

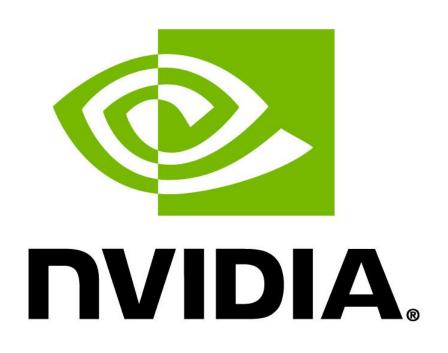
AGENDA

エヌビディアについて

ディープラーニングとは?

GPUディープラーニング/ディープラーニングSDK





創業1993年

共同創立者兼CEO ジェンスン・フアン (Jen-Hsun Huang)

1999年 NASDAQに上場(NVDA)

1999年にGPUを発明 その後の累計出荷台数は1億個以上

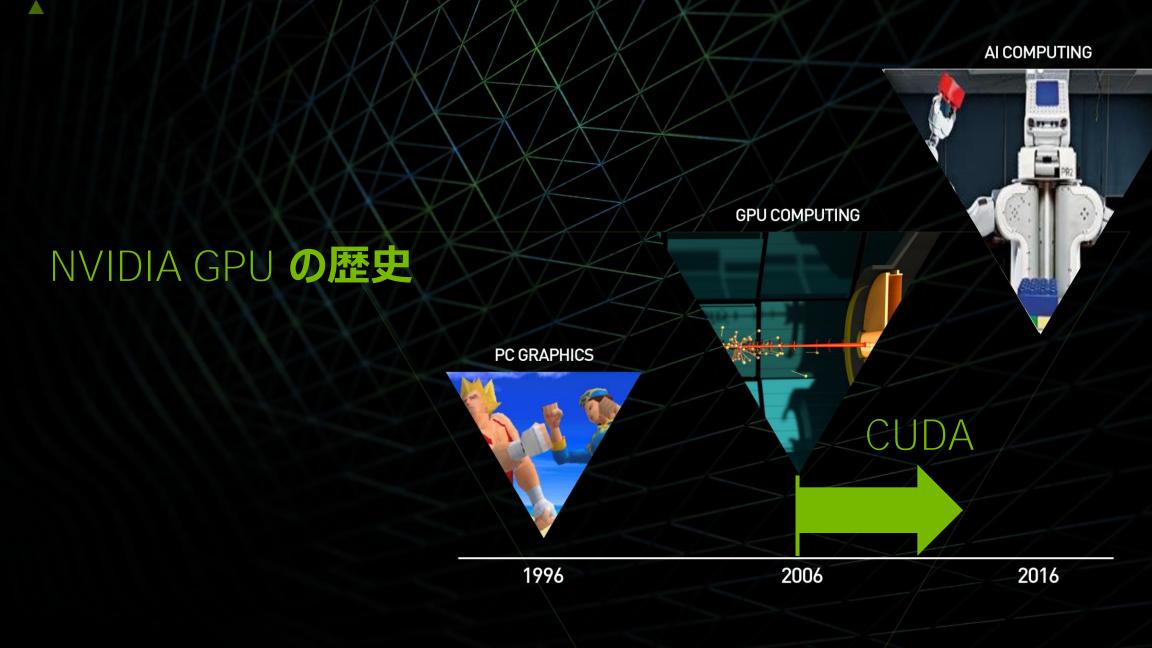
2016年度の売上高は50億ドル

社員は世界全体で9,227人

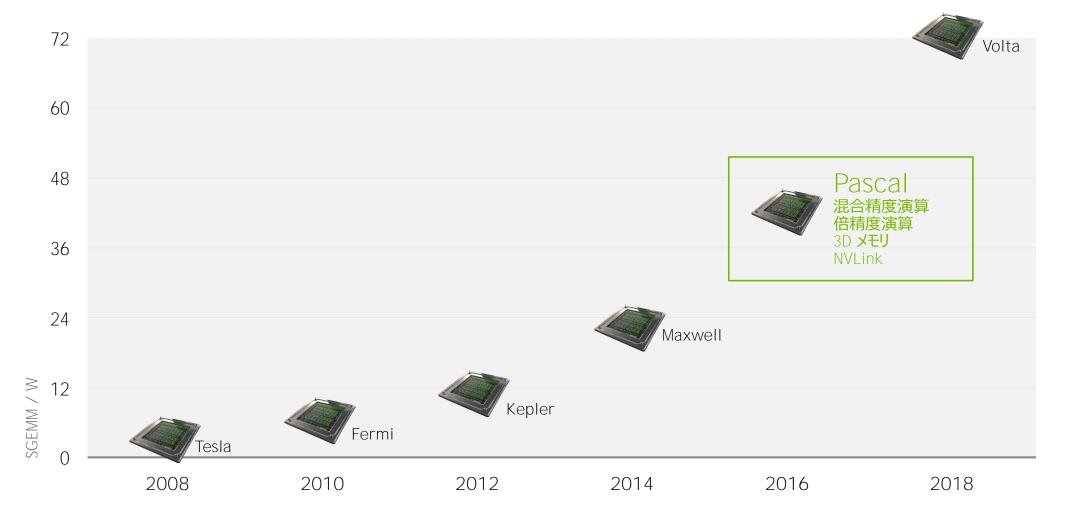
約7,300件の特許を保有

本社は米国カリフォルニア州サンタクララ



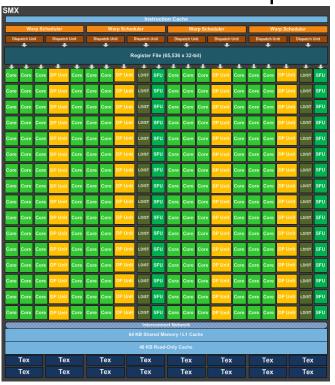


GPU ロードマップ



Compute Capability (CC)









Fermi CC 2.0 32 cores / SM

Kepler CC 3.5 192 cores / SMX

Maxwell CC 5.0 128 cores / SMM

Pascal CC 6.0 64 cores / SMM

https://developer.nvidia.com/cuda-gpus



Tesla K20x (CC 3.5)

2688 CUDA Cores

FP64: 1.31 TF

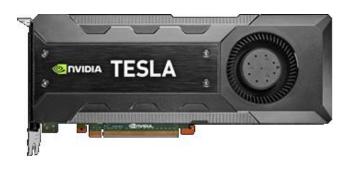
FP32: 3.95 TF

FP16: ...

INT8: ..

GDDR5

384 bit width 6GB 250 GB/s







Tesla P100 SXM2 (CC 6.0)

3584 CUDA Cores

FP64: 5.3 TF

FP32: 10.6 TF

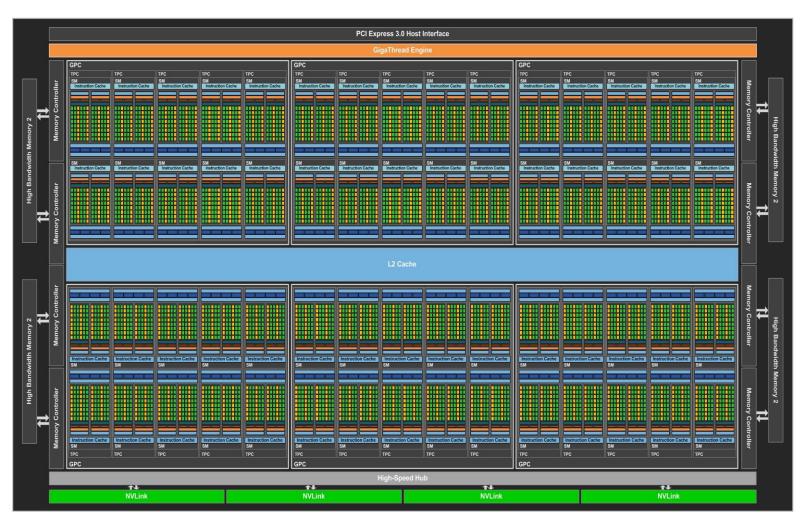
FP16: 21.2 TF

INT8: ...

HBM2

4096 bit width 16 GB 732 GB/s







Tesla P100 SXM2 (CC 6.0)

3584 CUDA Cores

FP64: 5.3 TF

FP32: 10.6 TF

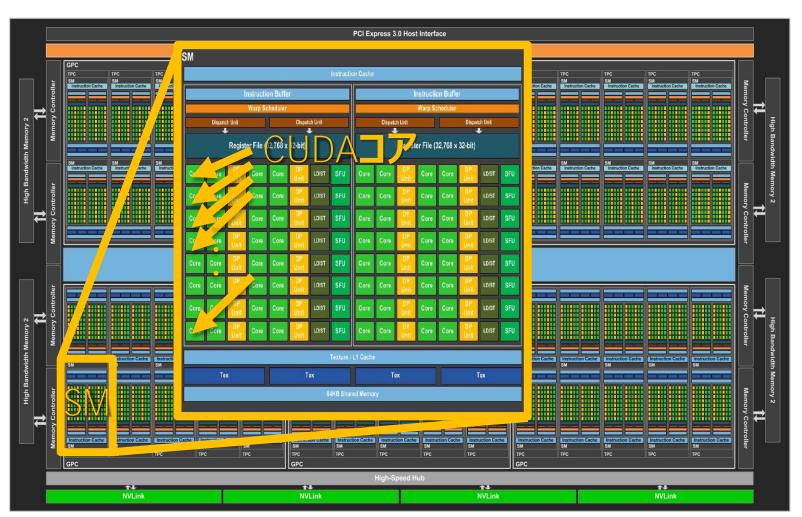
FP16: 21.2 TF

INT8: ...

HBM2

4096 bit width 16 GB 732 GB/s



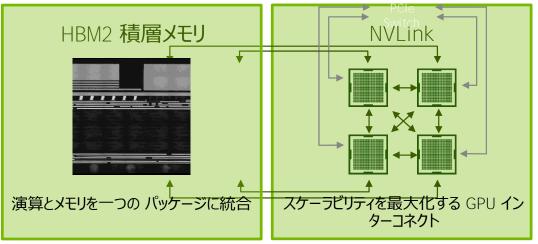


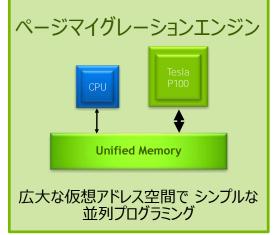


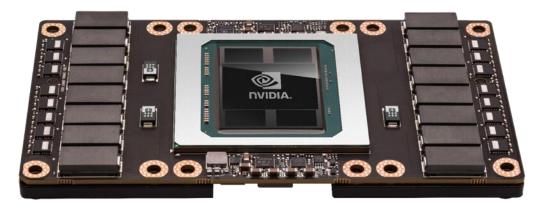
TESLA P100

最速の計算ノードを実現する新しい GPU アーキテクチャ









CUDA

GPU超並列コンピューティングプラットフォーム

CUDA8.0が最新

マルチOSで動作 (Windows, Linux(x86, power, arm), MacOS)

様々なライブラリや開発環境を提供



ライブラリ・ディレクティブ

OpenACC/NVIDIA CUDA Libraries/ 3rd Party Libraries



プログラム言語

C++/Fortran/Python/R/Matlab etc..



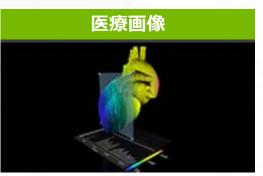
開発環境

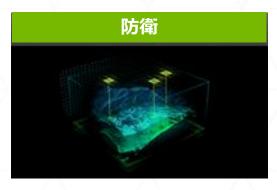
Debugger/Analyzer/IDE etc..

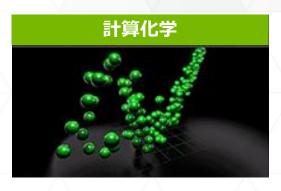
GPUアプリケーションの例

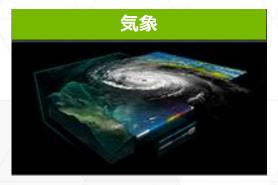
http://www.nvidia.com/object/gpu-applications.html

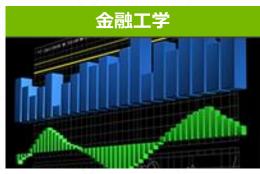




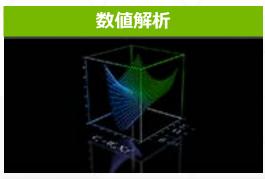












NVIDIA CUDA ライブラリ

cuFFT

フーリエ変換

cuBLAS

行列演算(密行列)

cuSPARSF

行列演算(疎行列)

cuSOLVER

行列ソルバ

cuDNN

ディープラーニング

cuRAND

乱数牛成

Thrust

C++テンプレート(STLベース)

NPP

画像処理・信号処理プリミティブ





GIE NVIDIA GPU Inference Engine is a high performance neural network inference library for deep learning

applications

NVIDIA Performance Primitives is a GPU accelerated library with a very large collection of 1000's of image processing primitives and signal processing primitives.



nvGRAPH Analytics Library is a GPU-

accelerated graph analytics library.

nvGRAPH



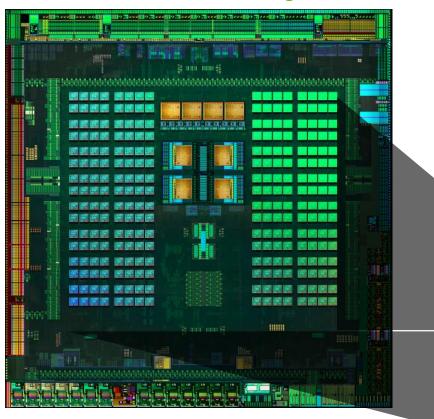


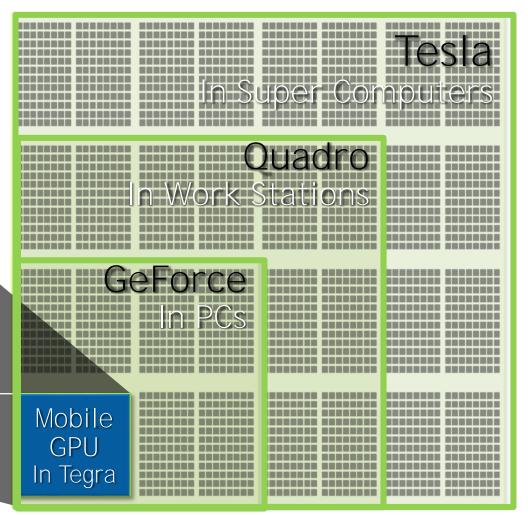


NVIDIA ONE-ARCHITECTURE

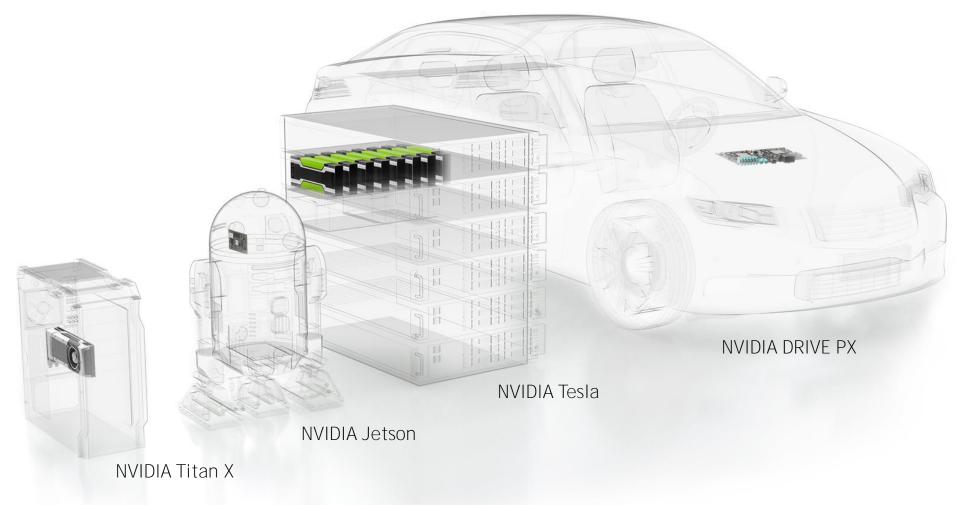
from Super Computer to Automotive SoC

Automotive Tegra



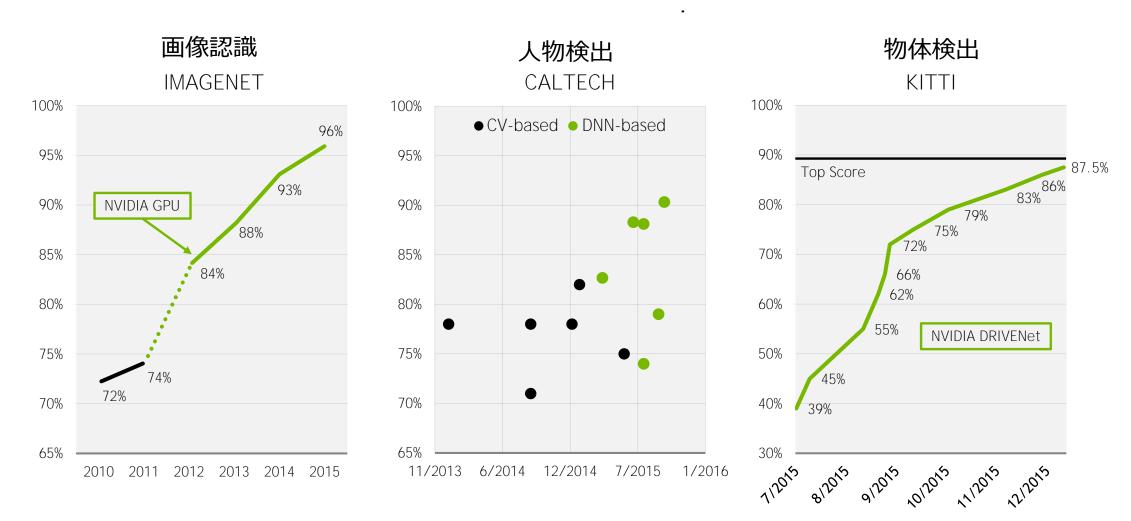


deep learning EVERYWHERE



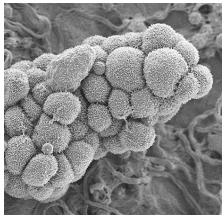


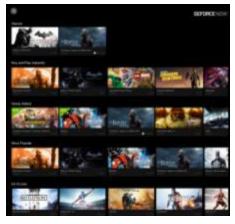
驚くべき認識精度の向上



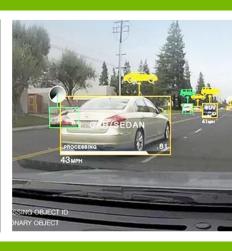
様々な分野でディープラーニングを応用











インターネットとクラウド

画像分類 音声認識 言語翻訳 言語処理 感情分析 推薦

医学と生物学

癌細胞の検出 糖尿病のランク付け 創薬

メディアとエンターテイメント

字幕 ビデオ検索 リアルタイム翻訳

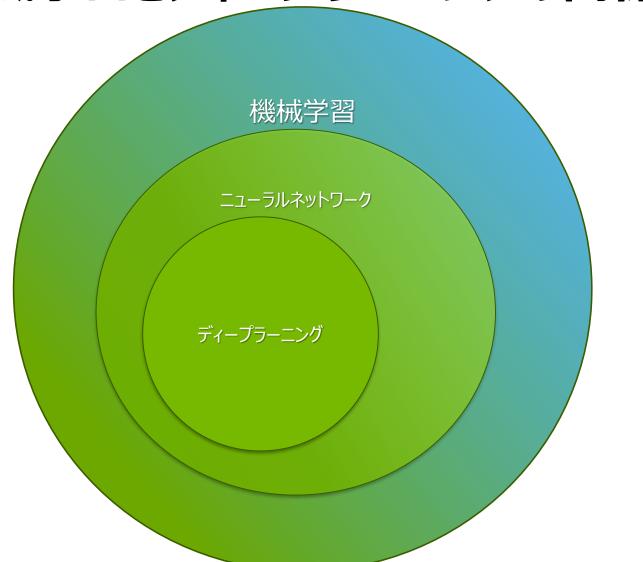
セキュリティと防衛

顔検出 ビデオ監視 衛星画像

機械の自動化

歩行者検出 白線のトラッキング 信号機の認識

機械学習とディープラーニングの関係

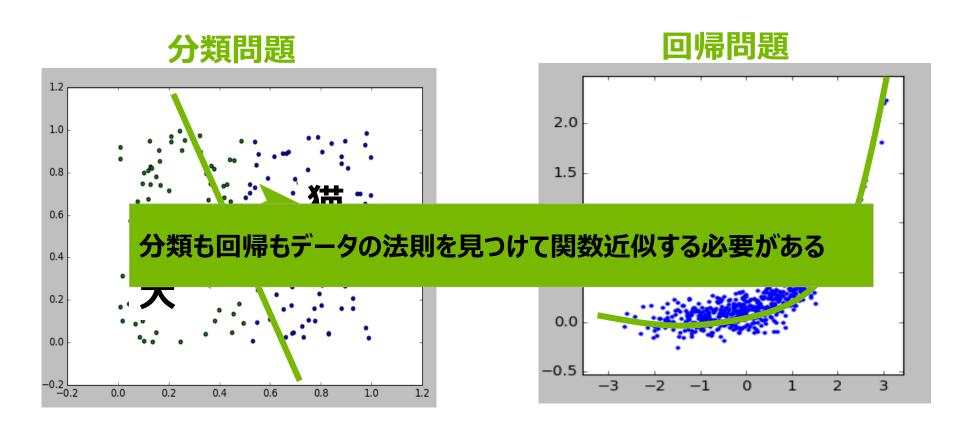


分類(Classification)と回帰(Regression)

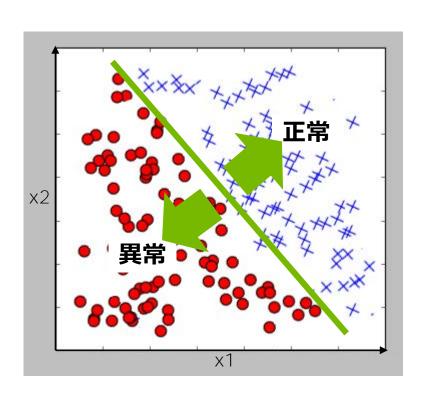
- 分類問題の例
 - 画像に写っているのが猫か犬の画像か推論する

- 回帰問題の例
 - N社の1ヶ月後の株価を予想する

分類(Classification)と回帰(Regression)



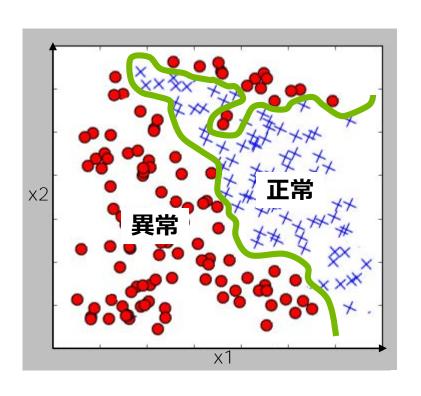
分類問題の例



例) 正常と異常に分類する事を考える 赤:異常 青:正常

$$Y = \alpha X + \beta$$

分類問題の例

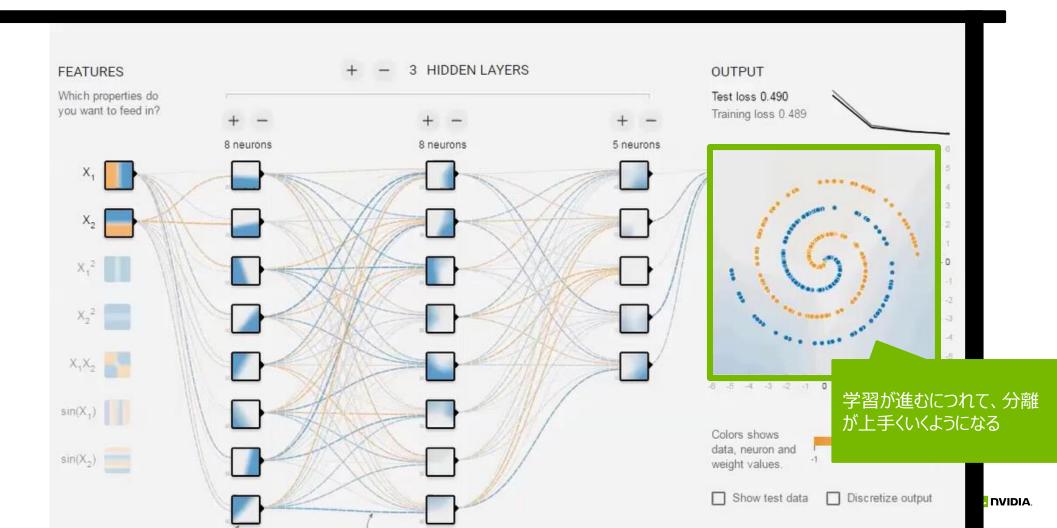


例) 正常と異常に分類する事を考える

赤:異常青:正常

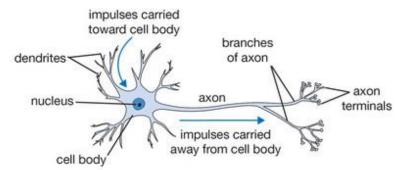
複雑な関数近似が必要になる

ディープラーニング手法による分類の例

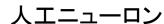


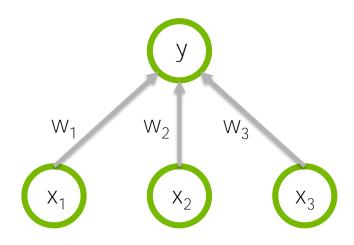
人工ニューロン 神経回路網をモデル化

神経回路網



スタンフォード大学cs231講義ノートより





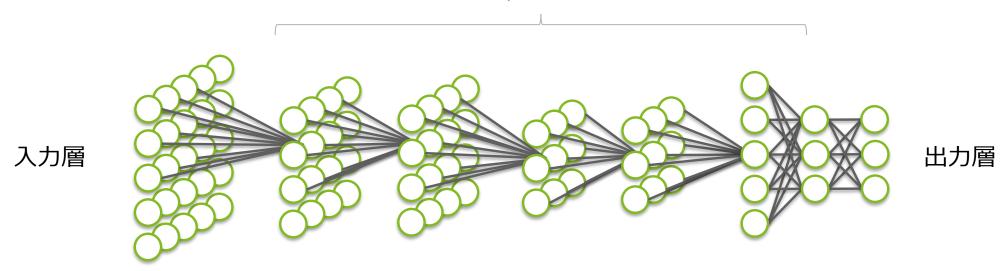
$$y=F(w_1x_1+w_2x_2+w_3x_3)$$

 $F(x)=max(0,x)$

人工ニューラルネットワーク

単純で訓練可能な数学ユニットの集合体 ニューラルネットワーク全体で複雑な機能を学習

隠れ層



十分なトレーニングデータを与えられた人工ニューラルネットワークは、入力データから判断を行う 複雑な近似を行う事が出来る。

ディープラーニングの恩恵

ディープラーニングとニューラルネットワーク

■ ロバスト性

特徴量の設計を行う必要がない。- 特徴は自動的に獲得される学習用データのバラつきの影響を 押さえ込みながら、自動的に学習していく

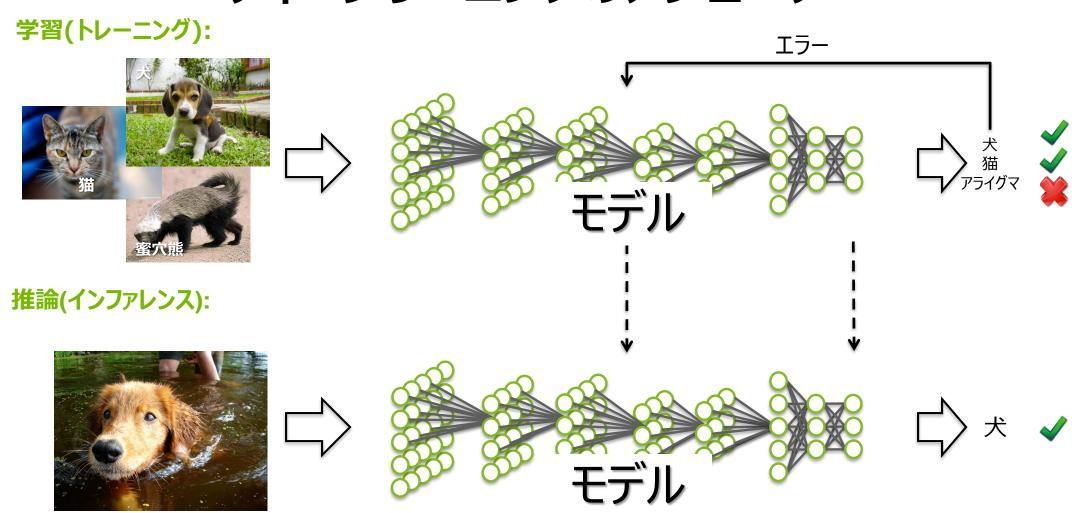
■ 一般性

同じニューラルネットワークのアプローチを多くの異なるアプリケーションやデータに適用する事が出来る

■ スケーラブル

より多くのデータで大規模並列化を行う事でパフォーマンスが向上する

ディープラーニングのアプローチ



畳込みニューラルネットワーク(CNN)

多量なトレーニングデータと多数の行列演算

目的

顔認識

トレーニングデータ

1,000万~1億イメージ

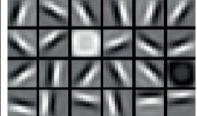
| ネットワークアーキテクチャ | ラーニングアルゴリズム

10層

10 億パラメータ

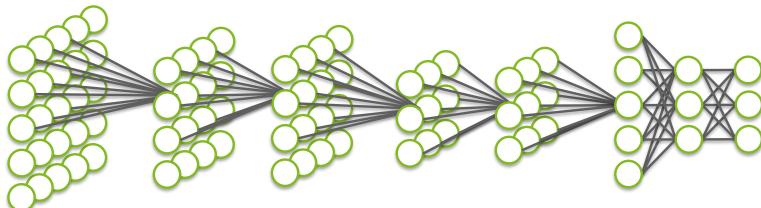
30 エクサフロップスの計算量 GPU **を利用して**30日











畳込みニューラルネットワーク(CNN)

画像認識・画像分類で使われる、高い認識精度を誇るアルゴリズム。畳込み層で画像の特徴を学習



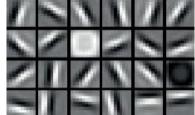


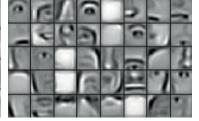


10層 10 億パラメータ

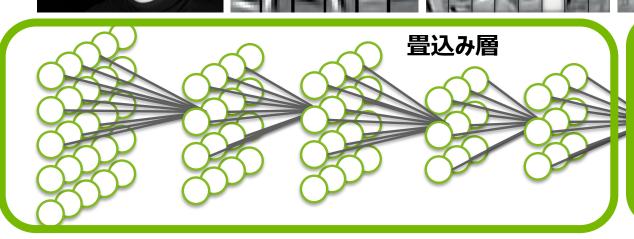
30 エクサフロップスの計算量 GPU **を利用して**30日





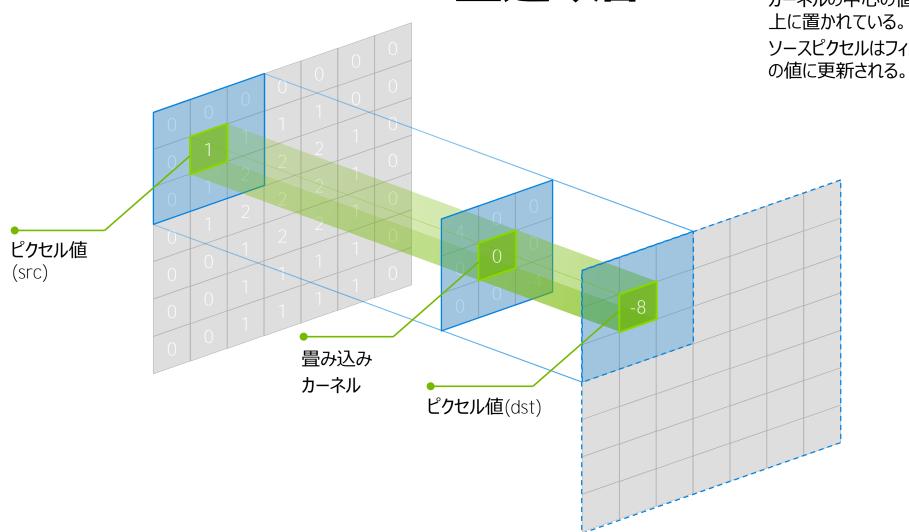








畳込み層



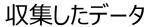
カーネルの中心の値はソースピクセル上に置かれている。 ソースピクセルはフィルタを自身の積

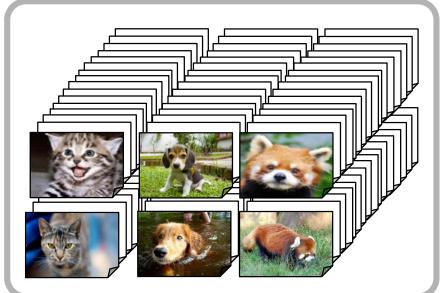


訓練データと検証データ

ディープラーニングの学習

■ データを訓練データ(training data)と検証データ(validation data)に分割する



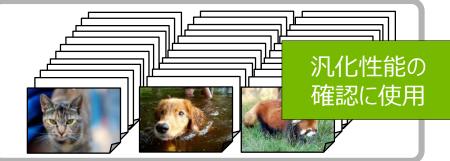








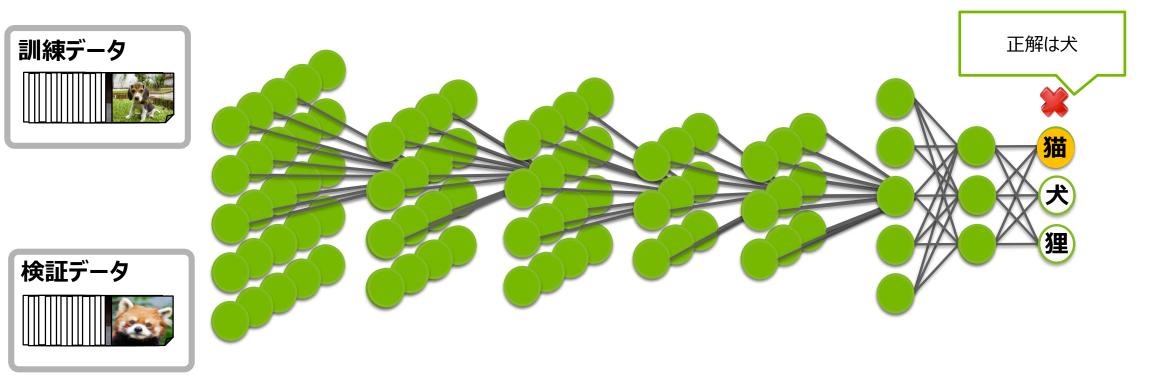




訓練データによる重みの更新

ディープラーニングの学習

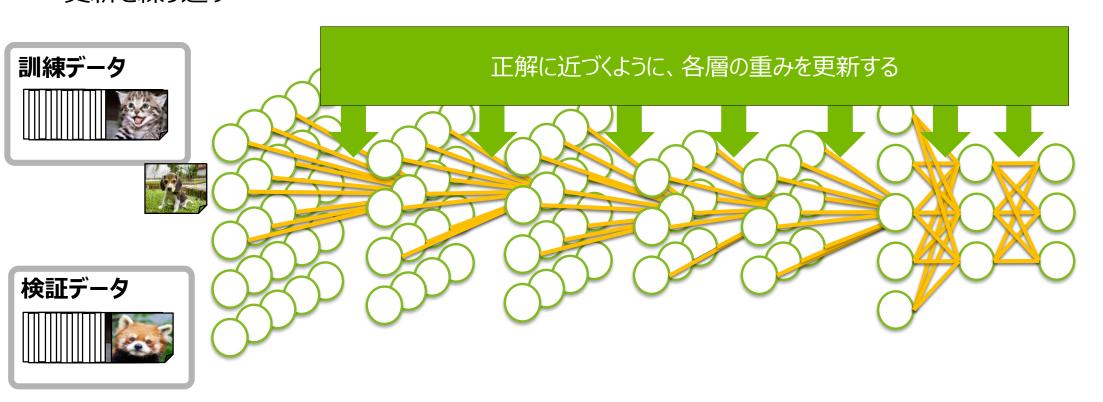
訓練データをニューラルネットワークに与え、正解ラベルと出力結果の誤差が無くなるように重みWの 更新を繰り返す



訓練データによる重みの更新

ディープラーニングの学習

訓練データをニューラルネットワークに与え、正解ラベルと出力結果の誤差が無くなるように重みWの 更新を繰り返す



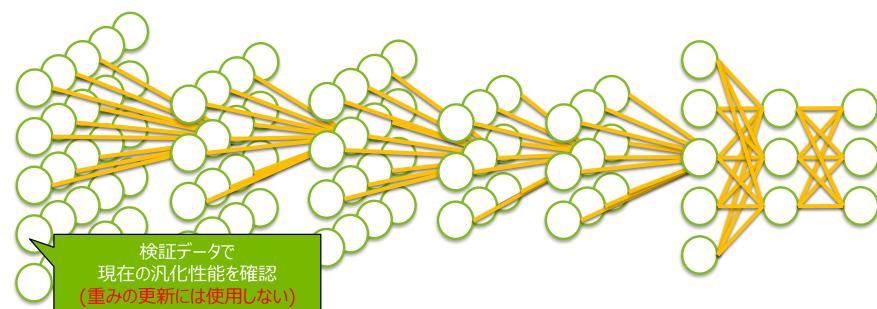
学習ループ

訓練データと検証データの役割

すべての訓練データを用いて重み更新を行う + すべての検証データで汎化性能を確認 ⇒ 1 **エポック**(epoch)**と呼ぶ**



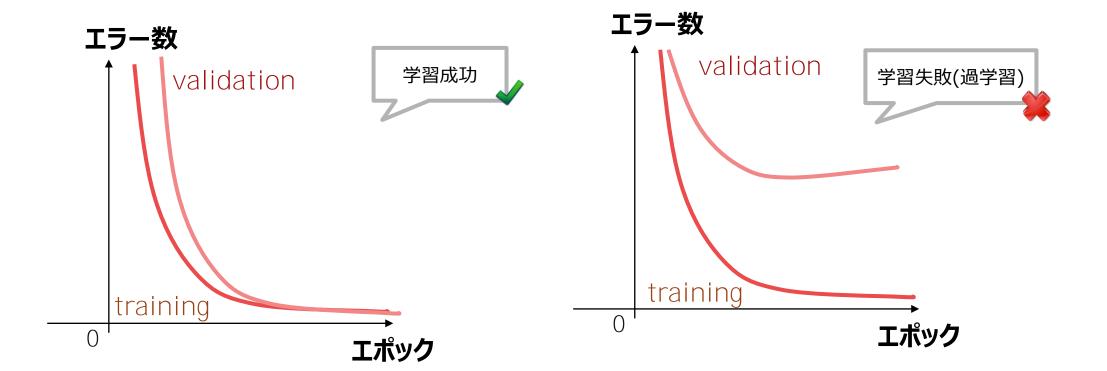




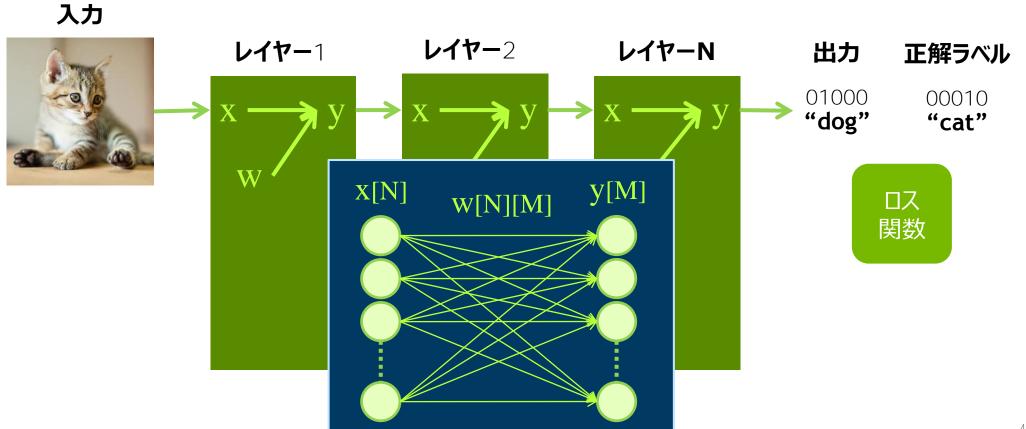
学習時の性能の確認

訓練データと検証データの役割

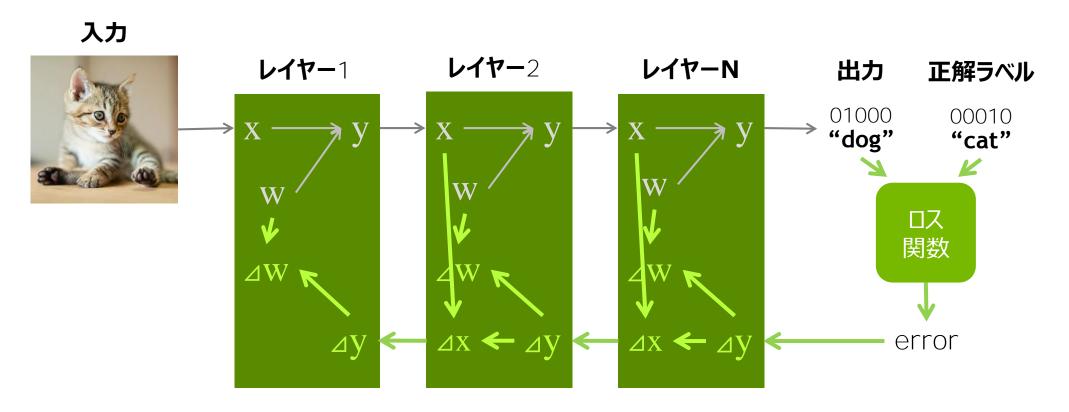
各エポックで訓練データをニューラルネットワークに与えた際の間違い率と検証データを与えた際の間違い率を確認しながら学習を進める必要がある



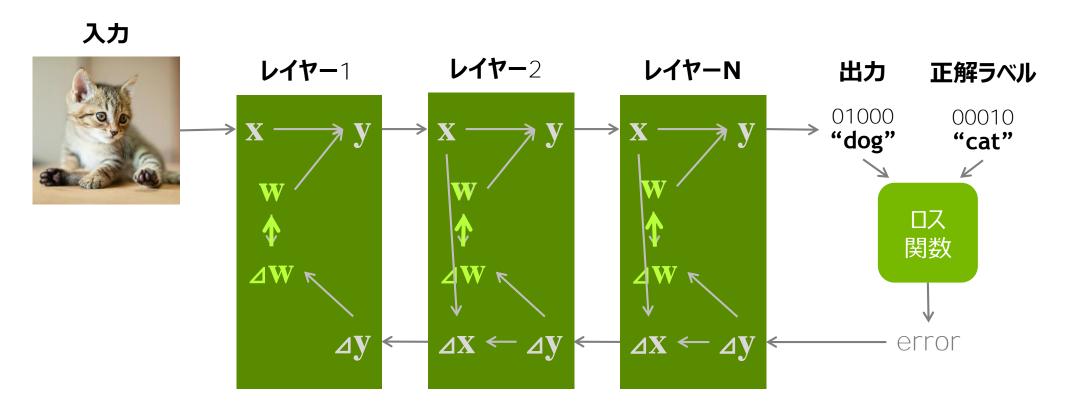
順伝播(フォワードプロパゲーション)



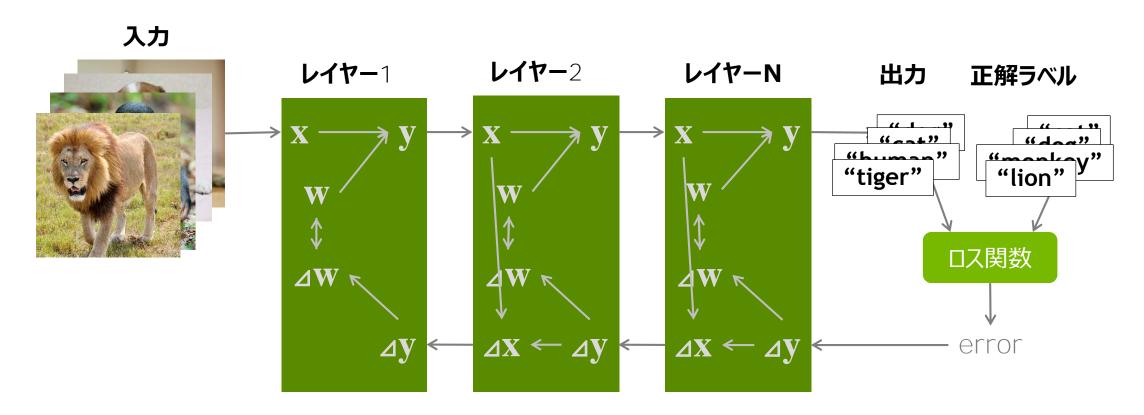
逆伝播(バックフ゜ロパゲーション)



重みの更新

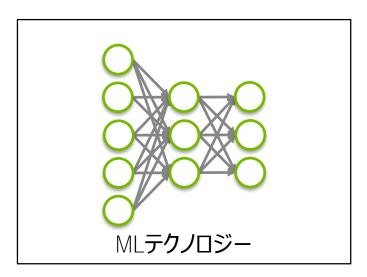


ミニバッチ





ディープラーニングを加速する3つの要因







TORCH	CAFFE	
NYU facebook.	Berkeley UNIVERSITY OF CALIFORNIA	
THEANO	MATCONVNET	
Université de Montréal	UNIVERSITY OF OXFORD	
MOCHA.JL	PURINE	
Massachusetts Institute of Technology	NUS National biowenity of singuise	
MINERVA	MXNET*	
MYU 🛞	W Carnegie UNIVERSITY of Mellon WASHINGTON University	

facebook.

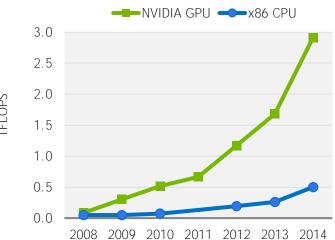
350 million images uploaded per day

Walmart > '<

2.5 trillion transactions per hour

You Tube

100 hour video uploaded per minute



ディープラーニング SDK

ディープラーニングを加速するディープラーニングライブラリ

















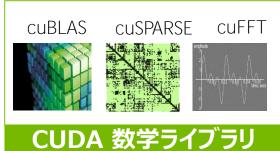


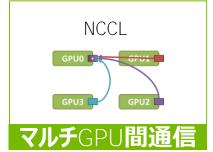




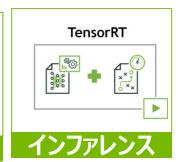












LINPACK

Benchmark Measures floating point computing power

Accelerated Features	Metric	Scalability
AII	GFLOPS	Multi-GPU, Multi Node

https://www.top500.org/project/linpack/

Linpack 2.1 Speedup Vs Dual-Socket CPU Server



CPU Server: Dual Xeon E5-2699 v4@2.2GHz (44-cores)
GPU Servers: Dual Xeon E5-2699 v4@2.2GHz (44-cores) with Tesla K80s or P100s PCle

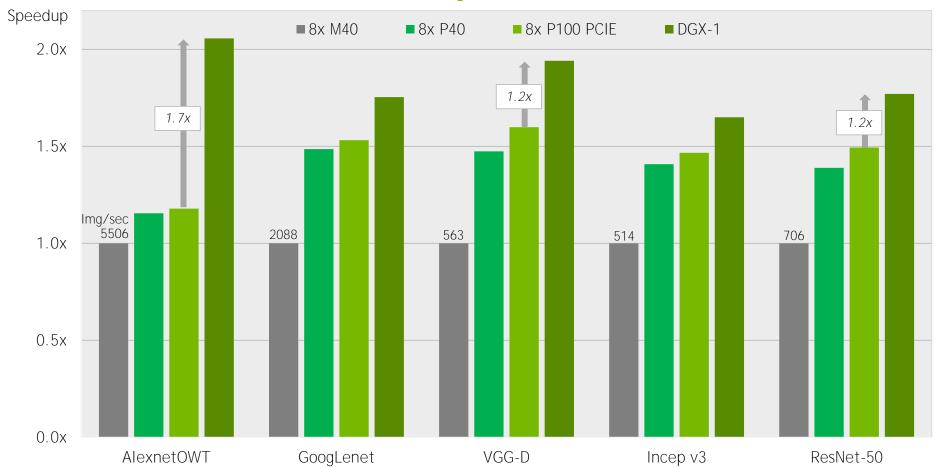
CUDA Version: CUDA 8.0.44

Dataset: HPL.dat



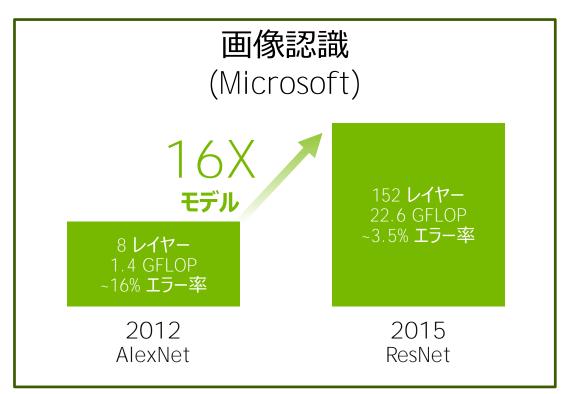
P100 SXM2**による高速な学習**

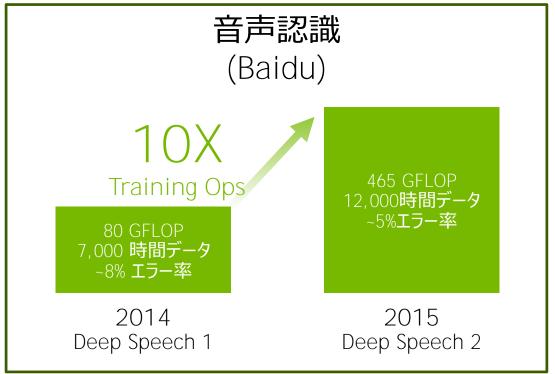
FP32 Training, 8 GPU Per Node



認識精度向上のため モデルはよりディープに、データはより大きく

強力な計算パワーが必要に





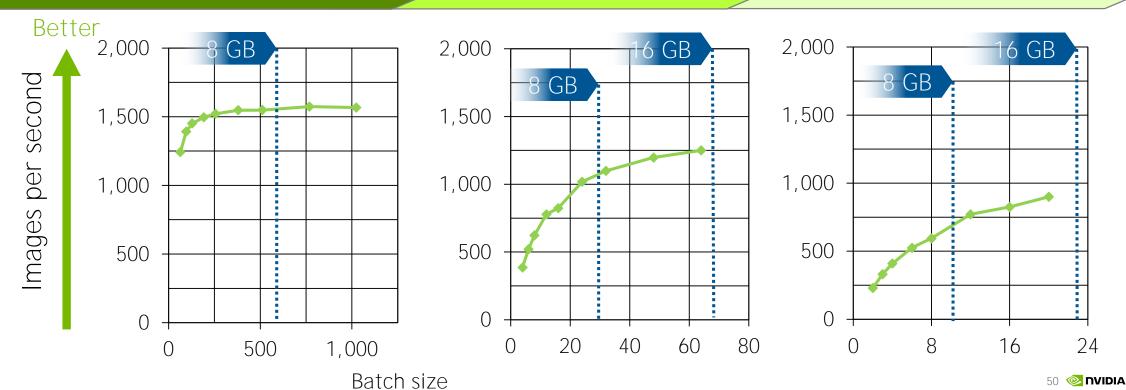
バッチサイズと計算性能

POOL NO CO

TESLA P100, Chainer 1.17.0



VGG-D (16 layers) 2014 ResNet (152 layers) 2015



FP16

[TESLA P100] FP16: 半精度浮動小数点

IEEE 754**準拠**

- 16-bit s e x p . m a n t i s s a
 - 符号部:1-bit, 指数部:5-bits, 仮数部:10-bit
- ダイナミックレンジ: 2⁴⁰
 - Normalized values: $2^{-14} \sim 65504 \ (\approx 2^{16})$
 - Subnormal at full speed

ユースケース

ディープラーニング

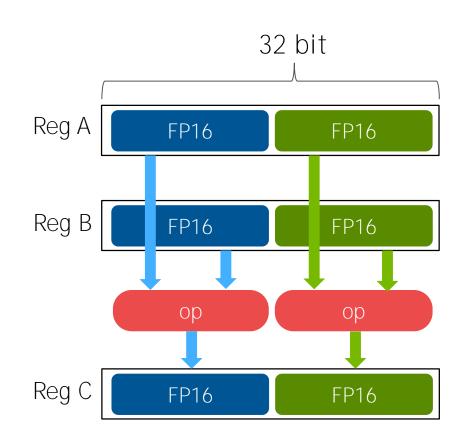
画像処理

信号処理



[TESLA P100] FP16: 半精度浮動小数点

- 2-way SIMD
 - •32ビットレジスタ内に2x FP16
 - メモリフットプリント: FP32の半分
- Instructions
 - Add, Sub, Mul
 - Fma (Fused Multiply and Add)
- •1サイクルあたり4つの算術
 - •FP32より2倍高速



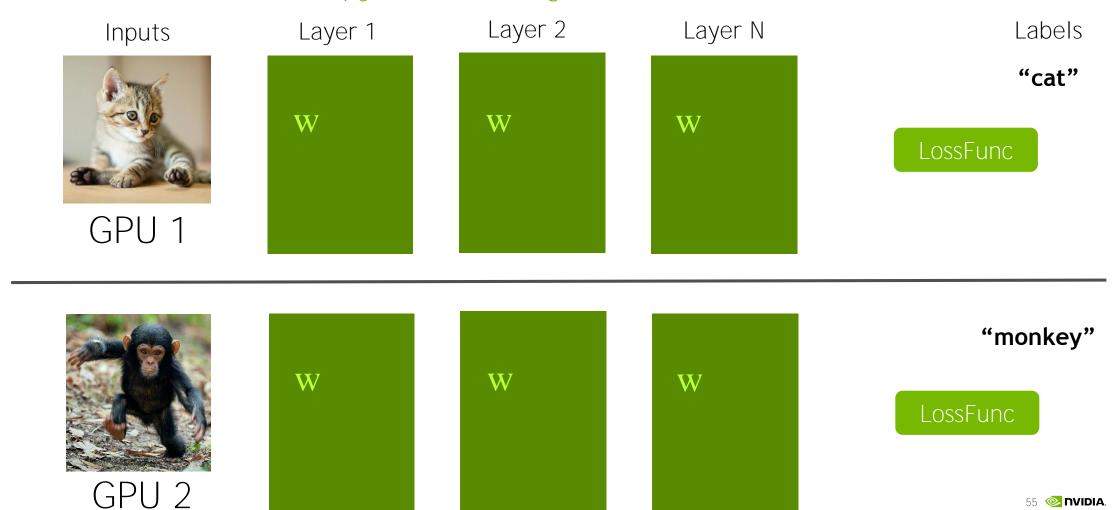
cuda_fp16.h

cuBLAS

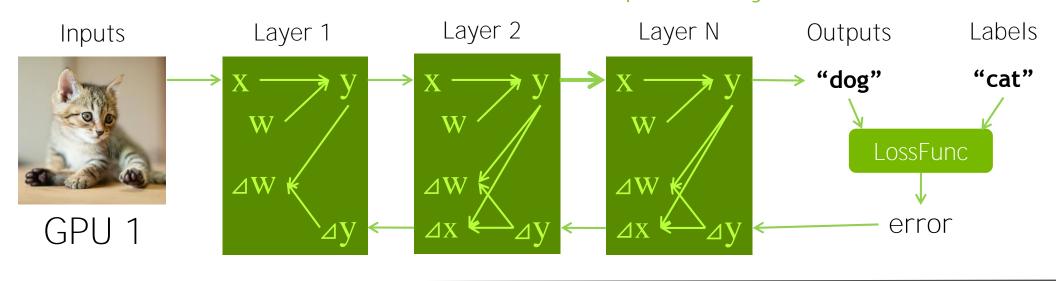


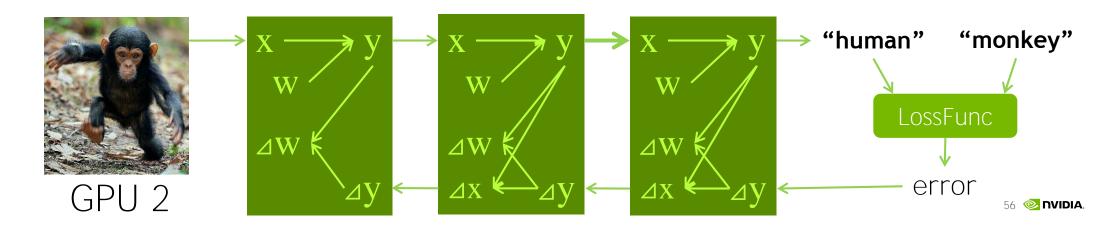
マルチGPU学習

Copy Model, Assigne different data

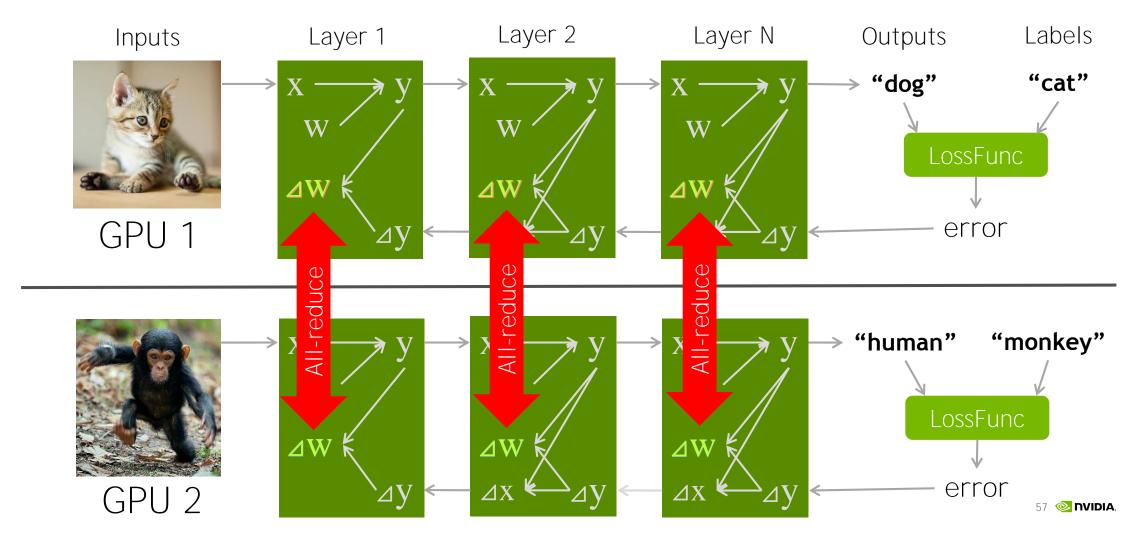


Forward & Backward Independently

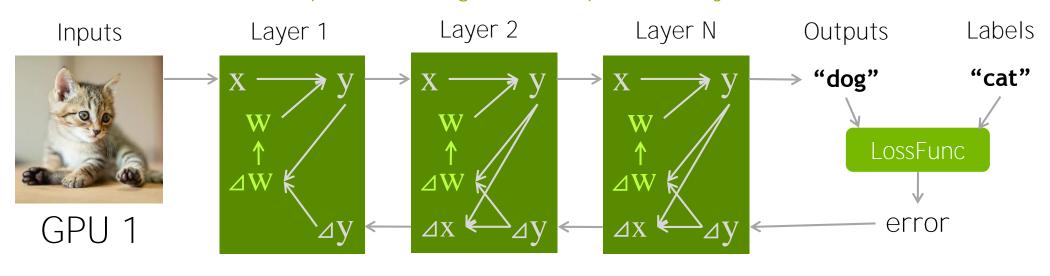


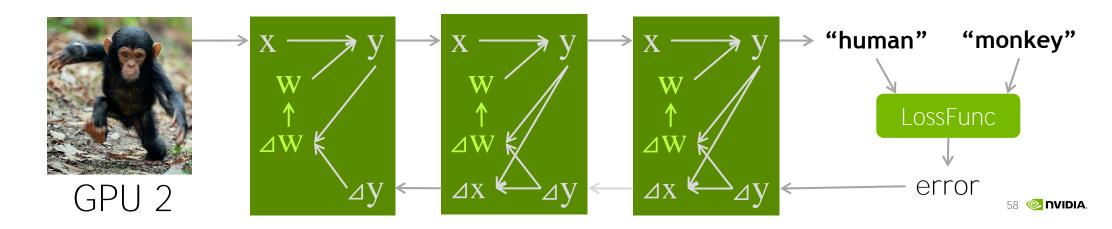


Combine ⊿w over multi-GPU



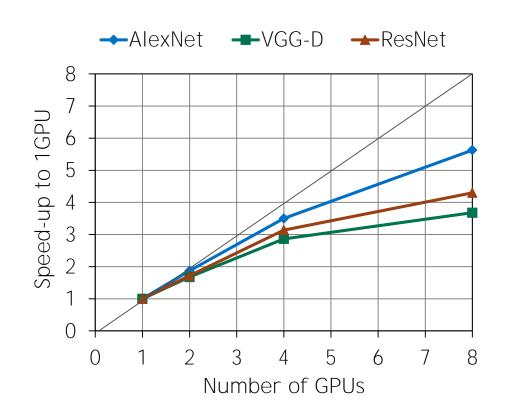
Update Weights Independently





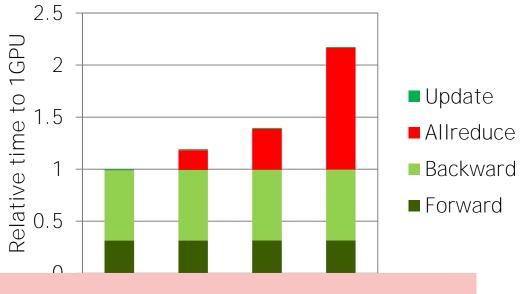
マルチGPU学習のパフォーマンス

NVIDIA DGX-1, Chainer 1.17.0 with multi-process patch



[Batch size per GPU] AlexNet:768, VGG-D:32, ResNet:12

Time per one batch (VGG-D)

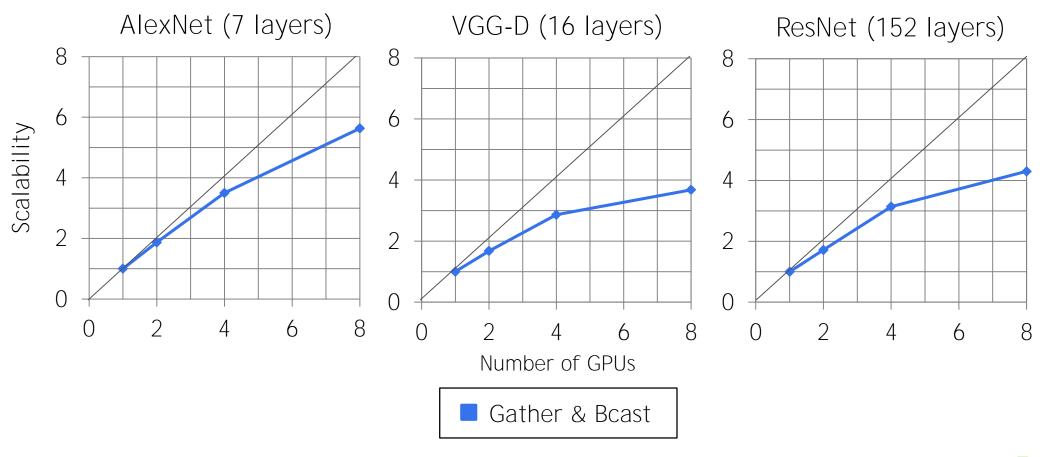


DGX-1's NVLink is not well utilized. Chainer's all-reduce implementation is naïve "gather and broadcat".

INVIDIA

マルチGPU学習のパフォーマンス(NCCL使用なし)

NVIDIA DGX-1, Chainer 1.17.0 with multi-process patch



NCCL

ディープラーニング SDK

ディープラーニングを加速するディープラーニングライブラリ



















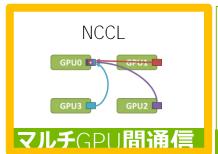




ディープラーニングフレームワーク











NCCL(NVIDIA Collective Collection Library) ディープラーニング SDK

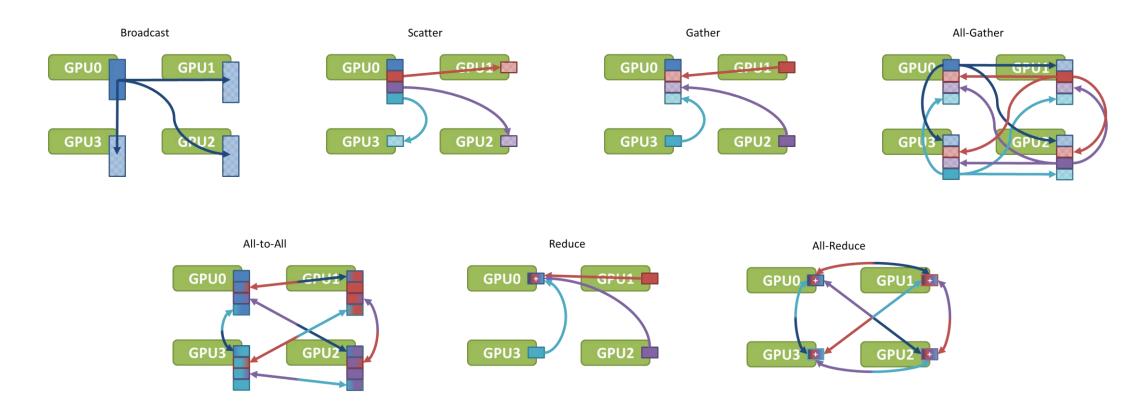
マルチGPU集合通信ライブラリ

- 最新リリースはv1.2.3
- https://github.com/NVIDIA/nccl

all-gather, reduce, broadcast など標準的な集合通信の処理をバンド幅が出るように最適化シングルプロセスおよびマルチプロセスで使用する事が可能

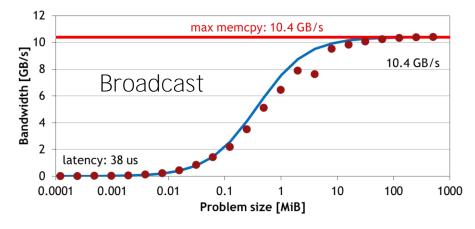
NCCL(NVIDIA Collective Collection Library)

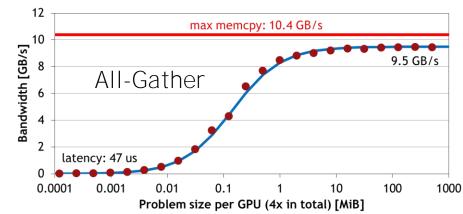
NCCLの集合通信処理

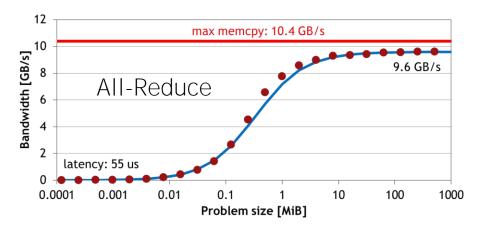


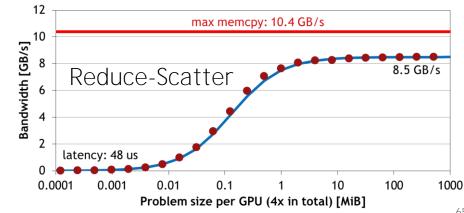
NCCL パフォーマンス

Bandwidth at different problem sizes (4 Maxwell GPUs)





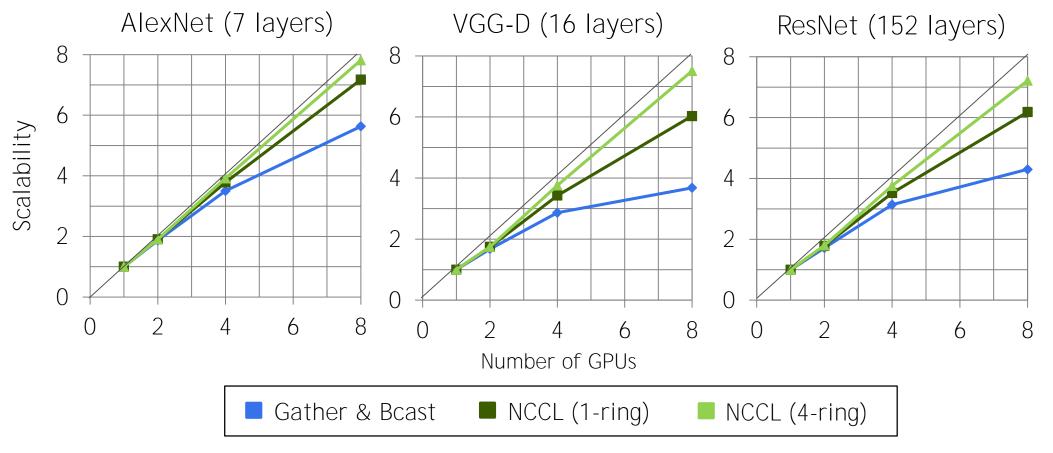




マルチGPU学習のパフォーマンス(NCCL使用)



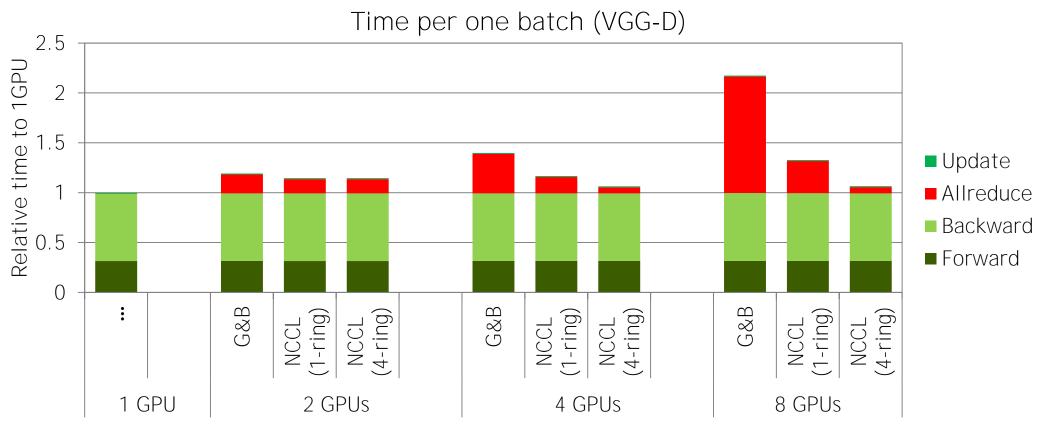
NVIDIA DGX-1, Chainer 1.17.0 with NCCL patch





マルチGPU学習のパフォーマンス(NCCL使用)

NVIDIA DGX-1, Chainer 1.17.0 with NCCL patch



ディープラーニング SDK

ディープラーニングを加速するディープラーニングライブラリ

























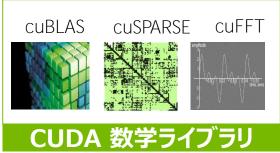
MINERVA

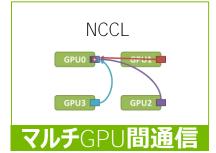


Pylearn2

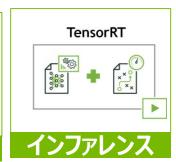
ディープラーニングフレームワーク











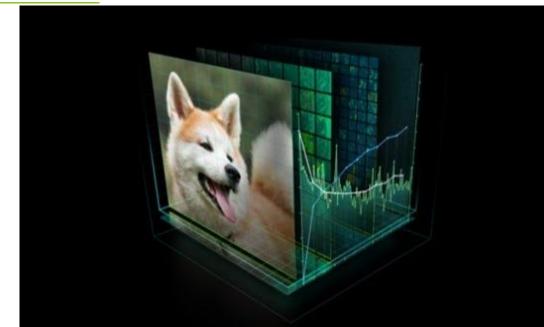
NVIDIA ディープラーニング学習コース

NVIDIA ディープラーニング・インスティチュート

ディープラーニングの為の自習型のクラス。ハンズオンラボ、講義資料、講義の録画を公開。ハンズオンラボは日本語で受講可能

https://developer.nvidia.com/deep-learning-courses

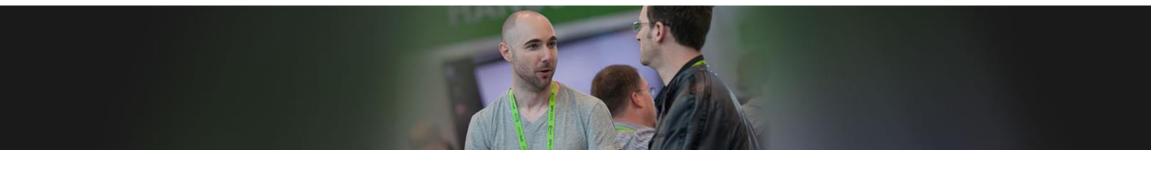
- 1. ディープラーニング入門
- 2. DIGITSによる画像分類入門
- 3. DIGITSによる物体検出入門
- 4. etc...





GTC 2017 参加登録受付中

2017/5/8 - 11 サンノゼで開催



基調講演

テクニカルセッション

ハンズオンラボ

ポスター展示

専門家との交流

スペシャルイベント

