

# tinyML® Talks

*Enabling Ultra-low Power Machine Learning at the Edge*

## “A TinyML Approach to Deploy Reduced-Order Model of Complex Systems on Microprocessor”

Brenda Zhuang – Engineering Manager, MathWorks

Greg Copenrath – Sr. Product Marketing Manager, MathWorks

July 18, 2023



[www.tinyML.org](http://www.tinyML.org)



Thank you, **tinyML Strategic Partners**,  
for committing to take tinyML to the next Level, together



T I N Y



TALKS  
*webcast*

# Executive Strategic Partners

**Qualcomm**  
AI research

# Advancing AI research to make efficient AI ubiquitous

## Power efficiency

Model design, compression, quantization, algorithms, efficient hardware, software tool

## Personalization

Continuous learning, contextual, always-on, privacy-preserved, distributed learning

## Efficient learning

Robust learning through minimal data, unsupervised learning, on-device learning

## A platform to scale AI across the industry



### Perception

Object detection, speech recognition, contextual fusion



### Reasoning

Scene understanding, language understanding, behavior prediction



### Action

Reinforcement learning for decision making



Edge cloud



Cloud



IoT/IIoT



Automotive



Mobile



Accelerate Your Edge Compute

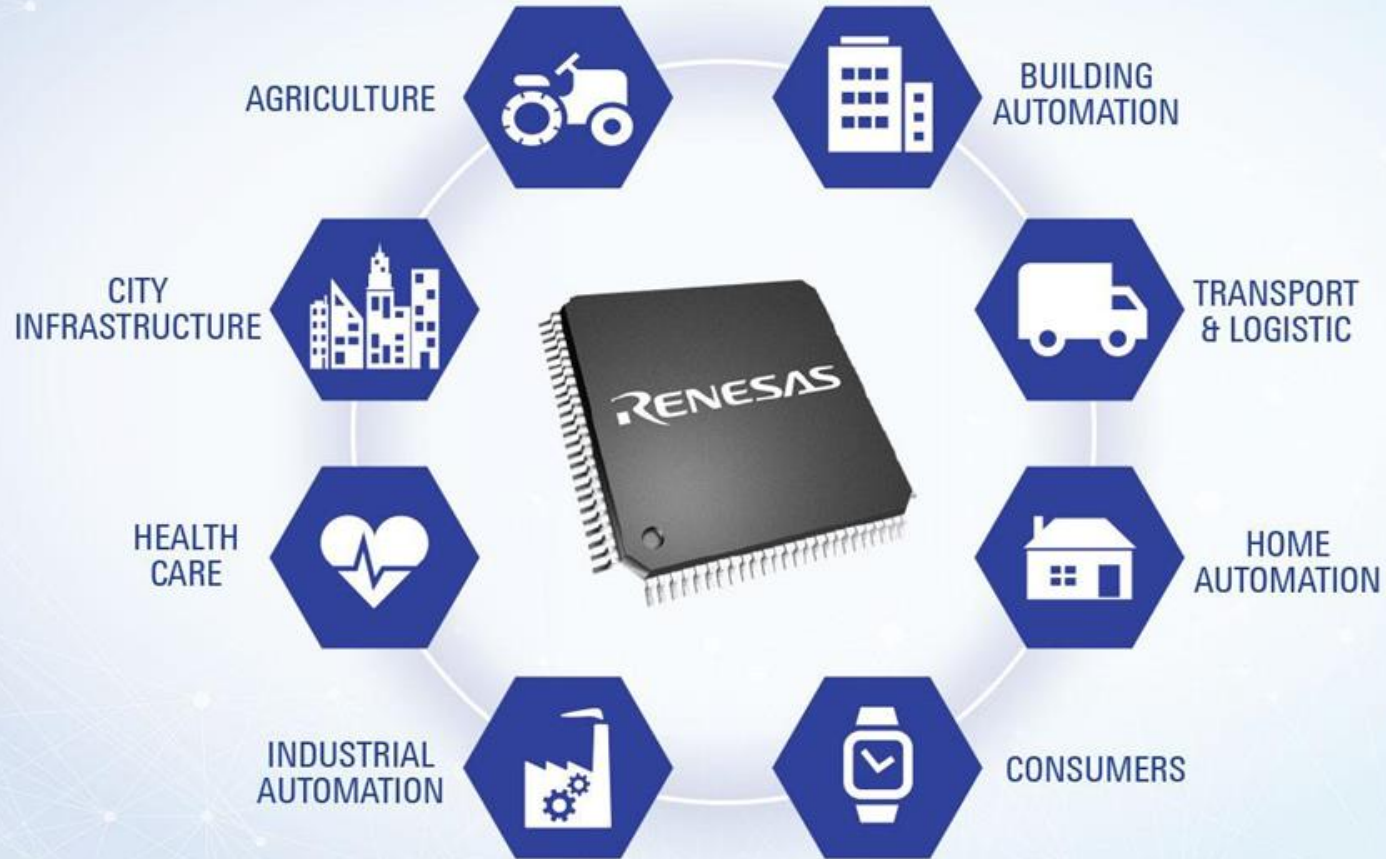
**SYNTIANT**

Making Edge AI A Reality

[www.syntiant.com](http://www.syntiant.com)

# Platinum Strategic Partners

**Renesas is enabling the next generation of AI-powered solutions that will revolutionize every industry sector.**



[renesas.com](https://www.renesas.com)



**DEPLOY VISION AI  
AT THE EDGE AT SCALE**

**SONY**



# Gold Strategic Partners



AHEAD OF WHAT'S POSSIBLE™



AHEAD OF WHAT'S POSSIBLE™

Where what if  
becomes what is.

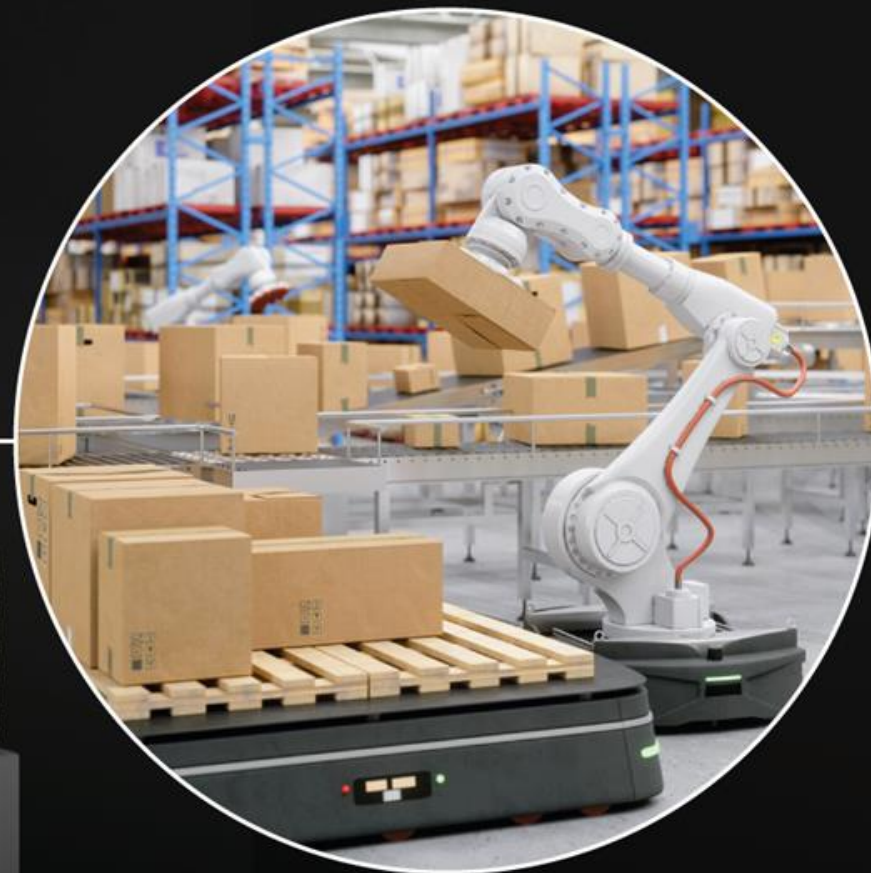
Witness potential made possible at [analog.com](http://analog.com).



**PRO**™

**Easily** deploy your  
**tinyML** solutions with  
**Arduino Pro**

[arduino.cc/pro](https://arduino.cc/pro)



**Made In Italy**

Build the  
Future of tinyML

on **arm**



T I N Y



TALKS  
*webcast*



**EDGE IMPULSE**

# The Leading Development Platform for Edge ML

[edgeimpulse.com](https://edgeimpulse.com)

Decarbonization

Digitalization



Driving decarbonization and digitalization. Together.

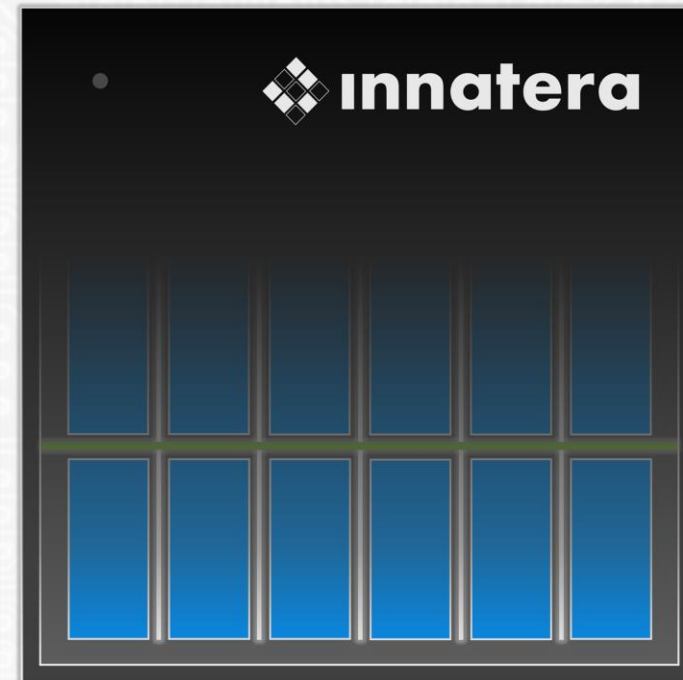
**Infineon serving all target markets as**  
**Leader in Power Systems and IoT**

[www.infineon.com](http://www.infineon.com)





# NEUROMORPHIC INTELLIGENCE FOR THE SENSOR-EDGE



[www.innatera.com](http://www.innatera.com)



Microsoft



The Right Edge AI Tools Can Make or Break Your Next Smart IoT Product



## Analytics Toolkit Suite





life.augmented

**STMicroelectronics provides extensive solutions to make tiny Machine Learning easy**



# ENGINEERING EXCEPTIONAL EXPERIENCES

We engineer exceptional experiences for consumers in the home, at work, in the car, or on the go.

[www.synaptics.com](http://www.synaptics.com)



T I N Y



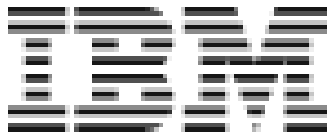
# Silver Strategic Partners



brainchip



Grovety Inc.



Nota AI





# Join Growing tinyML Communities:



15.8k members in  
49 Groups in 41 Countries

**tinyML - Enabling ultra-low Power ML at the Edge**

<https://www.meetup.com/tinyML-Enabling-ultra-low-Power-ML-at-the-Edge/>



3.8k members  
&  
12.7k followers

**The tinyML Community**

<https://www.linkedin.com/groups/13694488/>





Subscribe to  
**tinyML YouTube Channel**  
 for updates and notifications  
*(including this video)*

[www.youtube.com/tinyML](http://www.youtube.com/tinyML)



**tinyML**  
4.33K subscribers

**10k subscribers, 607 videos with 354k views**

HOME VIDEOS PLAYLISTS COMMUNITY CHANNELS ABOUT

13:24	33:27	32:39	36:41	34:03	34:58
On Device Learning Forum - Professors...	On Device Learning - Manuel Roveri: Is on-...	On Device Learning Forum - Warren Gros...	On Device Learning Forum - Yiran Chen...	On Device Learning Forum - Hiroku...	On Device Learning Forum - Song Han: O...
106 views · 4 days ago	138 views · 4 days ago	54 views · 4 days ago	47 views · 4 days ago	132 views · 4 days ago	137 views · 4 days ago
1:13	1:07:43	53:41	45:46	51:01	1:03:24
tinyML Smart Weather Station Challenge - ...	tinyML Talks Singapore...	tinyML Talks Shenzhen: Data...	tinyML Talks Singapore...	tinyML Smart Weather Station with Syntiant...	tinyML Trailblazers August with Vijay...
122 views · 4 days ago	262 views · 2 weeks ago	511 views · 3 weeks ago	229 views · 3 weeks ago	265 views · 3 weeks ago	286 views · 1 month ago
58:50	34:36	55:01	59:51	59:48	58:09
tinyML Auto ML Tutorial with SensiML	tinyML Auto ML Tutorial with Qeexo	tinyML Talks Germany: Neural network...	tinyML Trailblazers with Yoram Zylberberg	tinyML Auto ML Tutorial with Nota AI	tinyML Auto ML Tutorial with Neuton
351 views · 1 month ago	462 views · 2 months ago	374 views · 2 months ago	133 views · 2 months ago	287 views · 2 months ago	336 views · 2 months ago
1:02:30	34:31	1:00:30	1:06:44	1:53:07	42:13
tinyML Challenge 2022: Smart weather...	tinyML Talks South Africa - What is...	tinyML Talks: The new Neuromorphic Anal...	tinyML Talks Shenzhen: 分享主题...	tinyML Auto ML Forum - Paneldiscussion	tinyML Auto ML Forum - Demos
378 views · 2 months ago	214 views · 2 months ago	448 views · 2 months ago	159 views · 2 months ago	190 views · 2 months ago	545 views · 2 months ago

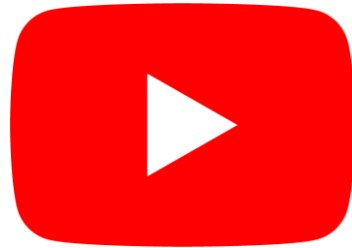


# Reminders

Slides & Videos will be posted tomorrow



[tinyml.org/forums](https://tinyml.org/forums)



[youtube.com/tinyml](https://youtube.com/tinyml)



Please use the Q&A window for your questions





## Brenda Zhuang



Dr. Brenda Zhuang is a consulting engineer and engineering manager at MathWorks, where she leads a team responsible for software tools for automatic deployment of embedded applications, such as motor controls and deep learning, to microprocessors and FPGAs. Brenda joined MathWorks in 2007. She received her PhD from Boston University in Systems Engineering. She serves on the technical program committee in control theory, modeling and simulation.





## Greg Coppentrath



Greg is the product marketing manager for Fixed-Point Designer and Deep Learning Toolbox Model Quantization Library. He has experience in the development of embedded systems and product development in the semiconductor industry. He received an MBA from Worcester Polytechnic Institute, an M.S. in Electrical Engineering from the University of Massachusetts Lowell, and received a B.S. in Electrical Engineering from WPI.



# A TinyML Approach to Deploy Reduced-Order Model on Microprocessor

# Meet the speakers today



**Brenda Zhuang, PhD**

Engineering Manager

MathWorks



**Greg Coppentrath**

Senior Product Manager

MathWorks

# MATLAB® & SIMULINK®



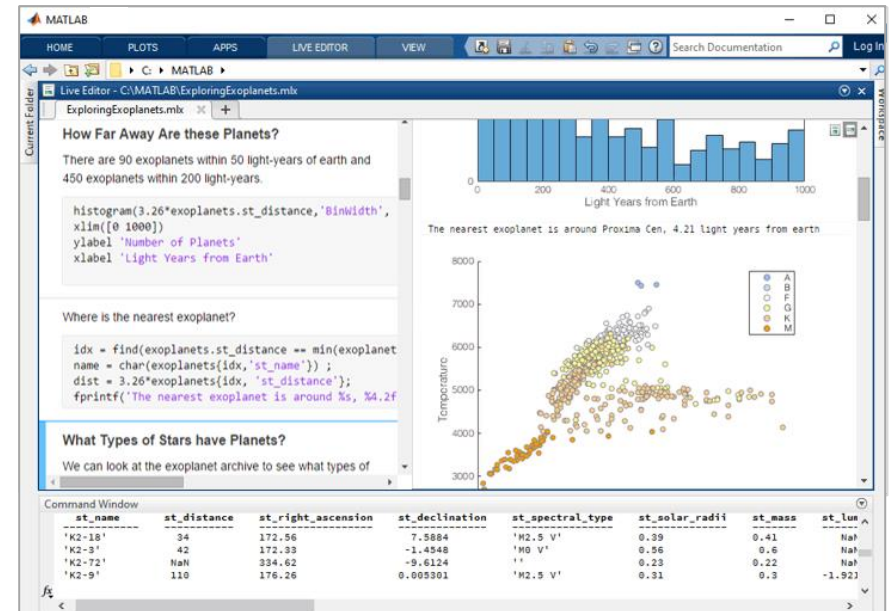
## Our Products

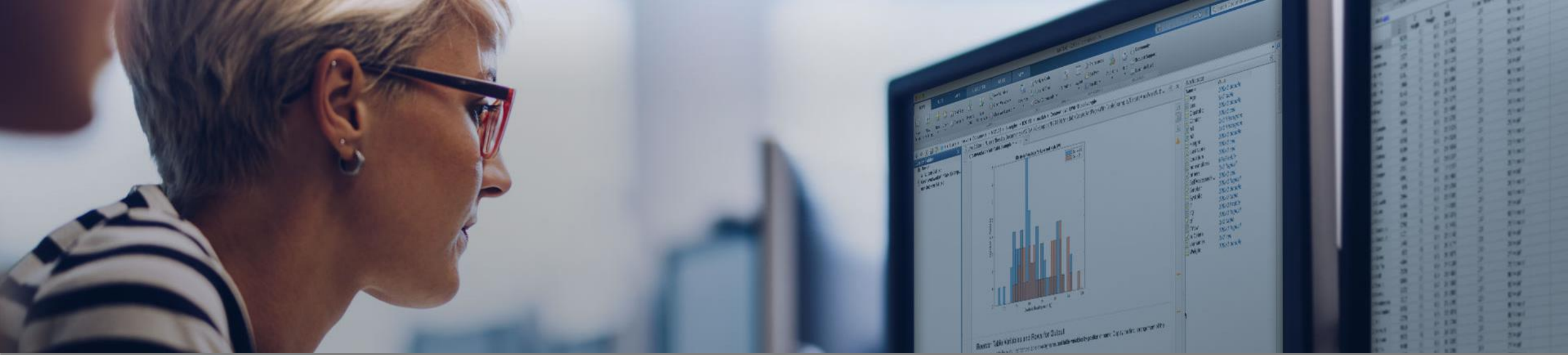
MATLAB is a programming environment for algorithm development, data analysis, visualization, and numeric computation.

Simulink is a block diagram environment for designing, simulating, and testing systems.

More than 120 add-on products for specialized tasks.

## Computer-Aided Design Toolbox





# Our Customers

Millions of engineers and scientists worldwide use MATLAB and Simulink.



**5 million+**

users in over 180 countries



**100,000+**

businesses, governments,  
and universities



All of the top 10  
automotive and  
aerospace companies

Fortune: 2021 Global 500 auto companies  
FlightGlobal: 2020 Top 100 aero companies\*

\*Excluding companies that are subject to embargos, sanctions, or other controls



**Headquarters**  
Natick, MA USA

**North America**  
United States

**Europe**

Finland  
France  
Germany  
Ireland  
Italy  
Netherlands  
Spain  
Sweden  
Switzerland  
UK

**Asia-Pacific**

Australia  
China  
India  
Japan  
Korea

## MathWorks Today



**6000+ staff**  
in 34 offices around  
the world



**\$1.25+ billion**  
in revenues



**Privately held**

# What's New for Our Products

## Megatrends changing our world

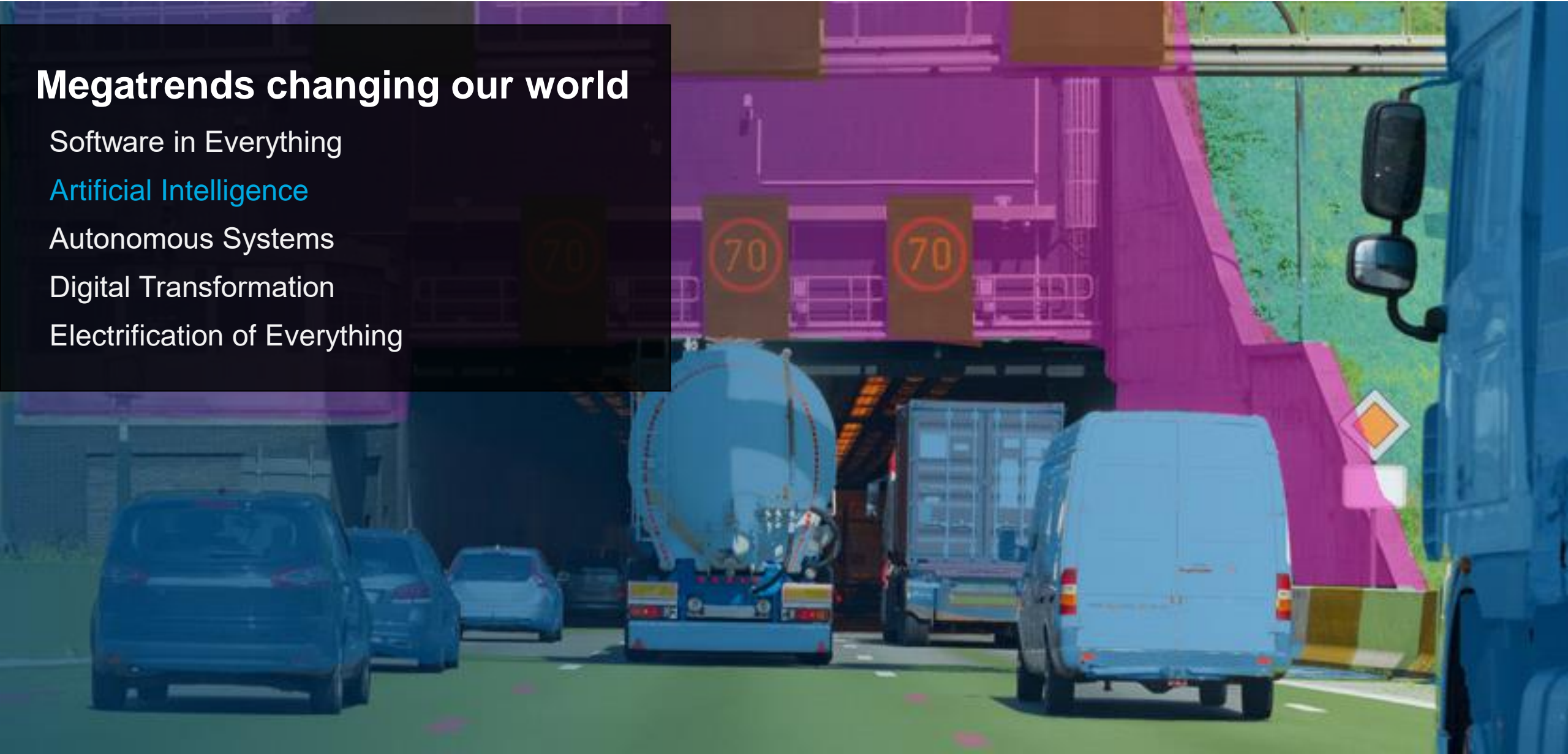
Software in Everything

Artificial Intelligence

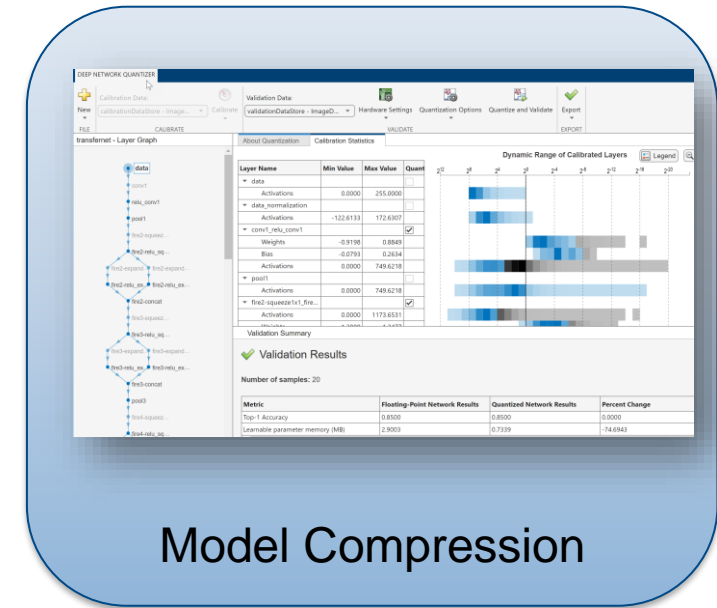
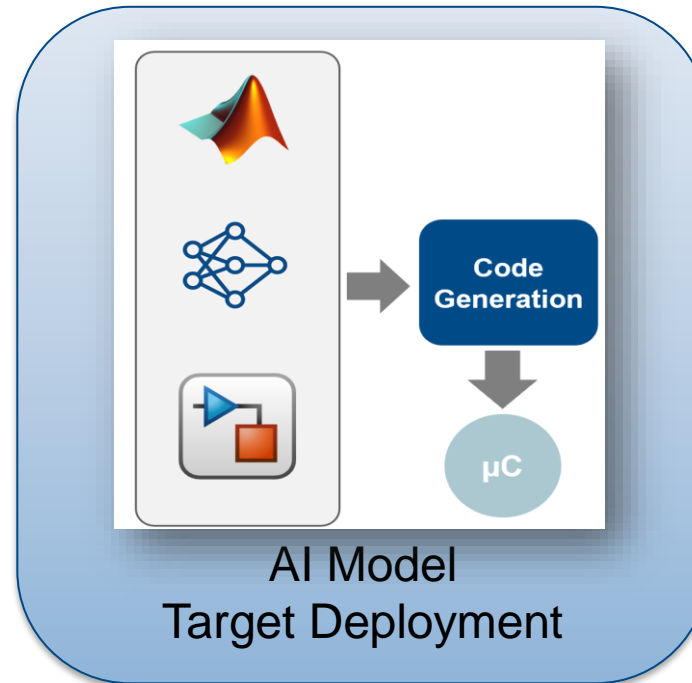
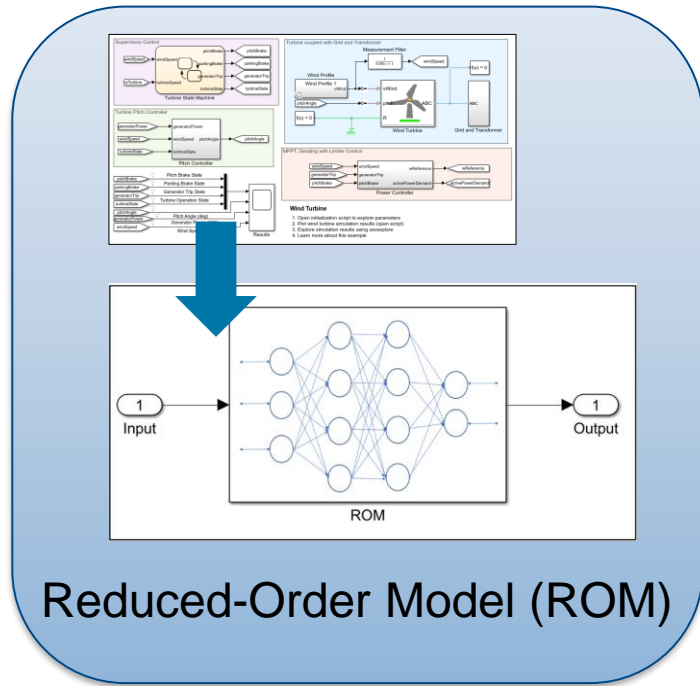
Autonomous Systems

Digital Transformation

Electrification of Everything



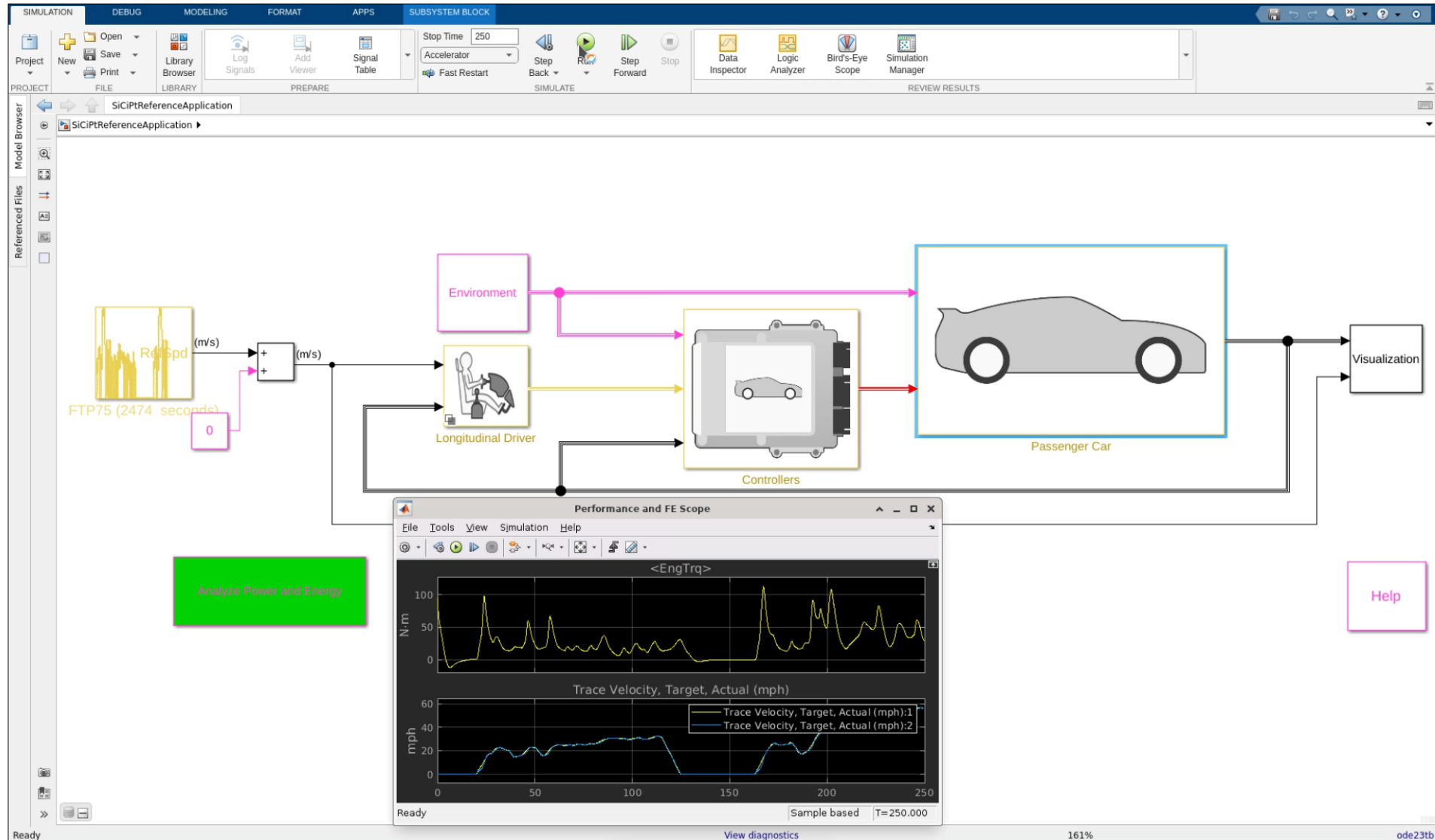
# Today, we cover the workflow steps from model development using AI-driven methodology to compression for target deployment



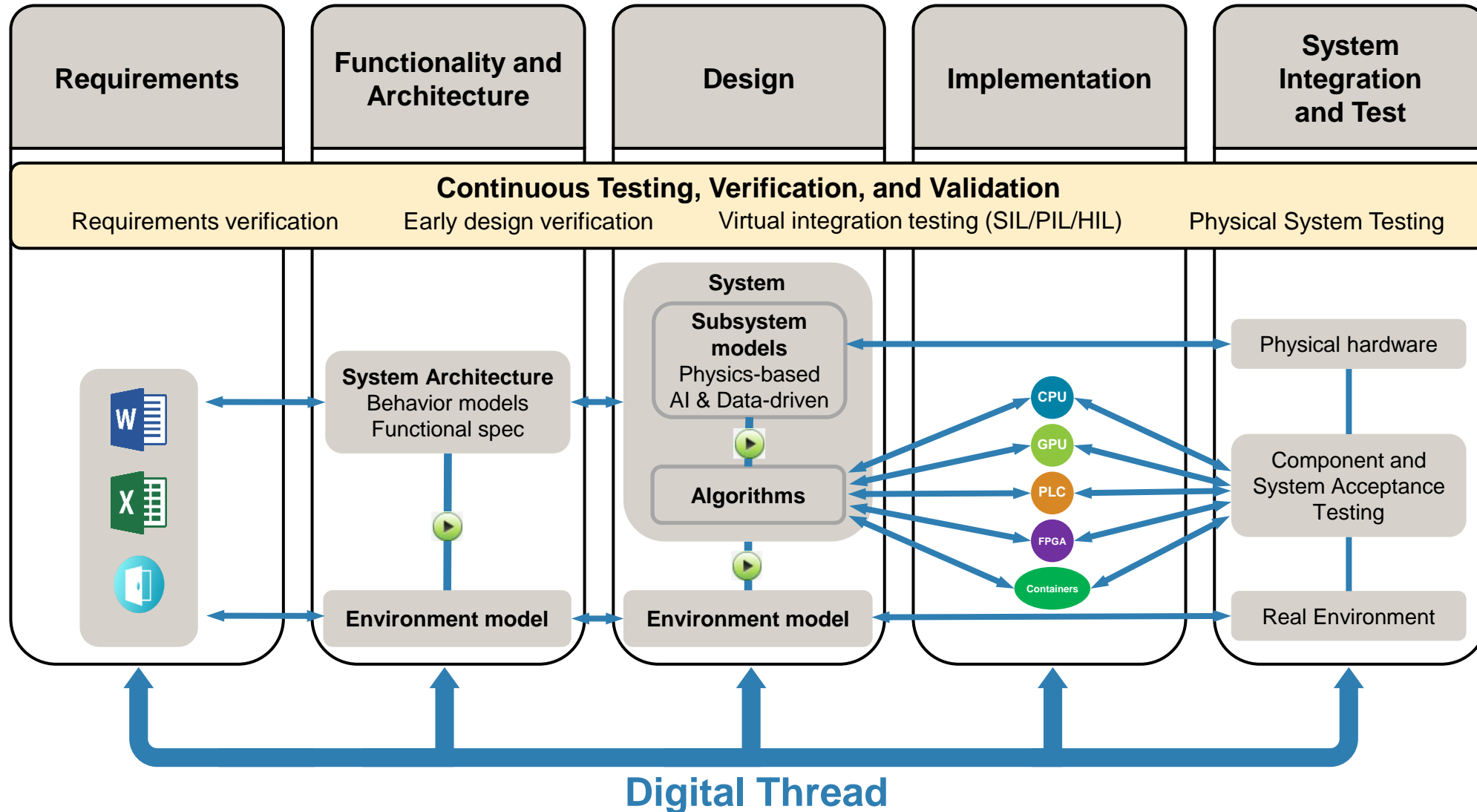


# Example overview

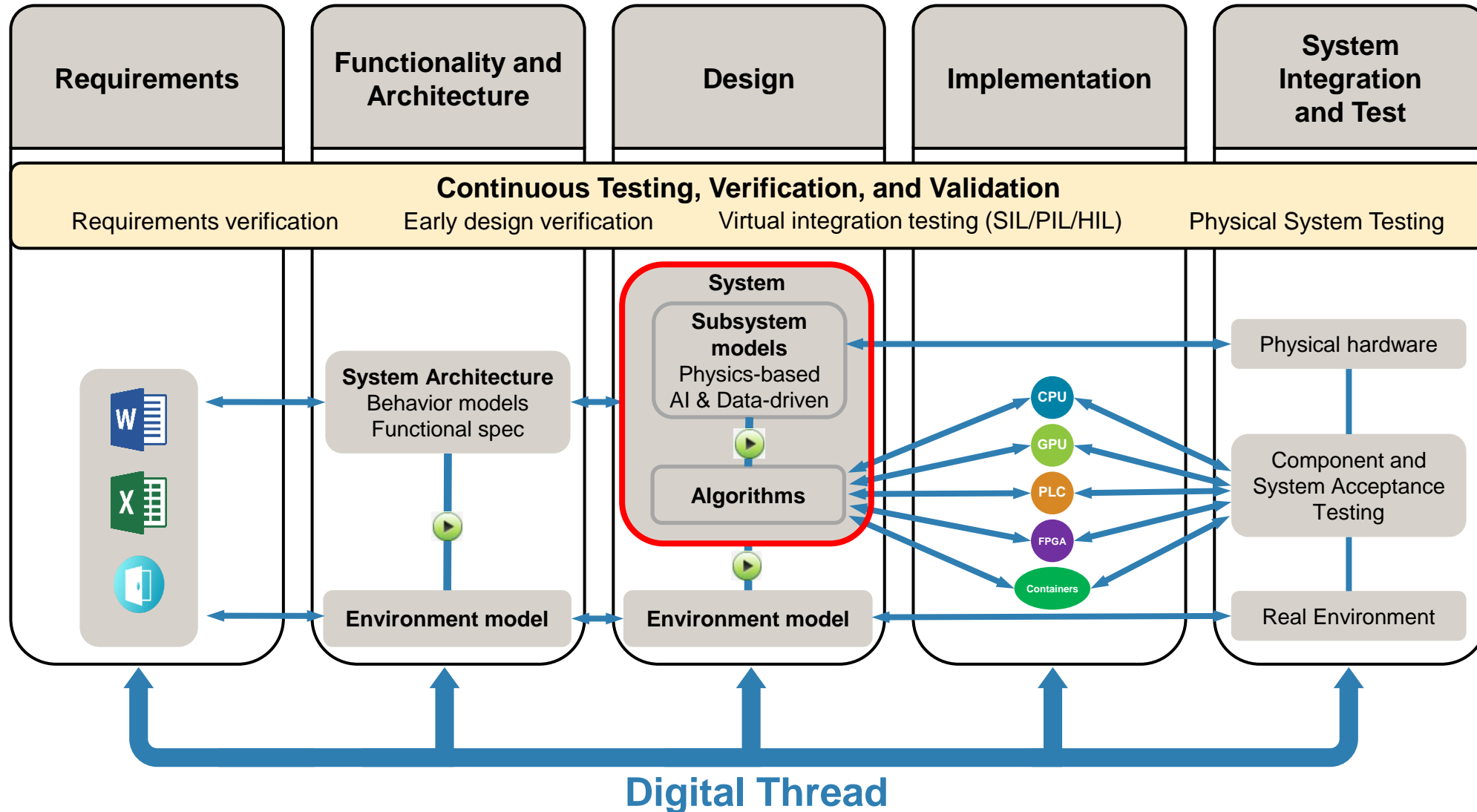
## System-level simulation



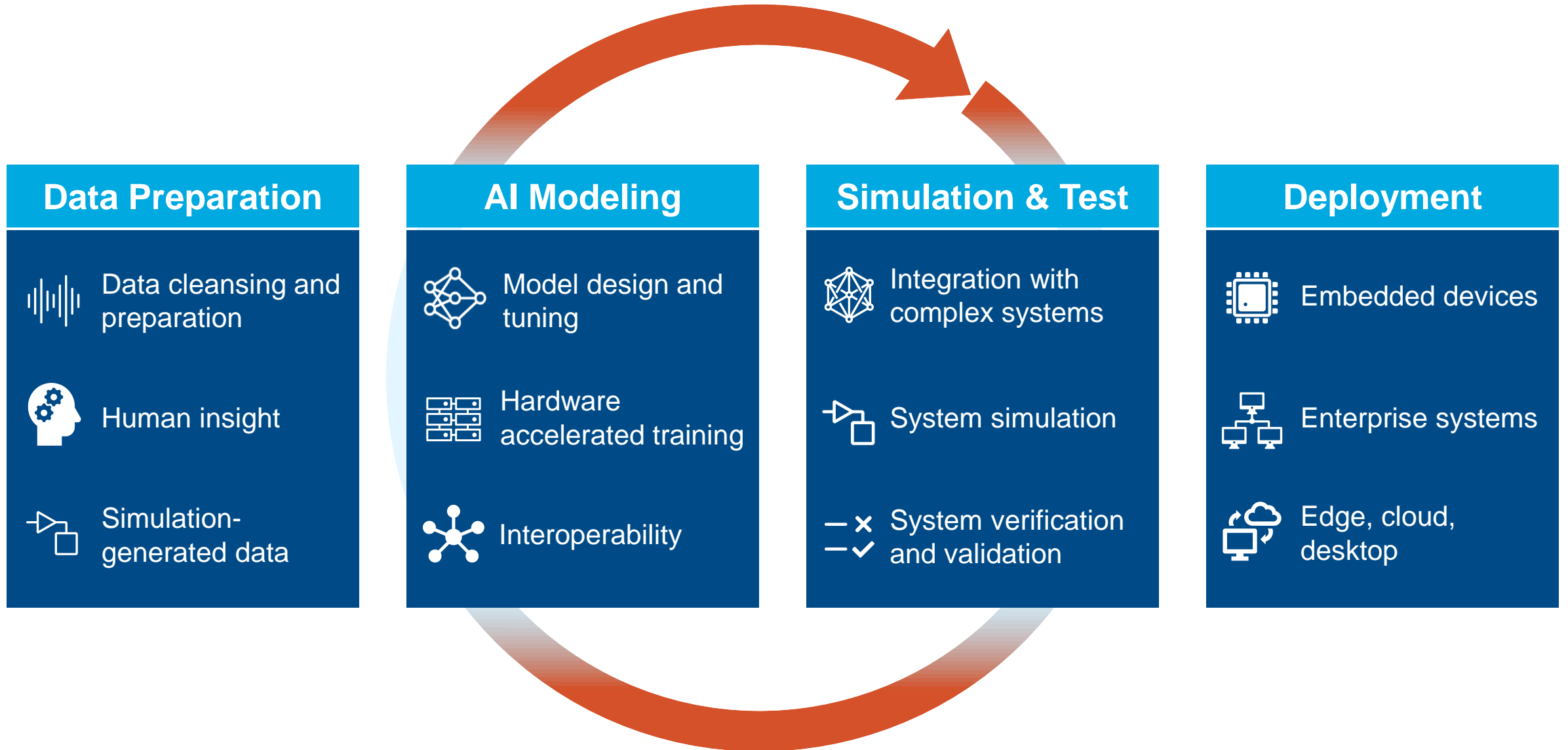
# Model-Based Design



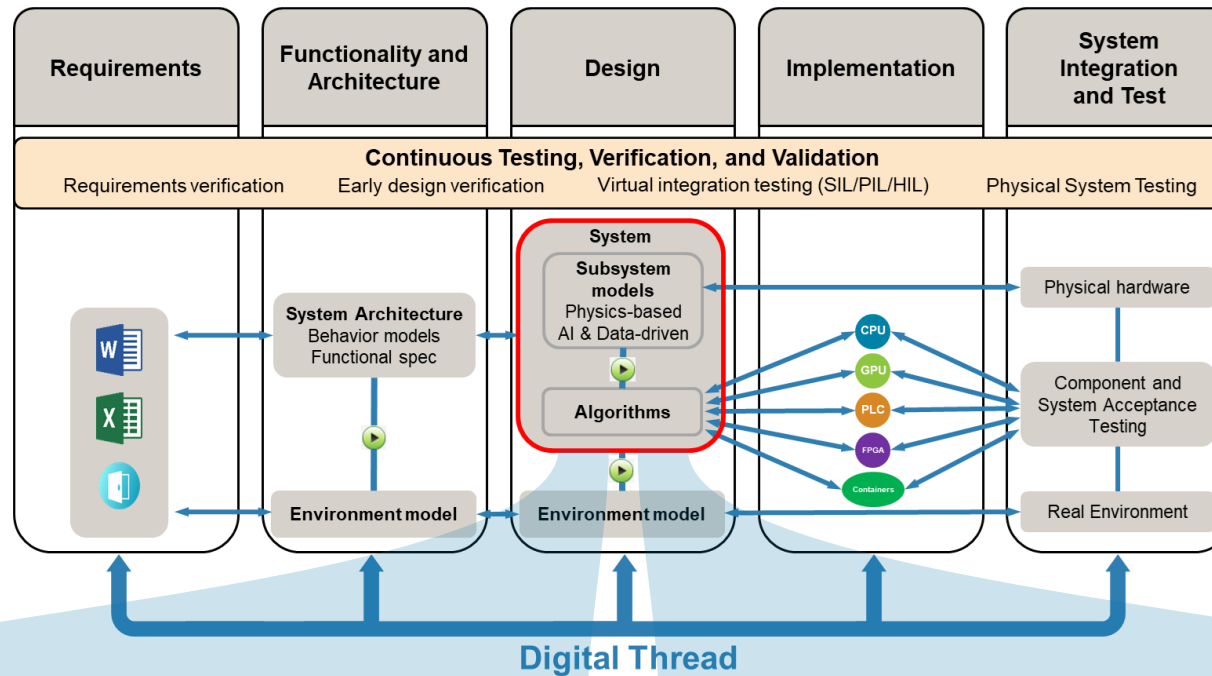
# Integrating AI into Model-Based Design



# AI-driven system design



# Integrate AI models into MBD for system-level simulation and code generation



## AI for component modeling

- Speeding up desktop and HIL simulations
- Modeling component dynamics from data when first-principles models cannot be obtained

## AI for algorithm development

- Virtual sensor modeling
- Sensor fusion
- Object detection

# Reduced Order Modeling

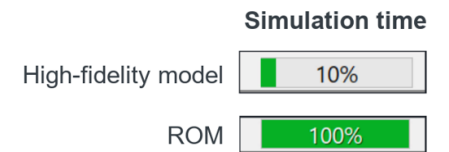
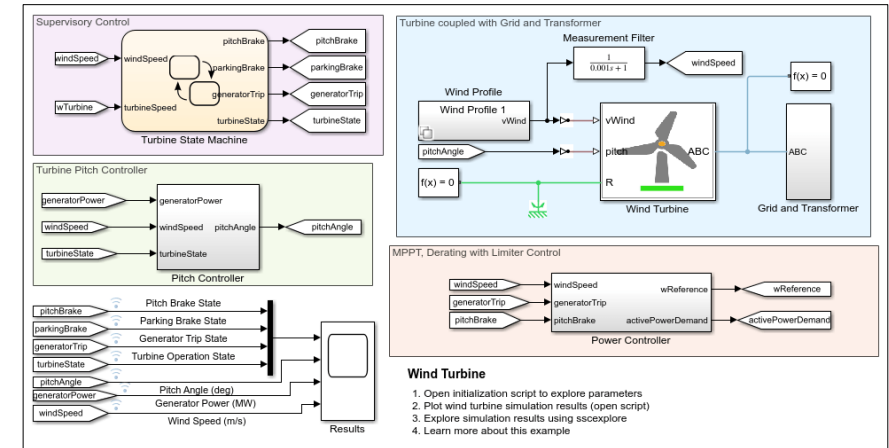
## What

- Techniques to **reduce the computational complexity or storage requirement** of a computer model
- Preserve the expected fidelity** within a controlled error

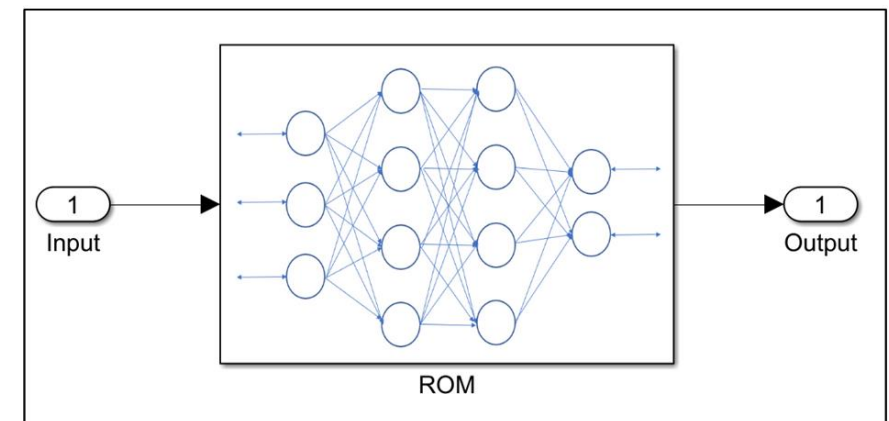
## Why

- Speed up system-level desktop simulation
- Hardware-in-the-loop testing
- Enable system-level simulation
- Develop virtual sensor, Digital twins
- Perform control design

## High-fidelity model

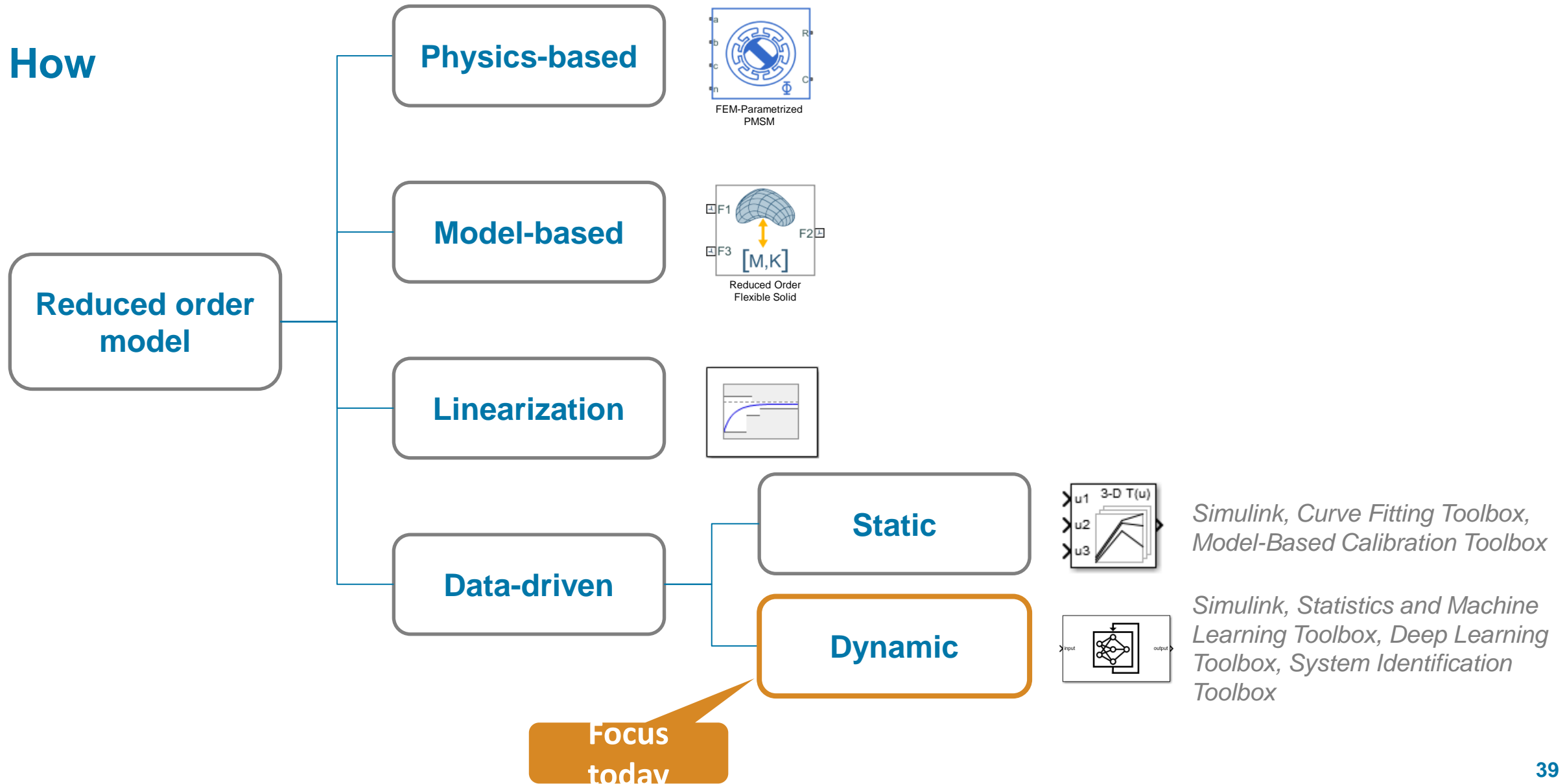


## Reduced-Order Model (ROM)



# Reduced Order Modeling

How



# Data-driven vs. first-principles modeling

Data-driven models and first-principles models can co-exist

## DATA-DRIVEN MODELS

Statistics, optimization, AI

## FIRST-PRINCIPLES MODELS

Physics, math, domain knowledge

**BLACK BOX**

**GREY BOX**

**WHITE BOX**

### Advantages

- May succeed when first-principles models are unavailable or challenging/impossible to find
- May reduce complexity, simulate faster
- Can leverage existing, measured data
- Do not require domain knowledge

### Challenges

- Require a lot of data
- Are often not
  - interpretable, explainable
  - easily parameterizable in a physically meaningful way
- Cannot extrapolate well beyond training data

### Advantages

- May capture (global) parameterizable behaviors with low/high fidelity
- Have clear (explainable) physical meaning
- Do not require data engineering

### Challenges

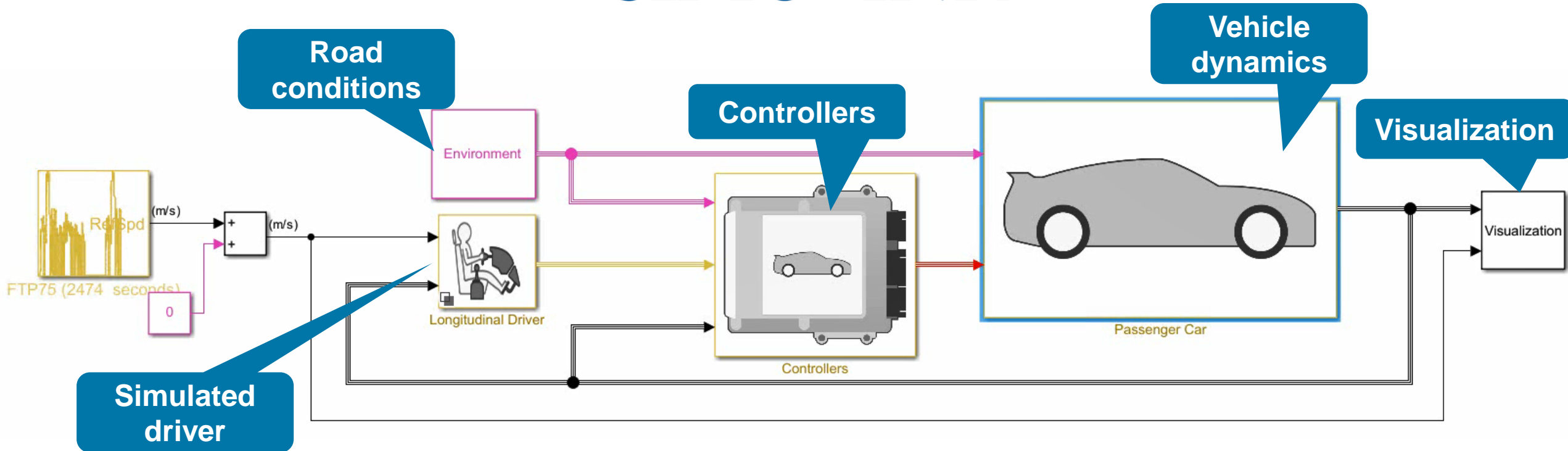
- Can be challenging/impossible to derive
- Require significant time for derivation
- Require expertise in the respective domain



# Example overview

*Replacing a first-principles engine model with an AI-based Reduced Order Model*

# SIMULINK®

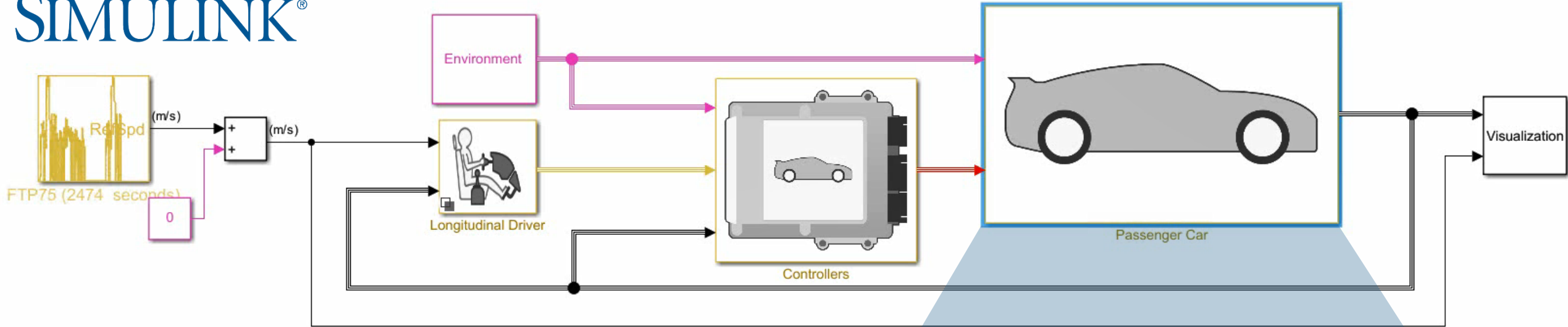


## Closed-loop control of vehicle speed

# Example overview

*Replacing a first-principles engine model with an AI-based Reduced Order Model*

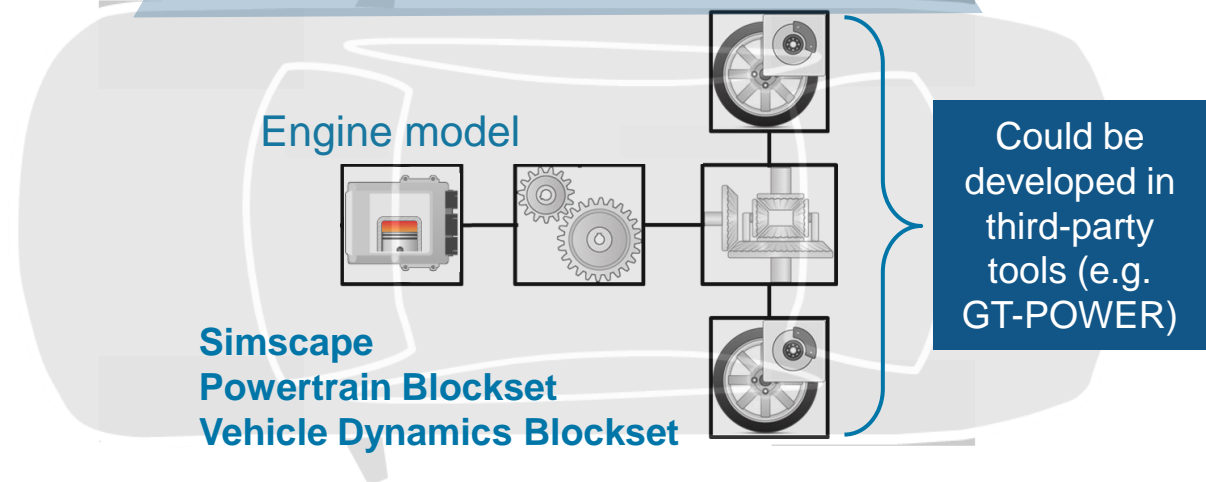
SIMULINK®



High fidelity

Complex model

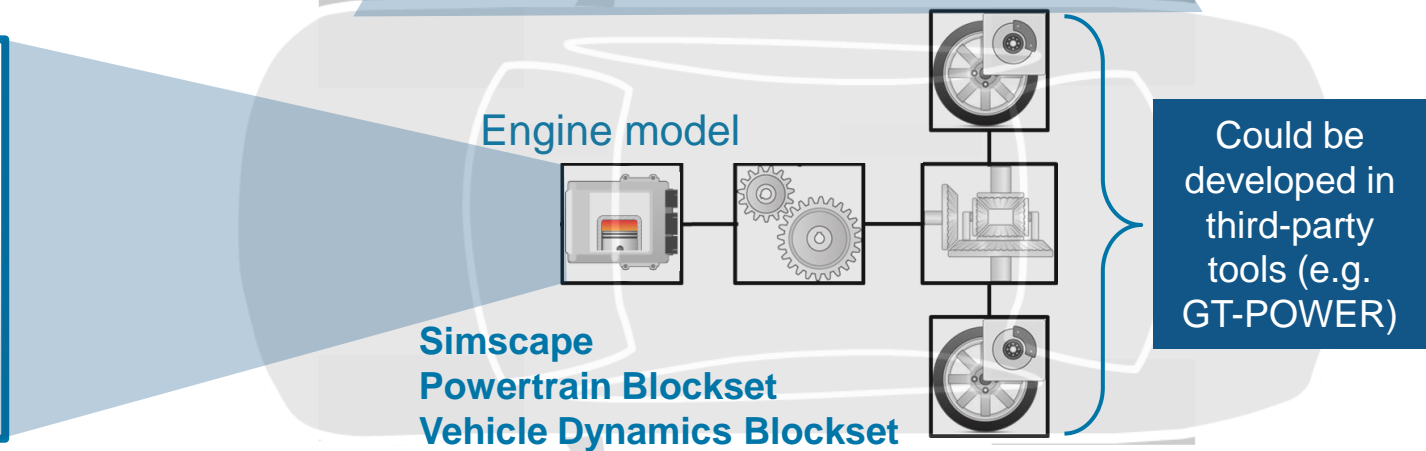
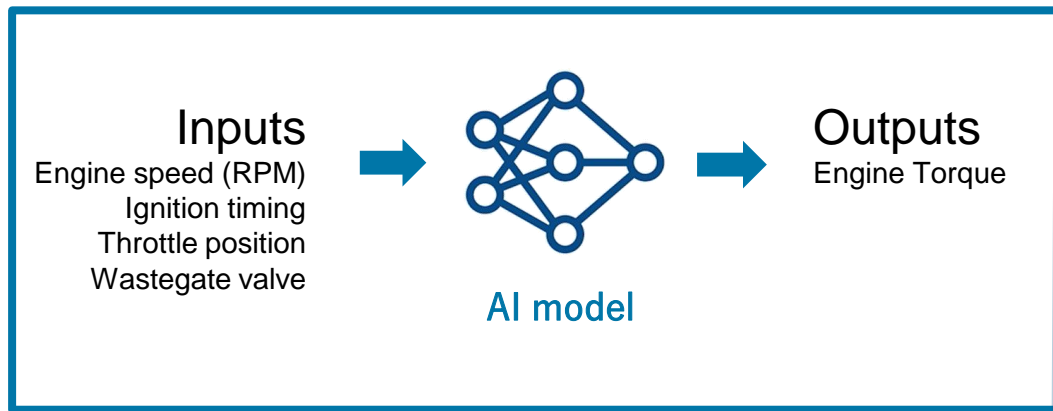
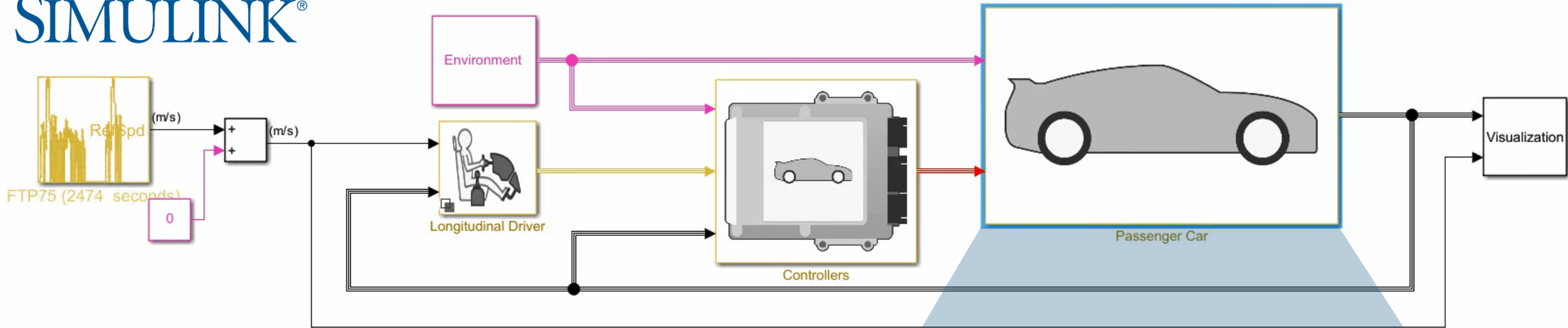
Slow simulation



# Example overview

*Replacing a first-principles engine model with an AI-based Reduced Order Model*

SIMULINK®

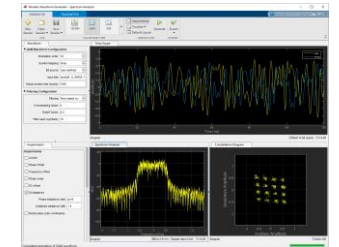


# Generate synthetic data for training

Other techniques:



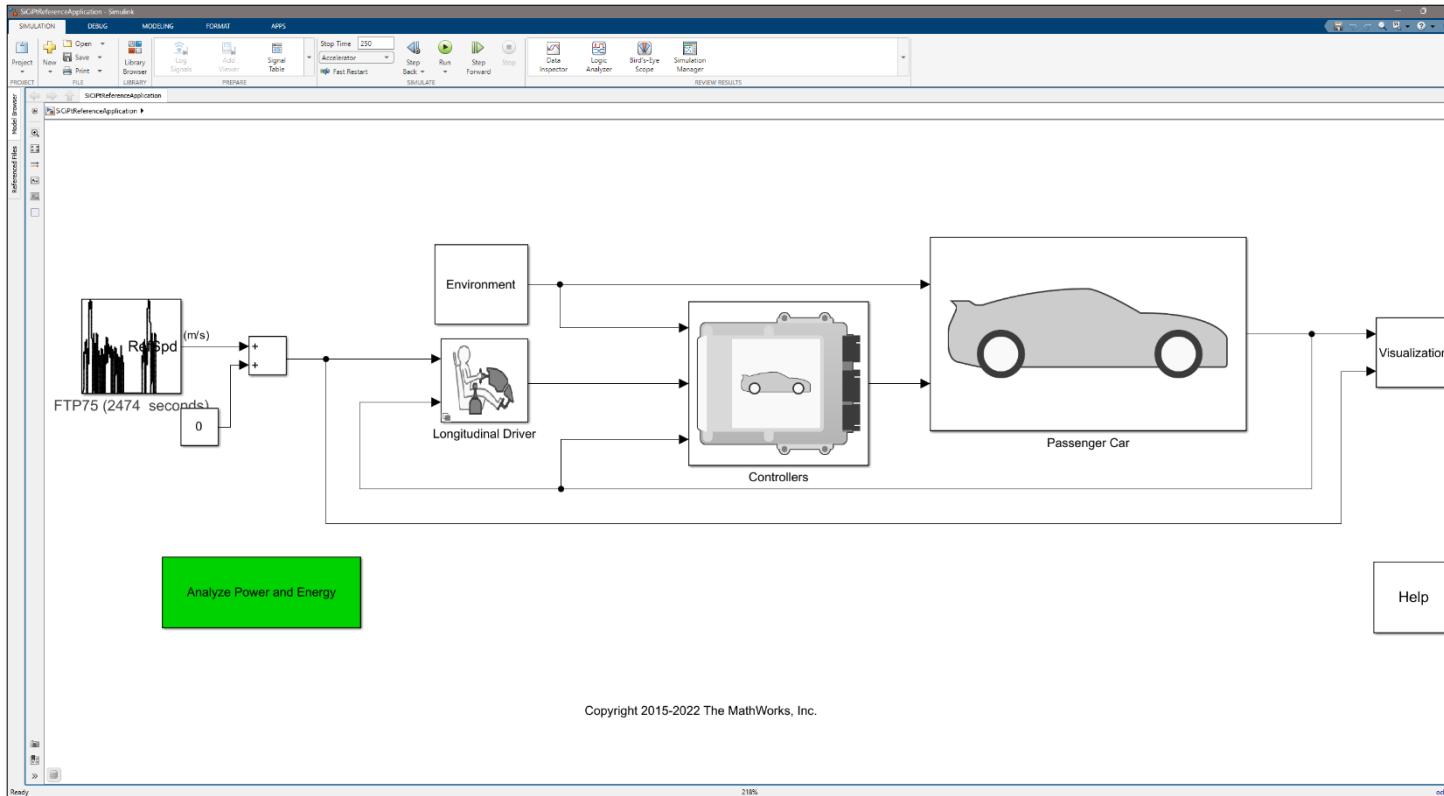
GANs



Wireless Waveform Generator



Unreal Engine®



Simulink/Simscape

Data Preparation

AI Modeling

Simulation & Test

Deployment

# Synthetic Data Generation

## Design of Experiments

DoE = 512x3 table

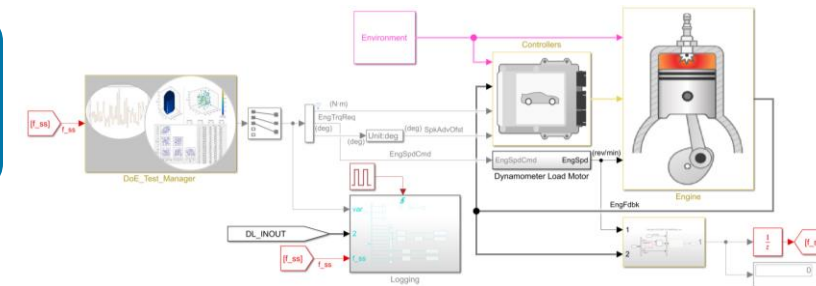
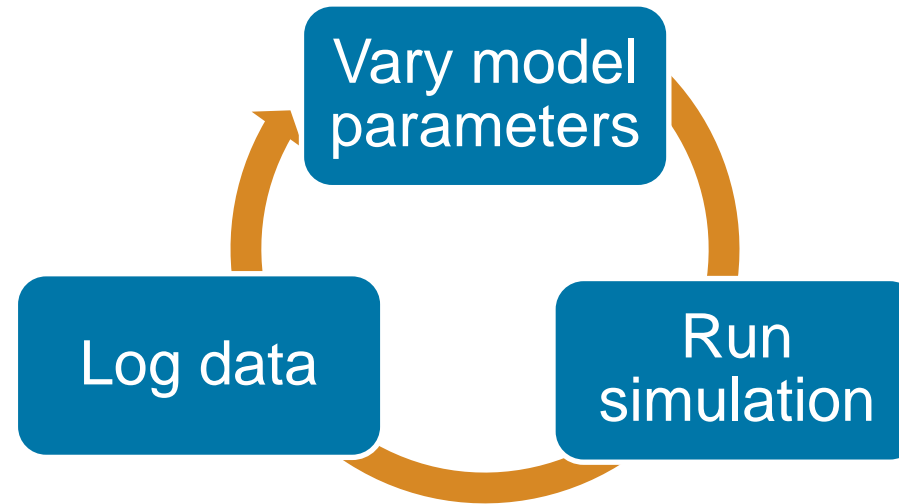
	EngTrqReq	EngSpdR...	SpkAdvOfst
1	60	2000	-30
2	128	2500	15
3	94	2750	8
4	111	2875	-19
5	77	2625	-11
6	144	2125	4
7	85	2563	-21
8	119	3313	-28
9	68	2938	21

### Input features

Engine speed (RPM)  
Ignition timing  
Throttle position  
Wastegate valve

### Response

Engine Torque



# Synthetic Data Generation

## Design of Experiments

**Execute Simulation**

For this demo, the engine model is brought from the workspace to the Simulink model.

- <https://www.mathworks.com/help/aut>
- <https://www.mathworks.com/help/aut>

Default setting: Simulation Type => "Rapid Acceleration"

```

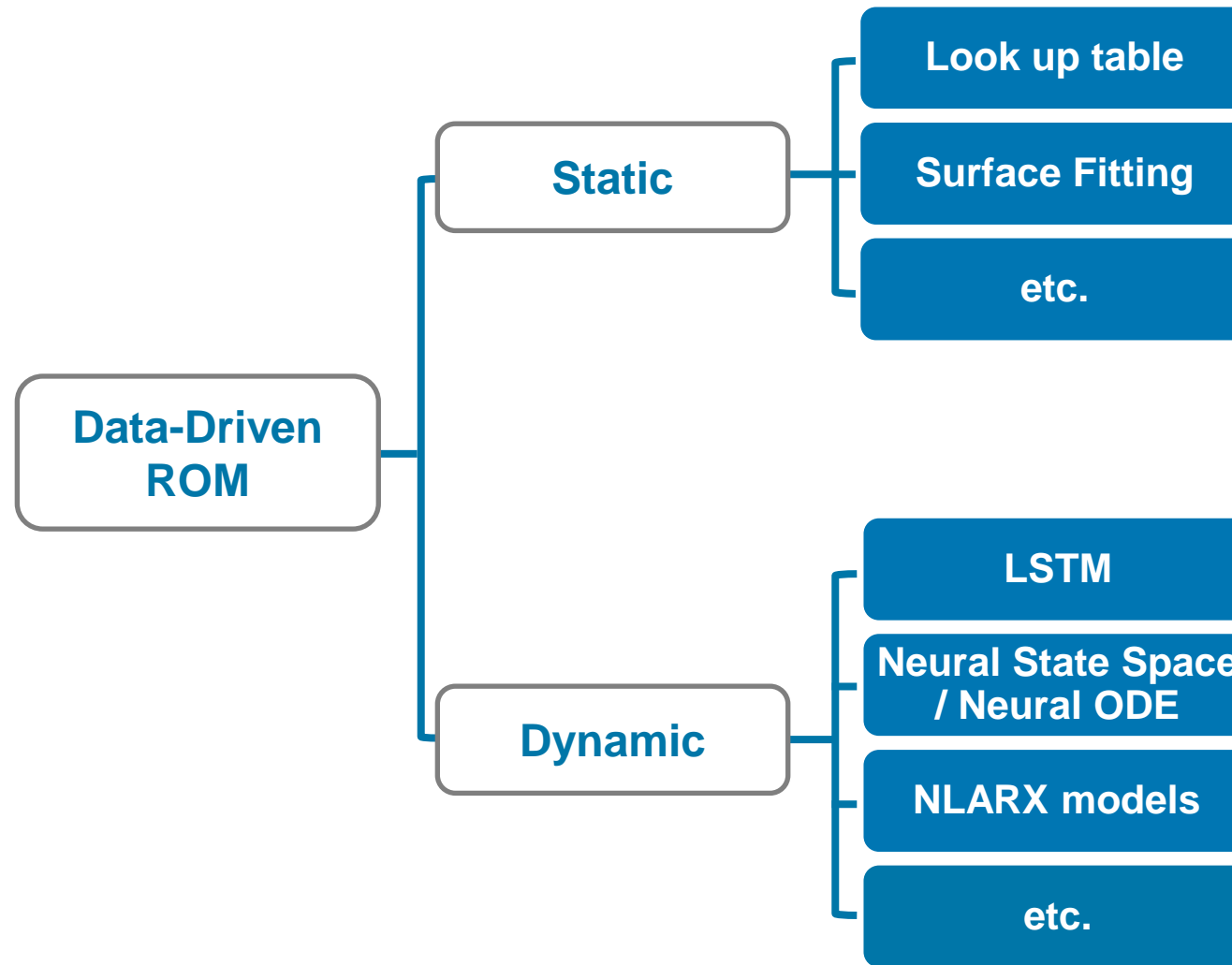
36 open_system('DoE_Engine_Test.slx')
37 % Run parallel simulations, 3000
38 simout = helper.simulation(DoE, 3000);

```

Starting parallel pool (parpool) using local 'labwin' engine.  
 Connected to the parallel pool (number of workers: 32).  
 [10-Jun-2022 03:33:09] Checking for updates.  
 [10-Jun-2022 03:33:09] Starting Simulink on parallel workers...  
 [10-Jun-2022 03:33:56] Loading project on parallel workers...  
 [10-Jun-2022 03:33:56] Configuring simulation cache folder on parallel workers...  
 [10-Jun-2022 03:34:18] Transferring base workspace variables used in the model to parallel workers...  
 [10-Jun-2022 03:34:20] Loading model on parallel workers...  
 [10-Jun-2022 03:34:53] Running simulations...

**Check and Save Simulation Result**

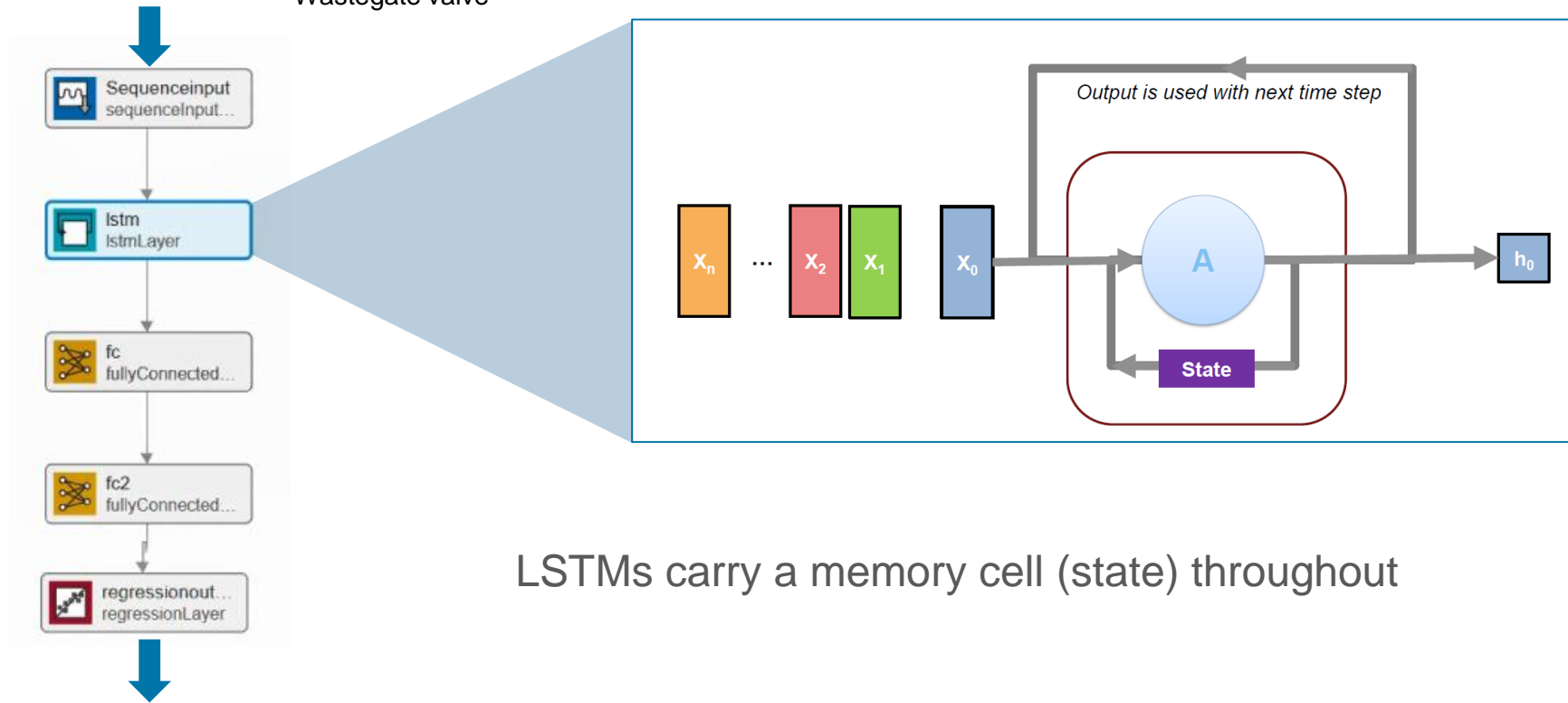
# Data-driven ROM



# AI-based ROM using LSTMs

*Capture time dependencies in time-series data*

- Engine speed (RPM)
- Ignition timing
- Throttle position
- Wastegate valve



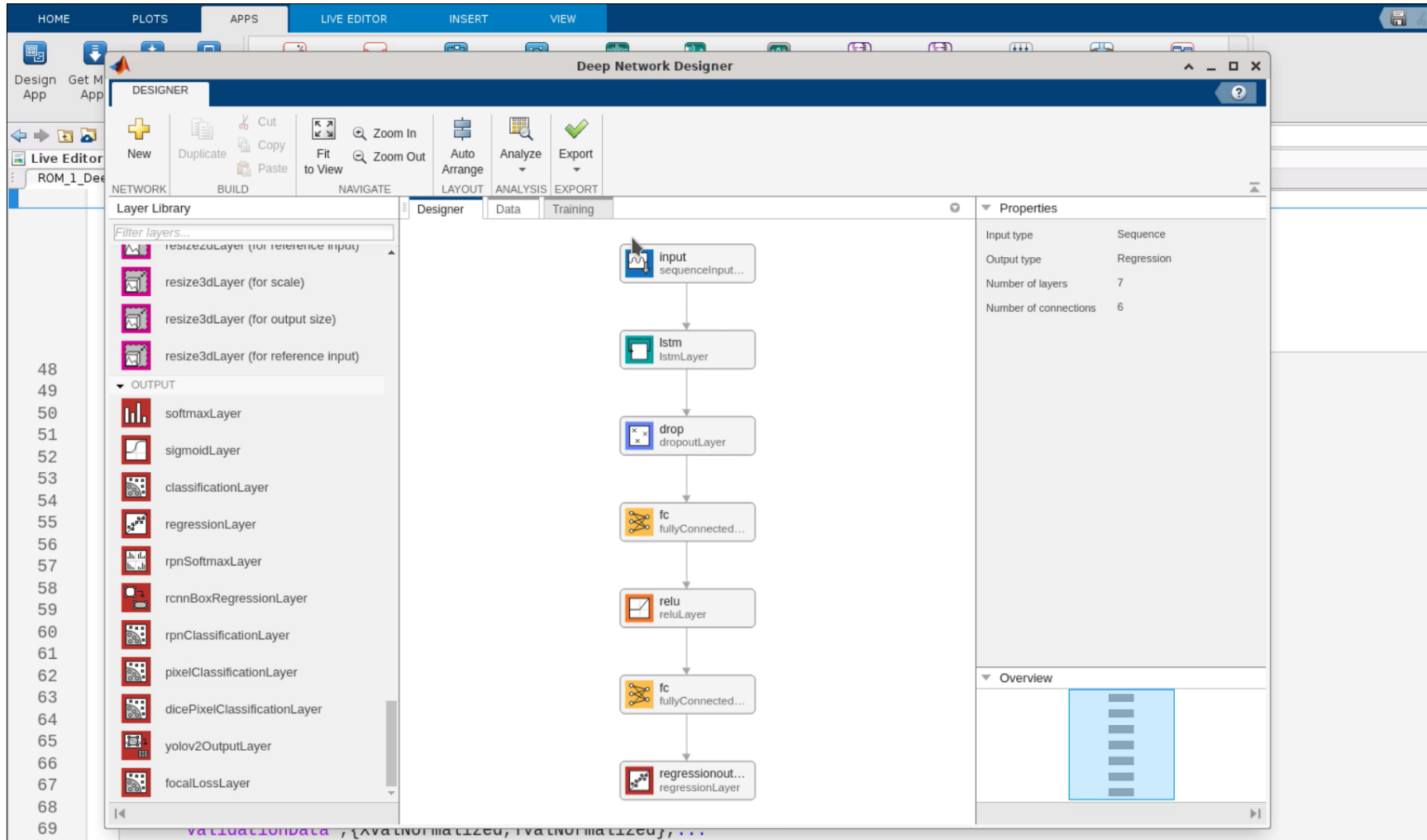
**Outputs** ▪ Engine torque

LSTMs carry a memory cell (state) throughout



# AI-based ROM using LSTMs

*Capture time dependencies in time-series data*



Data Preparation

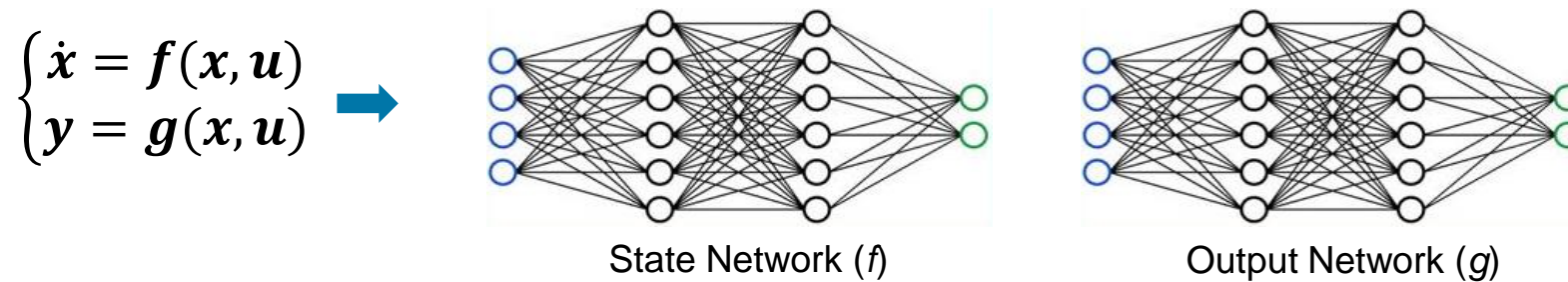
**AI Modeling**

Simulation & Test

Deployment

# AI-based ROM using Neural State Space

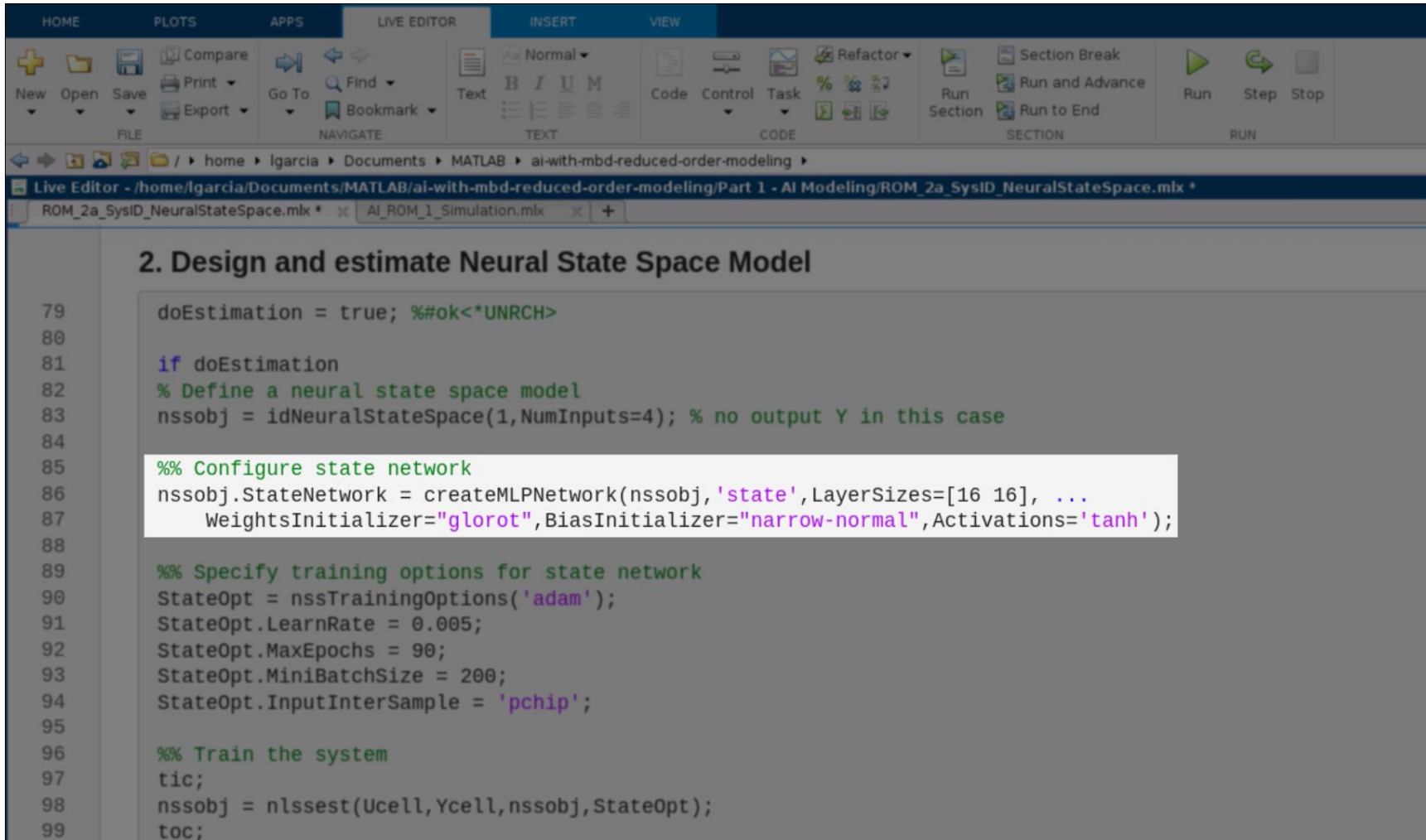
*Create DL-based nonlinear state-space models without having to be a deep learning expert*



- The nonlinear state function  $f$  and nonlinear output function  $g$  are feedforward neural networks that learn from data
- Popularly known as Neural ODE in deep learning community

# AI-based ROM using Neural State Space

*Create DL-based nonlinear state-space models without having to be a deep learning expert*



```
79 doEstimation = true; %#ok<*UNRCH>
80
81 if doEstimation
82 % Define a neural state space model
83 nssobj = idNeuralStateSpace(1,NumInputs=4); % no output Y in this case
84
85 %% Configure state network
86 nssobj.StateNetwork = createMLPNetwork(nssobj,'state',LayerSizes=[16 16], ...
87     WeightsInitializer="glorot",BiasInitializer="narrow-normal",Activations='tanh');
88
89 %% Specify training options for state network
90 StateOpt = nssTrainingOptions('adam');
91 StateOpt.LearnRate = 0.005;
92 StateOpt.MaxEpochs = 90;
93 StateOpt.MiniBatchSize = 200;
94 StateOpt.InputInterSample = 'pchip';
95
96 %% Train the system
97 tic;
98 nssobj = nlssest(Ucell,Ycell,nssobj,StateOpt);
99 toc;
```

Data Preparation

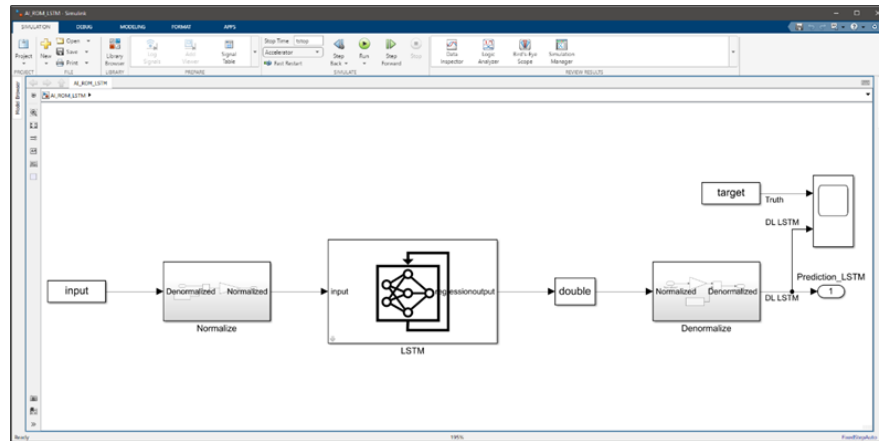
AI Modeling

Simulation & Test

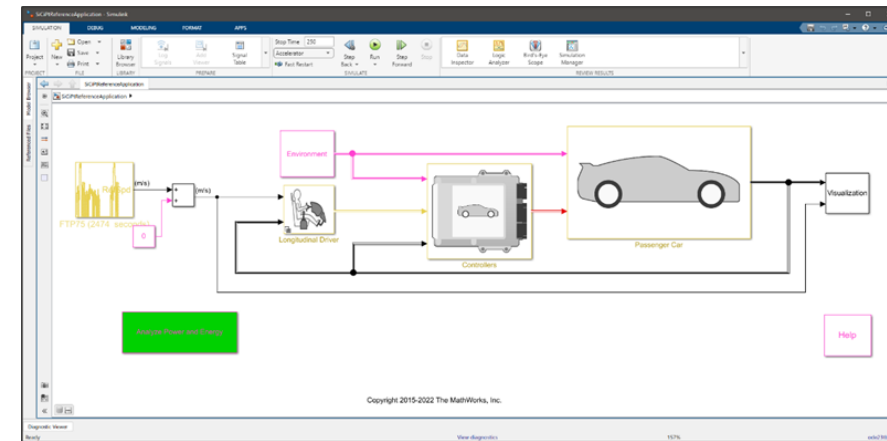
Deployment

# Integrate your AI model for system-level simulation and test

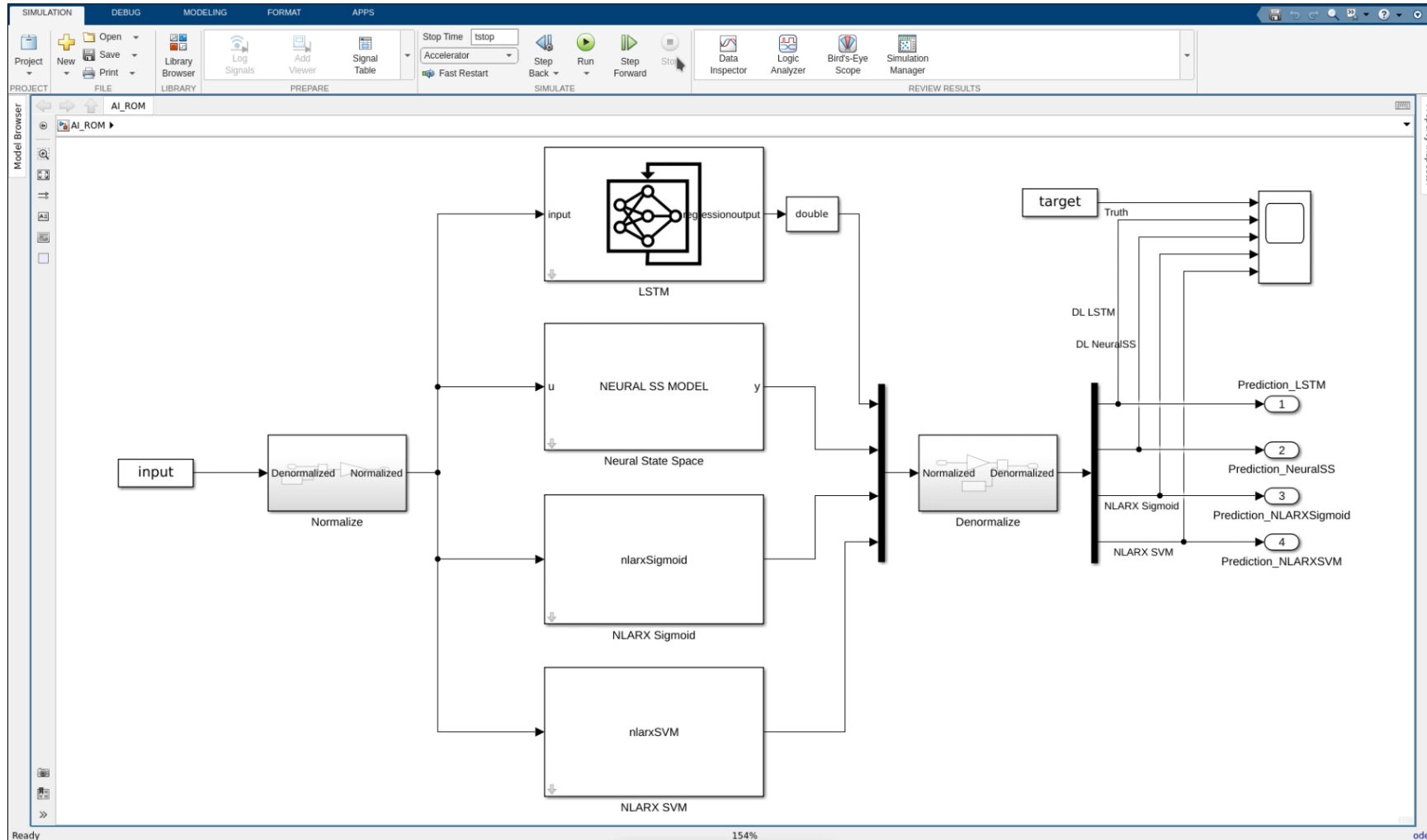
## Integration of trained AI model into Simulink



## System-level simulation



# Integration of trained AI models into Simulink



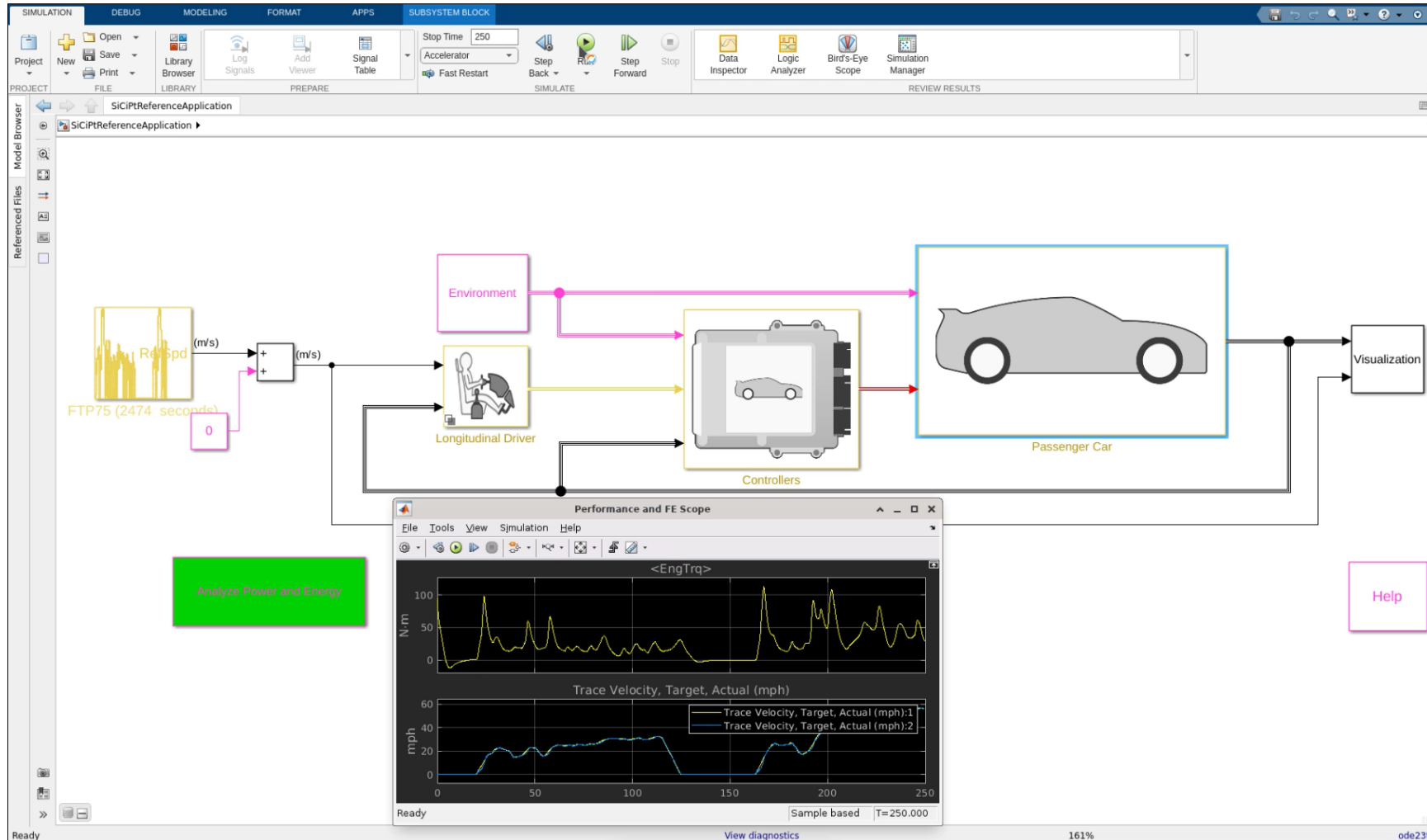
Data Preparation

AI Modeling

Simulation & Test

Deployment

# System-level simulation



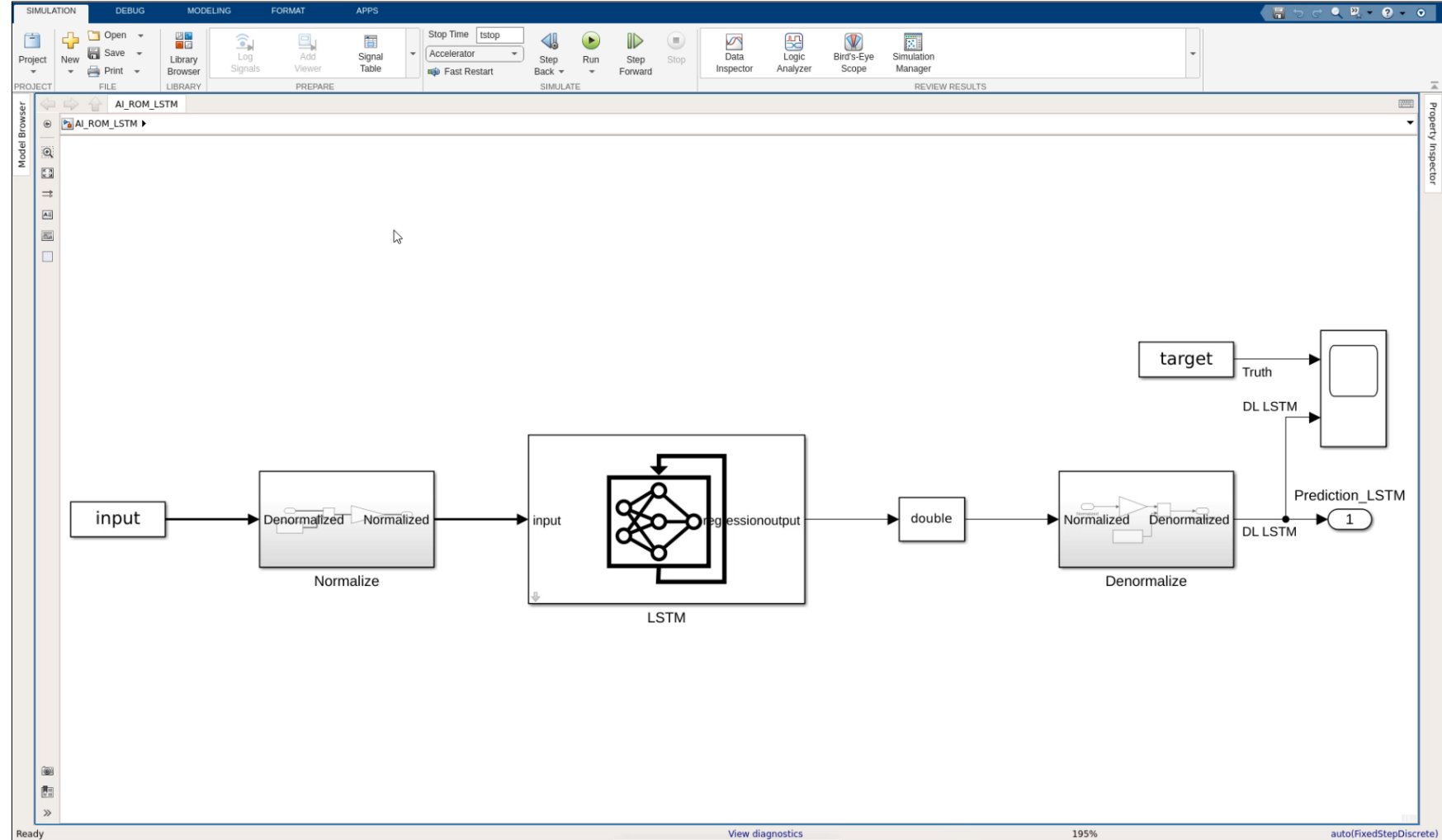
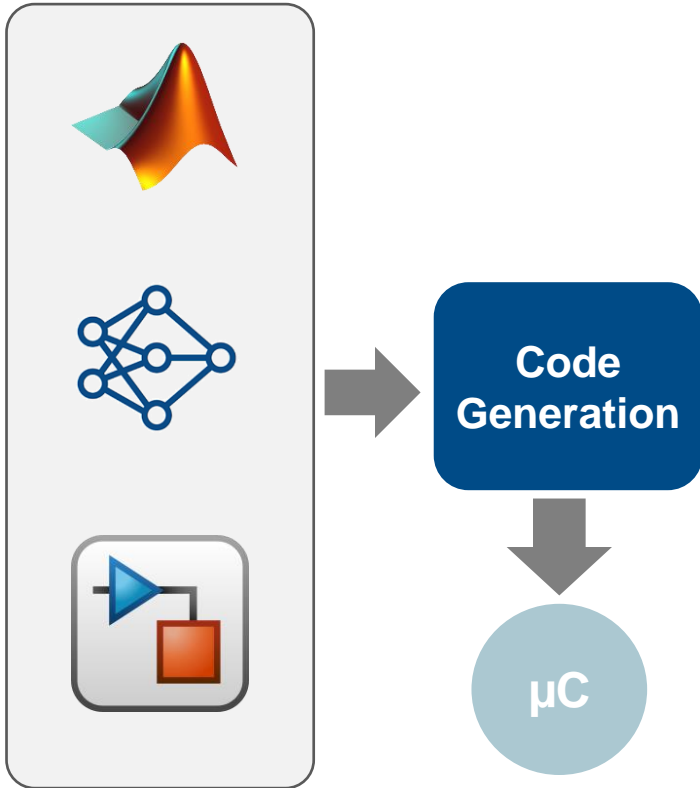
Data Preparation

AI Modeling

**Simulation & Test**

Deployment

# Generate Library-Free C Code for Deep Learning Networks



Data Preparation

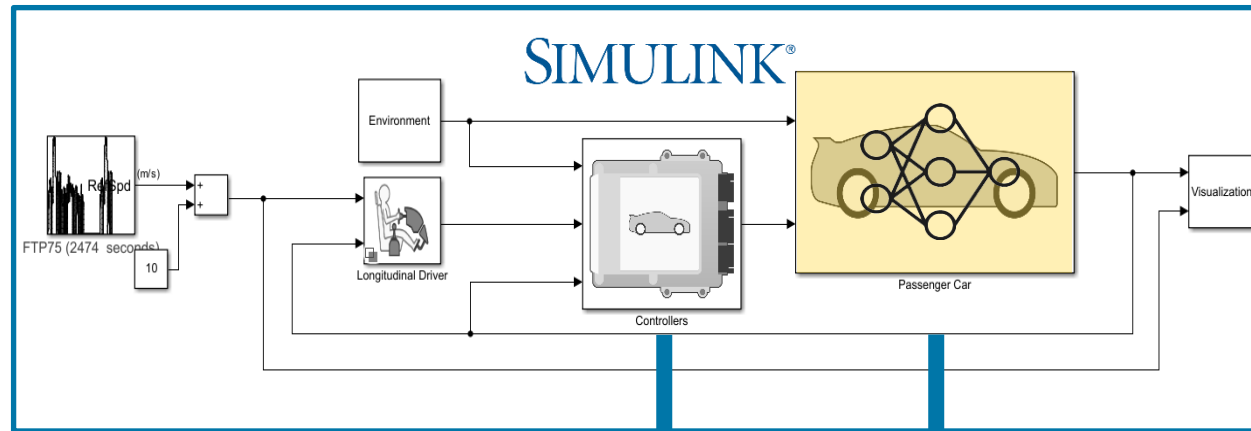
AI Modeling

Simulation & Test

Deployment

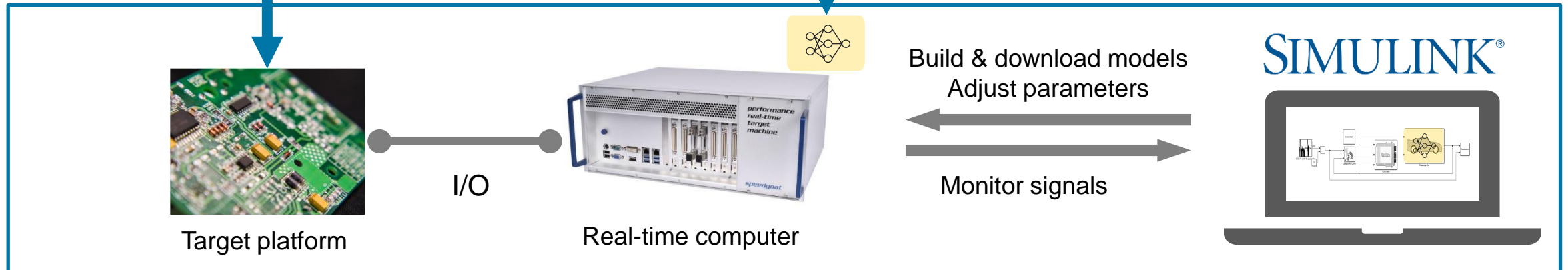
# Hardware-in-the-loop simulation

*System-level integration and test*



Code generation from algorithm

Code generation from plant model



Data Preparation

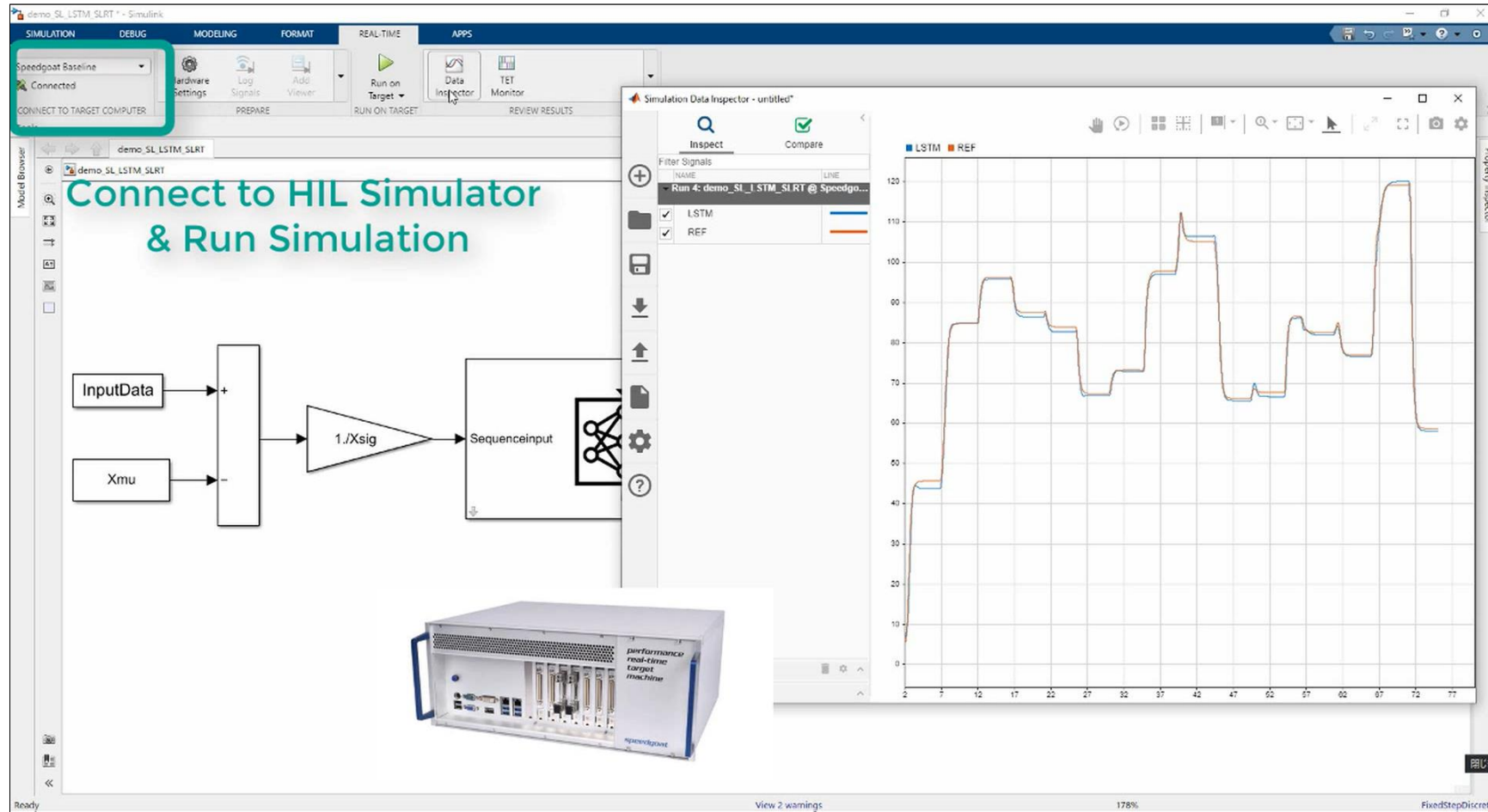
AI Modeling

Simulation & Test

Deployment



# Hardware-in-the-loop simulation













Data Preparation

AI Modeling

Simulation & Test

Deployment

# Manage AI tradeoffs for your system

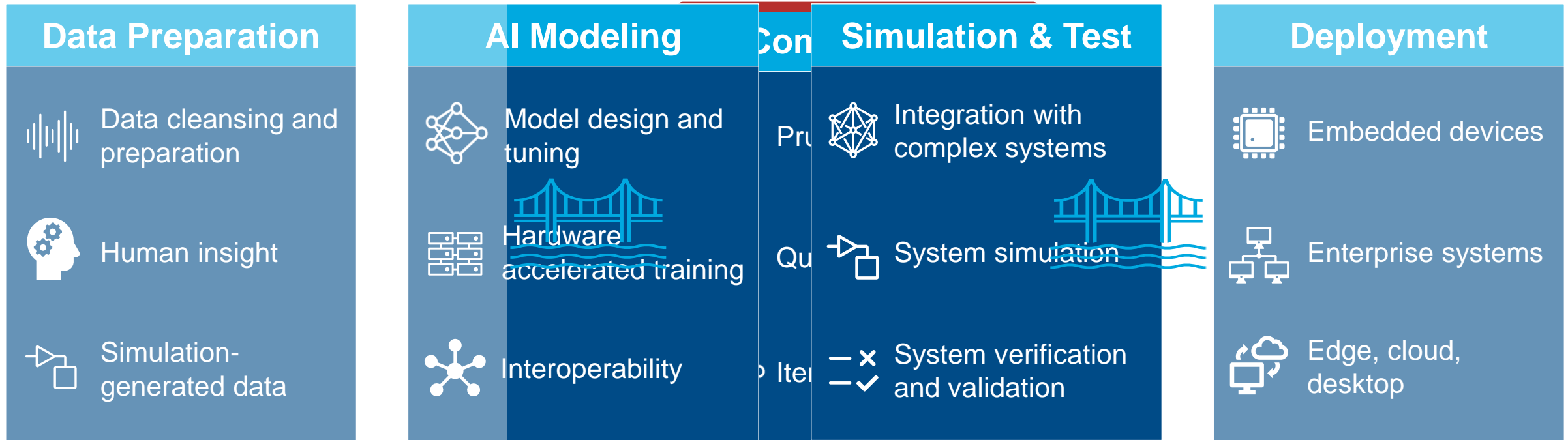
	<b>LSTM</b> Long Short-Term Memory Network	<b>Neural SS</b> Neural State Space (Neural ODE)
Training Speed	 *	
Interpretability		
Inference Speed		
Model Size		
Accuracy (RSME)		

*Results are specific to Vehicle Engine ROM example*

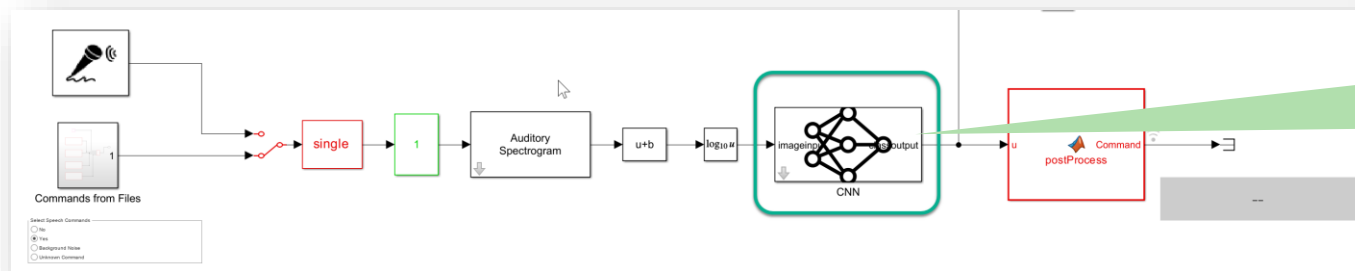
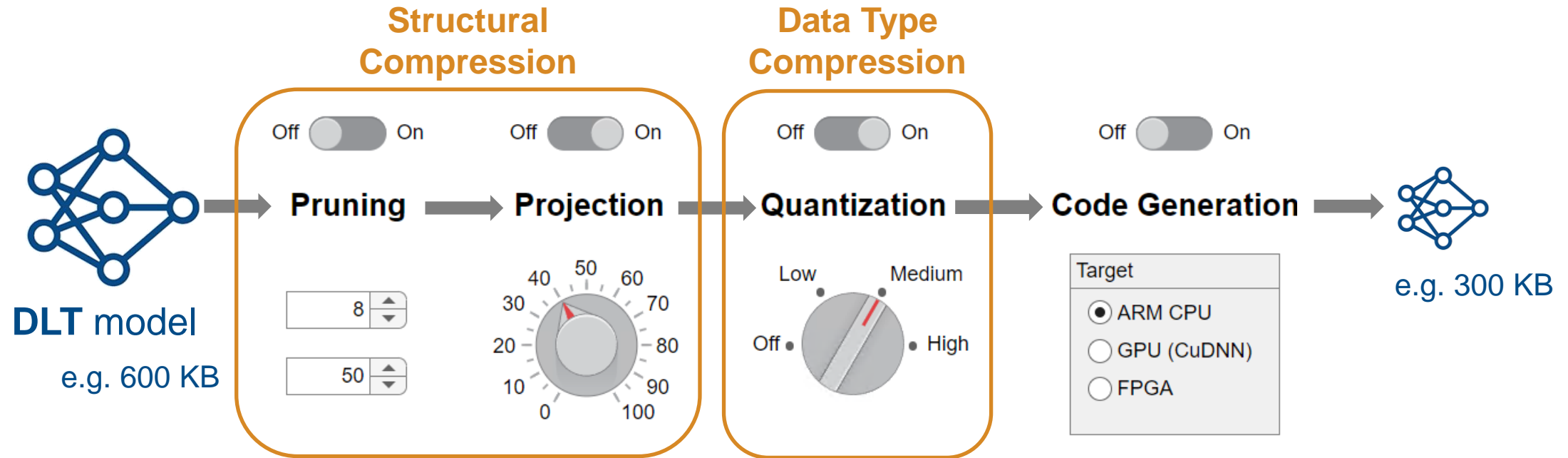


\*  if trained using a GPU. Testing made with GPU NVIDIA A100

# Model compression bridges the gap between AI modelling and embedded deployment.



# Problem Statement: Reduce model footprint and accelerate inference of DL models

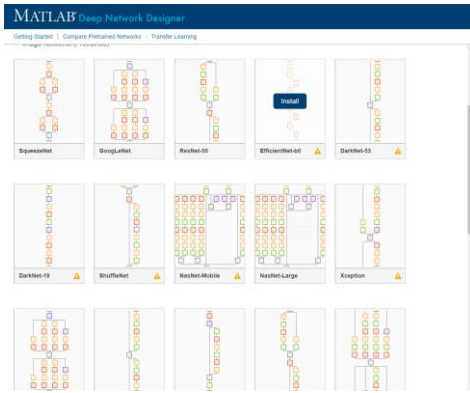


*“model is 600kb and want to reduce it to smaller. What to do?”*

# Workflow steps to compress Deep Neural Nets

1

Deep Network Designer



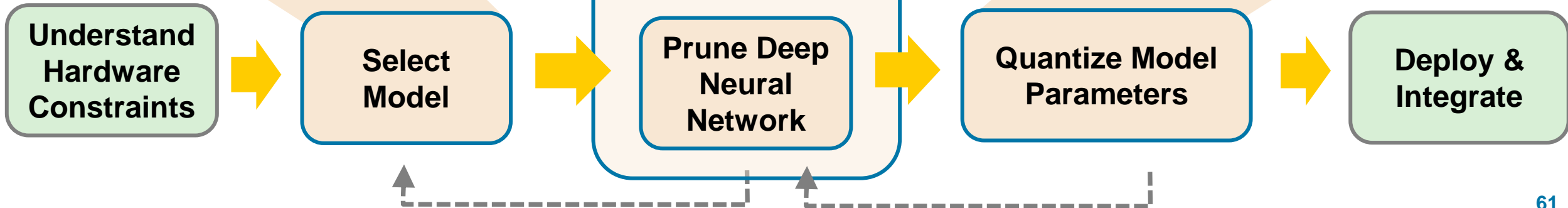
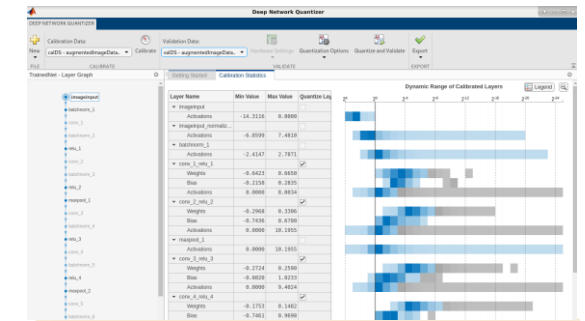
2

Pruning

```
taylorPrunableNetwork(net)
```

3

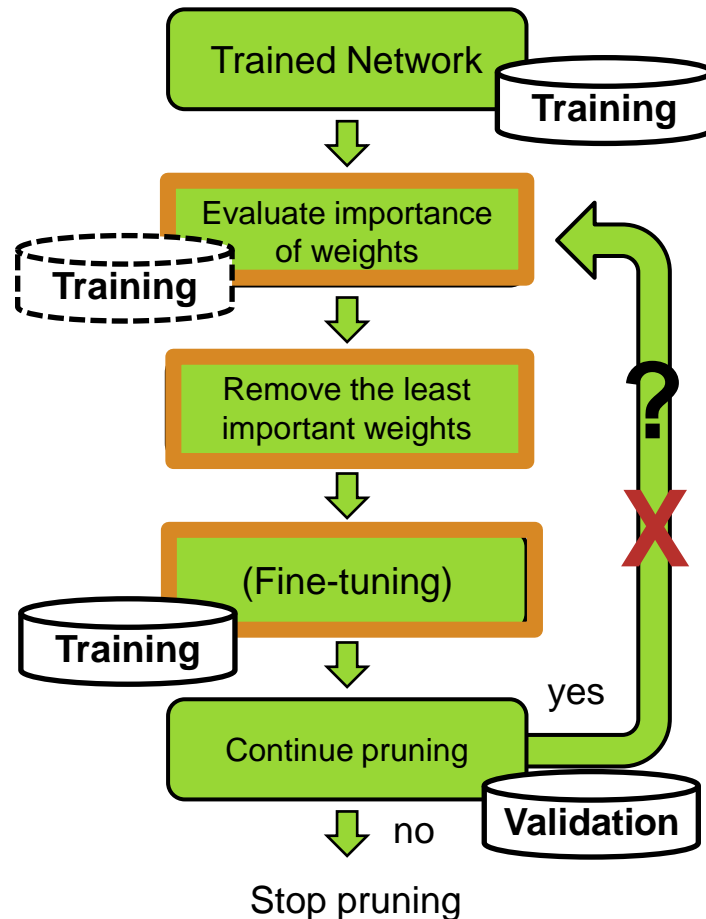
Deep Network Quantizer



# Pruning algorithms follow a common process but can have lots of small variations

No clear winner according to literature 😞

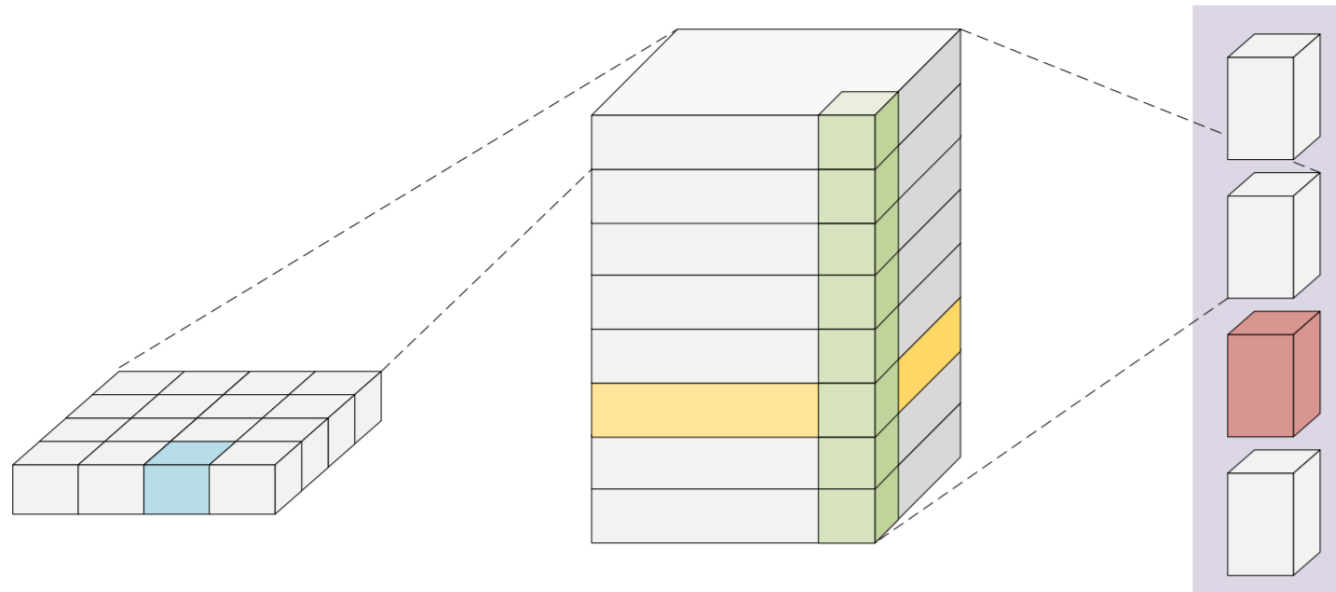
Best practices available 😊



Execution time vs. effect vs. data required

- Scoring
  - Absolute weight value
  - Gradient-based metric
  - Activations-based metric
- Pruning criteria
  - local (uniformly X% per layer)
  - global (X% across whole network)
- Fine-tuning, yes/no
- Scheduling
  - One-shot
  - Iteratively

# Which parts of the network can be pruned?



element-wise

individual connections

introduces sparsity

**UNSTRUCTURED**

channel-wise shape-wise filter-wise

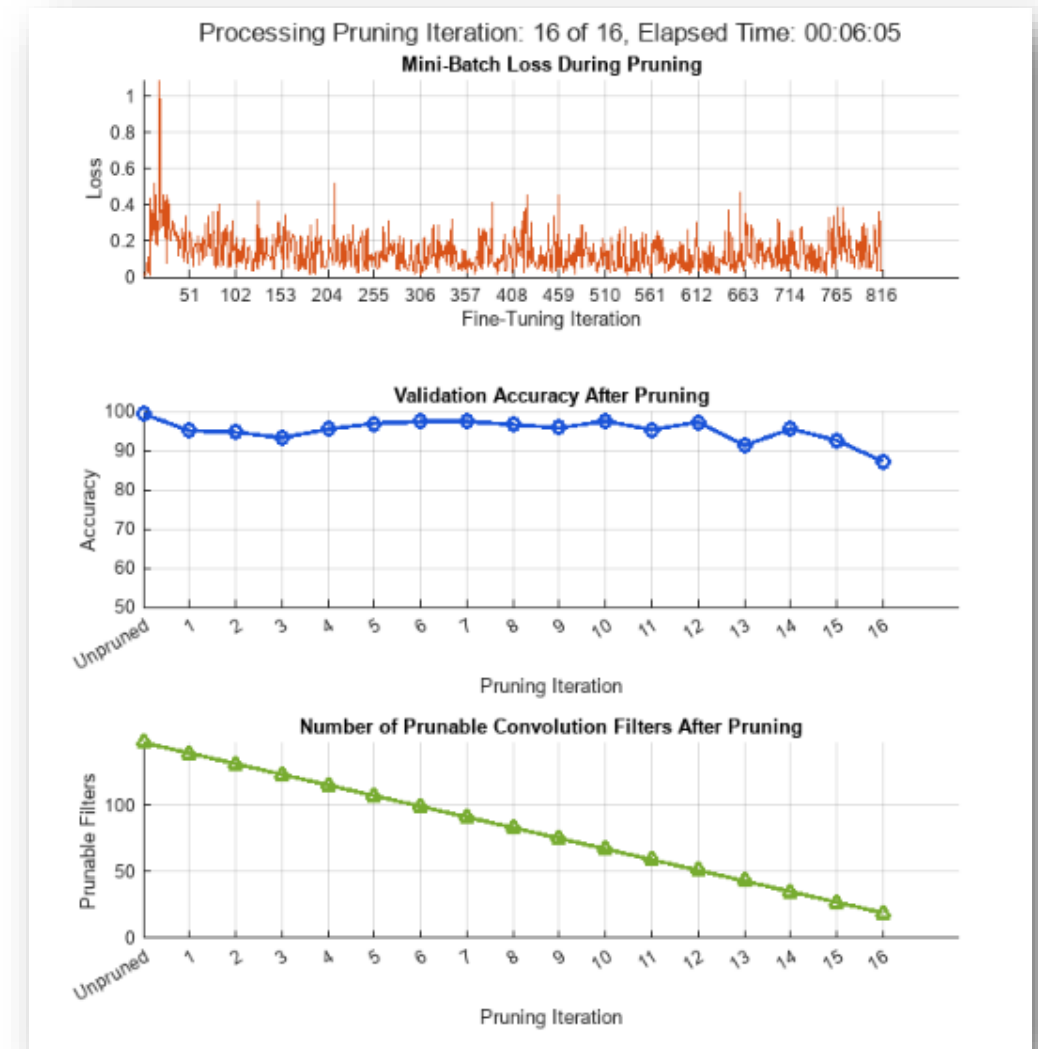
e.g. conv. filters,  
neurons in FC layer

**STRUCTURED**

layer-wise

# Taylor Pruning uses gradient score and eliminates number of filters in convolutional layers

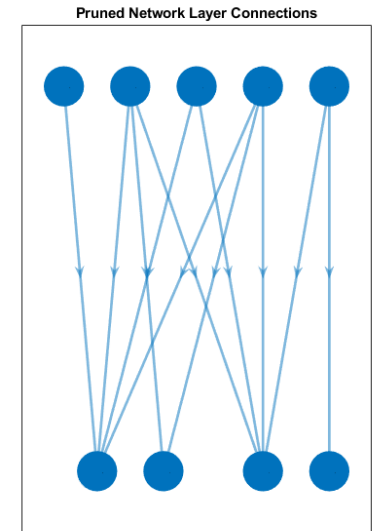
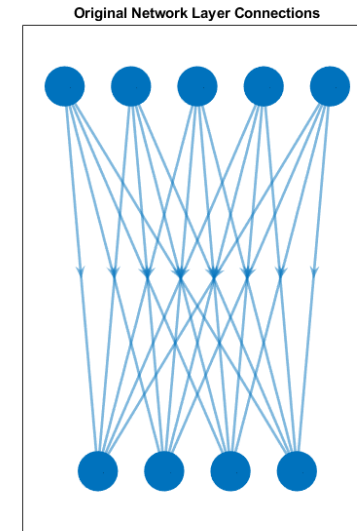
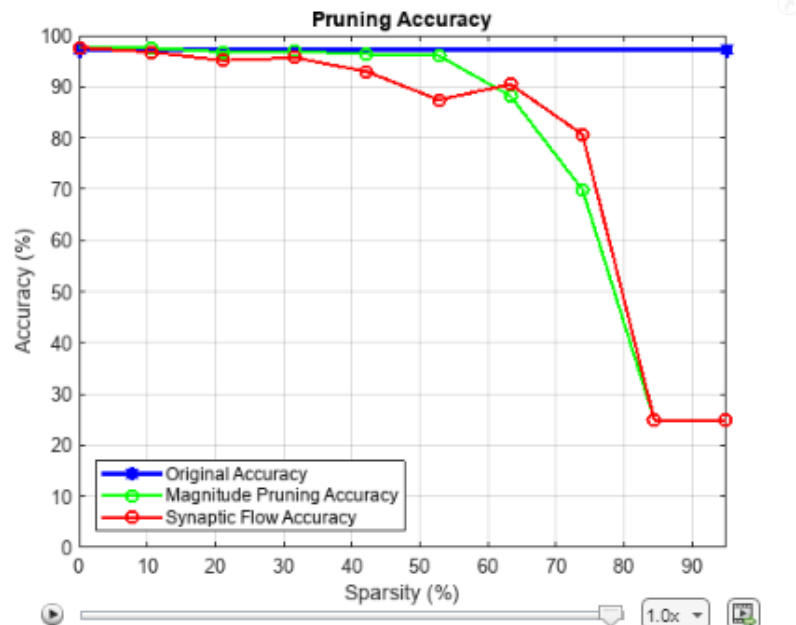
- Gradient-based method to estimate filter “importance” using first-order Taylor expansion
- Prune less important filters to reduce model size while maintaining predictive power
- **STRUCTURED** approach
- Fine-tune pruned model with data





# Parameter Pruning zeros out lower score connections

- Calculate numerical scores to rank the connections in the network
- Iteratively remove less “important” connections
- **UNSTRUCTURED** approach



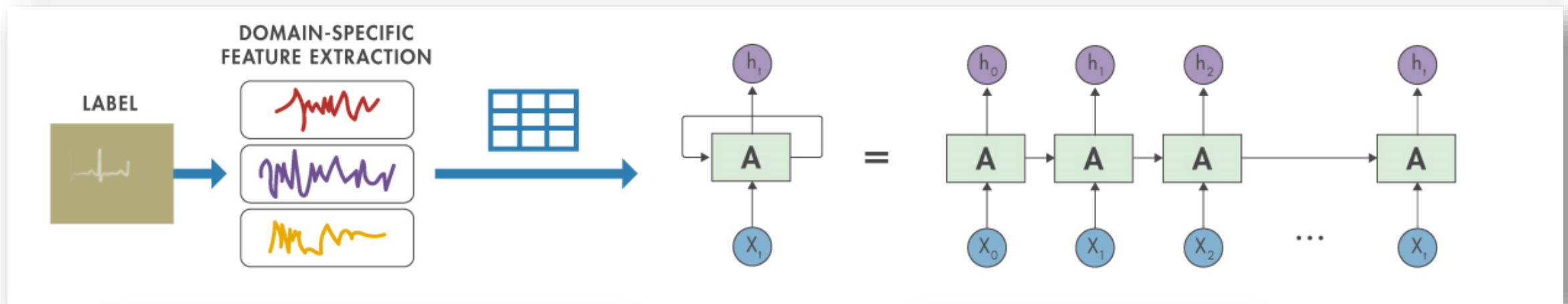
Examples are:

Magnitude Score

SynFlow Score: synaptic flow scores

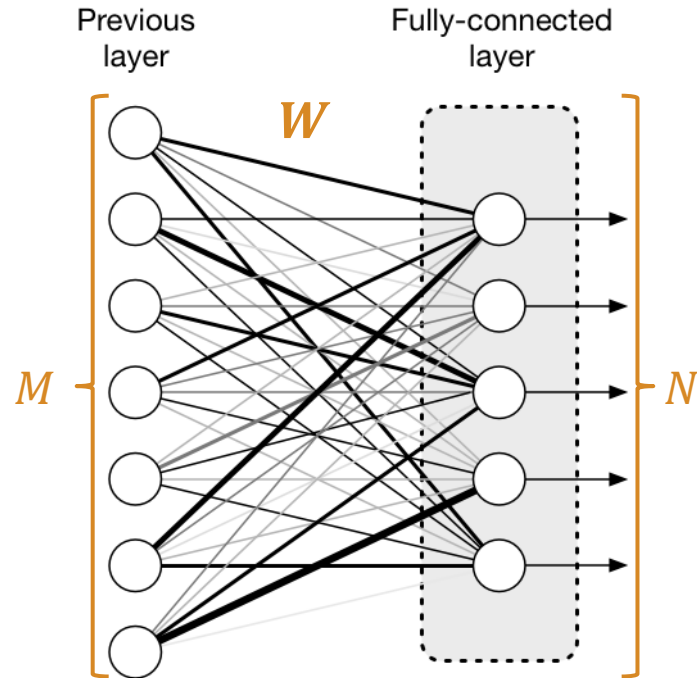
# Structural pruning reduces problem dimensions via projection into subspace

For example, a LSTM Networks



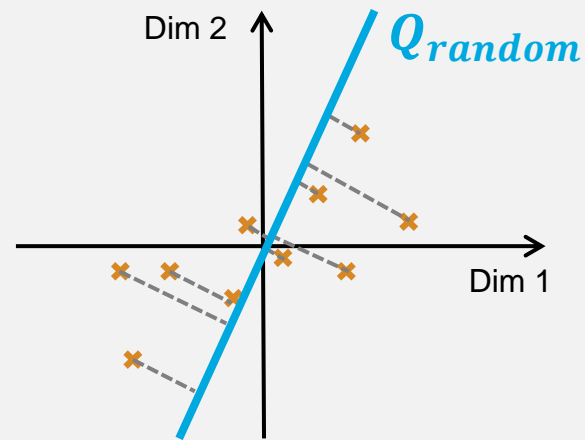
## Two-step approach: Projection compression with neuron PCA

High-dimensional space of input and output neurons is underutilized



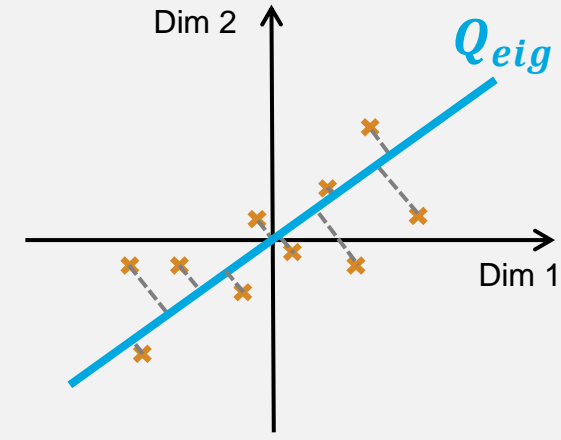
$W$  is an  $N$ -by- $M$  matrix

Dimensionality reduction via projection into subspace



projected layer

Minimize projection error via principal component analysis (PCA) of neurons



neuronPCA

# Structural compression of LSTM layers to reduce model size

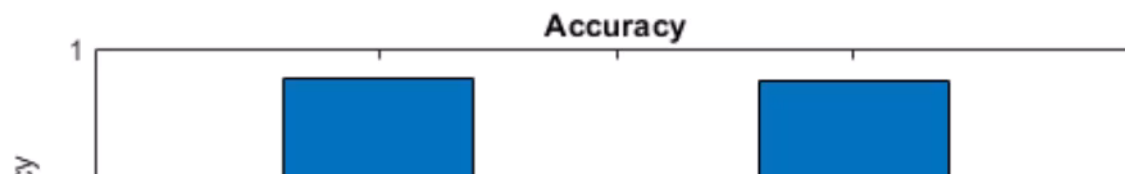
## Compress Neural Network Using Projection

This example shows how to compress a neural network using projection and principal component analysis.

To compress a deep learning network, you can use *projected layers*. The layer introduces learnable projector matrices  $Q$ , replaces multiplications of the form  $Wx$ , where  $W$  is a learnable matrix, with the multiplication  $WQQ^T x$ , and stores  $Q$  and  $W' = WQ$  instead of storing  $W$ . Projecting  $x$  into a lower dimensional space using  $Q$  typically requires less memory to store the learnable parameters and can have similarly strong prediction accuracy. A projected deep neural network can also exhibit faster forward passes when run on the CPU or deployed to embedded hardware using library-free C or C++ code generation.

The `compressNetworkUsingProjection` function compresses a network by projecting layers into smaller parameter subspaces. For optimal initialization of the projected network, the function projects the learnable parameters of projectable layers into a subspace that maintains the highest variance in neuron activations. After you compress a neural network using projection, you can then fine-tune the network to increase the accuracy.

This chart shows the effect of projection and fine tuning on a trained network. In this case, the projected network has significantly fewer learnable parameters at the cost of classification accuracy. The fine-tuned projected network yields similar classification accuracy to the original network.



# Results from compression of LSTM layers to reduce model size

Compress the network.

```
netProjected = compressNetworkUsingProjection(net,mbq);
```

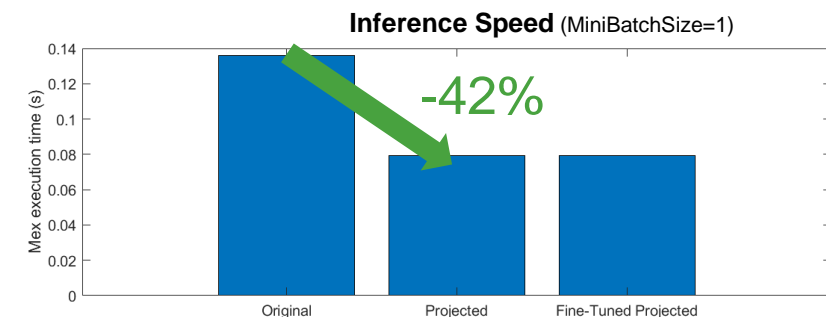
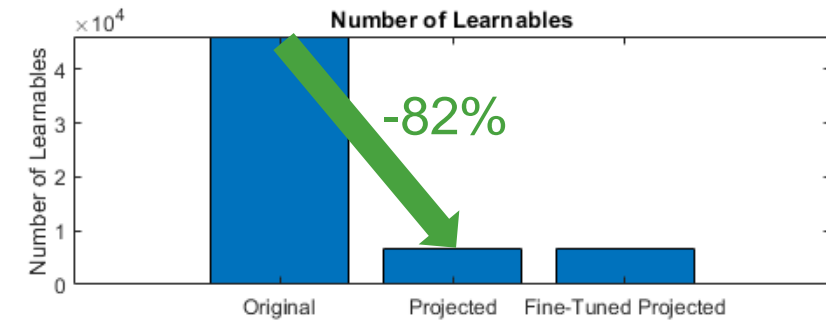
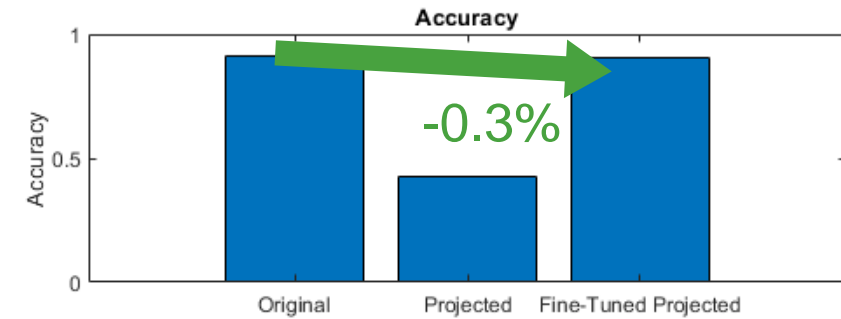
Compressed network has **82.4% fewer learnable parameters.**  
 Projected layers explain on average 96.6% of layer activation variance.

Optionally, precompute neuronPCA analysis for efficient experimentation:

```
npca = neuronPCA(netOriginal,mbqTrain,VerbosityLevel="steps");
```

```
Computing layer activations and covariance matrices...
Computing eigenvalues and eigenvectors...
neuronPCA analyzed 1 layers: "lstm"
```

Doc Example: [Compress Neural Network Using Projection](#)  
 (Seq-2-One Classification on Japanese Vowels data set)



# Deep Network Quantizer transparently applies quantization



DEEP NETWORK QUANTIZER

Calibration Data: calibrationDatastore - augme... Calibrate

Validation Data: validationDatastore - augmen... Hardware Settings Quantization Options Quantize and Validate Export

FILE CALIBRATE VALIDATE EXPORT

Net - Layer Graph

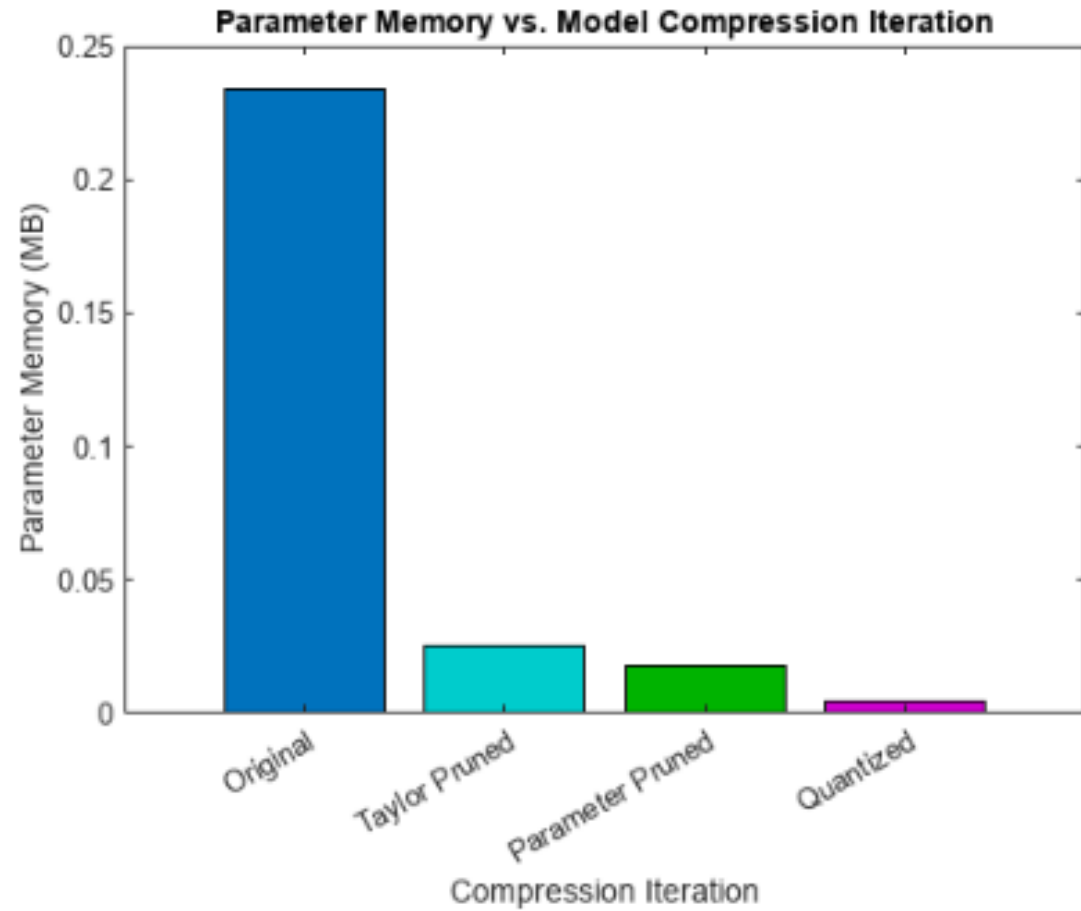
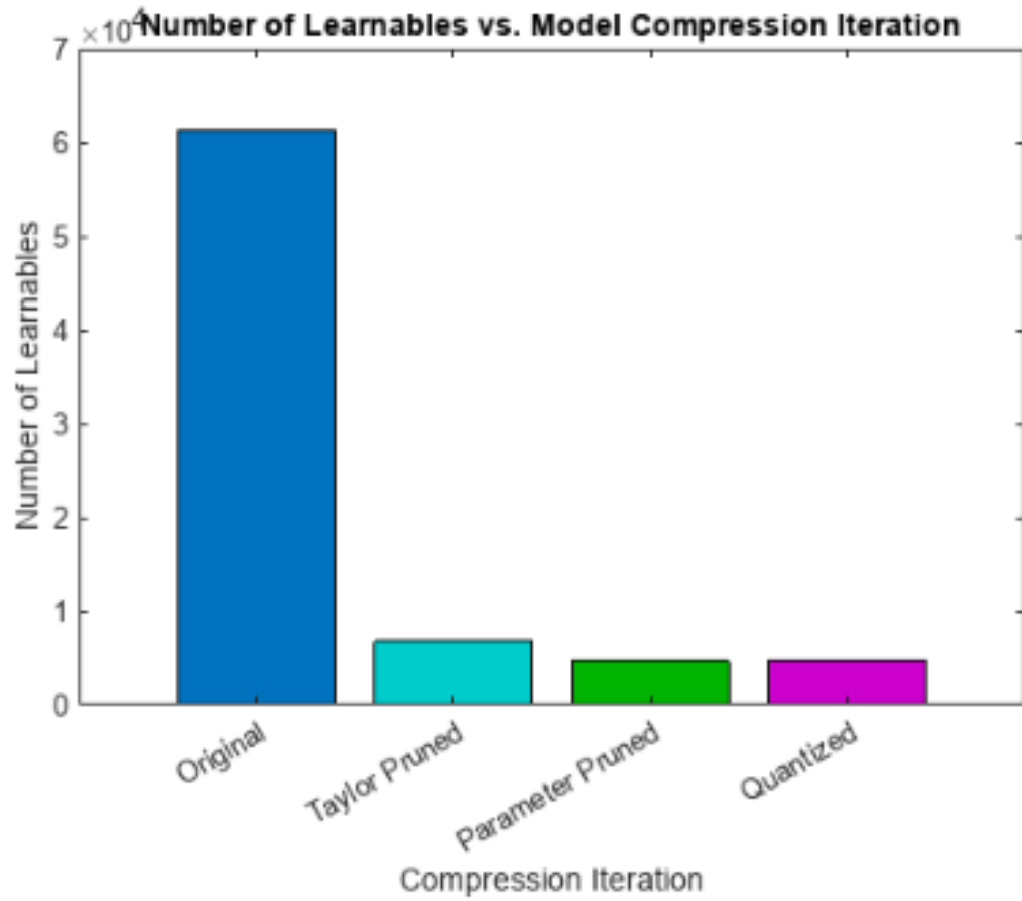
Getting Started Calibration Statistics

Layer Name	Min Value	Max Value	Quantize La
imageinput			<input type="checkbox"/>
batchnorm_1			<input type="checkbox"/>
conv_1			<input type="checkbox"/>
batchnorm_2			<input type="checkbox"/>
relu_1			<input type="checkbox"/>
conv_2			<input type="checkbox"/>
batchnorm_3			<input type="checkbox"/>
relu_2			<input type="checkbox"/>
maxpool_1			<input type="checkbox"/>
conv_3			<input type="checkbox"/>
batchnorm_4			<input type="checkbox"/>
relu_3			<input type="checkbox"/>
maxpool_1			<input type="checkbox"/>

Dynamic Range of Calibrated Layers

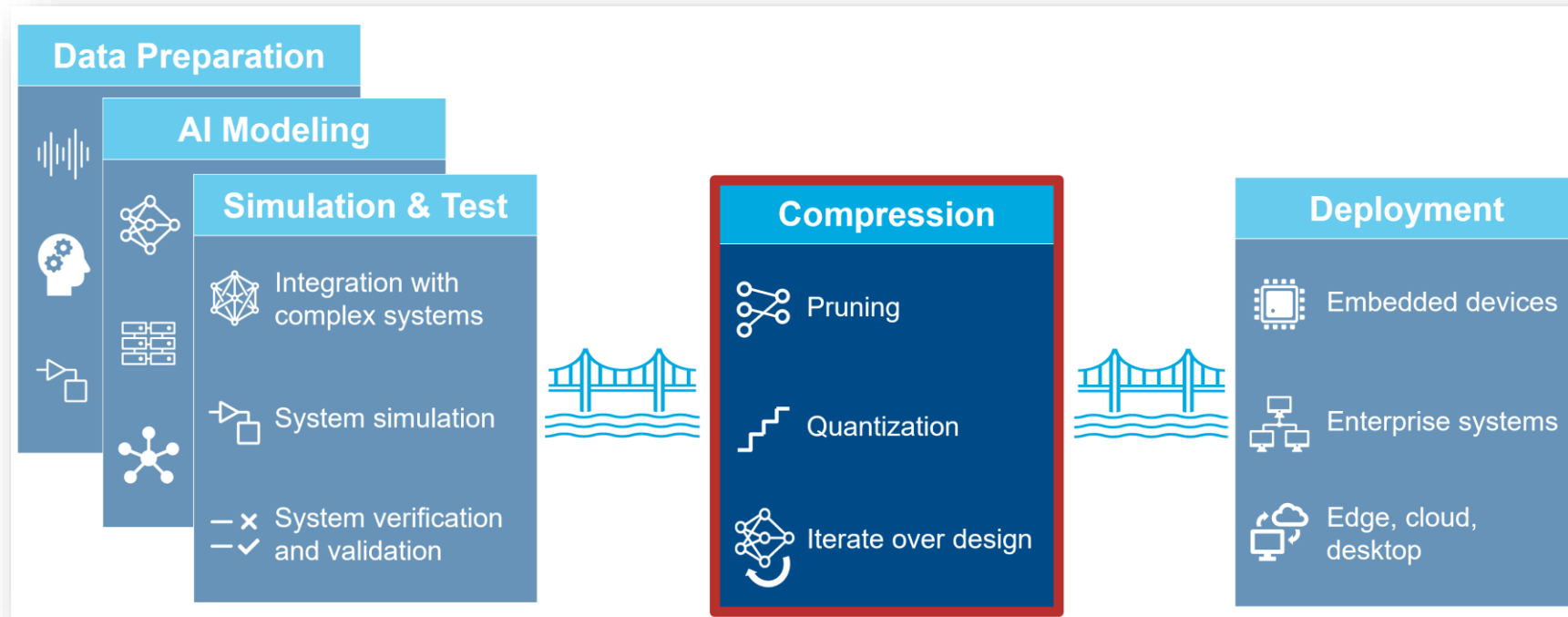
Legend

# Reduce learnable parameters



# Conclusion

*Compress AI-based reduced-order engine model for deployment*



- **Integrate trained AI model into Simulink** for system-level simulation together with first-principles components
- **Generate C code and perform HIL tests**
- **Deploy compressed TinyML model** to embedded target



# Thank you!

## Want to learn more?

### Quantization

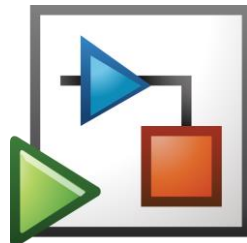
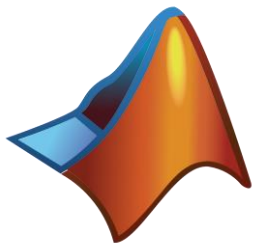
<https://www.mathworks.com/discovery/quantization.html>

### MATLAB for Deep Learning

<https://www.mathworks.com/solutions/deep-learning.html>

### More applications and capabilities

<https://www.mathworks.com/solutions.html>



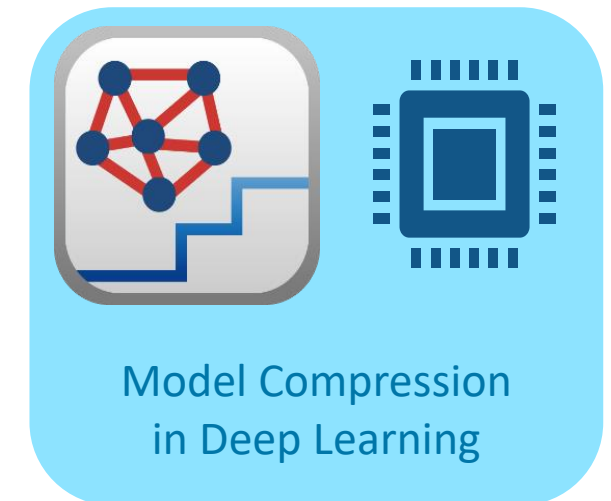
Brenda Zhuang

 @Brenda-Zhuang



Greg Coppenrath

 @GregCoppenrath





# Copyright Notice

This multimedia file is copyright © 2023 by tinyML Foundation. All rights reserved. It may not be duplicated or distributed in any form without prior written approval.

tinyML<sup>®</sup> is a registered trademark of the tinyML Foundation.

[www.tinyml.org](http://www.tinyml.org)



# Copyright Notice

This presentation in this publication was presented as a tinyML® Talks webcast. The content reflects the opinion of the author(s) and their respective companies. The inclusion of presentations in this publication does not constitute an endorsement by tinyML Foundation or the sponsors.

There is no copyright protection claimed by this publication. However, each presentation is the work of the authors and their respective companies and may contain copyrighted material. As such, it is strongly encouraged that any use reflect proper acknowledgement to the appropriate source. Any questions regarding the use of any materials presented should be directed to the author(s) or their companies.

tinyML is a registered trademark of the tinyML Foundation.

**[www.tinyml.org](http://www.tinyml.org)**