# (Big) Data, (Deep) Learning and AI

## When big data hits machine learning

Phạm Thành Lâm | Founder @ SaigonApps

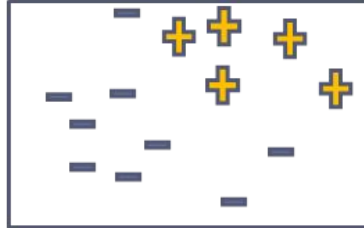17.09.2016

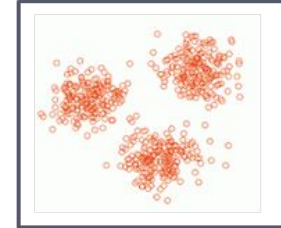**Big Picture: Big Data - Machine Learning/ Data Mining**

# History of AI, Machine Learning and Deep Learning



Image taken from:blogs.nvidia.com

**From Program To Machine/Deep Learning**



[Honglak Lee]

## Learned Feature Hierarchy



Image taken from the Internet

# Gartner Hype Cycle Emerging Technology 2015



Image taken from: http://bit.ly/1i4e8oL, kdnuggets.com

# In 2016, is Big Data still a "thing"?

- Enterprise Technology = building a data-driven culture, where Big Data is not "a" thing, but "the" thing
- The Ecosystem is Maturing (let see the picture)
- Big Data infrastructure:  Still Plenty of Innovation
- Big Data Analytics: Now with AI

Credited by http://bit.ly/1UIgzeJ

# Big Data Landscape 2016 (Version 2.0)

Last Updated 2/12/2016 © Matt Turck (@mattturck), Jim Hao (@jimrhao), & FirstMark Capital (@firstmarkcap)

FIRSTMARK

vietnamworks
Move Up!

# AI: **A**rtificial **I**ntelligence → **A**pplications and **I**nnovations

# Top acquirers AI startups MA Timeline



Apple Adds Startup 'Turi' To AI Arsenal, Pays $200M Sources Say

Siri.Vasile Cotovanu/Flickr

Race For AI: Most Active Acquirers In Artificial Intelligence

Jan 2012    Jan 2013    Jan 2014    Jan 2015    Jan 2016

Date of acquisition

CB INSIGHTS

www.cbinsights.com

Image taken from: cbinsights.com, econotimes.com

# Tech giants(FAGA) embracing AI

| Google | Facebook | Microsoft | Other |
|---|---|---|---|
| - **TensorFlow** DL framework and Tensor Processing Unit (TPU), a custom ASIC chip built specifically for machine learning<br><br>- 100+ different teams working on Google Today, Street View, Inbox Smart Reply, voice search, Google Play, etc.<br>- **Magenta** to play music<br>- **DeepDream** for creative pictures<br>- WaveNets: speech synthesis, music creator | - Fblearner **Flow** the tool, designed to help engineers build, test and execute machine learning assembly lines, is available to every engineer within the organisation like **Deep Text**<br><br>- Messenger platform, allowing businesses to create AI-powered chatbots to interact with their customers | - **Tay**, an artificial intelligence Twitter chatterbot, released by Microsoft in March<br>- **Cortana** – its equivalent to Apple's Siri and Android's Google Now – an artificial intelligence-powered personal assistant and knowledge navigator for Windows' Phones<br>- London-based AI startup **Swiftkey** is acquired in February | **Amazon**: unveiling DSSTNE, an open-source AI framework developed to run its recommendation system<br><br>**IBM**: Watson/Connie, IBM's AI computer system is able to answer questions posed in natural language, Bluemix apis.<br><br>**Sony**: undisclosed investment in Cogitai, a one year old California-based AI startup |

Info is curated from: techcitynews.com

vietnam**works**
*Move Up!*

# The pioneers of AI/ML/DL: (my bias)

**GodFather of DL, IEEE awarded 2016**

Geoffrey Hinton - Google

Yann Lecun – FB

Bengio Yoshua - Montreal University

Xavier Amatriain – Quora/Netflix

Demis Hassabis – DeepMind

Andrew Ng- Baidu

vietnam**works**
*Move Up!*

# Real world AI/DL applications



Image taken from: Luong's Machine Translation slide

# CẦM KỲ THI HOẠ

Prisma

Alpha Go

**THE ULTIMATE GO CHALLENGE**
GAME 1 OF 5

9 MARCH 2016

AlphaGo  VS  Lee Sedol

RESULT: W+Res  NUMBER OF MOVES: 186  TIME WHITE: 1h 55m  TIME BLACK: 1h 32m

Google Brain

Magenta

Image taken from:tuoitre.com

# Sample poetry

No.

he said.
"no," he said.
"no," i said.
"i know," she said.
"thank you," she said.
"come with me," she said.
"talk to me," she said.
"don't worry about it," she said.

Image taken from: Andrew Ng twitter

# Limits and challenges of DL/ML



Image taken from Internet: wsj.com, twitter.com

# Training DL is painful

- Tuning hyperparameters
- Network architecture: layers/nodes
- Some data preprocessing
- Weight initialization: $\sim N(0,1)$
- Learning rate, optimization algos
- Slowness
- Overfitting
- More ...

Yann LeCun

We know now that we don't need any big new breakthroughs to get to true AI
That is completely, utterly, ridiculously wrong.
As I've said in previous statements: most of human and animal learning is unsupervised learning. If intelligence was a cake, unsupervised learning would be the cake, supervised learning would be the icing on the cake, and reinforcement learning would be the cherry on the cake. We know how to make the icing and the cherry, but we don't know how to make the cake. We need to solve the unsupervised learning problem before we can even think of getting to true AI. And that's just an obstacle we know about. What about all the ones we don't know about?

# How to build ML/DL from scratch



Image taken from Internet

# Open source/Frameworks

## Top libraries by Github issues opened

| # | Count | Library |
|---|-------|---------|
| #1: | 2908 | BVLC/caffe |
| #2: | 2530 | fchollet/keras |
| #3: | 2456 | tensorflow/tensorflow |
| #4: | 1801 | dmlc/mxnet |
| #5: | 1705 | Theano/Theano |
| #6: | 1067 | deeplearning4j/deeplearning4j |
| #7: | 693 | Microsoft/CNTK |
| #8: | 505 | mila-udem/blocks |
| #9: | 498 | pfnet/chainer |
| #10: | 494 | NVIDIA/DIGITS |
| #11: | 394 | Lasagne/Lasagne |
| #12: | 342 | torch/torch7 |
| #13: | 233 | NervanaSystems/neon |
| #14: | 206 | tflearn/tflearn |
| #15: | 82 | IDSIA/brainstorm |
| #16: | 41 | karpathy/convnetjs |
| #17: | 39 | amznlabs/amazon-dsstne |
| #18: | 27 | torchnet/torchnet |

## Top libraries by Github stars

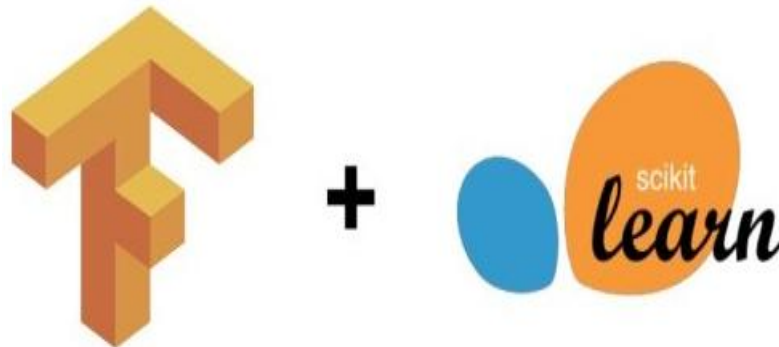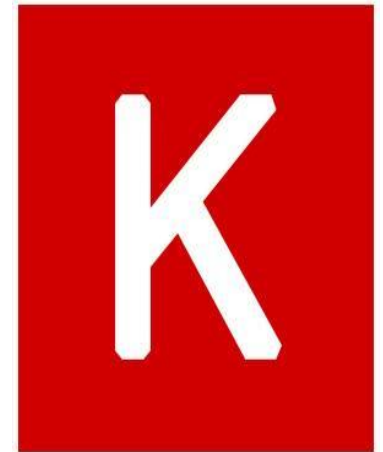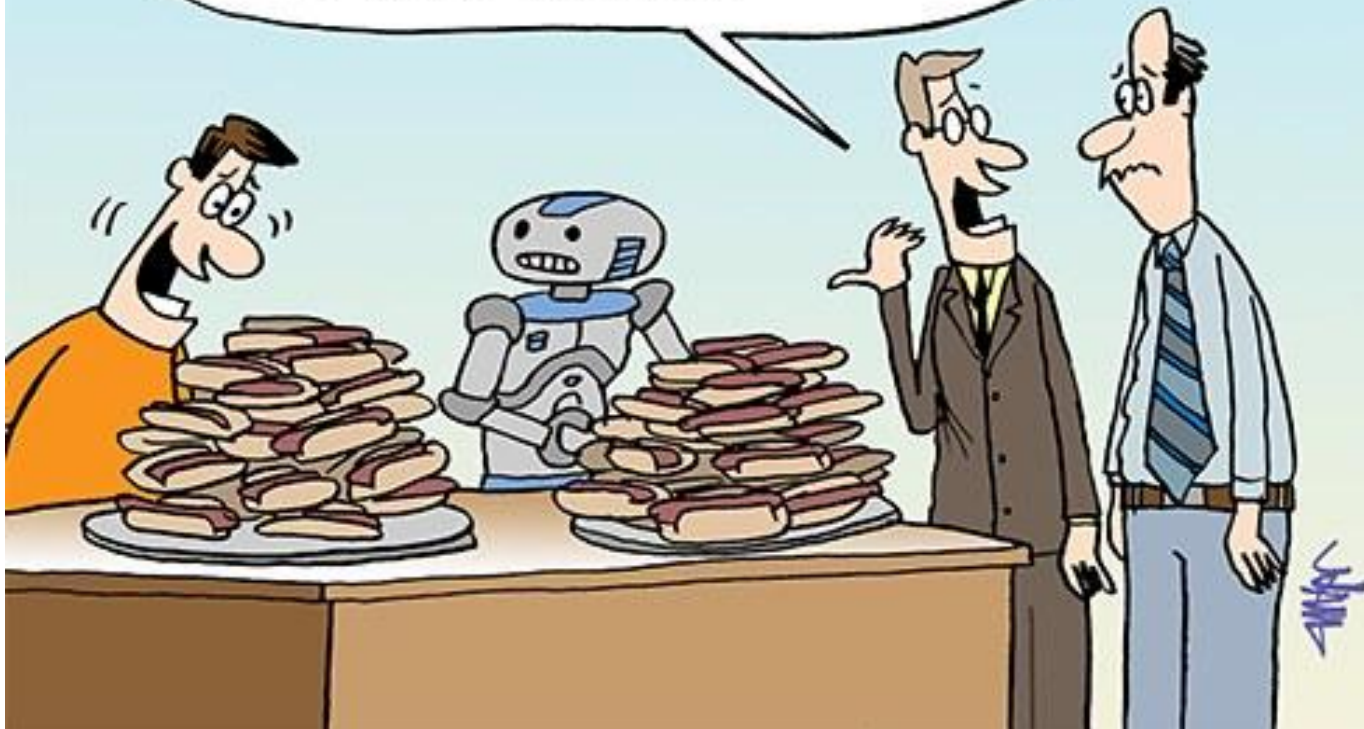| # | Count | Library |
|---|-------|---------|
| #1: | 29967 | tensorflow/tensorflow |
| #2: | 11914 | BVLC/caffe |
| #3: | 7595 | fchollet/keras |
| #4: | 5985 | Microsoft/CNTK |
| #5: | 5263 | karpathy/convnetjs |
| #6: | 5160 | torch/torch7 |
| #7: | 4740 | dmlc/mxnet |
| #8: | 4316 | Theano/Theano |
| #9: | 3723 | deeplearning4j/deeplearning4j |
| #10: | 3420 | tflearn/tflearn |
| #11: | 3162 | amznlabs/amazon-dsstne |
| #12: | 2372 | Lasagne/Lasagne |
| #13: | 2149 | NervanaSystems/neon |
| #14: | 1577 | pfnet/chainer |
| #15: | 1371 | NVIDIA/DIGITS |
| #16: | 1147 | IDSIA/brainstorm |
| #17: | 870 | mila-udem/blocks |
| #18: | 787 | torchnet/torchnet |

## Top libraries by Github contributors

| # | Count | Library |
|---|-------|---------|
| #1: | 348 | tensorflow/tensorflow |
| #2: | 244 | Theano/Theano |
| #3: | 234 | fchollet/keras |
| #4: | 202 | BVLC/caffe |
| #5: | 169 | dmlc/mxnet |
| #6: | 102 | torch/torch7 |
| #7: | 84 | deeplearning4j/deeplearning4j |
| #8: | 75 | Microsoft/CNTK |
| #9: | 72 | pfnet/chainer |
| #10: | 50 | Lasagne/Lasagne |
| #11: | 48 | mila-udem/blocks |
| #12: | 42 | NervanaSystems/neon |
| #13: | 39 | tflearn/tflearn |
| #14: | 28 | NVIDIA/DIGITS |
| #15: | 16 | amznlabs/amazon-dsstne |
| #16: | 15 | IDSIA/brainstorm |
| #17: | 14 | karpathy/convnetjs |
| #18: | 10 | torchnet/torchnet |

## Top libraries by Github forks

| # | Count | Library |
|---|-------|---------|
| #1: | 12506 | tensorflow/tensorflow |
| #2: | 7194 | BVLC/caffe |
| #3: | 2275 | fchollet/keras |
| #4: | 1777 | dmlc/mxnet |
| #5: | 1540 | Theano/Theano |
| #6: | 1484 | torch/torch7 |
| #7: | 1291 | Microsoft/CNTK |
| #8: | 1264 | deeplearning4j/deeplearning4j |
| #9: | 1024 | karpathy/convnetjs |
| #10: | 662 | Lasagne/Lasagne |
| #11: | 482 | amznlabs/amazon-dsstne |
| #12: | 450 | NervanaSystems/neon |
| #13: | 412 | NVIDIA/DIGITS |
| #14: | 377 | pfnet/chainer |
| #15: | 336 | tflearn/tflearn |
| #16: | 267 | mila-udem/blocks |
| #17: | 161 | torchnet/torchnet |
| #18: | 108 | IDSIA/brainstorm |

# Some Demos

Find me: @laampt | Github: lampts

THANK
YOU!