



JPSS Cloud Data Processing Future and Machine Learning for the Weather Value Chain

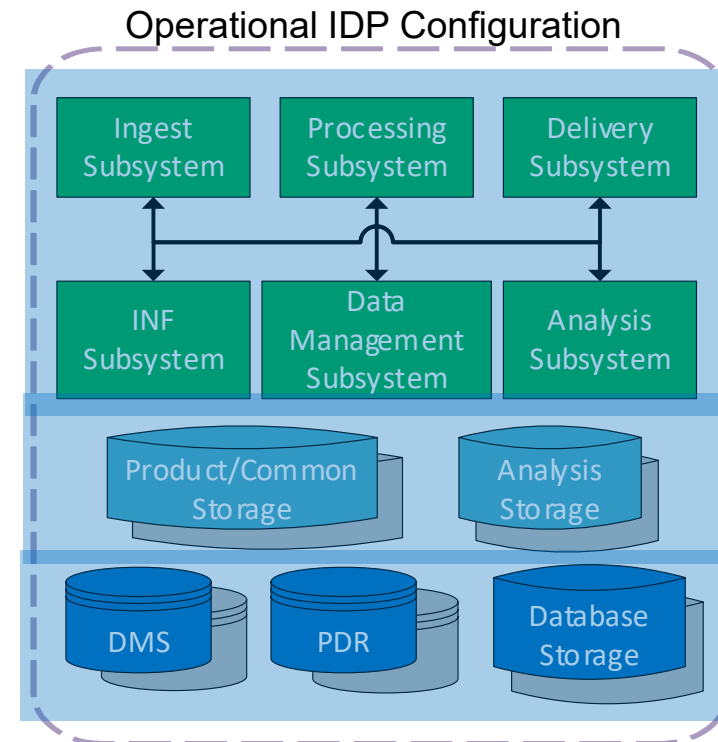
Scott Kern
Emily Greene
Shawn Miller

Agenda

- IDPS in the Cloud for NESDIS
- Optimization Plans
- Cloud enabled Opportunities
- Machine Learning for Operational Weather

Initial Implementation – Phase 1

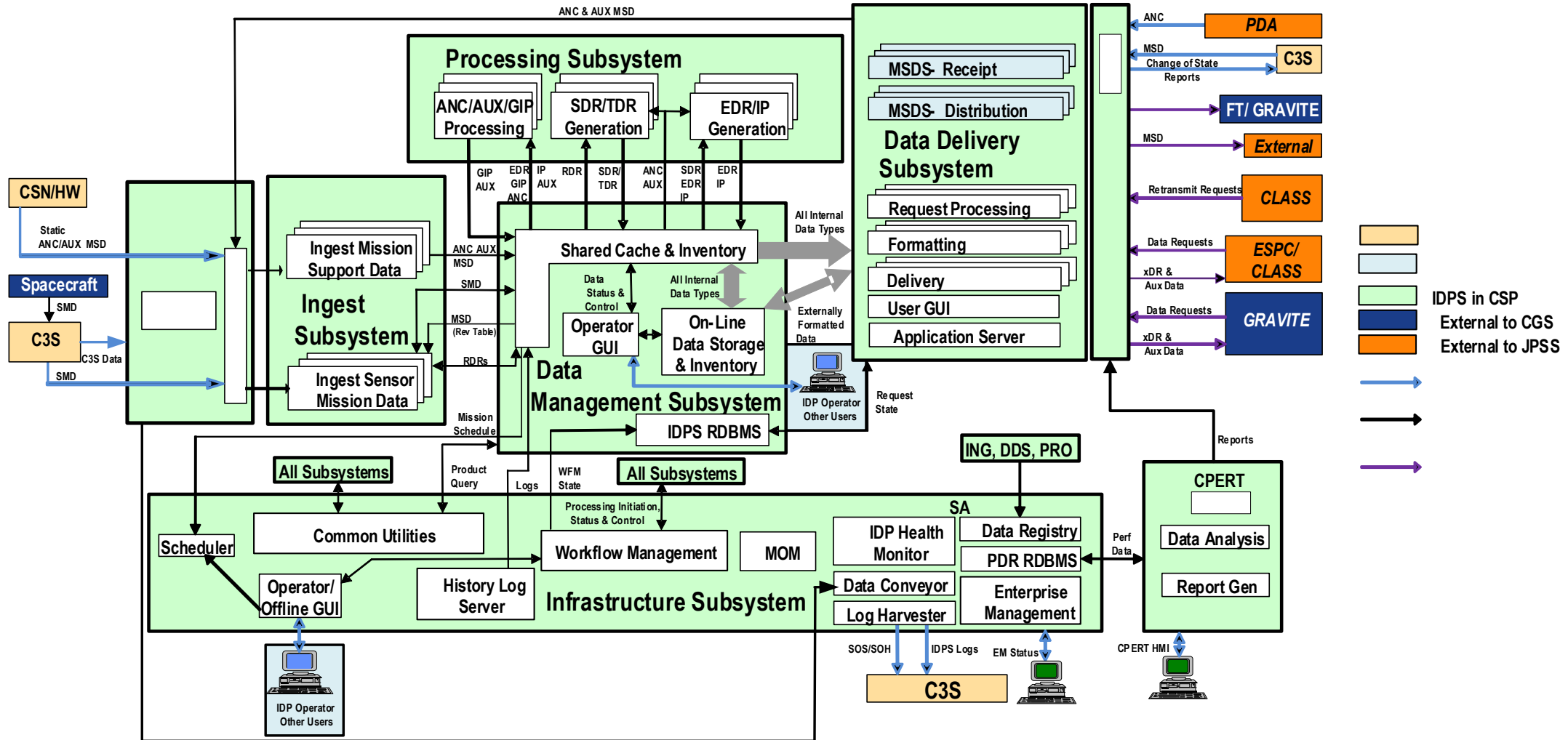
- Transition to Operations in Cloud must occur NLT EOY 2020 (Lenovo HW waiver expiration)
- NOAA direction to migrate current operational baseline to Cloud with minimal baseline changes
 - Only changes to baseline that are explicitly necessary to operate in the cloud
 - Moving primary IDP DB from Oracle to PostgreSQL to save Oracle licensing costs
- HOT backup of primary Operations IDP
 - Monthly Security Patching requires transition to backup IDP
 - 3rd IDP necessary to accommodate monthly patches and baseline upgrades while maintaining resiliency to failures
- Primary change is new Common Environment :
 - Route data to multiple IDPS systems from a single on-prem data source
 - Management of security solutions
 - On-prem IDPS is a “user” of security solutions from the C3S segment
- Leveraging DevOps Tools/Processes:
 - Environments 100% managed using Infrastructure-as-Code (Packer, Terraform, Chef)
 - Faster/Frequent algorithm releases to PRO subsystem decreases Research-to-OPS (R2O) cycle
- ~60 EC2 VMs and 500 TB EBS storage per Ops-capable IDPS



- Database Layer (EC2 and EBS)**
- Oracle Dataguard installed to EC2
 - Backup DB instance
 - EBS storage attached to EC2
 - DMS: Data Management
 - PDR: Performance Data Repo

IDPS in the Cloud Architecture Overview

CGS-018



AWS Services In Use For Initial Implementation

AWS Service	Purpose
EC2/EBS	<ul style="list-style-type: none"> • Processing VMs, significant tuning in Task Order prototypes to define the configuration • SIGNIFICANT volume of EBS storage required for GPFS, 250 TB writable data per IDP. GPFS installed to EC2 however needs full replication for resiliency to a single EC2 failure, 500TB total storage.
PostgreSQL Relational Database Service (RDS)	<ul style="list-style-type: none"> • Hosts primary data management database (DMS) • Deployed in multi-AZ configuration for redundancy, but other IDPS components are single-AZ
Simple Storage Service (S3)	<ul style="list-style-type: none"> • SMD and MSD archive hosted in cloud DPC • Factory use for testdata and other large datastores • Storage for Artifactory COTS in deployment pipeline • Drop-box for algorithm changes in DP-AE
Simple Notification Service (SNS) Simple Queue Service (SQS)	<ul style="list-style-type: none"> • Used by new Mission Data Distribution function to deliver one data source to many IDPs • New SMD/MSD product arrives in S3, SNS sends message to a SQS queue assigned to each IDP, guarantees each IDP receives all SMD/MSD even when not active
CloudWatch	<ul style="list-style-type: none"> • Monitoring and aggregation of cloud logs (from CloudTrail, VPC Flow Logs and CloudWatch agents) • Delivers security relevant events to on-prem Qradar
CloudTrail	<ul style="list-style-type: none"> • Logs API calls to AWS services
VPC Flow Logs	<ul style="list-style-type: none"> • Logs traffic in/out of each VPC
Direct Connect Gateway (DXG) Transit Gateway (TXG)	<ul style="list-style-type: none"> • DXG enables direct connect at NWAVE • TXG routes data from DXG to each VPC

Optimization – Phase 2

- Optimization Phase Updates the IDPS cloud design to take better advantage of cloud capabilities
- Provides significant cost savings over initial-implementation
 - Savings for Infrastructure, COTS, O&M
- Implements a better foundation for science/forecast product driven changes during Modernization Phase

Optimization	Description
Transition to Highly Available (HA) IDPS	<ul style="list-style-type: none"> • Deploy single HA IDPS spanning 2 Availability Zones <ul style="list-style-type: none"> • Subsystems deployed across AZs in auto-scaling groups • “Live” security patching on dynamic instances to eliminate OPS/Non-OPS transitions for monthly security patching
Dynamic Allocation of Processing Capacity	<ul style="list-style-type: none"> • Elastic processing capacity to dynamically respond to changing throughput needs in responding to anomalies
Complete migration of all databases to PDR	<ul style="list-style-type: none"> • COTS licensing savings • Reduces DBA support needs and security patching overhead
Modernize IDPS Storage Layer	<ul style="list-style-type: none"> • Product storage moved from GPFS to cloud-native blob storage (AWS S3) <ul style="list-style-type: none"> • Significant cost savings • Initial prototyping shows satisfactory performance with minimal code modifications • Common storage migrates to cloud-native shared file system (AWS Elastic File Service EFS) <ul style="list-style-type: none"> • Provides HA without overhead required to manage large replicated storage cluster
Utilize Clustered Messaging Service	<ul style="list-style-type: none"> • Develop HA messaging system or utilize “Messaging-as-a-Service from AWS (Amazon MQ)
Utilize Cloud-Native Monitoring and Alerting	<ul style="list-style-type: none"> • Initial-Implementation using legacy design of monitoring agents deployed on IDPS VMs delivering messages to operations.

BLUE: Denotes successfully prototyped and demonstrated capability

Modernization – Phase 3

- The modernization phase leverages IDPS’ proven data production platform
 - Provide an expanded number of enterprise data products
 - Decrease algorithm process overhead accelerating R2O cycle
- Data Delivery capability to expanded user base while minimizing data egress costs
 - Prioritize Real-time products critical to NWP delivered with IDPS’ proven low-latency and stability
 - Non-Real-time critical products have packaging and delivery processing

Optimization	Description
Modernize Processing Subsystem using Containerized Algorithms	<ul style="list-style-type: none"> • Science teams will directly develop algorithms and include dependencies in versioned containers • Run multiple algorithm versions in parallel, dependencies reside in container • Enterprise data product generation • Real-time Processing: Operational algorithms generating products • Off-line Processing: “Algorithm Sandbox” Evaluate updates to algorithms <ul style="list-style-type: none"> • Executed during “back-orbits”, spot-instances or serverless • Eliminates need for full IDPS dedicated for dedicated I&T and provides faster R2O cycles
Modernize Data Delivery via Cloud-based Content Delivery Network	<ul style="list-style-type: none"> • Data products delivered to single cloud location (S3) <ul style="list-style-type: none"> • Eliminate delivery of products through C3S facility to Mission Partners • Real-Time Delivery: Products delivered to S3 location <ul style="list-style-type: none"> • NWP products delivered in directly ingestible format (HDF, BUFR, NetCDR, etc) • Consumers who need real-time products will receive notification of new products and API to pull the data directly down to their system (S3 => SNS => SQS pipeline) • Off-Line Delivery: <ul style="list-style-type: none"> • Non-Real-Time consumers will be able to request aggregation and/or packaging of products which will create a new product in S3 and notification delivered to consumer
“Lights Out” IDPS decreases reliance on dedicated operations staff	<ul style="list-style-type: none"> • IDPS is highly stable system requiring almost no human interaction to function <ul style="list-style-type: none"> • Decreases reliance on 24x7 dedicated operators • Remove Java based GUIs and replace with simplified web GUI with APIs to drive IDP functions <ul style="list-style-type: none"> • Significantly improves security posture

Cloud Enabled Opportunities

- Data Delivery via Cloud Storage
 - Data products delivered to single cloud location (S3)
 - Eliminate delivery of products through C3S facility to Mission Partners
 - Real-Time Delivery: Products delivered to S3 location with no aggregation/packaging
 - Significant simplification of delivery subsystem
 - Consumers who need real-time products will receive notification of new products and API to pull the data directly down to their system
 - S3 => SNS => SQS pipeline
 - Off-Line Delivery:
 - Non-Real-Time consumers will be able to request aggregation and/or packaging of products which will create a new product in S3 and notification delivered to consumer

Cloud Enabled Opportunities

- Containerize Processing Algorithms

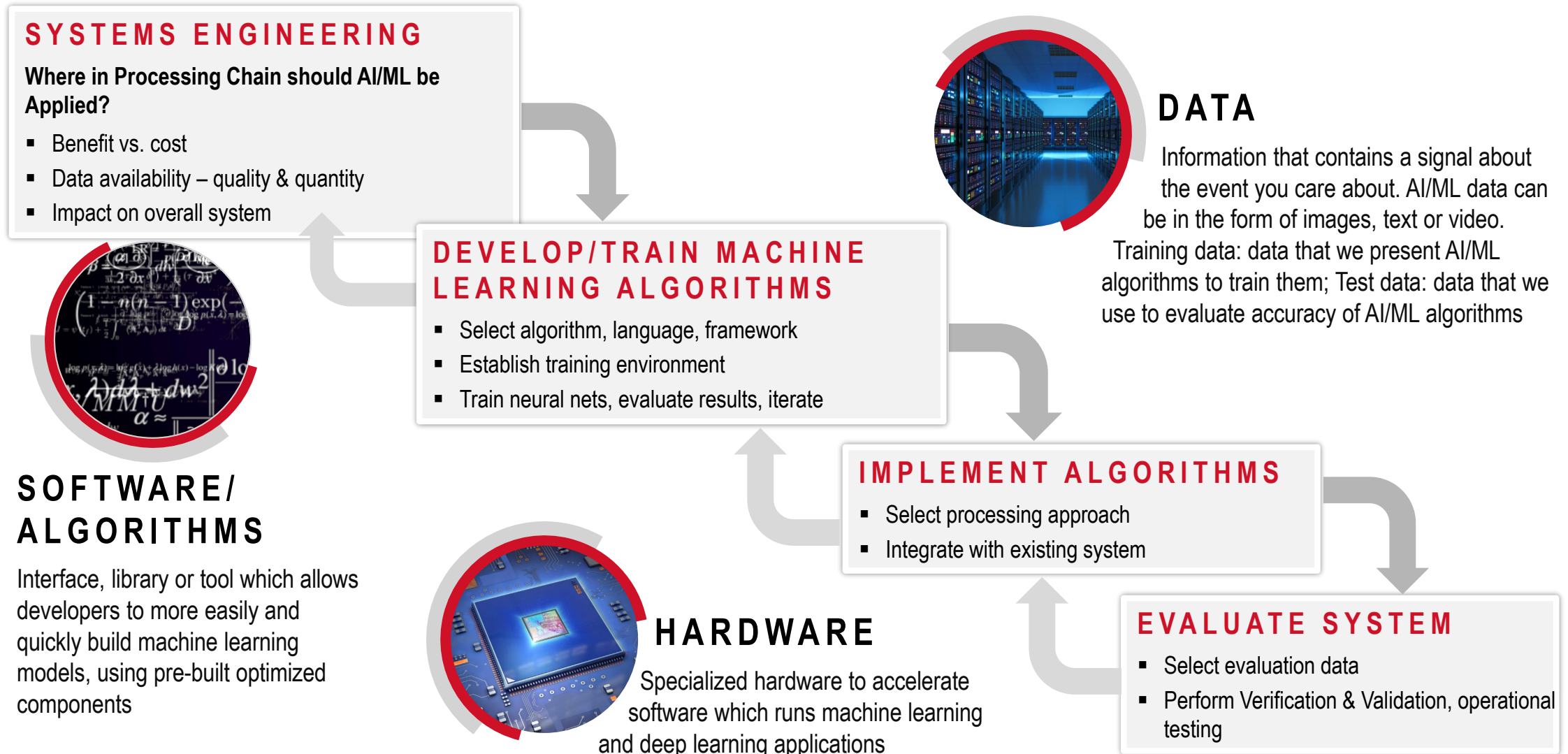
- Science teams will directly develop algorithms and include dependencies in versioned containers
- Run multiple algorithm versions in parallel, dependencies reside in container
- Real-time Processing: Operational algorithms generating products
- Off-line Processing: Evaluate updates to algorithms, executed during “back-orbits” or serverless
 - Less capacity required than real-time
 - Eliminates need for full IDPS dedicated for dedicated I&T and provides faster Science-to-OPS cycles

- Rapid Science Container updates

- 4 weeks cycles to reach Operations, 2 cycles in development at all times
- SW with 2 week sprints
- Extra time for Performance check and Science Quality evaluation

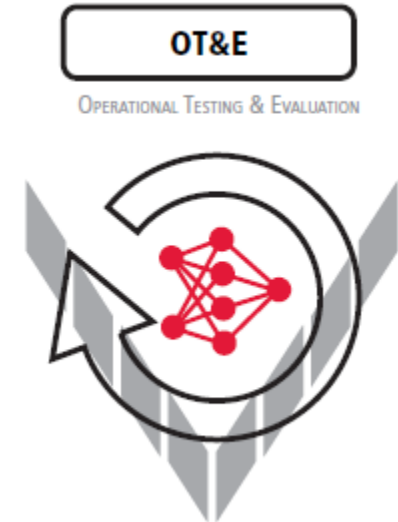
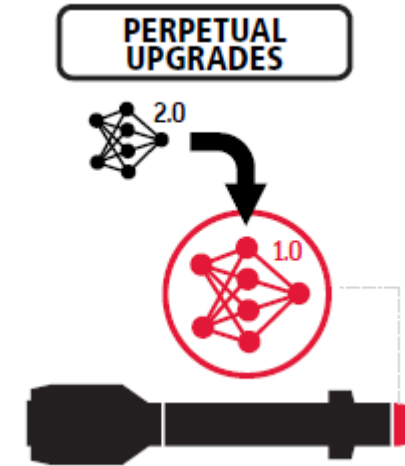
Week 1					Week 2					Week 3					Week 4					
M	T	W	R	F	M	T	W	R	F	M	T	W	R	F	M	T	W	R	F	
SW										Build	Int BUCO & SOL Delpoy		Regression and Perf.			Board Approv.		I&T Deploy & STAR co		Deploy to Ops
										Science Quality Check										

General Application of Machine Learning



Challenges to Machine Learning

- **Opacity (i.e., “the black box”):**
 - Explaining why ML got an answer is just as important as getting the answer
- **Perpetual Upgrades:**
 - Evolving requirements can cause models to have a life cycle of mere seconds; key discriminator is automation
- **Operational Test and Evaluation:**
 - ML models are inherently complex, non-deterministic systems; potential exists for unanticipated emergent behavior, indeterminate test results, Black Swan events (another fertile ground for innovative solutions)



Improving Machine Learning Itself

Strategies for Rapid Prototyping Machine Learning

- Difficulty in pattern recognition is the sheer volume of source data that must be analyzed – time required to acquire, label, and train the model
- We have developed a number of tools and approaches to maximize training efficiency and mitigate effects of limited training exemplars, bad labels, and noisy data
- Example 1: integration of Generative Adversarial Networks (GANs) into the training process (joint probability between inputs and outputs); shown in one study to reduce training iterations by factor of 10
- Example 2: Pseudo-Labels (labels that are created automatically for unlabeled data using a partially trained network) can significantly improve classification accuracy without changing network architecture, based on theory known as Entropy Regularization

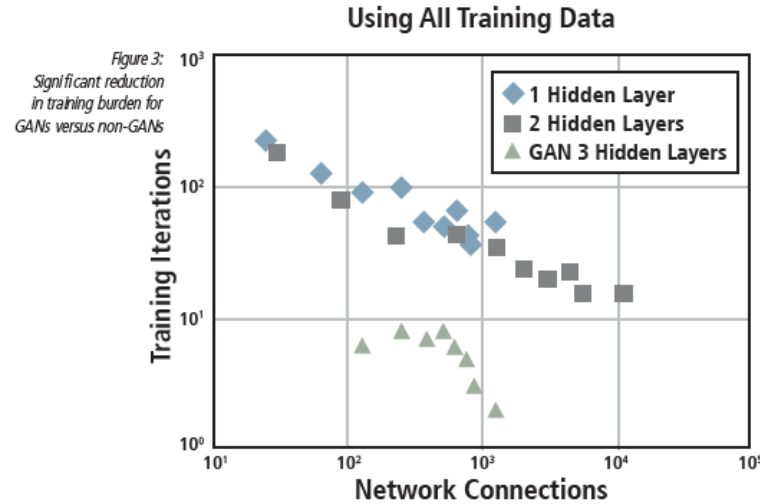
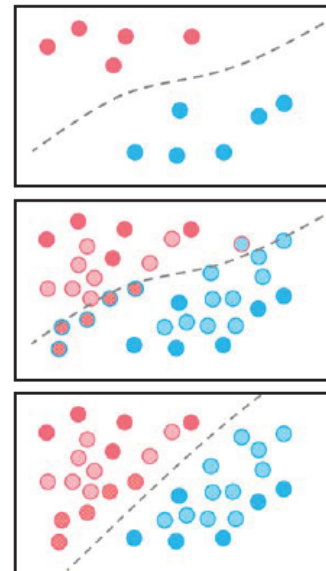
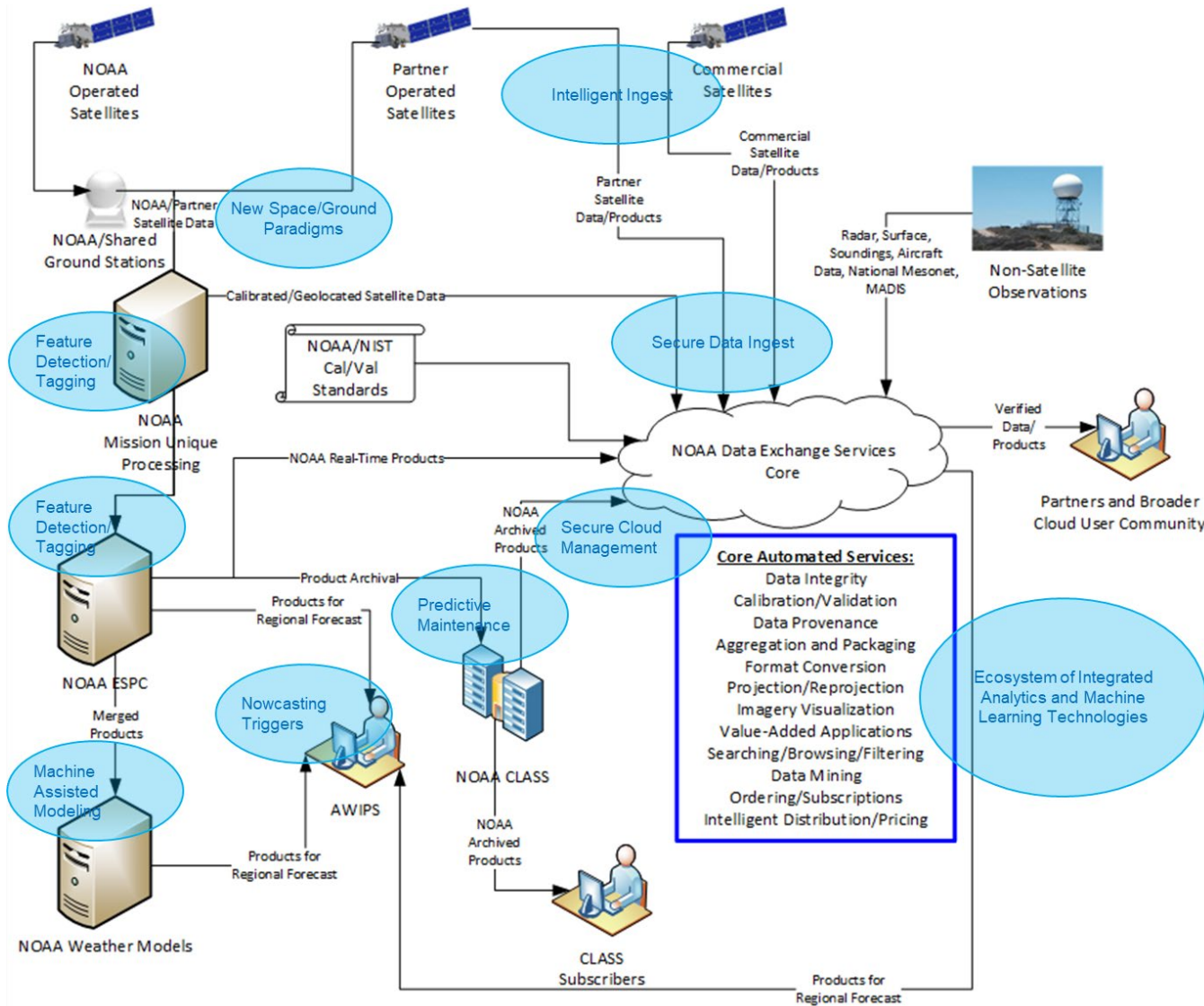


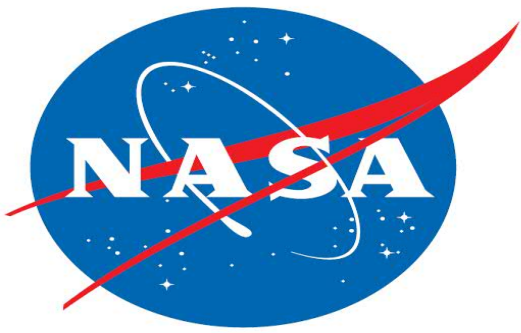
Figure 4:
 TOP: Sparse labeled data has ambiguous class boundaries,
 MIDDLE: Unlabeled data with pseudo-labels (outlines) added to dataset,
 BOTTOM: Re-training with pseudo-labels corrects decision boundaries based on data population density.



Analytics, ML, and the Weather Value Chain



Impact Area	Forecast Benefit
Accelerated Data Ingest	<ul style="list-style-type: none"> Commercial data available for forecast models Reduced data latency
Identify high-value data for forecast models during sensor processing	<ul style="list-style-type: none"> Improve initial conditions for model runs Improved Satellite Data Assimilation
Training models for weather anomalies	<ul style="list-style-type: none"> Increase speed and accuracy of extreme weather forecasts Provide forecaster tools to train regional datasets
Federate Weather data to process disparate data sets in the cloud	<ul style="list-style-type: none"> Faster R2O2R due to removal of data stovepipes “multi-int” data fusing across commercial/govt data types
Cloud technologies (massively parallel systems, GPUs, cloud computing)	<ul style="list-style-type: none"> ML for all classes of users, with elasticity for cost control Better collaboration across weather enterprise community



Raytheon

