

【エッジAIとは何か】

[最先端のAIテクノロジーとTDKの磁性技術のコラボレーション | Tech-eye | TDK Techno Magazine](https://www.tdk.com/ja/tech-mag/tech-eye/01)

<https://www.tdk.com/ja/tech-mag/tech-eye/01>

エッジAIは計算量を抑えつつ高度な知的機能を実現するアーキテクチャ

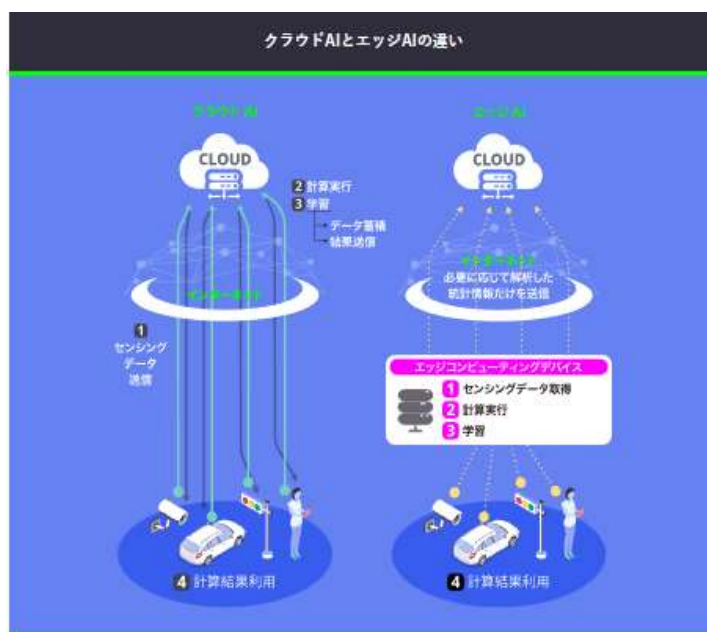
IoTなどの普及により、今後センサの数が増えると、データの量も計算の量も膨大になり、大量のデータをクラウドとの間で送受信すると、演算を行うための時間も膨大になり、電力消費の増大が大きな課題となります。

これを解決するために考えられているのが、エッジAI*というアーキテクチャ。

クラウドAIとエッジAIの違い

クラウドAI：ネット経由でクラウド側のAIでデータ処理して学習や予測などを行う従来型のAIアーキテクチャ。大量のデータ処理には有利だが、自動運転など、リアルタイム性が求められる用途ではデータの送受信にともなう遅延問題が避けられない。

エッジAI：無線通信技術を利用して、ユーザの端末の周辺部(エッジ)にAIデバイスを配置してデータ処理する分散型のAIアーキテクチャ。IoTデバイスからのデータの高速処理や、リアルタイム性が求められる用途に有利。

**リザーコンピューティングとは？**

エッジAIという言葉は、最近よく見聞きするようになったが、まだ端末にAIは入っていないで、基本的にはデータセンター側で処理をしている。また、クラウドAIからエッジAIに切り替えても、今度はエッジ側が学習機能を担うことになり、エッジ側での負担が増える。いかに学習の計算量を抑えて、かつAIの高い知的機能は維持するかが、エッジAIにおけるきわめて重要な課題となる。

そこで、最近着目されているのが、リザーコンピューティング*という新しい概念。ディープラーニング(深層学習)*などを可能にするニューラルネットワークの学習モデルは人間の脳を模したものだが、これに対してリザーコンピューティングは小脳を模したモデルとされている。

リザーバーコンピューティング：小脳の機能をモデル化したもので、入力層、リザーバー層(フィードバック結合するニューロン群)、出力層の3層からなるシンプルな情報処理構造が特長。リザーバーとは“貯水池”という意味で、リザーバーに蓄えられた過去のデータをもとに、計算量を抑えつつ、未来予測などの高度な知的機能を実現するのがリザーバーコンピューティング。リアルタイム性が求められるエッジAI、IoTデバイスのセンサデータなど、時系列データを低消費電力で高速処理する用途に向いている。

ディープラーニング：深層学習。大脳の機能をモデル化した多層構造のニューラルネットワークにより、膨大なデータから自動的に特徴を抽出する高度な機械学習を可能にした。画像認識や音声認識などに利用されるが、計算量が多いため処理時間が長く、電力消費が大きくなるのが難点。

大脳は記憶や論理的な考え方に関わるが、小脳は運動機能に関わる。たとえば、モノが落ちそうになったとき、反射的に手をのばすのは、大脳で考えてそうしているわけではなくて、「この状態になったものは落下する」という経験に基づいて身体が動く。そのときの計算量というのは、それほど多くない。

ニューラルネットワーク(ディープラーニング)とリザーバーコンピューティングによる情報処理の違い

ニューラルネットワーク(ディープラーニング)

- 大脳の機能をモデル化。●入力層と出力層の間に、複数の中間層をもつ階層的ニューラルネットワーク。
- 計算量が多く、処理時間が長くなるため、電力消費が大きい。

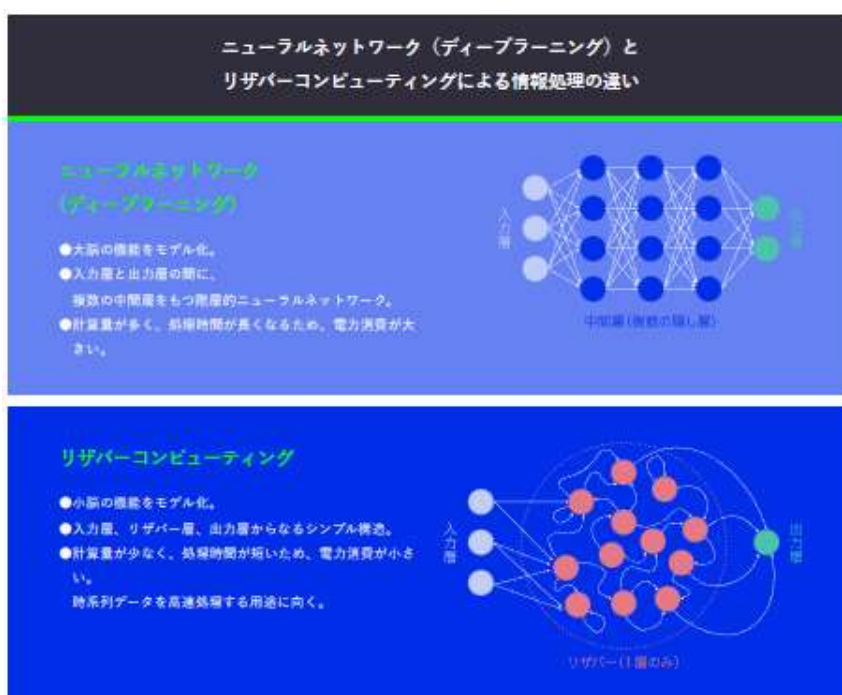
リザーバーコンピューティング

- 小脳の機能をモデル化。●入力層、リザーバー層、出力層からなるシンプル構造。
- 計算量が少なく、処理時間が短いため、電力消費が小さい。

時系列データを高速処理する用途に向く。過去の経験をもとに、少ない計算量で効率的に未来予測。

リザーバーコンピューティングができることは、ざっくり言うと、過去の経験をもとに、少ない計算量で未来予測すること。それは我々が日常的に行っていることで、勘のいい人なら3つ、4つ前の事象から次を予測する。

リザーバーコンピューティングの仕組みを取り入れると、10も20も前の、深い過去の事象をもとに、少ない計算量で未来予測できるようになる。計算量が少ないということは、処理時間が短く、消費電力が少なくなることを意味する。



リザーコンピュティングは、いくつかの特殊な **リカレントニューラルネットワークモデル** から派生した計算の枠組みで、主に時系列データの機械学習に利用される。典型的なリザーコンピュティングモデルは、**入力層**、**リザー**、**出力層** から構成される。入力層は時系列データを受け取って適当な重み付けを行い、リザーへ情報を送る。リザーは、**スパースでランダムな結合** をもつリカレントニューラルネットワークで与えられ、入力層からの情報を高次元時系列データに非線形変換する。時系列入力データの時間方向の関係性（依存性）を考慮するために、リザーは過去の入力情報を蓄積して記憶する役割を持ち、**リードアウト** では、その高次元時系列データを用いて、線形回帰などの簡便な学習アルゴリズムにより回帰や分類などのパターン解析を行う。リザーコンピュティングモデルの特徴は、リザーと出力層の間の結合重みだけを学習アルゴリズムで決定し、入力層とリザーの間の結合重みとリザー内のフィードバック結合重みはあらかじめ固定しておく点。この工夫によって、すべての結合重みを学習する一般のリカレントニューラルネットワークに比べて高速な学習が可能となる。ただし、高い計算性能を実現するには、**あらかじめ固定する結合重みの値を適切に設定しておく必要がある**。

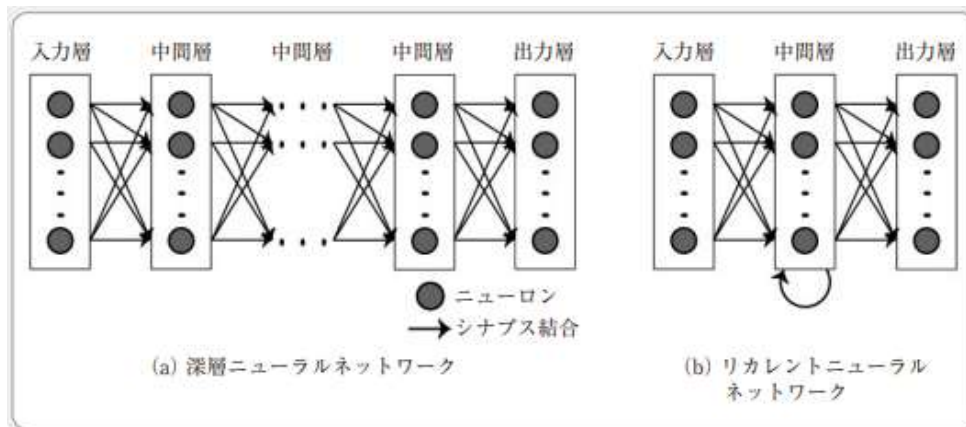


図1 深層学習を適用するニューラルネットワーク

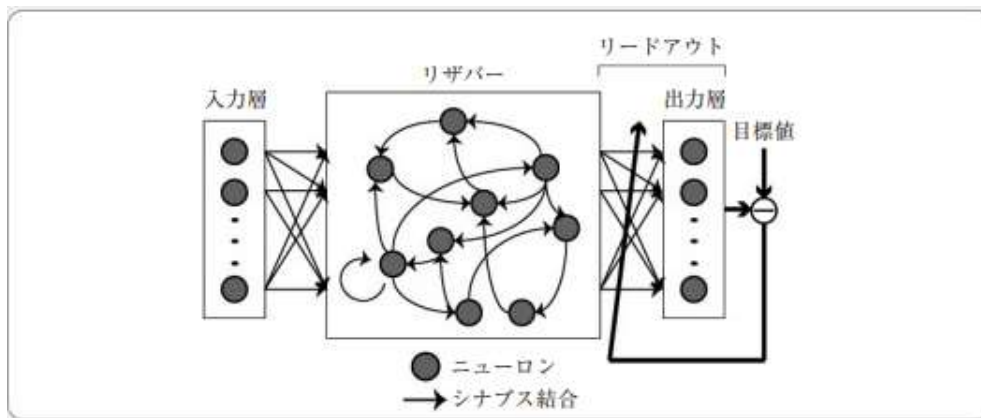
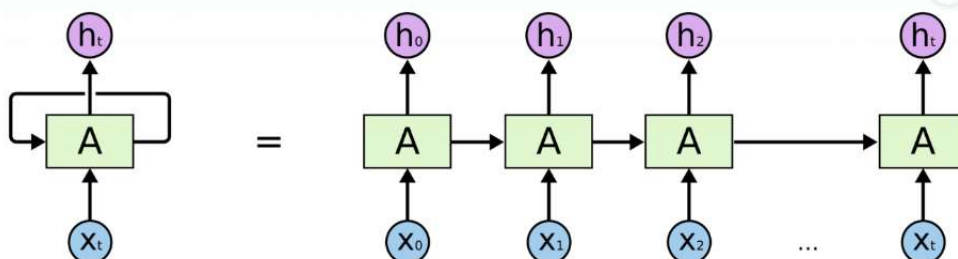


図2 リザーコンピュティングモデル



リザーバーコンピューティングでは時系列データを扱う。時系列データは、センサや計測器を通じて時々刻々と変化する一連の値として取得される。こうした時系列データに対する主な機械学習タスクには、時系列生成、時系列分類、時系列予測などがあり、図3 (a) は時系列生成タスクの例。この例では入力は定数入力が時間とともに切り替わる時系列データで、出力は入力値を周波数とするような正弦波。学習したモデルは、正弦波発生器（ファンクションジェネレータの一種）として動作する。このようなタスクは、例えばロボットの動的パターン生成などに応用される。図3 (b) は時系列分類タスクの例。この例では入力は正弦波または三角波の時系列データで、出力は正弦波または三角波に対応するラベルになる。学習したモデルは時系列入力の定性的な違いを判別することができる。このようなタスクは、例えば、音声信号から発話者や発話内容を推定する音声認識などに応用される。図3 (c) は時系列予測タスクの例。この例では入力は乱数の時系列データで、出力は入力のある非線形システムによって非線形変換した時系列データになる。学習したモデルは上記の非線形システムを近似するので、時系列データの将来を予測することができ、例えば、気象予測や経済予測などに応用される。

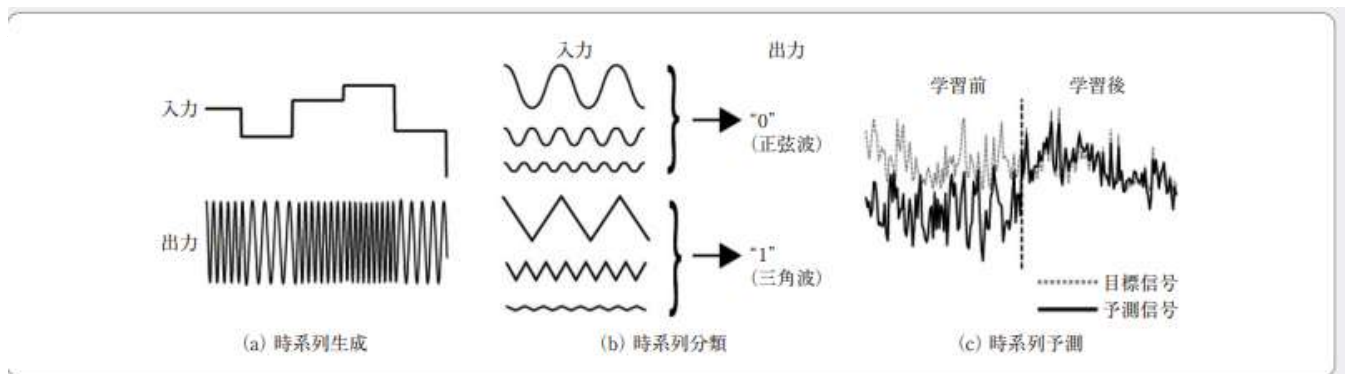


図3 代表的な時系列パターン認識

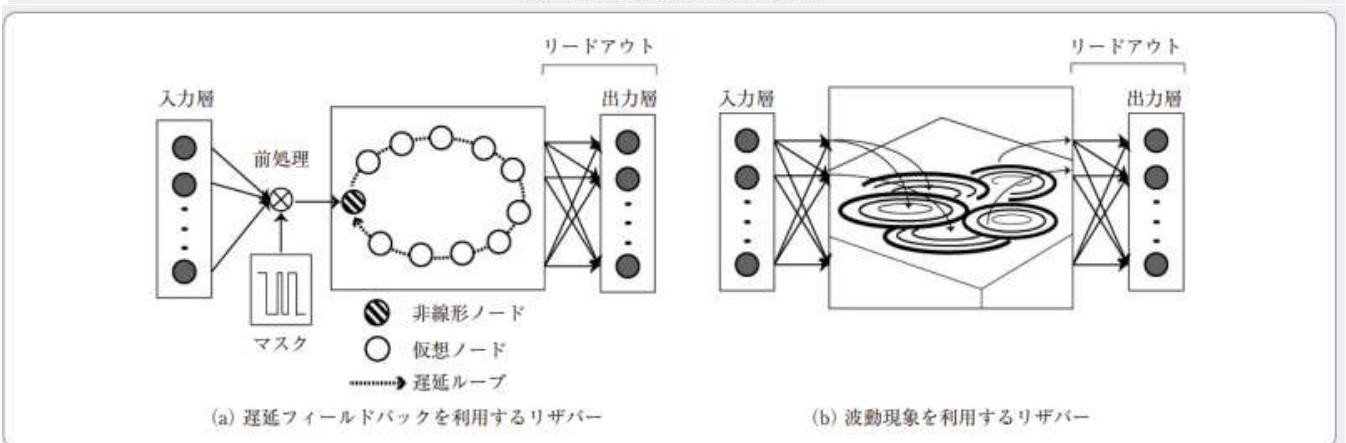


図4 物理リザーバーコンピューティング

[引用] 図1~4：「知っておきたいキーワード リザーバーコンピューティング」、田中剛平、映像情報メディア学会誌 74, 3 (2020)

<https://www.ite.or.jp/contents/keywords/2005keyword.pdf>