# Interpretability of Deep Learning Models

**Devinder Kumar**
PhD Candidate, UWaterloo & Vector Inst. for AI
Lead AI Scientist in Residence, NextAI
July 12th, 2018

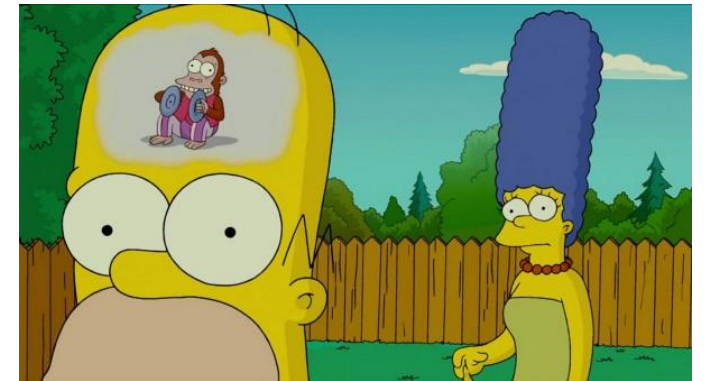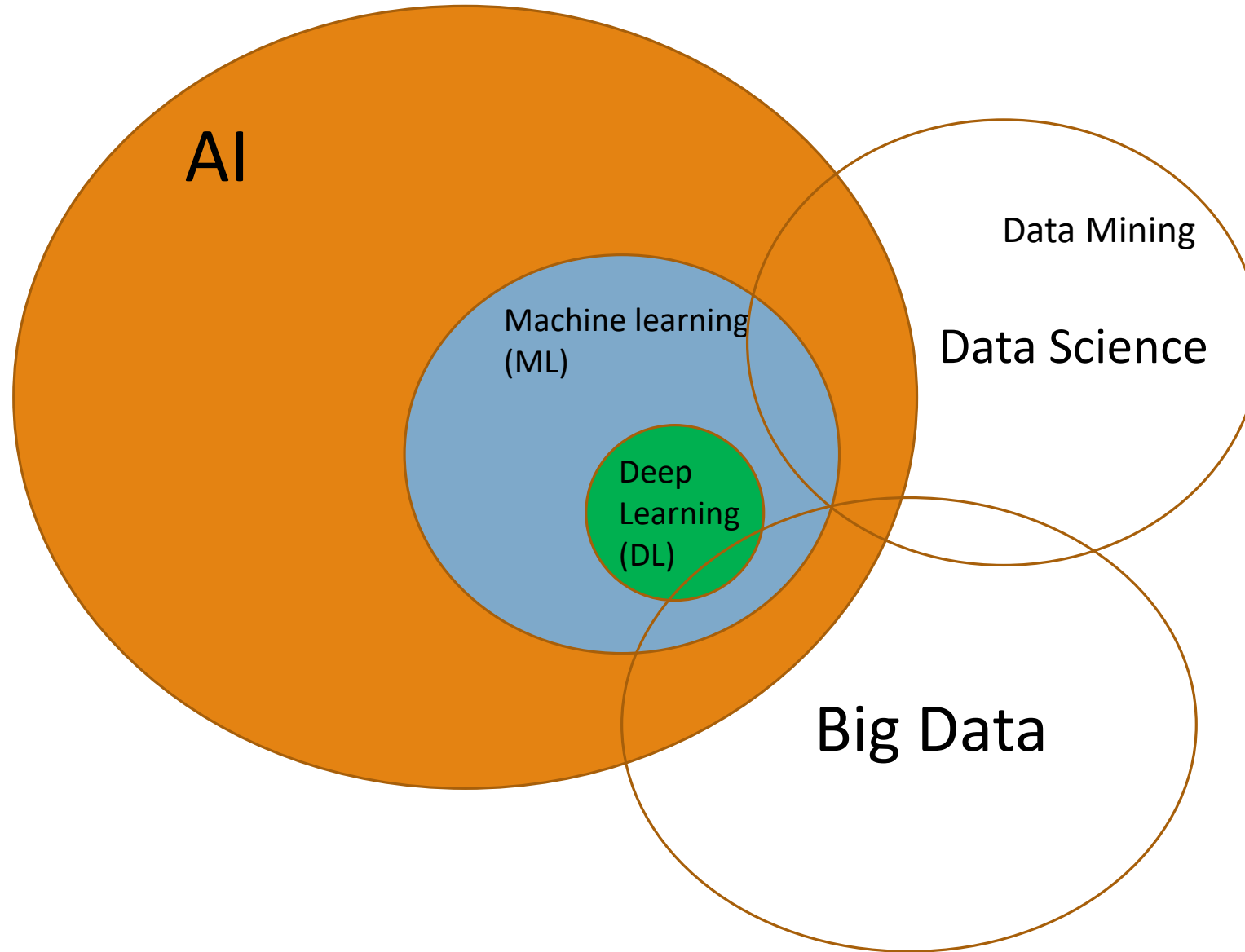# Today's Session Andrew Ng's words

- ~~Most of today's material is not very mathematical.~~

- Key ideas:
  - 1. What do we mean by Interpretability?
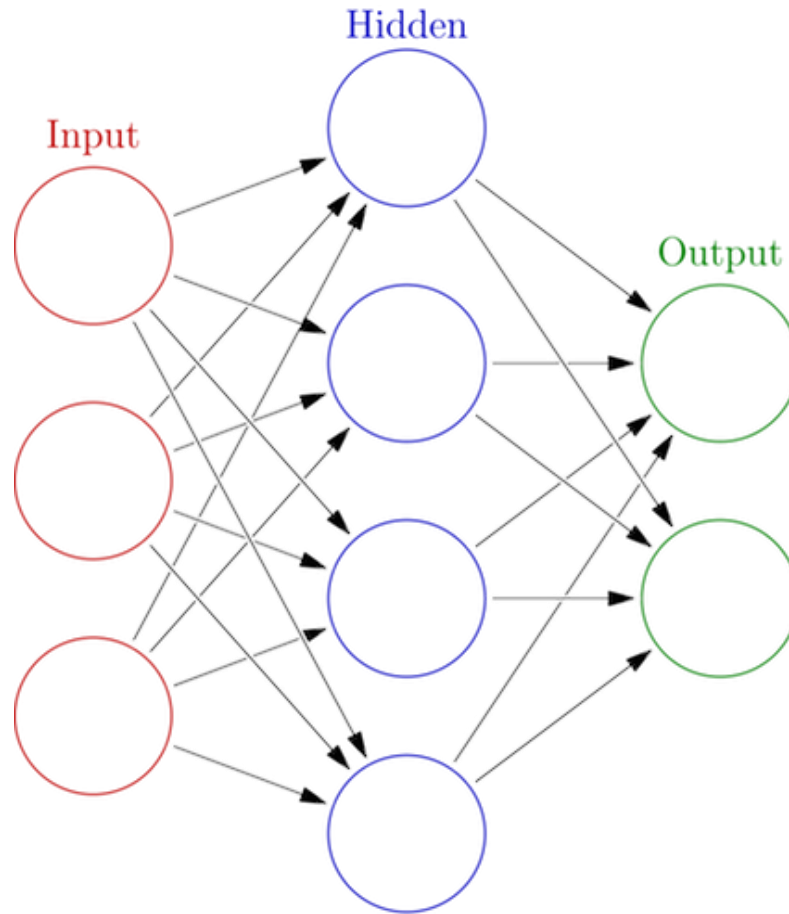  - 2. Why do we need Interpretability?
  - 3. How it is useful?

# What is AI?
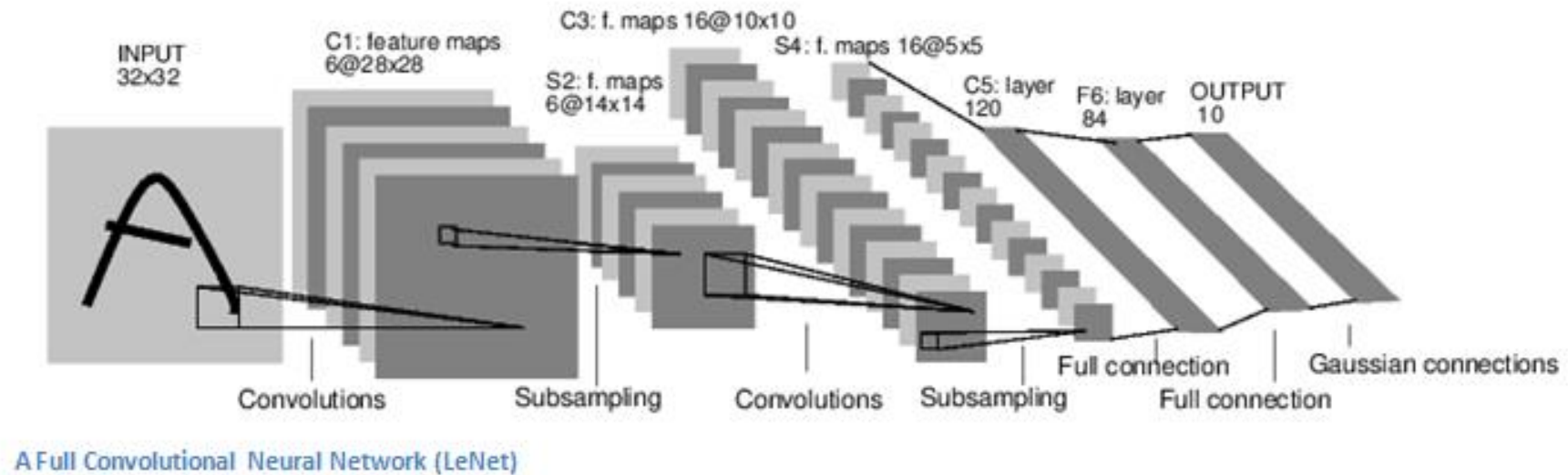
# Evolution of Neural Networks

Simple Neural Net
1980s

# Evolution of Convolutional Neural Networks!

LeNet 1989 (LIP-6, Paris)
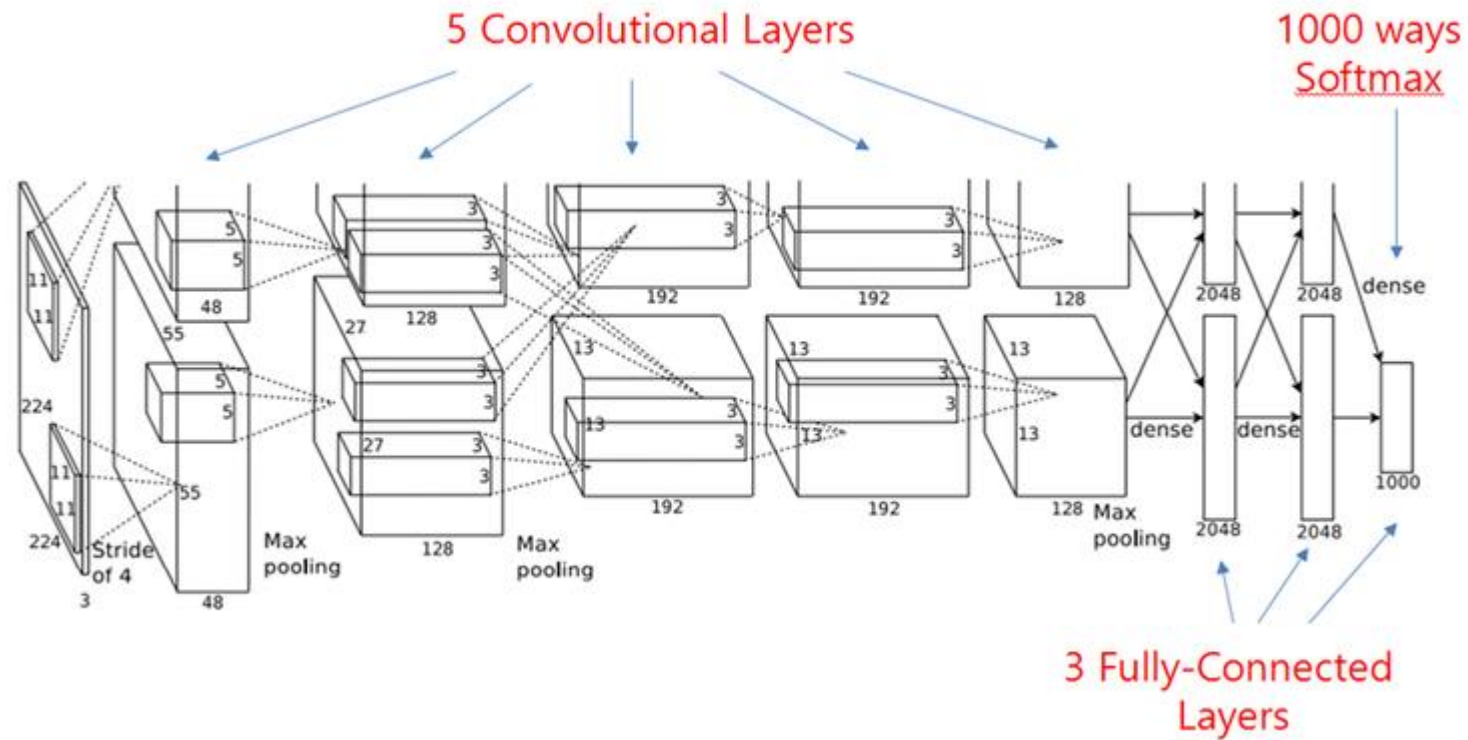


A Full Convolutional Neural Network (LeNet)

# Evolution of Convolutional Neural Networks!

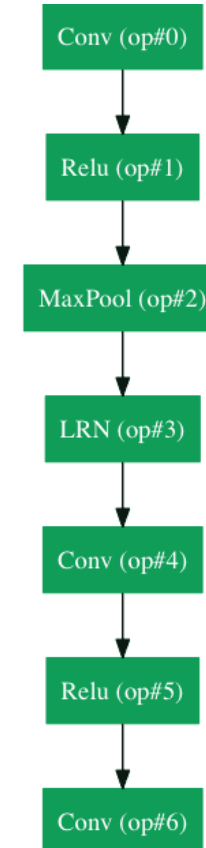## AlexNet 2012

60M parameters

WE NEED TO GO DEEPER

CHALLENGE ACCEPTED

quickmeme.com

# Current Models

**GoogleNet / Inception, 2015/16**

**6.9 M** parameters in the model

# Future AI : Promises

# Future AI : Reality

**Scalable Oversight: How can we efficiently ensure that a given AI system respects aspects of the objective that are too <span style="color:red">expensive to be frequently evaluated</span> during training?**

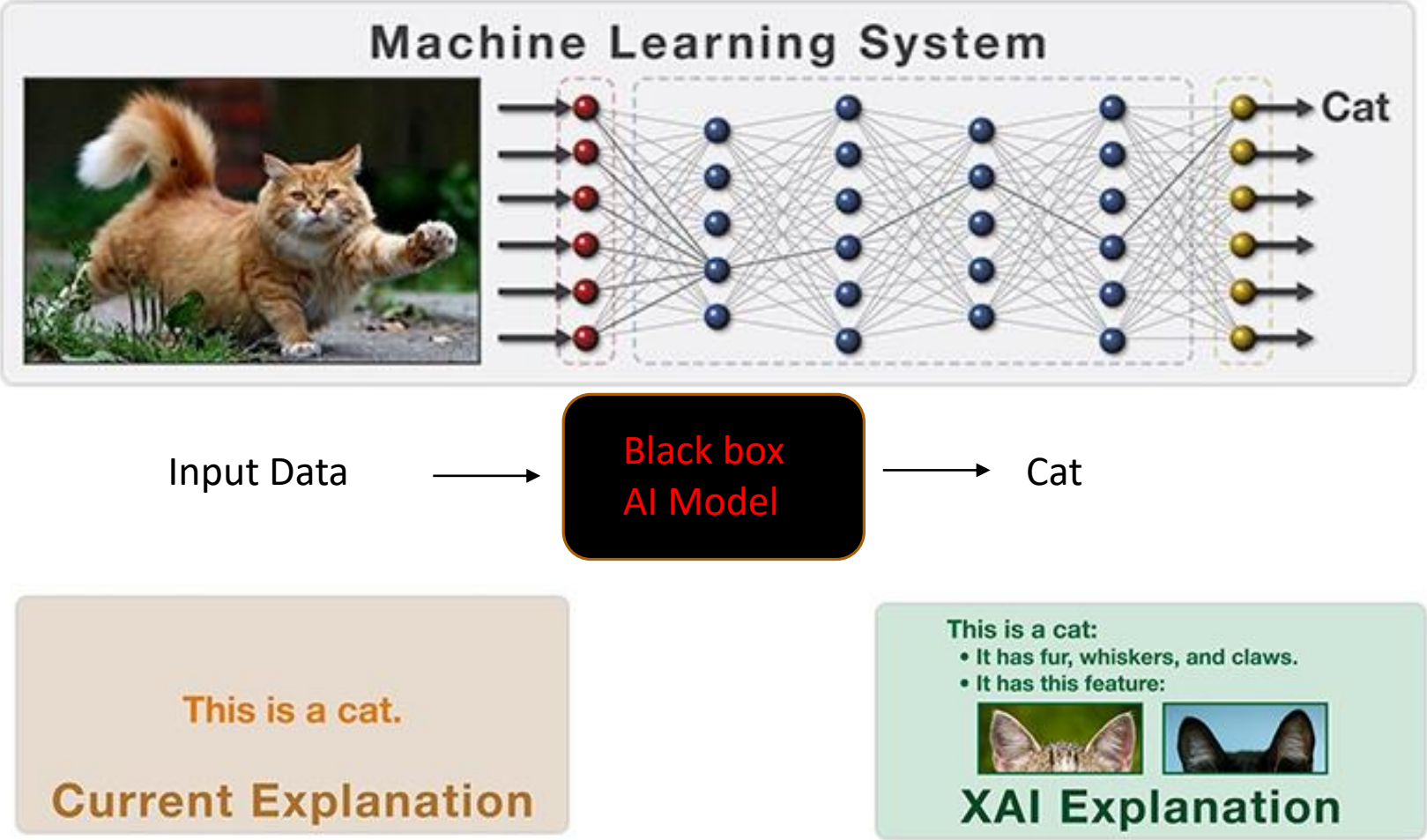**- Google on its challenges for AI, Dailymail UK, June 2016**

*"There is no neural network in the world, and <span style="color:red">no method</span> right now that can be trained to identify objects and images, play Space Invaders, and listen to music."*

*- R.Hadsel (Google DeepMind), The Verge, Oct 2016*

*"We can build these models, but we <span style="color:red">don't</span> know how they work."*

*- Deep Patient , MIT Review (April,2017)*

# Future AI: Explainable



Via XAI: DARPA

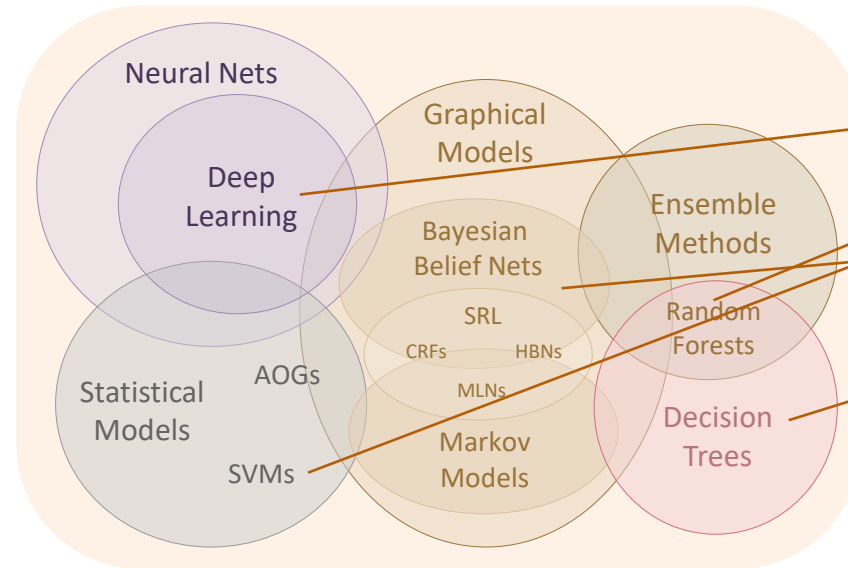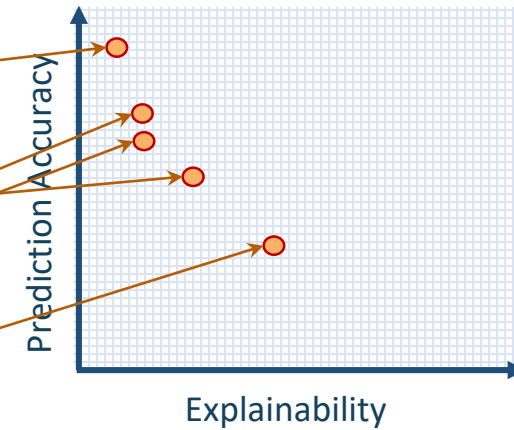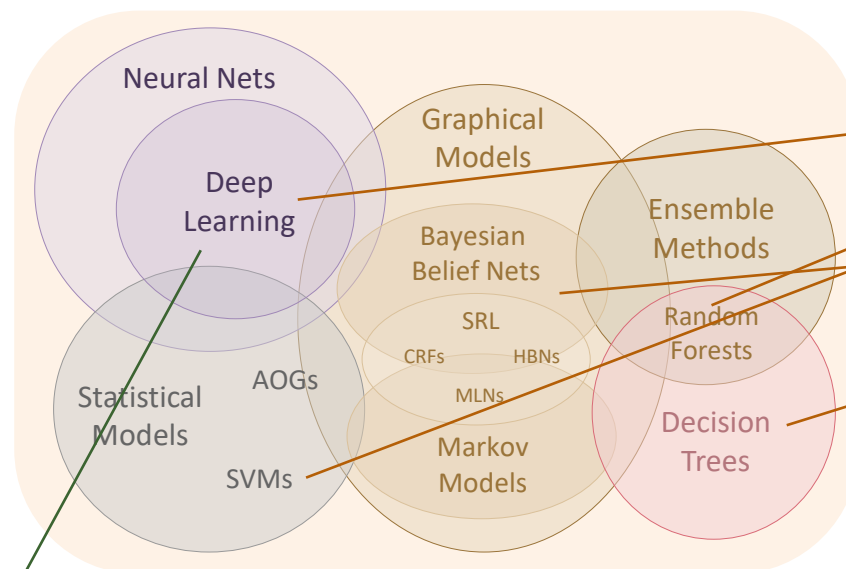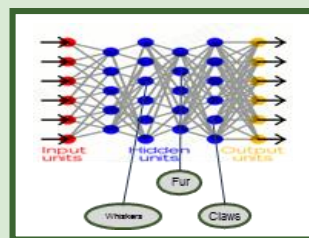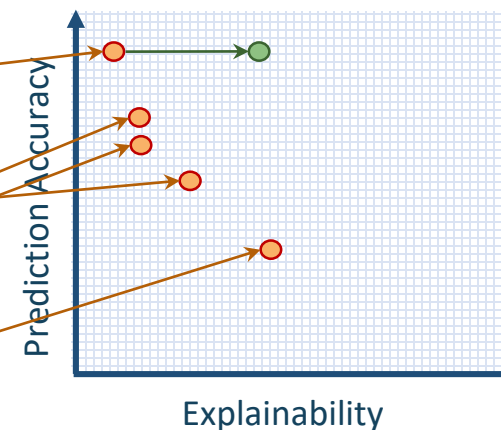# Why do we care about X-AI?

# Explainability

# Explainable Models

New
Approach

Create a suite of
machine learning
techniques that
produce more
explainable models,
while maintaining a
high level of learning
performance



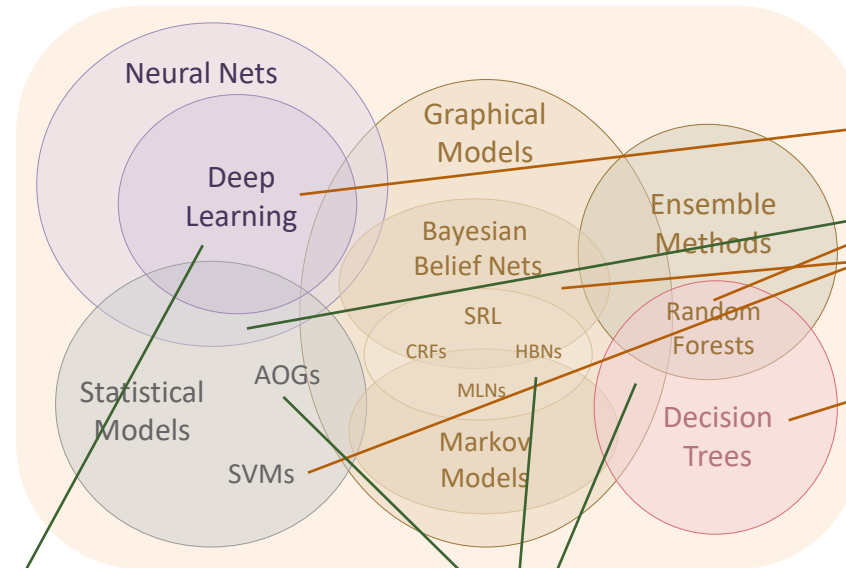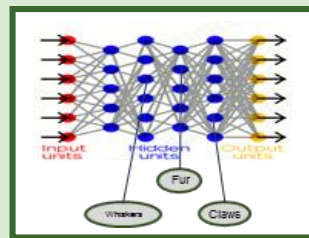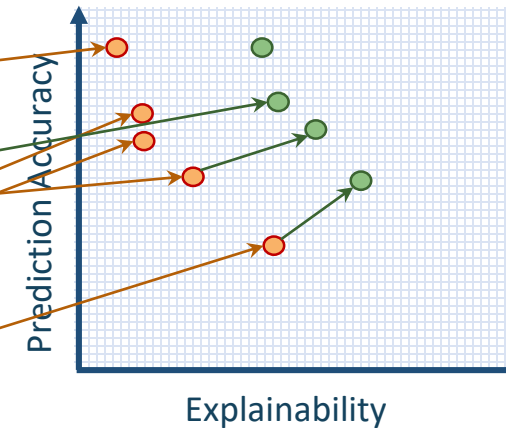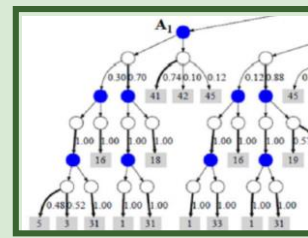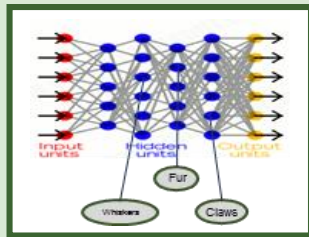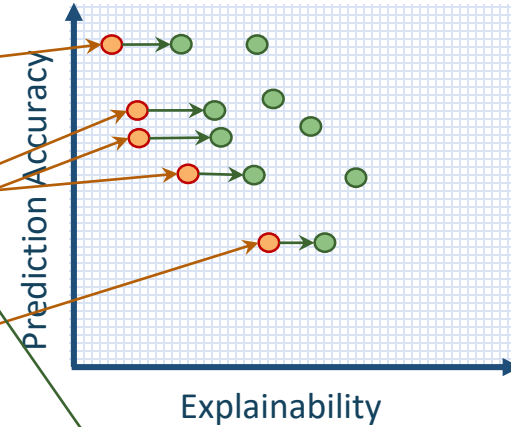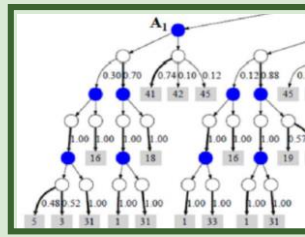Learning Techniques (today)

Neural Nets

Graphical Models
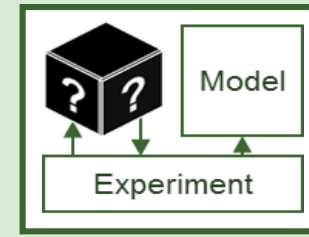
Deep Learning

Ensemble Methods

Bayesian Belief Nets

SRL

CRFs    HBNs

Random Forests

MLNs

Statistical Models

AOGs

Decision Trees

SVMs

Markov Models

Explainability
(notional)

Prediction Accuracy

Explainability

# Explainable Models



**New Approach**

Create a suite of machine learning techniques that produce more explainable models, while maintaining a high level of learning performance

**Learning Techniques (today)**

- Neural Nets
- Deep Learning
- Graphical Models
- Ensemble Methods
- Bayesian Belief Nets
- SRL
- CRFs
- HBNs
- MLNs
- Random Forests
- Statistical Models
- AOGs
- SVMs
- Markov Models
- Decision Trees

**Explainability (notional)**

Prediction Accuracy

Explainability

**Deep Explanation**
Modified deep learning techniques to learn explainable features
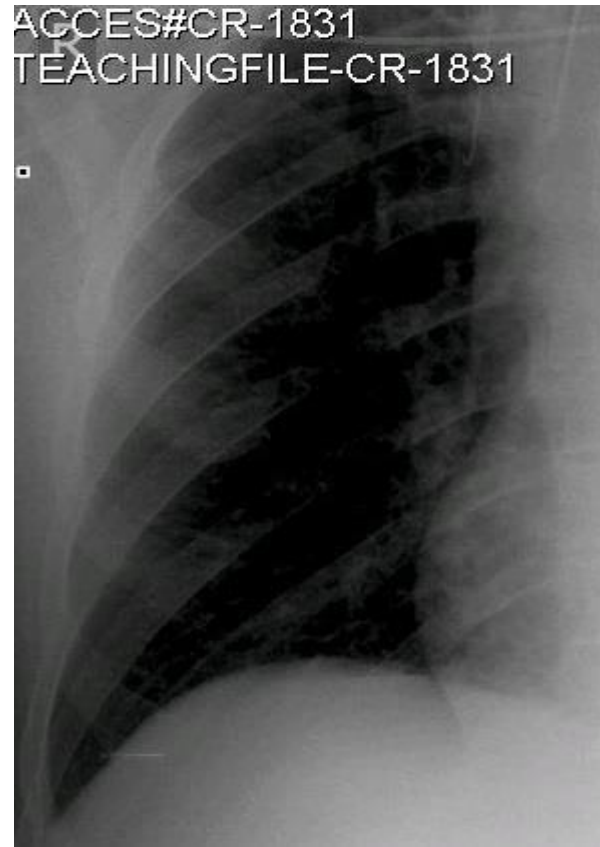
# Explainable Models



New Approach

Create a suite of machine learning techniques that produce more explainable models, while maintaining a high level of learning performance

Learning Techniques (today)

Neural Nets
Deep Learning
Graphical Models
Ensemble Methods
Bayesian Belief Nets
SRL
CRFs    HBNs
MLNs
Random Forests
Statistical Models
AOGs
SVMs
Markov Models
Decision Trees

Explainability (notional)

Prediction Accuracy

Explainability

**Deep Explanation**
Modified deep learning techniques to learn explainable features

**Interpretable Models**
Techniques to learn more structured, interpretable, causal models

# Explainable Models



**New Approach**

Create a suite of machine learning techniques that produce more explainable models, while maintaining a high level of learning performance

**Learning Techniques (today)**

- Neural Nets
- Deep Learning
- Graphical Models
- Ensemble Methods
- Bayesian Belief Nets
- SRL
  - CRFs
  - HBNs
  - MLNs
- Random Forests
- Statistical Models
- AOGs
- SVMs
- Markov Models
- Decision Trees

**Explainability (notional)**

Prediction Accuracy

Explainability

**Deep Explanation**
Modified deep learning techniques to learn explainable features

**Interpretable Models**
Techniques to learn more structured, interpretable, causal models

**Model Induction**
Techniques to infer an explainable model from any model as a black box

# What is Happening?



Breast X-ray



ACCES#CR-1831
TEACHINGFILE-CR-1831

# General Data Protection Regulation (GDPR)

**25th of May**

**2018**

*"more and clearer information about processing"*

# T-SNE: Visualization!

Input Data → Black box AI Model → Cat



Deconvolution: Zeiler et.al. ECCV 14

Guided backpropagation: ICLR 2015

Saliency: Simonyan et.al. CVPR 2013

Deep Taylor Decomp. Montavon et. al.

PR journal 2017

Class Activation Maps (CAM) Bolei Zhou MIT, 2016

Input Data → Black box AI Model → Cat

Prediction Difference:
Zintgraf et. al. ICLR 2017

# Stanford Dog Dataset Results

# Application in Healthcare

# Diabetic Retinopathy:
# Leading Cause of Blindness in the world



**Normal** Retina

**Diabetic** Retina

(a)

| Negative | Mild | Moderate | Severe | Proliferative |

Color Map

| No DR | Mild DR | Mod. DR | Severe DR | Proliferative DR |

# Application in Finance

# Correct Binary Stock Prediction

# Wrong Binary Stock Prediction

# Application in Science (the real one!)

**a** Incident radiation — Structure to classify — (Simulated) Diffraction pattern

**b**
Body-centered-tetragonal (bct$_{139}$) structure spgroup=139

Body-centered-tetragonal (bct$_{141}$) structure spgroup=141

Rhombohedral (rh) structure spgroup=166

Hexagonal (hex) structure spgroup=194

Simple cubic (sc) structure spgroup=221

Face-centered-cubic (fcc) structure spgroup=225

Diamond (diam) structure spgroup=227

Body-centered-cubic (bcc) structure spgroup=229

**d** Pristine structure — $D(\mathbf{q})^{\text{pristine}}$

**e** Defected structure: random displacements — $D(\mathbf{q})^{\text{disp}}$

**f** Defected structure: 25% vacancies — $D(\mathbf{q})^{\text{vac}}$

**g** $D(\mathbf{q})^{\text{disp}} - D(\mathbf{q})^{\text{pristine}}$ [%]

**h** $D(\mathbf{q})^{\text{vac}} - D(\mathbf{q})^{\text{pristine}}$ [%]
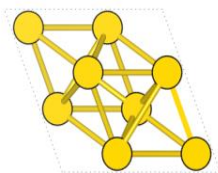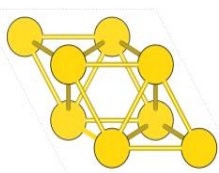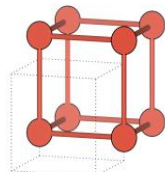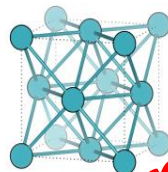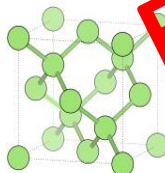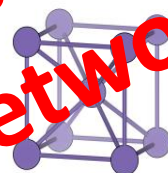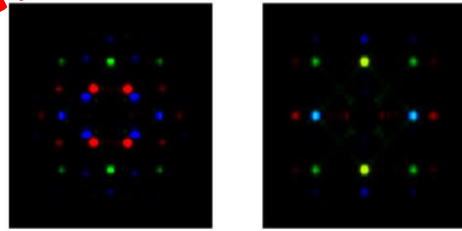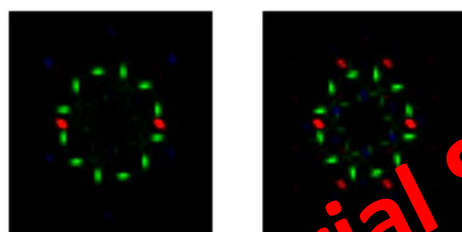
# Ground Truth

# Predicted



**b**

Body-centered-tetragonal
(bct$_{139}$) structure
spgroup=139

Body-centered-tetragonal
(bct$_{141}$) structure
spgroup=141

Rhombohedral
(rh) structure
spgroup=166

Hexagonal
(hex) structure
spgroup=194

Simple cubic
(sc) structure
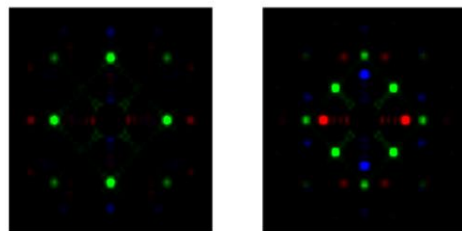spgroup=221

Face-centered-cubic
(fcc) structure
spgroup=225

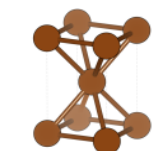Diamond
(diam) structure
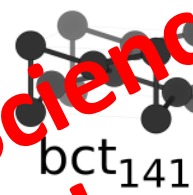spgroup=227
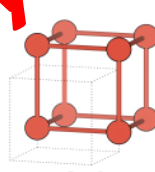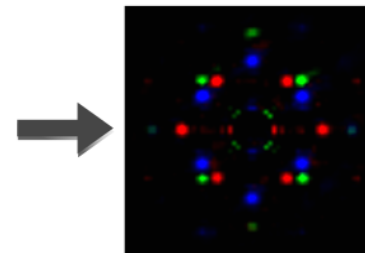
Body-centered-cubic
(bcc) structure
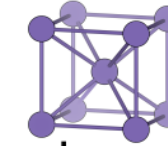spgroup=229

**c**

**b**

bct$_{139}$
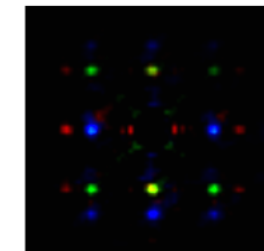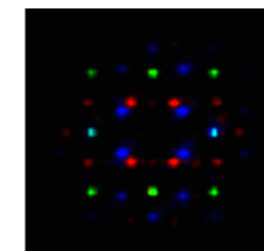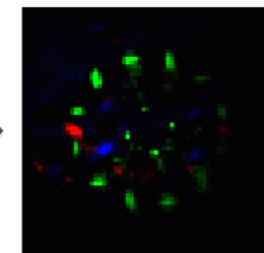
bct$_{141}$

sc

fcc

hex/rh

diam

bcc

First paper in Material Science that shows Neural Network Interpretability!

# Thank You!

devinder.kumar@uwaterloo.ca

http://devinderkumar.com