Inference At the Edge: A Case Study at the Amazon Spheres

Miro Enev, PhD Sr. Solution Architect NVIDIA WenMing Ye Specialist Solution Architect Amazon

Introduction

Agenda

Introduction: AI @ Amazon Spheres

Video: Welcome to the Amazon Spheres [living wall video]

Approach:

Anomaly Detection using DL on Time-Series Sensor Streams

Architecture:

Training (Amazon SageMaker) Inference (NVIDIA Jetson Xavier, Amazon SageMaker Neo)

Results:

Improved alerting

Future Work:

Computer vision based plant stress



Our Goal = Help the Caretakers

Claire



Ben







Temperature





Challenge 1: Lots of Systems to Manage





Challenge 2: Too Many Suspicious Values

Amazon Rufus 🗸			•	Miro Enev V SkySpark				
Site Spark				ii ★ 🕹 🖬				
Targets View Timeline Rules Select Info								
All Block 19								
Group	Rules	dur	Timelines	Targets				
Group () AC-46-2-1	Rules O Suspicious Values	dur 4.5hr	Timelines	Targets				
Group () AC-46-2-1 1 sparks	Rules O Suspicious Values	dur 4.5hr	Timelines	Targets				
Group () AC-46-2-1 1 sparks () AC-46-4-2 ()	Rules Image: Supplementary	dur 4.5hr	Timelines	Targets				

flt_suspiciousValue

flt_suspiciousValue(thePoint, sparkDate, duration: 30min, negativeThreshold: 0, airTempHighLimit: 140, waterTempHighLimit: 180, steamTempHighLimit: 250)

This point has data values that are less than zero, suspicious temperature values if the point is a temperature sensor or null data values.

Parameters (Default):

- duration (30min): The duration of the spark condition before firing a spark.
- negativeThreshold (0): The buffer of allowed negative numbers. Any negative numbers above this value are not analyzed. Used to disregard sensors that show slightly negative values when off.
- airTempHighLimit (140): The maximum allowed air temperature. Any air temp value above this limit will produce a spark.
- waterTempHighLimit (180): The maximum allowed water temperature. Any water temp value above this limit will
 produce a spark.
- steamTempHighLimit (250): The maximum allowed steam temperature. Any steam temp value above this limit will
 produce a spark.

Pseudocode: Find periods of time a number point with history data, that is not calculated and not outside air temperature is null and/or less than 0 for 30 minutes. In addition, find periods of time when a temperature sensor returns a value higher than 140 for "air" tags, 180 for "water" tags, or 250 for "steam" tags for 30 minutes.

When Issues Occur, They Go Unnoticed

Example 1: During a product launch (Alexa microwave integration), event organizers requested that the temperature be lowered for media and the air velocity reduced for better acoustics.

Problem: Incorrect temp. and air velocity for 4th floor plants for a week

Example 2: Building automation staff suspended the irrigation for the living wall to update/repairs several sensors.

Problem: 24 hours without water for living wall [low irrigation pressure warning was ignored]

Approach

Al to Assist the caretakers

- Accurate Alerts [low false alarm rate]
- Real-time & Low Cost
- Enable Current/Future Science
- Scalability & Availability of Technology







Why DL



Deep Learning @ Spheres [AutoEncoder Network]



Deep Learning @ Spheres [AutoEncoder Network]





Correlated Sensors [Weekday & Weekend Behaviors]



Detecting Anomalies



Multi Sensor Models

[AutoEncoder Network]







Sensor Zones

Living Wall [4 floors]

t,rh,d,co2 light level DLI-46-4-DG5] [X, AC-46-2-2, AC-46-3-2, AC-46-4-3]

[DLI-46-1-DG1, DLI-46-2-DG2, DLI-46-3-DG3,

North Conservatory [1st floor]

t,rh,d,co2 [AC-46-1-1, AC-46-1-2, AC-46-1-3, AC-46-1-4] light level [DLI-46-1-DS1, DLI-46-1-DM2, DLI-46-1-DM3]

South Conservatory [2nd floor]

t,rh,d,co2 [AC-46-2-3, AC-46-2-4, AC-46-2-5, AC-46-2-6] light level [DLI-46-2-DM1, DLI-46-2-DM2, DLI-46-2-DM3]

Canopy [3 floors above N. Conservatory]

t,rh,d,co2 [AC-46-2-1, AC-46-3-1, AC-46-4-2] light level [DLI-46-4-DL1, DLI-46-4-DL2, DLI-46-4-DL3, DLI-46-4-DL4, DLI-46-4-DL5, DLI-46-4-DL6, DLI-46-4-DL7, DLI-46-4-DL8, DLI-46-4-DL13, DLI-46-4-DL14]









Architecture



Amazon SageMaker Neo



Amazon SageMaker Neo

compiled_model.json - Visual Studio Code						
n D de	mo_documentation.txt D demo_workflow.txt C compiled_model.json	🗅 compiled.so 🌒 🛛 😶 Compiled_model.json 🗙				
D def 1 1 2 2 3 3 5 6 5 6 5 6 5 6 5 6 5 6 5 6 5 6 5 7 7 8 9 10 11 13 13 14 16 16 16 16 17 19 200 21 222 23 24 25 26 29 31 32 33 34 355 36 36 37 389 9 40 41 41 43	<pre>no_documentation.txt</pre>	Compiled_model.jon - Visual Statute Code Compiled_so • •••• D compiled_model.jon × 64 ** op": "null", 65 ** op": "tensorrt_subgraph_op", ** op": "num_inputs": "1" 76 77 ** op": "num_outputs": "1" 76 77 ** op": "num_outputs": "1" 76 77 ** op!, 80 81 83 10, 0, 0], 83 11, 0, 0], 84 11, 0, 0], 85 11, 0, 0], 85 11, 0, 0], 87 19, 0, 0], 88 110, 0, 0], 90 112, 0, 0] 91 91 91 92 ** subgraphs": [** op": "null", ** inputS*: [] 93 100 101 ** op": "null", ** inputS*: [] 95 106 107 ** op": "null", ** inputS*: [] 77 ** inputS*: [] 77 ** inputS*: [] 77 ** op": "null", ** op": op!				
45 46		108 "inputs": [] 109 }, 110 { "est "forces				
O 0 0 O			Ln 76, Col 10 Spaces: 2 UTF-8 LF JSON 🙂 🌲			

Training Architecture @ p3.4xlarge



Inference Architecture @ Jetson Xavier



Notebook Demo

Results

Sample Reconstructions





[Synthetic] Anomaly Detection





Real Anomaly Detection





" Nice catch. We altered the climate to encourage the blooming of our Amorphophallus titanum plant. The corpse flower is more accustomed to warmer temps and higher humidity than the normal spheres operating parameters. "

Future Work

Multi-spectral Imaging





Discussion & Q/A

Thank you!

WenMing Ye - wye@amazon.com

Miro Enev - menev@nvidia.com

Scheduled Lambdas Trigger Training and Batch Inference

CloudWatch Dashboards	Step 1: Create rule			
Alarms	Create rules to invoke Targets based on Events happening in	your AWS environment.		
ALARM 6 INSUFFICIENT 0 OK 0 Billing	Event Source Build or customize an Event Pattern or set a Schedule to invoke Targets.	Targets Select Target to invoke when an event matches your Event Pattern or when schedule is triggered.		
Events Rules Event Buses Logs Insights Metrics Alpine	 Event Pattern Schedule Fixed rate of 	Lambda function	<pre>def lambda_handler(event, context): ec2 = boto3.client('ec2', region_name=region) ec2.start_instances(InstanceIds=instances) print 'started your instances: ' + str(instance)</pre>	
Favorites	{ "version": "0", "id": #89d1a02d-5ec7-412e-82f5-13505f849b41", "idetail-type": "Scheduled Event", "source": "aws.events", "account": "122456789012", "time": "2016-12-30178-44.492", "region": "us-east-1", "resources": ["arn:aws:events:us-east-1:123456789012:rule/SampleRule"], "detail": {}		<pre>def lambda_handler(event, context): ec2 = boto3.client('ec2', region_name=region) ec2.stop_instances(InstanceIds=instances) print 'stopped your instances: ' + str(instances)</pre>	
	* Required	Cancel Configure details		

https://aws.amazon.com/premiumsupport/knowledge-center/start-stop-lambda-cloudwatch/

Multi-spectral Imaging & Computer Vision



Edge Processing + TensorRT



