Al/Autonomy Solution Architecting Workflow Modeled in Digital Engineering

Ground System Architectures Workshop (GSAW)

October 2021

Approved for public release. OTR 2022-00081

© 2021 The Aerospace Corporation

Overview

Digital engineering supports integrating trusted Al/autonomy from inception through demonstration



Al/Autonomy Solution Architecting is part of the larger architecting and systems engineering effort that integrates the necessary Al/autonomy from inception to reduce cost, schedule, and technical risk

Architecting, system engineering, and modeling and simulation (M&S) iterate and inform each other in a spiral process

A solution team includes all relevant specialties serving iteratively throughout the process

Algorithm complexity drives/can be driven by software and hardware requirements and thus architectures within trade space and constraints

Conducting co-design within a digital engineering platform fosters rapid iteration and integration to develop and test prototypes

Demonstrating autonomous systems progressively reduces risk

Once an algorithm is selected, evidence must be collected to demonstrate trust

Digital engineering supports this entire process, from design to demonstrating trust

9

Integrating Al/autonomy from Inception Requires Building in Trust

Digital engineering supports building in trust, which requires careful design, development, & testing

- Trusted Al/autonomy is a nascent area
 - Researchers are working to develop best practices & standards, but no set of practices & standards has risen to widespread adoption yet
- Digital engineering tools, such as Cameo or GENESYS, support integrated co-design of trusted Al/autonomy
- The Aerospace Corporation has developed a trusted AI framework that maps to & includes other researchers' trust frameworks
 - Helps illustrate which information to collect for trust
- Digital engineering can instantiate the Aerospace trust framework—or any other trust framework—to expand its tenets into test routines & metrics
- This presentation introduces the Aerospace Al/autonomy solution process and its digital engineering model
 - The design use case is a convolutional neural network (CNN) used for pose estimation in on-orbit refueling within rendezvous proximity operations (RPO)
 - Digital engineering supports concrete data, tests, & metrics

Budgeting design, development, & testing for trust up front saves time & money



The Aerospace Al/autonomy Solution Architecture Workflow

Digital engineering can model this process, which allows complexity requirements to drive design



The Aerospace Al/autonomy Solution Architecture Workflow in GENESYS

The model instantiates the inner workings of the illustrated workflow concept



1: Refine CONOPS

The complexity required for the CONOPS and use cases drives the hardware/software design



1.2: Outline High-level States, Behaviors



1.3: Create Activity Diagrams of Workflow, Roles, Responsibilities

Activity diagrams help illustrate the what, who, when, and why of Al/autonomy decisions & action



1.4: Document Detailed CONOPS, Including Assumptions & Unhappy Path



2: Assess Considerations

Which factors constrain the solution? What does the trade space look like?



2.1: Assess Update Considerations

These considerations influence design re which Al/autonomy, why, & where within the solution



3: Identify Type(s) of AI, Autonomy, ML with TRLs



3.1: Identify Update Type(s) of AI, Autonomy, ML with TRLs



4: Perform Trusted AI ("MA for AI") Process

This stage of the solution architecting process links to the Aerospace Trusted AI Framework



Stage AIF.0, the Aerospace Trusted AI Framework, is captured in a separate GENESYS model/project



The Aerospace Trusted AI Framework, Modeled in GENESYS

15

5: Develop and Size Software and Hardware Architecture



5.1: Update Development & Sizing of Software & Hardware Architectures

Trade studies inform this stage of the co-design process,

in which Al/autonomy software is part of the larger software architecture

and Al/autonomy hardware is part of the larger hardware architecture

The overall hardware architecture must support the entire software-hardware co-design



6: Develop & Test Prototypes using Digital Engineering



6.1: Develop Update & Test Prototypes using Digital Engineering





6.1.1.2: Co-design Non-AI Software and Hardware, Including Sensors for maneuvering in RPO



7: Demonstrate Prototype to Customer; Iterate as Needed



Simulation Based on the GENESYS Model

10 iterations

Trust.1 Refine CONOPS		B	۵	8	0 8		0	1	Q	٥	0	
Trust.1.1 Identify players, roles, responsibilities	Simulation	Complete			×		1	- 1		1	1	
Trust.3.0.1 Develop Model		complete						I	I	I !		1
Trust.3.0.2 Select Model												
Trust.1.2.1 Outline/update high-level states, behaviors	🛛 💼 s	imulation run co	ompleted. All	events w	ere processed.		1	1	1	1	1	
Trust.1.2 Outline high-level states, behaviors							1	- 1	1	1	1	
Trust.1.3.1 Create/Update activity diagrams of workflow, roles, responsibilities							1	1	1	1	1	
Trust.1.3 Create activity diagrams of workflow, roles, responsibilities					ОК		1	1	1	1	1	
Trust.1.4 Document detailed CONOPS, including assumptions & unhappy path							1	8	1	1	1	
Trust.1.4.1 Document/Update detailed CONOPS, including assumptions & unha	ppy path	1	1	1	- I I		1	1	1	1	1	
Trust.2 Assess Considerations		8		B	0	3	0	0	8	0	E	1
Trust.2.1 Assess/Update Considerations		L 0	0 8 0	Ð	đ	3	0	0	Q	Ô	0	8
Trust.2.1.1 Begin Assessment			111	I I	I		1	1	1	1	- 1	1
Trust.1 Refine CONOPS			20								Ľ	1
Trust.2.1.2.1 Answer 5Vs of Data: Volume, Velocity, Variety, Veracity Value?			111	I I	1	I	1	- 1	1	1	1	1
Trust.2.1.2.2 Answer Questions Without Sub-questions			111	I I	I	I	1	1	1	1	1	1
Trust.2.1.3 Consolidate Assessments			111	1	I	I	- I	1	- I	1	1	1
Trust.3.1 Identify/Update Type(s) of AI, Autonomy, ML with TRLs		0	Ø			0] [1 2		0
Trust.3.1.1 Analyze Decision			1	I I	I	I	1	1	I	- I		1
Trust.3 Identify Type(s) of AI, Autonomy, ML with TRLs				\sim		0	E		1 🛛	1 0		2
Trust.3.1.2.1 Identify AI/Autonomous Appropriate Decisions				I I		I		1	I I			
Trust.3.1.2.2 Identify Decisions Needing Fully Manual			1		I		1			1		1
Trust.3.1.2.1.1 Determine Level of Automation				I I		I		1	I I			
Trust.3.1.3 Design or Select Algorithm for automation, planning or Al		1	1	Ш	II	I	II		I II	1		1
Trust.3.1.4 Assess Algorithm TRLs		1	- I	Ш	11	I I	I		I II	I		1
Trust.3.1.5 Assess Processor TRL			1	Ш	11	I	I			. I		1
Trust.3.1.5.1 Continue AI Process			1	I.	I	1			I	l		
Trust.3.1.6 Confirm Customer Timeframe				1	I				1	I		1

23

Summary

- Digital engineering supports the process of designing integrated Al/autonomy as well as demonstrating that such solutions are worthy of trust
 - Digital engineering can be used to model The Aerospace Corporation's Al/autonomy Solution Architecting workflow
 - This process uses the customer's complexity requirements to drive the hardware-software co-design
 - The trusted AI part of the process uses The Aerospace Corporation's Trusted AI Framework
 - This framework can be broken down into a Goal Structuring Notation, which can be instantiated in digital engineering for the design and demonstration of tests & metrics for trust
- The overall approach reduces cost, schedule, and technical risk by promoting a tailored Al/autonomy solution and helping customers and suppliers understand which information and evidence to collect to demonstrate that the resulting solution is trustworthy



Additional Information

- Introduction to The Aerospace Corporation's Trusted AI Framework
- How to begin collecting evidence for trusting a CNN to perform pose estimation

A Trusted AI Framework from Multiple Sources for Everyone

Explicit guidance on how to accomplish trust in relevant applications

The Aerospace Trusted AI Framework maps to and covers key industry frameworks as well as additional trust frameworks, including the IDA Roadmap to Assurance (May 2020) and the NIST Workshop on AI Trustworthiness (Aug 2020) (not shown in the table below)

Aerospace's Trusted Al Framework		National AI Initiative Office Characteristics of Trust	DoD's Principles of Al Ethics		AI Ethics Framework for the Intelligence Community		Deloitte's Trustworthy Al Framework	IBM - Trusting AI		Microsoft Responsible and Trusted Al	U.S. Air Force Research Laboratory
Thread 1	Objective Specification	Ethical use of AI			Testing, Version Control (builds, models, data), Stewardship	Governing the Al and the data; Documentation of		Transparency	Value Alignment		Ethics
	Data Specification	Privacy				Purpose, Parameters, Limitations, and Design Outcomes	Privacy	and Accountability		Privacy and Security	
Thread 2	Stability		Responsible, Traceable, Reliable				Robust/Reliable				
	Confidence and Uncertainty	 Accuracy, Reliability 								Te	Testing, Formal Methods
	Adversarial Robustness	Robustness, Resilience					Robust/Reliable		Robustness		
	Interpretability	Explainability and Interpretability, Transparency				Transparency: Explainability and Interpretability	Transparent/ Explainable		Explainability	Transparency	Transparency
	Familiarity			Equitable							
	Fairness	Fairness, Bias Mitigation				Mitigating Undesired Bias and Ensuring Objectivity	Fair/Impartial		Fairness	Fairness, Inclusiveness	
Thread 3	Monitoring	Security		Governable		Periodic Review	Safe/Secure	- Transparency and Accountability		Reliability and Safety	Reliability and Security
	Control	Safety				Human Judgement and Accountability	Responsible/ Accountable			Accountability	

This framework facilitates development of concrete requirements, test routines, and performance metrics that provide evidence for trust

Digital Engineering Supports Designing a Trusted Software Architecture

AI-ML models, such as CNNs, require specific trade analyses and trust

Follow the AI Solution Architecture Process	Perform the model trade analyses	Select or design the model	Use the Trusted AI Framework to build trust in the model in from inception			
CONOPS Considerations Identify Al/ Autonomy, TRLs	 Model architecture requires optimization: Model size, Data sources Languages, Libraries Accuracy 	Use AI Solution Architecture inputs, e.g., trades, to select/design a feasible, appropriately sized model:	 Assess current capabilities Identify risks & degree of autonomy Define Needs * Specify objectives * Specify data * Specify data * Specify data 			
Software, Hardware Architectures Develop, Test Prototypes Demos	 More complex model Harder to train More data, compute resources Simplicity Smaller model 	 CONOPS Considerations AI/Autonomy needs Model size Data needs and availability 	 Monitoring Control Thread 3: How can trust be maintained? Monitoring Control Stability Stability Stability Confidence & uncertainty Adversarial robustness Interpretability 			
Software and hardware architectures require iterative co-design to build in tests, monitoring, and control for trust	 Fewer parameters Less data, compute resources Might reduce overfitting Transfer learning vs learning from scratch 	 Data sources, quality, etc. (5Vs) Design of and 5V requirements for simulated data Ceiling analysis Hardware constraints 	• Fairness; Familiarity Trusted AI is a nascent field requiring explicit definitions into meaningful, generalizable, measurable, and testable attributes. High consequence environments often entail high risk in mission criticality, algorithm complexity to meet mission criticality and complexity, and level of autonomy to meet issues like communications latency. data volume, etc.; technical, cost, and schedule risks must be quantified so they can be mitigated			

The notional selection of a model, such as a CNN, for this use case is only the beginning of a solution

How do you know you can trust an artificial neural network like a CNN?

Trusting neural networks is not a solved problem—there is no established set of tests or metrics

- Convolutional neural networks (CNNs) are inspired by the mammalian brain and vision system
 - CNNs use multiple hidden layers & thus qualify as deep neural networks (DNNs)
 - DNNs are non-procedural and non-linear
- CNNs are well established and accepted for applications like computer vision, including pose estimation
- CNN architecture is feed forward, but testing for trust is not intuitively obvious
 - Trade analyses are required to validate the best model for an application
 - Metrics & tests must be developed
- The Aerospace trusted AI framework is a helpful starting point for developing these metrics & tests



Trusting the Model via the Trusted AI Framework in Digital Engineering

Goal Structuring Notations help define what's needed for **actionable confidence** that a model meets our objectives **quantifiably and understandably over the system lifetime**



Threads of Trusted AI: Feature Tests

The example in this presentation focuses on a subset of trust features listed under Thread #2

For this demonstration, the GSN expounds on three features, or subthreads

- Stability
- Confidence and Uncertainty
- Interpretability

for modeling and simulation in GENESYS, a digital engineering environment M2.1 Stability: The AI system's predictions are consistent when provided inputs fall within routine parameter ranges.

Trust has been explicitly defined for the AI

M2.2

M2/G1. Thread 2:

Confidence and Uncertainty: Levels of confidence and uncertainty bounds can be discerned for the AI system.

M2.3

<u>Adversarial Robustness</u>: The AI system is not only robust to a wide variety of known inputs but also robust to purposefully misleading inputs.

M2.4

Interpretability: The AI system is instrumented in a way for users to easily understand the underlying causes of how responses were formulated.

M2.5

Familiarity: Users are able to understand when to trust and when not to trust the Al system.

M2.6

Fairness: The AI system ensures that decisions made are not unfair or do not cause unintentional negative consequences due to bias. **Stability:** predictions are consistent when provided inputs fall within routine ranges

Confidence and Uncertainty: levels of confidence and uncertainty bounds can be discerned

Interpretability: users can easily understand the underlying causes of how responses were formulated

30

Integrating the Satellite Pose CNN into Digital Engineering

GSN elements of interest



This is the GENESYS instantiation of the GSN that was developed using a drawing tool like Visio

Digital Engineering: MBSE AI-ML Evaluation Flow

Elements of a GSN instantiated in GENESYS can access external code, like Python scripts



The Solution Team selected GENESYS for this use case, but other digital engineering environments should work

Trusted AI Framework Example: Thread 2 Pose Estimation Trust Features

Can a ResNet-50 CNN estimate the pose of a client spacecraft well enough for trust?

- Satellite Pose Estimation
 - Estimate the 3D position and orientation of a client spacecraft from 2D imagery
 - Deep Learning with a ResNet-50 Convolutional Neural Network (CNN) from the literature
- Spacecraft Pose Estimation Dataset (SPEED)
 - Created for Satellite Pose Estimation Challenge [Kisantal 2020]
 - Stanford University, Space Rendezvous Laboratory (SLAB)
 - Images of Tango spacecraft from PRISMA mission
 - First publicly available dataset for spacecraft pose estimation
 - Includes both synthetic and real images (mostly synthetic images)
 - 12,000 training images (labeled); 2,998 synthetic test images (not labeled)











Prototype Satellite Pose CNN

Overview

- Training a Regression Model in Python (train_dnn.py) (speed_resnet50.pth)
 - ResNet-50 architecture [He 2015]
 - Pose estimation is a regression problem (not classification)
 - Estimate orientation q and position r
 - Training objective (loss): minimize L_1 error between truth and predictions
 - Orientation L₁ loss: $\theta_q = |\hat{q} q|$
 - Position L₁ loss: $\theta_r = |\hat{r} r|$
 - Combined L₁ loss: $\theta = 15 \times \theta_q + \theta_r$
 - Implemented in Python with PyTorch machine learning library
 - Model pre-trained with ImageNet weights [Russakovsky 2015]



Stability

Predictions are consistent when provided inputs fall within routine ranges

- Evaluation (test_dnn.py)
 - Evaluate performance statistics against an independent and identically distributed dataset
 - Training dataset split: 80% training, 20% validation
 - Model trained with 9,600 images from training set
 - Model evaluated against 2,400 unseen images from training set
- Metrics
 - $-L_1$ loss for both orientation and position
 - Error metrics from Satellite Pose Estimation Challenge [Kisantal 2020]
 - Orientation error: $e_q = 2 \times \arccos(\langle \hat{q}, q \rangle)$
 - Normalized position error: $e_r = \frac{|\hat{r} r|_2}{|r|_2}$
 - Combined pose error: $e = e_q + e_r$

Stability

How stable is the prototype CNN?

- Prototype Results
 - Orientation L_1 loss (θ_q): 0.3180
 - Position L_1 loss (θ_r): 0.3977
 - Orientation error (e_q) : 1.1732
 - Normalized position error (e_r) : 0.0936
 - Combined pose error (e): 1.2667
- Comparison
 - State-of-the-art (2020) comparison appears below:

DETAILED RESULTS OF THE TOP TEN SUBMISSIONS COMPARED TO THE SLAB'S BASELINE PERFORMANCE									
Team	$E_{ m syn}$	$E_{\rm real}$	e_{q} [°]	e _t [m]	PnP				
1. UniAdelaide [47]	0.0094	0.3752	0.41 ± 1.50	$0.032~\pm~0.095$	Yes				
2. EPFL_cvlab	0.0215	0.1140	0.91 ± 1.29	$0.073~\pm~0.587$	Yes				
3. pedro_fairspace [48]	0.0571	0.1555	2.49 ± 3.02	0.145 ± 0.239	No				
SLAB Baseline [49]	0.0626	0.3951	$2.62~\pm~2.90$	0.209 ± 1.133	Yes				
4. Team_Platypus	0.0500	1 2001	0 4 4 1 4 0 4	0.001 1.0 700	No				
5. motokimura1					No				
6. Magpies	Performance falls within Top 3 for								
7. GabrielA					No				
8. stainsby	Satellite Pose Estimation Challenge								
9. VSI_Feeney				U	No				
10. jblumenkamp	0.0037	7.2710	00.04 L T0.14	2.000 ± 2.110	Yes				

Best results for each metric are highlighted with bold fonts. The mean and the standard deviation of the orientation errors (e_q) as in (3) and position errors (e_t) as in (1) are measured on the synthetic test set.

Confidence and Uncertainty

Levels of confidence and uncertainty bounds can be discerned

- Determining Out-of-distribution Inputs [Gawlikowski 2021]
 - Performance metrics are only known for training/validation data
 - Performance for data that deviate from training/validation is uncertain
 - Method needed for determining if given images are close enough to training images
- Generative Adversarial Network (GAN) (train_gan.py) (speed_dcgan.pth)
 - Deep Convolutional Generative Adversarial Network (DCGAN) [Radford 2016]
 - Separate DNN with two sub-networks: a Generator & Discriminator
 - Generator: creates similar (but fake) images
 - Discriminator: determines if an image is real or fake

Confidence and Uncertainty: the Role of Generators and Discriminators

The Discriminator helps identify Uncertainty by determining out-of-distribution inputs

- DCGAN Generator
 - Attempts to create images similar to training set from random noise
 - Note: these are low fidelity images



- DCGAN Discriminator
 - Can classify images that are real (matching distribution) or fake (out of distribution)
 - Evaluation results (test_gan.py)
 - Identifying *real training* images (79% accuracy)
 - Identifying *fake generated* images (66% accuracy)

Quantifying Confidence is challenging for regression problems like this: there's no general assessment that a neural network's prediction is "correct." It is necessary to define margins of error and instead measure confidence relative to those margins of error.

Interpretability

Users can easily understand the underlying causes of how responses were formulated

- Visualizing Focal Regions (viz_gradcam.py)
 - Gradient-based Localization (Grad-CAM) [Selvaraju 2019]
 - Highlights image regions that most significantly impacted predictions
 - Enables users to reject predictions when focal regions seem suspect



Tools like this illustrate to human users where within the imagery the CNN focuses to make its determinations

Demonstrating the Prototype Satellite Pose CNN in Python

Now the Solution Team can integrate the prototype CNN into digital engineering to illustrate the trusted AI framework and connect its tenets with requirements and other program artifacts

- Graphical Interface (run_demo.py)
 - Enables user to process images in real time with Satellite Pose CNN
 - Displays both "true" and "estimated" pose (computed in real time)
 - Enables processing with DCGAN to determine "accept" / "reject"
 - Enables Grad-CAM projection for *interpretability*
 - Enables image transformations to observe robustness to change



The dashboard with its slider bars facilitates Interpretability by showing the user how performance changes based on manual parameter adjustments