



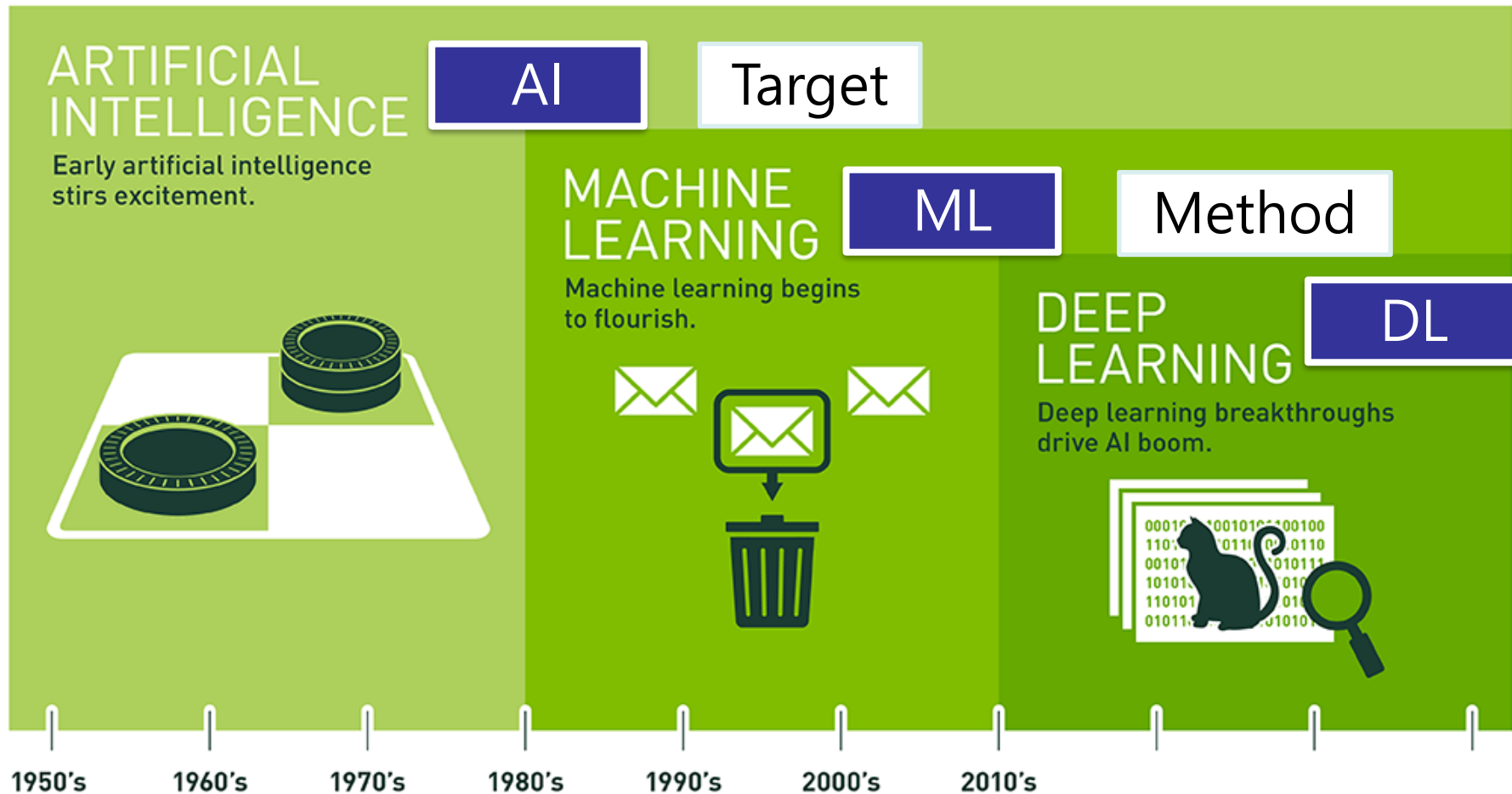
ECE 661

COMP ENG ML & DEEP NEURAL NETS

1. INTRODUCTION

HAI “HELEN” LI, SPING 2024

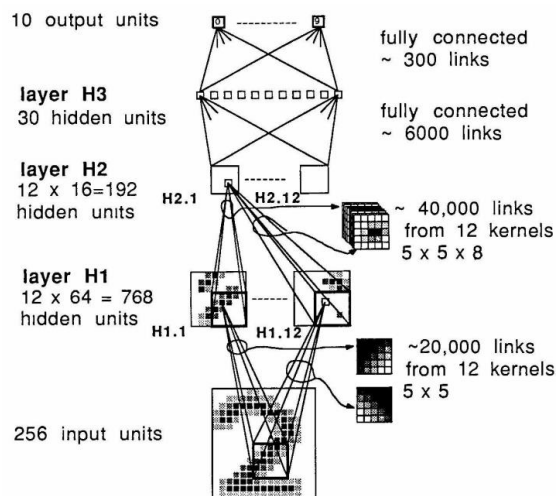
AI ↔ ML ↔ DL



Since an early flush of optimism in the 1950s, smaller subsets of artificial intelligence – first machine learning, then deep learning, a subset of machine learning – have created ever larger disruptions.

Historical Overview

Convolutional Network (1980s)

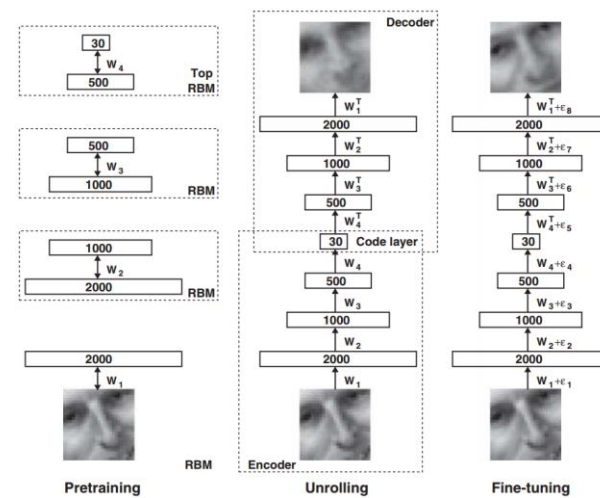


Dark period (1990s)

- Serious problem: Vanishing gradient
- No benefits observed by adding more layers
- No high-performance computing devices



Renaissance (2006 ~ Present)



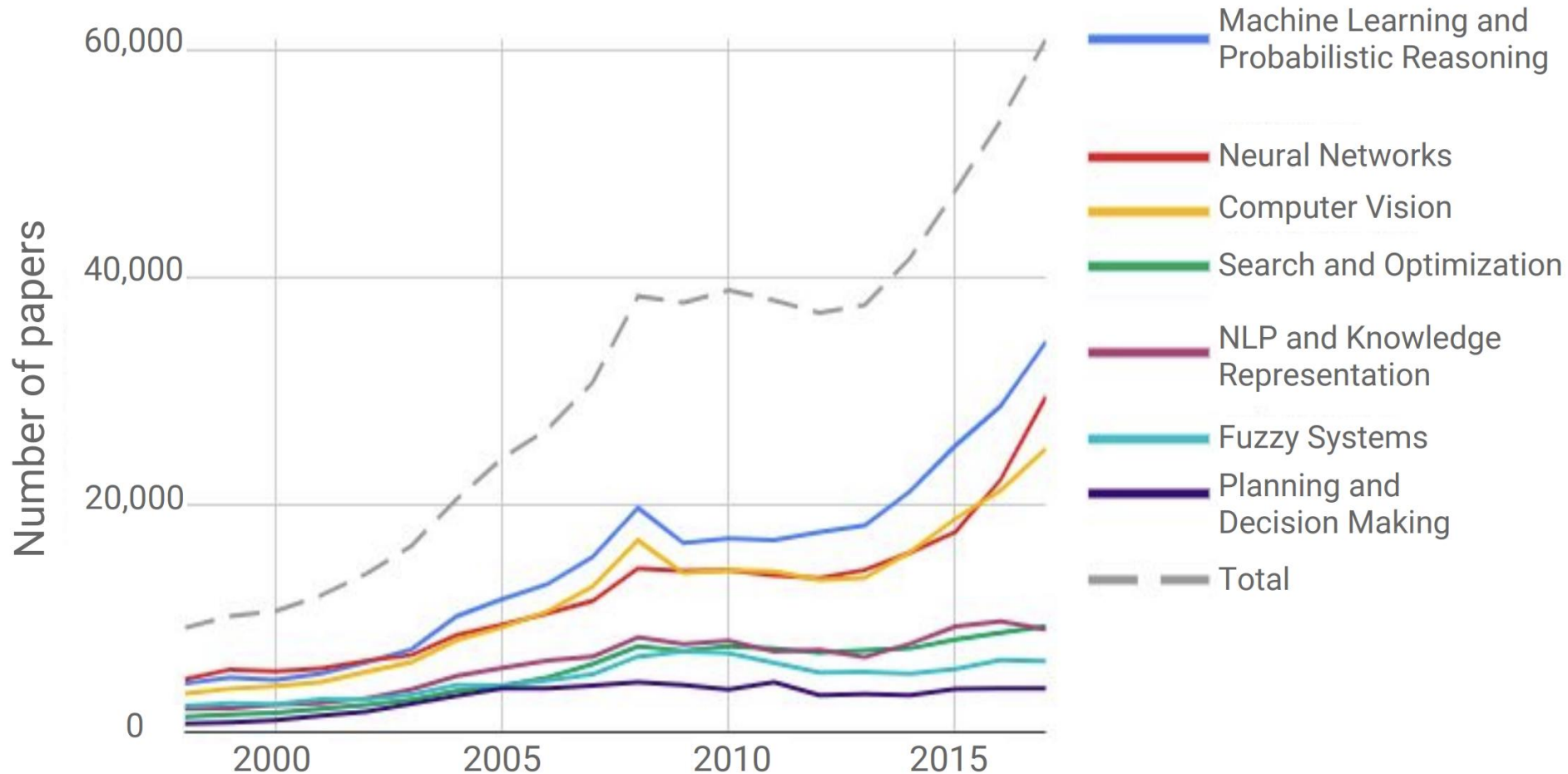
Y. Lecun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard and L. D. Jackel. Backpropagation Applied to Handwritten Zip Code Recognition. 1989.

J. Schmidhuber. Deep Learning in Neural Networks: An Overview. arxiv, 2014.

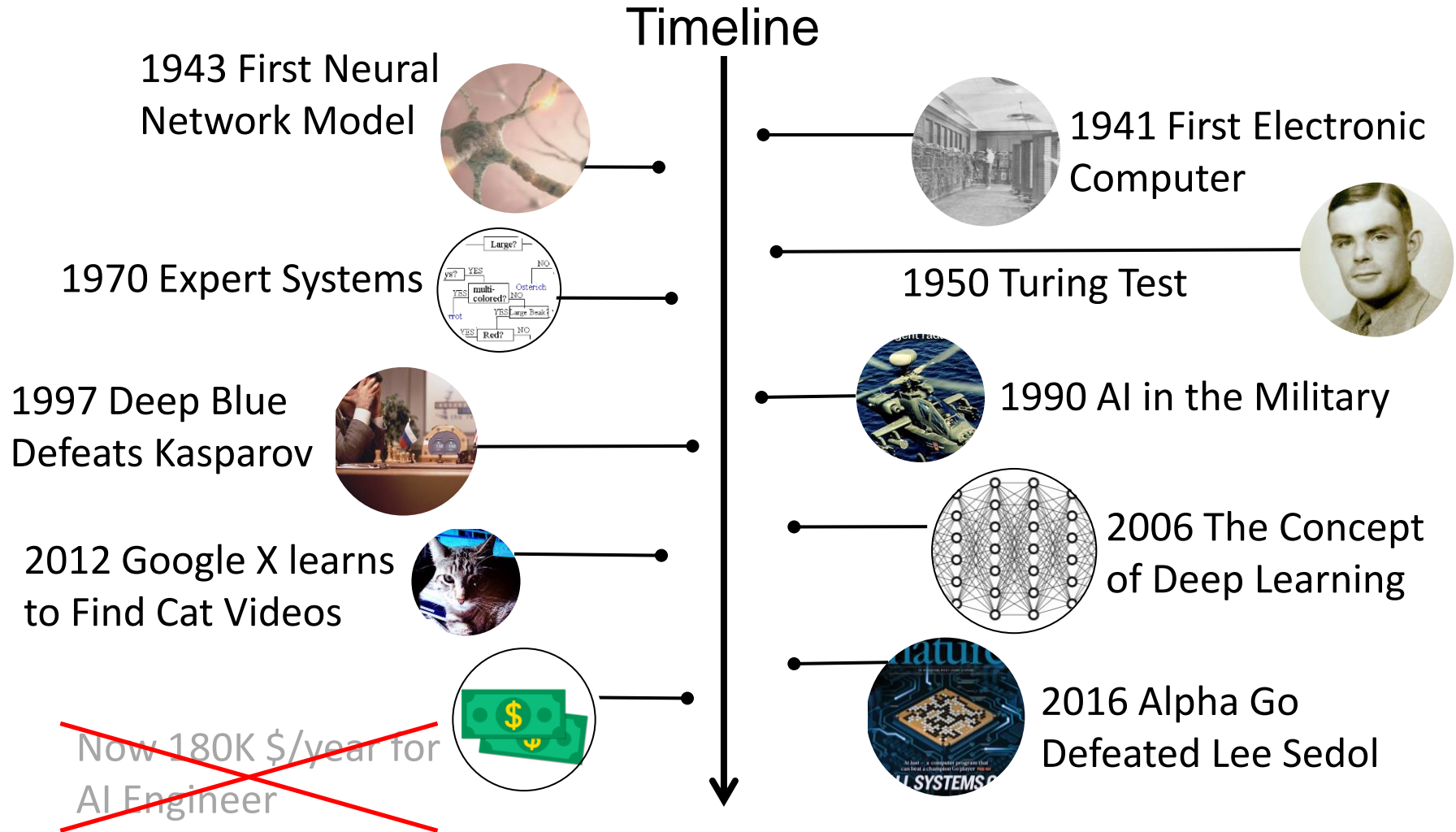
G. E. Hinton and R. R. Salakhutdinov. Reducing the Dimensionality of Data with Neural Networks. Science, 2006.

AI Papers

The number of AI papers on Scopus

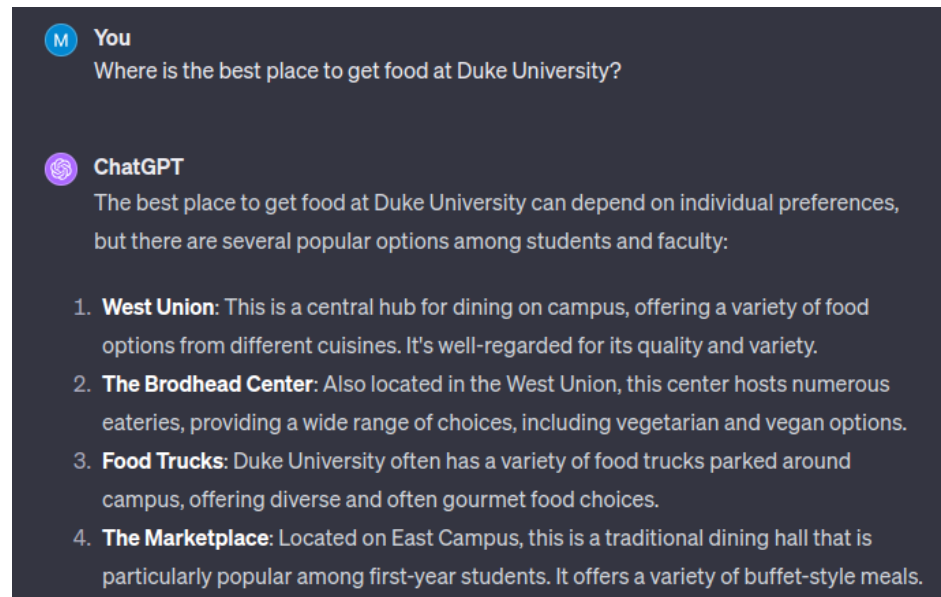


Historical Milestones



Recent Development: ChatGPT and LLMs

- The recent creation of ChatGPT has drawn massive attention to AI.
- Large Language Models (LLMs) have been a developing research area for years, but ChatGPT was the first highly competent chatbot.
- Useful for simple automation tasks, coding, teaching, etc leading to a renaissance in LLM-based products and research.
- Prone to hallucination and not human-level at reasoning/planning, but very knowledgeable across a wide range of topics.



Outline

- Course Introduction
- Machine Learning & Deep Neural Networks
 - Applications
 - Categories
 - Important Metrics
 - Platforms & Frameworks

Course Objectives

- For MS/MEng students and undergraduate students who want to learn computer engineering methods commonly performed in developing and using machine learning and deep neural network models.
- For PhD students who want to learn and practice a wide variety of ML topics that are beyond any single focus area. The breadth of knowledge covered may spur new ideas to be used in your own research.
- **Practice** will be the focus of this course, while **theoretical understanding** is essentially important.

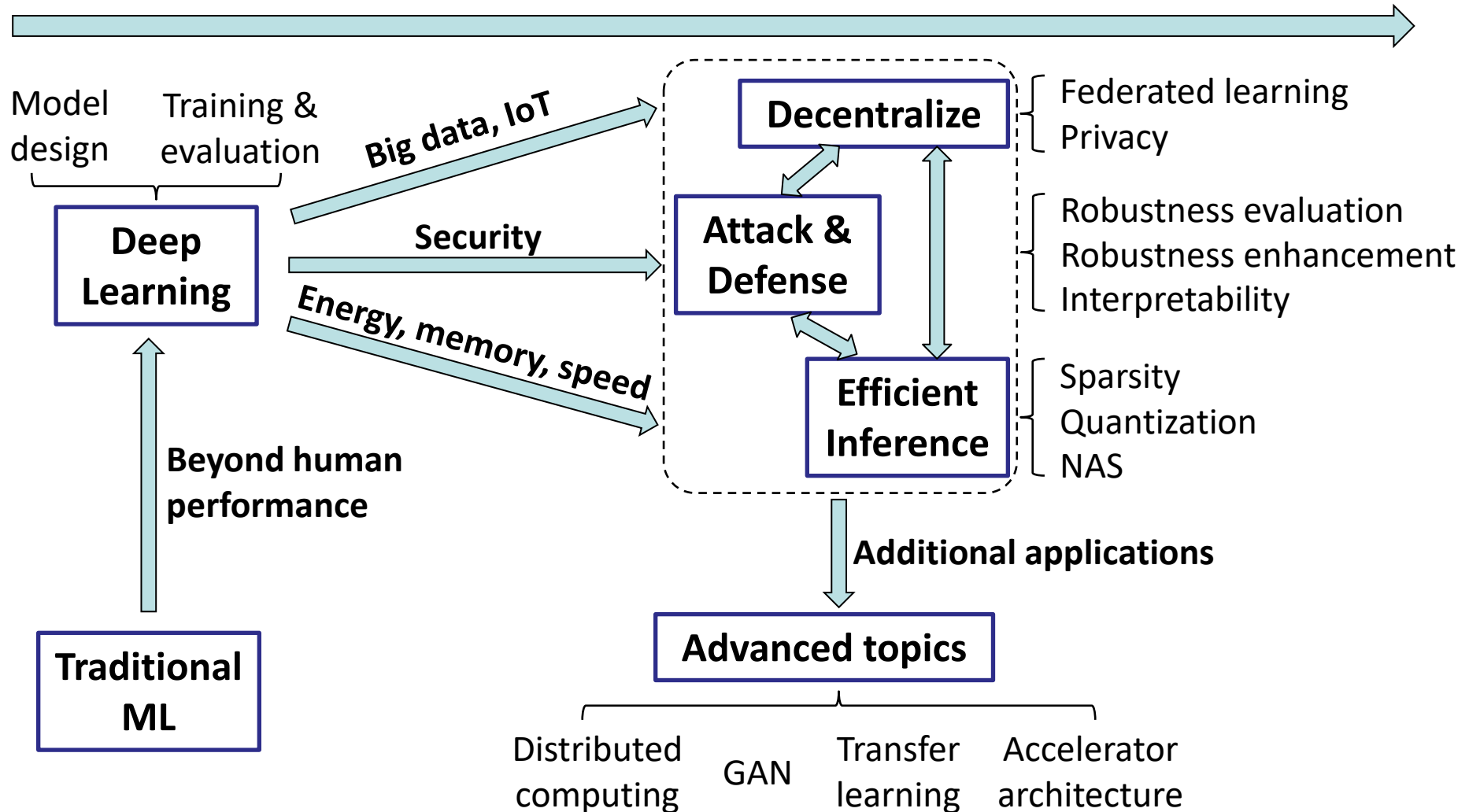
Course Objectives

This course is designed to improve your ability to:

1. **Comprehend** the mechanisms, applications, and limitations of techniques commonly used in training and inference of machine learning and deep neural networks algorithms;
2. **Formulate** hypotheses and conduct experiments employing these techniques;
3. **Analyze** experimental results obtained by these techniques and your own practices and **derive** the conclusions that are supported or not supported by your data;
4. **Synthesize** and **communicate** the experimental results and data through oral narrative, graphs, figure legends, and result narratives;
5. **Utilize** proper engineering techniques for novel machine learning algorithms and deep neural network models;
6. **Propose** new engineering approaches and techniques to further enhance machine learning and deep neural network training and inference execution.

Course Roadmap

Applying machine learning into the real world



Course Overview

- DNN fundamentals
 - Forward/backward propagation, training, convolutional neural network (CNN), network architecture, recurrent neural network (RNN), language models, ...
- DNN acceleration
 - Compact neural architecture, model compression, pruning, quantization, sparsification, ...
- Machine learning security
 - Adversarial attack, robust learning method, ...
- Advanced topics
 - Distributed computing, neural architecture search (NAS), transfer and reinforcement learning, generative adversarial network (GAN), decentralization and privacy, DNN accelerator, ...

Spring 2024 Tentative Schedule

(Subjective to change)

Week	Date	Lecture	Content	Assign	Due	Assignment
1	1/10	Course Introduction	Lec01			
2	1/15	Martin Luther King Jr. Day holiday; No Class				
	1/17	Perceptron and back propagation	Lec02			
3	1/22	Image feature and 2D convolution	Lec03	HW 1		Gradient computation (logistic regression), CNN weight observation
	1/24	Convolutional Neural Network (CNN)	Lec04			
4	1/29	CNN Training - Basic	Lec05			
	1/31	CNN Training - Basic & Advanced	Lec05/06			
5	2/05	CNN Training - Advanced	Lec06	HW 2	HW 1	CIFAR-10 training: exploring DNN architectures and parameter tuning
	2/07	CNN Architectures	Lec07			
6	2/12	Compact Neural Architecture Design	Lec08			
	2/14	RNN and Language Models	Lec09			
7	2/19	Attention Model & LLM	Lec10	HW 3	HW 2	RNN
	2/21	Deep Learning Hardware Systems	Lec11			
8	2/26	Project Introduction	Lec12	Project idea		
	2/28	Deep Compression	Lec13			
9	3/04	Sparse Regularization	Lec14			
	3/06	Sparse Optimization	Lec15	HW 4	HW 3	Sparsification and quantization
10	3/11	Spring recess; No Class				
	3/13	Spring recess; No Class				
11	3/18	Fixed-point Quantization	Lec16		Project Proposal	
	3/20	Machine Learning Security	Lec17			
12	3/25	Adversarial Attack - Attacks	Lec18	HW 5		
	3/27	Adversarial Attack - Defenses	Lec19		HW 4	Adversarial attack and adversarial training
13	4/01	Transfer learning / Generative Adversarial Network	Lec20			
	4/03	AutoML	Lec21			
14	4/08	Neural Architecture Search	Lec22		HW 5	
	4/10	Distributed/Federated Learning	Lec23		Project Mid-point Check-in	
15	4/15	Lab Q&A				
	4/17	Lab Q&A				
16	TBD	Project Poster Session			Poster & Report	

Related Topics and Courses

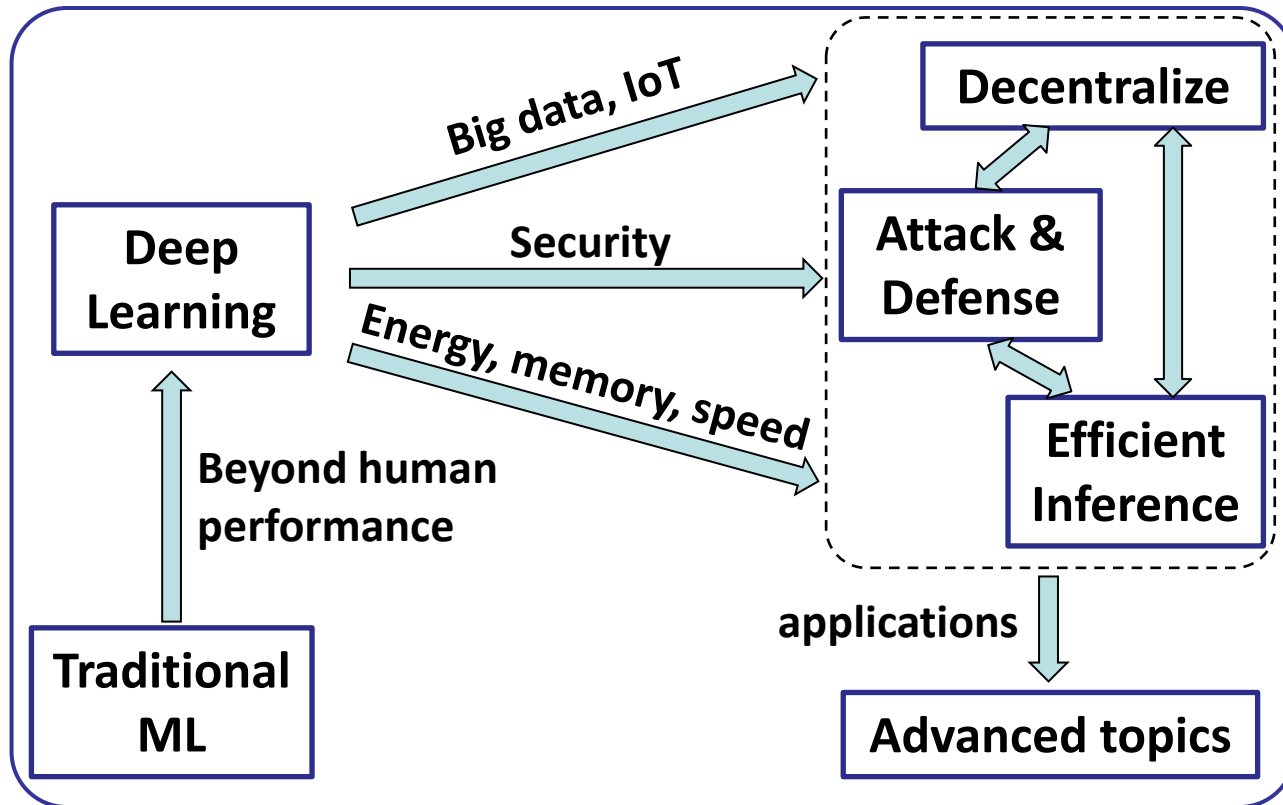
NLP: **ECE 684**

Deep learning: **ECE 685**

Cloud computing: **ECE 563**; Smart sensor: **ECE 590-04**

Security: **ECE 590-03**; Image processing & denoise: **ECE 588**

Information theory: **ECE 587**; Compressed sensing: **ECE 741**



Math basics: **ECE 581, 586**

Machine learning: **ECE 681, 687**

Implementation: **ECE 550, 551, 650**

Architecture design: **ECE 552, 590-24**

System optimization: **ECE 558, 563, 565**

Hardware: **ECE 538, 539, 559, 590 (ML accel.)** ¹³

Prerequisites

- We expect that students to have basic object-oriented programming experience (e.g., C++, Python) and be familiar with linear algebra and computer hardware fundamentals prior to taking this course, such as
 - For graduate students: ECE 551D
 - For undergraduate students: CS201
- If you are familiar with a topic that we are covering ...
 - You may learn something new
 - I may present it slightly differently than you are used to
 - You may be able to help other students learn it
- If you do not have these pre-requisites and are unfamiliar with these topics
 - We will **NOT** be slowing down to cover them
 - Please come talk to me or a TA sooner rather than later!

“Learn by Doing”

- Five (5) assignments (conceptual questions + labs in PyTorch)
 - 1: Building and understanding CNN modules
 - 2: CIFAR-10 training
 - 3: RNN model
 - 4: Advanced sparse optimization and quantization techniques
 - 5: Adversarial attack and adversarial training
- One (1) project

Logistics

ECE 661 COMP ENG ML & DEEP NEURAL NETS		
Faculty:	Dr. Yiran Chen	yiran.chen@duke.edu
Lectures:	Mondays/Wednesdays 10:05 AM – 11:20 AM Wilkinson 021 In person only. No recording	
Office Hours:	By Appointment (please send email to Dr. Yiran)	
Teaching Assistants:	Haoxuan Shan (Lead TA)	haoxuan.shan@duke.edu
	Junyao Zhang	junyao.zhang2@duke.edu
	Mark Horton	mark.horton@duke.edu
	Martin Kuo	martin.kuo@duke.edu
Office Hours:	Tuesday 5:00 PM – 7:00 PM, Location: TBD Thursday 3:00 PM – 5:00 PM, Location: TBD	

*TAs are **NOT** under obligation to bail you out at 3am or debug your code.
Your best bet is to get help in a timely and reasonable manner!*

Getting Info

- **Sakai:**
 - Syllabus, schedule, slides, assignments, rules/policies, prof/TA info, office hour info
 - Links to useful resources
 - Just assignment submission and gradebook
- **Slack workspace:** questions/answers
 - Use you Duke email to sign up at the following link
https://join.slack.com/t/ece661-2024sp/shared_invite/zt-2a0myy7pi-lbZit8dL6TqKVQHlKvpx5Q
 - Post all your questions here
 - Questions must be “public” unless good reason otherwise
 - No code in public posts!

Getting Answers to Questions

- What do you do if you have a question?
 - Check Sakai
 - Check Slack
 - If you have questions about homework, use Slack – then everyone can see the answer(s) posted there by me, a TA, or your fellow classmate
 - Contact TA directly if need additional background materials for prerequisite knowledge
 - Contact professor directly if issue that is specific to you and that can't be posted publicly (e.g., regrade)

Textbook & Software

- There are no designated textbooks for this course.
- The related reading materials (e.g., papers, webpages, etc.) will be distributed through Sakai before the classes.
- We use PyTorch (<https://pytorch.org/>) in this course



Grading

Assignment	%
Assignments (5)	55%
Project	25%
Quiz	20%

- Completion of all assignments is required in order to earn a passing grade of D- or better in this course.
- Course grades are determined using an absolute, but adjustable scale (i.e., there is no curve). A final course average (rounded to the nearest 0.1 point) of at least 93.3 = A, 90.0 = A-, 86.7 = B+, 83.3 = B, 80.0 = B-, etc.
- Note: the professor reserves the rights to scale the grades.

Homework Submission

- Homework assignments and lab reports will be submitted as **PDF files** through the Assignments tool in Sakai/Gradescope. **The code** of lab assignments will be submitted to our servers. The details will be given in assignments.
- Late policy
 - < 24 hours late: deduct 10% credits
 - < 48 hours late: deduct 25% credits
 - No credit for late work after 48 hours
- Strict cutoff time will be enforced based on submission timestamp on Sakai
- Consider a small margin in case of system/internet issue

Grade Appeals

- All re-grade requests must be in writing
 - your assignment in question, with
 - a brief written description of the error, and
 - your Duke NetID.
- I will respond to your regrade request by email and make arrangement to return your work to you.
- As a matter of policy, when you request a regrade you are agreeing that the grading of the entire assignment may be re-evaluated.
- All regrade requests must be submitted no later than 1 week after the assignment was returned to you.

Academic Misconduct

- Academic Misconduct
 - Refer to **Duke Community Standard**
 - Homework/lab is individual – you do your own work
 - Common examples of cheating:
 - Running out of time and using someone else's output
 - Borrowing code from someone who took course before
 - Using solutions found on the Web
 - Having a friend help you to debug your program
- **We will not tolerate any academic misconduct!**
 - We use software for detecting cheating
- “But I didn’t know that was cheating” is not a valid excuse

Academic Integrity: General

Some general guidelines

- If you don't know if something is OK, please ask me.
- If you think “I don't want to ask, you will probably say no” that is a good sign its NOT acceptable.
- If you do something wrong, and regret it, please come forward—I recognize the value and learning benefit of admitting your mistakes. (Note: this does NOT mean there will be no consequences if you come forward).
- If you are aware of someone else's misconduct, you should report it to me or another appropriate authority.

Course Problems

- Struggling in course
 - Come to see me/TAs: We are here to help
- Other problems:
 - Feel free to talk to the instructor, who generally understands and will try to work with you
 - Some problems may extend well beyond my course
 - Academic Advisor
 - DGS Team

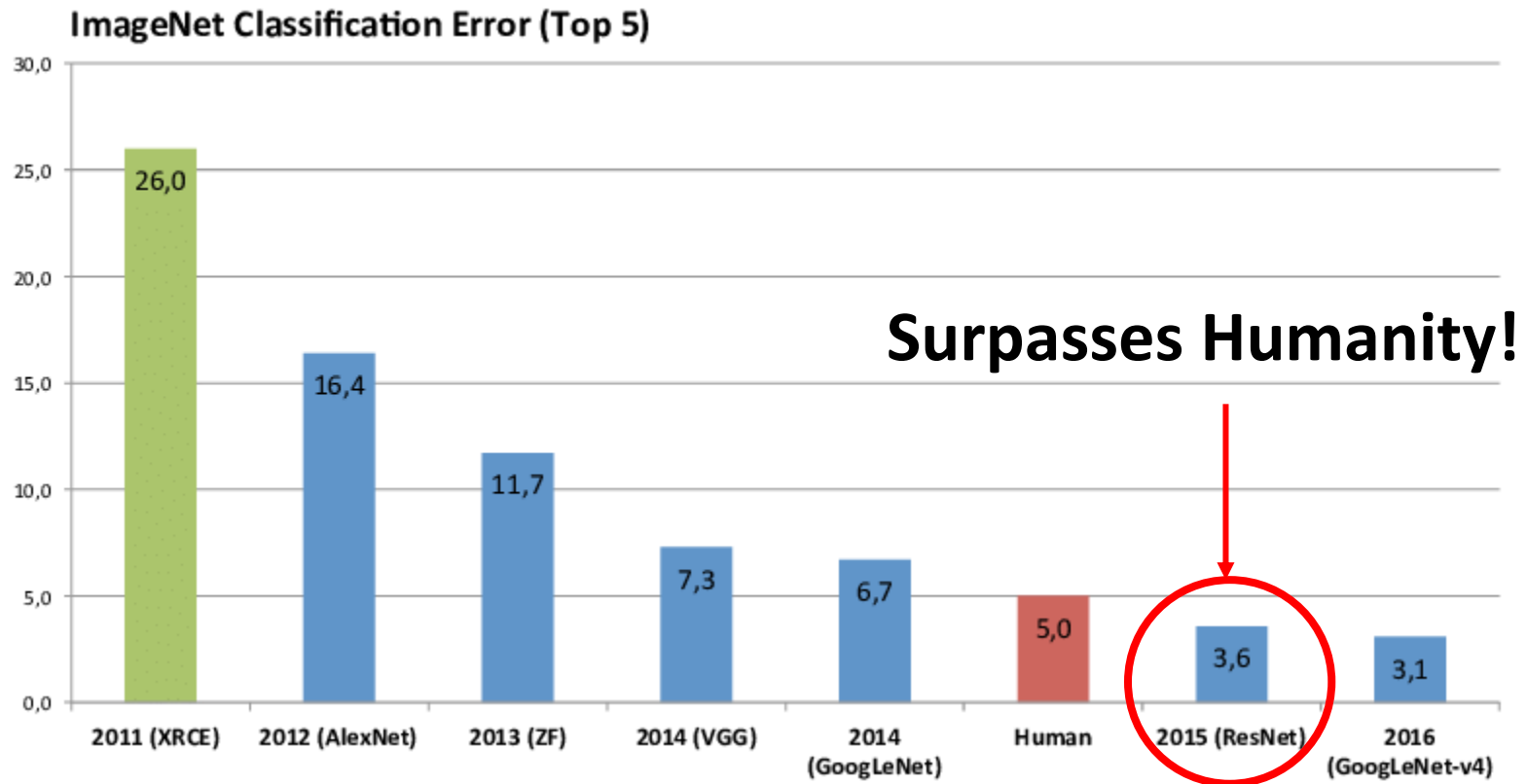
Our Responsibilities

- The instructor and TAs will...
 - Provide lectures/office hours at the stated times
 - Set clear policies on grading
 - Provide timely feedback on assignments
 - Be available out of class to provide reasonable assistance
 - Respond to comments or complaints about the instruction provided
- Students are expected to...
 - Receive lectures/recitations at the stated times
 - Turn in assignments on time
 - Seek out of class assistance in a timely manner if needed
 - Provide frank comments about the instruction or grading as soon as possible if there are issues
 - Assist each other within the bounds of academic integrity

Outline

- Course Introduction
- Machine Learning & Deep Neural Networks
 - Applications
 - Categories
 - Important Metrics
 - Platforms & Frameworks

Applications: Images



Can you tell
what kind of
turtle this is?



- A. *Dermochelys coriacea*
- B. *Caretta caretta*
- C. *Lepidochelys kempii*
- D. *Lepidochelys olivacea*

Applications: Images

- However, surprisingly weak...



Panda

57.7% confidence

+ 0.007 ×



noise

=



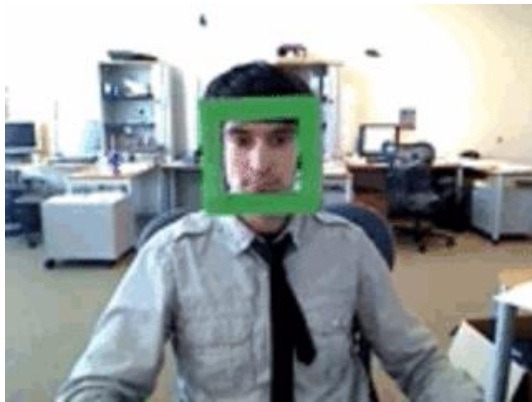
gibbon

99.3% confidence

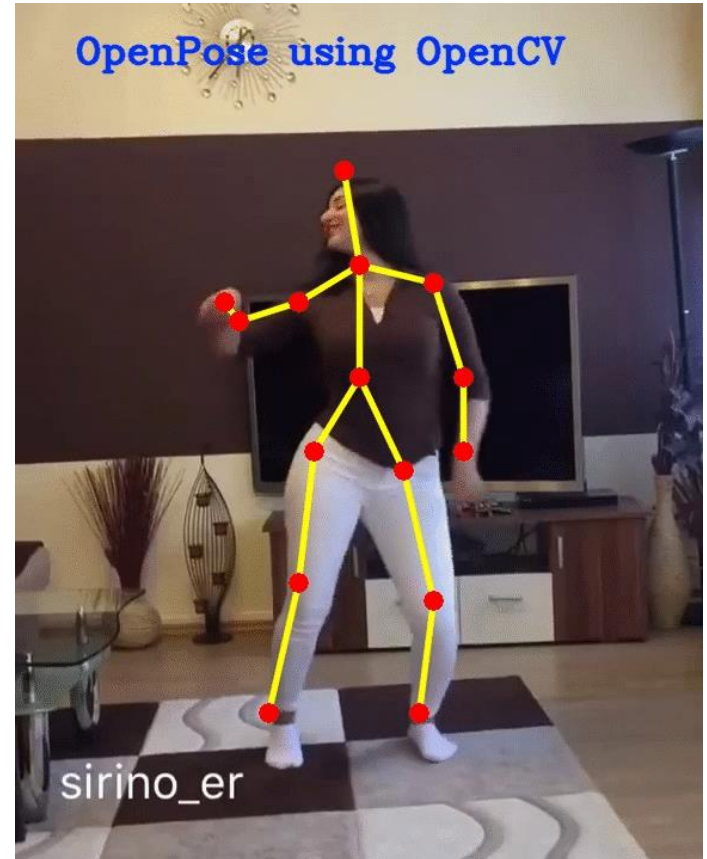
Applications: Videos



Object Detection

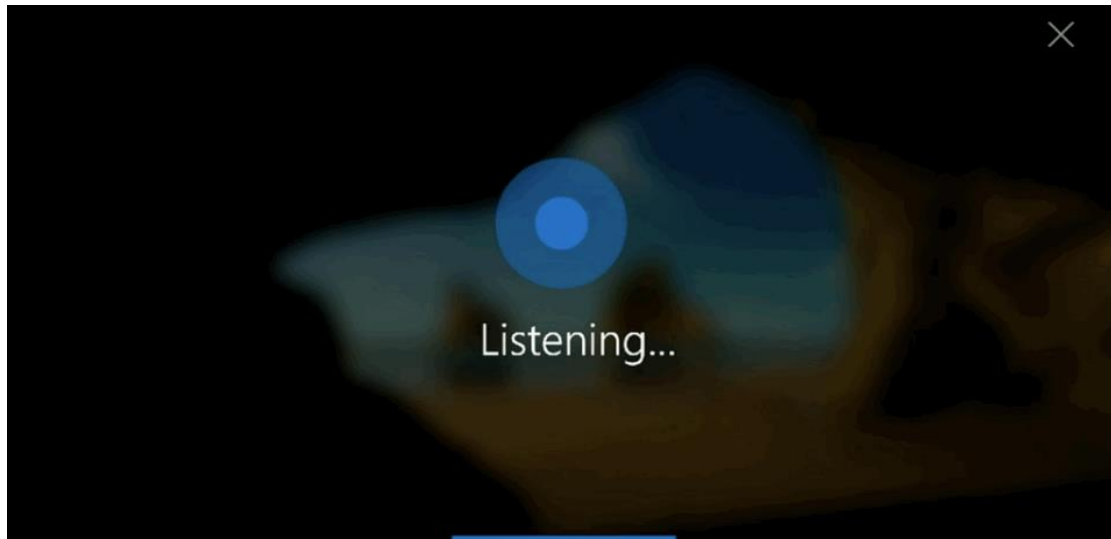


~~The Perfect Real Time
Face Tracking~~

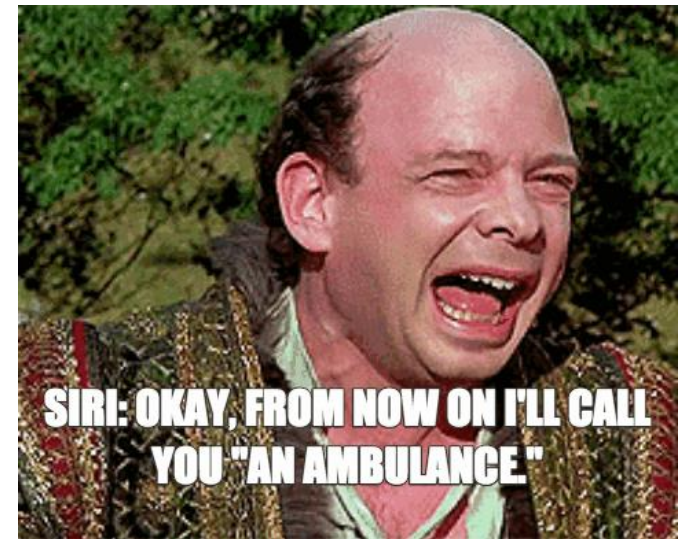


Human Pose Estimation

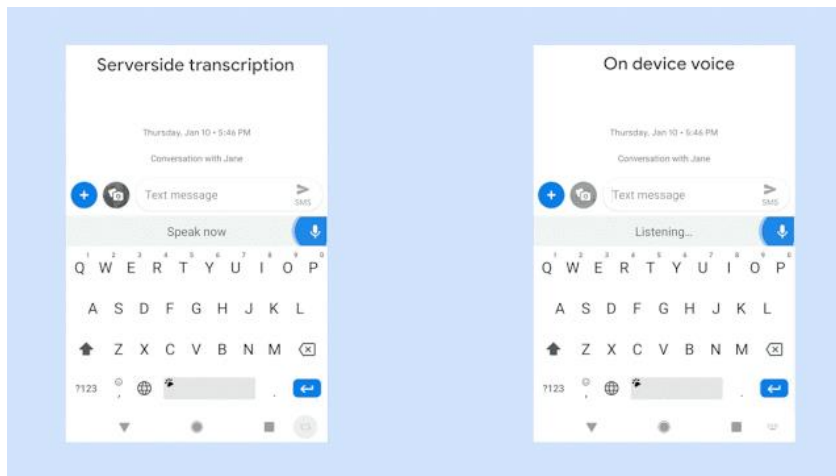
Applications: Speech



Cortana



"Siri, Call me an ambulance"



Speech To Text

"Remember when people typed with two fingers? My voice is faster."

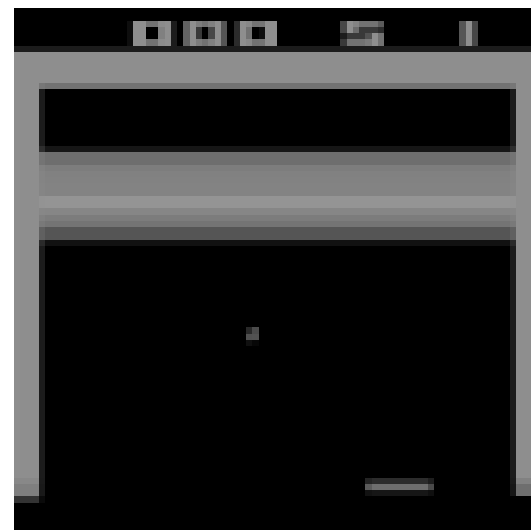
Applications: Game; Strategy



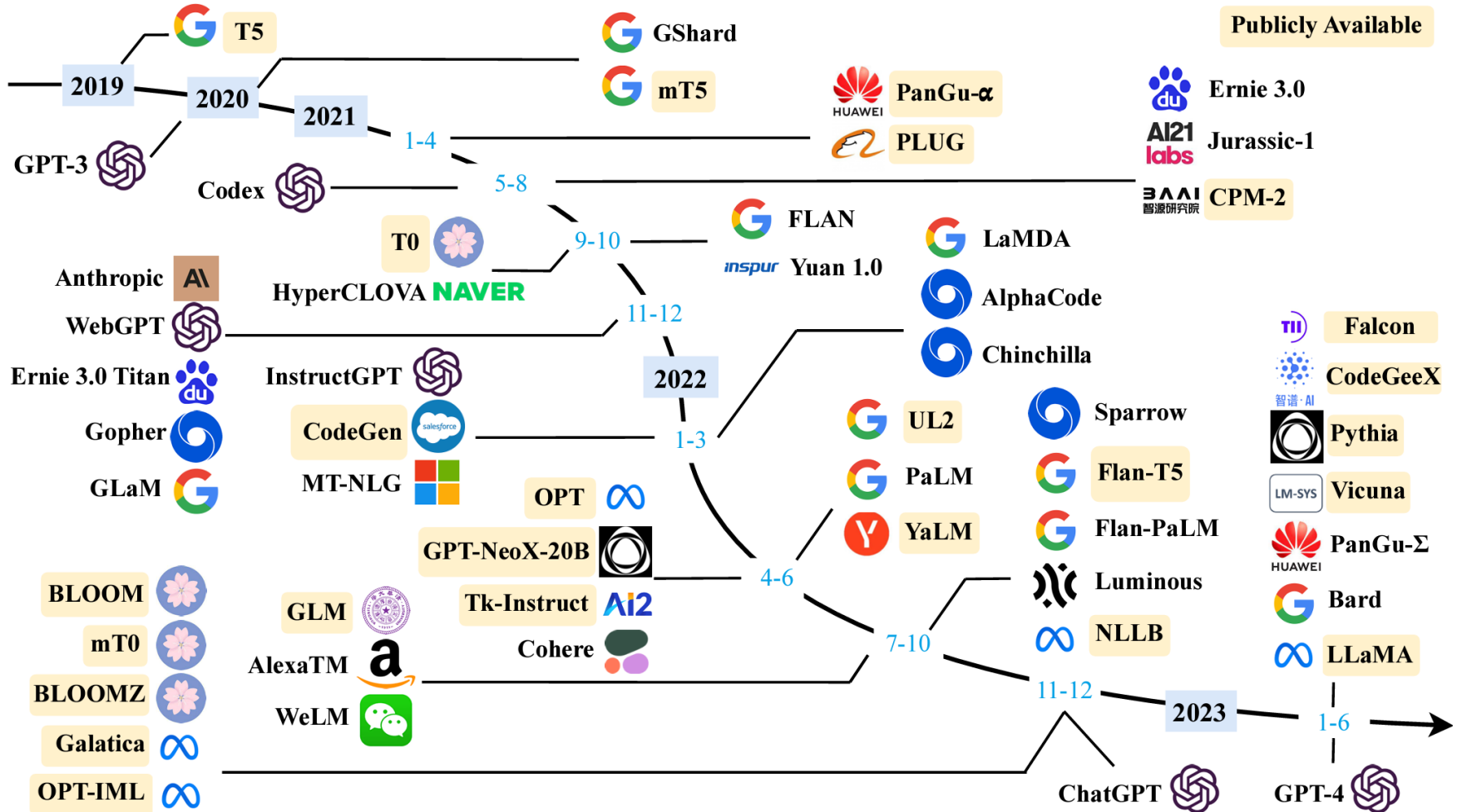
AlphaStar: StarCraft II



Alpha Go



Large Language Models (LLM)



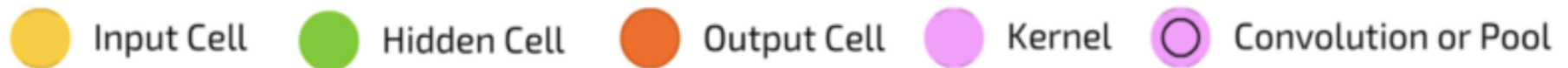
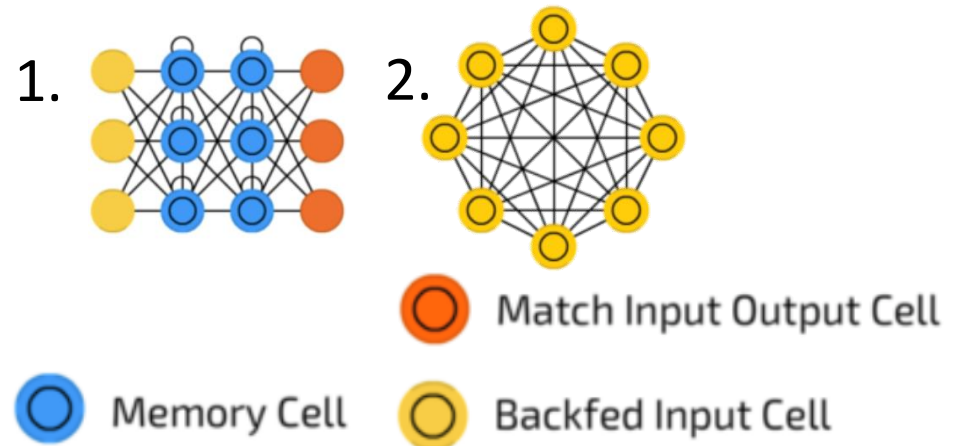
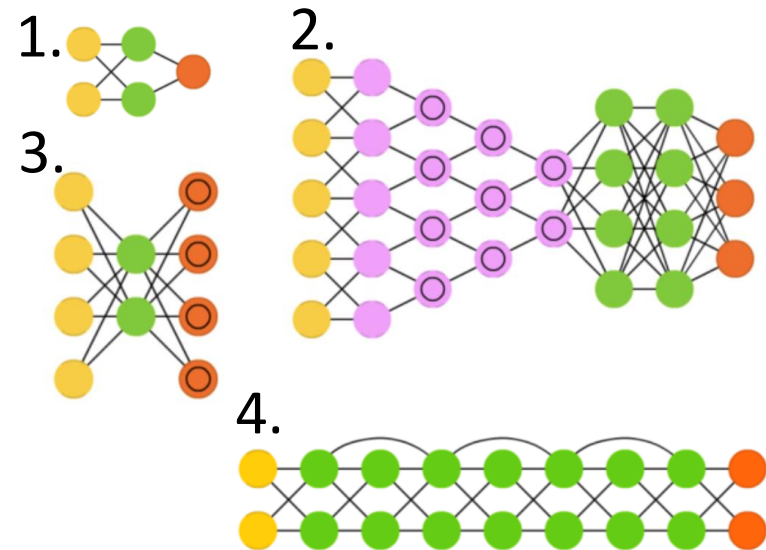
A Survey of Large Language Models, [Zhao et al., 2023]

Outline

- Course Introduction
- Machine Learning & Deep Neural Networks
 - Applications
 - Categories
 - Important Metrics
 - Platforms & Frameworks

Structures

- Feedforward neural network:
 1. Multilayer perceptron
 2. Convolutional neural network
 3. Autoencoder
 4. Deep residual network
- Recurrent neural network:
 1. Long short-term memory
 2. Hopfield
 3. ...
- Spiking neural network

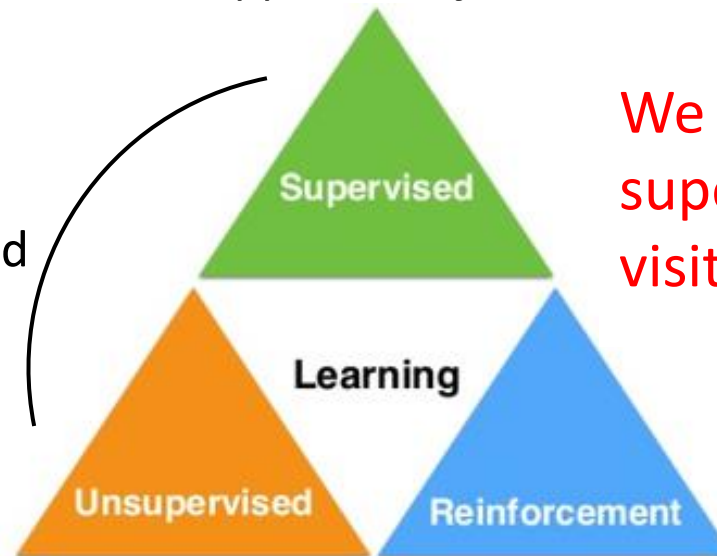


Learning Types

- Supervised
- Semi-supervised
- Unsupervised
- Reinforcement

Labeled data
Direct feedback
Predict outcome/future
Apps: Classification

Semi-supervised
Weakly-supervised



We will start with supervised learning and visit other topics later

No labels
No feedback
Find hidden representations
Apps: Reconstruction

Decision process
Reward system
Learn series of actions
Apps: Decision-making

Outline

- Course introduction
- Machine Learning & Deep Neural Networks
 - Applications
 - Categories
 - Important Metrics (**LASER**)
 - **L**atency
 - **A**ccuracy
 - **S**ize of Model
 - **E**nergy Efficiency
 - **R**obustness
 - Platforms & frameworks

Latency

- Latency is a measure of delay.
 - The length of time it takes for the data that you feed into one end of your network to emerge at the other end.
- Better accuracy? Longer latency!
- VGG-16 needs ~3s to process a single image on your smart phone, which is **unacceptable**.

Going Deeper!

Deep = Many hidden layers

Error Rate

16.4%

AlexNet
(2012)

7.3%

VGG
(2014)

6.7%

GoogleNet
(2014)

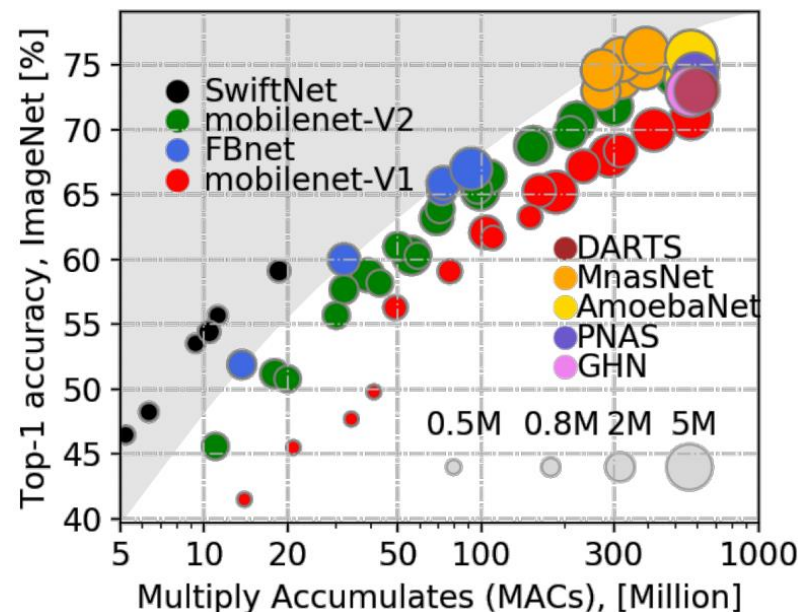
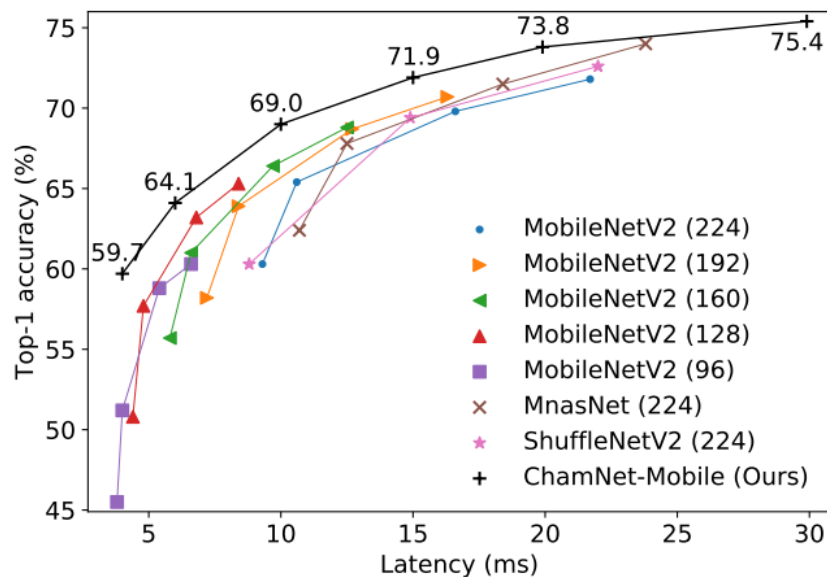
3.57%

152 layers

Residual Net
(2015)

Accuracy

- Accuracy is a metric for classification problem
- We call it: “Top-K Accuracy”
- Higher accuracy is good, but we need to pay for it
 - Everything is a trade-off.



Source:

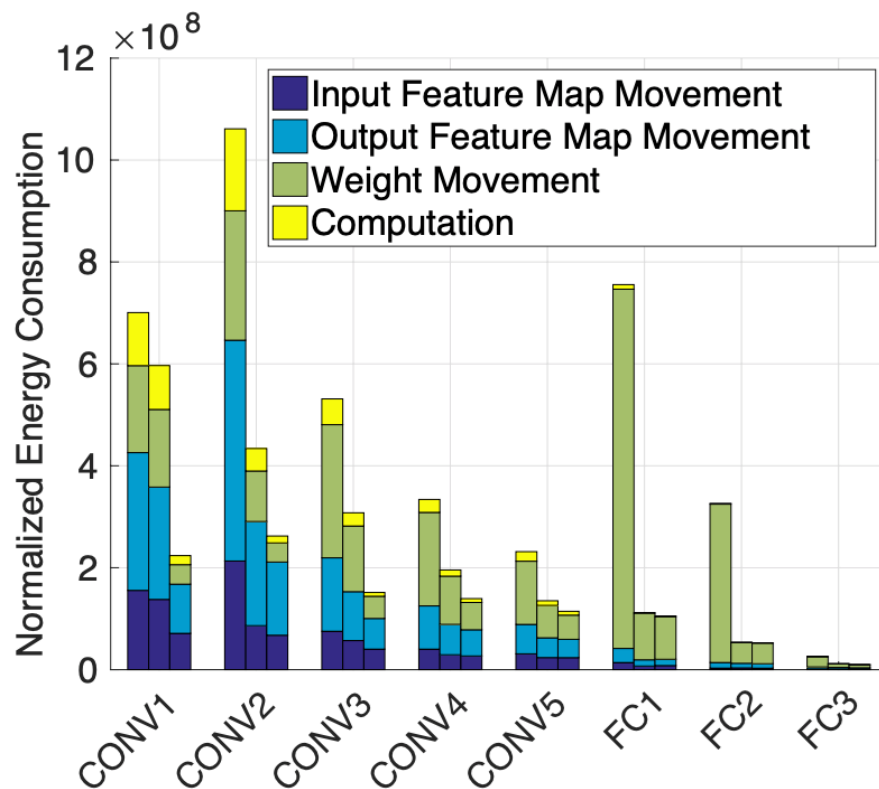
1. Dai, Xiaoliang, et al. "Chamnet: Towards efficient network design through platform-aware model adaptation." (2019)
2. Cheng, Hsin-Pai et al. "SwiftNet: Using Graph Propagation as Meta-knowledge to Search Highly Representative Neural Architectures" (2019)

Size of Model

- # FLOP: Number of floating point operations.
- # MAC: Number of multiply-and-accumulate operations
 - Usually, 1 floating-point multiply-and-accumulate is considered equivalent to 2 FLOPs.
- # Parameters
- Area [mm²]

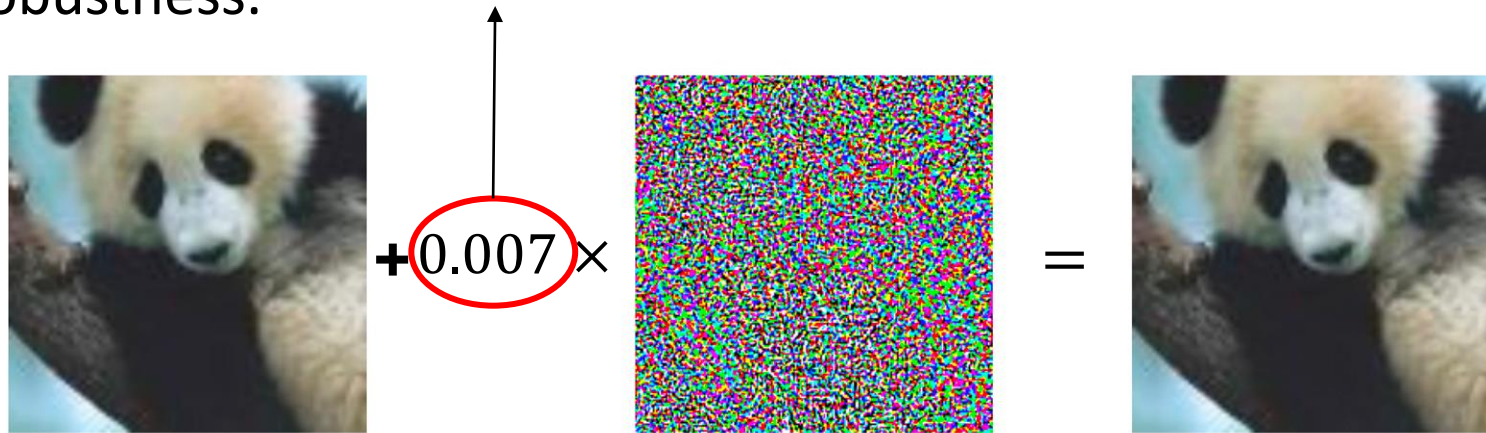
Energy Efficiency

- Power consumption [mW]
- Energy is mainly used for
 - Calculation
 - Data movement
- Energy is a different thing:
 - A lower number of MACs **does not** necessarily lead to lower energy consumption.
 - Convolutional layers **consume more** energy than fully-connected layers.
 - Deeper CNNs with fewer weights **do not** necessarily consume less energy than shallower CNNs with more weights.



Robustness

- This parameter, is used to evaluate a neural network's robustness.



Panda

57.7% confidence

noise

gibbon

99.3% confidence

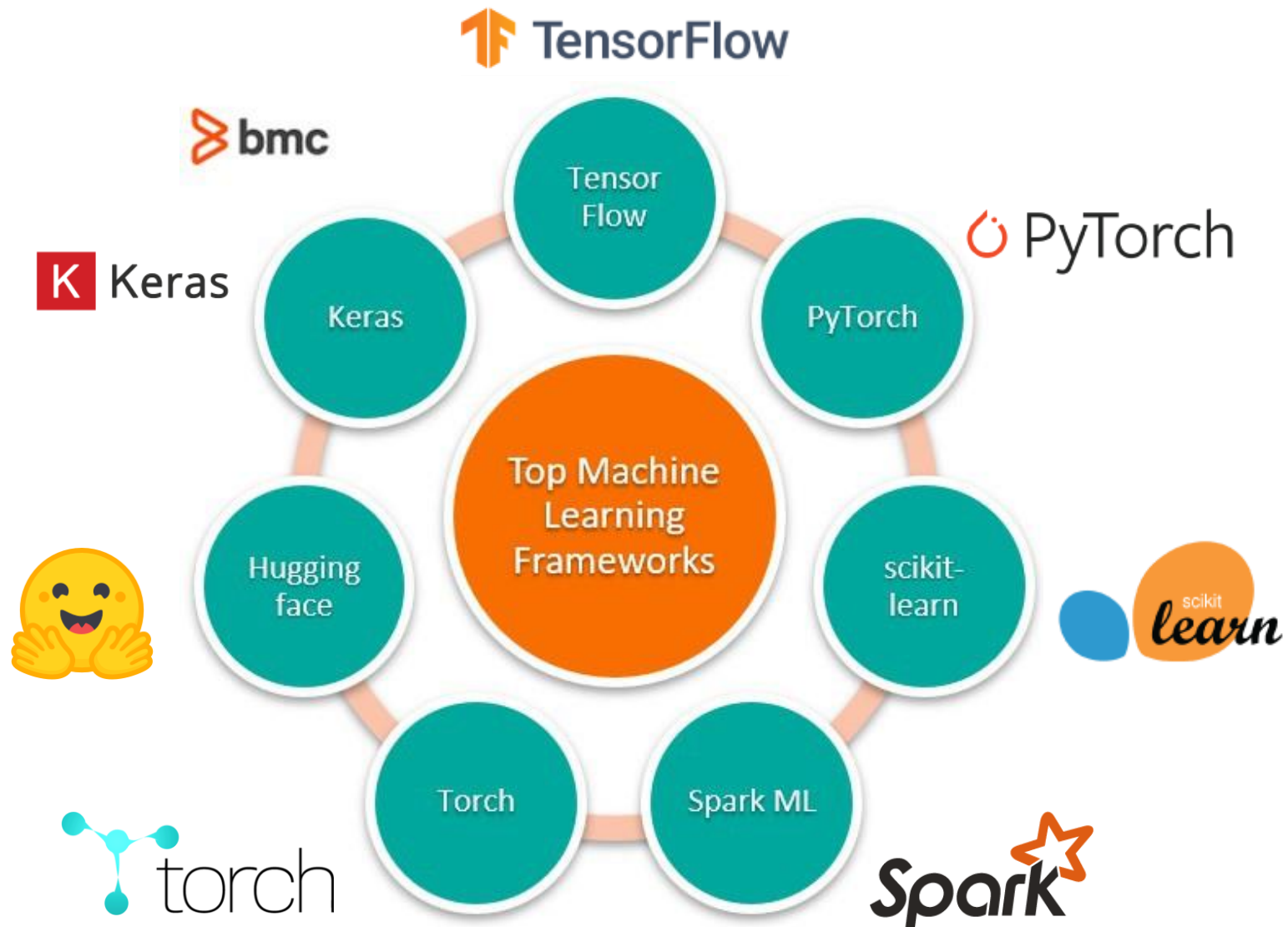
- Usually, a high accuracy model is not robust.
- Compared to the size of a neural network, the structure has more impact towards robustness.

Everything is
a trade-off

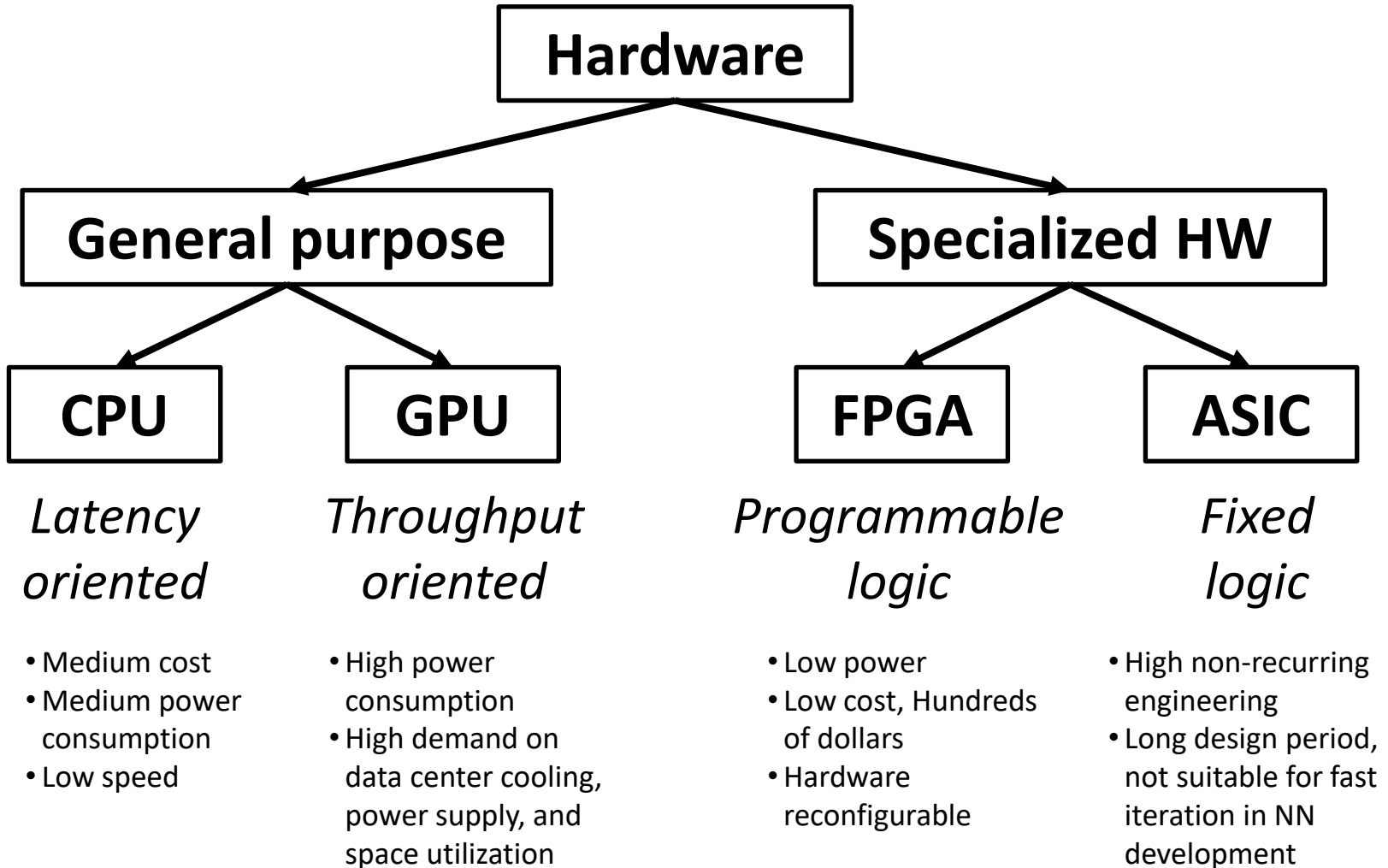
Outline

- Course introduction
- Machine learning & Deep Neural Networks
 - Applications
 - Categories
 - Important Metrics
 - Platforms & Frameworks
 - Software Platforms
 - Hardware Computing Devices

Top ML Frameworks



Hardware Computing Devices



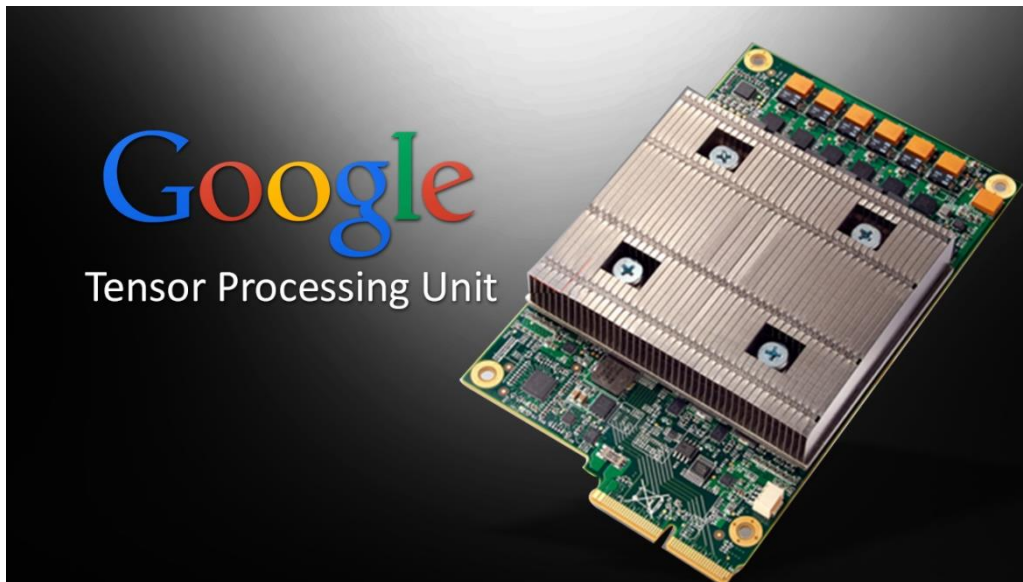
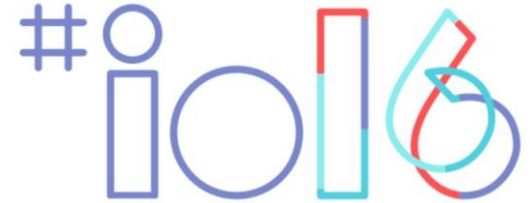
ImageNet-1K Classification Performance

Platform	Inference Throughput	Peak TFLOPs	Effective TFLOPs	Power	Power Efficiency GOPs/J
Intel Xeon E5-2450	53 images/s	0.27T	0.074T (27%)	~225W	~0.3
Altera Arria 10 GX1150	369 images/s	1.366T	0.51T (38%)	~40W	~12.8
NVIDIA Titan X	4129 images/s	6.1T	5.75T (94%)	~250W	~23.0

Neural network is usually trained in back-end GPU clusters, while FPGA is very suitable for low-power real-time inference job

Tensor Processing Unit (TPU)

- Unveiled during Google I/O Conference, Mountain View, CA (May 2016).
- Tensor Processing Unit (TPU): a custom ASIC built specifically for machine learning — and tailored for TensorFlow.
- This unit is designed for dense matrices, sparsity will have higher priority in the future.



Tensor Processing Unit (TPU)

Applications:

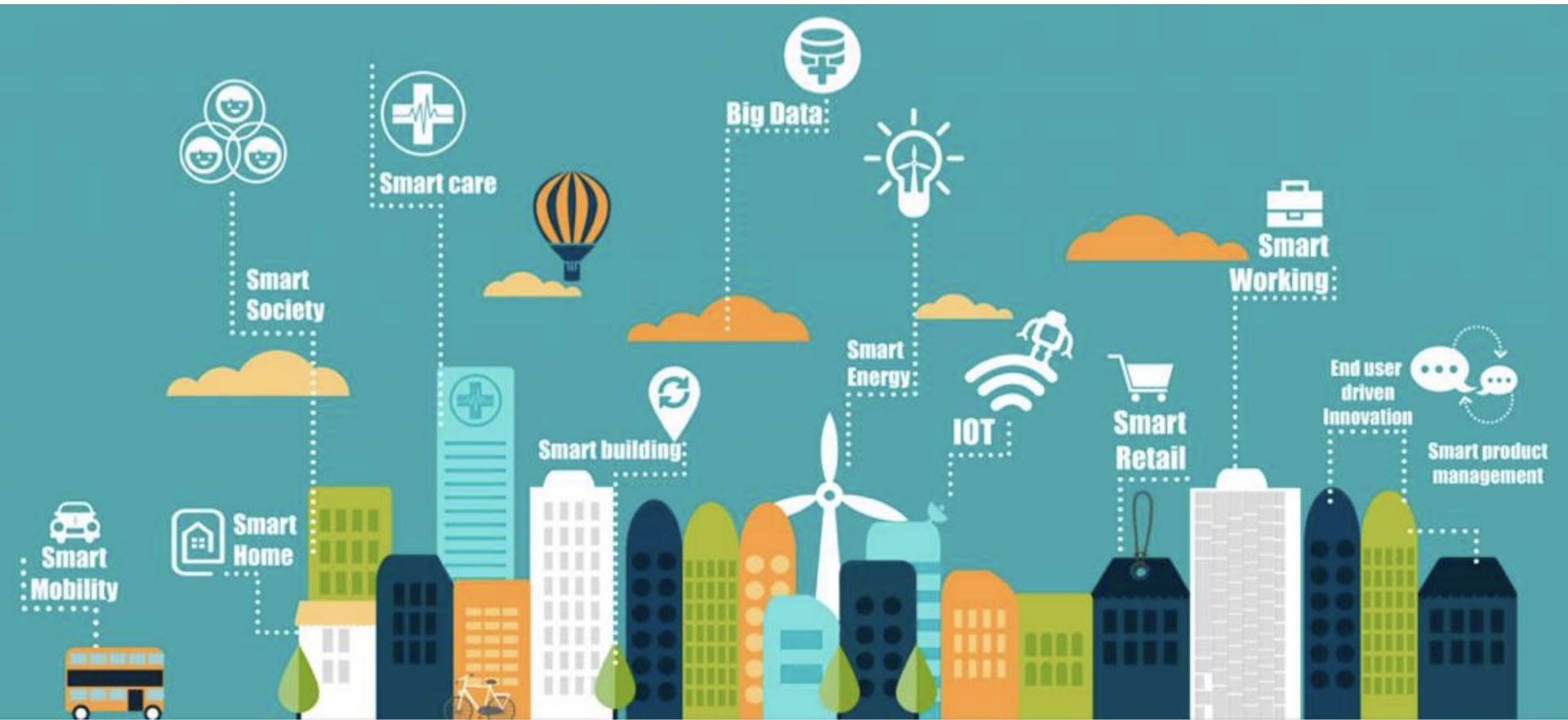
1. RankBrain: improve the relevancy of search results.
2. Street View: improve the accuracy and quality of our maps and navigation.
3. AlphaGo: “think” much faster and look farther ahead between moves.



Server racks with TPUs used in the AlphaGo matches with Lee Sedol



Future



Smart

Low latency

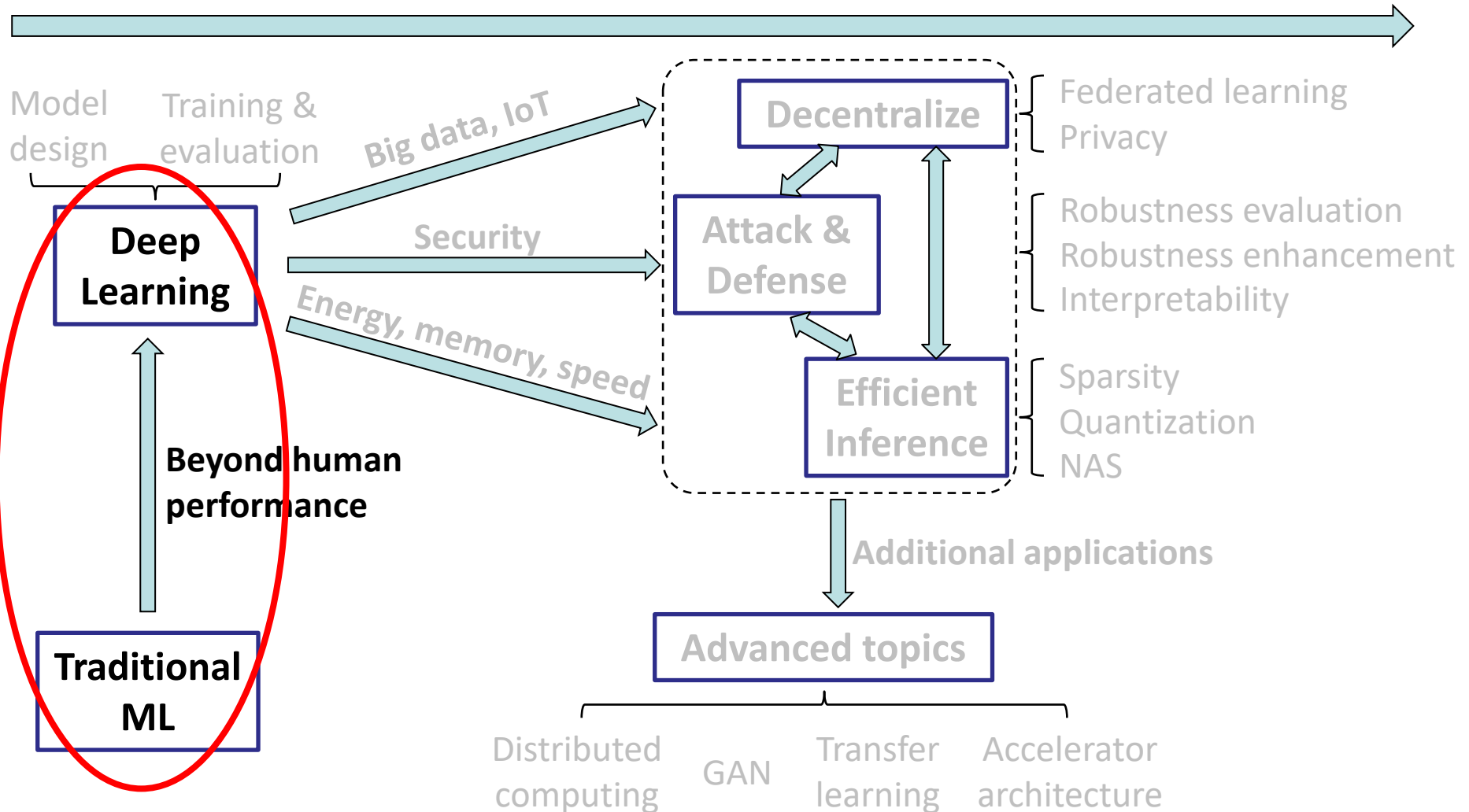
Privacy

Mobility

Energy efficient

Next Lecture

Applying machine learning into the real world



Reading Material

- Deep Learning (2016), Ian Goodfellow and Yoshua Bengio and Aaron Courville
<http://www.deeplearningbook.org/>
– Chapter “Introduction”

