# Probabilistic Graphical Models & Probabilistic AI

Ben Lengerich

Lecture 15: Deep Learning from a GM Perspective

April 1, 2025

Reading: See course homepage

# Entering Module 3: Modern Probabilistic AI

- Outstanding graded material:
  - Exam (20%, grades TBD)
  - Project midterm report (5%, 4/11)
  - Project presentation (5%, 4/31, 5/1)
    - Sign up here!
  - Project final report (15%, 5/5)
  - Extra credit (3%, sign-up)
- Module 3:
  - Papers > Textbooks

| Weeks | Lecture Dates | Topic | Assignments |
|-------|--------------|-------|-------------|
| **Module 1: Foundations of PGMs, Exact Inference** | | | |
| 1-4 | Jan 21- Feb 13 | Course Introduction, Foundations of PGMs, Exact Inference | HWs 1, 2 |
| 4 | Feb 13 | **Quiz** | |
| **Module 2: Learning** | | | |
| 5-9 | Feb 18 - Mar 18 | Parameter Learning, Structure Learning, Approximate Inference | HWs 3,4,5 |
| 9 | Mar 20 | **Midterm Exam** | |
| 10 | Mar 21 - Mar 30 | Spring Recess | |
| **Module 3: Modern Probabilistic AI** | | | |
| 11-14 | Apr 1 - Apr 24 | Deep Learning, LLMs from a GM perspective | Project Midway Report |
| 15 | Apr 29 - May 1 | Project Presentations | Project Final Report |

# A note on research papers
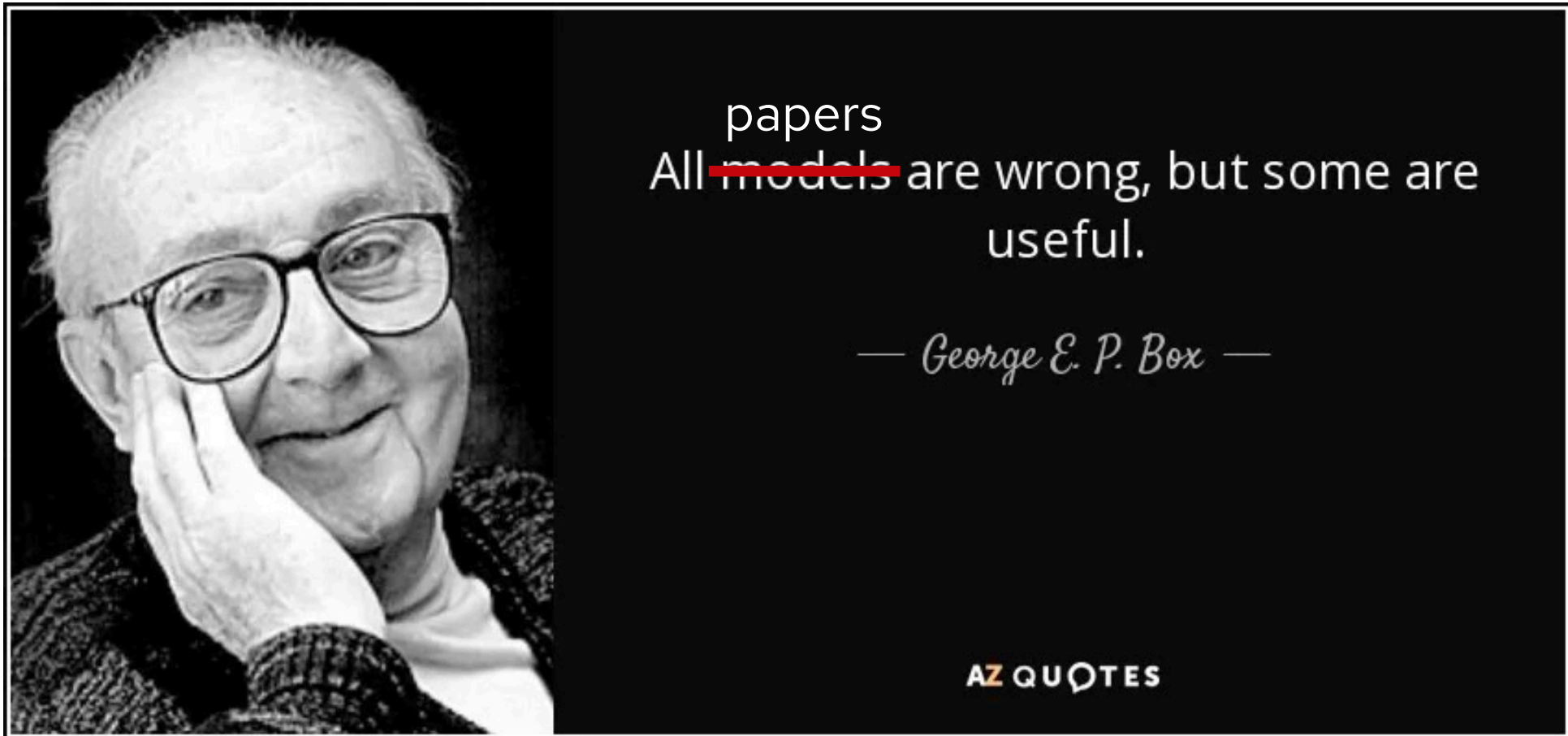
How we imagine research papers:

How research papers **actually** are:



Holes big enough to drive a car through!
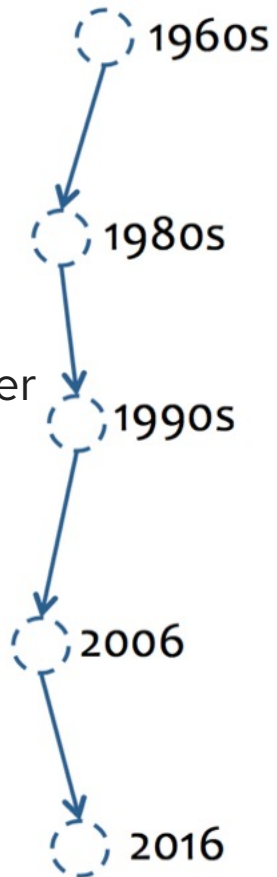
# A note on research papers → let's be optimists.



papers
All ~~models~~ are wrong, but some are useful.

— George E. P. Box —

AZ QUOTES

# Deep Learning from a GM Perspective

# History - Motivation



1973 – Pres. Gerald Ford viewing computer translation

1960s
1980s
1990s
2006
2016

## Deep learning:

– Has won numerous pattern recognition competitions
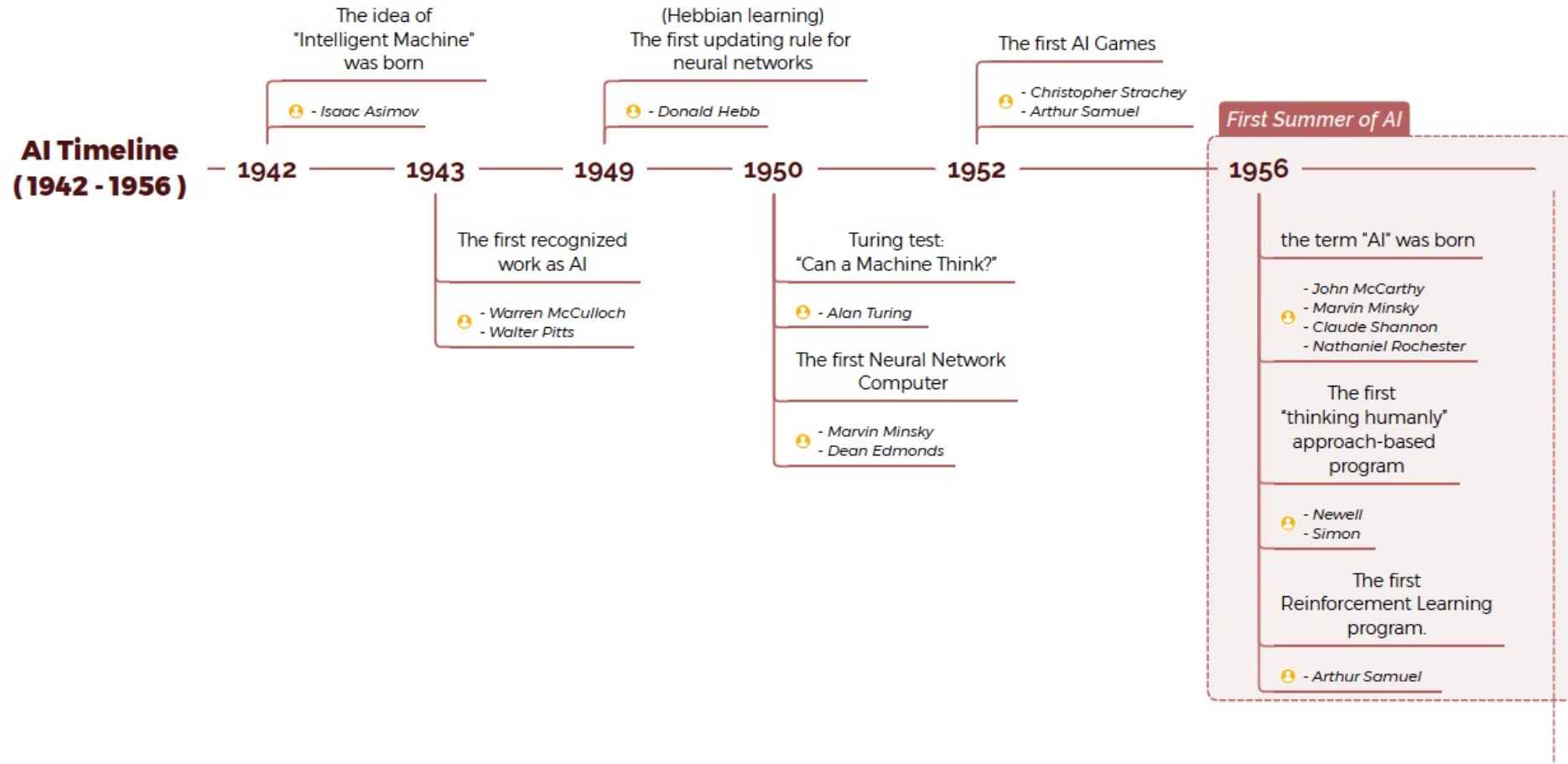
– Does so with minimal feature engineering

**This wasn't always the case!**
Since 1980s: Form of models hasn't changed much, but lots of new tricks…
– More hidden units
– Better (online) optimization
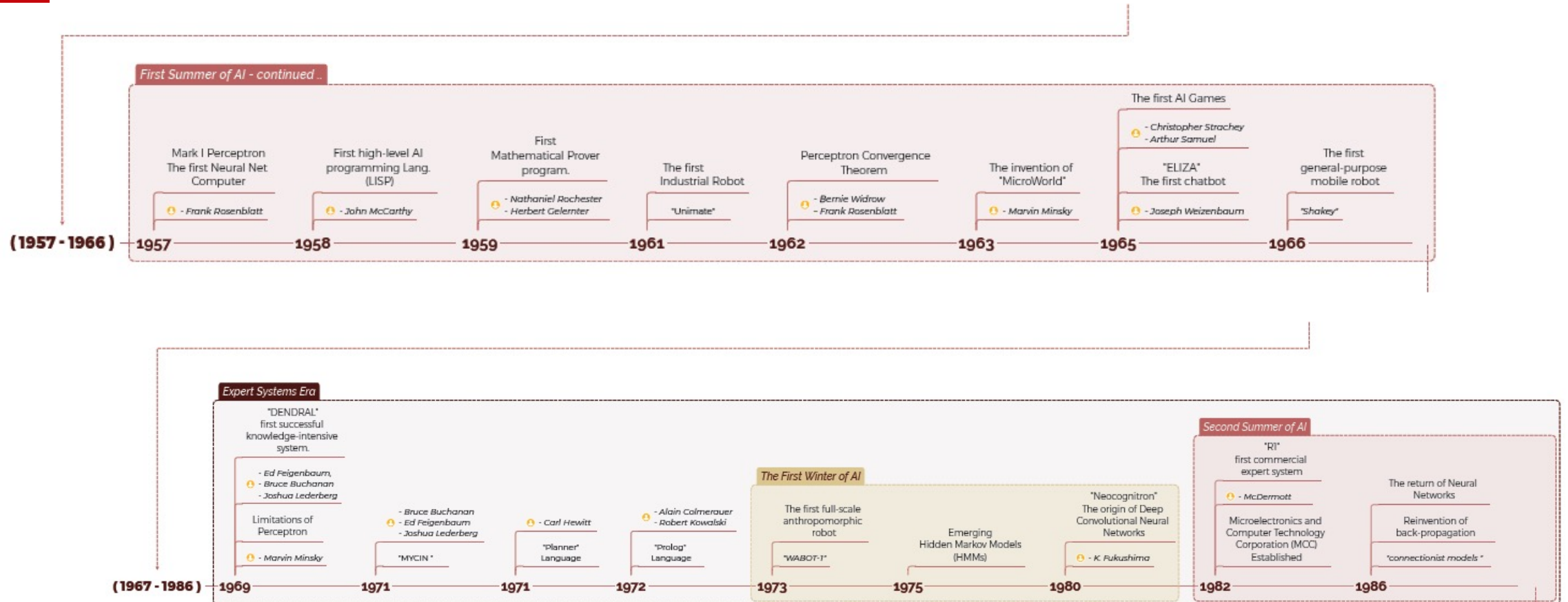– New nonlinear functions (ReLUs)
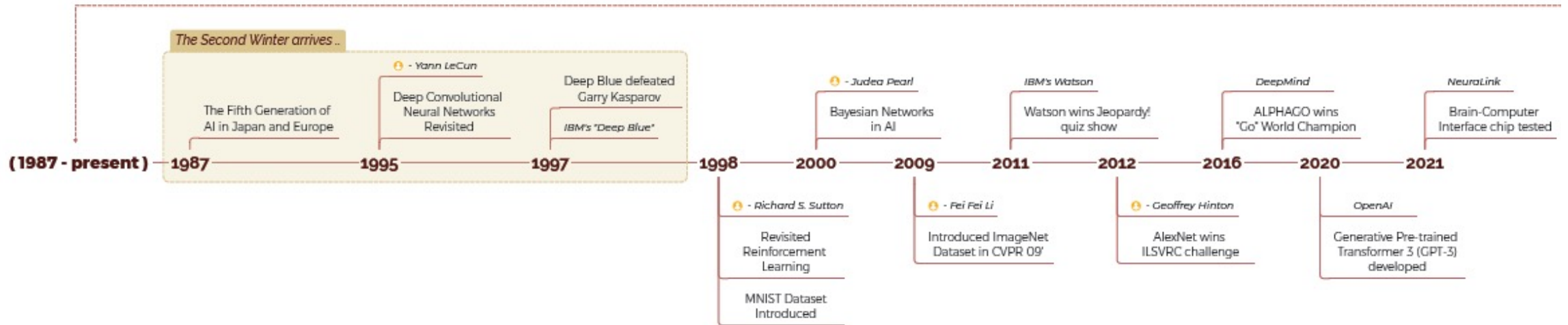– Faster computers (CPUs and GPUs)

# A brief history of AI



**AI Timeline (1942 - 1956)**

**1942** — The idea of "Intelligent Machine" was born
- Isaac Asimov

**1943** — The first recognized work as AI
- Warren McCulloch
- Walter Pitts

**1949** — (Hebbian learning) The first updating rule for neural networks
- Donald Hebb

**1950** — Turing test: "Can a Machine Think?"
- Alan Turing

The first Neural Network Computer
- Marvin Minsky
- Dean Edmonds

**1952** — The first AI Games
- Christopher Strachey
- Arthur Samuel

**1956** — *First Summer of AI*

the term "AI" was born
- John McCarthy
- Marvin Minsky
- Claude Shannon
- Nathaniel Rochester

The first "thinking humanly" approach-based program
- Newell
- Simon

The first Reinforcement Learning program.
- Arthur Samuel

[Toosi et al 2021]

# A brief history of AI



First Summer of AI - continued ...

**The first AI Games**
- Christopher Strachey
- Arthur Samuel

**Mark I Perceptron**
The first Neural Net Computer
- Frank Rosenblatt

**First high-level AI programming Lang. (LISP)**
- John McCarthy

**First Mathematical Prover program.**
- Nathaniel Rochester
- Herbert Gelernter

**The first Industrial Robot**
"Unimate"

**Perceptron Convergence Theorem**
- Bernie Widrow
- Frank Rosenblatt

**The invention of "MicroWorld"**
- Marvin Minsky

**"ELIZA"**
The first chatbot
- Joseph Weizenbaum

**The first general-purpose mobile robot**
"Shakey"

(1957 - 1966)   1957   1958   1959   1961   1962   1963   1965   1966

Expert Systems Era

**"DENDRAL"**
first successful knowledge-intensive system.
- Ed Feigenbaum,
- Bruce Buchanan
- Joshua Lederberg

**Limitations of Perceptron**
- Marvin Minsky

- Bruce Buchanan
- Ed Feigenbaum
- Joshua Lederberg
"MYCIN"

- Carl Hewitt
"Planner" Language

- Alain Colmerauer
- Robert Kowalski
"Prolog" Language

The First Winter of AI

**The first full-scale anthropomorphic robot**
"WABOT-1"

**Emerging Hidden Markov Models (HMMs)**

**"Neocognitron"**
The origin of Deep Convolutional Neural Networks
- K. Fukushima

Second Summer of AI

**"R1"**
first commercial expert system
- McDermott

**Microelectronics and Computer Technology Corporation (MCC) Established**

**The return of Neural Networks**
Reinvention of back-propagation
"connectionist models"

(1967 - 1986)   1969   1971   1971   1972   1973   1975   1980   1982   1986

[Toosi et al 2021]

# A brief history of AI



The Second Winter arrives ..

**(1987 - present)**

- **1987** — The Fifth Generation of AI in Japan and Europe
- **1995** — 👤 - Yann LeCun — Deep Convolutional Neural Networks Revisited
- **1997** — Deep Blue defeated Garry Kasparov — IBM's "Deep Blue"
- **1998** — 👤 - Richard S. Sutton — Revisited Reinforcement Learning; MNIST Dataset Introduced
- **2000** — 👤 - Judea Pearl — Bayesian Networks in AI
- **2009** — 👤 - Fei Fei Li — Introduced ImageNet Dataset in CVPR 09'
- **2011** — IBM's Watson — Watson wins Jeopardy! quiz show
- **2012** — 👤 - Geoffrey Hinton — AlexNet wins ILSVRC challenge
- **2016** — DeepMind — ALPHAGO wins "Go" World Champion
- **2020** — OpenAI — Generative Pre-trained Transformer 3 (GPT-3) developed
- **2021** — NeuraLink — Brain-Computer Interface chip tested

[Toosi et al 2021]

# From biological neuron to artificial neuron



- McCulloch & Pitts neuron – **Threshold only**



Warren McCulloch



Walter Pitts

# From biological neuron to artificial neuron

- McCulloch & Pitts neuron – **Threshold only**

- Can represent "AND", "OR"

- But not "NOT", "XOR"

# Perceptrons generalize MP neurons



**Weighting!**

**Activation function**

$$net = \sum_{i=0}^{n} w_i x_i$$

$$o = \sigma(net) = \frac{1}{1 + e^{-net}}$$

CORNELL AERONAUTICAL LABORATORY, INC.

Report No. 85-460-1

THE PERCEPTRON
A PERCEIVING AND RECOGNIZING AUTOMATON
(PROJECT PARA)

January, 1957

# Perceptrons generalize MP neurons



**Weighting!**

**Activation function**

$$net = \sum_{i=0}^{n} w_i x_i$$

$$o = \sigma(net) = \frac{1}{1 + e^{-net}}$$

- Consider regression problem f: X→Y for scalar Y
  - Let $Y \sim N(f(x), \Sigma^2)$
  - Then $\text{argmax}_w \log \prod_i P(y_i \mid x_i; w) = \text{argmin}_w \sum_i \frac{1}{2}(y_i - f(x_i; w))^2$

# Perceptron learning algorithm

- Recall the nice property of sigmoid: $\frac{d\sigma}{dt} = \sigma(1-\sigma)$

$$
\begin{aligned}
\frac{\partial E_D[\vec{w}])}{\partial w_j} &= \frac{\partial}{\partial w_i} \frac{1}{2} \sum_d (t_d - o_d)^2 \\
&= \frac{1}{2} \sum_d 2(t_d - o_d) \frac{\partial}{\partial w_i}(t_d - o_d) \\
&= \sum_d (t_d - o_d)\left(-\frac{\partial o_d}{\partial w_i}\right) \\
&= -\sum_d (t_d - o_d) \frac{\partial o_d}{\partial net_i} \frac{\partial net_d}{\partial w_i} \\
&= -\sum_d (t_d - o_d) o_d (1 - o_d) x_d^i
\end{aligned}
$$

$x_d$ = input

$t_d$ = target output

$o_d$ = observed output

$w_i$ = weight i|

Batch mode:

Do until converge:

1. compute gradient $\nabla E_D[w]$
2. $\vec{w} = \vec{w} - \eta \nabla E_D[\vec{w}]$

Incremental mode:

Do until converge:

- For each training example $d$ in $D$
  1. compute gradient $\nabla E_d[w]$
  2. $\vec{w} = \vec{w} - \eta \nabla E_d[\vec{w}]$

where

$\nabla E_d[\vec{w}] = -(t_d - o_d) o_d (1 - o_d) \vec{x}_d$

# Can a Perceptron represent XOR?



**Weighting!**

**Activation function**

$$net = \sum_{i=0}^{n} w_i x_i$$

$$o = \sigma(net) = \frac{1}{1 + e^{-net}}$$

- No!
- If there were, then there would be constants $w_1$ and $w_2$ such that:
  - When $x_1 = x_2$, then $\sigma(w_1 x_1 + w_2 x_2) < \theta$
  - When $x_1 \neq x_2$, then $\sigma(w_1 x_1 + w_2 x_2) \geq \theta$
  - Let $x_1 = 1, x_2 = 0$
    - Eq. (1): $\sigma(w_1) \geq \theta$
  - Let $x_1 = 0, x_2 = 1$
    - Eq. (2): $\sigma(w_2) \geq \theta$
  - Let $x_1 = 1, x_2 = 1$:
    - Eq. (3): $\sigma(w_1 + w_2) < \theta$

  **Eq. (1) + Eq. (2) contradicts Eq. (3)**

# An XOR Logic Gate



**Multi-layer Perceptron?**

$$net = \sum_{i=0}^{n} w_i\, x_i$$

$$o = \sigma(net) = \frac{1}{1 + e^{-net}}$$

https://byjus.com/jee/basic-logic-gates/

# Neural Network Model (MLP)

# "Combined Logistic Models"...

# "Combined Logistic Models"...

# "Combined Logistic Models"...

# ...But no target for hidden units

# Backpropagation

- Neural networks are function compositions that can be represented as computation graphs:



$x$ — Input variables

Intermediate computations

$f(x)$ — Outputs

$$\frac{\partial f_n}{\partial x} =$$

- By applying the chain rule, and working in reverse order, we get:

$$\frac{\partial f_n}{\partial x} = \sum_{i_1 \in \pi(n)} \frac{\partial f_n}{\partial f_{i_1}} \frac{\partial f_{i_1}}{\partial x} = \sum_{i_1 \in \pi(n)} \frac{\partial f_n}{\partial f_{i_1}} \sum_{i_2 \in \pi(i_1)} \frac{\partial f_{i_1}}{\partial f_{i_2}} \frac{\partial f_{i_1}}{\partial x} = \ldots$$

# Model building blocks

- **Activation functions**
  - Linear and ReLU



$x_1$ $w_1$

$x_2$ $w_2$ $\quad$ f $\quad$ $f(Wx + b)$

$x_3$ $w_3$

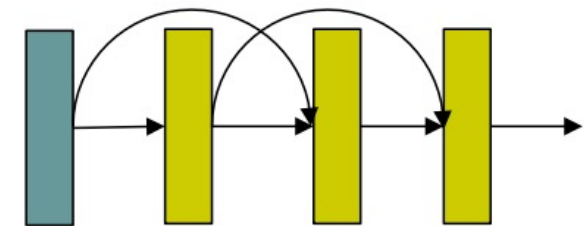Linear

Rectified linear

# Model building blocks

- ## Activation functions
  - Linear and ReLU
  - Sigmoid and tanh
  - Etc.



$x_1 \quad w_1$

$x_2 \quad w_2 \quad f \quad f(Wx + b)$

$x_3 \quad w_3$

**Sigmoid**

**Hyperbolic tangent**

# Model building blocks

- **Activation functions**
  - Linear and ReLU
  - Sigmoid and tanh
  - Etc.

- **Layers**
  - Fully connected
  - Convolutional & pooling
  - Recurrent
  - ResNets
  - Etc.

# Model building blocks

- Activation functions
  - Linear and ReLU
  - Sigmoid and tanh
  - Etc.
- Layers
  - Fully connected
  - Convolutional & pooling
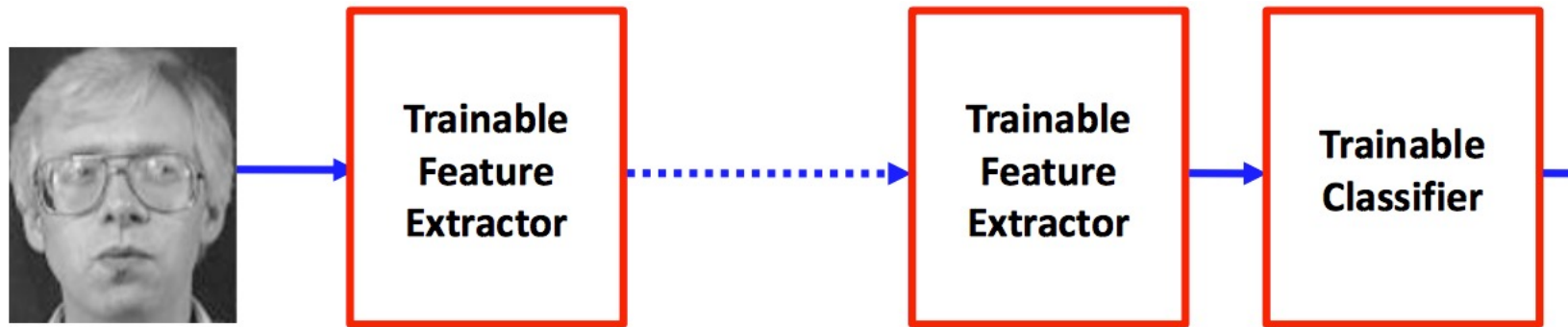  - Recurrent
  - ResNets
  - Etc.
- Loss functions
  - Cross-entropy loss
  - Mean squared error
  - Etc.



- Arbitrary combinations of the basic building blocks
- Multiple loss functions – multi-target prediction, transfer learning, and more
- Given enough data, deeper architectures just keep improving
- Representation learning: the networks learn increasingly more abstract representations of the data that are "disentangled," i.e., amenable to linear separation.

# Using DNNs for hierarchical representations



- In Language: hierarchy in syntax and semantics
  - Words → Parts of Speech → Sentences → Text
  - Objects, Actions, Attributes... → Phrases → Statements → Stories
- In Vision: part-whole hierarchy
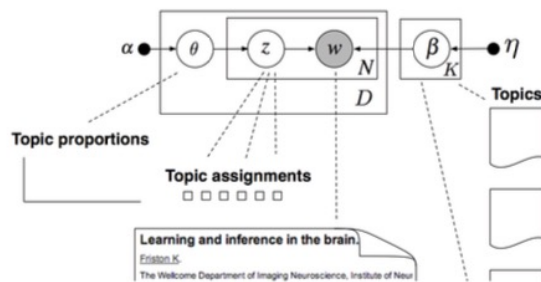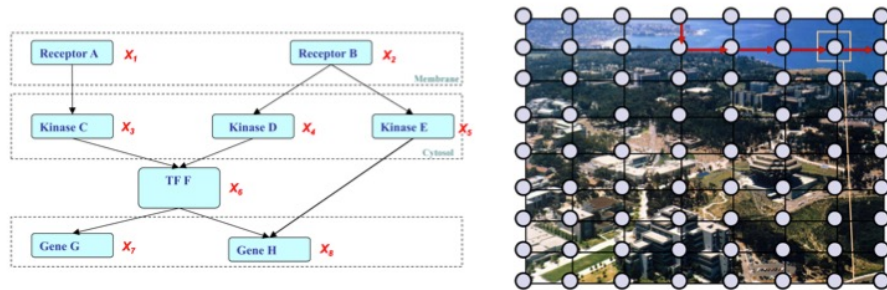  - Pixels → Edges → Textons → Parts → Objects → Scenes

| | DL | ⪌ ? ML (e.g., GM) |
|---|---|---|
| Empirical goal: | e.g., classification, feature learning | e.g., latent variable inference, transfer learning |
| Structure: | Graphical | Graphical |
| Objective: | Something aggregated from local functions | Something aggregated from local functions |
| | | |
| Vocabulary: | Neuron, activation function, … | Variable, potential function, … |
| Algorithm: | A single, unchallenged, inference algorithm – Backpropagation (BP) | A major focus of open research, many algorithms, and more to come |
| Evaluation: | On a black-box score – end performance | On almost every intermediate quantity |
| Implementation: | Many tricks | More or less standardized |
| Experiments: | Massive, real data (GT unknown) | Modest, often simulated data (GT known) |

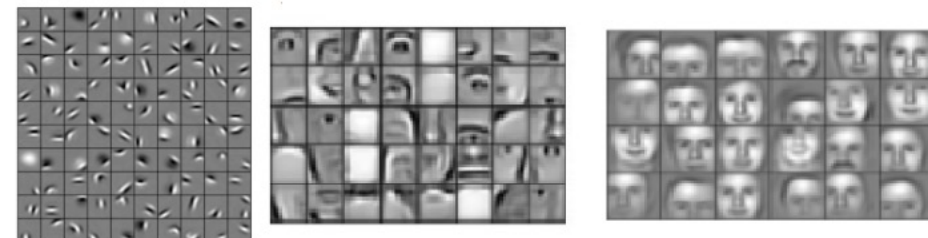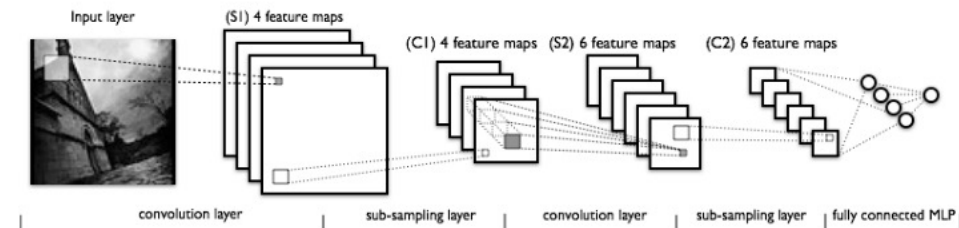# Graphical models vs Deep nets

## Graphical models

- Representation for encoding meaningful knowledge and the associated uncertainty in a graphical form



## Deep neural networks

- Learn representations that facilitate computation and performance on the end-metric (intermediate representations may not be meaningful)

# Graphical models vs Deep nets

## Graphical models

- <u>Representation</u> for encoding meaningful knowledge and the associated uncertainty in a graphical form

- <u>Learning and inference</u> are based on a rich toolbox of well-studied (structure-dependent) techniques (e.g., EM, message passing, VI, MCMC, etc.)

- Graphs <u>represent models</u>

## Deep neural networks

- <u>Learn representations</u> that facilitate computation and performance on the end-metric (intermediate representations may not be meaningful)

- <u>Learning</u> is predominantly based on the gradient descent method (aka backpropagation); <u>Inference</u> is often trivial and done via a "forward pass"

- Graphs <u>represent computation</u>
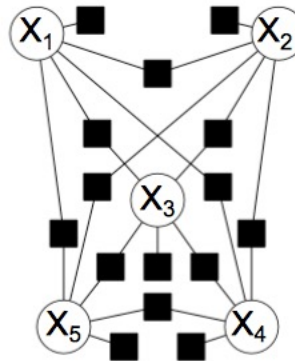
# Graphical models vs Deep nets

## Graphical models
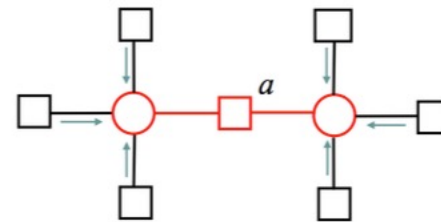
### Utility of the graph

- A vehicle for synthesizing a global loss function from local structure
  - potential function, feature function, etc.
- A vehicle for designing sound and efficient inference algorithms
  - Sum-product, mean-field, etc.
- A vehicle to inspire approximation and penalization
  - Structured MF, Tree-approximation, etc.
- A vehicle for monitoring theoretical and empirical behavior and accuracy of inference

### Utility of the loss function

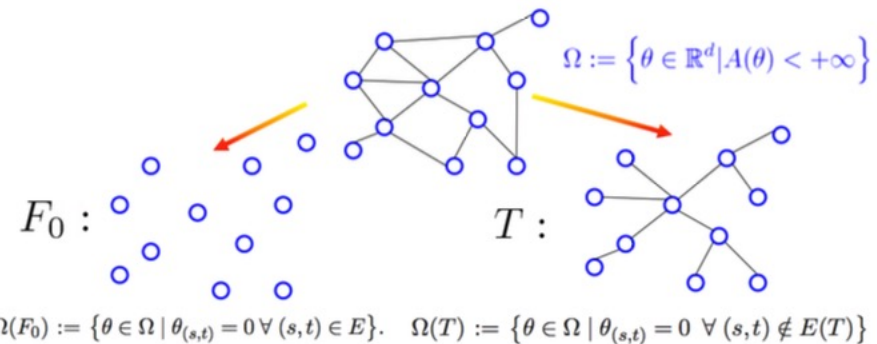- A major measure of quality of the learning algorithm and the model

$$\log P(X) = \sum_i \log \phi(x_i) + \sum_{i,j} \log \psi(x_i, x_j)$$

$$m_{i \to a}(x_i) = \prod_{c \in N(i) \setminus a} m_{c \to i}(x_i)$$

$$b_a(X_a) \propto f_a(X_a) \prod_{i \in N(a)} m_{i \to a}(x_i)$$

$$m_{a \to i}(x_i) = \sum_{X_a \setminus x_i} f_a(X_a) \prod_{j \in N(a) \setminus i} m_{j \to a}(x_j)$$

$$\Omega := \left\{ \theta \in \mathbb{R}^d \mid A(\theta) < +\infty \right\}$$

$F_0:$ $\qquad$ $T:$

$$\Omega(F_0) := \left\{ \theta \in \Omega \mid \theta_{(s,t)} = 0 \, \forall \, (s,t) \in E \right\}. \qquad \Omega(T) := \left\{ \theta \in \Omega \mid \theta_{(s,t)} = 0 \, \forall \, (s,t) \notin E(T) \right\}$$

# Graphical models vs Deep nets

## Graphical models

### Utility of the graph

- A vehicle for synthesizing a global loss function from local structure
  - potential function, feature function, etc.
- A vehicle for designing sound and efficient inference algorithms
  - Sum-product, mean-field, etc.
- A vehicle to inspire approximation and penalization
  - Structured MF, Tree-approximation, etc.
- A vehicle for monitoring theoretical and empirical behavior and accuracy of inference

### Utility of the loss function

- A major measure of quality of the learning algorithm and the model

## Deep neural networks

### Utility of the network

- A vehicle to conceptually synthesize complex decision hypothesis
  - stage-wise projection and aggregation
- A vehicle for organizing computational operations
  - stage-wise update of latent states
- A vehicle for designing processing steps/computing modules
  - Layer-wise parallelization
- No obvious utility in evaluating DL inference algorithms
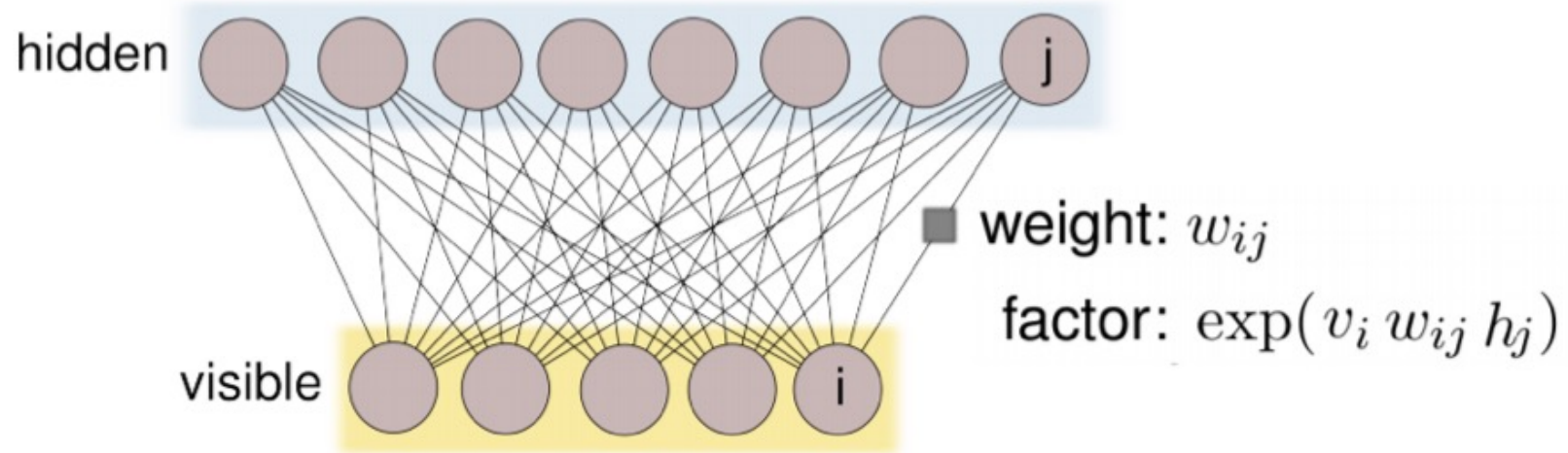
### Utility of the Loss Function

- Global loss? Well it is complex and non-convex...

# Sometimes nets are proposed as true GMs:

- Boltzmann machines (Hinton & Sejnowsky, 1983)
- Restricted Boltzmann machines (Smolensky, 1986)
- Learning and Inference in sigmoid belief networks (Neal, 1992)
- Fast learning in deep belief networks (Hinton, Osindero, Teh, 2006)
- Deep Boltzmann machines (Salakhutdinov and Hinton, 2009)

# Restricted Boltzmann Machines

- Assume visible units are one layer, and hidden units are another.
- Throw out all the connections within each layer.



weight: $w_{ij}$

factor: $\exp(v_i\, w_{ij}\, h_j)$
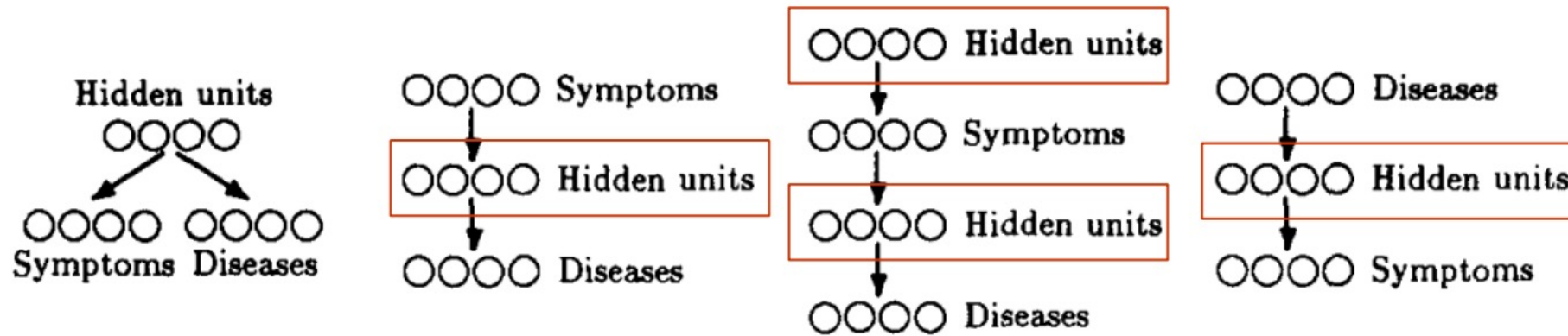
# Restricted Boltzmann Machines

$$\frac{\partial}{\partial w} \log L \;\propto$$

$$\underbrace{\frac{1}{N} \sum_{\mathbf{v} \in \mathcal{D}}}_{\text{data}} \underbrace{\sum_{\mathbf{h}} P(\mathbf{h} \mid \mathbf{v})}_{\text{av. over posterior}} \frac{\partial}{\partial w} \log P^\star(\mathbf{x}) \;-\; \underbrace{\sum_{\mathbf{v},\mathbf{h}} P(\mathbf{v},\mathbf{h})}_{\text{av. over joint}} \frac{\partial}{\partial w} \log P^\star(\mathbf{x})$$

Both terms involve averaging over $\frac{\partial}{\partial w} \log P^\star(\mathbf{x})$.

Contrastive Divergence estimates the second term with a Monte Carlo estimate from 1-step of a Gibbs sampler!
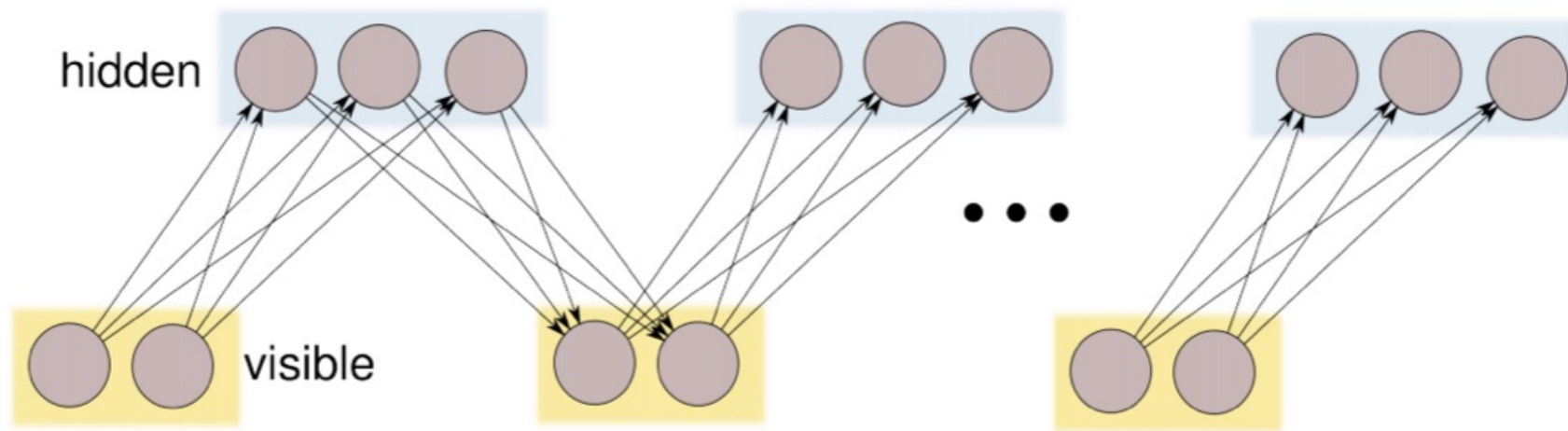
# Sigmoid Belief Networks



Sigmoid belief nets are simply Bayes networks conditionals represented in a particular form:

$$P(S_i = x \mid S_j = s_j : j \neq i)$$

$$\propto P(S_i = x \mid S_j = s_j : j < i)$$

$$\cdot \prod_{j > i} P(S_j = s_j \mid S_i = x, \; S_k = s_k : k < j, \; k \neq i)$$

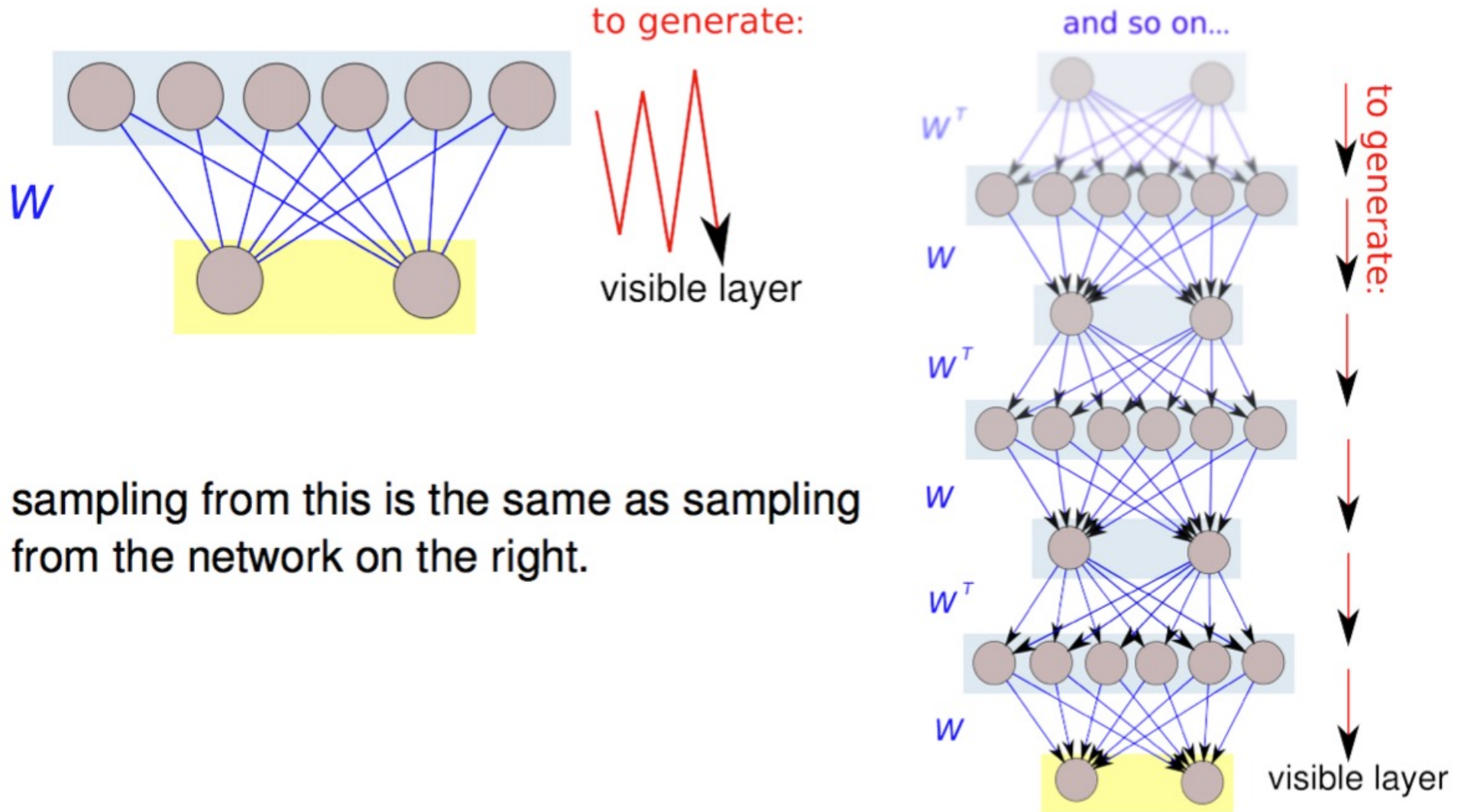$$P(S_i = s_i \mid S_j = s_j : j < i) = \sigma\left( s_i^* \sum_{j < i} s_j w_{ij} \right)$$

# RBMs are infinite belief networks (with tied weights)

Since none of the units within a layer are interconnected, we can do Gibbs sampling by updating the whole layer at a time.



(with time running from left ⟶ right)
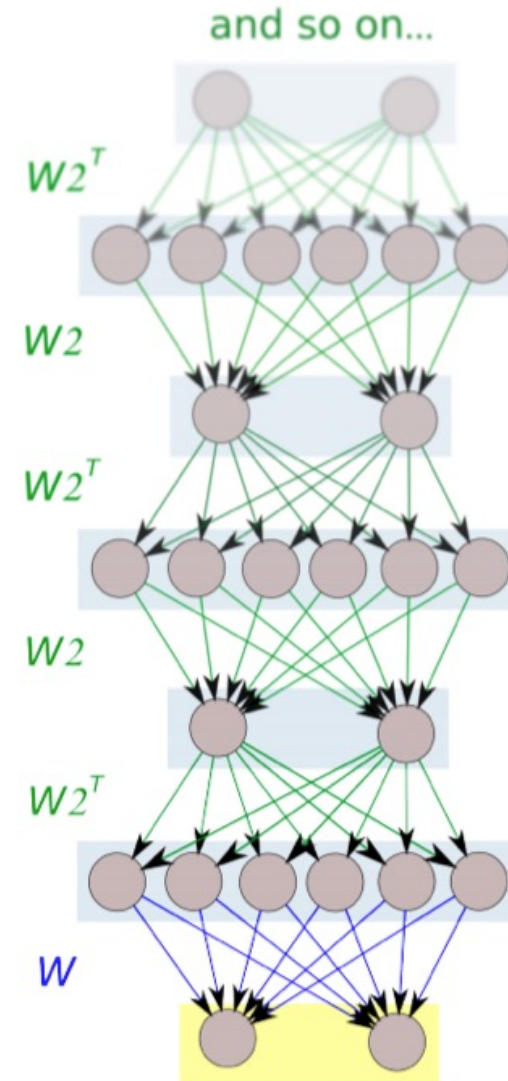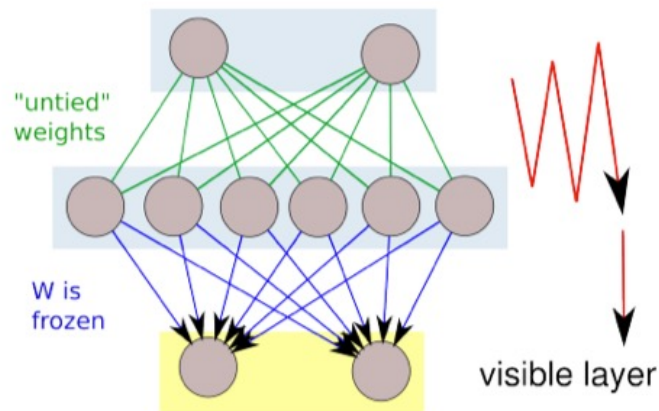
# RBMs are infinite belief networks (with tied weights)



to generate:

visible layer

sampling from this is the same as sampling from the network on the right.

$W$

and so on...

$W^T$

$W$

$W^T$

$W$

$W^T$

$W$

to generate:

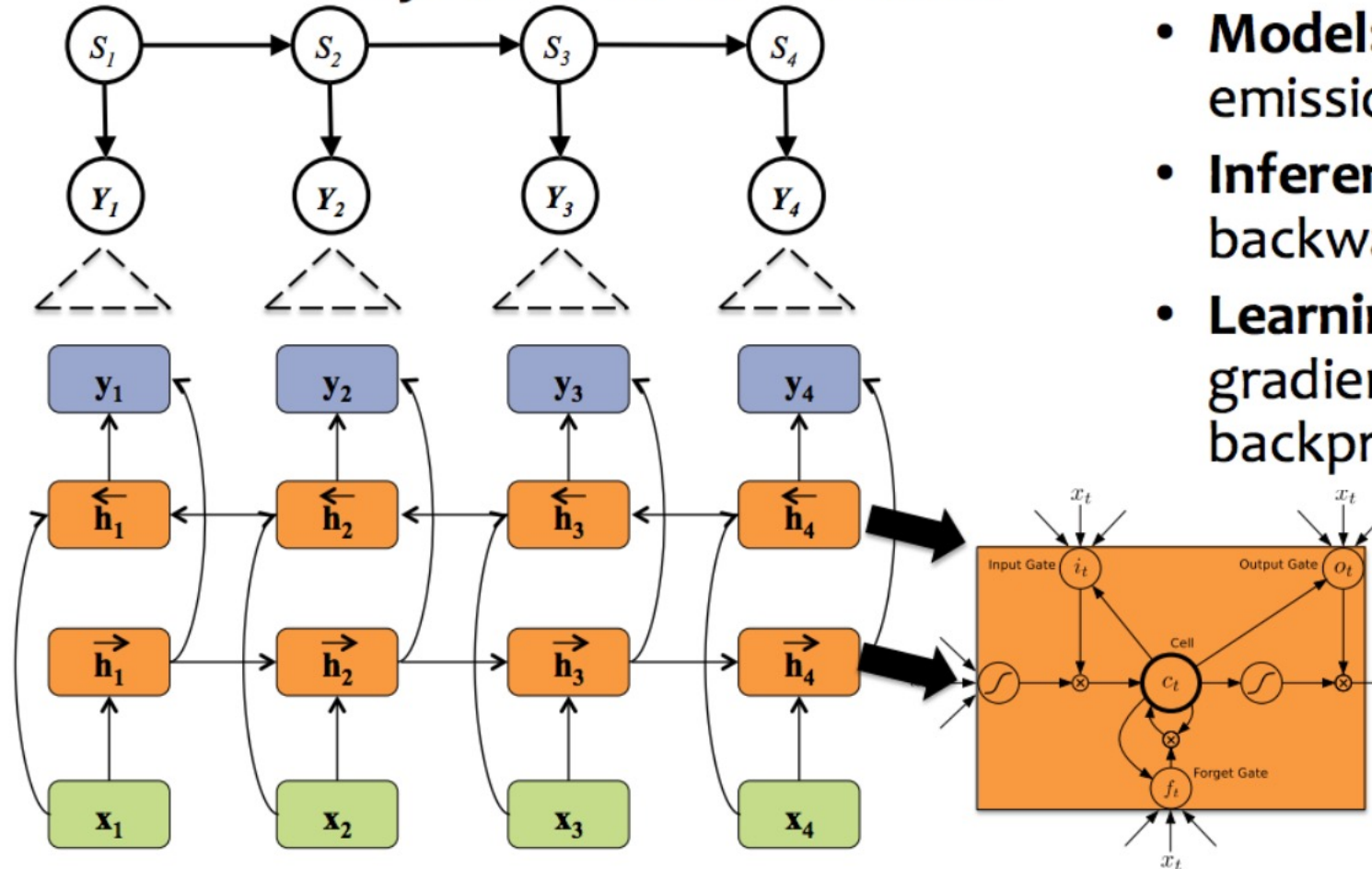visible layer

# Deep Belief networks: layer-wise pre-training

Un-tie the weights from layers 2 to infinity

If we freeze the first RBM, and then train another RBM atop it, we are untying the weights of layers 2+ in the ∞ net (which remain tied together).

and so on...
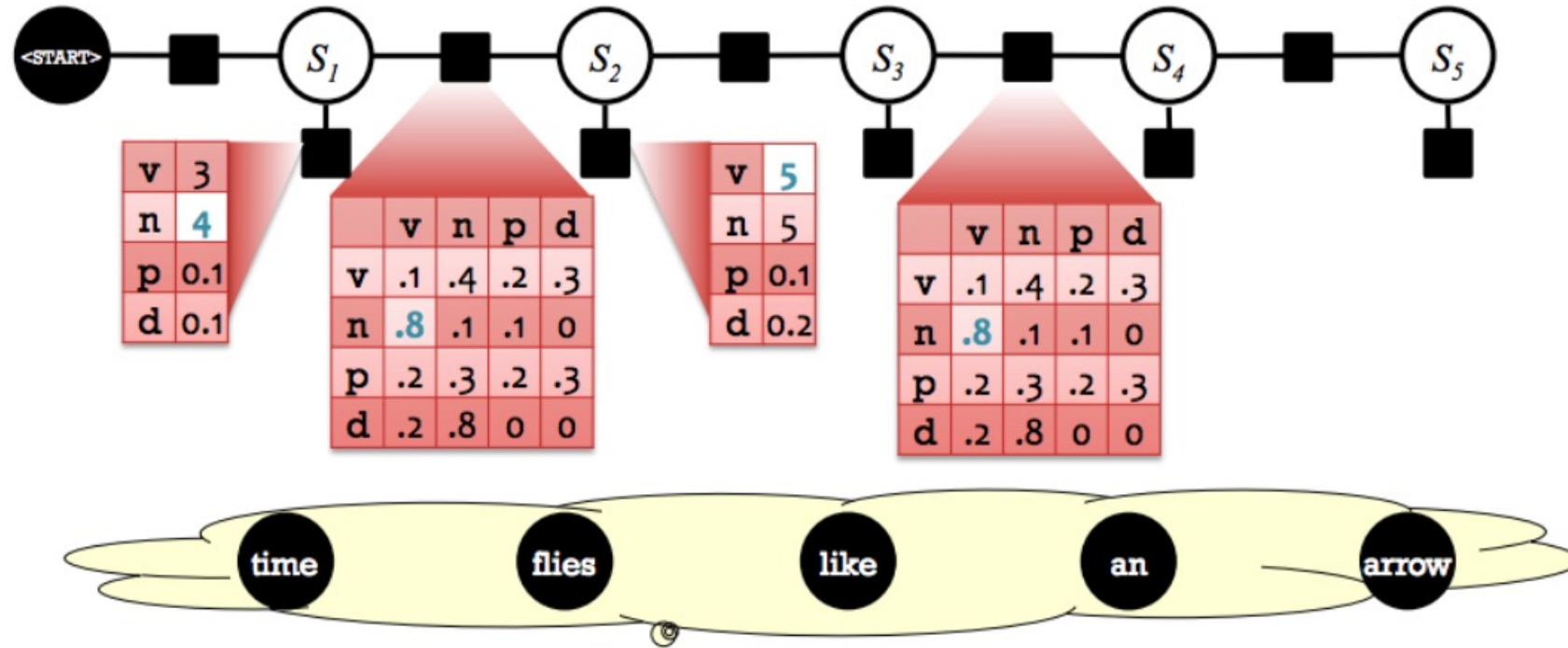
# NNs and GMs: Natural Complements



Hybrid: RNN + HMM

- **Objective:** log-likelihood
- **Model:** HMM/Gaussian emissions
- **Inference:** forward-backward algorithm
- **Learning:** SGD with gradient by backpropagation

[Graves et al. 2013]

# NNs and GMs: Natural Complements



- In a standard CRF, each of the factor cells is a parameter (e.g. transition or emission)
- In the hybrid model, these values are computed by a neural network with its own parameters

[Collobert & Weston 2011]

# Looking ahead

| | | Module 3: Modern Probabilistic AI | |
|---|---|---|---|

| 4/1 | Lecture #16 (Prof. Lengerich): **Deep Learning from a GM Perspective** [ slides | notes ] | • Goodfellow et al., Deep learning book, Ch. 6.2-5, 20.3-4 <br> • Salakhutdinov and Hinton, Deep Boltzmann Machines <br> • Ranganath et al., Deep exponential families |
|---|---|---|
| 4/3 | Lecture #17 (Prof. Lengerich): **CNNs, RNNs, Autoencoders** [ slides | notes ] | • Pascanu, Mikolov, Bengio, On the difficulty of training recurrent neural networks |
| 4/8 | Lecture #18 : **Deep Generative Models: GAN, VAEs** [ slides | notes ] | • Goodfellow et al., Deep learning book, Ch. 20.9-10 <br> • Kingma and Welling, Variational Autoencoders <br> • Goodfellow et al., Generative Adversarial Nets <br> • Arora., Generative Adversarial Networks (GANs), Some Open Questions |
| 4/10 | Lecture #19 (Prof. Lengerich): **Attention and Transformers** [ slides | notes ] | • Vasvani et al., Attention is all you need <br> • Devlin et al., BERT - Pre-training of Deep Bidirectional Transformers for Language Understanding. <br> • Raschka, Build an LLM from Scratch 3 (video) <br> • Sanderson, Visualizing transformers and attention (video) |

| 4/15 | Lecture #20 (Prof. Lengerich): **LLMs from a Probabilistic Perspective 1: Implementing a GPT from Scratch** [ slides | notes ] | • Radford et al., Improving Language Understanding by Generative Pre-Training (the GPT-1 paper) <br> • Radford et al., Language Models are Unsupervised Multitask Learners (the GPT-2 paper) <br> • Brown et al., Language Models are Few-Shot Learners (the GPT-3 paper) <br> • Raschka, Build an LLM from Scratch 4 (video) <br> • Karpathy, Let's Build GPT from Scratch (video) |
|---|---|---|
| 4/17 | Lecture #21 (Prof. Lengerich): **LLMs from a Probabilistic Perspective 2: Training on Unlabeled Data** [ slides | notes ] | • Bi. et al, DeepSeek LLM Scaling Open-Source Language Models with Longtermism <br> • Liu et al., DeepSeekV2 A Strong, Economical, and Efficient Mixture-of-Experts Language Model <br> • Liu et al., DeepSeekV3 Technical Report |
| 4/22 | Lecture #22 (Prof. Lengerich): **LLMs from a Probabilistic Perspective 3: Fine-tuning on Labeled Data** [ slides | notes ] | • Raffel et al., Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer <br> • Ouyang et al., Training language models to follow instructions with human feedback <br> • Li & Liang, Prefix-tuning - Optimizing continuous prompts for generation |
| 4/24 | Lecture #23 (Prof. Lengerich): **Context-Adaptive Graphical Models** [ slides | notes ] | • Lengerich et al., Contextualized Machine Learning |
| 4/29 | | Project Presentations |
| 5/1 | | Project Presentations |

Questions?