

# **Creating a Digital Twin of an Insider Threat Detection Enterprise using Model Based Systems Engineering**

JAMES LEE, GEORGE MASON UNIVERSITY AHMAD ALGHAMDI, GEORGE MASON UNIVERSITY DR. ABBAS ZAIDI , GEORGE MASON UNIVERSITY

![](_page_0_Picture_3.jpeg)

![](_page_1_Picture_0.jpeg)

### Agenda

Motivation

Benefits of Using a Digital Twin (DT)

Insider Threat Detection Evaluation Challenges

Inference Enterprise Modeling as a DT Technique

Solution Method – Model Based Digital Twin Generation

Discussion

![](_page_2_Picture_0.jpeg)

### Motivation

Constant threats -> need for continuous monitoring

#### Operational testing is a challenge

- Lack of ground truth
- Data Privacy
- Blackbox models (COTS)
- Sensitive to disruptions

Solution: Use digital twin techniques to create real-time virtual counterpart of the physical system

![](_page_3_Picture_0.jpeg)

![](_page_3_Picture_1.jpeg)

Provides virtual and dynamic digital representation of the system

Enables real-time monitoring of the system

Reduces cost associated with testing and verification

Integration with MBSE reduces ambiguity through unifying system structures and behaviors

![](_page_4_Picture_0.jpeg)

### Insider Threat

Malicious threat that comes from people within the organization

#### Examples:

- 2019 Capital One: \$150M
  - Data breach via former employee of vendor
- 2018 Google and Uber: \$245M
  - IP theft via former employee

#### 2016 DoD 5220.22-M:

• Insider Threat Program Mandate

![](_page_4_Picture_10.jpeg)

![](_page_5_Picture_0.jpeg)

Limited understanding of detection accuracy / performance

Confusion Matrix

![](_page_5_Figure_4.jpeg)

Overwhelming number of false alarms / limited ability to reduce the population to manageable subpopulations for more detailed analysis

![](_page_5_Picture_6.jpeg)

Inference Enterprise (IE): enterprises that make inferences based on incomplete information (lack of ground truth)

### Methodology Development

![](_page_6_Picture_1.jpeg)

19 Challenge Problems: description of Inference Enterprise (IE) and data

- Range of threat behavior: accidental compromise to deliberate sabotage
- Missing and/or incomplete data
- Down-select algorithms used: classification / clustering algorithms
- Forecast future performance, evaluate hypothetical upgrades to system

### Inference Enterprise Modeling (IEM) methodology

Build models and evaluate performance

![](_page_6_Picture_9.jpeg)

![](_page_7_Picture_0.jpeg)

## Inference Enterprise Modeling (IEM)

![](_page_7_Figure_2.jpeg)

### Relevant Methods - Ontology

Formal description of knowledge as a set of concepts within a domain and the relationships between those concepts

Benefits:

- Provides coherent and easy navigation between one concept to another - classes and instances
- By having essential relationships between concepts, enables automated reasoning
- Easy to connect with other ontologies
- Can represent any data format, enabling smoother data integration

![](_page_8_Figure_8.jpeg)

![](_page_8_Figure_9.jpeg)

![](_page_8_Picture_10.jpeg)

### Relevant Methods – Process Modeling

![](_page_9_Picture_1.jpeg)

Graphical representation of business processes or workflows

Business process: collection of tasks an organization performs to create products, reach goals, provide value

Process Modeling Languages: Business Process Modeling Notation (BPMN), Event Driven Process Chain (EPC), UML Activity Diagrams

![](_page_9_Figure_5.jpeg)

![](_page_10_Picture_0.jpeg)

### Solution Components

![](_page_10_Figure_2.jpeg)

![](_page_11_Picture_0.jpeg)

### IEM Process Ontology Overview

![](_page_11_Figure_2.jpeg)

# **ProblemRequirements**: characteristics that define problem

- Simulate future data based on data provided
- Use certain ML algorithm to predict
- Calculate precision, recall, and false pos. rate

**Solutions:** software assets and/or manual activities that fulfill requirements

- Population assumptions based on SME input
- Population simulation and/or ML algorithms

#### **ProcessTemplate**

- Step A: Simulate Populations
- Step B: Down-select
- Step C: Calculate estimates

### Small Example

![](_page_12_Picture_1.jpeg)

![](_page_12_Figure_2.jpeg)

#### **Input New Problem:**

- Correlate data sources
- Use decision tree to predict
- Calculate precision, recall, and false pos. rate

#### **Output ProcessTemplate:**

StepA	StepB	StepC						
Sol2A	Sol1B	Sol1C						

![](_page_13_Picture_0.jpeg)

1. Use formal process modeling to document the solution workflows

![](_page_13_Figure_3.jpeg)

Notation name	Concrete syntax	note
Initial Node	•	Starting point of an activity.
Activity Final Node	۲	Marked the end or completion of an activity.
Fork/Join Node	Ι	Fork node has a single incoming edge and many outgoing edges, and it is the reverse for the join node. These elements help modelers in expressing the occurrence of concurrent sequences in an activity.
Decision/Merge Node	$\diamond$	Used to show the starting/ending of an alternative sequence in an activity. Usually merge nodes are used in conjunction with decision nodes to model loops within an activity.
Wait Time Action	$\boxtimes$	Implies wait until certain action occurs.
Opaque Action	(1)	Introduced for implementation-specific actions. Allow modeler to implement action by running user defined script in languages such as Java, C++ and others.
Accept Event Action		Action that waits for the occurrence of an Event that meets specified conditions.
Action	0	An Action is a named element that is the fundamental unit of an executable functionality. Regular actions are used to represent manual tasks in this work.

![](_page_14_Picture_0.jpeg)

 $\times$ 

### IEM Process Ontology Development

1. Use formal process modeling to document the solution workflows

Specification of Opaque Action RCP17 Install Dependencies

#### Specification of Opaque Action properties

Specify properties of the selected Opaque Action in the properties specification table. Choose the Expert or All options from the Properties drop-down list to see more properties.

![](_page_14_Picture_6.jpeg)

 $\times$ 

Specification of Opague Action <>

🗏 te 🔽 🏾 🎗	RCP17 Install Dependencies											
RCP17 Install Dependencies		Properties: All	~									
Documentation/Comments	Opaque Action		^									
Navigation/Hyperlinks	Name	RCP17 Install Dependencies										
Usage in Diagrams	Qualified Name	Model::RCP17 Install Dependencies										
Pins	Owner	🔁 Model										
Inner Elements      Deletione	Applied Stereotype											
- Entracons - Entracons - Constraints - Traceability	Body and Language	<pre>String[] cmd = {"c:/Program Files/K/R-4.0.3/bin/Rscript.exe", "C:/RCP17/install_dependencies.R"}; Runtime.getRuntime().exec(cmd)</pre>										
- Allocations	Body	String[] cmd = {"c:/Program Files/R/R-4.0.3/bin/Rscript.exe", "C:/RCP17/install_dependencies.R"}; Runtime.getRuntime().exec(cmd)										
	Active Hyperlink											
	Redefinition Context											
	Redefined Element											
	Is Leaf	false										
	Applied Stereotype Instance											
	Owned Comment											
	Owned Element											
	Handler											
	In Partition											
	In Interruptible Region		~									
	Name The name of the NamedElement.	Name The name of the NamedElement.										
	Q Type nere to filter properties											
		Close Back Ecoward	Heln									

Specification of documentation and commen Write documentation for the selected Opaque A	ts ction and create new comments.	
= 1: 🖸 <i>2</i>	Documentation/Comments	
<sup>(27)</sup> RCP12 Estimate Correlation <sup>(17)</sup> Documentation/Comments <sup>(17)</sup> Navigation/Hyperlinks <sup>(17)</sup> Navigation/Hyperlinks <sup>(17)</sup> Navigation/Hyperlinks <sup>(17)</sup> Navigation/Hyperlinks <sup>(17)</sup> Pins <sup>(17)</sup> Pins <sup>(17)</sup> Pins <sup>(17)</sup> Relations <sup>(17)</sup> Constraints <sup>(17)</sup> Traceability <sup>(17)</sup> Allocations	HTML      B      C:\Users\James\Desktop\Dissertation\RCP12\Updated Instructions.docx      C:\Users\James\Desktop\Dissertation\RCP12\Correlation Study.xlsm	
	Delete	
	Comments	
	Close Back Forward Help	

![](_page_15_Picture_1.jpeg)

2. Generalizing solution methods by creating a process template and categorizing each part of the solution as a section of the template

![](_page_15_Figure_3.jpeg)

![](_page_16_Picture_0.jpeg)

#### 3. Defining traceability links from problem requirements to each solution method – Problem Side

	Available Data								Data Characteristic						Down-Select Algorithm Types											Addition	Performance Evaluation							
	Correlation IPE ObsDist Target Behavior			et Behavior	Correlation Amt Num Obs Num TB Observable Data Types				Algorithms without Training Set Algorithms with Training Set																									
	None	Obs	ObsTime	Yes	No	Yes	Yes	: No	Full	Partial	<= 16	>16	>2	Both	ContOnly	DiscOnly	Alerts>=2	Cate≻=2	Days>=8	PCADBS	NB	RF	DT	LR	GBM	SVM	HMM	NN	HPTune	MultTrain	MultTest	CMR	SP	SPOQ
Q1(RCP 1, 2, 3)		8			х	8	8		8		ж					х	×																	х
RCP 4, 6		ж		×		8		×		×	×					×	×																	х
RCP 5		×		×		×		×	8		×				8					×														×
RCP 7		8			8	8		8	8		×				8					х														х
RCP 8		×			×	×		×	×			×				×			х															х
RCP 9		ж			8	8		8	8		х					8		х																х
RCP 10		×			х	×	×		8			×		×									х								х			х
RCP 11		х			х	8	×		8		×			×									×								х			х
RCP 12	х				8	8	8				8		х			8							×								х	х		
Q5 (RCP 13, 14)			х		х	8	×		×			×		×														8	×					х
RCP 15			8		8	х	8			8		8		8									×			×	×		×	х	×		х	
RCP 16		х			х	8	×		×			×		×							х	×									х		х	
Q7 (RCP 17, 18, 19)			х		х	8	×		8			×		х																	8		8	
IRCP 1		х			х	×		×	8			×			8									х						х	×			х
IRCP 2		×			×	×		×	8			×		×										х						×	×			×
IRCP 3		×			×	×	8		8			×		8										×						х	×			×

IRCP1	Property assertions: IRCP3
IRCP2	
IRCP3	Object property assertions
LessThanEqualTo16	hasTargetBehavior TargetLabelInfoProvided
LogisticRegression	hasDSCharacteristics MultTrain
MoreThan16	hasDownSelectAlgorithm LogisticRegression
MultTest	hasCorrAmount FullCorrelations
MultTrain	hasDSCharacteristics MultTest
NaiveBayes	hasObsDistribution ObservableDistributions
NeuralNetworks	
NewDownSelectNoTrain	hasNumObs MoreThan16
NewDownSelectTrain	hasCorrelation ObservableCorrelations
NoCorrelations	hasPerfEvalQuestions SysPerfObsQual
NoTBInfo	hasObsDataVarType BothTypesObs
<ul> <li>ObservableAndTemporalCorrelations</li> <li>ObservableCorrelations</li> </ul>	hasIndepPerfEst IPENo

![](_page_17_Picture_1.jpeg)

3. Defining traceability links from problem requirements to each solution method

![](_page_17_Figure_3.jpeg)

![](_page_18_Picture_0.jpeg)

v... with a construction of the constructio

hasReq

implementedBy

#### 3. Defining traceability links from problem requirements to each solution method – Solution Side

														4
MIEM Problem Specification	Discrete Event Activity Counts Complete	Stochastic Optimization Complete	Discrete Event Activity Counts with IPE	Stochastic Optimization with IPE	Inverse CDF Copula Discrete with SME Correlation	Factorized Stochastic Optimization with DEAC	Inverse CDF Copula	Inverse CDF Copula and Logistic Regression with IPE	Inverse CDF Copula Continuous and Discrete	Rank Correlat ion Copula	2D Rank Correlat ion Copula	2D Rank Correlation Copula with SME Partial Correlation	Hierarch ical Bayes Net	Tree Augmented Naïve Bayes
no correlations	[				x									
observable correlations	x	x	x	x			x	x	x	x				
observable and temporal correlations											x	x		
IPE			x	x				x						
no IPE	x	x			x	x	x		x	x	x	x	х	x
observable distributions	x	х	x	x	x	x	x	x	x	x	x	x	x	x
target labels no target label info														x
full correlations	x	х					x	x	x	x	х			
partial correlations			x	x								x		
num obs restriction <= 16	x	x	x	x	x									
any observables continuous observables only discrete observables only	v	v	v	v	v	v	x	x	x	x	x	x	x	x
discrete observables only	1^	^	^	^	^	^								

![](_page_18_Figure_4.jpeg)

### Small Example

![](_page_19_Picture_1.jpeg)

![](_page_19_Figure_2.jpeg)

![](_page_20_Picture_0.jpeg)

### Simulation Tool

![](_page_20_Picture_2.jpeg)

![](_page_21_Picture_0.jpeg)

### New Challenge Problems

Both problems use the GMU phishing experiment data set

- A: There are two models that use different sets of data to predict susceptibility to phishing.
   Since the two datasets are disjoint, correlation data between the sets do not exist. Forecast the performance of a model that uses all the data from both models.
  - Has continuous and discrete variables
  - Has binned data
  - Does not have full correlations
  - Has random forest
- B: There is a model that predicts staff members' susceptibility to phishing using staff data.
   Faculty data is not yet available. Forecast the performance of a model trained on staff data and tested on simulated faculty data.
  - Has full correlations
  - Has continuous variables
  - Has binned data
  - Has logistic regression

### Suggested Workflows

![](_page_22_Picture_1.jpeg)

![](_page_22_Figure_2.jpeg)

1.Understand the overall flow of the suggested solution

2. Identify the software and documentation location by inspecting the activity elements of the activity diagram

3.Read the documentation, run the software modules as-is, inspect the input and output files, and understand its behavior.

4.Create a copy of the software modules that point to the new dataset of the challenge problems

5.Compare the input data file formats of the new problem to the input and output of the suggested module and identify places were appropriate changes are necessary
6.Make the appropriate code changes to the modules and execute

![](_page_23_Picture_0.jpeg)

![](_page_23_Picture_1.jpeg)

Introduced IEM as digital twin technique for insider threat detection enterprises

Developed knowledge base of IEM expertise and model-based solution that can be used to rapidly prototype digital twins of new IEM scenarios

![](_page_24_Picture_0.jpeg)

### Future Works

Accurate problem requirement generation

Flexible template generation

Full simulation capabilities of systems modeling software

IEM – Ransomware protection

![](_page_25_Picture_0.jpeg)

### Thank you

Submitted full manuscript at SysCon 2022

Email: jlee194@gmu.edu