

## In-Network Computing による遠隔 GPU リソースを活用した 低遅延 AI 映像解析の実証に成功

～6G 時代の AI・ロボットがその能力を最大限発揮するネットワークの実現に向けて前進～

発表のポイント:

- ◆ 分散配備された GPU リソースと 5G ネットワークを IOWN APN<sup>※1</sup> で接続する INC エッジ<sup>※2</sup> を実装し、通信の制御に加えて AI 推論処理<sup>※3</sup> をネットワーク側で制御する技術を確認しました。
- ◆ 上記技術の実証実験を実施し、遠隔 GPU リソースを活用した In-Network Computing<sup>※4</sup> によって、6G 時代の遠隔でのロボット制御に十分な低遅延を実現できる見通しをえました<sup>※4</sup>。

株式会社 NTTドコモ(以下、ドコモ)と NTT 株式会社(以下、NTT)は、分散配備された遠隔 GPU リソースと 5G ネットワークを IOWN APN で接続する INC エッジを活用した In-Network Computing<sup>※5</sup>(以下、INC)により、低遅延な AI 映像解析の実証実験に成功しました。

本実証実験では、5G コアネットワーク上に実装した INC エッジにより、通信の制御に加えて AI 推論処理をネットワーク側で制御しました。これにより、IOWN APN を介して接続された遠隔 GPU リソースを用い、端末から転送された映像データを低遅延に解析できることを確認し、6G 時代の遠隔でのロボット制御に十分な低遅延を実現できる見通しをえました。

なお本成果は、2026 年 3 月 2 日(月曜)から 5 日(木曜)にかけてスペイン・バルセロナで開催される GSMA 主催「Mobile World Congress Barcelona 2026<sup>※6</sup>」の NTT グループブースにて展示いたします。ドコモおよび NTT は、今後も機能が簡素化された端末の普及に向けて INC の技術検討・実証および国際標準化を推進し、6G 時代の AI・ロボットがその能力を最大限発揮するネットワークの実現をめざします。

### 1. 背景

6G 時代に向けて、没入型 XR や AI/ロボットを活用した新たなサービスの展開が進むと言われています。これらのサービスでは、従来に比べ大容量・低遅延のデータ転送や大規模なデータ処理を必要とする場合があります。例えば、ロボットが自律的に動作する場合において、ロボット周囲の映像やセンサーデータを取得し、AI を用いてロボットの移動先の障害物などを解析し、ロボット制御に即時にフィードバックすることなどが考えられます。特に、小型のロボットや簡素なウェアラブル端末などで AI における学習や推論を用いるアプリケーションを利用する場合、お客さま体感を落とすことなくサービスを提供するためには、端末以外の処理リソースにおいてもリアルタイムに大量のデータを処理する能力が求められ、6G 時代のネットワークには通信の処理だけでなく、サービスのデータ処理も含めた制御を実施し、品質を担保することが期待されています。

一方で、AI 推論処理の分散実行は、従来、アプリケーションやサーバ側で制御されることが一般的であり、ネットワークは主にデータ転送を担う役割にとどまっていた。そのため、推論処理に用いる GPU リソースの配置や通信遅延がサービス品質に大きく影響し、通信遅延の面で有利である地理的に近い場所にある計算リソースの利用が前提となるなど、柔軟なリソース活用には課題がありました。

このような期待と課題からドコモと NTT は、6G 時代のネットワークに必要な要素技術として、INC の研究開発を進めています。INC では、ネットワークの中に GPU をはじめとしたさまざまなリソースが分散配備され、通信だけではなくサービスの計算処理もネットワークで制御し、AI などのサービスを高品質で提供します。

## 2. 実証実験の概要

本実証実験では、ネットワーク内に分散配備された遠隔 GPU リソースと 5G ネットワークを INC エッジを用いて IOWN APN を介して接続し、5G ネットワークに接続された端末から送信された映像データの AI 推論処理の検証を実施しました。

一般に AI 推論処理を各リソースの処理負荷軽減のための分散実行するケースでは、GPU リソース間の通信遅延が推論処理全体の遅延に大きく影響するため、同一拠点内など地理的に近い場所にある GPU リソースの利用が前提とされています。本実証実験では、INC エッジと IOWN APN を活用し、通信の制御に加えて AI 推論処理をネットワーク側から制御することにより、地理的に離れた遠隔 GPU リソースを用いた場合でも、高い推論性能を維持できるかを検証しました。

本実証実験にあたっては、INC エッジとして、新たに IOWN APN とモバイル網の接続機能に加え、AI 推論処理を推論の前段にあたる処理と推論の実行部分に分け、前段処理後のデータを遠隔 GPU リソースへ低遅延に転送・振分けするための仕組みをネットワーク機能として実装しました。また、映像データの転送には、AWS 上に構築した商用 5G コアネットワークの優先制御機能を適用し、INC エッジの役割と組み合わせることで、5G ネットワークおよび IOWN APN を活用した広帯域・低遅延な AI 映像解析を実現しました。また、今回の実験において、通信と AI 映像解析の合計処理遅延は、人間の周囲でロボットが自律制御に基づいて動作する場合に想定される要求遅延と比較して、要求遅延以内であることを確認し<sup>※4</sup>、6G 時代の遠隔でのロボット制御に十分な低遅延を実現できる見通しをえました。

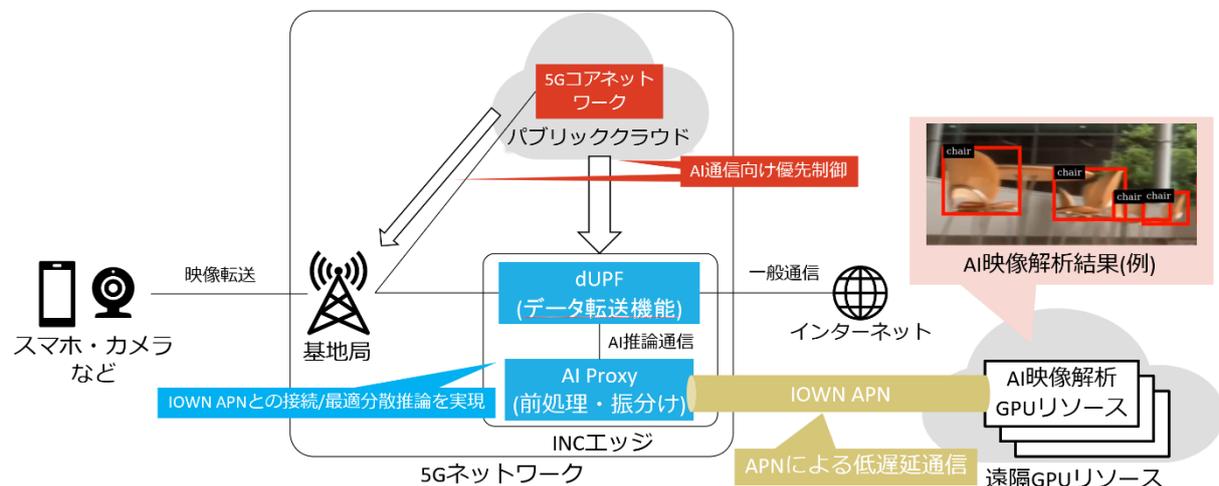


図 1. 実証実験のシステム構成

### 3. 各社の役割

#### ドコモ

- ・ 実証実験全体の計画策定、全体管理
- ・ コアネットワークや無線基地局装置などの 5G SA 商用環境およびノウハウの提供
- ・ 実証実験における IOWN APN の設計検討・構築
- ・ 実現方式の検討およびネットワーク構成の設計

#### NTT

- ・ INC 基盤の提供
- ・ 5G コアネットワークと INC を IOWN APN を介して接続・融合し分散推論を実現するエッジ機能 INC エッジの提供
- ・ 実現方式の検討およびネットワーク構成の設計

### 4. 今後の展開

本実証実験から得られた結果は、6G 時代の AI やロボット向けのデータ転送・処理にも応用できることが期待されます。ドコモ、NTT は今後も 6G ネットワークの要素技術として、機能が簡素化された端末の普及に向けて通信とデータ処理を包括的に提供する INC の技術検討・実証および国際標準化を推進していき、6G 時代の AI・ロボットがその価値を最大限発揮するネットワークの実現をめざします。

### 5. 関連する過去の報道発表

・2025 年 3 月 3 日「6G 時代の高機能サービスの利用に向け、ネットワークとサービスの連携によるコンピューティングサービスのオンデマンド一括制御の実証に成功 -In-Network Computing による 6G 時代の AI 活用に向けて前進-

([https://www.docomo.ne.jp/info/news\\_release/2025/03/03\\_01.html](https://www.docomo.ne.jp/info/news_release/2025/03/03_01.html))

・2026 年 3 月 2 日「国内で初めて AWS 上に構築した 5G コアの商用サービス展開を開始するとともに、世界で初めての AI を用いたコアネットワークの自動構築に成功」

([https://www.docomo.ne.jp/info/news\\_release/2026/03/02\\_00.html](https://www.docomo.ne.jp/info/news_release/2026/03/02_00.html))

※1 IOWN 構想に基づく光ネットワーク基盤であり、超低遅延・広帯域・低消費電力を特長とします。本実証実験では、5G ネットワークと接続することで、分散配備された遠隔 GPU リソースを低遅延かつ安定的に接続する基盤として活用しました。

※2 5G コアネットワーク上に実装され、5G ネットワークと IOWN APN を接続するとともに、通信の制御に加えて AI 推論処理をネットワーク側から制御するエッジ機能です。INC エッジは、5G ネットワークと IOWN APN および INC を接続する UPF である dUPF(DPU offloaded dUPF)と、AI 推論処理を推論の前段にあたる処理と推論の実行部分に分け、前段処理後のデータを IOWN APN を介して遠隔 GPU リソースへ転送・振分けすることで分散推論を実現する AI Proxy から構成されます。

※3 AI が、事前に学習した知識やパターンを用いて、新しいデータ(画像、音声、テキストなど)に対して分析を行い、予測、分類、判断などの結果を導き出す一連の処理のことです。

※4 2026年3月2日時点。ドコモ・NTT調べ。協働ロボットの安全要求事項を定義するISO/TS 15066内で定められた特定の条件(ロボットと人間の距離、人間の移動速度)における要求遅延の数値内で、本実証に成功。

※5 アプリケーションレイヤの処理機能を、ネットワークのデータ転送制御と一体で扱うことで、遅延や端末の消費電力を低減しつつ、高性能・高機能なサービスを実現する技術コンセプトです。ネットワークが通信だけでなく、情報処理の配置や実行の制御にも主体的に関与する点が特徴であり、ネットワーク内に配備されたアクセラレータや計算リソースへ情報処理をオフロードすることで、端末の負荷を低減することが期待されます。

※6 <https://www.mwcbarcelona.com/>

■本件に関する報道機関からのお問い合わせ先  
株式会社 NTTドコモ  
ブランドコミュニケーション部 広報担当  
Tel:03-5156-1366

NTT 株式会社  
情報ネットワーク総合研究所 広報担当  
[問い合わせフォームへ](#)