetaData





Logical Design:

–LinearReferencingSystem (LRS)







Ascending (with measure) Descending (against measure)

IMPORTANT

The LRS tables defined above are based on the ESRI ArcGIS for Pipeline Referencing (APR) Location Model tables.

Implementing the LRS Location Model tables within PODS Next Gen is optional.

Implementing LRS within PODS Next Gen does not require adoption or use of the ESRI APR software tools. However, PODS next Gen has been designed to work seamlessly with APR.



- It is all about data, the systems that store and present data
 - I either don't have the data or I have too much data or I don't trust my data (for sure I don't trust your data)
 - Data is Siloed (lost, not digital, within organizations, between organizations)
- Data Accuracy is Missing
 - Focus on data accuracy for the most important part of the system The Centerline and the Assets
- Geography becomes a beautiful thing (planar geometry and coordinates)
 - Allows for spatial interactions
 - But not always business interactions
- Location and Tags are the common reference systems
 - If I have a coordinate I can locate something in space
 - If I have a tag then I can tie two seemingly disparate pieces of information together



- True Operational and Data Excellence
 - Build and maintain the system of record
 - Turn everyone into data hunters, gatherers and consumers
 - Notify people when things change
 - Provide constant access to 'what is here' information
- Commoditize of Information Technology (Availability)
 - "Put your stuff in the cloud because we are more tired of your IT people than you are"
 - "Our IT department would never allow us to do that ..."
 - "More people work on security (for any cloud provider) then you have working for your entire company"



- Form a Data Sharing Alliance
 - Leverage the cloud to store and provide service based access to commonly shared and publically available data
 - Provide uncertainty level indexes to these data based on spatial position
 - Provide concurrency and provenance (pedigree) metadata for attribute data
 - Crowd-source data in shared corridors (line crossing, structures)
 - Work with industry groups (API, INGAA, GTI, NACE, CGA) to support and fund (host) these initiatives
 - -This is where standards come into play
 - Charge modicum of cost for use of data services



- Doubt, there is a one-size-fits-all solution
- All, some or most of the parts, form the solution

BIG Data	Machine Learning	Business Intelligence
Cloud Technology	Services	Uncertainty Management
Data Warehouse	SME	Qualitative
PODS	Quantitive	GIS





Questions?

Thank you

Peter Veenstra

P: 816-820-7841

E: pveenstra@trcsolutions.com

TRC's Guiding Principles



Our Mission

We understand our clients' goals and embrace them as our own, applying creativity, experience, integrity and dedication to deliver superior solutions to the world's energy, environment and infrastructure challenges.

Our Vision

We will solve the challenges of making the Earth a better place to live – community by community and project by project.

