

CIFAR CANADIAN INSTITUTE FOR ADVANCED RESEARCH

> Université 🇰 de Montréal

Yoshua Bengio

November 4, 2016

PLUG: Deep Learning, MIT Press book in presale, chapters online for feedback OSDC West, Santa Clara



Cars are now driving themselves...

(far from perfectly, though)



Unusual...

March 2016: World Go Champion Beaten by Machine

AI: The Upcoming Industrial Revolution

First industrial revolution:

 Machines extending humans' mechanical power

Upcoming industrial revolution:

- Machines extending humans' cognitive power
 - From the digital economy to the Al economy
 - Predicted growth at least 25%/yr
 - All sectors of the economy



A new revolution seems to be in the work after the industrial revolution.

Devices are becoming intelligent.

And Deep Learning is at the epicenter of this revolution.



Breakthrough in deep learning

A Canadian-led trio at CIFAR initiated the deep learning AI revolution

• Fundamental breakthrough in 2006:

first successful recipe for training a deep supervised neural network

- Second major advance in 2011, with rectifiers
- Breakthroughs in applications since then



Al Needs Knowledge

- Failure of classical AI: a lot of knowledge is not formalized, expressed with words
- Solution: computer gets knowledge from data, learns from examples

MACHINE LEARNING



Machine Learning, Al & No Free Lunch

- Five key ingredients for ML towards AI
 - 1. Lots & lots of data
 - 2. Very flexible models
 - 3. Enough computing power
 - 4. Powerful priors that can defeat the curse of dimensionality
 - 5. Computationally efficient inference

Bypassing the curse of dimensionality

We need to build compositionality into our ML models

Just as human languages exploit compositionality to give representations and meanings to complex ideas

Exploiting compositionality gives an exponential gain in representational power

Distributed representations / embeddings: feature learning

Deep architecture: multiple levels of feature learning

Prior assumption: compositionality is useful to describe the world around us efficiently

Non-distributed representations



- Clustering, n-grams, Nearest-Neighbors, RBF SVMs, local non-parametric density estimation & prediction, decision trees, etc.
- Parameters for each distinguishable region
- # of distinguishable regions is linear in # of parameters

 \rightarrow No non-trivial generalization to regions without examples

The need for distributed representations

- Factor models, PCA, RBMs, Neural Nets, Sparse Coding, Deep Learning, etc.
- Each parameter influences many regions, not just local neighbors
- # of distinguishable regions grows almost exponentially with # of parameters
- GENERALIZE NON-LOCALLY TO NEVER-SEEN REGIONS



Hidden Units Discover Semantically Meaningful Concepts

- Zhou et al & Torralba, arXiv1412.6856 , ICLR 2015
- Network trained to recognize places, not objects



Each feature can be discovered without the need for seeing the exponentially large number of configurations of the other features

Consider a network whose hidden units discover the following features:

Feature maps

- Person wears glasses
- Person is female
- Person is a child
- Etc.

If each of *n* feature requires *O(k)* parameters, need *O(nk)* examples

Non-parametric methods would require $O(n^d)$ examples

The Depth Prior can be Exponentially Advantageous

Theoretical arguments:



2 layers of - Logic gates Formal neurons RBF units

= universal approximator

RBMs & auto-encoders = universal approximator

Theorems on advantage of depth:

(Hastad et al 86 & 91, Bengio et al 2007, Bengio & Delalleau 2011, Braverman 2011, Pascanu et al 2014, Montufar et al **NIPS 2014**)

Some functions compactly represented with k layers may require exponential size with 2 layers



subroutine1 includes subsub1 code and subsub2 code and subsubsub1 code

subroutine2 includes subsub2 code and subsub3 code and subsub3 code and ...

"Shallow" computer program

mair



"Deep" computer program

Exponential advantage of depth

- Expressiveness of deep networks with piecewise linear activation functions: exponential advantage for depth (Montufar et al, NIPS 2014)
- Number of pieces distinguished for a network with depth L and n_i units per layer is at least

$$\left(\prod_{i=1}^{L-1} \left\lfloor \frac{n_i}{n_0} \right\rfloor^{n_0}\right) \sum_{j=0}^{n_0} \binom{n_L}{j}$$

or, if hidden layers have width n and input has size n_0

$$\Omega\left((n_{n_0})^{(L-1)n_0} n^{n_0}\right)$$

A Myth is Being Debunked: Local Minima in Neural Nets → Convexity is not needed

- (Pascanu, Dauphin, Ganguli, Bengio, arXiv May 2014): On the saddle point problem for non-convex optimization
- (Dauphin, Pascanu, Gulcehre, Cho, Ganguli, Bengio, NIPS' 2014): *Identifying and attacking the saddle point problem in high- dimensional non-convex optimization*
- (Choromanska, Henaff, Mathieu, Ben Arous & LeCun AISTATS 2015): *The Loss Surface of Multilayer Nets*

Saddle Points

- Local minima dominate in low-D, but⁴
 saddle points dominate in high-D
- Most local minima are close to the bottom (global minimum error)







2010-2012: breakthrough in speech recognition



Source: Microsoft

2012-2015: breakthrough in computer vision

- Graphics Processing Units (GPUs) + 10x more data
- 1,000 object categories,
- Facebook: millions of faces
- 2015: human-level performance



ImageNet Accuracy Still Improving

Top-5 Classification task



IT companies are racing into deep learning



From computer vision to self-driving cars: 2016

Holmdel, New Jersey February 2016

Ongoing progress: combining vision and natural language understanding



A woman is throwing a <u>frisbee</u> in a park.



A dog is standing on a hardwood floor



A <u>stop</u> sign is on a road with a mountain in the background

With a lot more data... visual question answering



.....



Tap here to ask a question about the photo



Recurrent Neural Networks

 Selectively summarize an input sequence in a fixed-size state vector via a recursive update

$$s_{t} = F_{\theta}(s_{t-1}, x_{t})$$

$$s_{t-1} \xrightarrow{F_{\theta}} \underbrace{s_{t+1}}_{f_{\theta} \text{ shared over time}} \xrightarrow{F_{\theta}} \underbrace{s_{t+1}}_{f_{\theta} \text{ shared over time}} \xrightarrow{F_{\theta}} \underbrace{s_{t+1}}_{x_{t-1} x_{t}} \xrightarrow{F_{\theta}} \underbrace{s_{t+1}}_{x_{t+1}}$$

$$s_t = G_t(x_t, x_{t-1}, x_{t-2}, \dots, x_2, x_1)$$

Generalizes naturally to new lengths not seen during training

Generative RNNs

 An RNN can represent a fully-connected directed generative model: every variable predicted from all previous ones.



Attention Mechanism for Deep Learning

(Bahdanau, Cho & Bengio, ICLR 2015; Jean et al ACL 2015; Jean et al WMT 2015; Xu et al ICML 2015; Chorowski et al NIPS 2015; Firat, Cho & Bengio 2016)

- Consider an input (or intermediate) sequence or image
- Consider an upper level representation, which can choose « where to look », by assigning a weight or probability to each input position, as produced by an MLP, applied at each position



End-to-End Machine Translation with Recurrent Nets and Attention Mechanism

(Bahdanau et al ICLR 2015, Jean et al ACL 2015, Gulcehre et al 2015, Firat et al 2016)

• Reached the state-of-the-art in one year, from scratch

	NMT(A)	Google	P-SMT	
NMT	32.68	30.6*		
+Cand	33.28	_	37.03 •	
+UNK	33.99	32.7°		
+Ens	36.71	36.9 °		

(a) English → French (WMT-14)

(b) EnglishightarrowGerman (WMT-15)

(c) English→Czech (WMT-15)

Model	Note	Model	Note	
24.8	Neural MT	18.3	Neural MT	
24.0	U.Edinburgh, Syntactic SMT	18.2	JHU, SMT+LM+OSM+Sparse	
23.6	LIMSI/KIT	17.6	CU, Phrase SMT	
22.8	U.Edinburgh, Phrase SMT	17.4	U.Edinburgh, Phrase SMT	
22.7	KIT, Phrase SMT	16.1	U.Edinburgh, Syntactic SMT	

Neural MT Contributions from Montreal

- Soft attention (Bahdanau et al ICLR 2015)
- Minibatch fast training with large vocabulary (Jean et al ACL 2015)
- Combining with neural language model (Gulcehre et al 2015)
- Subword and character-level NMT (Chung et al 2016)
- Multi-lingual NMT (Firat et al 2016)

Google-Scale NMT Success (Wu et al & Dean, Nature, 2016)

- After beating the classical phrase-based MT on the academic benchmarks, there remained the question: will it work on the very large scale datasets like used for Google Translate?
- Distributed training, very large model ensemble
- Not only does it work in terms of BLEU but it makes a killing in terms of human evaluation on Google Translate data

	PBMT	GNMT	Human	Relative		
				Improvement		
$English \rightarrow Spanish$	$3.594{\pm}1.58$	$5.031{\pm}1.09$	$5.140{\pm}1.04$	93%		
English \rightarrow French	$3.518{\pm}1.70$	$5.032{\pm}1.22$	$5.215{\pm}1.03$	89%		
English \rightarrow Portuguese	$3.675{\pm}1.64$	$4.856{\pm}1.29$	$4.973{\pm}1.17$	91%		
English \rightarrow Chinese	$2.457{\pm}1.48$	$4.154{\pm}1.42$	$4.580{\pm}1.26$	80%		
$\text{Spanish} \to \text{English}$	$3.410{\pm}1.65$	$4.921{\pm}1.16$	$4.930{\pm}1.12$	99%		
$\mathrm{French} \to \mathrm{English}$	$3.639{\pm}1.63$	$5.000{\pm}1.07$	$5.016{\pm}1.09$	99%		
$\text{Portuguese} \to \text{English}$	$3.471{\pm}1.74$	$5.029{\pm}1.05$	$5.040{\pm}1.03$	99%		
$\text{Chinese} \to \text{English}$	$1.994{\pm}1.47$	$3.884{\pm}1.37$	$4.334{\pm}1.20$	81%		

Table 10: Side-by-side scores on production data

Deep Learning: Beyond Pattern Recognition, towards Al

- Many researchers believed that neural nets could at best be good at pattern recognition
- And they are really good at it!
- But many more ingredients needed towards AI. Recent progress:
 - REASONING: with extensions of recurrent neural networks
 - Memory networks & Neural Turing Machine
 - PLANNING & REINFORCEMENT LEARNING: DeepMind (Atari and Go game playing) & Berkeley (Robotic control)

The next frontier: to reason and answer questions

Sam walks into the kitchen. Sam picks up an apple. Sam walks into the bedroom. Sam drops the apple.

Q: Where is the apple? A: Bedroom Brian is a lion. Julius is a lion. Julius is white Bernhard is green

Q: What colour is Brian? A: White The Biggest Challenge: Unsupervised Learning & Learning Commonsense Autonomously

- Recent progress mostly in supervised DL
- Real technical challenges for unsupervised DL
- Potential benefits:
 - Exploit tons of unlabeled data
 - Answer new questions about the variables observed
 - Regularizer transfer learning domain adaptation
 - Easier optimization (local training signal)
 - Structured outputs
 - Necessary for RL without given model or domain simulator

Learning « How the world ticks »

- So long as our machine learning models « cheat » by relying only on surface statistical regularities, they remain vulnerable to outof-distribution examples
- Humans generalize better than other animals by implicitly having a more accurate internal model of the underlying causal relationships
- This allows one to predict future situations (e.g., the effect of planned actions) that are far from anything seen before, an essential component of reasoning, intelligence and science

Invariance and Disentangling

- Invariant features
- Which invariances?



- Alternative: learning to disentangle factors
- Good disentangling →
 avoid the curse of dimensionality

Learning Multiple Levels of Abstraction

- The big payoff of deep learning is to allow learning higher levels of abstraction
- Higher-level abstractions disentangle the factors of variation, which allows much easier generalization and transfer



GAN: Generative Adversarial Networks

Goodfellow et al NIPS 2014





• 40% of samples mistaken *by humans* for real photos



- Sharper images than max. lik. proxys (which min. KL(data|model)):
- GAN objective = compromise between KL(data|model) and KL(model|data)

Convolutional GANs

(Radford et al, arXiv 1511.06343)

Strided convolutions, batch normalization, only convolutional layers, ReLU and leaky ReLU



GAN: Interpolating in Latent Space

If the model is good (unfolds the manifold), interpolating between latent values yields plausible images.











43



man without glasses

ALI: Adversarially Learned Inference

(Dumoulin et al 2016)

 Combines ideas from VAE and from GAN



Figure 1: The adversarially learned inference (ALI) game.





More Technical Challenges

- Learning long-term dependencies in recurrent neural networks
- Optimization challenge of training deep neural networks
- Taking advantage of feedback connections for attention, iterative inference & learning
- Incorporating "general knowledge" or commonsense (mostly from unsupervised learning) in RL

Applications on the horizon



Computer Interaction







Robotics

How to Attract the Best Researchers in Industry

- Extreme current demand for deep learning expertise, crazy salaries and acquisitions
- Not enough trained PhDs, too much industry demand
- Long-term open research
 - Necessary to attract and retain the strongest researchers
 - Success stories: DeepMind, FAIR, OpenAI
 - Need a pipeline & portfolio of different horizons
- Focused research: strategic, targeted choices
- Untying research org. from product-driven R&D

AI Corporate Research Strategy & Execution

- Difficult to reconcile
 - short-term pressure to deliver products and sales
 - creative & leading-edge AI research aiming at 5-10 year horizon

because the short-term guys have the money

- Need to have BOTH
 - 1. a firewall between the research organization and R&D
 - 2. a fluid path for people and ideas between the two

→ need to have independent funding for research and make it easy for (2) to happen, e.g., physical proximity, multiple "layers" in the pipeline.

Open Science & Open Source

- Best deep learning researchers (even in industry) demand open science →
 - Open and early publications (arXiv)
 - Accessible open source code (github)
- Both are
 - Reputation building (attracts more scientists)
 - Reproducible science
 - Generate follow-ups, citations & impact
 - Responsible: contribute to the community

Machine Learning Patents?

- ML scientists do not like ML patents because
 - work done at one company cannot be continued when the author of the work moves to another
 - algorithm is not available to the community, reducing the probability of follow-up by others, thus reducing the scientific impact (citations)
- ML scientists go to places with less IP constraints (OpenAI, FAIR)
- ML patents can easily be bypassed (different implementation) or are abusive (math patent)
- Patents only used for legal defense → the same can be achieved by arXiv posting

Montreal Institute for Learning Algorithms

Université

de Mont

MILA Faculty



Yoshua Bengio Director Aaron Courville

Pascal Vincent Roland Memisevic Christopher Pal



Laurent Charlin



Simon Lacoste-Julien



Doina Precup Joelle Pineau