



# Distributed AI for Intelligence at Edge

Presenter: Dinesh Verma  
IBM Fellow, Distributed AI  
IBM T. J. Watson Research Center  
Email: [dverma@us.ibm.com](mailto:dverma@us.ibm.com)

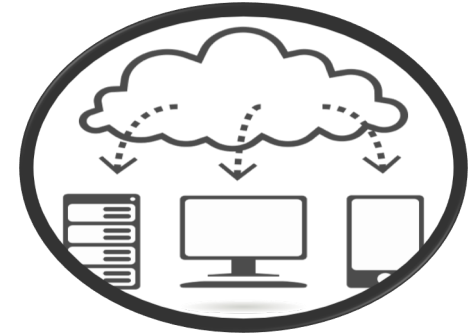


*Many promising technologies  
are converging together to  
transform the world within the  
next decade*

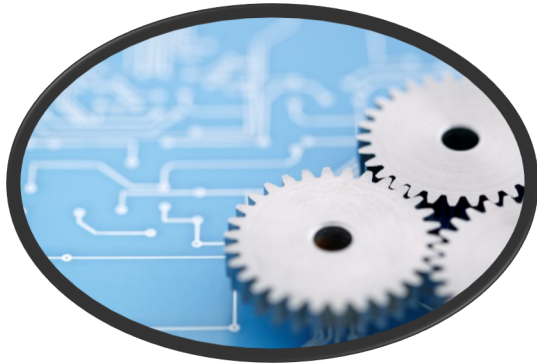
Cloud Computing



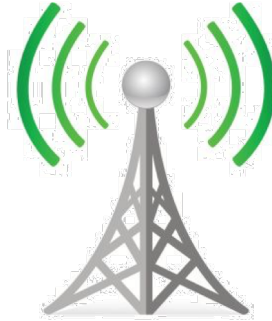
Edge



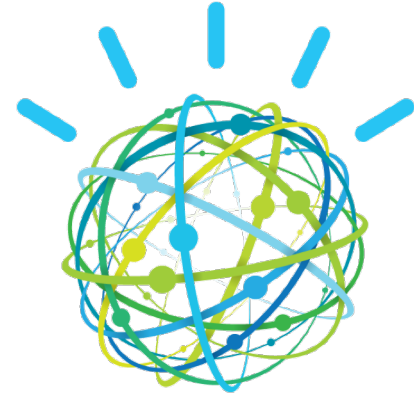
Internet of Things



5G



Artificial Intelligence



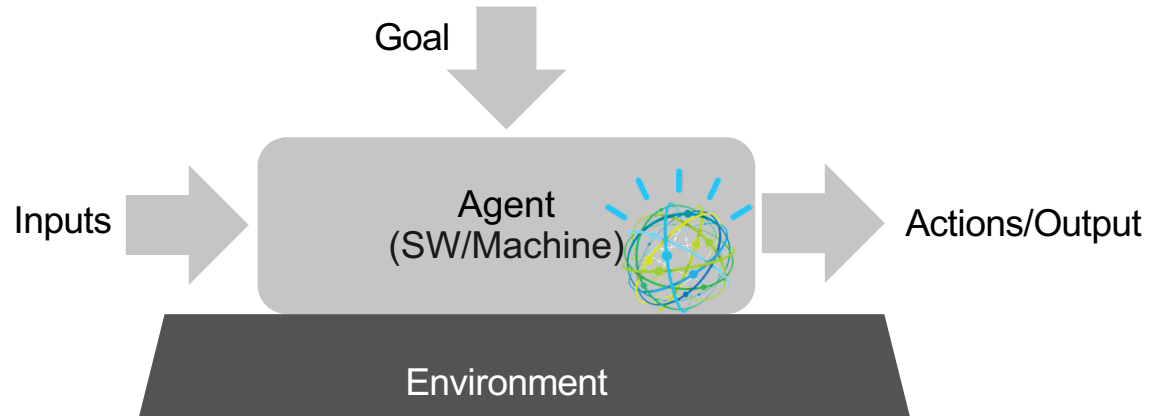
# What is AI (really)

## Definitions of Intelligence

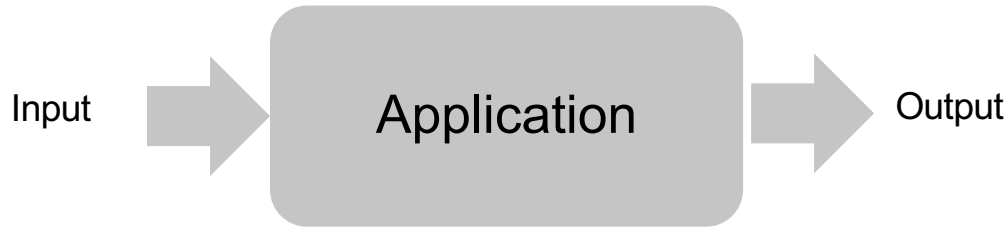
- Too many – For a collection of definitions, see Legg, Shane, and Marcus Hutter. "A collection of definitions of intelligence." *Frontiers in Artificial Intelligence and applications* 157 (2007): 17

## Informally

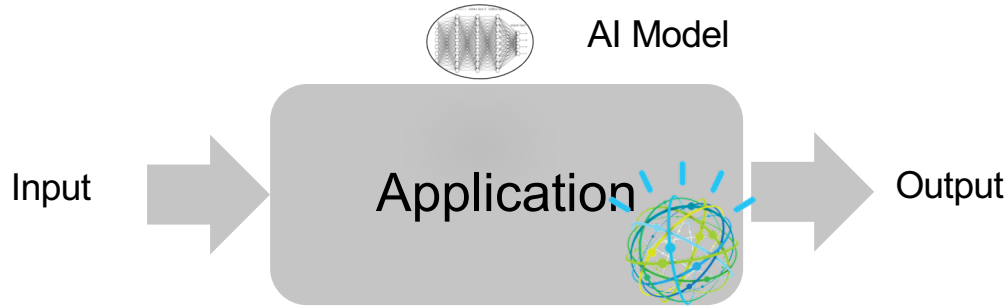
- “Intelligence is a measure of an agent’s ability to achieve goals in a wide range of environments.
- Artificial Intelligence is the ability of a machine to display intelligence, i.e. an ability to achieve its goal in a variety of environments



# Any Computer Application converts an input to an output

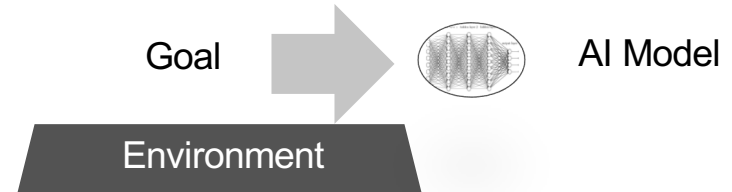


Traditional App:  
Conversion Logic coded in software



*Inference*

AI enabled App:  
Conversion Logic coded in a model  
Model defines how to attain goal in an environment



*Training*



# AI versus not AI

From the conversion of input to output, an AI enabled model and a non-AI model are equivalent

- Every Software Program can be converted to a equivalent Turing Machine Program
- Every AI model can be converted into a equivalent “Lookup table”

AI does not enable us to do anything we could not already do using a software encoded logic

What AI provides is non-functional attributes

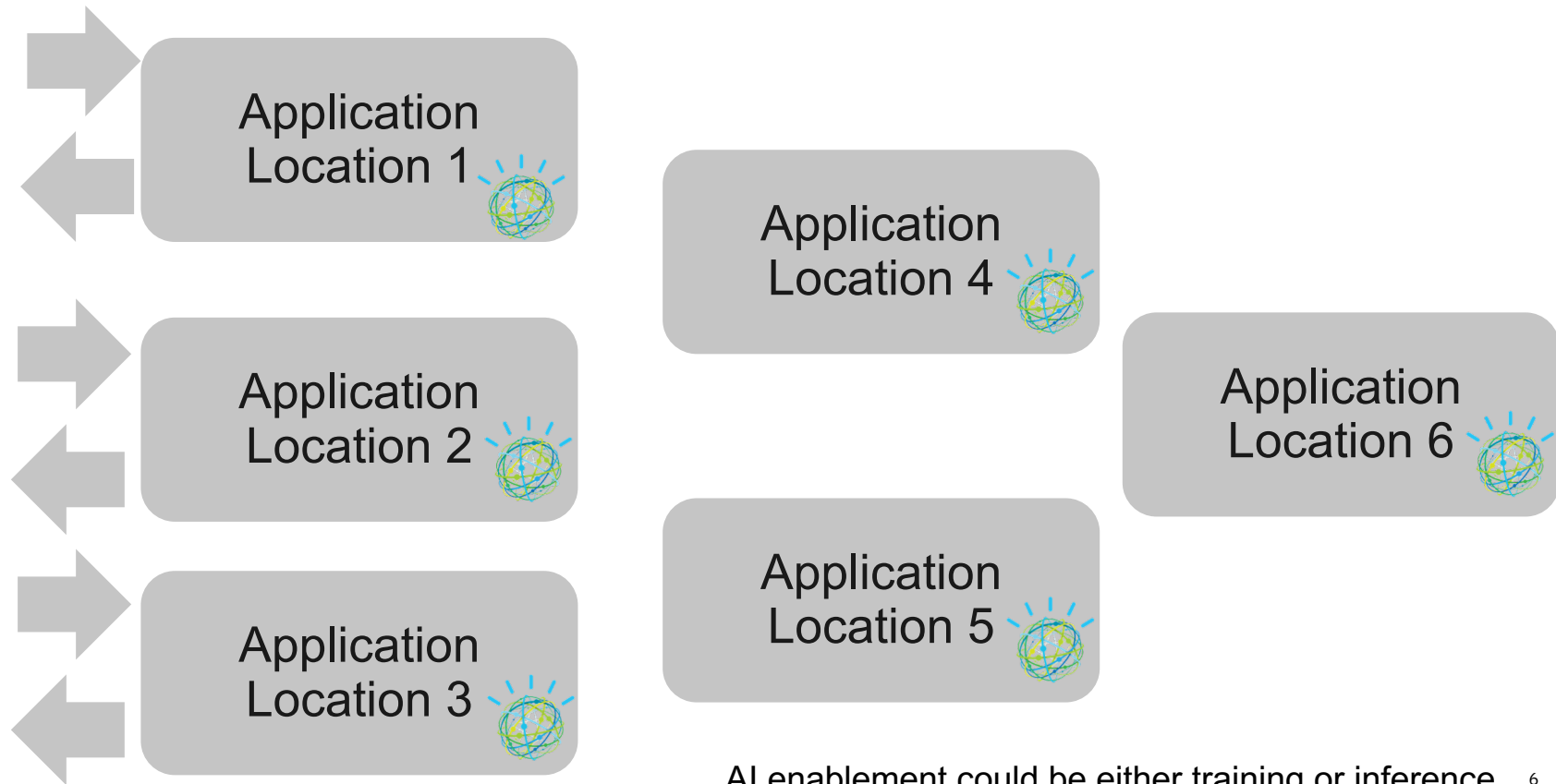
## Cheaper

- Use different skills, e.g. data scientists instead of programmers
- Uses less of an expensive resource (network/storage/processor)

## Better

- More flexible and adaptable
- Faster time to decision making
- More maintainable

# Distributed AI – application where AI enablement happens at many different locations

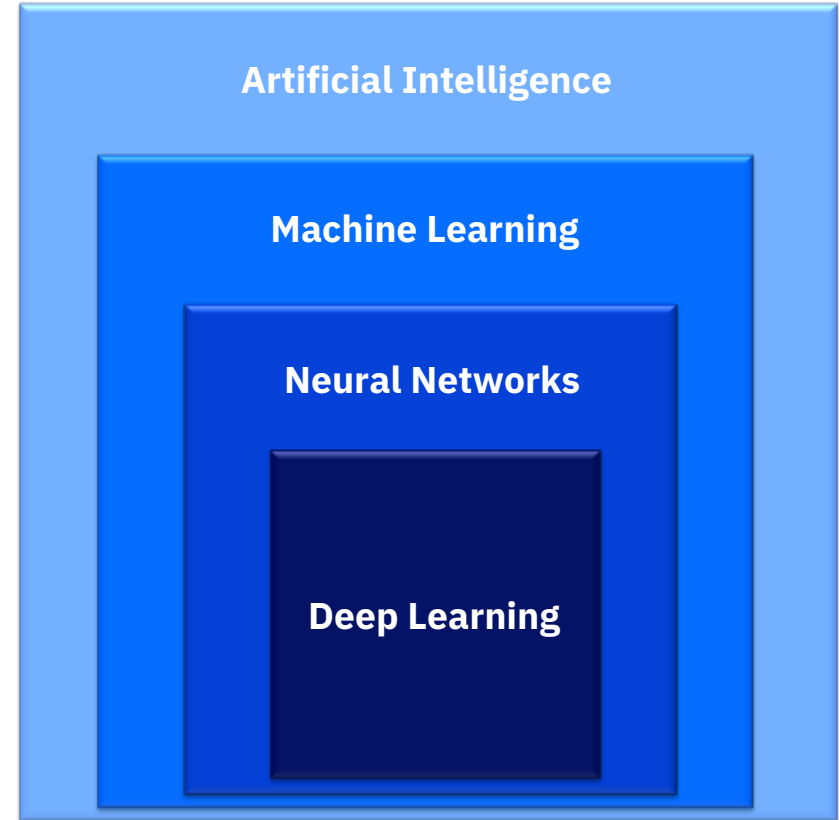
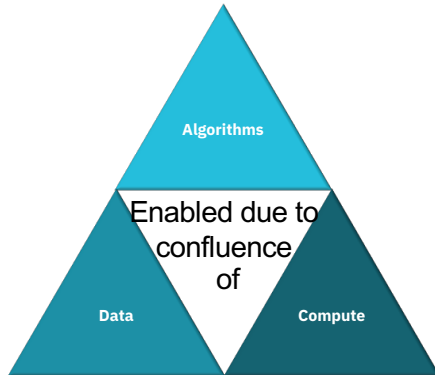


AI enablement could be either training or inference 6

# Creating an AI model

AI model can be created in a variety of ways

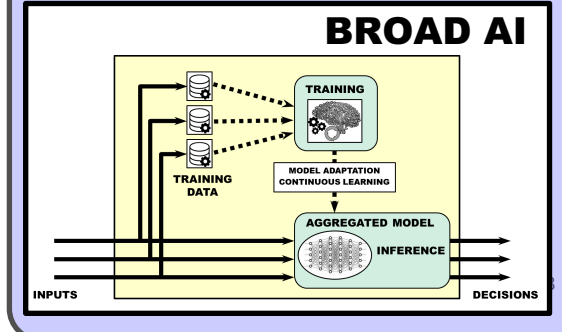
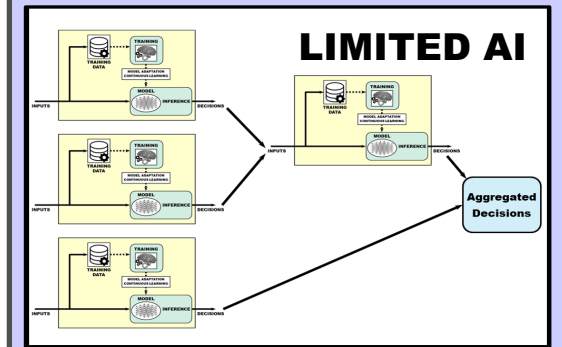
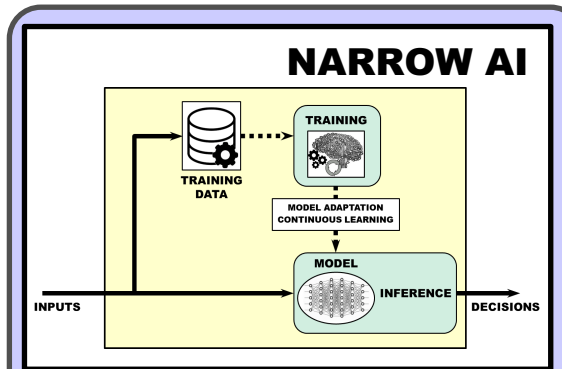
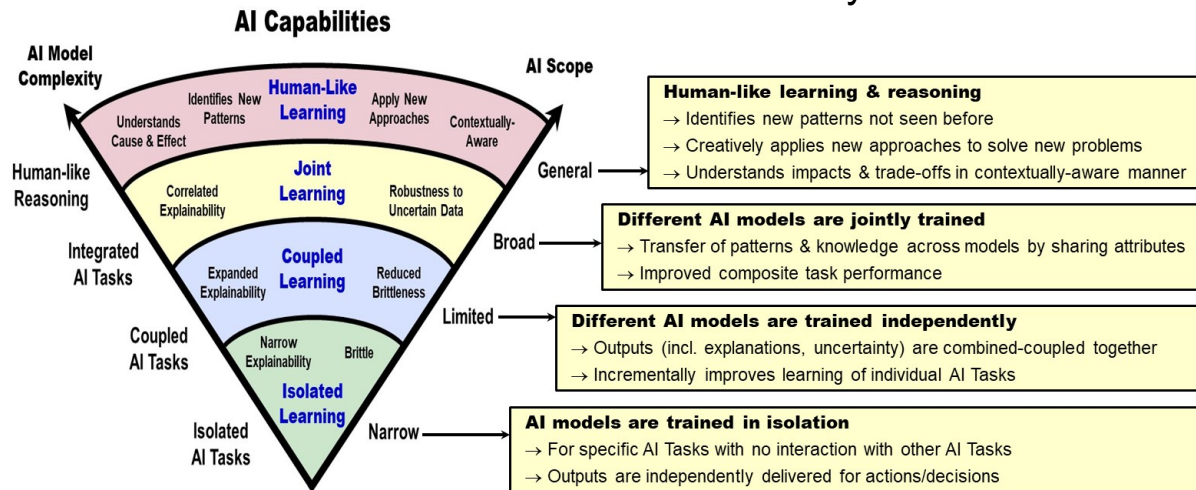
- Symbolic AI: A human being provides the model
- Machine Learning: A machine extracts the model from training data
- Neural Networks: A specific type of AI model
- Deep Learning: AI model using neural networks with many layers



# AI model encodes a spectrum of capabilities

The capabilities inherent in AI are related to:

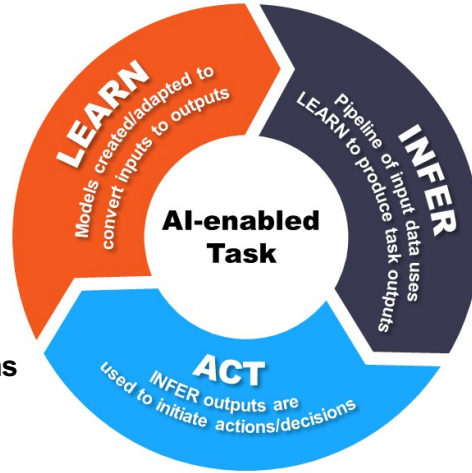
- Scope of the learning
- Complexity of the AI learning model
- Potential level of explainability
- Potential for robustness under uncertainty



# The process for creating an AI-enabled application

## Systematic approach to create AI:

1. Design AI/ML to meet goals
2. Obtain/curate relevant training/validation data
3. Train models under a variety of conditions  
→ dinky, dirty, dynamic, deceptive
4. Validate AI/ML in realistic conditions
5. Analyze to understand & predict its behavior/performance for safety
6. Determine allowable use cases including learning during deployment & autonomy

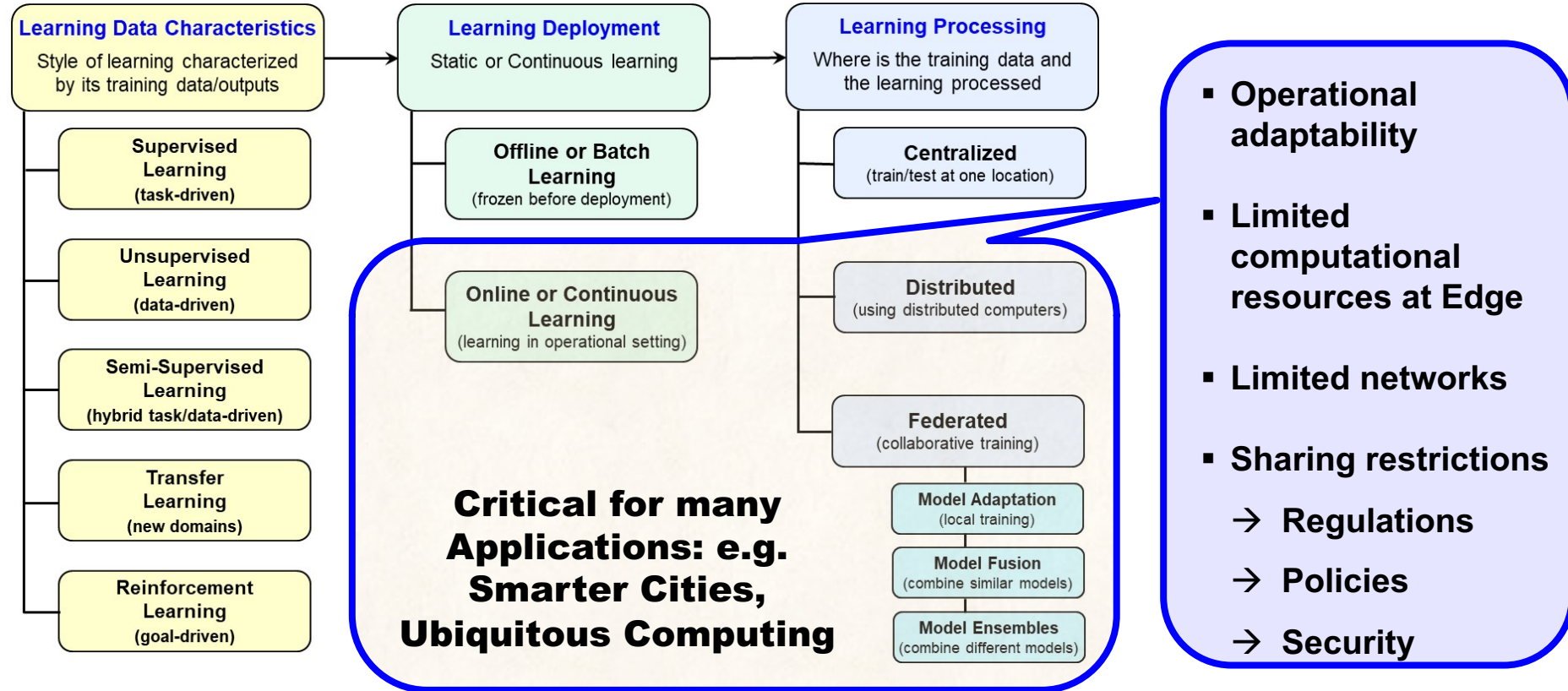


## Design Inputs for

### Learn → Infer → Act Processes

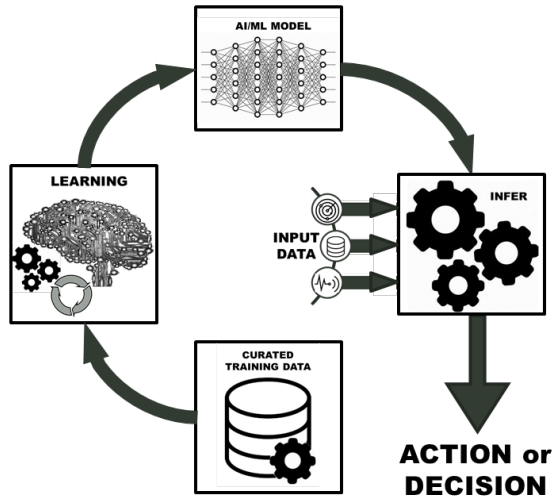
1. Where does learning, inferring, acting occur? (not necessarily co-located)
2. What are the performance & resiliency requirements of the action?
3. What level of autonomy is required?
4. What constraints will exist in the operational setting?
5. What is the availability & location of the training & input data?
6. How large & complex is the model to be used?
7. Where is computational capacity to support learning & inferencing?

# Factors involved in enabling an AI-based application



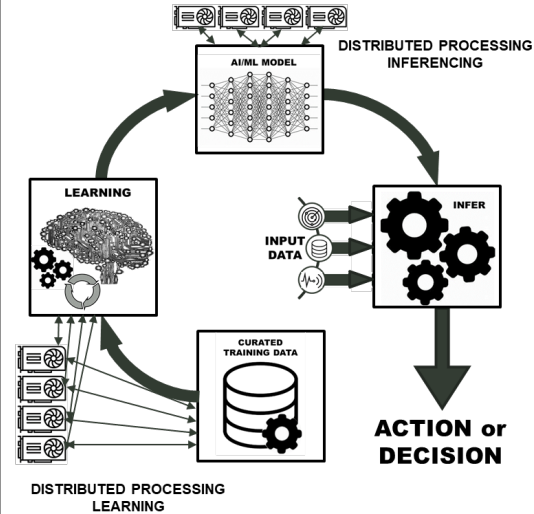
# Different options for Learning Processing

## Centralized Machine Learning



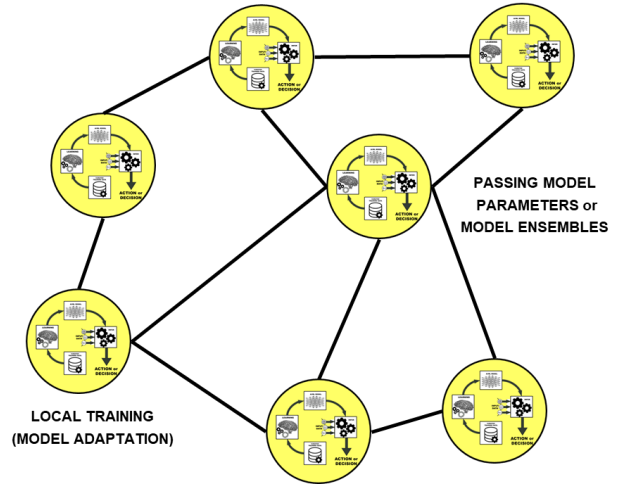
**When the Network is Always Reliably Available**

## Distributed Machine Learning



**When the Network is Reliable but a single processor is not adequate**

## Federated Machine Learning

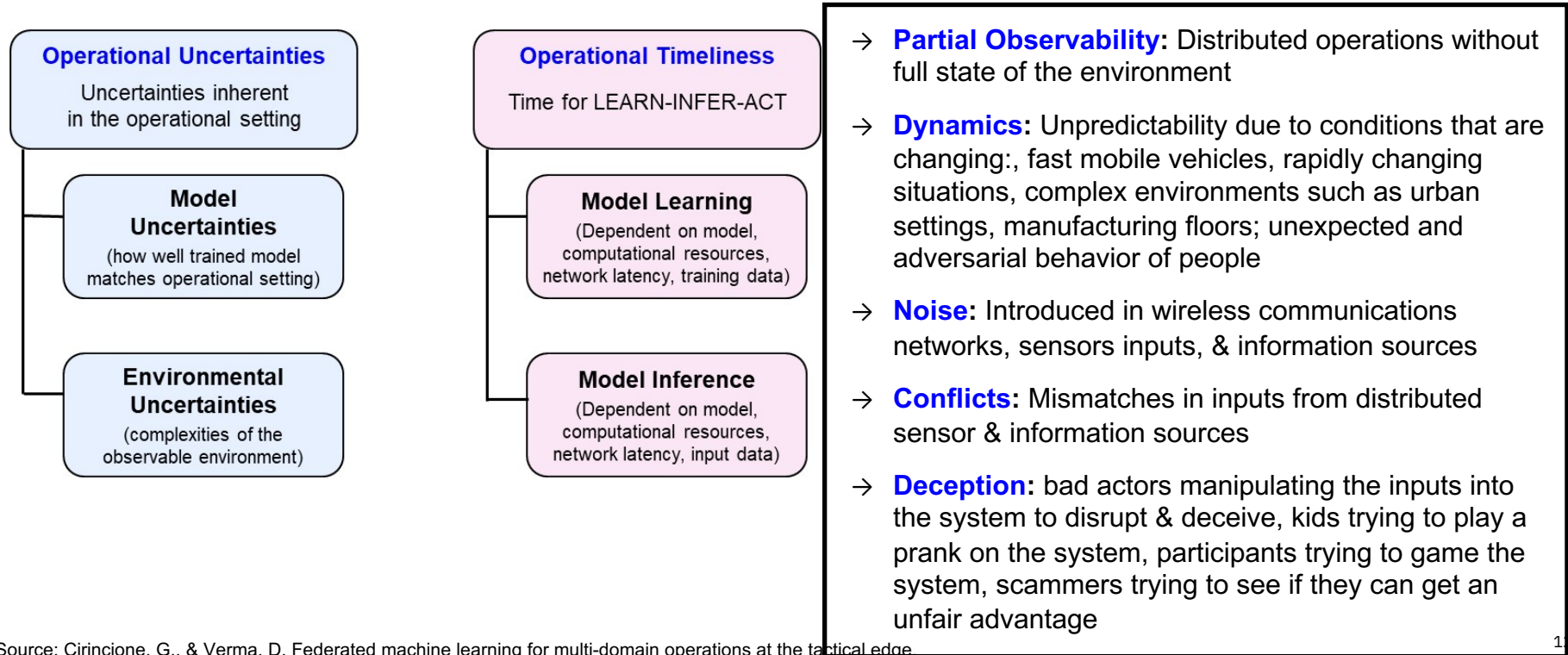


**When the Network is unavailable or crosses a regulatory domain**



# Any resilient application has to deal with the inherent uncertainties in the environment

## Operational factors impacting AI/ML effectiveness

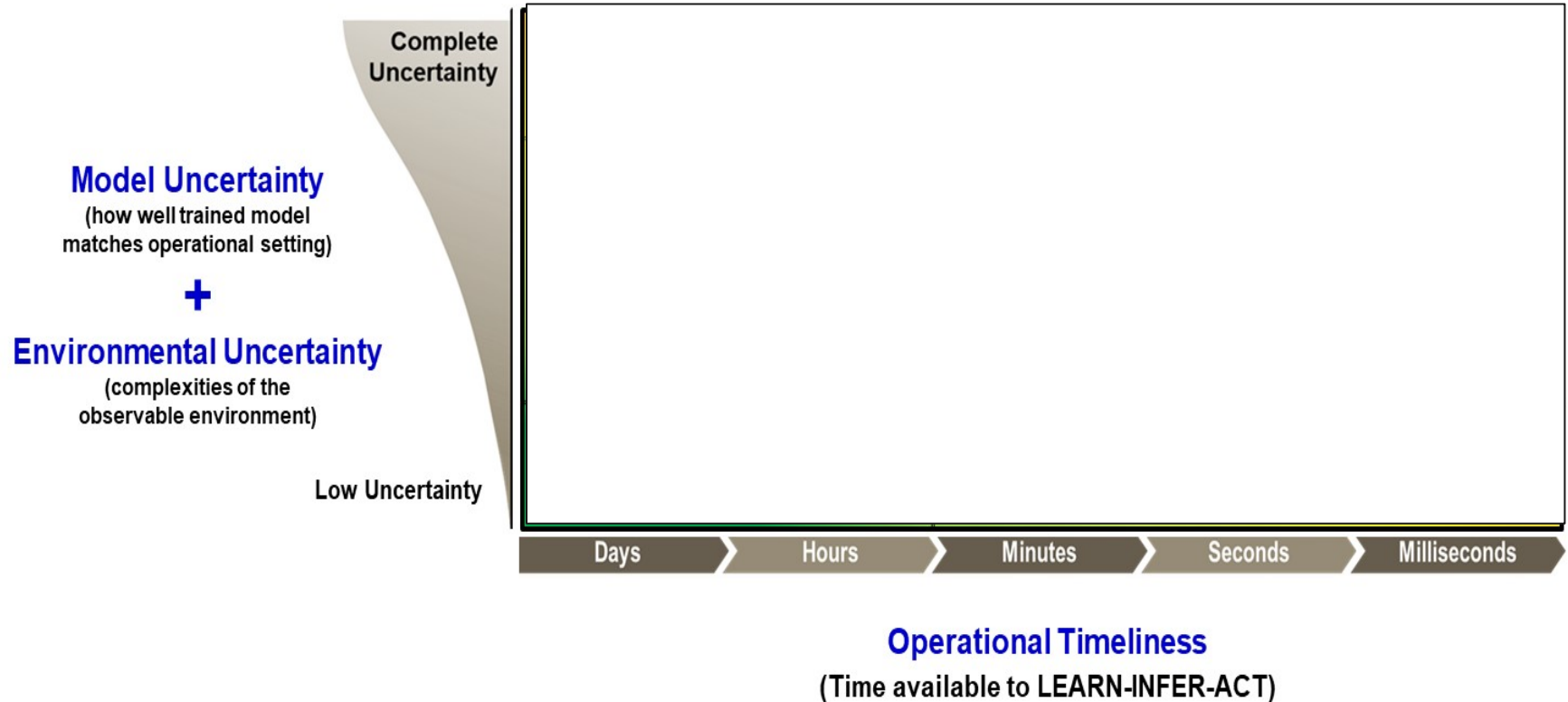


# Applications need to meet operational timeliness requirements

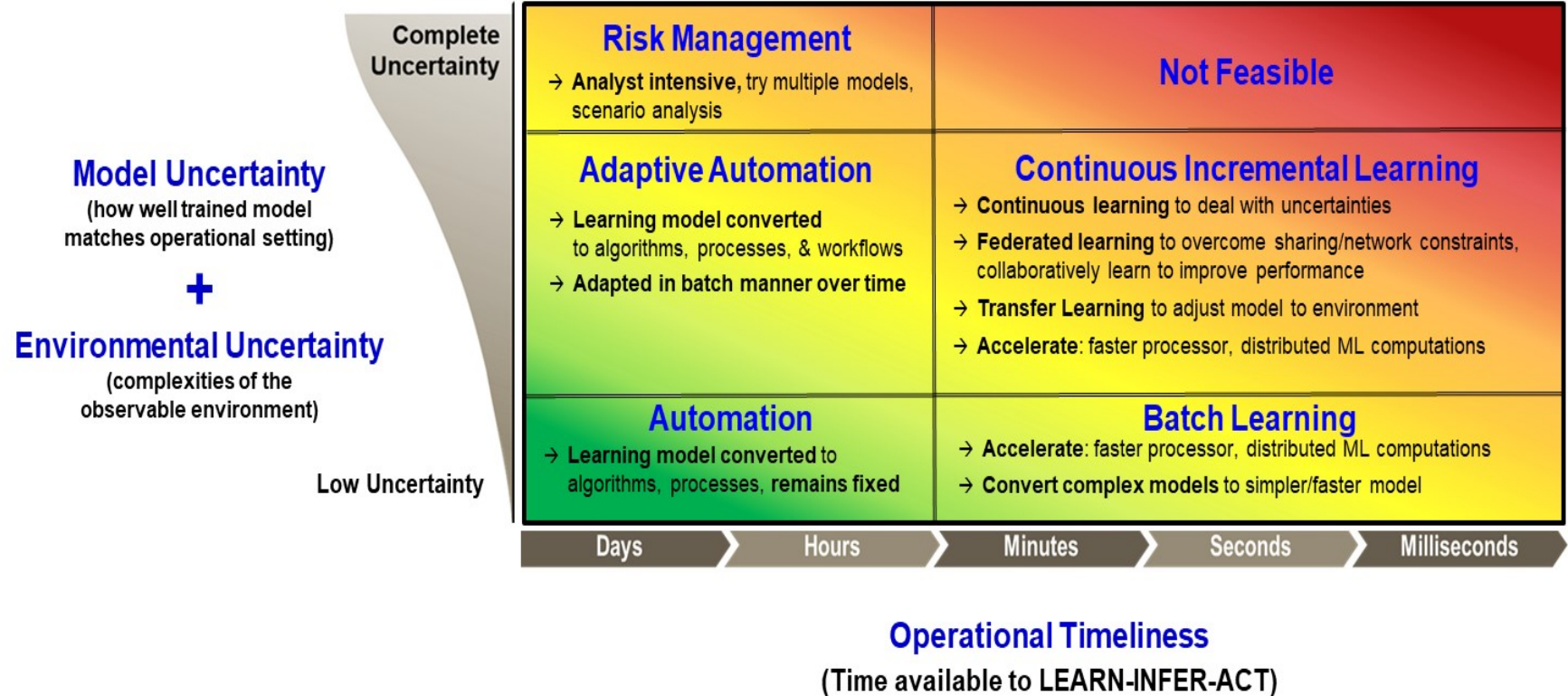
## Attributes of Operational Timeliness requirements in some smart computing environments

	Attributes	Operational Timeliness	Examples
Non Real-Time	<ul style="list-style-type: none"><li>• Complex Analytics from heterogenous and structured/unstructured sources</li><li>• Strategic decision making</li></ul>	Hours to Days	<ul style="list-style-type: none"><li>• Supply Chain</li><li>• Asset Maintenance</li><li>• Business Process Optimization</li><li>• Customer Journey Analysis</li><li>• Customer Retention Analysis</li></ul>
Near Real-Time	<ul style="list-style-type: none"><li>• Speed is important, but some delays are acceptable</li><li>• Quick response with soft or hard deadlines</li></ul>	Seconds to Minutes	<ul style="list-style-type: none"><li>• Voice Assistants</li><li>• Safety of Crowds and Cities</li><li>• Premises Access with Biometrics</li><li>• Industrial Worker Safety</li><li>• Bank Fraud Detection</li></ul>
Real-Time	<ul style="list-style-type: none"><li>• Constant input with a steady data output requirement</li><li>• Frequently hard deadlines</li></ul>	Seconds to milliseconds	<ul style="list-style-type: none"><li>• Manufacturing Quality Control</li><li>• Autonomous Vehicle Control</li><li>• Production Robotic Control</li><li>• Industrial Safety and Controls</li></ul>

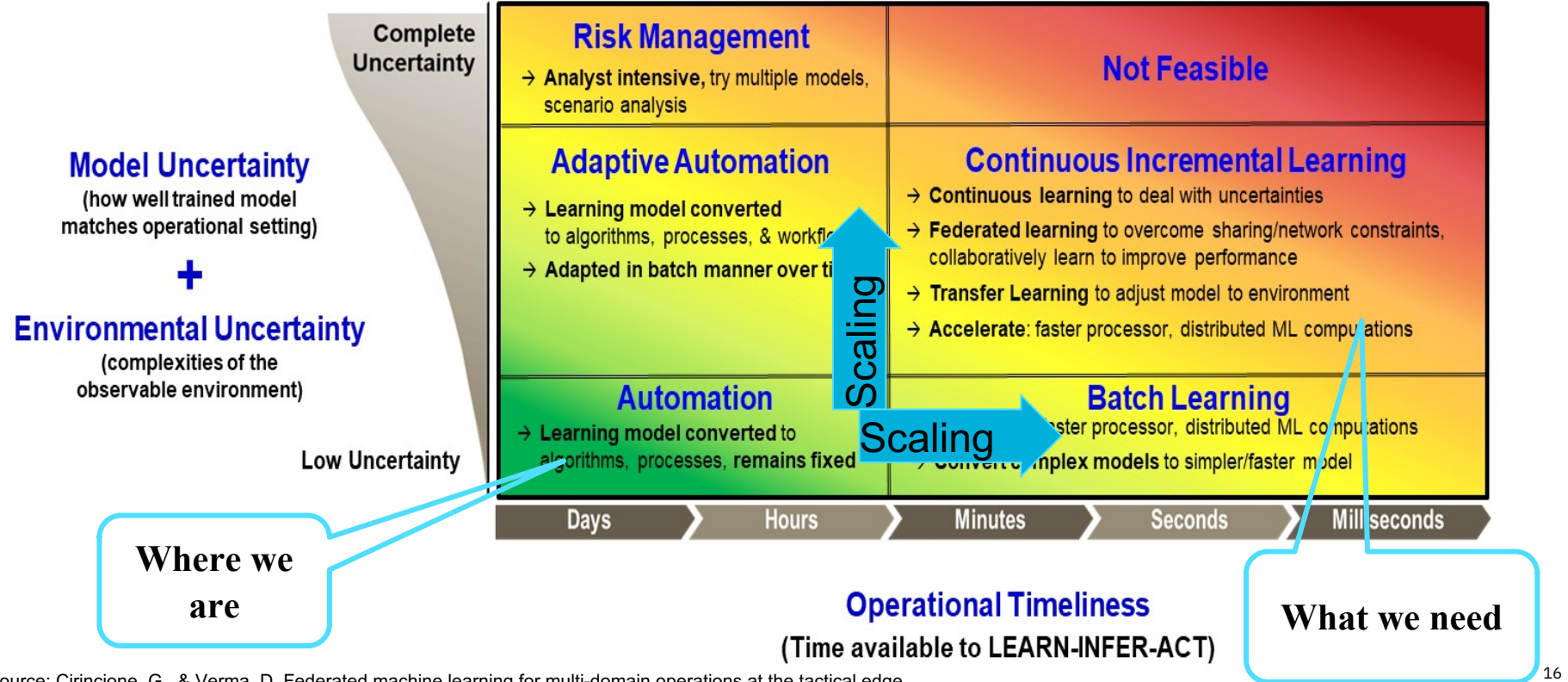
# Any Solution addresses a point in uncertainty-timeliness space



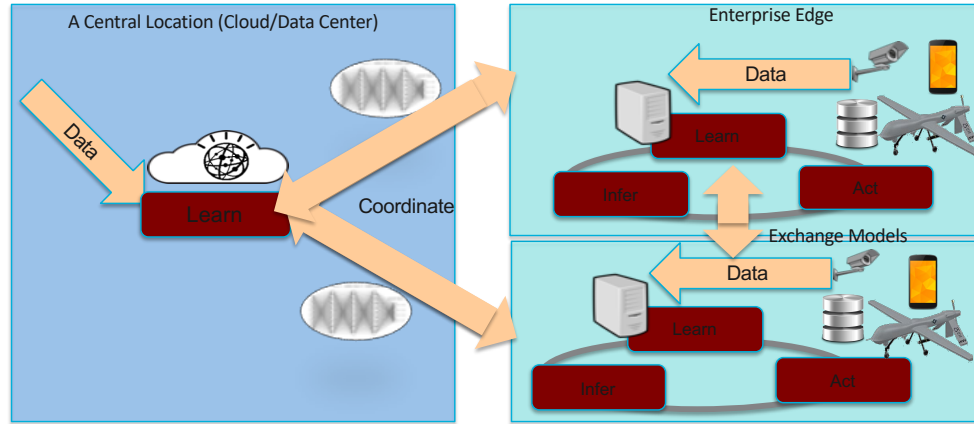
# Different types of approaches required at different regions



# Distributed AI is required to address the operational timeliness requirements

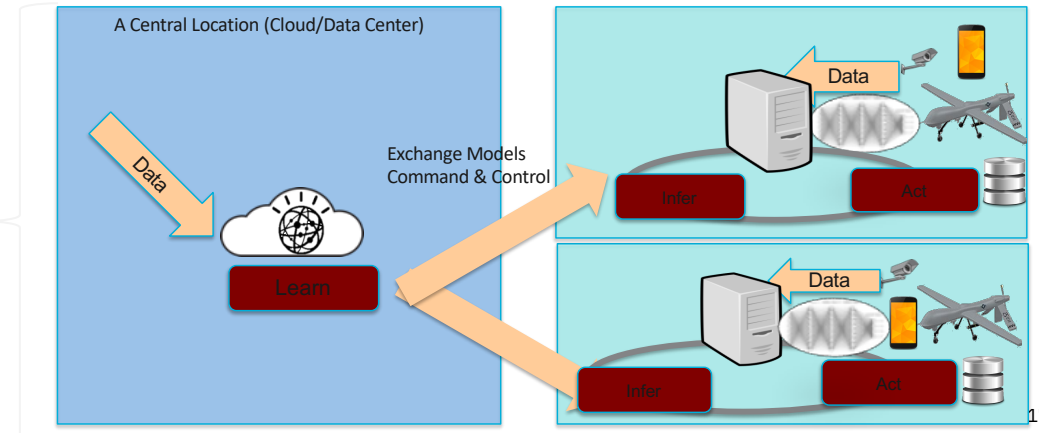


# Two Approaches for Distributed AI during model training phase

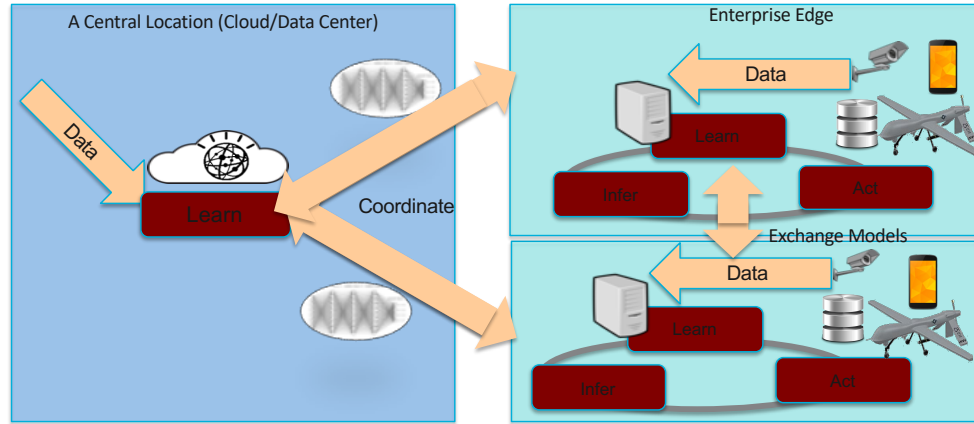


- Each edge site learns the AI model from their local data
- Models are exchanged and combined into a composite.

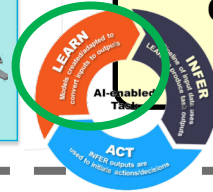
- Cloud learns the model and sends to each edge site
- Each edge site adapts the model to match its environment.



# Two Approaches for Distributed AI during model training phase



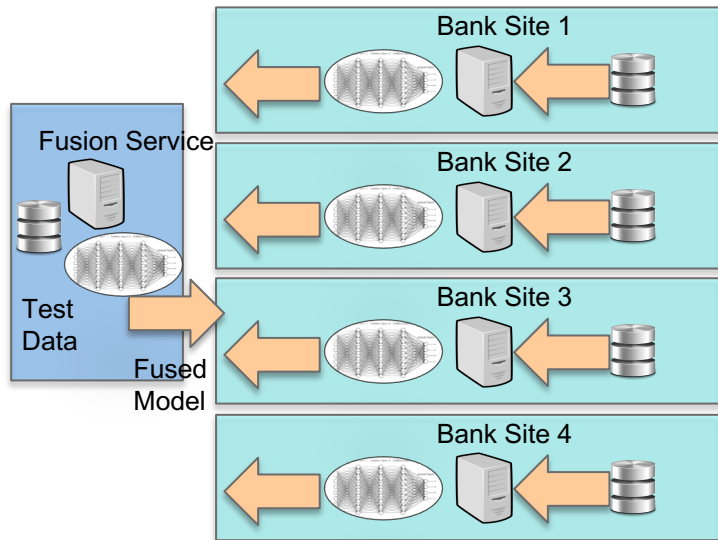
- Each edge site learns the AI model from their local data
- Models are exchanged and combined into a composite.



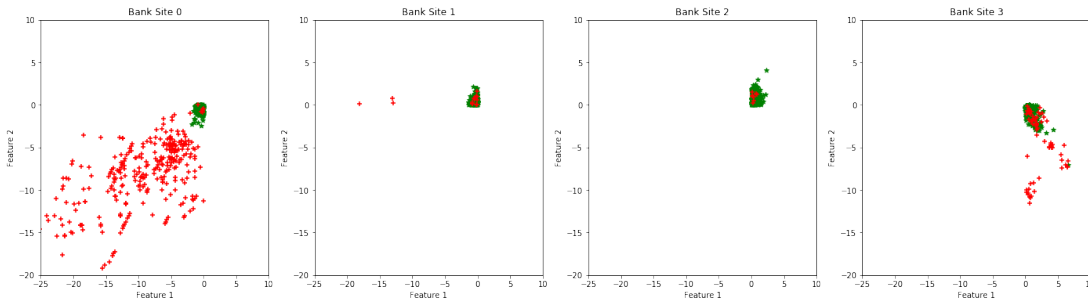


# Improving credit card fraud models using Federated Learning

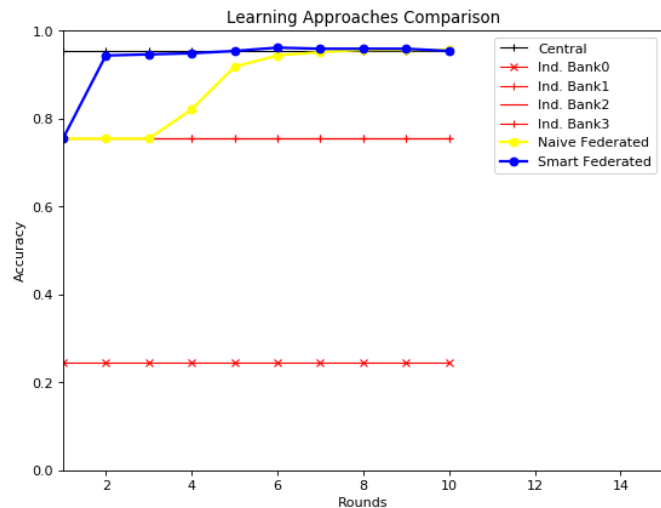
## Scenario



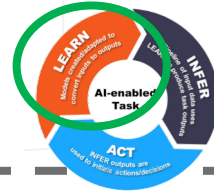
## Data



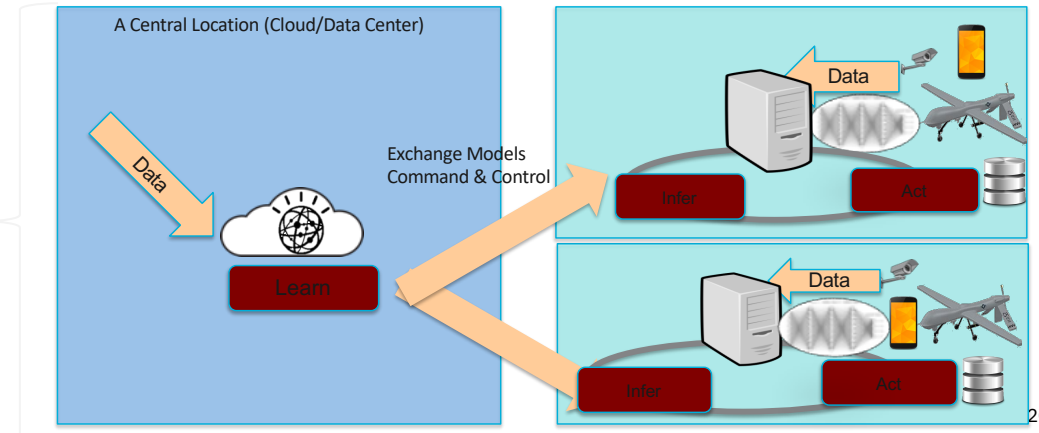
## Results



# Two Approaches for Distributed AI during model training phase

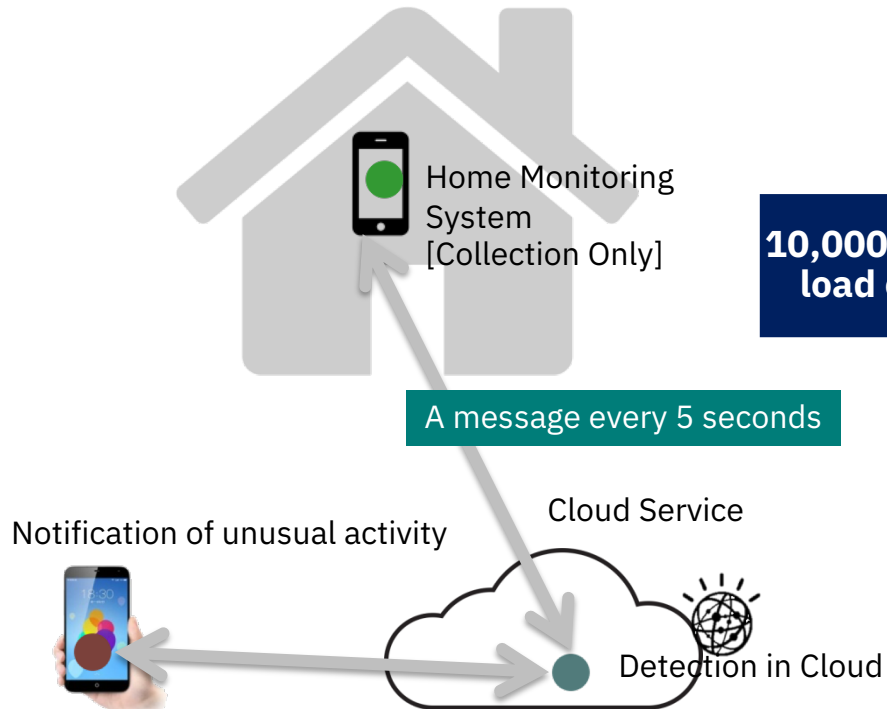


- Cloud learns the model and sends to each edge site
- Each edge site adapts the model to match its environment.



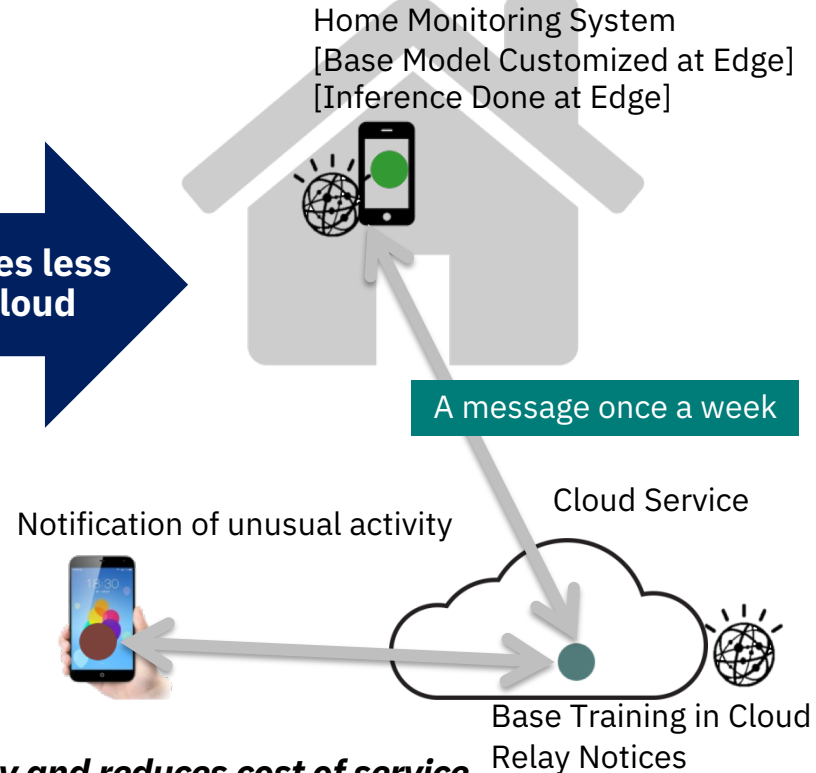
# A Specific Example: Comparing two Smart Computing Implementations

## Cloud-based AI



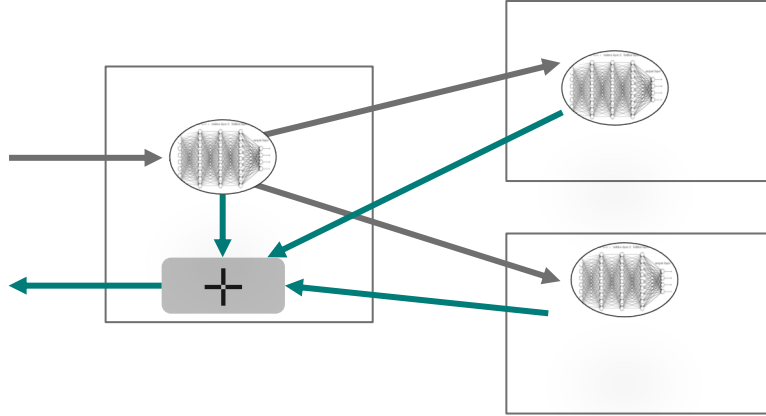
**10,000 times less  
load on Cloud**

## Edge Inference

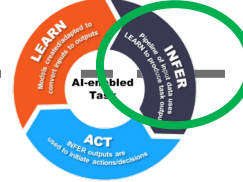


***Distributed AI improves Scalability and reduces cost of service***

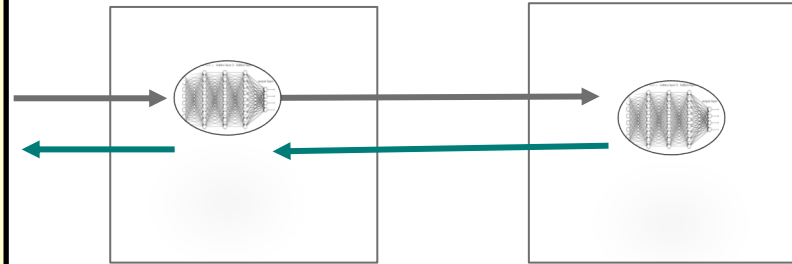
# Two Approaches for Distributed AI during the infer/act phase



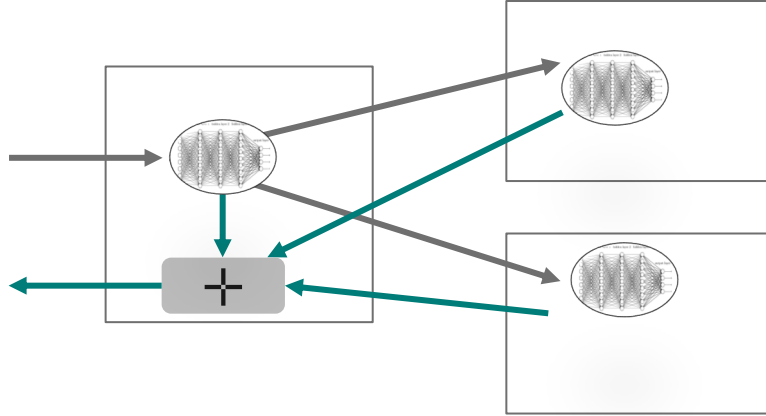
- Each site sends the inference request to one or more peer sites
- Results are combined into a aggregate answer by requesting site.



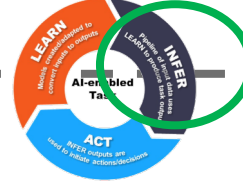
- Edge Site gets the model from Cloud Site and adapts it locally
- It checks its ability to perform the inference and if not, consults the cloud site



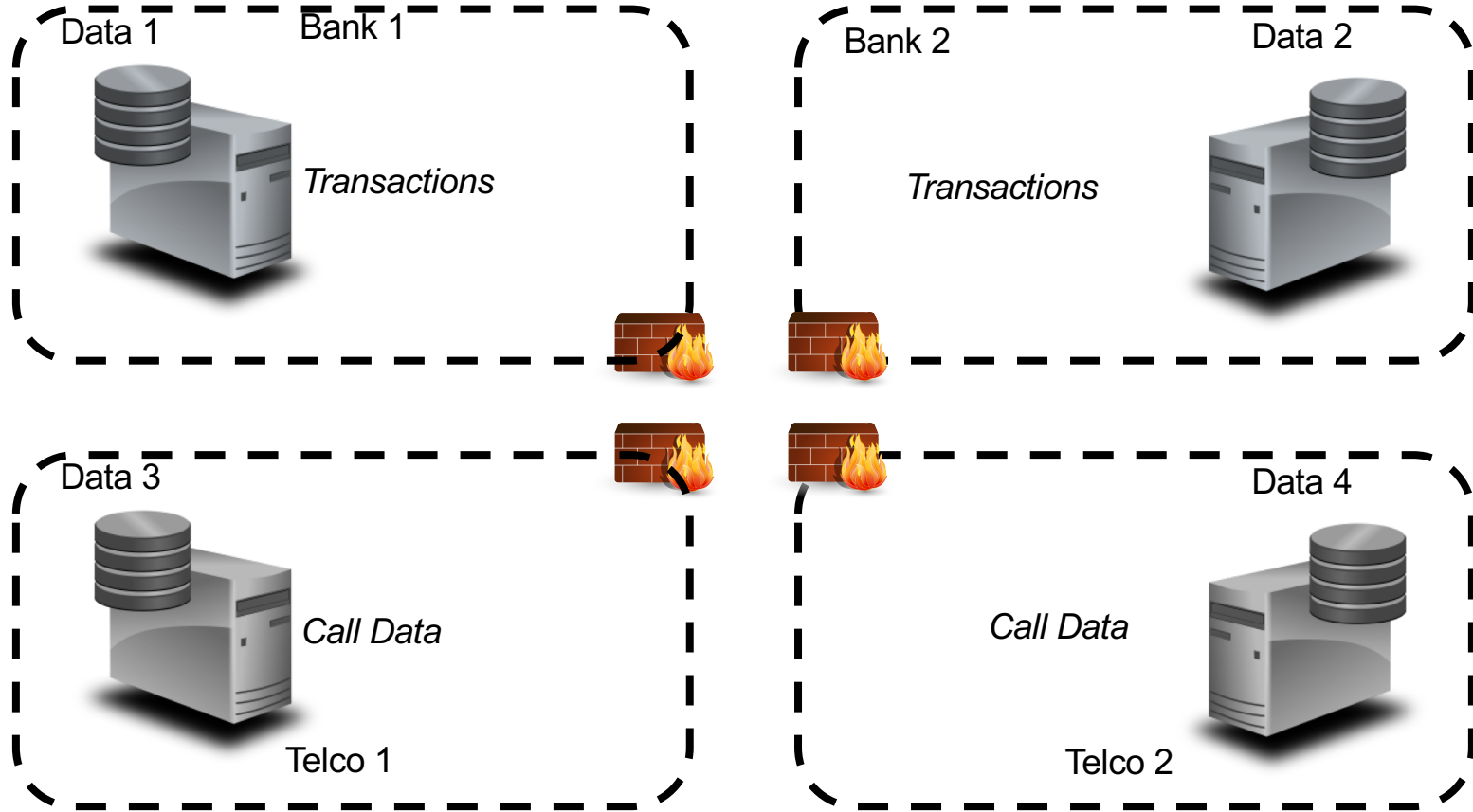
# Two Approaches for Distributed AI during the infer/act phase



- Each site sends the inference request to one or more peer sites
- Results are combined into a aggregate answer by requesting site.



# Example: The World Environment



2 Banks and 2 Telcos want to detect scammers

# Typical Results

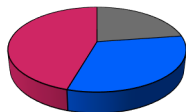
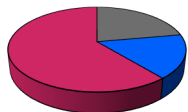
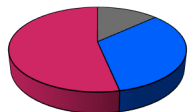
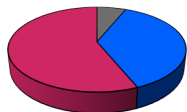
■ True Positives ■ False Positives ■ False Negatives

Bank 0

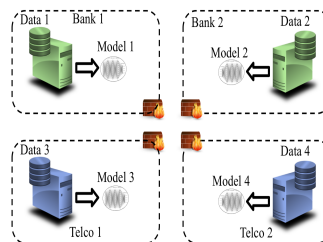
Bank 1

Telco 0

Telco 1



**Independent Detection**



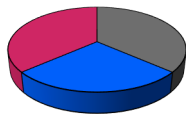
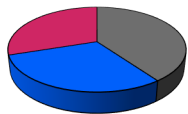
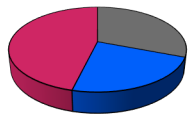
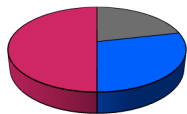
Each organization  
Works independently  
to detect scammers

Bank 0

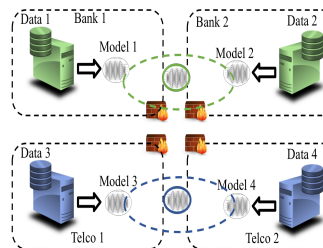
Bank 1

Telco 0

Telco 1



**Federated Learning**



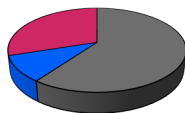
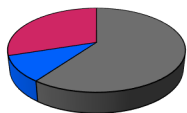
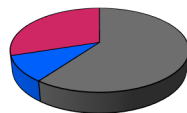
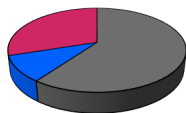
Banks and Telcos  
use Federated Learning  
to detect scammers

Bank 0

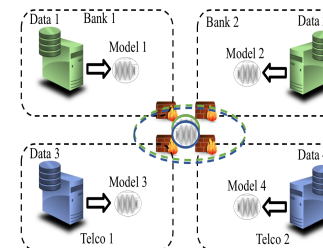
Bank 1

Telco 0

Telco 1



**Federated Learning + Federated Inference**

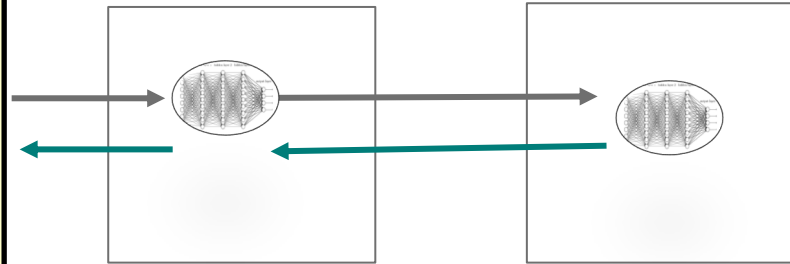
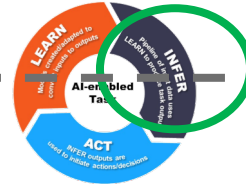


Banks and Telcos  
use Federated Learning  
and Federated Inference  
to detect scammers



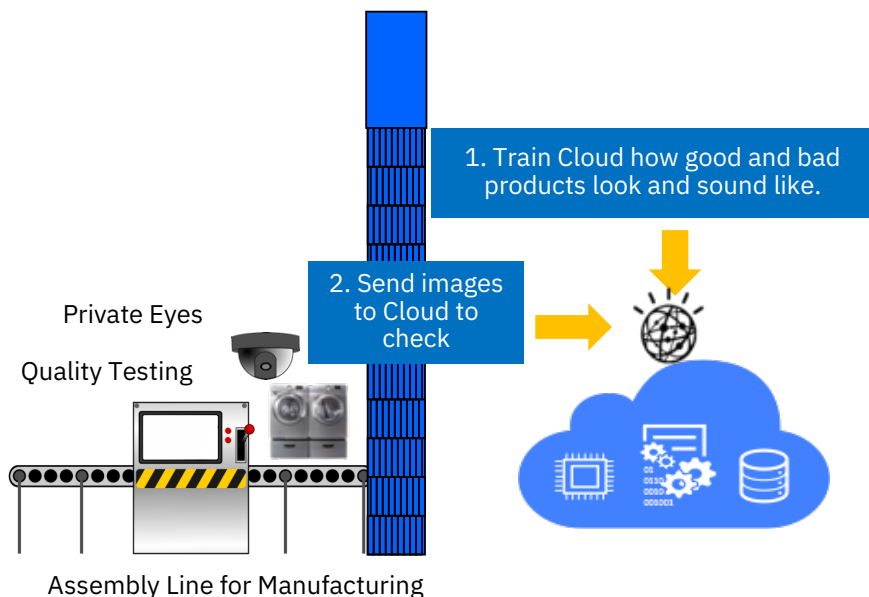
# Two Approaches for Distributed AI during the infer/act phase

- Edge Site gets the model from Cloud Site and adapts it locally
- It checks its ability to perform the inference and if not, consults the cloud site



# A Specific Example: Meeting Requirements of Manufacturing

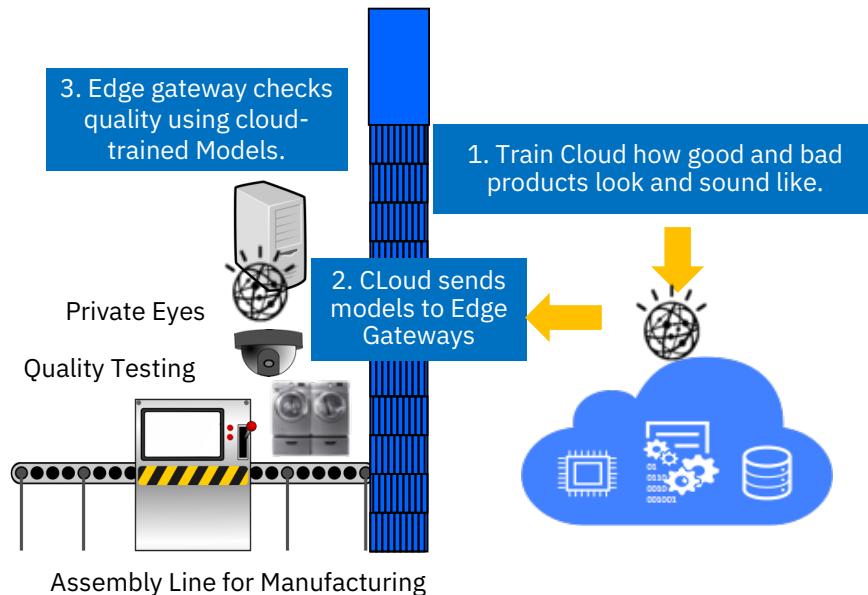
## Cloud-based AI



Time per check = 1000 ms:

- Network ( $10 \times 100\text{ms}$ ) + Image Time (2ms) + Inference Time (1ms)
- Operational data goes to the cloud

## Edge Inference



Time per check = 50 ms:

- Network ( $10 \times 5\text{ms}$ ) + Image Time (2ms) + Inference Time (2ms)
- Plus No Operational data leaves premises

***Distributed AI improves operational timeliness and privacy risks***