



CloudEdge Fusion Project

Jason H. HAGA (on behalf of the project members)

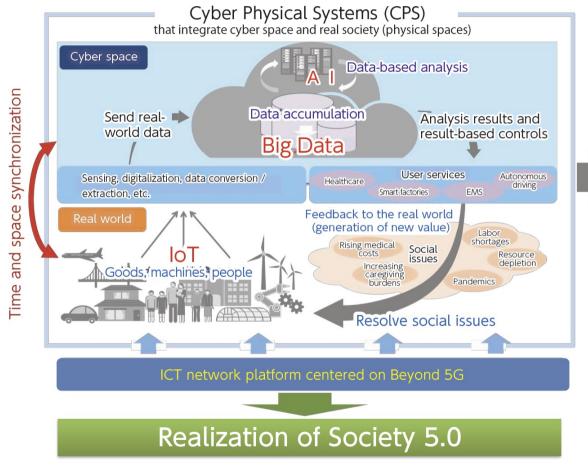
Chief Senior Research Scientist,

Digital Architecture Research Center,

National Institute of Advanced Industrial Science and Technology (AIST)

Cyber Physical Systems / Society 5.0

Figure 4-1-1-1 Vision of society expected in the 2030s



Vision of Society in the 2030s A resilient and vigorous society Inclusive A society in which everyone can be active in all places, with the removal of barriers and gaps between urban and rural areas, national borders, age, and physical abilities Sustainable A society that is convenient and achieves sustainable growth without loss of sociability Dependable A human-centered society that assures safety and security, even in the event of contingencies, with unwavering bonds of trust

In the 2030s, cyber space and physical spaces will become even more tightly integrated.

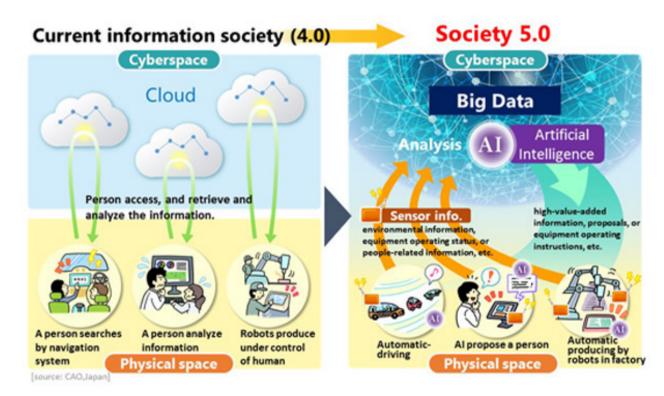
Moreover, a resilient and vigorous society will emerge in which cyber space not only extends the functions of physical spaces but also maintains the smooth functioning of people's lives and economic activities even when contingencies occur in physical spaces

Information and Communications in Japan 2020 white paper Ministry of Internal Affairs and Communications, Japan



How Society 5.0 works

- A huge amount of information from sensors (Real world data; RWD) in physical space is accumulated in cyberspace.
- In cyberspace, analysis results are fed back to humans in physical space in various forms.
- Different applications have different requirements for analysis processing capacity, latency, etc.



Source: https://www8.cao.go.jp/cstp/society5_0/



CPS and **IT** infrastructure

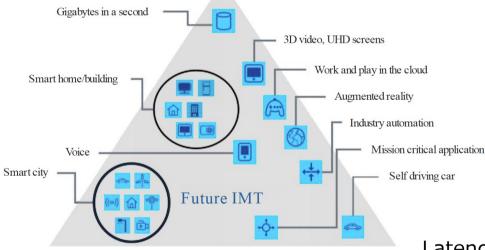
- In some applications, data from the physical world (Real World Data: RWD) must be processed within a certain time frame
 - Autonomous driving
 - Factory control
 - VR/AR
- Requirements for IT infrastructure
 - Reliable low-latency processing
 - Can process data from a huge number of IoT devices
 - Processing and communication throughput
 - Security and privacy



5G/Post-5G Communication

Anticipated Post-5G technologies promise to simultaneously provide low latency and high reliability, massive connections, large capacity, and security.

Enhanced mobile broadband (eMBB)



Latency <= 1ms (wireless section)

Massive machine type communications (mMTC) Ultra-reliable and low latency communications (URLCC)

Recommendation ITU-R M.2083-0



All we need is 5G/P5G?

- Cyberspace is enabled by networking and computing
- In addition to wireless networks, wired networks and computing must be able to achieve:
 - Low latency
 - High reliability
 - Support a huge number of connections
 - Security.
- Edge computing at the network edge is one solution, but…
 - Computational resources at the network edge are small, so statistical multiplexing effects cannot be expected
 - Lower utilization, higher cost (CAPEX, OPEX)
 - Difficulties in sharing data in wide-area.



Low latency and availability required for 5G/post-5G use cases

Use case		Latency (ms)	Reliability	Availability (%)	Device density	Traffic density	User throughput	Mobility (km/h)
Smart Factory [3], [4]	Manufacturing cell [3]	5	10-9	> 99.9999	0.33-3/m ²	N/S ²	N/S	<30
	Machine tools [3]	0.5	10-9	> 99.9999	$0.33-3/m^2$	N/S	N/S	<30
	Printing machines [3]	2	10-9	> 99.9999	$0.33-3/m^2$	N/S	N/S	<30
	Packaging machines [3]	1	10-9	> 99.9999	0.33-3/m ²	N/S	N/S	<30
	Cooperative motion control [4]	1	10-9	> 99.9999	0.33-3/m ²	N/S	N/S	<30
	Video-operated remote control [4]	10-100	10-9	> 99.9999	0.33-3/m ²	N/S	N/S	<30
	Assembly robots or milling machines [4]	4-8	10 ⁻⁹	> 99.9999	0.33-3/m ²	N/S	N/S	<30
	Mobile cranes [4]	12	10-9	> 99.9999	0.33-3/m ²	N/S	N/S	<30
	Process automation - Monitoring [81]	50	10^{-3}	99.9	10000/plant	10 Gbps/km ²	I Mbps	<5
	Process automation - Remote control [81]	50	10 ⁻⁵	99.999	1000/km ²	100 Gbps/km ²	<100 Mbps	<5
Smart Grids	Electricity distribution - Medium Voltage [81]	25	10^{-3}	99.9	1000/km ²	10 Gbps/km ²	10 Mbps	0
	Electricity distribution - High Voltage [81]	5	10-6	99.9999	1000/km ²	100 Gbps/km ²	10 Mbps	0
Smart Vehicle [81]	Autonomous driving [84]	5	10 ⁻⁵	99.999	500-3000/km ²	N/S	0.1-29 Mbps	urban < 10 highway< 50
	Collision warning [87]	10	$10^{-3} - 10^{-5}$	99.999	500-3000/km ²	N/S	0.1-29 Mbps	urban < 10 highway< 50
	High-speed train [81]	10	N/S	N/S	1000/train	12.5-25 Gbps/train	25-50 Mbps	<500
IIS [49]	Road safety urban [49]	10-100	$10^{-3} - 10^{-5}$	99.999	3000/km ²	10 Gbps/km ²	10 Mbps	<100
	Road safety highway 49	10-100	$10^{-3} - 10^{-5}$	99.999	500/km ²	10 Gbps/km ²	10 Mbps	<500
	Urban intersection [49]	<100	10-5	99.999	3000/km ²	10 Gbps/km ²	10 Mbps	<50
	Traffic efficiency [49]	<100	10^{-3}	99.9	3000/km ²	10 Gbps/km ²	10 Mbps	<500
	Traffic jam [85]	8	N/S	95.0	N/S	480 Gbps/km ²	20-100 Mbps	N/S
0	Large outdoor event [85]	N/S	10-2	99.0	4/m ²	900 Gbps/km ²	30 Mbps	N/S

Shopping mall [8:

Media on demand [85]

Stadium [8

 10^{-2}

10-2

N/S

N/S

99.0

99.0

N/S

4/m²

200000/km

4000/km²

N/S

0.1-10

Mbps/m²

700 Gbps/km2

60-300 Mbps

0.3-20 Mbps

60-300 Mbps

15 Mbps

N/S

N/S

TABLE IV: The summary of 5G requirements for CAV use cases [3], [4], [49], [80]-[87].

RWD applications require large capacity, many simultaneous connections, high reliability and low latency, but the requirements vary by application

Data volume, type and amount of processing required, processing time (latency) requirements, etc.

For automated driving and smart factory use cases, low latency (<100ms) and high availability (>99.9%) must be met.

(source) T.M.Ho, et. al., Next-generation Wireless Solutions for the Smart Factory, Smart Vehicles, the Smart Grid and Smart Cities, arXiv:1907.10102

How small can the latency of a wired network be?

• SINET: An academic research network operated by the National Institute of Informatics (NII)

 A mesh network connecting prefectures with a bandwidth of 100 Gbps or more (Upgraded to 400 Gbps in SINET6)

• Interconnection including local areas is realized with short transmission lines

SINET5 (\sim 2022) Topology



Los Angeles

New York

(As of June 1, 2021)

	Participation (%)
National Universities	86 (100%)
Public Universi	ties 90 (96%)
Private Universities	430 (70%)
Junior College	es 84 (26%)
Technology Colleges	56 (98%)
Inter-Universit Research Institu Corporation	
Other	217
Total	979

Courtesy of NII

: SINET DC : Domestic line (400Gbps) : Domestic line (100Gbps) : International line (100Gbps)



What we can learn from latency on SINET5

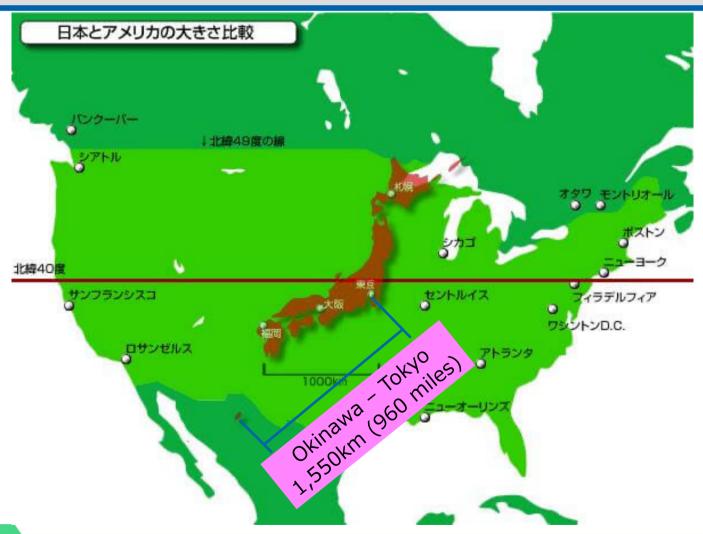
		lwate	Kochi	Mie	Oita	Osaka	Tokyo	Toyama	Yamaguchi	Okinawa
	←one way latency(msec)	6.65	15.50	13.00	16.70	12.40	9.70	9.97	15.50	22.20
	→one way latency(msec)	6.56	15.50	13.00	16.60	12.50	9.77	9.93	15.50	22.10
Kitami	distance (km)	508.00	1443.00	1189.00	1578.00	1238.00	970.00	967.00	1505.00	2436.00
	est. fiber length/distance	2.60	2.15	2.19	2.11	2.01	2.01	2.06	2.06	1.82
	←one way latency(msec)		9.47	6.94	11.30	7.03	3.63	5.07	10.60	17.30
Iwate	→one way latency(msec)		9.46	6.94	11.30	7.12	3.74	5.07	10.60	17.30
	distance (km)		963.40	688.60	1115.00	748.10	463.80	480.20	1056.60	1951.00
	est. fiber length/distance		1.96	2.02	2.03	1.89	1.59	2.11	2.01	1.77
	←one way latency(msec)			3.31	2.04	2.10	5.51	4.67	2.67	9.25
Kochi	→one way latency(msec)			3.32	2.00	2.18	5.61	4.65	2.68	9.23
	distance (km)			303.70	182.00	222.00	612.20	483.20	202.80	990.80
	est. fiber length/distance			2.18	2.22	1.93	1.82	1.93	2.64	1.87
	←one way latency(msec)				5.18	1.45	3.39	2.56	5.03	11.40
Mie	→one way latency(msec)				5.14	1.52	3.48	2.55	5.02	11.40
Ivile	distance (km)				481.70	90.70	308.70	227.10	466.70	1267.70
	est. fiber length/distance				2.14	3.27	2.23	2.25	2.15	1.80
	←one way latency(msec)					3.94	7.35	6.51	1.50	7.43
Oita	→one way latency(msec)					4.05	7.48	6.52	1.55	7.43
	distance (km)					395.20	790.20	639.00	105.90	866.50
	est. fiber length/distance					2.02	1.88	2.04	2.88	1.71
	←one way latency(msec)						3.59	2.74	3.75	7.40
Osaka	→one way latency(msec)						3.60	2.64	3.67	10.00
Osaka	distance (km)						395.90	270.40	376.20	1202.90
	est. fiber length/distance						1.82	1.99	1.97	1.45
	←one way latency(msec)							3.45	7.44	10.10
Tokyo	→one way latency(msec)							3.31	7.34	13.70
TORYO	distance (km)							249.40	769.00	1553.60
	est. fiber length/distance							2.71	1.92	1.53
	←one way latency(msec)								5.74	13.80
Toyama	→one way latency(msec)								5.75	12.40
	distance (km)								590.70	1471.80
	est. fiber length/distance								1.95	1.78
	←one way latency(msec)									6.85
Yamaguchi	→one way latency(msec)									6.83
	distance (km)									956.10
	est. fiber length/distance									1.43

- The fiber length can be estimated to be roughly twice as long as the distance between the cities
 - The delay of fiber is about 5µs/km (5ms/1000km)
- The latency of the wired network can be constrained if a SINET-like topology is used.
 - Tokyo-Kitami 970km: RTT 19.47ms
 - Tokyo-Oita 790km: RTT 14.83ms
 - Tokyo-Okinawa 1,550km: RTT 23.80ms

The latency between SINET DCs Source: NII



Japan is not so small



Resources to support CPS

Resource	type	Guaranteed latency and throughput	Deployment and maintenance cost		
Network	Wide area network	✓ Can support if properly configured	-		
	5G / P5G	√ Will support	-		
Compute	Cloud / Data centers	?	√ Relatively low		
	Network edge	✓ Can support by using dedicated resources for each application	x High		

Challenges in computing

- Conventional technologies prioritize processing efficiency, often focus on throughput and give little consideration to latency,
 - Memory hierarchy, batch processing, etc.
- Cloud emphasizes cost and throughput
- Most PoC of 5G low latency use cases uses dedicated compute infrastructure
 - Major obstacles to practical deployment of services
- From the perspective of network-connected computing, research and development from the perspective of latency and reliable processing is lacking
 - Technology to appropriately combine and utilize geographically dispersed computing resources is needed, especially in terms of latency and security.

CloudEdge Fusion (CEF) Project

- We will develop technologies to support the realization of a cyberinfrastructure across the cloud-edge continuum to provide optimal processing power (latency, bandwidth, security) in response to service requirements.
- The developed technologies will be integrated as a system towards commercialization and the effectiveness will be verified through practical demonstrations.
- NEDO will fund 77M USD over 5 years (until Mar. 2028)

CEF Project Organization







Information Science



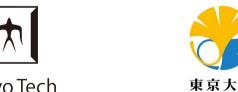


Date Labs













Tatebe Labs, Amagasa Labs



Sakamoto Labs, Endo Labs



Nakamura Labs, Sekiya Labs





Challenges in cloud-edge continuum platform

Time-sensitive service execution environment provisioning

• Short end-to-end connections enabled in part by geo-positioning resources, so that worst-case latencies in the data pipeline do not exceed service request response requirements.

Declarative application deployment

- Even if the service provider does not know the details of the infrastructure, applications will be deployed in the right places according to the characteristics.
- A mechanism is needed to predict and adjust the situation in a timely manner to prevent service interruptions or other SLA violations.

Secure service federation

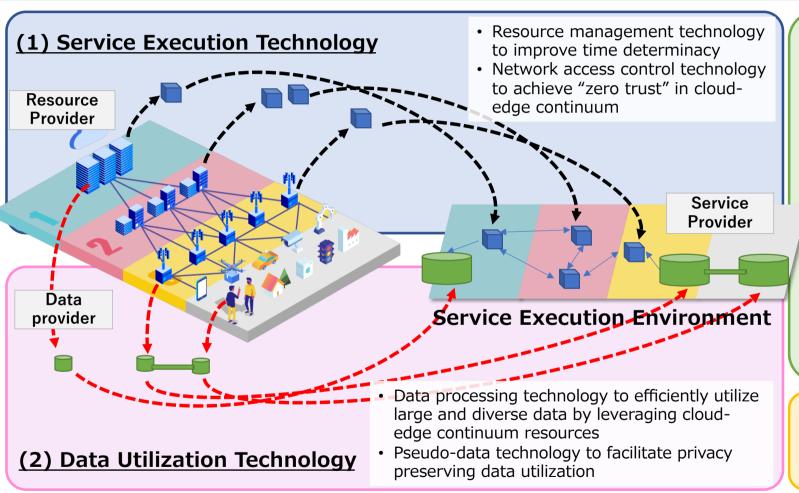
 An application can communicate securely with internal/external services, depending on the context and the data it uses.

> CloudEdge Fusion

Project items of CEF

- Project Item (1) Computing continuum service execution infrastructure technology
 - Reduce request processing time fluctuation in one order of magnitude or more compared to existing technology (about 100 ms).
- Project item (2) Fundamental technologies for continuum data
 - Establish elemental technologies for data processing, search, pseudoization, storage management, etc.
- Project item (3) System Integration
 - integrate the technologies researched and developed in project items (1) and (2) with existing technologies as a system.
 - Confirm that the functionality of the computing continuum infrastructure technology satisfies the design and that the performance is sufficient to realize services that require an end-to-end response time of about 100 ms.
- Project Item (4): System Demonstration
- Project item (5) Promotion of commercialization

CEF Project Working Packages



(3) System Integration Technology

System integration of the R&D results from (1) and (2) with existing software

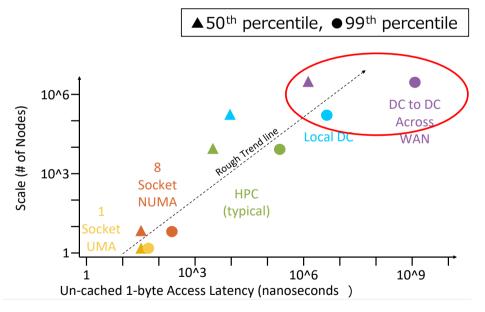
(4) System Demonstra tion

Demonstrati on of the proposed platform through industrial application services

(5) Social Implementation and Standardization

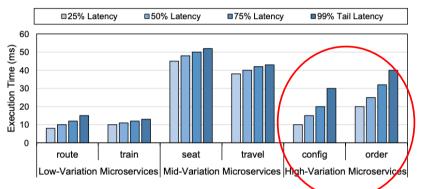
Performance fluctuations in distributed computing

- Real-world data processing requires that the response time to be within the requirements.
- In general, the larger the system, the larger the performance fluctuation.



(出典) Rob Sherwood (Intel) [Exacomm2022]

 The execution time of microservices fluctuates by several tens of milliseconds.



(出典) X.Wang, et al., Exploring Efficient Microservice Level Parallelism [IPDPS2022]

Real service consists of several dozen microservices.

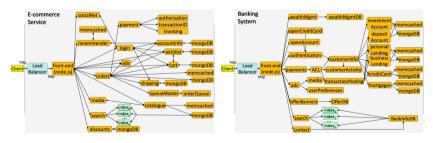


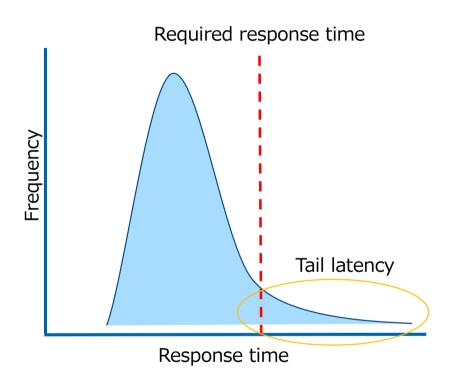
Figure 6. The architecture of the E-commerce service. Figure 7. The architecture of the Banking end-to-end service.

(出典) U.Gan, et al., An Open-Source Benchmark Suite for Microservices… [ASPLOS2019]

18

Probabilistic Service Reliability

- Ideally, processing latency and throughput should be fixed, but in reality, this is difficult to achieve in a shared computing infrastructure.
- The response time in distributed system is known to follow "tail latency," caused by various factors:
 - Contention for cache, memory and storage access,
 - Process scheduling
 - Network congestion
- We introduce the concept of "**probabilistic** service reliability" to satisfy trade-offs between performance and economic cost.
 - A platform probabilistically guarantees service levels based on the assumption that there is a certain amount of fluctuation in processing capacity.

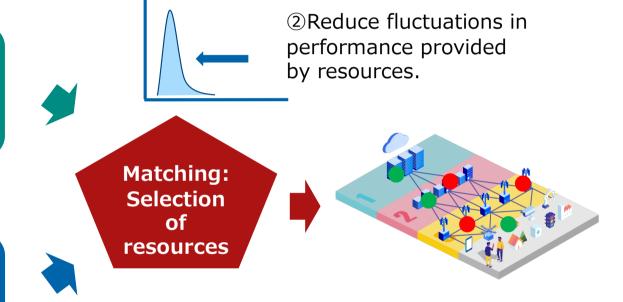


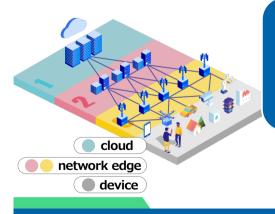
Matching and resource allocation

③ Combine resources to provide a service execution environment that meets application requirements.



Application requirements (Manifestos)

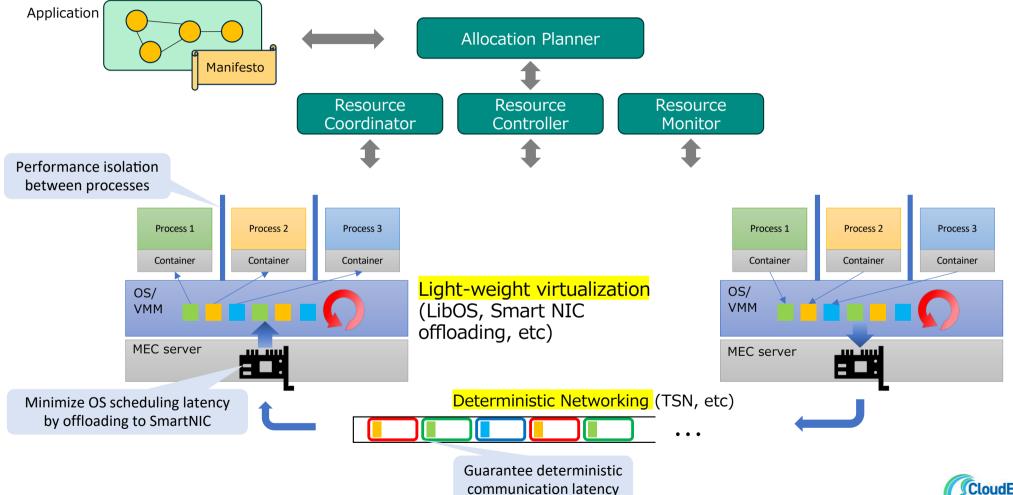




Performance of various infrastructure

①Parameterize the fluctuations in performance provided by resources and define the quality of service that can be provided.

Highly time-sensitive service execution mechanism



Summary

- In a smart society where cyber space and physical space converge, <u>a cyberinfrastructure that satisfies the requirements for processing</u> <u>capacity (throughput) and processing latency (latency)</u> for each service is essential.
- Post-5G communications are designed with these requirements, but today's computing infrastructures such as cloud computing are not designed to address them.
- <u>The CEF project</u> aims to satisfy these requirements and support the realization of a smart society <u>by developing a cyber-infrastructure</u> access cloud-edge continuum to provide optimal processing power in response to service requirements.

Thank you for your attention!



ポスト5G情報通信システム 基盤強化研究開発事業

This presentation is based on results obtained from the project "Research and Development Project of the Enhanced Infrastructures for Post-5G Information and Communication System" (JPNP20017), commissioned by the New Energy and Industrial Technology Development Organization (NEDO).