



SEMINAR REPORT

## ディープラーニング（深層学習）とハードウェア技術 ～富士通におけるディープラーニングへの取り組み～



富士通の吉山と申します。今日は、「ディープラーニングとハードウェア技術～富士通におけるディープラーニングへの取り組み～」ということで、お話しさせていただきたいと思っております。私は、富士通のAI サービスを提供する部門に所属しております。実際にお客さまとお話させていただいて、業務へのディープラーニングの適用を行っております。

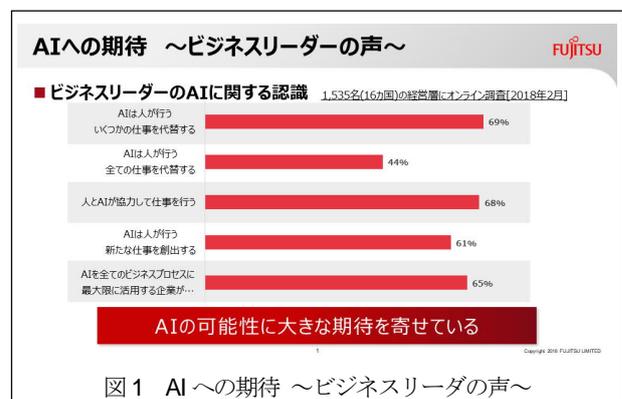
このディープラーニングについて説明しようとする、どうしても微分とか偏微分とか行列とか、数式を出さないといけません。なるべく簡単に説明することに努めたのですが、一部数式が出てくるところはご容赦願いたく存じます。

さらにもう1つ。私はこの事業部に来て、まだ1年ほどです。ある意味、初心者に近いところもあるので、用語とか理解とかに皆さんと違うところがあるかもしれませんが、そのところはちょっと目をつぶっていただければと思っております。

### はじめに

AIへの期待ということで、16か国の経営層に「AIをどのように認識をしていますか?」というアンケートを取りました。図1をご覧くださいますと、「AIはいつか仕事を代替する」とか、「ビジネスプロセスに最大限活用する」とかのご意見をいただいております。スコアの低いところでは、「人が行う全ての仕事を代替する」というのがあります。ニュースとかメディアとかによると、AIが発展すると人のする仕事が全て無くなってしまわないのか、全てAIに奪われてしまわないのかと言われていますが、さすがビジネスリーダーの方は冷静に見えて、全て代替することはないだろうと思っておられます。こ

が1つ重要なポイントだと思っています。そうはいつても、他の項目は全てスコアが高いので、やはりAIの可能性に大きな期待を寄せているというのが現状ではないかと思っています。



### 人工知能(AI)とは

まず、「人工知能、AIとは何でしょうか」というところからお話しします。現在のテレビCMを見ている、事あるごとにAIという言葉が出てきていますが、実際問題として、AIとは何かについてきちんと押さえておく必要があると思っています。もっとも、AIという言葉には曖昧なところがあって、人工知能学会で定義されているのですが、正直なところ知性とか知能とかに係わる問題であって、そもそも知性・知能には定義がないのです。なので、AIは人によって色々な定義があるので、そのところを踏まえた上で見ていただければよいと思います。

図2に示すように、人工知能の範疇には機械学習というものがあります。これは何かといいますと、機械で学習するというよりも、人が色々情報を集めてきて、Aが来たらBを返す、Cが来たらDを返すというようなルールを決めて、それをコンピュータに覚え込ませる手法となります。それをさらに狭義で捉えるなら、ニューラルネットワークというものがあり、神経細胞をモデル化したものです。このモデルを使って学習させるもので、その中の1つがディープラーニングです。神経細胞がスマートフォンを見て「スマートフォンだ」と分かるのは、脳の中で何かの反応が起きていて、それでスマートフォンと判断しているのであって、それと同じようなモデルをコンピュー

タにやらせて、色々なもの、例えば、画像認識とか、音声認識とかをやらせようというのがディープラーニングです。

**人工知能(AI)、機械学習、ディープラーニング**

**人工知能**

**機械学習** 開発者が予めすべての動作を決めておく従来型のプログラムとは異なり、与えられた情報を元に学習し、自律的に法則やルールを見出す手法やプログラムのこと

**ニューラルネットワーク** 機械学習の一種でニューラルネットワークと呼ばれる、主に生物の神経系の挙動を模して学習できるようにデザインされたもの

**ディープラーニング**

- ✓ 画像認識
- ✓ 音声認識
- ✓ 自然言語処理 など

図2 人工知能、機械学習、ディープラーニング

AIは万能のように思われていますが、実は、ディープラーニングの世界は万能ではありません。普通のコンピュータ処理では、Aの処理をさせてBの結果となって、必ず正しい答えが出せるのですが、ディープラーニングでは、常に正しい答えを出せるかという、そうではなくて、8, 9割方が正しいという世界です。

コンピュータとAIでは何が違うのかというと、コンピュータは人がすると難しいようなことを、とても速く処理してくれます。例えば、10万行のデータの中からあるデータを探し出すとき、人はとても時間がかかるのに、コンピュータはすぐに見つけてくれる。逆に、AI特にディープラーニングは、人のできることがようやくやれる程度です。そここのところが違うと思います。先ほどのスマートフォンの話ではないですが、これはスマートフォンだと人の目にはすぐにわかりますが、これをAIに判断させようとする、長い時間をかけて勉強させて、ようやく判断できるのです。そのような例を、次から具体的にお見せしたいと思います。

写真を見て犬と分かるのは何故か。古典的な機械学習では、犬の特徴を人が教え込みます。犬はどのような特徴があるかという、簡単な例では、図3に示すように、①4本足で、②鼻がシュッとして、③毛で覆われている。この3つぐらいが犬の特徴ではないかと思えます。他に何か特徴があるのかとなって、特徴を次々と覚え込ませていかないといけませんが、従来の機械学習のロジックです。

**写真を見て「犬」と分かるには?**

■ (例) 犬を認識する画像認識処理

■ 古典的な機械学習

● 画像から犬を判別する「特徴」を人が教え込み、コンピュータが判断

① 4本足      ② 鼻がシュツ      ③ 毛で覆われている

では、キツネとの違いがわかる特徴は?  
↓  
特徴をすべて言葉にするのは難しい

犬には他にどのような特徴がありますか?

図3 写真を見て犬と分かるには?

それで、犬は分かっただけでも狐はどうなのか。犬とどこが違うのか。確かに狐も4本足で、鼻がシュツとしていて、毛で

覆われています。なので、特徴を全て言葉にして覚え込ませるとなると非常に難しいです。機械学習では、昔はこれを一生懸命行っていました。プログラム全体が「If~then else~」の形で書き連ねられていて、最終的には分からないという経験をされた方もおられるかと思いますが、そういうことが起きます。機械学習で教え込ませることは非常に大変で、どうしても漏れが生じてしまい、認識率が低くなってしまっているのが現状です。

図3では、4本足の話をしました。図4では、これは何本足ですかという話です。①飛んでいるとわからないし、2本足で立たれると難しい。②鼻がシュツと言われても、そうでない犬もいる。③全身が毛で覆われているかとなると、最近では色々な犬がいて、覆われていないものもいる。人が全ての犬種、動き、状態を加味して特徴で教えることは、非常に難しいのです。なので、ディープラーニング手法が出現する前は、人が一生懸命覚えさせていたのが、ディープラーニングが出てきて、状況ががらりと変わってきたというのが現在の状況です。

**人が特徴量を教え込むことの限界**

■ 人が特徴量を考え出し、教え込むことは**すごく大変**

■ 漏れが生じやすく、認識率が低くなってしま

① 4本足?      ② 鼻がシュツ?      ③ 毛で覆われている?

人がすべての犬種、動き、状態などを加味した特徴量を教えることは質・量とも難しい

図4 人が特徴量を教え込むことの限界

先ほどの犬の話では4本足、鼻がシュツというのが特徴であったと思いますが、人が教えずとも機械が画像から特徴量を自分で見つけ出すということが、ディープラーニングの特徴となります。なので、まずは入力として犬の画像をたくさん用意して、一生懸命覚え込ませます。そうすると、画像の中で特徴量を見つけて、これが犬だと分かる学習済のモデルを作ります。そうすると、次にこのモデルを使って、例えば、猫と入力して認識済モデルにかけると、これは犬ではないということで、犬の確率を0.0001とはじき出します(図5)。これが、ディープラーニングの特徴です。なので、機械学習では人が特徴を覚え込ませていたものを、ディープラーニングでは画像を大量に用意すれば、何か知らない内に犬と分かるモデルを作ることができるのです。

**ディープラーニングの特長 - 特徴量を見つける**

■ 人が機械に教えずとも、大量の画像から**特徴量を自ら見つけ出す**

入力: 数万~数千万単位の大量の犬画像

学習フェーズ: ニューラルネットワーク

出力: 犬認識学習済モデル

大量の犬画像(教師データ)から、犬の特徴量を見つけ出す

これを使うことで...

入力: 猫の画像

推論フェーズ: 犬認識学習済モデル

出力: 犬の確率: 0.0001

<https://www.what-dog.net/>

図5 ディープラーニングの特長 - 特徴量を見つける

図6は、従来の機械学習とディープラーニングの違いを比較したものです。従来の機械学習では、人がデータを分析して特徴量を抽出して、学習方法を設計しています。なので、どのように学習しているかは、作っている人には分かっているのです。ディープラーニングでは何が違うのかというと、機械が自動抽出・設計して特徴量を見つけています。なので、抽出から学習のところはブラックボックスになっています。ここで何が起きているのか分からないが、結果は分かるということになります。

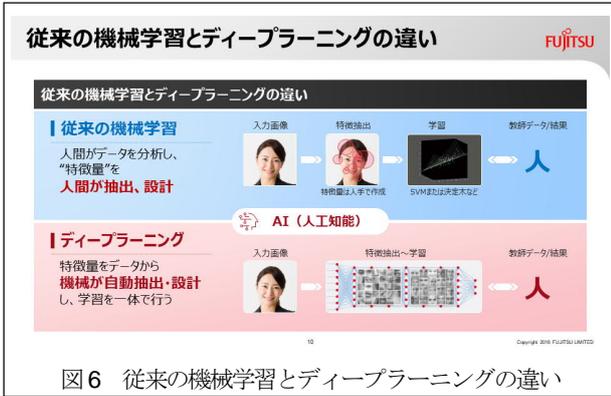


図6 従来の機械学習とディープラーニングの違い

機械学習とディープラーニング、少し前までは機械学習の方が成績がよかったのですが、あるときブレークスルーが起きて、ディープラーニングが脚光を浴びるようになりました。ILSVRC (ImageNet large scale visual recognition challenge) という画像認識のコンテストがあります。2012年に、トロント大の SuperVision というディープラーニングを使ったプログラムが出展されました。エラー率は16%です。これは相当間違っているような気がしますが、2位は20%を超えています。20%以上間違っているというのは、人の能力と比較するとほど遠い世界です。人の間違い率は5%ぐらいです。それでも、圧勝したのがディープラーニングということで、脚光を浴びたのです。確かにぶっちぎりでしたが、そうは言っても、16%というのが2012年当時の実力でした。

では、今はどうなっているかということ、図7の棒グラフ黄色がエラー率16%ですが、1回脚光を浴びると誰もがそれに群がってやり始めて、今はどこまで向上しているかということ、2017年で2.3%です。なので、ディープラーニングはほぼ間違えない、人ほど間違えないレベルまで到達しています。劇的に進化しているのです。2012年より前までは、28%ぐらいの性能でけっこう機械学習が使われていたのですが、現在は、ディープラーニングに置き換わってきています。

図8は、AIと富士通の歴史を表わしたものです。富士通は昔からAIに取り組んでいます。日本初のAI搭載コンピュータとか、知識情報システムとかを世に送り出しています。ディープラーニングには計算機パワーが必要ということで、昔からスパコン「京」に取り組んでいたのが、AIコンピュータにも取り組んでいるのです。第1次AIブームがあって、冬の時代があって、第2次AIブームがあって、また冬の時代があって、今は第3次AIブームです。今のブームを牽引しているのは、やはりディープラーニングということになるかと思えます。

なにゆえ、このように脚光を浴びるようになったかということ、やはり、ビッグデータと計算機パワーとアルゴリズムの3点セットが揃って、ようやくできるようになったのではないかと思います。

っています。1つ目は、データが圧倒的に多数揃えられるようになったことです。第2次AIブームのときに失敗したのは、ビッグデータを自動的に集められなかったのが大きかったようです。極論すると、当時は町中歩いて、猫の写真を撮りまくるといったような状況でした。今はGoogleで検索すれば、猫の画像がいっぱい出てきます。2つ目は、計算機パワーが圧倒的に上がったことです。後ほどお話ししますが、GPGPU (General purpose computing on graphics processing units) が出現して、計算能力が圧倒的によくなったことです。3つ目が、アルゴリズムのニューラルネットワークの出現です。色々なことに対応できるようになったことが非常に大きいと思っています。

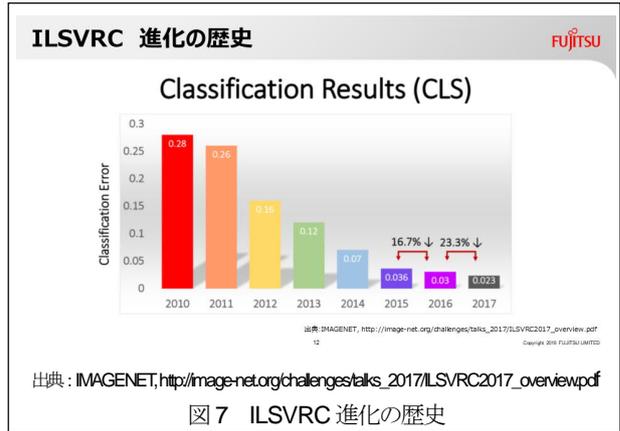


図7 ILSVRC 進化の歴史

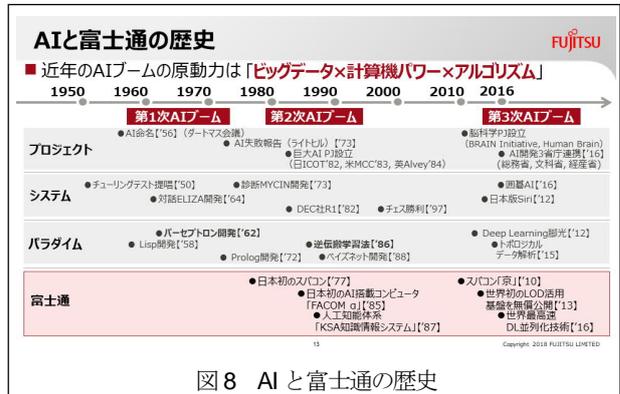


図8 AIと富士通の歴史

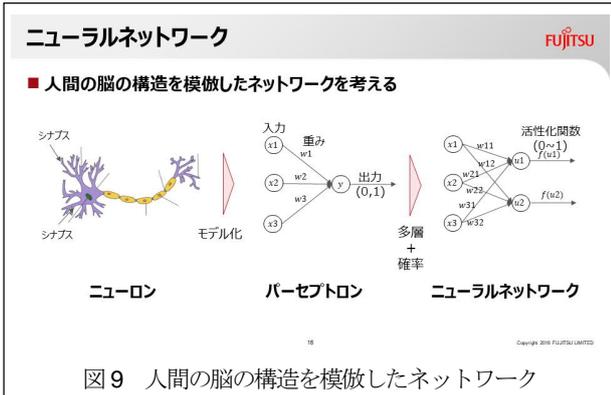
## アルゴリズム

それでは、3点セットの内、ビッグデータはスキップして、アルゴリズムと計算機パワーについてお話をさせていただきたいと思っています。

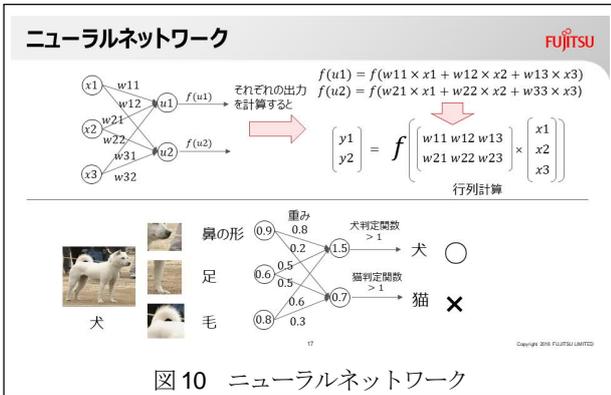
ニューラルネットワークは、人の脳の構造を模倣したネットワークをコンピュータ上で構築しようというのがスタートポイントです。シナプスが外部からの刺激を受けて違う刺激を出して次々と伝えていくことで、スマートフォンを認識したり、言葉を読んだりできるというのが人の脳の構造です。

図9中央に示すように、これを模して考えられたのがパーセプトロンです。入力x1, x2, x3を何かしら入れて、それに重み付けをして計算した結果を出力します。この出力したものを、また次の入力にする。要するに、このような脳の単純なモデルの組み合わせと考えたのがパーセプトロンです。このパーセプトロンの出力は、0と1、真か偽かしかありません。これでは使いにくいということで、もう少し使いやすいのを考えよう

ということになって、脳の構造に近いもの、多層で組み合わせることにします。先ほどの図5の犬と猫ではないですが、犬は80%ぐらいで100%ではないということで、0/1 がはっきりしないものは、ある関数を使って0から1までの確率分布で出力するようにします。このように組み合わせたものがニューラルネットワークになります。なので、それぞれの入力に対して重み付けした計算結果を踏まえて、次の出力に出すか出さないかを組み合わせてやるのが、このニューラルネットワークです。これを組み合わせれば、脳みそと同じような構造ができて、先ほどのような画像認識ができるのではないかとというのが、ニューラルネットワークのコンセプトです。



それで、計算は実際にはどのようにするかというと、例えば  $u_1$  の出力は、それぞれの入力に対して重み付けして足していきます。 $f(u_1)$  は  $w_{11}x_1 + w_{12}x_2 + w_{13}x_3$  と計算します。同様に  $u_2$  の出力  $f(u_2)$  も同じような計算式となります。これを行列計算で表すと、図10のような行列計算式に置き換えられます。

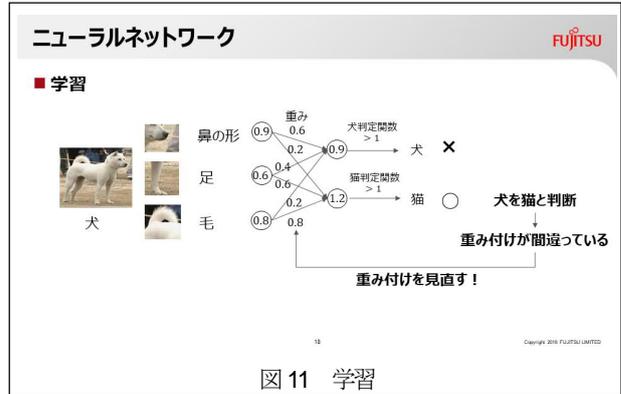


先ほど、図3で示した特徴、鼻の形、足、毛の組み合わせで犬と猫を判断します。犬は重み付けした計算結果が1.5になりました。1以上は犬であるというフラグが立ちます。猫判定関数は0.7となって、1以下なので猫ではない。ここでは単純な例で示したのですが、仕組み的にはこのような重み付け計算して、犬か猫かを判断していくことになります。

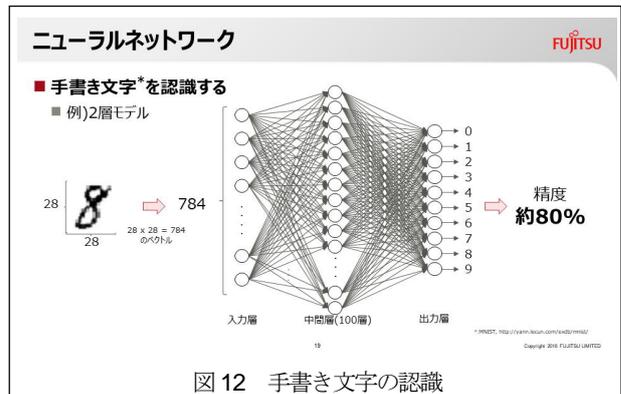
図11に示すように、最初は適当に重み付けしているので、実際に計算してみると、犬と猫を正しく判定してくれません。犬判定関数が1以上にならないので、犬判定にならない。判定が間違っているので、重み付けを見直さないといけない。重みの0.6, 0.2, 0.4を直していく、すなわち学習することになります。こうして重み付けの見直しを繰り返すことで、最終的に犬

か猫かが正しく判定できるようになるというのが、ディープラーニングの重要なポイントとなります。

この重み付けの見直しを人がするのはとても無理で、プログラミングで自動的に行うのがディープラーニングです。この例では、入力が3つで出力が2つなので、重みは全部で6つしかありませんが、層が深くなると、重みを次々と変えていかなければいけない。これがうまく調整できると、犬は犬、猫は猫と判断できるようになります。



そこで、実際に図12のようなモデルを作ってみました。入力は  $28 \times 28$  のデータなので、全部で784。これをビット列に直して、784入力のモデルを作りました。中間層に計算した結果を1回保存しておいて、もう1回計算して出力するモデルになっています。0から9までの数字を判定するもので、実際に計算してみました。デフォルト状態で計算して、80%の認識率で数字を認識することができました。家にあるパソコンで計算したもので、80%ならまずまずといったところです。計算時間は1分もかかっていません。



この80%の認識率はどうかというと、これでは容認できなくて、90%台までもって行かないと業務には使えません。そこで次のアクションとしては、80%をいかにして上げるかがポイントとなります。ディープラーニングでいうところのチューニング作業になります。

色々やり方はありますが、図13に例として書かせていただいたのは、中間層を変えてみることです。入力層784、中間層100を1回、出力層10にしているので、中間層が多い方が認識率が上がるのではないかと考えて、200にしてみたら81%になりました。それでは、減らすとどうなるかということで、50にしてみたら75%になりました。次に、活性化関数変えて

みました。デフォルトはシグモイド関数ですが、ランプ関数に変えてみました。そうすると、認識率が88%になりました。その次は、学習回数を変えてみようということで、2万回にしてみたら83%になりました。

チューニング作業でこれら3つを組み合わせると、91%まで認識率が向上しました。デフォルトのままでは80%なのが、チューニングを施すことで91%まで上がりました。このような単純なモデルでも、チューニングすれば91%まで上げられるということです。

このチューニングも、AIでできたらよいのですが、ここのところは人がしているのです。中間層の数を増やしてみるとか、中間層の回数をもう1つ増やしてみるとか、人が行って認識率を向上させているのが実情です。

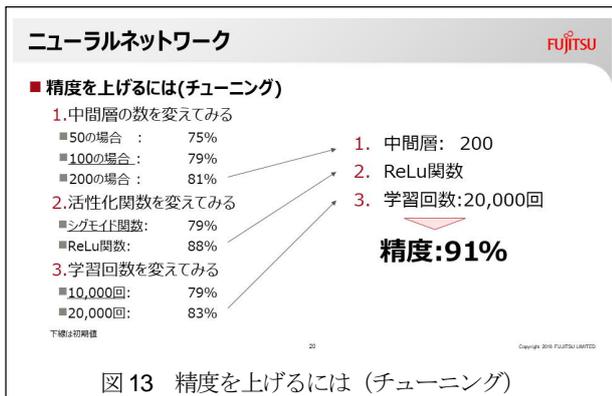


図 13 精度を上げるには (チューニング)

当事業部では、年間40件ぐらいの案件を扱っているのですが、先ほどのようなチューニングを繰り返して行っています。実ビジネスでは、顧客からデータをもって、最初は50~60%ぐらいだったものを80~90%まで上げていく作業をしている状況です。これを2~3か月で終えなければいけないので、現状のAIは結構泥臭い作業だと思います。

図14は川崎地質様の実際の例です。今までは、撮ってきたレーザー画像を人が認識して、「ここは空洞だ」と見つけていましたが、これにAIを適用しました。認識率も非常に高く、探索時間を減らすことができたということです。

これも泥臭い作業で、1回の学習に10時間以上かかっています。精度を上げていくとなると、数か月はかかります。モデルを作って夜中に動作させておくと、帰った後にエラーになっていることもあって、チューニングに結構時間がかかるのです。



図 14 道路陥没を防ぐ路面下空洞調査

実際のディープラーニングは、すべきことが非常に多くて、

利用までの道のりはとても長いです。図15は、当社のプロセスモデルを示したのですが、まず、顧客の要求が結構曖昧です。とにかくAIを導入したいという顧客は多くいるのですが、課題設定、要件定義のところを誤ると後が大変です。例えば、カメラ画像でトラックを識別したいという話であれば、データの確認、データの準備、データの分析のところは非常に重要です。このトラックの写った画像にトラックというタグを付けていますが、これを誤ると違うものを認識してしまいます。データの数をいかに綺麗に揃えるかも非常に重要で、ここにも大それた時間がかかっています。

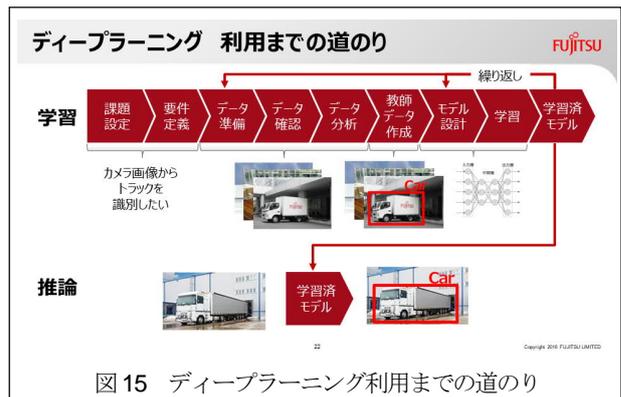


図 15 ディープラーニング利用までの道のり

これらのことを踏まえてモデル設計して、学習させる。結果がよくないなら、モデルを見直したりデータを見直したりして、ようやく学習済のモデルができ上がります。それで推論をして、実際にトラックの写った画像を見せて、これはトラックだと認識させる。このような手順を踏みます。ですので、通常のプログラミングでは、ロジックが決まったらそのとおりに組めばよく、後はバグがあるかないかの問題ですが、ディープラーニングはバグがあるかどうか分からないブラックボックスなので、やってみないとまよくいくかどうか分からない。やってみても、その正解率でよいのかどうか分からないということで、判断が非常に難しいのです。

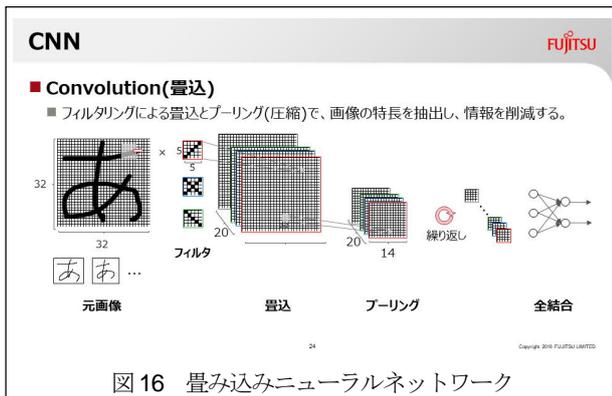
ニューラルネットワークには、様々なモデルがあります。RNN (Recurrent Neural Network) とか、LSTN (Long Short-Term Memory) とか、CNN (Convolution Neural Network) など色々なモデルがあって、これらのモデルの中から顧客の要望に沿うモデルを見つけてきて、組み合わせて構築する。これをベースにさらに改善をしていくのが、現状のビジネスではないかと思っています。

この中で、現在、画像認識でもっとも使われているのがCNN (畳み込みニューラルネットワーク) です。これは、ILSVRCのコンテストでも使われているモデルです。文字を識別するとき、人は全てを見なくても分かります。なぜなら、人は特徴を見つけて文字を判断しているのです。文字の特徴を抽出すれば、画像認識がうまくいくのではないかとというのが、CNNのコンセプトになります。要するに、人の見方と同じような考え方を導入しようということです。

どのような処理をしているかという、図16に示すように、例えば、「あ」という文字を32x32で区切ります。このデータは、単に飛び飛びの点で組み合わせられているのではなく、結構線状に連続していたりして、それぞれの部分部分で特徴があるのです。なので、特徴を探し出すために、一定の区画ごとにフ

フィルタリングをかけて抽出します。例えば、5x5 で斜めの特徴を見つけたいなら、これで一旦フィルタリングをかけます。これ以外の特徴に対しても、フィルタリングをかけることで特徴だけを次々と抽出していきます。このような処理のことを畳み込みと言います。

色々な特徴をたくさん抽出して、そのままではデータが大きくなるので、さらにプーリング、すなわち圧縮処理を施します。このプーリングしたものをフィルタリングにかけて、特徴を導き出します。このように処理を繰り返して、最後に全結合のニューラルネットワークが書けるのです。CNN は、このように処理することで画像の特徴を抽出して、「あ」とか「い」とかの文字を抽出するアルゴリズムなのです。これらのモデルは一般に公開されているので、簡単に使うことができます。



当社は、DeepTensor というグラフデータ学習技術、人やモノのつながりを表すグラフデータから新たな知見を導くディープラーニングの新技术を提供しています。一般的な技術と当社独自技術を使って、顧客の様々なニーズに対応できるようにしています。

ディープラーニングをプログラミングするときには、フレームワークを使います。行列計算式を全て人が書いていたのでは大変なので、また、誤差が生じたときの重み付けを人が直すのも大変なので、これらを全てしてくれるフレームワークが用意されています。オープンソースになっていて、無償で使えます。なので、これらオープンソースを使うことで、簡単にプログラムが組めるようになっています。

たくさんあって、どれを使ったらよいか悩みますが、図 17 はフレームワークの人気度<sup>1)</sup>を示したものです。一番人気は TensorFlow で次が Keras です。Keras は、TensorFlow をラッピングしているフレームワークで、さらに簡単に使えるモデルとなっていて、急激に伸びてきています。なので、TensorFlow を勉強すれば、比較的情報量も多いので、簡単にプログラムが組めるようになります。

先ほどの図 12 のモデルを TensorFlow で書いてみると、複雑な行列演算とか、間違っただけの重み付けを直すとかも含めて、わずか 73 行で書いてしまいます。

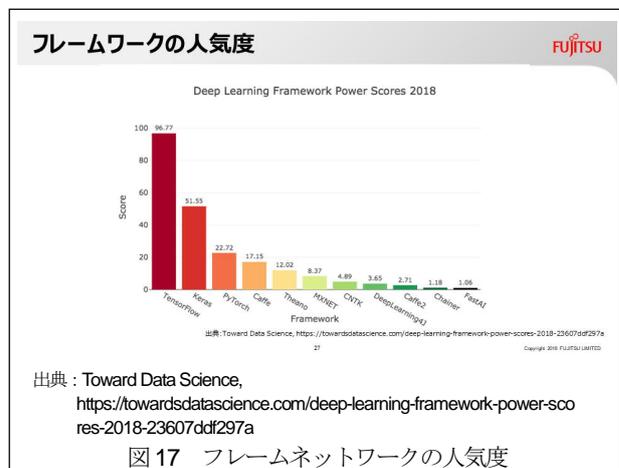
TensorFlow はコンソールで書きますが、最近は GUI が用意されていて、もっと簡単にディープラーニングができます。当社は、Zinrai という名称で体系化していて、クラウドでこれを提供しています。他には、Sony、NVIDIA などが GUI で提供しています。

当社の Zinrai をご紹介します。画像を認識して、それが何か

を当てるとい例です。教師データは、24 種類 3 万枚用意しました。このデータを使って、このモノが何かを推論するプログラムを実際に作ってみたいと思います。

まず、フォルダにデータを用意します。教師データはインターネットから拾ってきます。次に、学習ウィザードを立ち上げて、プログラムの名前、ネットワーク種別、教師データ、最終的に学習済みにするモデル、パラメータなどを設定して、プログラムを走らせます。学習が完了すると、学習済みモデルができ up します。これで画面の上の認識用画像を選んで、推論できるようになります。

Zinrai の特徴は、エッジモデルも作ることができます。スマートフォンに学習結果を送って、カメラで撮って、その場で推論ができる仕組みも用意しています。例えば、お茶の葉の開き具合をその場で確かめたいのであれば、葉の状態を学習させておいて、現場で写真を撮って、摘み取るタイミングかどうかの判定モデルを簡単に作ることができます。



## 計算機パワー

先ほどの図 14 の学習では 10 時間以上かかるという話をしましたが、これでは 1 日に 1 回しか学習できません。競合他社があるので、この時間をいかに短くできるかが重要なポイントとなります。それには、計算機パワーが必要だということです。なぜ GPGPU を使うとディープラーニングができるかというと、ニューラルネットワークは行列の計算です。同じ計算をするなら、GPGPU にやらせた方が圧倒的に速いのではないかとすることで、GPGPU がニューラルネットワークの計算に使われるようになりました。

汎用 CPU でも計算はできます。図 18 に示すように、汎用 CPU と GPGPU でクロック速度も消費電力も似かよったものですが、明らかに構造が違います。汎用 CPU は連続で複雑な処理が得意で、GPGPU は単純で大量な処理が得意です。画像の点を一気に書き換えるには一気に計算しないといけないので、並列処理が得意でないといけません。コア数を比べたら、汎用 CPU が 16 に対して GPGPU は 5,120 と数が大きく違います。なので、GPGPU は特定のことはできないが、行列演算だけを考えると、圧倒的に GPGPU の方が処理が速いです。それでは、普通の表計算を GPGPU にやらせたらどうかというと、それぞれ得手不得手があって、使い方に依って選択するというのが一般的な考え方です。

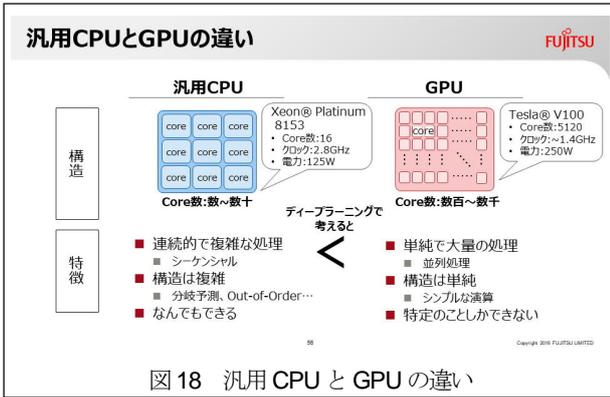


図 18 汎用 CPU と GPU の違い

計算時間は 27 倍ぐらい違って、圧倒的に GPGPU の方が速いです。汎用 CPU でもできないことはないですが、図 12 のモデルで 1 分ぐらいかかります。大きな画像を扱うとか、CNN で大きなネットワークを作るとなると、それなりの計算能力が必要になるということです。コンピューティングアーキテクチャは、汎用 CPU で何でも扱える時代は終わりを迎えていて、図 19 に示すように、ディープラーニング、スーパーコンピュータ、量子コンピューティングなどをそれぞれの用途に応じて使い分けるのが、今の流れだと思っています。

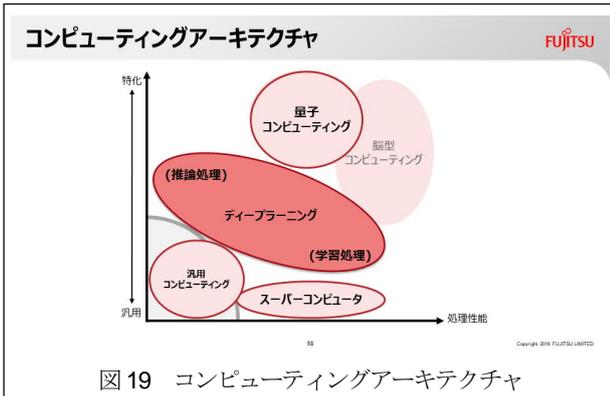


図 19 コンピューティングアーキテクチャ

ディープラーニング用プロセッサは、学習用と推論用で分けられます。学習には大きな CPU パワーが必要ですが、推論はそうでもないです。学習用で有名なところは、Google の TPU、Xeon の Phi、NVIDIA の Tesla@V100 です。今ほとんど、これらが使われていると思います。当社も参入しています。推論用は、FPGA チップに独自の回路を組み込んで使うのが、今の主力となりつつあります。

当社はメインフレームからスーパーコンピュータまで、ずっとプロセッサを独自で作っています。この技術を使ってディープラーニング用チップを作ろうということで、今は電力性能 10 倍を目指して取り組んでいるところです。

## DLU (Deep Learning Unit)

当社の DLU の基本的な考え方は、小さいコアをたくさん並べて、同時に行列演算させようというものです。図 20 に示すように、チップ内に DPU というユニットを作って、その中に DPE というさらに細かい演算ユニットを作って、同時に計算できるプロセッサを構成します。1 チップでできることは限られているので、脳というすごいネットワークと同じような構成に

するためには大規模にしなければいけないということで、「京」を作ったときのネットワークのように 6 次元のネットワークを組み込むことで、超並列処理ができるように今開発を進めているところです。

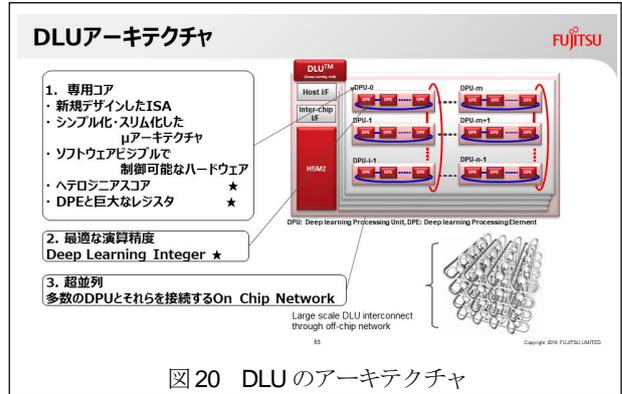


図 20 DLU のアーキテクチャ

同時に実行させるために Master コアを用意して、小さなコアである DPU に「計算して」と投げるようなアーキテクチャになっています。半導体チップ内にいかに多くの演算ユニットを詰め込めるかがポイントになると思っていて、不要なものを取り去るということで、ディープラーニングに特化した形でハードウェア設計を行っています。

一つ特徴的なのは、普通のプロセッサは演算器があって、レジスタがあって、キャッシュメモリがあるのですが、DLUにはキャッシュメモリがありません。なぜかという、重み付けの値はすぐ変えてしまうので、キャッシュに覚える必要がないのです。なので、キャッシュを全て取り払って、図 21 に示すように、大きなレジスタファイルを置いて、レジスタと演算器をなるべくたくさん並べるとというのが、ディープラーニングユニットの特徴となっています。

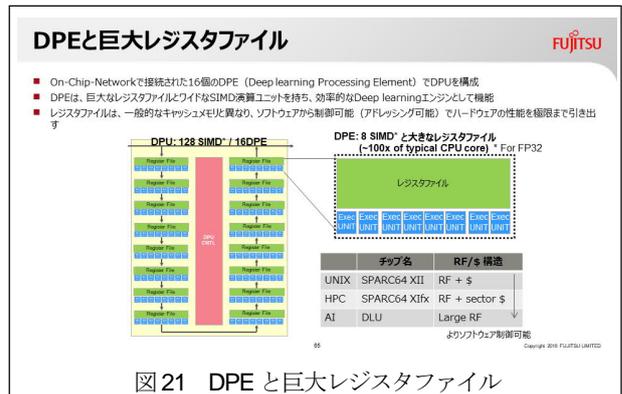


図 21 DPE と巨大レジスタファイル

2つ目の特徴は演算精度です。GPU で 3D 演算するには、8bit 計算では画像が崩れてしまうので、32bit 倍精度の浮動小数点演算で計算しています。そうすると、計算が重くなってしまいます。ディープラーニングでは、それほど精度はいらないということで、精度を落とした計算ができる仕組みを取り入れています (図 22)。

演算精度を落として、最終的に 8bit で計算するようにすると、32bit の 4 倍の計算能力になります。もともと、8bit 計算で表せるのは 256 とおろしかないので、演算精度がガクッと落ちてしまいます。そこで、演算精度は 16bit にしておいて、統計情報

