

Data Dependency Network Analysis in Systems of Systems

May 1, 2018

Reggie Cole Senior Fellow Lockheed Martin Corporation reggie.cole@Imco.com

Messy Data Has Been Around for a Long Time – But the Problem is Particularly Acute for SoS & IoT



In 1963, Hafley and Lewis first discussed the analysis of <u>messy data</u>, which they defined simply as data that cannot be effectively consumed and used in the intended way.

The last 50 years has seen a continual growth in data, across an increasingly complex set of systems and enterprises.

This inherent messiness is particularly problematic for SoS & IoT, where constituent systems produce data for purposes that might be different than the intended uses of consuming systems.

So, how can we start to get our hands around this problem?

Messy Data in SoS & IoT — Effects & Mitigation





	Data Perspective	User Perspective
Context	Spelling / Formatting Error	Inaccessible Information
Independent	Missing / Duplicated Data	Insecure Information
	Incorrect / Unallowable Value	Irretrievable Information
	Outdated Data	Difficult to Aggregate
	Syntax / Constraint Violation	Transformation Errors
	Domain Constraint Violation	Unreliable Information
	Business Rule Violation	Irrelevant Information
Context	Company Policy Violation	Inconsistent Information
Dependent	Regulatory Violation	Incomplete Information
	Local Constraint Violation	Difficult to Process
		Difficult to Understand

Two Key Mitigation Strategies



Messy Data is a Key Source of System Degradation

But... What is the right mix of cleansing and governance, and how should the strategies be applied to a particular SoS?

ELERNSING

DATA GOVERNANCE

Govern

Cleanse Incoming Data

How Can We Apply Modeling to the Problem?

<u>Modeling Question #1</u> — Given a specific SoS architecture, what are the effects of messy data on the SoS?



<u>Modeling Question #2</u> — How should we apply governance and cleansing strategies to optimize cost and performance?

Development of the DDNA Model

Garvey and Pinto established the foundation for functional dependency network analysis in their 2009 paper.



Guariniello and DeLaurentis applied FDNA to systems-of-systems operability in their 2013 paper.





Applies the model to <u>data</u> <u>dependency</u> and <u>data</u> <u>quality level of provider</u> This paper extends the FDNA approach to handle data dependencies and their effect on system-of-systems operability.



Model Details — Variables & State Transition



- SE = System Self Effectiveness (From 0 to 100)
- $0 = System \ Operability \ (From \ 0 \ to \ 100)$

DQL = Output Data Quality Level (From 0 to 100)

- $\rho = Incoming \ Data \ Cleansing \ Effectiveness \ (From 0 \ to \ 1.0)$
- $\mu = Data \ Governance \ Effectiveness \ (From 0 \ to \ 1.0)$

 $\varphi = Operability Strength (From 0 to 1.0)$

 $\alpha = Strength of Data Dependency (From 0 to 1.0)$

 $\gamma = Contextual Alignment Factor (From 0 to 1.0)$

 $\beta(x) = Criticality of Dependency$

$$\beta(x) = \begin{cases} 0 \text{ if Provider Operability} < x \\ 1 \text{ if Provider Operability} \ge x \end{cases}$$

Where:

x = Minimum Operability Threshold for Provider System

For systems with no dependencies: O = SE

Since the system has no external dependencies, the system is in complete control of its own operability. In the diagram, the operability for System C and System D are governed by this equation.

For systems with N dependencies (where N > 0):

$$O = \prod_{i=1}^{N} \beta(x)_{i} \cdot \left[\sum_{i=1}^{N} \left(\frac{\alpha_{i} \cdot [\gamma_{i} \cdot DQL_{i} + \rho \cdot (100 - \gamma_{i} \cdot DQL_{i})]}{N} \right) + SE \cdot \left(1 - \sum_{i=1}^{N} \frac{\alpha_{i}}{N} \right) \right]$$

In the diagram, the operability of System B is a factor of its own self effectiveness, the quality of data being provided, the strength of dependency on the data being provided, and whether any of its predecessors have an operability below their respective operability thresholds.

For each system the quality of data being provided (DQL) is a factor of its own operability and the effectiveness of its own data governance.

$$DQL = \varphi \cdot O + \mu \cdot (100 - \varphi \cdot O)$$

The *Total System Performance* of the system-of-systems with *M* constituent systems is calculated as follows:

$$Performance = \sum_{k=1}^{M} w_k \cdot O_k$$

Where:

Model Details — System Operability



 $\varphi = Operability Strength (From 0 to 1.0)$

 $\alpha = Strength of Data Dependency (From 0 to 1.0)$

 $\gamma = Contextual Alignment Factor (From 0 to 1.0)$

 $\beta(x) = Criticality of Dependency$

(0 if Provider Operability < x $\beta(x)$ 1 if Provider Operability $\geq x$

Where

x = Minimum Operability Threshold for Provider System

For systems with no dependencies:

Since the system has no external dependencies, the system is in complete control of its own operability. In the diagram, the operability for System C and System D are governed by this equation.

For systems with N dependencies (where N > 0):



System Operability is the key observable state. It is calculated for each node based on:

SE = System Self Effectiveness (From 0 to 100)

- DQL = Output Data Quality Level (From 0 to 100)
 - $\rho = Incoming Data Cleansing Effectiveness (From 0 to 1.0)$
 - $\alpha = Strength of Data Dependency (From 0 to 1.0)$
 - $\gamma = Contextual Alignment Factor (From 0 to 1.0)$

 $\beta(x) = Criticality of Dependency$

$$\beta(x) = \begin{cases} 0 \text{ if Provider Operability} < x \\ 1 \text{ if Provider Operability} \ge x \end{cases}$$

Where:

x = Minimum Operability Threshold for Provider System

$$O = \prod_{i=1}^{N} \beta(x)_{i} \cdot \left[\sum_{i=1}^{N} \left(\frac{\alpha_{i} \cdot [\gamma_{i} \cdot DQL_{i} + \rho \cdot (100 - \gamma_{i} \cdot DQL_{i})]}{N} \right) + SE \cdot \left(1 - \sum_{i=1}^{N} \frac{\alpha_{i}}{N} \right) \right]$$

Model Details — Output Data Quality Level



- SE = System Self Effectiveness (From 0 to 100) $0 = System \ Operability \ (From \ 0 \ to \ 100)$ DOL = Output Data Quality Level (From 0 to 100)
- $\rho = Incoming Data Cleansing Effectiveness (From 0 to 1.0)$ $\mu = Data Governance Effectiveness (From 0 to 1.0)$
- $\varphi = Operability Strength (From 0 to 1.0)$
- $\alpha = Strength of Data Dependency (From 0 to 1.0)$ $\gamma = Contextual Alignment Factor (From 0 to 1.0)$
- $\beta(x) = Criticality of Dependency$
 - (0 if Provider Operability < x $\beta(x)$ 1 if Provider Operability $\geq x$
- Where:
- x = Minimum Operability Threshold for Provider System

For systems with no dependencies:

O = SE

Where:

Since the system has no external dependencies, the system is in complete control of its own operability. In the diagram, the operability for System C and System D are governed by this equation.

For systems with N dependencies (where N > 0):

 $O = \prod_{i=1}^{N} \beta(x)_{i} \cdot \left| \sum_{i=1}^{N} \left(\frac{\alpha_{i} \cdot [\gamma_{i} \cdot DQL_{i} + \rho \cdot (100 - \gamma_{i} \cdot DQL_{i})]}{N} + SE \cdot \right) \right|$

In the diagram, the operability of System B is a factor of its own self effectiveness, the quality of data being provided, the strength of dependency on the data being provided, and whether any of its predecessors have an operability below their respective operability thresholds.

For each system the quality of data being provided (DQL) is a factor of its own operability and the effectiveness of its own data governance.



 $w_{i} = Performance Contribution of System k$

Then the output *Data Quality Level* is calculated based on:

- O = System Operability (From 0 to 100)
- $\mu = Data Governance Effectiveness (From 0 to 1.0)$
- $\varphi = Operability Strength (From 0 to 1.0)$

 $DQL = \varphi \cdot O + \mu \cdot (100 - \varphi \cdot O)$

Model Details — Overall SoS Performance



- SE = System Self Effectiveness (From 0 to 100) O = System Operability (From 0 to 100) DOL = Output Data Quality Level (From 0 to 100)
- $\rho = Incoming Data Quality Level (170m 0100 100)$ $<math>\rho = Incoming Data Cleansing Effectiveness (From 0 to 1.0)$ $<math>\mu = Data Governance Effectiveness (From 0 to 1.0)$
- $\varphi = Operability Strength (From 0 to 1.0)$
- $\alpha = Strength of Data Dependency (From 0 to 1.0)$
- $\gamma = Contextual Alignment Factor (From 0 to 1.0)$
- $\beta(x) = Criticality of Dependency$
- $\beta(x) = \begin{cases} 0 \text{ if Provider Operability} < x \\ 1 \text{ if Provider Operability} \ge x \end{cases}$
- Where:
- x = Minimum Operability Threshold for Provider System

For systems with no dependencies:

O = SE

Where:

Since the system has no external dependencies, the system is in complete control of its own operability. In the diagram, the operability for System C and System D are governed by this equation.

For systems with N dependencies (where N > 0):

 $O = \prod_{i=1}^{N} \beta(x)_i \cdot \left| \sum_{i=1}^{N} \left(\frac{\alpha_i \cdot [\gamma_i \cdot DQL_i + \rho \cdot (100 - \gamma_i \cdot DQL_i)]}{N} \right) + SE \cdot \left(1 - \sum_{i=1}^{N} \frac{\alpha_i}{N} \right) \right|$

In the diagram, the operability of System B is a factor of its own self effectiveness, the quality of data being provided, the strength of dependency on the data being provided, and whether any of its predecessors have an operability below their respective operability thresholds.

For each system the quality of data being provided (DQL) is a factor of its own operability and the effectiveness of its own data governance.

 $DQL = \varphi \cdot 0 + \mu \cdot (100 - \varphi \cdot 0)$ The Total System Performance of the system-of-systems with *M* constituent systems is calculated as follows:

 $Performance = \sum_{k=1}^{M} w_k \cdot O_k$ $w_k = Performance Contribution of System k$

Finally, the overall SoS performance is calculated based on the system operability of each node in the system

$$Performance = \sum_{k=1}^{M} w_k \cdot O_k$$

$$w_k = Performance \ Contribution \ of \ System \ k$$

NetLogo Implementation of the DDNA Model



How NetLogo Entities Are Used in DDNA:

- ✤ Turtles Systems
- ✤ Patches Not Used
- ✤ Links Interfaces

NetLogo is an open source agent-based modeling tool based on turtles, patches and links:

- Turtles Modeled Entities
- Patches Location Entities
- Links Relationship Entities



Putting it Together in the NetLogo Simulation



Calculation of System States



Identification & Mitigation of "Bad Actors"



As we adjust the <u>Self Effectiveness</u> of each system in the SoS, we are able to identify systems whose data quality level has the most profound effect on the SoS as a whole.

Then we can start to think about modeling mitigation strategies for addressing this SoS "architectural weakness."

By improving the data cleansing capabilities in just two systems, we see an overall mitigation of the problem.



Summary

- Messy data is a big problem for systems of systems
- There are strategies for mitigating its effects
- The DDNA model provides a tool for understanding the effect of messy data on systems of systems
- The DDNA model provides a tool for optimizing the strategy for mitigating the effect of messy data in systems of systems

