# How do artificial neural networks learn to predict protein structures?

Badri Adhikari

adhikarib @ umsl.edu

Assistant Professor of CS
Department of Mathematics & Computer Science
University of Missouri-St. Louis

**PRINCIPAL INVESTIGATOR**

Badri Adhikari, PhD
Assistant Professor of Computer Science
Department of Mathematics and Computer Science
University of Missouri-St. Louis

312 Express Scripts Hall
St. Louis, MO 63121

Phone: 314-516-7393
Email: adhikarib@umsl.edu
Homepage: https://badriadhikari.github.io/

**TEACHING**

- Artificial Intelligence - 2018 Fall, 2019 Spring, 2019 Fall
- Deep Learning - 2019 Spring
- Programming and Data Structures - 2018 Spring
- Hands-on Deep Learning Workshops - 2018 Fall
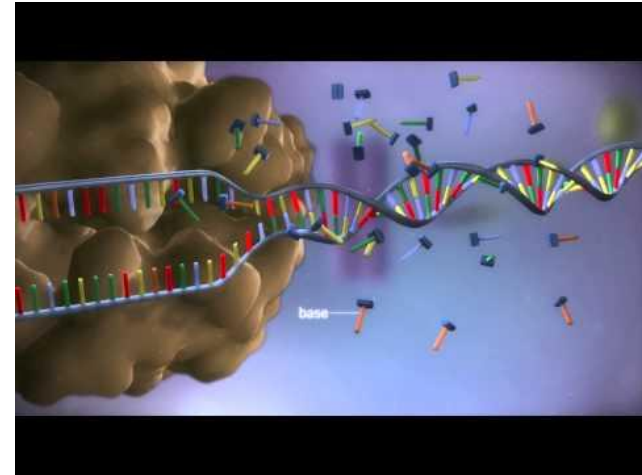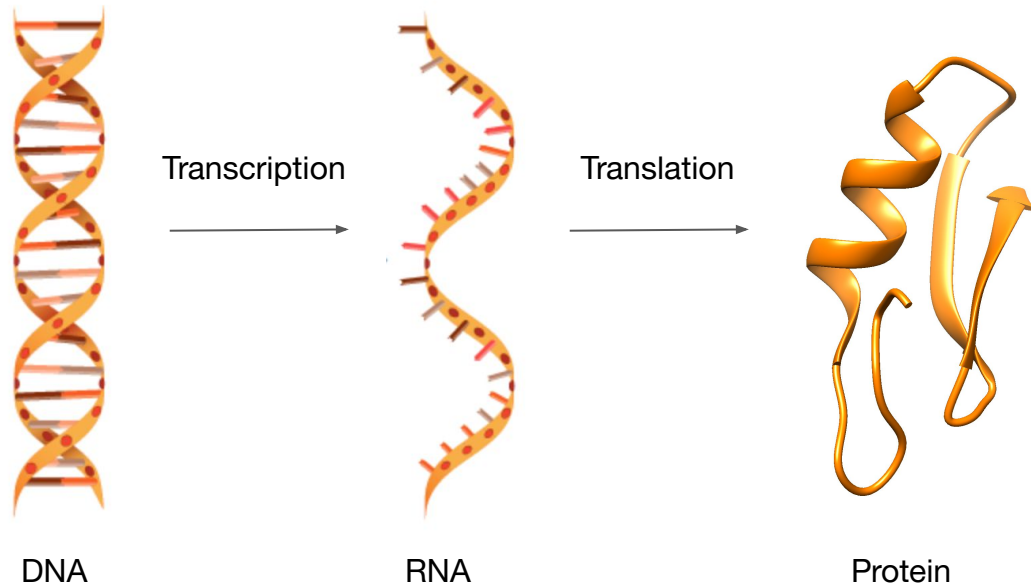- Advanced Data Structures and Algorithms - 2017 Fall

**RESEARCH**

- Investigate deep learning methods for protein structure prediction
- Collaboration with Robert Paul at MIMH (various mental health datasets)
- Collaboration with Lauren Salmimen at USC (UK biobank datasets, and various other datasets)
- Collaboration with researchers at International Maize and Wheat Development Center (CIMMYT)

How do artificial neural networks learn to predict protein structures?

# The central dogma in molecular biology

- The fundamental process in life is the flow of information from DNA to proteins
  - How does this happen?



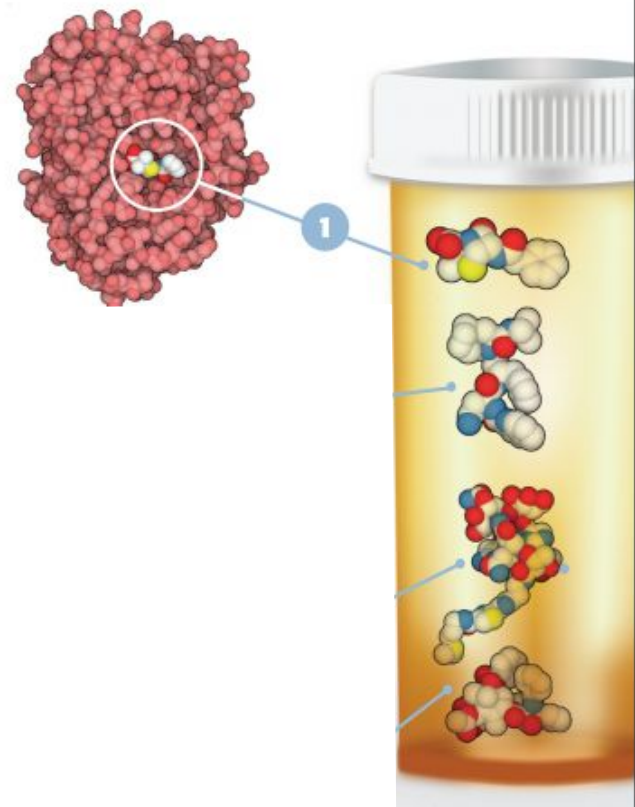Transcription

Translation

DNA

RNA

Protein

# Protein folding is important

- Proteins are fundamental to understanding their role within the body

- Many diseases are believed to be caused by misfolded proteins
  - Alzheimer's, Parkinson's, Huntington's, cystic fibrosis, etc.

- One of the top 100 questions selected by the Science magazine

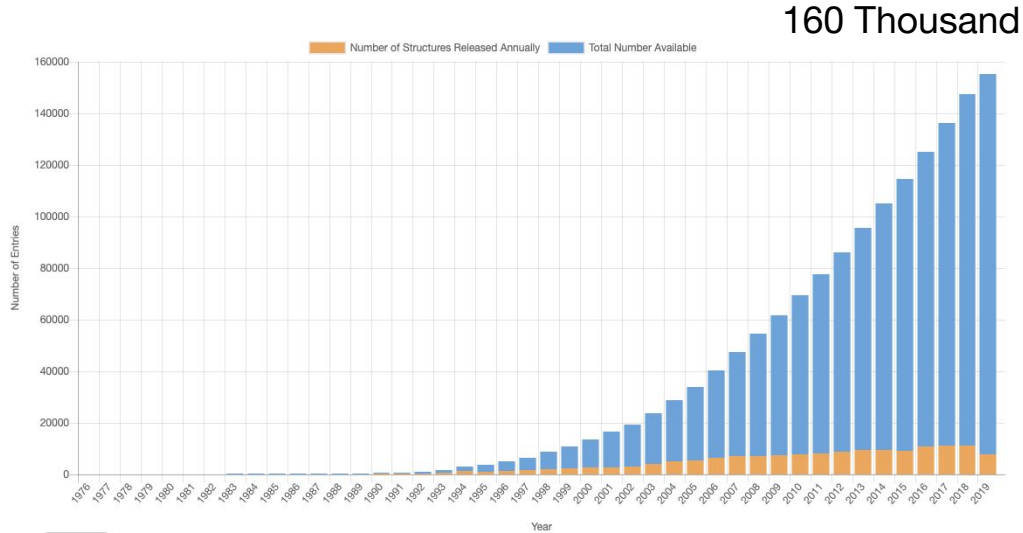# The need to obtain precise 3D structures

- Antibiotics need to kill pathogenic organisms like bacteria without poisoning the patient

- Often, these drugs attack proteins that are only found in the targeted bacterium which are crucial for their survival or multiplication

- For instance, penicillin attacks the enzyme that builds bacterial cell walls



https://cdn.rcsb.org/pdb101/learn/resources/how-do-drugs-work-flyer.pdf

# Only around 160K protein structures are solved so far



PDB Statistics: Overall Growth of Released Structures Per Year

160 Thousand

| | Total |
|---|---|
| UniRef100 | 199,397,329 |
| UniRef90 | 99,657,864 |
| UniRef50 | 37,541,209 |

200 Million

UniRef clusters per taxonomic group

# Can we predict protein structures accurately today?

http://predictioncenter.org/
Critical Assessment of Protein Structure Prediction



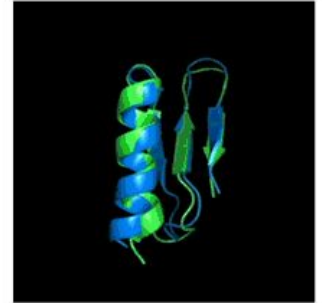World-wide competition
held every two years
(3 months long)



vs

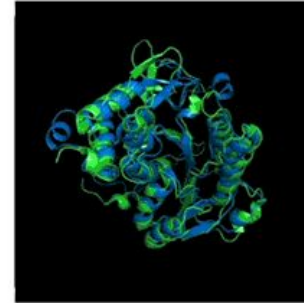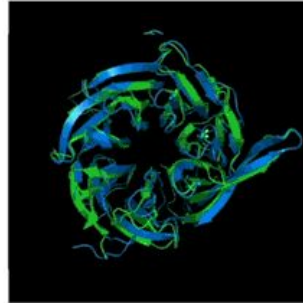Models predicted by DeepMind in CASP 13 (2018)



Structures:
Ground truth (green)
Predicted (blue)

T0954 / 6CVZ        T0965 / 6D2V        T0955 / 5W9F
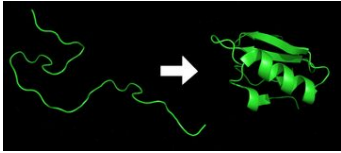
https://deepmind.com/blog/article/alphafold

# What does AI do exactly?

# This is how computer scientists predict protein structures

Anfinsen's Dogma (1973)
Native structure is determined **only**
by the protein's amino acid sequence
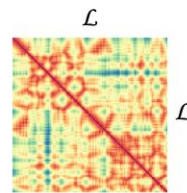
Experimentally determined
3D structure of a Protein

Sequence
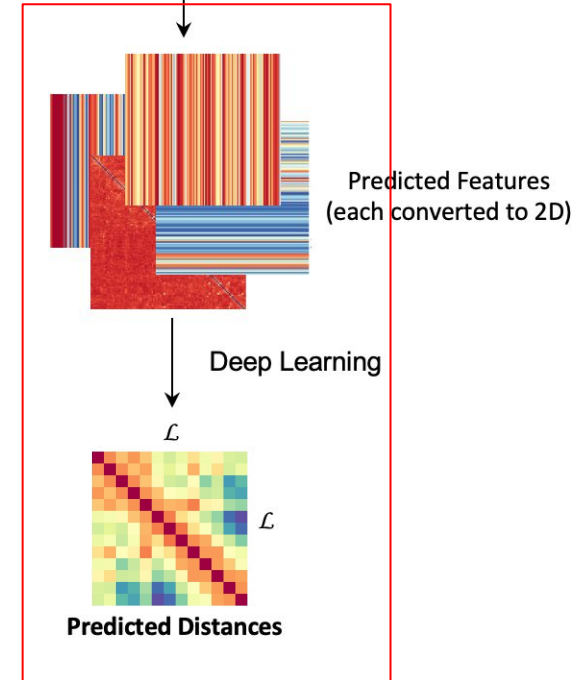
Protein Amino Acid Sequence of Length $\mathcal{L}$

MKTFLYFCLLFIVQTAFAADSIY...VREQ

Calculate pairwise distances
between carbon-β atoms

$\mathcal{L}$

$\mathcal{L}$

**True Distance Matrix**

Predicted Features
(each converted to 2D)

Deep Learning
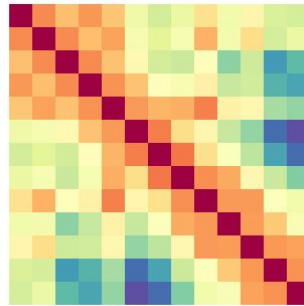
$\mathcal{L}$

$\mathcal{L}$

**Predicted Distances**

# The distance prediction problem is an AI problem



A Protein's Amino Acid Sequence

G F G C N G P W ... D E D D M

?

Prediction

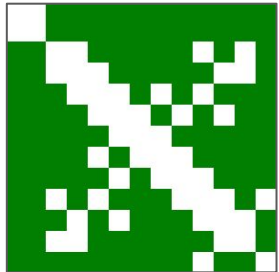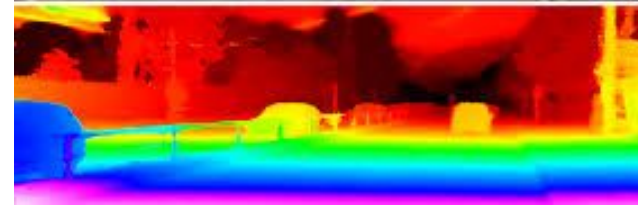| Pair | D |
|------|-----|
| $G_1 - F_2$ | 3.5 |
| $G_1 - G_3$ | 6.2 |
| $G_1 - C_4$ | 8.1 |
| ... | |
| $G_2 - E_9$ | 8.0 |
| ... | |

Distance Map

Threshold

Contact Map

Input

Output

**Depth Estimation** / Image Segmentation

# Attacking the protein distance prediction problem



Multiple Sequence Alignment

Input

GFGCNGPWDEDDM

1D and 2D Predictions

Sequence profiles, Sec struc., etc..

Input Feature Matrix

Deep Learning

Predicted Contact/Distance Map

3D Model

True Structure

Ground Truth for DL

1 Generating Large and **High-quality MSAs**

2 The Right Deep Learning **Architecture**

DL

3 **Feature Engineering** (encoding, 1D to 2D, covariance/precision matrix)

4 Appropriate **Loss Function** (besides binary_crossentropy, MSE, etc.)

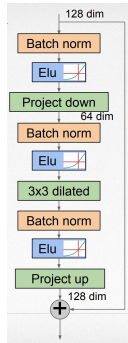# All top methods predict distances using residual neural networks

🏆

Top Methods in the most recent CASP Competition



De novo protein folding using statistical potentials from deep learning

R.Evans, J.Jumper, J.Kirkpatrick, L.Sifre, T.F.G.Green, C.Qin, A.Zidek, A.Nelson, A.Bridgland, H.Penedones, S.Petersen, K.Simonyan, D.T.Jones [UCL], K.Kavukcuoglu, D.Hassabis, A.W.Senior
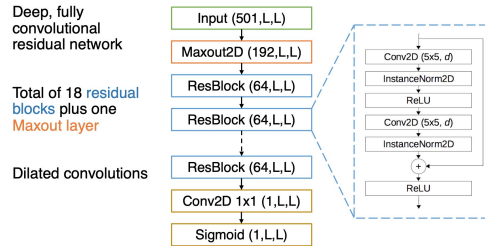
DeepMind
Group 043 / A7D / AlphaFold

**DeepMetaPSICOV (DMP) in CASP13**

Shaun M Kandathil
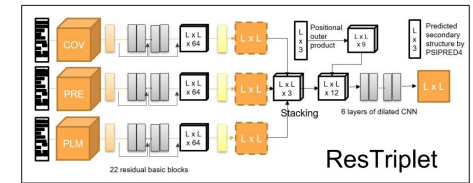University College London
&
The Francis Crick Institute

DeepMetaPSICOV model architecture

**ResTriplet/TripletRes:**
**Learning contact-maps from a triplet of coevolutionary matrices**

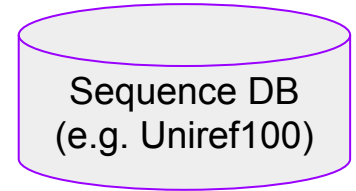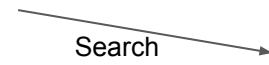Eric W. Bell, Yang Li, Chengxin Zhang, Dong-Jun Yu, Yang Zhang

Department of Computational Medicine and Bioinformatics, University of Michigan - Ann Arbor
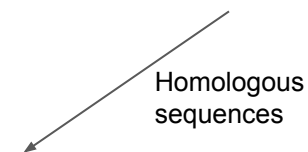
# All these results show that residual networks are best architectures (for this problem)

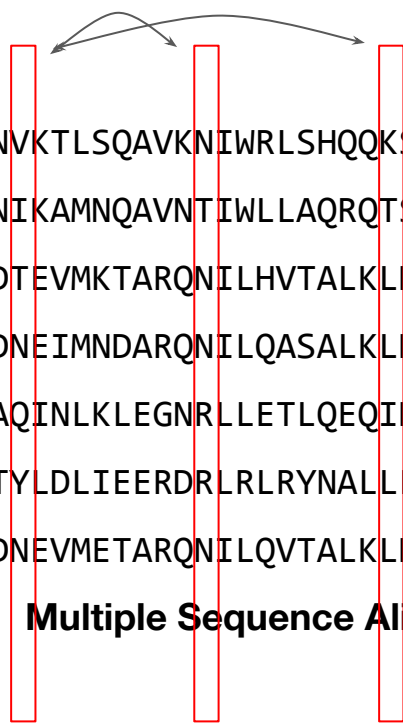# Can We Learn to Predict Contacts WITHOUT 'True' Contacts?

MSEIITFPQQTVVYPEINVKTLSQAVKNIWRLSHQQKSGIEIIQEKTLRISLYSRDLDEA

Search → Sequence DB (e.g. Uniref100)

Homologous sequences

MSEIITFPQQTVVYPEINVKTLSQAVKNIWRLSHQQKSGIEIIQEKTLRISLYSRDLDEA

----NTLSQKENMYPEINIKAMNQAVNTIWLLAQRQTSGIEIINDKVKRISLYSREFDE-

------------LTPPDTEVMKTARQNILHVTALKLDFLPVMKEKMRPLQDALISADK-

-----------ILTPPDNEIMNDARQNILQASALKLDFLPVMKEKMLPLQTALKRADKV

MVVRNSAKAIAEHSDDMAQINLKLEGNRLLETLQEQIDSITLRSAALESTMGEITA----

AGIARLGKLLDKVSSALTYLDLIEERDRLRLRYNALLEESRTAHQEEKATAAKLDELT--

------------LTPPDNEVMETARQNILQVTALKLDFLPVMKEKMLPLQAALMSADKV
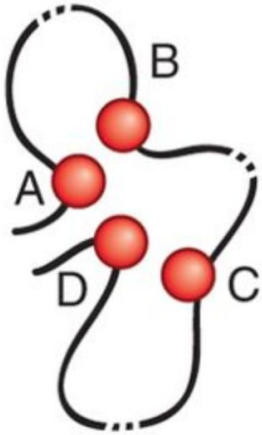
Covariance / Coevolution

**Multiple Sequence Alignment**

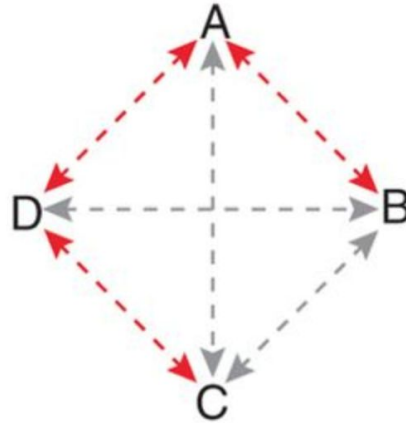**High covariance implies physical closeness!**

Does this mean we can write an algorithm to predict contacts?

# Can We Learn to Predict Contacts WITHOUT 'True' Contacts?



Physical contacts

Observed correlations

Predicted contacts

■ Causative   ■ Transitive
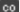
Perspective | Published: 08 November 2012

Protein structure prediction from sequence variation

Debora S Marks ✉, Thomas A Hopf & Chris Sander ✉

*Nature Biotechnology* **30**, 1072–1080 (2012) | Download Citation ⤓

# Can We Write Algorithms to Remove Transitive Noise?

## Protein 3D Structure Computed from Evolutionary Sequence Variation

Debora S. Marks co ✉, Lucy J. Colwell co, Robert Sheridan, Thomas A. Hopf, Andrea Pagnani, Riccardo Zecchina, Chris Sander

## FreeContact: fast and free software for protein contact prediction from residue co-evolution

László Kaján, Thomas A Hopf, Matúš Kalaš, Debora S Marks and Burkhard Rost ✉

## PSICOV: precise structural contact prediction using sparse inverse covariance estimation on large multiple sequence alignments FREE

David T. Jones ✉, Daniel W. A. Buchan, Domenico Cozzetto, Massimiliano Pontil

Author Notes

## CCMpred—fast and precise prediction of protein residue−residue contacts from correlated mutations 🔓

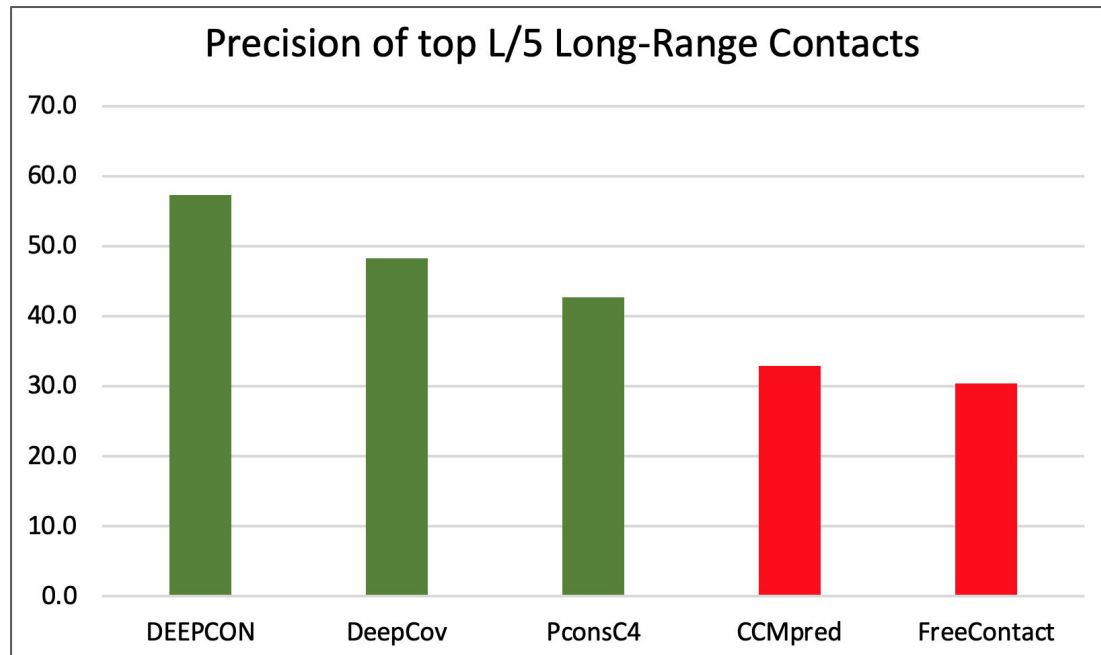Stefan Seemayer, Markus Gruber, Johannes Söding ✉     Author Notes

# Can Deep Learning Remove Transitive Noise?



Precision of top L/5 Long-Range Contacts

DEEPCON: protein contact prediction using dilated convolutional neural networks with dropout

Badri Adhikari ✉

*Bioinformatics*, btz593, https://doi.org/10.1093/bioinformatics/btz593
**Published:** 29 July 2019    **Article history ▾**

# What does AI do exactly?

It learns from input/output pairs and can outperform an algorithm based on theories!

# Neural plasticity

In 1953, Professor Theodor Erismann devised an experiment
- performing it upon his assistant and student, Ivo Kohler

He made Kohler wear a pair of hand-engineered goggles
- Specially arranged mirrors flipped the light that would reach eyes, top becoming bottom, and bottom top.



After 10 days, Kohler had grown accustomed to the invariably upside-down world
- everything seemed to him normal, rightside-up
- He could do everyday activities in public perfectly well: walk along a crowded sidewalk, even ride a bicycle

https://www.theguardian.com/education/2012/nov/12/improbable-research-seeing-upside-down

# Neural plasticity

**SHARE**

REPORTS

## Eye-specific termination bands in tecta of three-eyed frogs
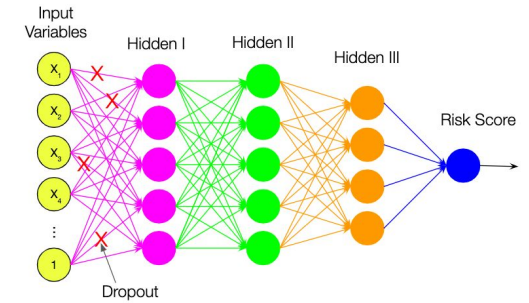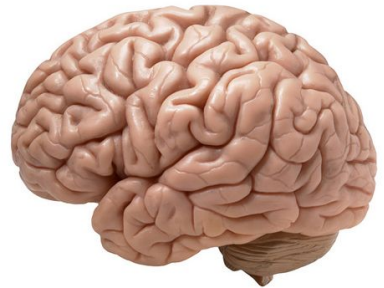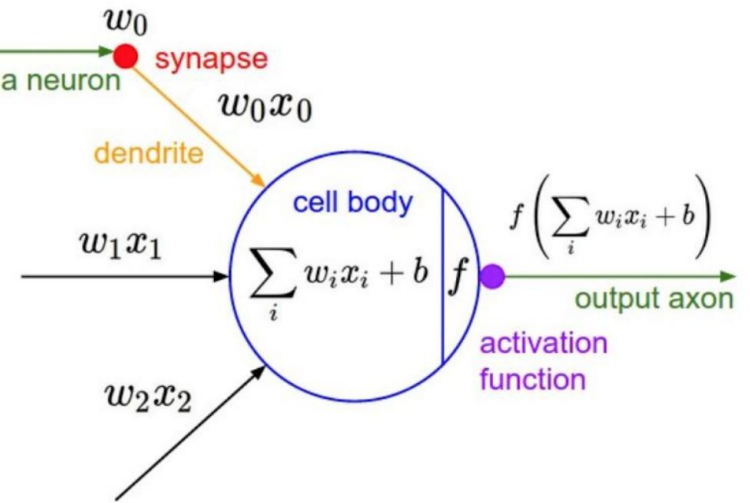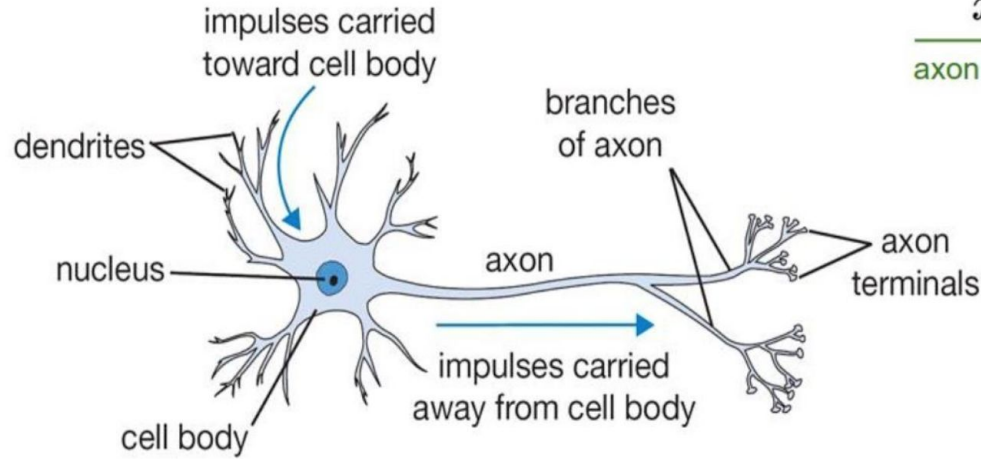
M Constantine-Paton, MI Law
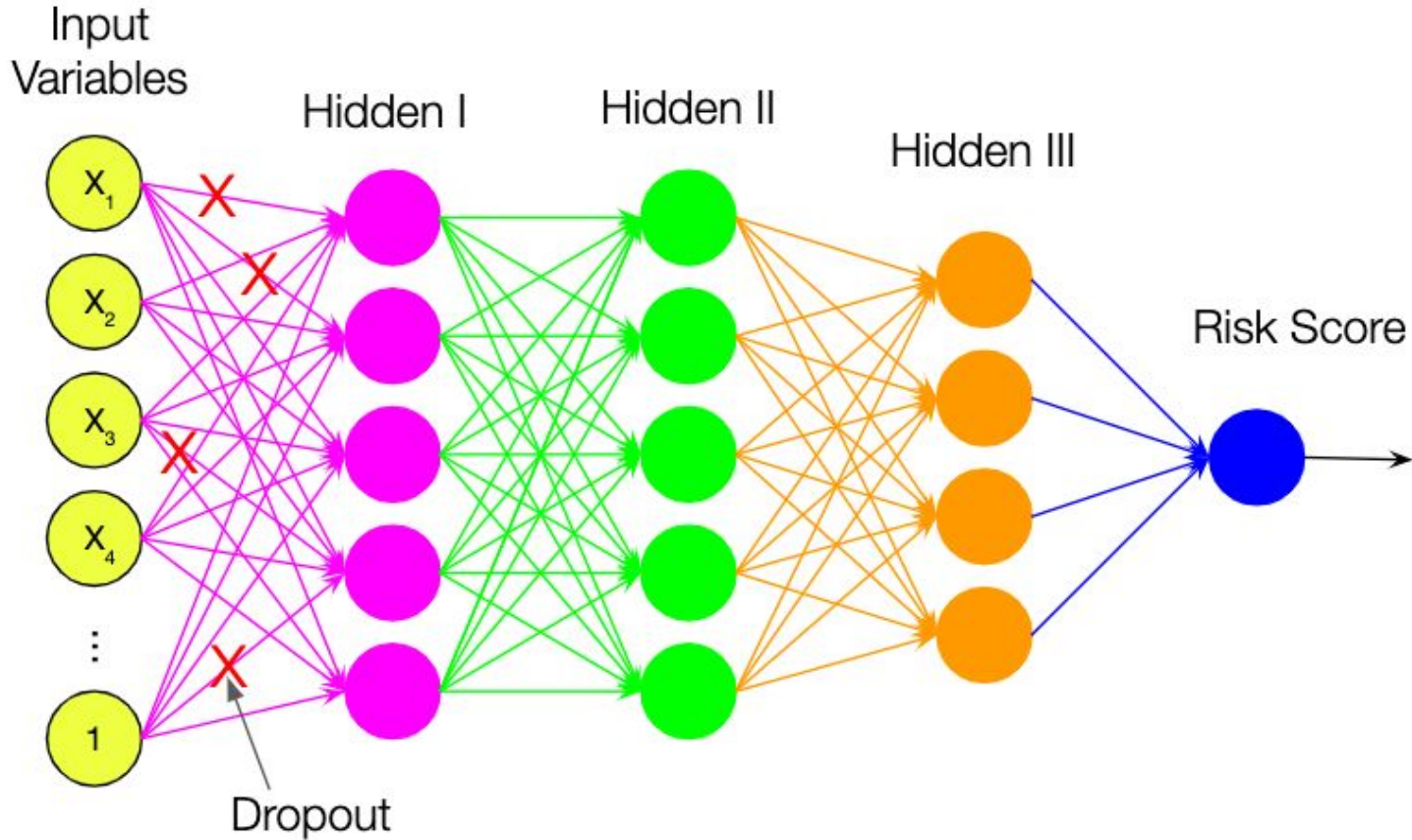+ See all authors and affiliations

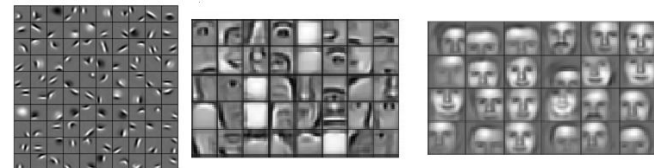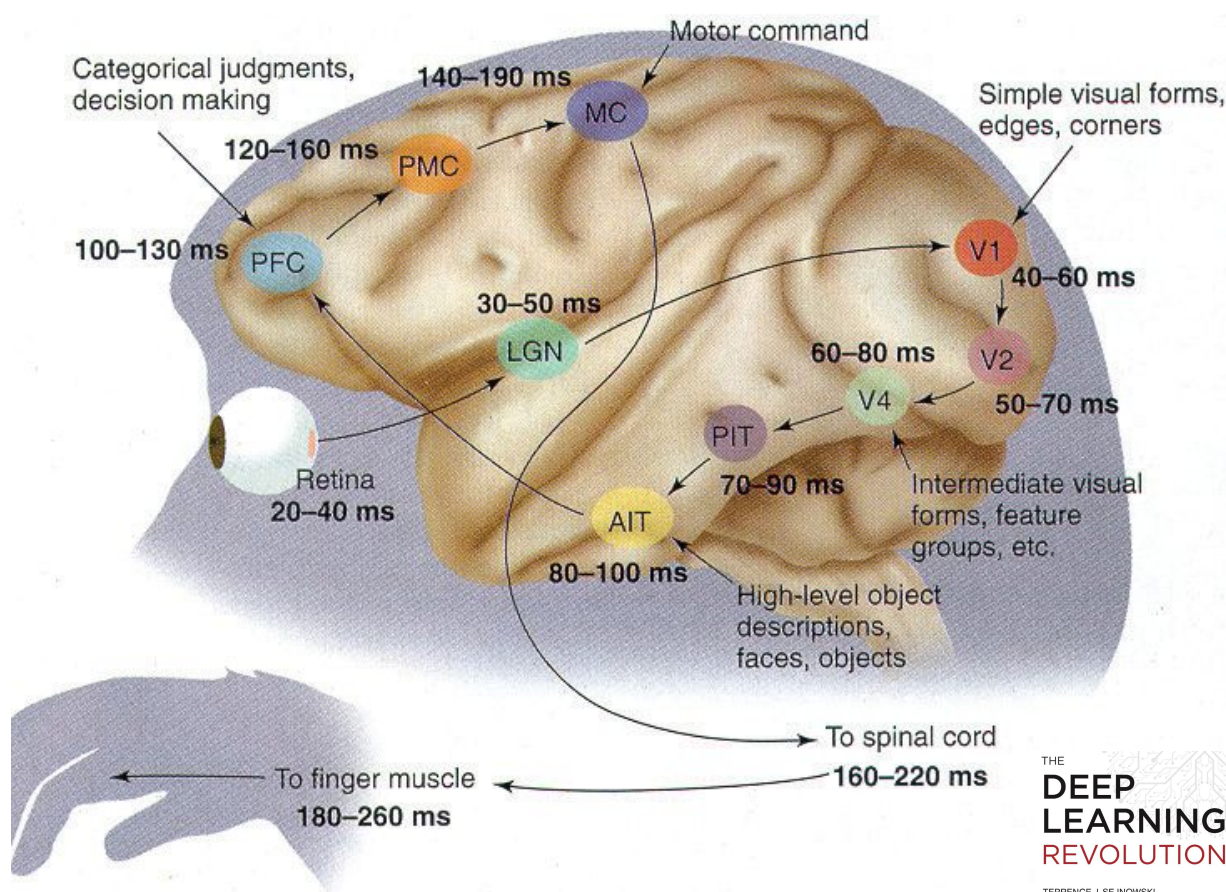Dr. Martha Constantine-Paton is a neuroscientist at MIT

# Biological vs. artificial neurons

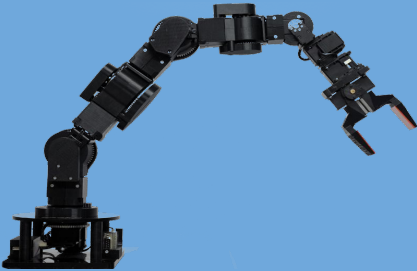# Feed-forward neural networks are very successful

# Human visual cortex is hierarchical



LGN: lateral geniculate nucleus
V1: primary visual cortex
V2: secondary visual cortex
V4: visual area 4
AIT and PIT: anterior and posterior inferotemporal cortex
PFC: prefrontal cortex
PMC: premotor cortex
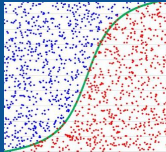MC: motor cortex

# Artificial Intelligence vs. Machine Learning vs. Deep Learning
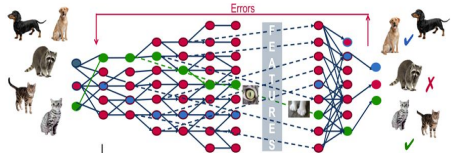


ARTIFICIAL INTELLIGENCE
a very broad field

MACHINE LEARNING
Fundamentals of "learning from data"

Deep Learning
Powerful trending ML methods

1950s

1980s

2010s

# The state of computer vision and AI: we are really, really far away



http://karpathy.github.io/2012/10/22/state-of-computer-vision/

**Some things "we" understand easily**

There are **3 mirrors** in the scene so some of those people are **"fake" replicas** from different viewpoints

**Recognize Obama** from the few pixels that make up his face

You recognize that there's a person **standing on a scale**, even though the scale occupies only very few white pixels that blend with the background

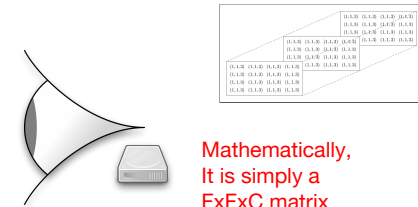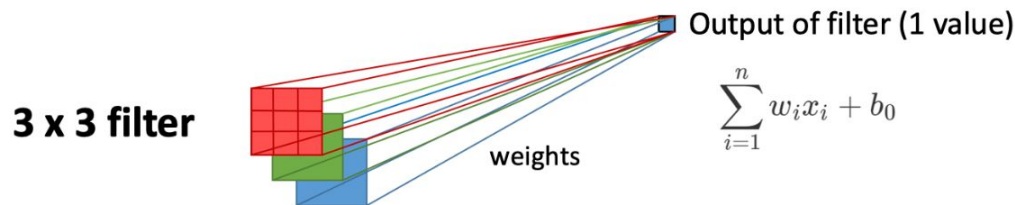Obama has his **foot positioned just slightly** on top of the scale
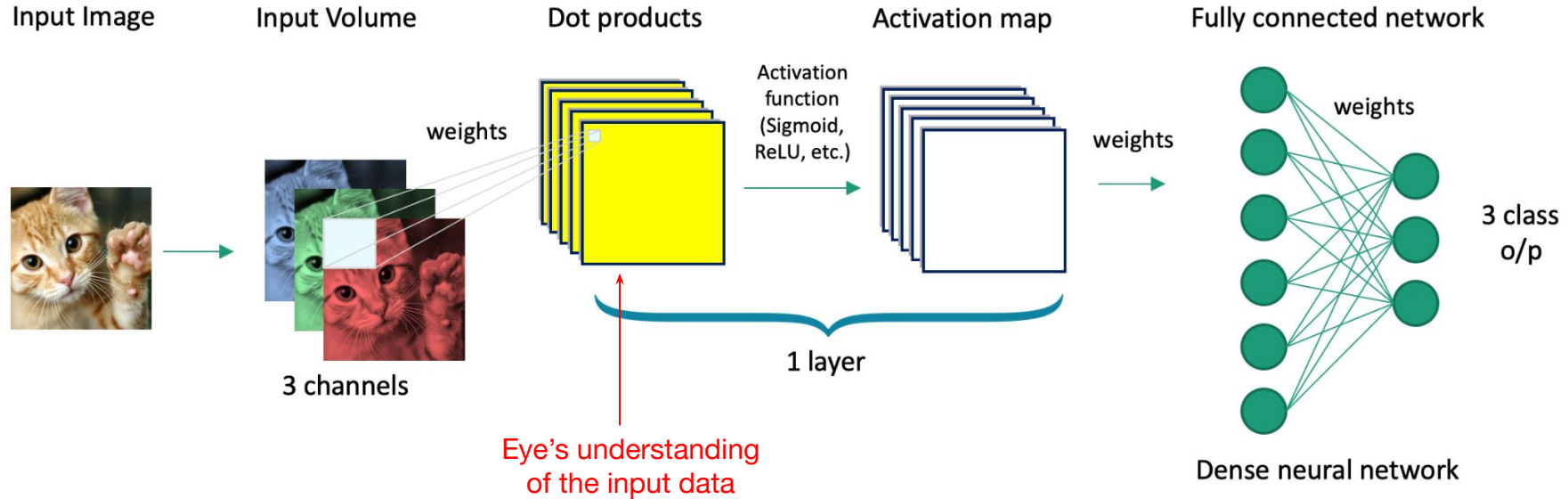
Working physics - **Obama is leaning in on the scale**, which applies a force on it. Scale measures force that is applied on it, that's how it works => it will over-estimate the weight of the person standing on it.

The **person** measuring his weight **is not aware** of Obama doing this

There are **people in the back** who **find** the person's imminent confusion **funny**
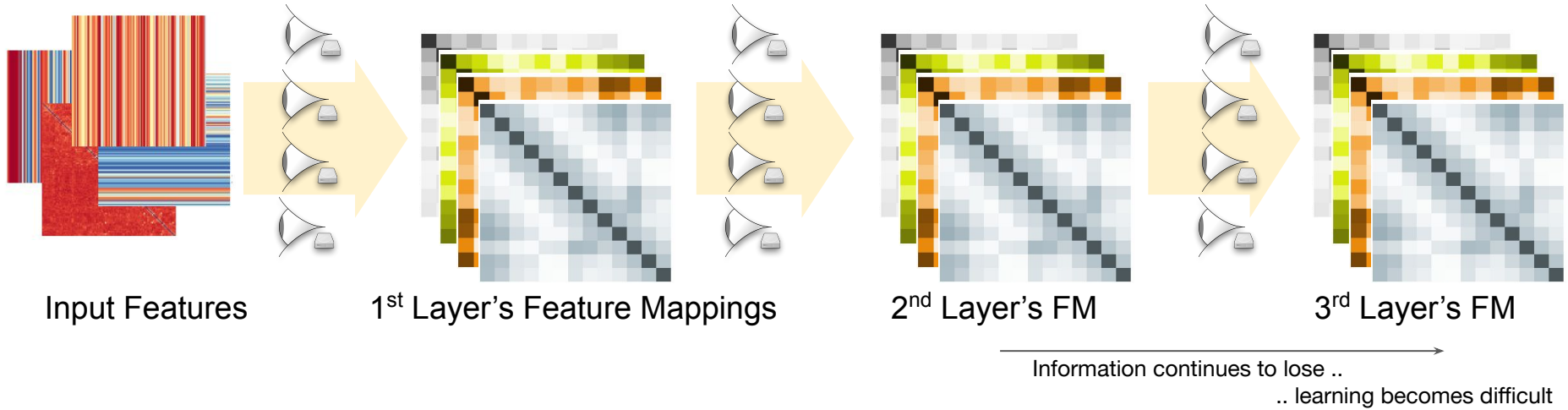
# Residual networks, used for distance prediction, are ConvNets



Input Image

Input Volume

3 channels

weights

Dot products

Activation function (Sigmoid, ReLU, etc.)

Activation map

weights

Fully connected network

weights

3 class o/p

Dense neural network

1 layer

Eye's understanding of the input data

**3 x 3 filter**

Output of filter (1 value)

weights

$$\sum_{i=1}^{n} w_i x_i + b_0$$

Mathematically, It is simply a FxFxC matrix

**Convolutional neurons** are like our **"eyes with memory"**..

# Residual networks



Input Features      1ˢᵗ Layer's Feature Mappings      2ⁿᵈ Layer's FM      3ʳᵈ Layer's FM

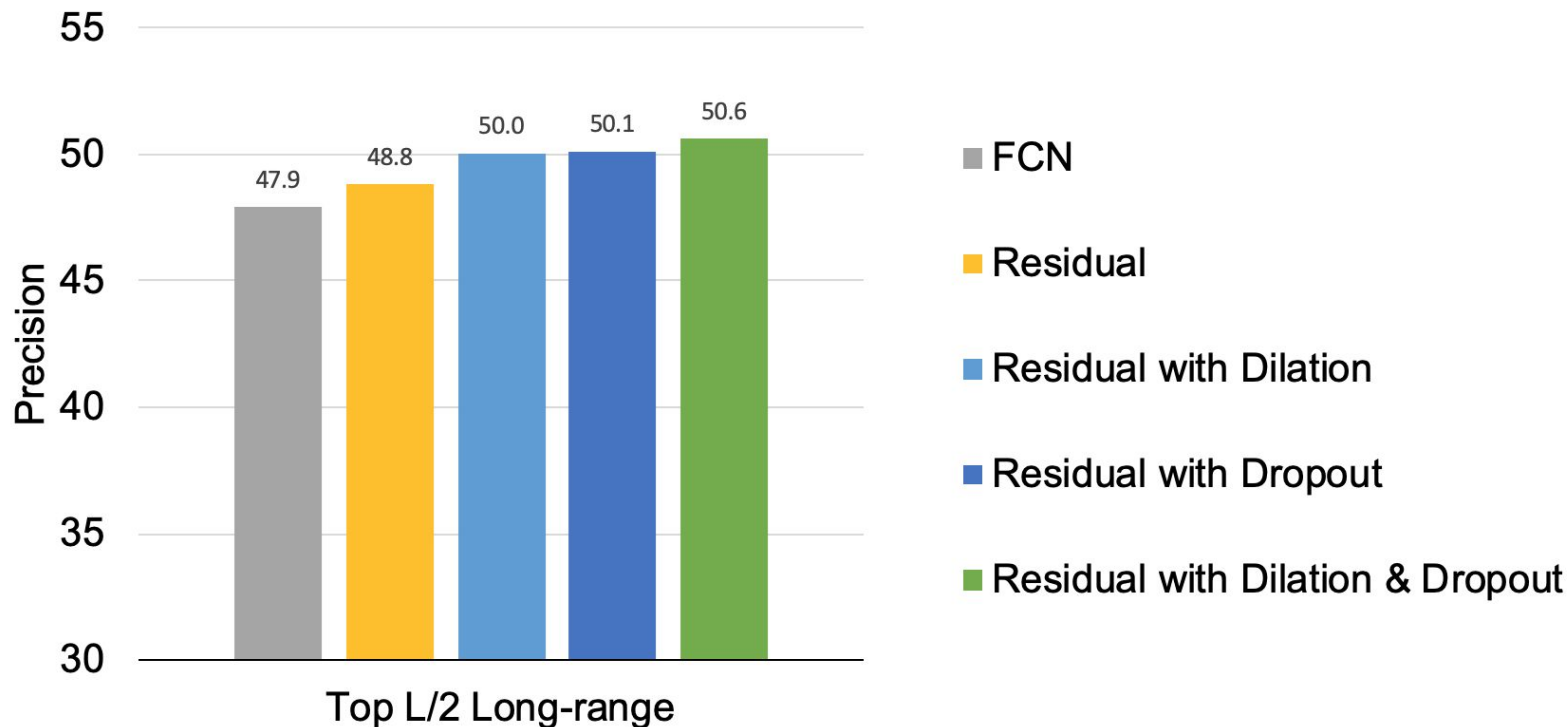Information continues to lose ..

.. learning becomes difficult

# Residual networks

# Variants of residual networks can perform even better



DEEPCON: protein contact prediction using dilated convolutional neural networks with dropout

Badri Adhikari ✉

*Bioinformatics*, btz593, https://doi.org/10.1093/bioinformatics/btz593
**Published:** 29 July 2019    **Article history** ▾

# Hardware for protein distance prediction

Most current deep learning experiments are performed on sample datasets:

- 3 K representative proteins [200 GB]

- Special SSDs known as M2s that directly attach to the motherboard

- Powerful GPUs such as V100 and P6000 are required

- One experiment (training) takes about 24 hours

Full dataset:

- 50 K proteins [around 10 TB]

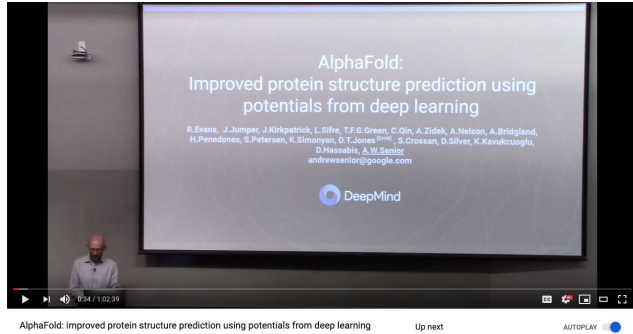- One experiment can take up to 10 days

Feature generation:

- 1000s of CPU time for a few days

# Present and Future Research in Protein Folding
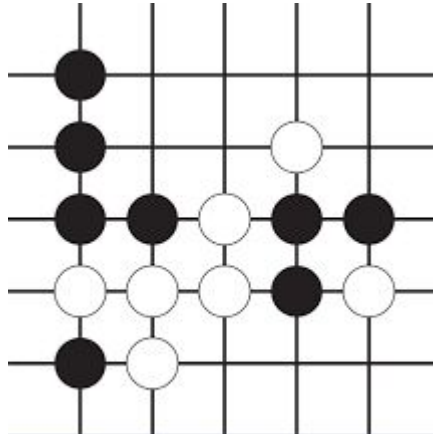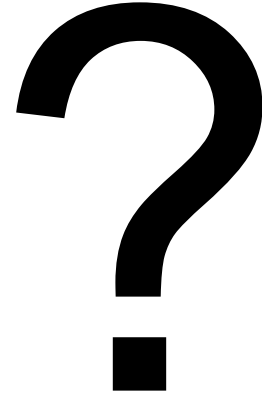
# Increased interest from industry

# Reinforcement learning.. a big hope..



Chess



Game of Go

# A lot of data + many powerful algorithms = much work to do



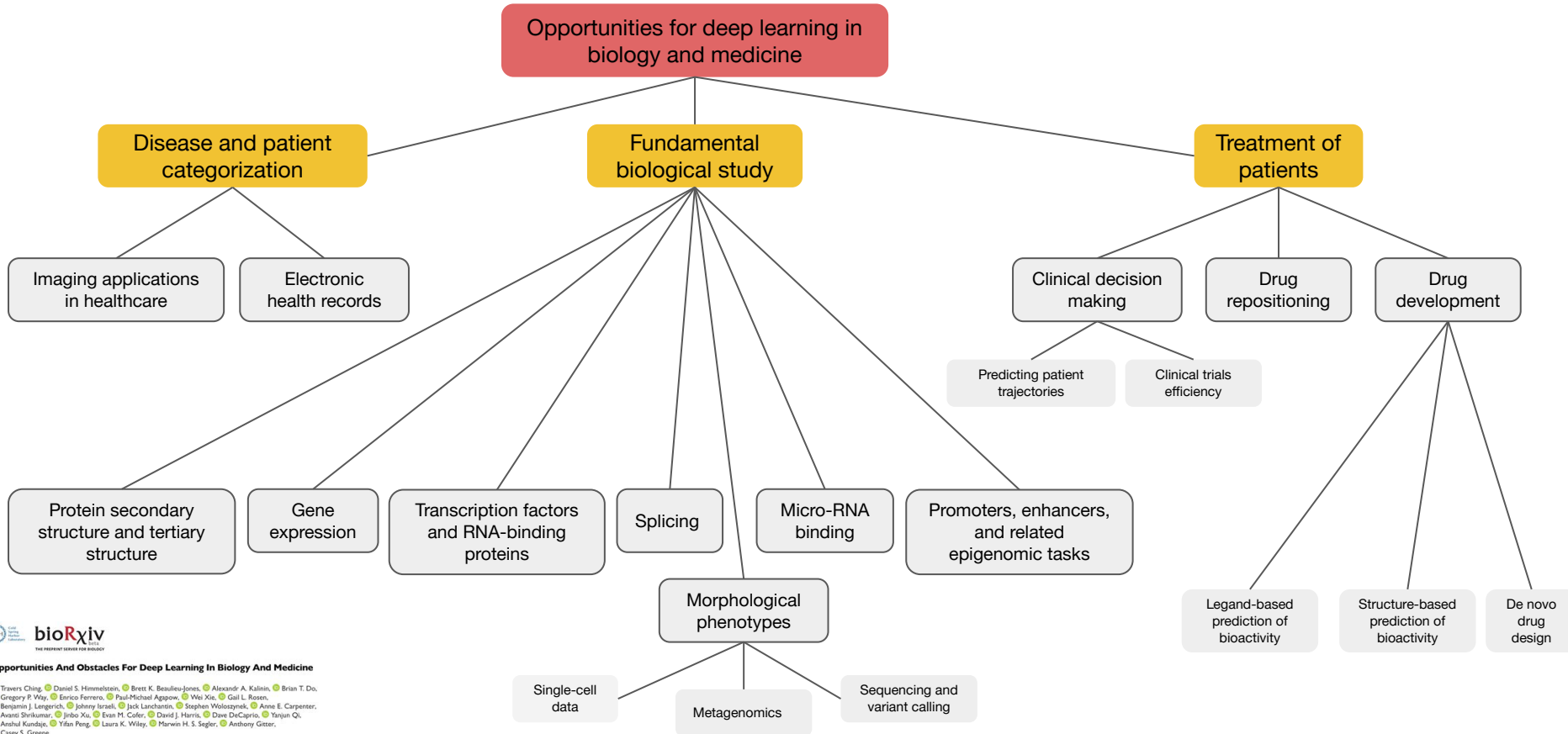| Model | Total |
|---|---|
| UniRef100 | 199,397,329 |
| UniRef90 | 99,657,864 |
| UniRef50 | 37,541,209 |

+ Deep Learning +

Multiple Sequence Alignment

Input
GFGCNGPWDEDDM

1D and 2D Predictions

Sequence profiles, Sec struc., etc..

Input Feature Matrix

Deep Learning

Predicted Contact/Distance Map

3D Model

True Structure

Ground Truth for DL

DL
1 Generating Large and **High-quality MSAs**
2 The Right Deep Learning **Architecture**
3 **Feature Engineering** (encoding, 1D to 2D, covariance/precision matrix)
4 Appropriate **Loss Function** (besides binary_crossentropy, MSE, etc.)

A unique problem

# Deep learning for biology and medicine

# Conclusions

1) Deep learning methods are full of promise but also have a lot of limitations

2) A key component of the protein folding problem, distance prediction, is largely a deep learning problem

3) Solving the problem of protein folding requires expertise from both domains - deep learning and bioinformatics

4) Protein folding problem will potentially unravel the limitations of AI and DL

# Acknowledgements

## Research Support & Contribution



Cezary Janikow



Sharlee Climer



Cynthia Jobe



Anthony Ackah-Nyanzu



Patrick Kong



Sri Harsha Akurathi

## IT Support



Philip Reiss



Kenneth Voss



Michael Remier



MU - Research Computing Support Services (RCSS)

## Computing Resources

THANK YOU