

ホワイトペーパー：  
2030年代に向けた  
コンピューティング  
インフラストラクチャー

LF Edge と IOWN GF DCI  
プラットフォームの融合

1st edition - April 2024  
[www.akraino.org](http://www.akraino.org)

# 目次

---

|  |    |
|--|----|
| 1 概要.....                              | 3  |
| 2 背景.....                              | 4  |
| 3 課題解決への取組.....                        | 6  |
| 4 IOWN GF が提案するインフラストラクチャー .....       | 7  |
| 5 IOWN GF/LF Edge ジョイント PoC について ..... | 9  |
| 6 まとめ .....                            | 10 |
| 7 略語.....                              | 10 |
| 参考文献.....                              | 11 |
| 著者.....                                | 11 |

## 図

|  |   |
|--|---|
| 図 1: AI モデルに必要な計算性能 .....                  | 4 |
| 図 2: アプリケーションのパフォーマンス要件 .....              | 5 |
| 図 3: IOWN GF/LF EDGE ジョイント POC コンセプト ..... | 6 |
| 図 4: ディスアグリゲーション .....                     | 7 |
| 図 5: アクセラレータ デバイス間の直接転送 .....              | 8 |
| 図 6: データ セントリック インフラストラクチャー (DCI) .....    | 8 |
| 図 7: IOWN GF/LF EDGE ジョイント POC の詳細.....    | 9 |

## 本訳文について

この日本語文書は、**Computing Infrastructure into the 2030s** の参考訳として、The Linux Foundation Japan が便宜上提供するものです。英語版と翻訳版の間で齟齬または矛盾がある場合（翻訳版の提供の遅滞による場合を含むがこれに限らない）、英語版が優先されます。

この日本語文書を引用する際には、下記の一文を記載してください。

引用：Computing Infrastructure into the 2030s 参考訳（The Linux Foundation Japan 提供）

翻訳協力：天満尚二

# 1 概要

エッジコンピューティングは、データからリアルタイムに意思決定を行うデータ駆動型社会にとって重要な技術です。なぜならば、エッジコンピューティングは、データの発生元に近いエッジが処理することによりレイテンシを削減できるからです。Linux Foundation のプロジェクトである LF Edge は、エッジコンピューティングの実用化に向けた課題となっているハードウェア、プロセッサ、クラウド、オペレーティングシステムに依存しないオープンで相互運用可能なエッジコンピューティングフレームワークの確立に取り組んでいます。また、LF Edge は、ファクトリーオートメーションや自動運転など、膨大な量のデータをリアルタイムに AI 処理するニーズの高まりを受け、エッジ AI の実現にも注力しています。しかしながら、データ量の増加、AI モデルの計算複雑さの増大、レイテンシやメモリ容量などのアプリケーション固有の要件への対応、エネルギー効率と柔軟性を兼ね備えたインフラストラクチャーの実現には課題があります。だからこそ、インフラストラクチャー技術のイノベーションが必要なのです。そこで、The Linux Foundation と IOWN GF (Innovative Optical and Wireless Network Global Forum) は、IOWN GF が提案するインフラストラクチャー上に Linux Foundation のソフトウェアを統合し、性能向上、レイテンシ低減、エネルギー効率向上を実現する共通のインフラストラクチャーを開発する基本合意書を 2023 年 6 月に締結しました。この合意に基づき、両技術の融合と性能向上を実証する IOWN GF/LF Edge 共同 PoC (Proof of Concept) を計画しました。本 PoC では、IOWN GF のインフラストラクチャーと LF Edge のプラットフォームおよびソフトウェアを用いて、エッジからクラウドまでのエンドツーエンドの環境を構築し、デモ用の環境で実際のアプリケーションを動作させます。本稿では、IOWN GF/LF Edge のジョイント PoC の内容について述べます。

## 2 背景

LF Edge は、オープンで相互運用可能なエッジ コンピューティング フレームワークを確立することに加えて、ファクトリーオートメーションや自動運転など、大量のデータのリアルタイム AI 処理のニーズの高まりに対して、エッジで AI を実現することにも重点を置いています。

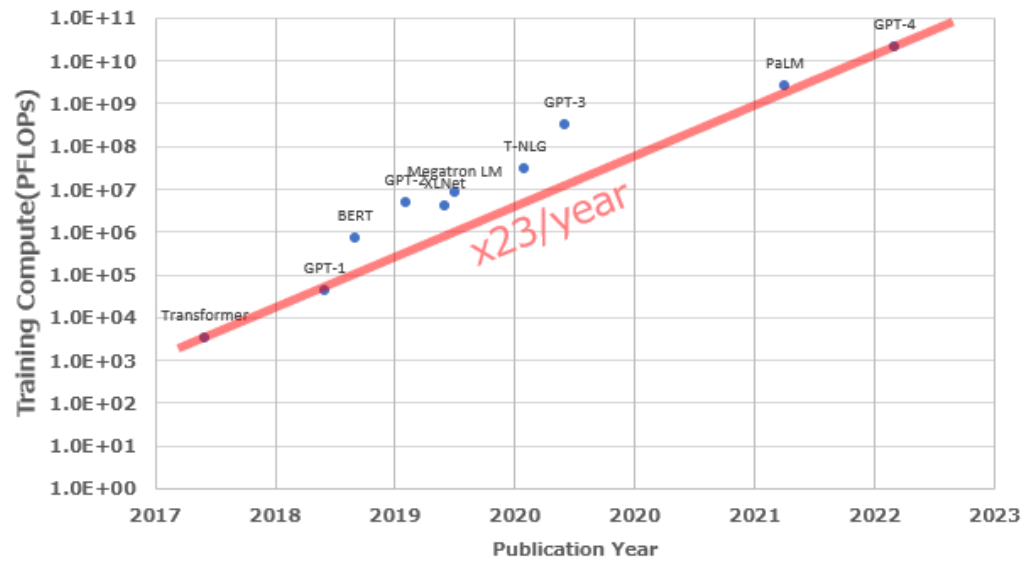


図 1: AI モデルに必要な計算性能

AI 分野では、図 1 に示すように、学習データ量の増加や生成 AI モデルの規模拡大により、コンピューティング インフラストラクチャーに求められる計算性能が急速に増加しています<sup>[1][2]</sup>。これは、AI の計算処理を担う CPU や GPU の性能進化である「ムーアの法則」を超えています。コンピューティング インフラストラクチャーは、CPU と GPU を備えたサーバーを並列に増やすことにより、この速度を追従しています。しかし、データセンターの電力消費量は 2030 年には現在の約 15 倍となり、電力不足につながります<sup>[3][4][5]</sup>。

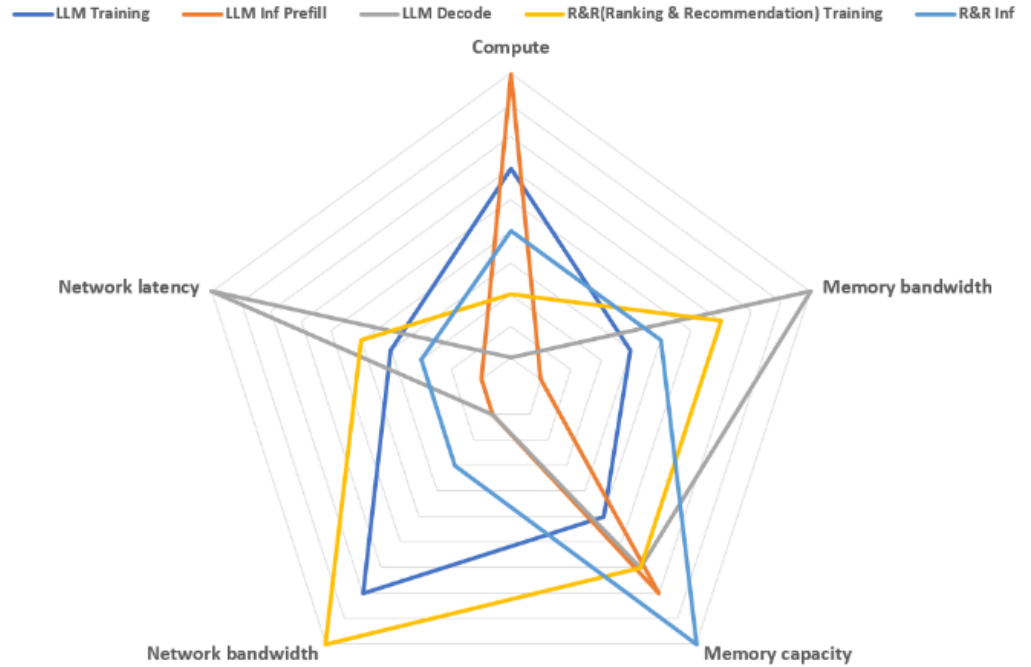


図 2: AI アプリケーションのパフォーマンス要件

また、図 2 に示すように、AI のアプリケーションごとに求められる性能は多様化しています<sup>[6]</sup>。たとえば、大規模言語モデル推論のプレフィル (LLM Inf prefill) は、メモリ帯域幅よりも計算パフォーマンスを重視しますが、大規模言語モデル推論のデコード (LLM Decode) はその逆です。現在、インフラストラクチャーは、コンピューティング、メモリ、およびネットワークリソースの中でアプリケーションが最も重視するパフォーマンスに基づいて、サーバーごとに構築されています。そのため、一部のリソースが十分に活用されず、インフラストラクチャーが拡張されています。

したがって、AI アプリケーションの規模と多様性の拡大を持続的に実現するためには、より電力効率が高く、よりリソースの割り当てが柔軟なコンピューティング インフラストラクチャーが必要です。この課題は、エッジのようなコンピューティング インフラストラクチャーに電力やスペースなどの制限がある環境では、より重要になるはずです。私たちは、コンピューティング インフラストラクチャーのアーキテクチャの変更が有効な解決策の一つと考え、試みました。

### 3 課題解決への取組

IOWN GF (Innovative Optical and Wireless Network Global Forum) では、光技術などの革新的技術を活用して、既存のインフラの限界を超える高速・大容量の通信と膨大な計算リソースを提供できるネットワークと情報処理インフラストラクチャー DCI (Data Centric Infrastructure) を提案しています<sup>[7]</sup>。Linux Foundation と IOWN GF (Innovative Optical and Wireless Network Global Forum) は、2023 年 6 月に IOWN GF が提案するインフラストラクチャーに Linux Foundation のソフトウェアを統合し、性能、レイテンシ、エネルギー効率を向上させる共通のインフラストラクチャーを開発するための基本合意書に署名しました<sup>[8]</sup>。LF Edge は、DCI をエッジ コンピューティングにおけるこれらの課題に対する可能性の高い解決策の 1 つと見なしています。そこで、DCI による両技術の融合と性能向上を実証する IOWN GF/LF Edge ジョイント PoC (Proof of Concept) を企画しました。本 PoC では、IOWN GF のインフラストラクチャーと、LF Edge のプラットフォーム、ソフトウェアを用いて、エッジからクラウドまでのエンドツーエンドの環境を構築し、デモ用環境において実際のアプリケーションを実行します。下図は、ジョイント PoC の概要を示しています。

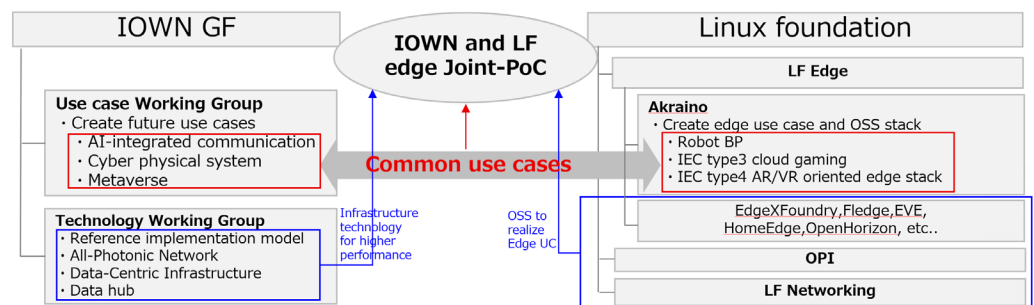


図 3: IOWN GF/LF EDGE ジョイント POC コンセプト

# 4 IOWN GF が提案するインフラストラクチャー

IOWN GF は、新たな効率的で柔軟なコンピューティング インフラストラクチャーとして、Data Centric Infrastructure (DCI) を提唱しています。現在のコンピューティング インフラストラクチャーと比較して、2つの主要な機能があります。

1つ目は、ディスクアグリゲーションです。これにより、デバイスの柔軟な再配置が可能になります。たとえば、次の図に示すように、各アプリケーションが異なるリソースを必要とする場合、現在のコンピューティング アーキテクチャである“Server-oriented”では、サーバー毎に各アプリケーションにリソースを割り当てます。そのため、一部のデバイスが未使用のままになり、利用効率が低下することがあります。一方、ディスクアグリゲーションは、アプリケーションに応じてリソース プールから必要なデバイスを選択し、PCIe スイッチを介して論理ノードを構成することで、さまざまなアプリケーションの要件を効率的に満たすことができます。

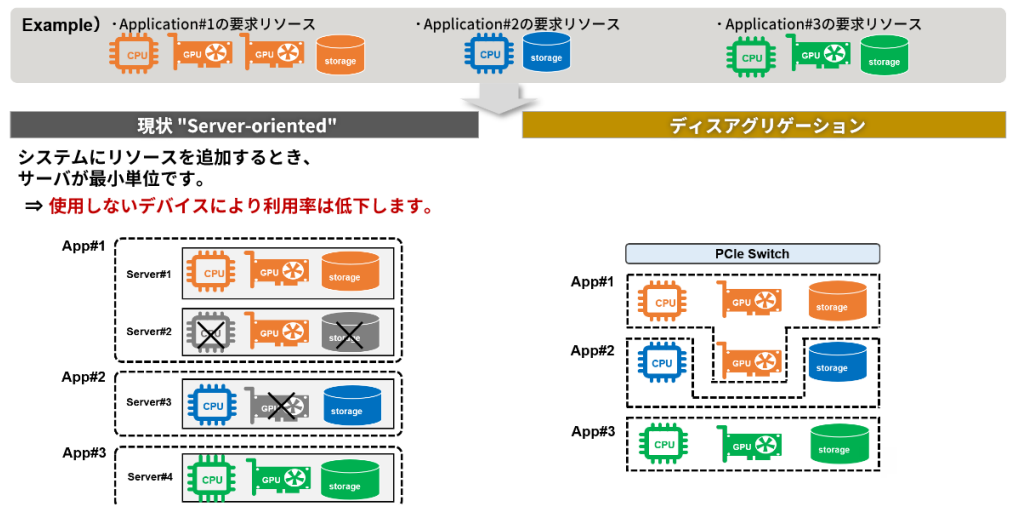


図 4: ディスクアグリゲーション

2つ目は、光接続または共有メモリを介したアクセラレータデバイス間の直接接続です。現在のコンピューティング アーキテクチャは CPU 中心で、データは CPU メモリを介して CPU によってアクセラレータ デバイス間を転送されます。一方、直接接続の場合は、共有メモリや光スイッチを介してアクセラレータ デバイス間で直接データが転送されます。これにより、データ転送に使用される CPU コアの使用率が削減され、コストと消費電力が削減されます。処理遅延も低減できます。

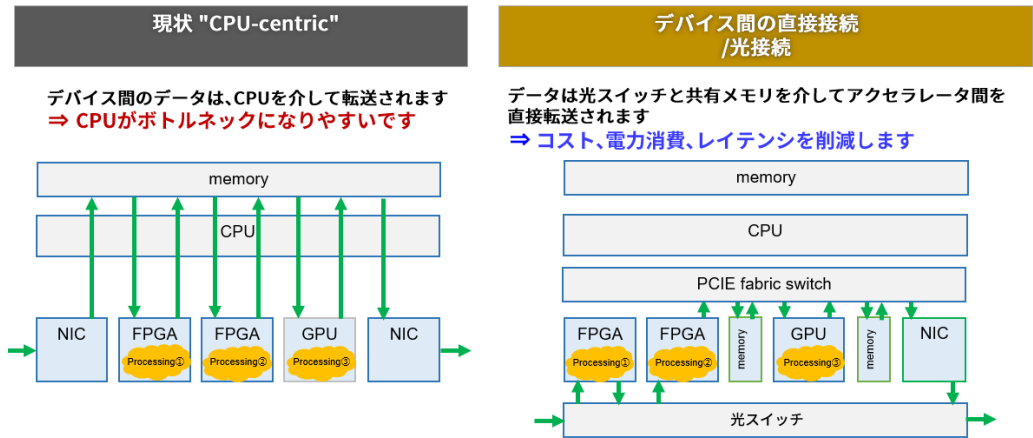


図 5: アクセラレータデバイス間の直接転送

DCI の概要と仕組みを以下の図に示します。まず、ユーザーがアプリケーションの実行に必要なリソースを DCI に要求し (1.Request)、DCI のコントローラーは、リソース プールと論理ノードを管理する Composable Disaggregated Infrastructure (CDI) 管理ソフトウェアに論理ノードの構成を指示することによって応答します (2.Allocate)。また、コントローラーは、アクセラレータ デバイス間の直接接続を担当する光スイッチの構成も行います。

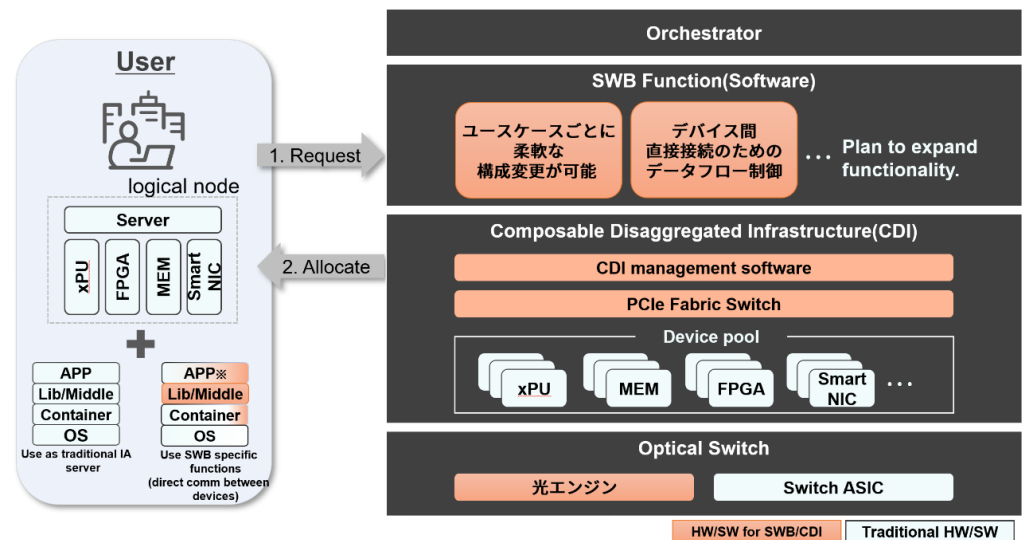


図 6: データ セントリック インフラストラクチャー (DCI)



# 5 IOWN GF/LF Edge ジョイント PoC について

この章では、IOWN GF/LF Edge ジョイント PoC の詳細について説明します。PoC では、以下の2点を実証します。1点目は、IOWN GF と LF Edge の融合です。LF Edge のプラットフォームとソフトウェアが DCI 上に統合できることを確認し、実際の環境構築を通じて現在のサーバーアーキテクチャから DCI に変更することにより、DCI をリソースとして管理・制御するために必要なメカニズムを抽出します。2点目は、DCI を適用することで処理性能が向上し、消費電力が削減されることを実証します。このデモでは、下図のように、エッジ、クラウドにエンドツーエンドにアクセスするアプリケーションとして、エッジ AI の代表的な処理である AI による動画推論を実装します。アクセスエリアから映像が入力され、5G RAN やコアネットワークを介してエッジクラウドに送信され、AI による推論処理が行われます。処理されたビデオはクラウドに送信されます。ここでは、エッジクラウドのコンピューティングリソースに、市販のコンポーザブル ディスアグリゲーション製品を使用して DCI をエミュレートして使用します。AI の推論処理は、デコード、フィルタリング / リサイジング、推論の3つのフェーズに分かれており、FPGA や GPU などのアクセラレータ デバイスに実装され、DCI の特徴の一つであるアクセラレータ デバイス間の直接接続によりデータの送受信が行われます。この実証においては、現在のサーバーアーキテクチャをコンピューティングリソースとして使用する場合と比較して、DCI を適用することによる処理性能の向上と消費電力の削減を測定します。加えて、エッジクラウド環境を構築するためのミドルウェアとして、LF Edge Akraino IEC (Integrated Edge Cloud) type2 Blueprint を DCI 上に実装し、DCI と LF Edge エッジコンピューティングプラットフォームが統合可能であることを検証します。

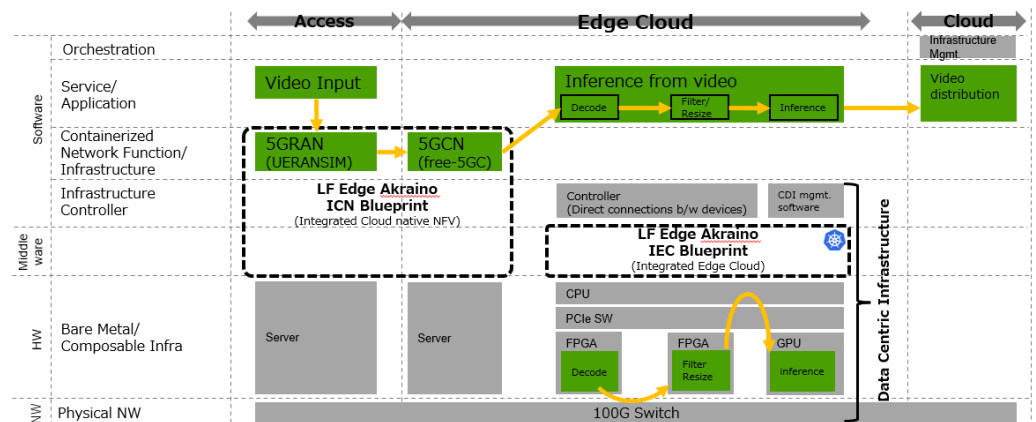


図 7: IOWN GF/LF EDGE ジョイント POC の詳細

## 6 まとめ

上記エンドツーエンド環境を今後の実証のベースとして構築し、2024年4月に本稿の第2弾としてご報告します。第2弾では、PoCの最初の実証ポイントであるIOWN GFが提案するDCIとLF Edgeのプラットフォーム、ソフトウェアの統合結果をレポートします。そして、引き続き、もう一つの実証ポイントであるDCI適用による処理性能の向上と消費電力の削減を、DCI適用前と適用後の比較で実証していきます。さらに、冒頭で述べたコンピューティングの課題のうち柔軟性についても、今後、実証実験の内容について検討していく予定です。例えば、IOWNが提案する光とディスアグリゲーションによってエッジサイトを直接つなぐAll-Photonic Network (APN)は、離れたエッジサイトのリソースを組み合わせ、より柔軟に論理ノードを構築することができます。このような例により、LF Edgeのプラットフォームとソフトウェアをどのように使用できるかを引き続き実証していきます。また、この活動をIOWN GF PoCレポートとして公開することも検討したいと考えています。これらの活動により、IOWN GFとLF Edgeの協業による技術の進化とシナジーの実現を加速し、産業の発展に貢献してまいります。

## 7 略語

| 用語      | 説明   |
|---------|--|
| IOWN GF | Innovative Optical and Wireless Network Global Forum |
| CPU     | Central Processing Unit                              |
| GPU     | Graphics Processing Unit                             |
| DCI     | Data-Centric Infrastructure                          |
| PCIe    | Peripheral Component Interconnect-Express            |
| FPGA    | Field Programmable Gate Array                        |

## 参考文献

- [1] A. Gholami, “AI and Memory Wall,” 30 3 2021. [Online]. Available: <https://medium.com/riselab/ai-and-memory-wall-2cb4265cb0b8>.
- [2] S. McAleese, “Retrospective on ‘GPT-4 Predictions’ After the Release of GPT-4,” 18 3 2023. [Online]. Available: <https://www.lesswrong.com/posts/iQx2eeHKLwgBYdWPZ/retrospective-on-gpt-4-predictions-after-the-release-of-gpt>.
- [3] S. R. a. F. Boshell, “The nexus between data centres, efficiency and renewables: a role model for the energy transition,” 26 6 2020. [Online]. Available: <https://energypost.eu/the-nexus-between-data-centres-efficiency-and-renewables-a-role-model-for-the-energy-transition/>.
- [4] Japan Science and Technology Agency, “Impact of Progress of Information Society on Energy Consumption,” 2 2020. [Online]. Available: <https://www.jst.go.jp/lcs/pdf/fy2020-pp-03.pdf>.
- [5] N. Willing, “New Technologies Are Needed to Curb Data Center Energy Use, Says the IEA,” 3 8 2023. [Online]. Available: <https://www.techopedia.com/new-technologies-are-needed-to-curb-data-center-energy-use-says-the-iea>.
- [6] M. Omar Baldonado, “Meta’ s evolution of network for AI,” 17 10 2023. [Online]. Available: <https://www.youtube.com/watch?v=5gOOtFySrQA>.
- [7] IOWN GLOBAL FORUM, “DCI Product Concept Paper,” 19 10 2023. [Online]. Available: [https://iowngf.org/wp-content/uploads/formidable/21/IOWN-GF-RD-DCI\\_PCP-1.1.pdf](https://iowngf.org/wp-content/uploads/formidable/21/IOWN-GF-RD-DCI_PCP-1.1.pdf). [Accessed 9 2 2024].
- [8] IOWN GLOBAL FORUM, “Linux Foundation and IOWN Global Forum to Collaborate for Future Smart Connected World,” IOWN GLOBAL FORUM, 14 6 2023. [Online]. Available: <https://iowngf.org/press-releases/linux-foundation-and-iown-global-forum-to-collaborate-for-future-smart-connected-world/>. [Accessed 9 2 2024].

## 著者

Haruhisa Fukano (Fujitsu)  
Toshimichi Fukuda (Fujitsu)  
Reo Inoue (Fujitsu)

