## **AI INTERCONNECT FOR 6G**

Sasu Tarkoma University of Helsinki

HELSINGIN YLIOPISTO HELSINGFORS UNIVERSITET UNIVERSITY OF HELSINKI

## CONTENTS

BACKGROUND
TECHNOLOGY MIX: 6G AND AI
AI INTERCONNECT AND NEURAL PUB/SUB
BUCATION: 6G MOOC



Finland's 6G network of excellence.





## NOKIA CENTER FOR ADVANCED RESEARCH (NCAR) OVERVIEW

#### **Publications Events Externally Funded Projects and** International Collaborations +70 scientific papers published Workshops, conferences and other Multiple national and international on high-ranking journals and research projects dissemination events co-organized conference since 2016 Key project: Academy-NSF Lean6G together with the FCAI and the 6G (https://tinyurl.com/5asazpbz) with professor Leandros Tassiulas at Flagship ൳൭ഁ EVENTS **Yale University** Patents and patent applications FCAI Finnish Center for Artifician FLAGSHIP UNIVERSITY Education People Research An international team of 5 key research themes 1-2 PhD degrees annually and +12 members including faculty, investigated: many students working on university researchers, and **Environmental Sensing** theses. Edge Computing and IoT postdocs currently involved in Online education (MOOCs) Low Latency Transport Protocols NCAR-related activities

Cloud Computing and NFV Distributed and Cognitive Processing

https://www.helsinki.fi/en/networks/6g-research





- Four partners:
- University of Helsinki (UH), lead
- KTH Royal Institute of Technology (KTH)
- Aarhus University (AU)
- Norwegian University of Science and Technology (NTNU).
- Research focus on Edge Intelligence of each partner (shown in Figure 1).
- Coherent, joint strategic approach to take a good advantage of joint cross-border collaboration.

#### **KUMPULA CAMPUS AS** "PLAYGROUND" FOR OUR RESEARCH



## CONTENTS

> BACKGROUND
> TECHNOLOGY MIX: 6G AND AI
> AI INTERCONNECT AND NEURAL PUB/SUB
> EDUCATION: 6G MOOC

## **KEY DEVELOPMENTS THAT ARE CONVERGING**





**'ට** 

**Cellular networks:** cloud RAN, O-RAN, network slicing, Mobile Edge Computing (MEC), private networks



**AI:** distributed techniques (federated learning, split learning), transfer learning, differential privacy, LLMs, regulation

## Proceedings of LEEE

#### Edge Intelligence: Empowering Intelligence to the Edge of Network

Video Summarization Using Deep Neural Networks: A Survey

Point of View: How Fast Do Algorithms Improve? Scanning Our Past: Swedish Cryptology I



Edge Intelligence: Empowering Intelligence to the Edge of Network D Xu, T Li, Y Li, X Su, S Tarkoma, T Jiang, J Crowcroft, P Hui Proceedings of the IEEE 109 (11), 1778-1837. (60 pages)

## **Centralized cloud for Al**

## **Edge Intelligence**



Edge Intelligence: Empowering Intelligence to the Edge of Network. D Xu, T Li, Y Li, X Su, S Tarkoma, T Jiang, J Crowcroft, P Hui Proceedings of the IEEE 109 (11), 1778-1837.



#### A Taxonomy





# FOG COMPUTING FOR DEEP LEARNING WITH PIPELINES

We demonstrate split learning and inference pipelines on device, cloudlet and cloud with a vision based use case (CIFAR-10).

If the fog hardware can be fully utilized, it will improve the throughput of heavy workloads.



- Fog is a hierarchical network structure.
- Resources in fog can be pooled with pipelines.
- Sparse: Stream Processing Architecture for Resource-Subtle Environments<sup>a</sup>

<sup>&</sup>lt;sup>a</sup>https://github.com/AnteronGitHub/sparse



## **KEY TASKS: LLMS AS CONTROLLERS**



HuggingGPT: Solving AI Tasks with ChatGPT and its Friends in HuggingFace. <u>https://huggingface.co/papers/2303.17580</u> Large Language Models as Tool Makers. <u>https://arxiv.org/pdf/2305.17126.pdf</u>

#### HELSINGIN YLIOPISTO HELSINGFORS UNIVERSITET UNIVERSITY OF HELSINKI

## **Distributed GPTs in the Network**

#### Drones and other autonomous sensors



Distributed AI based on foundation models that are incrementally finetuned (stacking property) for application domains

Orchestrator places and synchronizes models and state

Interconnect enables orchestration

## CONTENTS

> BACKGROUND
> TECHNOLOGY MIX: 6G AND AI
> AI INTERCONNECT AND NEURAL PUB/SUB
> EDUCATION: 6G MOOC



#### AI INTERCONNECT – LARGE LANGUAGE MODELS FOR 6G SYSTEMS

Al-native Interconnect Framework for Integration of Large Language Model Technologies in 6G Systems Sasu Tarkoma, Roberto Morabito, Jaakko Sauvola https://arxiv.org/abs/2311.05842

> Al-native Interconnect Framework for Integration of Large Language Model Technologies in 6G Systems

> > Sasu Tarkoma\*, Roberto Morabito\*, Jaakko Sauvola

Abstract—The evolution towards 6G architecture promises a transformative shift in communication networks, with artificial intelligence (AI) playing a pivotal role. This paper delves deep into the seamless integration of Large Language Models (LLMs) and Generalized Pretrained Transformers (GPT) within 6G systems. Their ability to grasp intent, strategize, and execute intricate commands will be pivotal in redefining network functionalities and interactions. Central to this is the AI Interconnect framework, intricately woven to facilitate AI-centric operations within the network. Building on the continuously evolving current state-of-the-art, we present a new architectural perspective for the upcoming generation of mobile networks. Here, LLMs and GPTs will collaboratively take center stage alongside traditional pre-generative AI and machine learning (ML) algorithms. This union promises a novel confluence of the old and new, melding tried-and-tested methods with transformative AI technologies. Along with providing a conceptual overview of this evolution, we delve into the nuances of practical applications arising from such an integration. Through this paper, we envisage a symbiotic integration where AI becomes the cornerstone of the next-generation communication paradigm, offering insights into the structural and functional facets of an AI-native 6G network.

Index Terms—6G, Generative AI, Large Language Model (LLM), Generative Pre-trained Transformers (GPT), AI Interconnect, Edge Intelligence, Open RAN (O-RAN).

HELSINGIN YLIOPISTO HELSINGFORS UNIVERSITET UNIVERSITY OF HELSINKI



#### AI INTERCONNECT – LARGE LANGUAGE MODELS FOR 6G SYSTEMS

The paper addresses **opportunities**, **challenges**, and possible **solutions** arising from the integration of Large Language Models, such as GPT, in the context of 6G systems

**On-Device Cognitive Computing** 

#### Edge Intelligence with LLMs

Delving into the intricacies of deploying LLMs in edge environments, addressing the challenges and innovations within 6G mobile edge computing (MEC) systems.

#### Generative AI in 6G

A holistic view exploring the synergy between generative large language models (LLMs) and 6G technologies, outlining the transformative potential for future wireless networks. Beyond edge infrastructure: exploring the feasibility and implications of integrating LLMs directly into devices, paving the way for truly decentralized intelligence.

> LLMs for 6G State-of-the-Art

#### Demystifying Telco Through LLMs

Leveraging the natural language processing capabilities of LLMs to interpret, manage, and simplify complex telco terminologies and operations. This encapsulates the potential to make telecommunication operations and maintenance (OAM) accessible to non-experts.

#### Diverse Applications in the 6G Ecosystem

Highlighting a spectrum of use cases from hardware acceleration, software optimization, IoT robotics, to dynamic network topology management, emphasizing the versatility of LLMs in the evolving 6G landscape.

HELSINGIN YLIOPISTO HELSINGFORS UNIVERSITET UNIVERSITY OF HELSINKI







## NEURAL PUB/SUB AND LARGE LANGUAGE MODELS

- Large Language Models (LLMs) like OpenAI's GPT at the forefront of 6G technological evolution.
- Enhancing network performance through intelligent management and optimization.
- Development of a proof-of-concept using GPT for advanced network slicing management.

Roberto Morabito et al. work in progress



## NEURAL PUB/SUB FOR FUNNEL PATTERN

- Builds upon the mapping pattern for advanced data handling.
- Involves subscribing to one or more publications and applying a specific function (F).
- Function (F) determines the order and timing of data processing.
- Results in a single emitted publication based on the output of function F.



Funnel pattern workflow of a Neural Pub/Sub broker system designed to aggregate CQI data from three different slices.

- Real-time monitoring of Channel Quality Indicator (CQI) across network slices.
- Automated generation of Slices Status Reports in natural language, offering actionable insights for network operators.
- Dynamic adjustments to network slices based on CQI data analysis to maintain optimal service quality.



Funnel pattern workflow of a Neural Pub/Sub broker system designed to aggregate CQI data from three different slices.

HELSINGIN YLIOPISTO HELSINGFORS UNIVERSITET UNIVERSITY OF HELSINKI



Funnel pattern workflow of a Neural Pub/Sub broker system designed to aggregate CQI data from three different slices.

#### HELSINGIN YLIOPISTO HELSINGFORS UNIVERSITET UNIVERSITY OF HELSINKI

```
def query openai(data):
    start time = time.time()
    response = clientOpenAI.chat.completions.create(
        model="gpt-3.5-turbo-1106", #gpt-3.5-turbo-1106 gpt-4-32k-0314 gpt-4
        response_format={ "type": "text" },
        #stream=True,
        messages=[
            {"role": "system", "content": "The system receives CQI data for each network slice in real-time or at scheduled intervals.
             It analyzes the CQI data to detect any values that fall below a predefined threshold which indicates suboptimal performance.
             Upon detecting low CQI values, the system determines the best course of action, such as suggesting a reallocation of resources
             or adjusting other slice parameters like bandwidth or priority levels. The system could also be responsible for generating alerts for network operators,
             providing detailed reports, or executing pre-approved actions to remedy the situation."},
            {"role": "user", "content": f"We have a network slicing system where each slice is dedicated to a specific service type.
             The system continuously monitors the Channel Quality Indicator (CQI) to maintain high service quality.
             When a slice's CQI falls below the threshold of 10, it's an indicator that the service quality is degrading.
             As the system, analyze the incoming CQI data and suggest actions to reallocate resources dynamically and adjust the slice parameters to improve the CQI.
             Consider factors such as current network load, priority of the service, and historical performance data of the slice when making suggestions.
             Generate a response that outlines potential adjustments to the network slice configuration to address the low CQI values. Consider the following {data}"}
    end time = time.time()
    processing latency = end time - start time
```

return response, processing\_latency

The Query Function at the Broker: The interconnect design accommodates multimodal subscriptions and requests, complex subscriptions/tasks can be decomposed



#### DEMO

#### Phase 1:

The Neural Pub/Sub Broker manages the publication of network slices data



HELSINGIN YLIOPISTO HELSINGFORS UNIVERSITET UNIVERSITY OF HELSINKI



#### DEMO

#### Phase 2:

The Network Operator subscribes to a GPT-generated Slices Status Report



HELSINGIN YLIOPISTO HELSINGFORS UNIVERSITET UNIVERSITY OF HELSINKI



#### DEMO

**Phase 3:** Network Slices data is published by the three network slices

y) robertomorabito@RobertoMacBook-Pro neural_pub_sub_bro n publisher1.py	(gpt_latency) robertomorabito@RobertoMacBook-Pro neural_pub_sub_broker % python neural_pub_sub_broker_agg.py	(gpt_latency) robertomorabito@RobertoMacBook-Pro neural_pub_ % python subscriber.py
py:14: DeprecationWarning: Callback API version 1 is dep	neural_pub_sub_broker_agg.py:92: DeprecationWarning: Callback API versi	subscriber.py:17: DeprecationWarning: Callback API version 1
date to latest version matt (light(matt callbackAPTVersion VERSTON1 "COT SLICE	broker client - matt Client(matt CallbackAPTVersion VERSION1 "Neural	sub client - matt Client(matt CallbackAPTVersion VERSTON1
	Broker")	eSubscriber")
data published	Connected with result code 0	Connected with result code 0
		Subscribing to GPT-generated Slices Status Report
		0
	Data Received from all the Slices	
	The Neural Pub/Sub Broker is set to query the GPT-based Slice Analysis	
	Service to examine the CQI data of various network slices, identifying	
	any values falling below a predefined threshold that indicates suboptim	
Θ	al performance.	
y) robertomorabito@RobertoMacBook-Pro neural_pub_sub_bro		
n publisher2.py		
py:14: DeprecationWarning: Callback API version 1 is dep		
matt Client(matt Callback&PTVersion VERSTON1 "COT SLICE		
inderretenedinderrenteneden interstorit, edizette		
data published		
Θ		
y) robertomorabito@RobertoMacBook-Pro neural_pub_sub_bro		
n publisher3.py		
py:14: DeprecationWarning: Callback API version 1 is dep	8	
matt.Client(matt.CallbackAPIVersion.VERSION1, "COI SLICE	A	
data published		

HELSINGIN YLIOPISTO HELSINGFORS UNIVERSITET UNIVERSITY OF HELSINKI



#### DEMO

#### Phase 4:

The Neural Pub/Sub Broker query the **GPT-based Slice Analysis Service** 



HELSINGIN YLIOPISTO HELSINGFORS UNIVERSITET UNIVERSITY OF HELSINKI



#### DEMO

#### Phase 5:

The Neural Pub/Sub Broker publishes to the Subscriber the **GPT-generated Slices Status Report** 



HELSINGIN YLIOPISTO HELSINGFORS UNIVERSITET UNIVERSITY OF HELSINKI





#### **FULL DEMO**

HELSINGIN YLIOPISTO HELSINGFORS UNIVERSITET UNIVERSITY OF HELSINKI

- Demonstrated effectiveness using real data.
- Comparison between GPT-4 and GPT-3.5 models:
  - GPT-3.5 shows lower processing latency and higher verbosity, highlighting trade-offs in deployment.
- Private models (llama3, mixtral etc) are also suitable



Comparative analysis of GPT-4 and GPT-3.5, illustrating the sum of processing latency and network latency against the average completion tokens, represented by red dots for visual distinction.



## **NEURAL PUB/SUB ROUTER**



HELSINGIN YLIOPIS HELSINGFORS UNIV UNIVERSITY OF HE Deliver to clients

Forward subscriptions and advertisements to external systems in the distributed environment



#### LLM-Based Content Router

#### Input Subscriptions:

```
prompt : "Net the whether a transfer update tails of any for device, especially for our smart thermostats and lighting systems, to ensure an devices are turning the
latest software for optimal performance and security."
},
{
    "identifier": 4,
    "prompt": "Provide real-time alerts if any industrial IoT sensors in our manufacturing plant lose connectivity for more than 5 minutes, as this could indicate critical
production issues or equipment failures."
},
{
    "identifier": 5,
    "prompt": "Send an alert if any environmental sensors, like those monitoring temperature, humidity, or air quality, report readings outside of our set safety thresholds in our
storage facilities, to prevent damage to sensitive goods."
```

#### Input Events:

```
{
    "event": "Battery level of security camera in downtown clinic falls to 8%",
    "matching_subscription_ids": [1]
    },
    {
        event": "Weekly data usage report shows a smart meter exceeded 600MB in a single day",
        "matching_subscription_ids": [2]
    },
    {
        event": "Firmware update failed on smart thermostats in multiple residential buildings",
        "matching_subscription_ids": [3]
```

Start the routing process.

#### Aggregated Subscriptions: Aggregated Subscriptions:

[{"identifier": [1, 3], "Ilm\_output": "Notify me immediately when critical IoT devices, like healthcare monitors or security cameras, have less than 10% battery remaining or whenever a firmware update fails on any of these devices, to ensure we can proactively replace or recharge them and that all devices are running the latest software for optimal performance and security."}, {"identifier": [2, 5], "Ilm\_output": "Send a weekly summary of data usage and real-time alerts for environmental sensors for all IoT devices, highlighting any device that exceeds 500MB of data in a day or reports readings outside of our set safety thresholds in our storage facilities, to ensure we stay within our data plan limits, identify any unusual activity, and prevent damage to sensitive goods."}, {"identifier": [4], "Ilm\_output": "Provide real-time alerts if any industrial IoT sensors in our manufacturing plant lose connectivity for more than 5 minutes, as this could indicate critical production issues or equipment failures."}] JSON difference to original 152 and percentage 12.958226768968457Number of subs in the merged 5 Number of subs in the merged 3

#### Matching Events Per Subscription:

Matching: [{ "event": "Battery level of security camera in downtown clinic falls to 8%", "matching\_subscription\_ids": [1] }, { "event": "Weekly data usage report shows a smart meter exceeded 600MB in a single day", "matching\_subscription\_ids": [2] }, { "event": "Firmware update failed on smart thermostats in multiple residential buildings" "matching\_subscription\_ids": [3] } { "event":



36

#### ABHISHEK KUMAR, LAURI LOVÉN, SUSANNA PIRTTIKANGAS AND SASU TARKOMA. "DATA FABRIC FOR INDUSTRIAL METAVERSE." ACCEPTED TO 44TH IEEE INTERNATIONAL CONFERENCE ON DISTRIBUTED COMPUTING SYSTEMS (ICDCS 2024), NEW JERSEY, US.





Saleh, Alaa, Praveen Kumar Donta, Roberto Morabito, Naser Hossein Motlagh, and Lauri Lovén. "Follow-Me AI: Energy-Efficient User Interaction with Smart Environments." arXiv preprint arXiv:2404.12486 (2024).

## CONTENTS

> BACKGROUND
> TECHNOLOGY MIX: 6G AND AI
> AI INTERCONNECT AND NEURAL PUB/SUB
> EDUCATION: 6G MOOC

## WE NEED 6G EDUCATION AND SKILLS

Upskilling Europe and the world. Explaining 5G and beyond. Raising interest in technology. For creating a safer, more sustainable and connected planet.

## **Core 5G and Beyond**

In this University of Helsinki open online course, you learn the fundamental concepts behind mobile networks, focusing in particular on 5G and the envisioned future 6G networks.

START COURSE

0



5

MINISTRY OF EDUCATION AND CULTURE FINLAND

0





University of Helsinki together with 6G Finland invites everyone to learn about 5G and 6G.

## **6G.MOOC.FI**



MINISTRY OF EDUCATION AND CULTURE FINLAND

