

AI Application Deployment: From Cloud & Edge to the FarEdge Dr. John Soldatos, Senior Innovation Delivery Specialist OpenDay-CTM 2021, May 4th, 2021

INTRASOFT in a nutshell







Over 20 years of successful R&D projects implementation and...



- System Integration & Architecture Specification
- Technical Coordination
- Application **Development**
- Communication

Efficient/effective **coordination** and team-leading EC co-funded R&D practices:

- FP5/FP6/FP7/H2020
- CIP/ECSEL/CEF
- Studies/Impact assessment/Policy support

...innovation



Pre-commercialisation of Information and Communication Technologies **research results**



Innovation-related activities:

- Internal Innovation
 Competition
- Innovation radar/Technology transfer
- Joint ventures

Strategy positioning:

- Feasibility study of cutting-edge technologies (Big Data, Data Analytics, Blockchain, ...)
- Opening up new innovative fields (agrifood, AI, maritime security, ...)



Product Development and cross - collaboration:

- Extend/New product offerings
- Cross-departmental business and commercial collaboration
- Spin-off (SmartPACT plockchain in insurance)

Edge Computing for Industrial Use Cases: When Cloud is Not



- Rationale for Edge Computing in Industrial Use Cases:
 - Latency & Real-Time Field Operations (e.g., Real-Time Actuation and Control)
 - Energy Efficiency (e.g., Sustainable Manufacturing)
 - Data Privacy and Protection (e.g., Industrial Data)
- Different Edge Computing Approaches:
 - Edge Computing with Local Clouds
 - Fog Computing
 - FarEdge Approaches (e.g., Embedded Machine Learning)
 - Federated Machine Learning
 - Decentralization based on Distributed Ledger Technologies

Edge Computing VS. Standard Industrial Architectures (source: https://www.iiconsortium.org/)





OpenFog Reference Architecture

Actuator Sensor with Routing Actuator with Routing Actuator with Routing



FEATURES	CLOUD (e.g., IOT/Cloud Integration (WAN))	EDGE (e.g., Local Clouds, Fog Computing)	FarEdge (e.g.,AIoT, Embedded Devices, Microcontollers)
Data Points Availability	High	Medium	Low
Energy Efficiency	Low	Medium-to-High	Very High
Privacy	Low-Medium	Medium-to-High	High
Real-Time Opportunities	Low	Medium-to-High	Very High

Cloud vs. Edge vs. FarEdge



	Transfer Learning								
INDUSTRY USE CASE	CLOUD	EDGE	FarEdge						
Asset Management & Predictive Maintenance	RUL Calculation at Cloud	RUL Calculation at Edge Cluster or Edge Devices	Fault Detection Inside the Machinery						
Quality Management & Zero Defect Manufacturing (ZDM)	Quality Inspection – Defect Prediction	Defect Prediction & Near Real-Time Control	Defect Detection & Real- Time Control						
Value Proposition	ML Model Accuracy	Speed vs. Accuracy Balance	Speed and Power Efficiency						
	Predictive, Preventiv Strategies to Asset Mar Manager								



- Is AI/ML in the Cloud Always the Best Option?
 - Siri, Alexa, OK Google etc.: benefit from an instant response
 - Industrial maintenance: More timely detection of failures and abnormalies
 - Agriculture: Instant disease detection using a plant's image
 - In several cases ML/AI «at the edge» dramatically reduces costs & complexity, and limits potential data privacy leaks.
- TinyML: An Alternative Form of Machine Learning and AI at the Far Edge



TinyML Hardware Examples



- ARDUINO NANO 33 BLE
 SENSE WITH HEADERS, Cost
 ~ 30€
- SparkFun Edge Development Board -Apollo3 Blue DEV-15170, Cost ~ 15€

 STM32F746G-DISCO discovery board (32F746GDISCOVERY) - complete platform for STMicroelectronics ARM[®] Cortex[®]-M7 corebased STM32F746NGH6 microcontroller, Cost ~ 70€



Building TinyML Models and Applications

- Standard Tools Work for Data Science & Model Development:
 - Python, Jupyter Notebooks, Arduino IDE,...
- Machine Learning Framework for TinyML:
 - <u>**TensorFlow</u>**: Suite of tools for building and running ML models</u>
 - <u>Keras</u>: TensorFlow's high-level API focused on building and training Deep Learnin applications
 - <u>TensorFlow Lite</u>: Specifically designed for inference on devices with limited computing capacity (e.g., phones, tablets, embedded devices).
 - <u>TensorFlow Lite Micro</u>: Deploy models on microcontrollers and other devices with only few kilobytes of memory e.g.,
 - Core runtime just fits in 16 KB on an Arm Cortex M3
 - Doesn't require operating system support, any standard C or C++ libraries, or dynamic memory allocation



TensorFlow

 Θ



INTRASOFT's DataCrop IIOT Platform







Attributes



Processing Jobs are of three types:

- Pre-processing
- Storage
- Analytics (e.g., ML, calculation)

Nikos Kefalakis, Aikaterini Roukounaki, John Soldatos: Configurable Distributed Data Management for the Internet of the Things. Inf. 10(12): 360 (2019)

ML Framework Integration over DataCrop - QARMA

- Different ML Algorithms run over DataCrop
 - Recurrent Neural Networks (RNN)
 - Long Short-Term Memory Networks (LSTM)
 - Attention-based Networks
- QARMA4Industry:
 - Rules Mining Approach
 - Identify Rules that hold on data and quantify them
 - Used in H2020 QU4LITY Project



Ioannis T. Christou, Nikos Kefalakis, Andreas Zalonis, John Soldatos, Raimund Bröchler, End-to-End Industrial IoT Platform for Actionable Predictive Maintenance, IFAC-PapersOnLine,

Volume 53, Issue 3, 2020, Pages 173-178, ISSN 2405-8963,

-									and the second second							_
4865	4799725	0.556	0.589	0.647	0.964	0.414	1.212	0.044	0.018	0.04	0.038	0.341	1.24	0.724	0.497	0.6
4895	4179254	0.546	0.57	0.923	1.016	0.197	1.306	0.085	0.038	0.033	0.053	0.197	1.242	0.743	0.536	0.5
9922	4181234	0.52	0.569	0.77	1.023	0.19	1.254	0.073	-0.005	0.033	0.028	0.252	1.259	0.87	0.564	0.6
4956	4261778	0.498	0.755	0.854	1.016	0.415	1.288	0.068	-0.028	0.027	0.037	0.273	1.298	0.952	0.497	0.7
976	4262645	0.491	0.672	0.871	1.044	0.156	1.299	0.085	0.022	0.046	0.028	0.26	1.328	0.987	0.497	0.6
5078	4287840	0.539	0.705	0.785	1.057	0.491	1.25	0.107	0.029	0.047	0.012	0.343	1.323	0.972	0.488	0.6
5081	4151867	0.517	0.634	0.88	1.028	0.384	1.295	0.07	0.056	0.031	-0.001	0.065	1.245	0.774	0.529	0.5
5093	4873303	0.526	0.587	0.76	1.172	0.332	1.242	0.198	0.093	0.126	0.08	0.679	1.321	0.901	0.569	0.6
\$151	4732013	0.517	0.651	1.713	1.115	1.5	1.562	0.152	0.085	0.094	0.027	0.54	1.291	0.802	0.527	0.6
5182	4766937	0.537	0.68	0.814	1.018	0.42	1.287	0.064	-0.018	0.034	0.007	0.443	1.39	1.102	0.461	0.7
\$204	4280359	0.586	0.684	0.892	1.031	0.696	1.192	0.079	-0.012	0.031	-0.014	0.485	1.398	1.094	0.502	0.7
5210	4741697	0.507	0.674	1.317	1.084	0.865	1.45	0.148	0.061	0.07	0.019	0.571	1.359	0.997	0.536	0.6
211	4830223	0.577	0.663	0.761	1.153	0.527	1.243	0.182	-0.026	0.04	0.015	0.521	1.323	0.948	0.547	0.6
215	4773427	0.441	0.612	1.135	1.078	0.404	1.404	0.115	0.024	0.052	0.07	0.558	1.315	0.887	0.472	0.6
\$233	4261770	0.497	0.725	0.952	1.044	0.273	1.298	0.1	0.055	0.055	0.013	0.579	1.274	0.787	0.535	0.6
5236	4737824	0.577	0.659	0.901	1.095	0.821	1.292	0.12	0.027	0.035	0.076	0.516	1.32	0.977	0.564	0.7
5266	4128382	0.536	0.582	0.937	1.033	0.063	1.288	0.075	0.059	0.058	0.014	0.267	1.235	0.702	0.55	0.6
5287	4918946	0.516	0.682	0.857	1.15	0.466	1.266	0.188	0.099	0.046	0.047	0.574	1.283	0.823	0.517	0.6
5337	4750905	0.551	0.756	0.992	1.083	0.431	1.33	0.137	0.091	0.084	-0.117	0.543	1.346	0.964	0.47	0.6
488	4843862	0.585	0.592	1.005	1.256	0.649	1.361	0.313	0.203	0.18	0.02	0.442	1.295	0.845	0.518	0.7
5554	4836785	0.486	0.594	1.297	1.283	0.722	1.451	0.319	0.304	0.384	0.054	0.289	1.236	0.75	0.511	0.6
9614	3738214	0.2	0.535	0.656	1.014	0.467	1.202	0.057	-0.022	0.038	0.041	0.289	1.279	0.843	0.203	0.5
9621	4221970	0.508	0.605	1.055	1.063	0.266	1.338	0.098	0.087	0.069	0.096	0.353	1.244	0.701	0.494	0.6

QARMA4Industry outperforms popular algorithms





UCode:45.714285714285715

CS1 AVG

0.981 0.841 1.015 0.763 0.863 0.918

1.005

H2020 FAR-EDGE : Distributed Ledger Technologies for Orchestrating & Configuring Processes at the Edge



- Edge Computing Functions:
 - Edge Automation & Edge Analytics
- Edge Automation:
 - Low-overhead, low-latency connectivity
- Two way interactions
 - Edge Analytics:
 - Real-time analytics
 - One-way data savvy operations



Distributed ledger technology for decentralization of manufacturing processes. ICPS 2018: 696-701



DLT Technologies for Edge Computing Functions (see: www.Edge4Industry.eu)





Edge Simulation Use Case (WhiteAppliances Factory) (source: FAR-EDGE Project)



5G: New Capabilities & Services - Impact on Vertical Sectors

20 BILLION

₽ 90%



Source: 5G Infrastructure Association: Vision White Paper, February 2015



RASOFT

Source: 5G Infrastructure Association: 5G Empowering vertical industries. White Paper, 2016

Why 5G/6G for Edge Intelligence?





• Flexibility through Slicing and NFVs

Enabling

 Seamless connectivity across multiple IoT technologies (e.g., RFID, short-range communications, UWB, Blue Tooth,

• Tactile Internet • Edge Al



Source: Why 5G is the way to go for IoT: https://www.linkedin.com/pulse/why-5g-way-go-internet-things-john-soldatos/

Conclusions



- Alternative Edge Computing Approaches in Industrial Environments
 - With 5G/6G as key Enabler
- No Silver Bullet: Selection Depends on Requirements:
 - Latency / Real-Time Concerns
 - Power Efficiency
 - Privacy
 - Decentralization
 - Security & Robustness
 - Federated Learning (not addressed in this presentation)
- An "Edge Operating System" should take into account relevant trade-offs
 - Zero Defects Manufacturing Approach Combines Reactive, Preventive and Predictive Strategies
- Infrastructures like GAIA-X can boost a modular approach



Thank you for your attention



