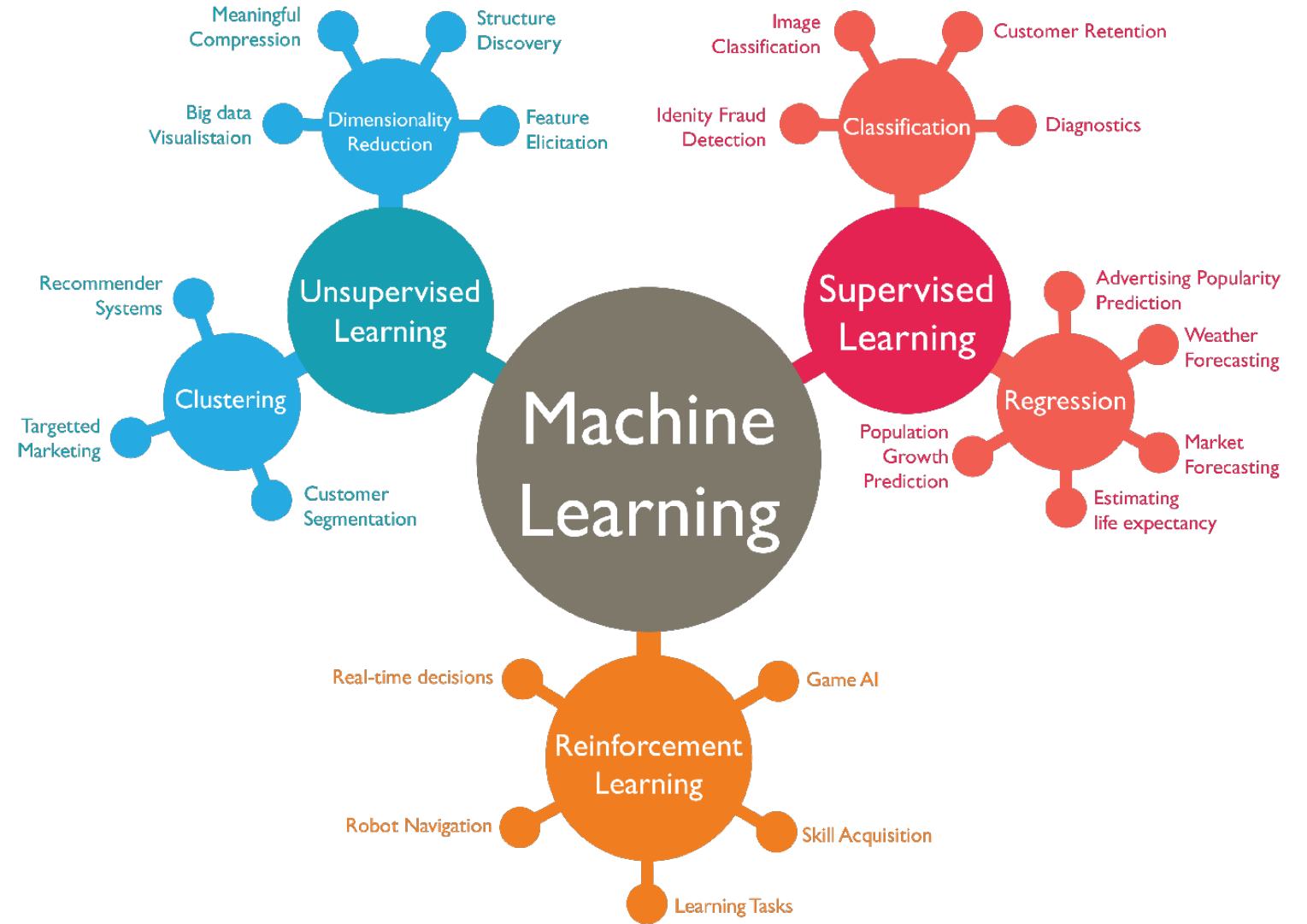


Matthew Fricke

Machine Learning at CARC



Other things I do... autonomous robotics



Melanie Moses BCLab

UNM Newsroom

About Us News Topics Schools and Colleges For News Media

UNM Newsroom / News / UNM's VolCAN team makes history in Canary Islands



UNM's VolCAN team makes history in Canary Islands

Researchers use drones to sample gases from active La Palma volcano

By Kim Deller @ January 05, 2023

Related



SCIENCE ADVANCES | RESEARCH ARTICLE

APPLIED SCIENCES AND ENGINEERING

Aerial strategies advance volcanic gas measurements at inaccessible, strongly degassing volcanoes

Other things I do... autonomous robotics

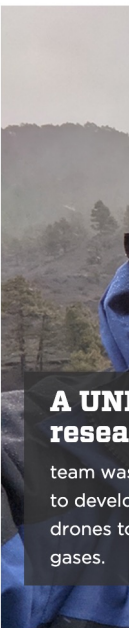
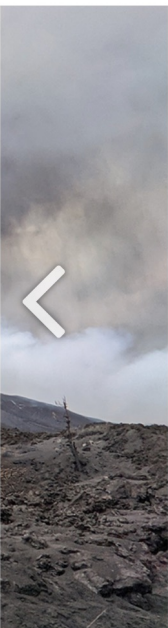
UNM Newsroom

About Us News Topics Schools and Colleges For News Media 



Ilanie Moses BCLab

UNM Newsroom / News / UNM's Volcanic Islands



A UNM research team was to develop drones to sample gases.

UNM's Volcanic Islands

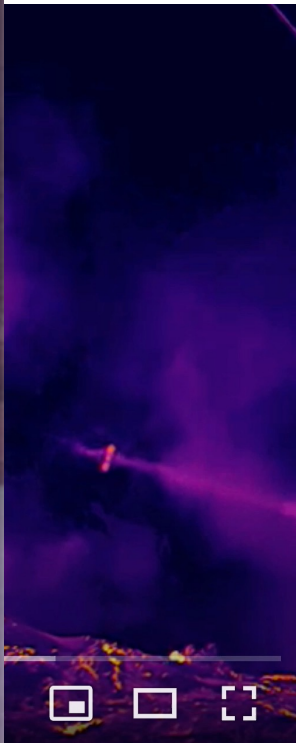
Researchers use drones to sample gases from active La Palma volcano
By Kim Deller @ January 05, 2023



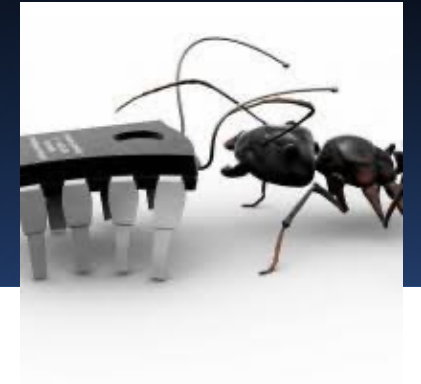
NASA MINDS ChiliHouse – Lunar Gardener



Aerial strategies advance volcanic gas measurements at inaccessible, strongly degassing volcanoes



Computational Immunology



Melanie Moses BCLab

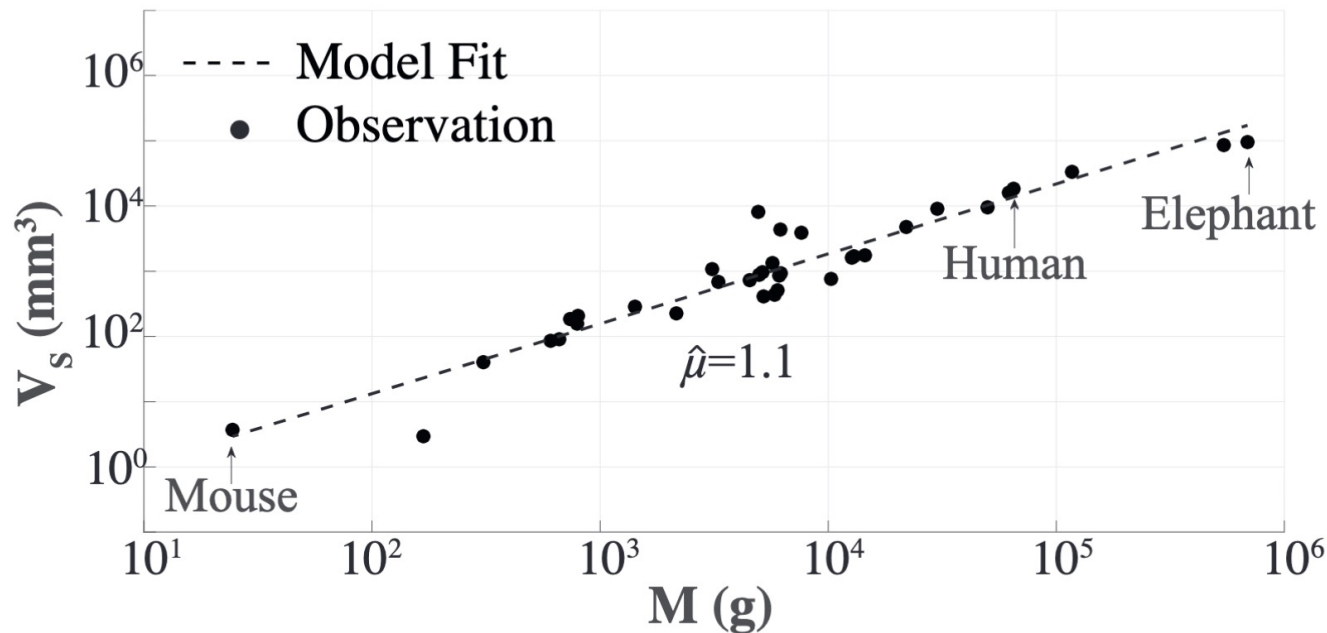


Fig. 1 Spleen Volume Scaling with Animal Mass. Species shown on log-log axes. The dashed line indicates the best model fit. The best-fit exponent for spleens (36 species) is $\hat{\mu} = 1.1$ 95% CI [0.95, 1.2]. $R^2 = 0.91$.

The Sizes and Numbers of Lymph Nodes Support a Scalable Immune Response in Mammals

Jannatul Ferdous^{1*}, G. Matthew Fricke^{1,2†}, Judy L. Cannon^{3†} and Melanie E. Moses^{1,4,5†}

¹Department of Computer Science.

²Center for Advanced Research Computing.

³Molecular Genetics and Microbiology.

⁴Biology Department, The University of New Mexico, Albuquerque, USA.

⁵Santa Fe Institute, Santa Fe, USA.

*Corresponding author(s). E-mail(s): jannat@unm.edu;

Contributing authors: mfricke@unm.edu;

jucannon@salud.unm.edu; melaniem@unm.edu;

[†]These authors contributed equally to this work.

Or... How to Build a Scalable Immune System

Machine Learning



INTERFACE

rsif.royalsocietypublishing.org

Research



Article submitted to journal

Subject Areas:

Agnostic Polymer Detection using Mass Spectrometry using Machine Learning

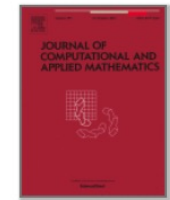
Victoria Da Poian^{*,1,2}, Lu Chou^{*,1,3}, Natalie Grefenstette^{*,4,5},
G. Matthew Fricke⁶, and Christopher P. Kempes⁴

¹NASA Goddard Space Flight Center, Greenbelt, USA, ²Microtel LLC, Greenbelt, USA, ³Georgetown University, Washington DC, USA, ⁴Santa Fe Institute, Santa Fe NM, USA, ⁵Blue Marble Space Institute of Science, Seattle WA, USA, ⁶Department of Computer Science, University of New Mexico, Albuquerque NM, USA.



Journal of Computational and Applied Mathematics

Volume 395, 15 October 2021, 113451



Machine learning feature analysis illuminates disparity between E3SM climate models and observed climate change



Melanie Moses BCLab

All these projects involved CARC computations.

Topics

- What is machine learning anyway?
 - Why should you care?
 - Why CARC?
-
- An Example ML Run at CARC, Monitoring GPU usage
(and a case for Jupyter Notebooks)

Just from today's talks...

- Jeremy Hogeveen (Neural Networks)
- Susan Atlas (Neural Networks)
- Sam McKenzie (K-means clustering)
- Tameem Talbash (Neural Networks)
- And more ...

Just from today's talks...

- Jeremy Hogeveen (Neural Networks)
- Susan Atlas (Neural Networks)
- Sam McKenzie (K-means clustering)
- Tameem Talbash (Neural Networks)
- And more ...

But the GPUs on our Machine Learning Cluster are often idle.

We need CARC ML users to capitalize on those GPUs.

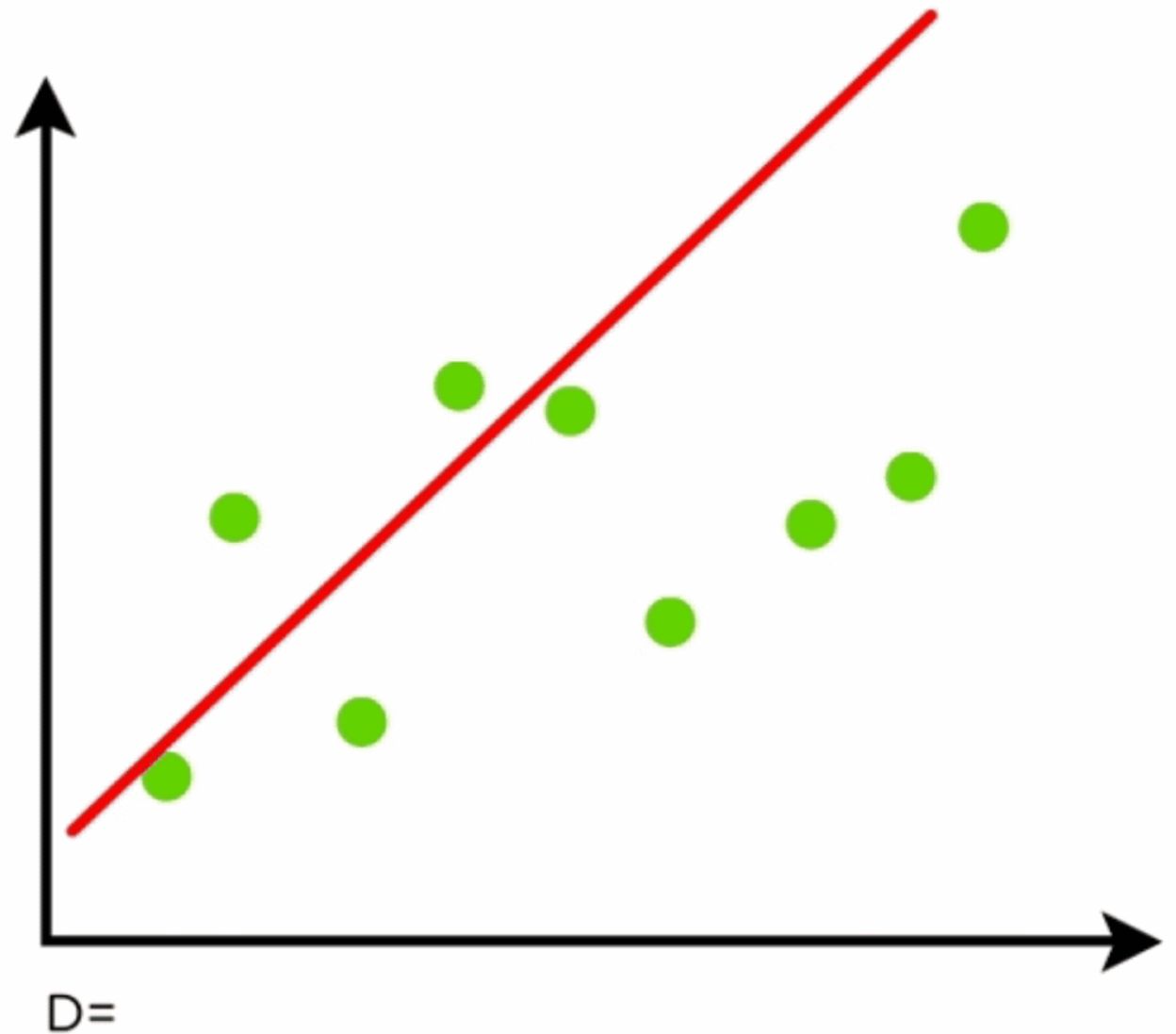
Good for you! Lot's of free GPU resources!



Herbert Simon

- “Learning is any process by which a system improves performance from experience”
- “Machine learning is concerned with computer programs that automatically improve their performance through experience”
- Professor of Psychology and Computer Science at Carnegie Mellon. Turing award and Nobel winner.

Regression is Machine Learning



Procedure Driven vs Data Driven

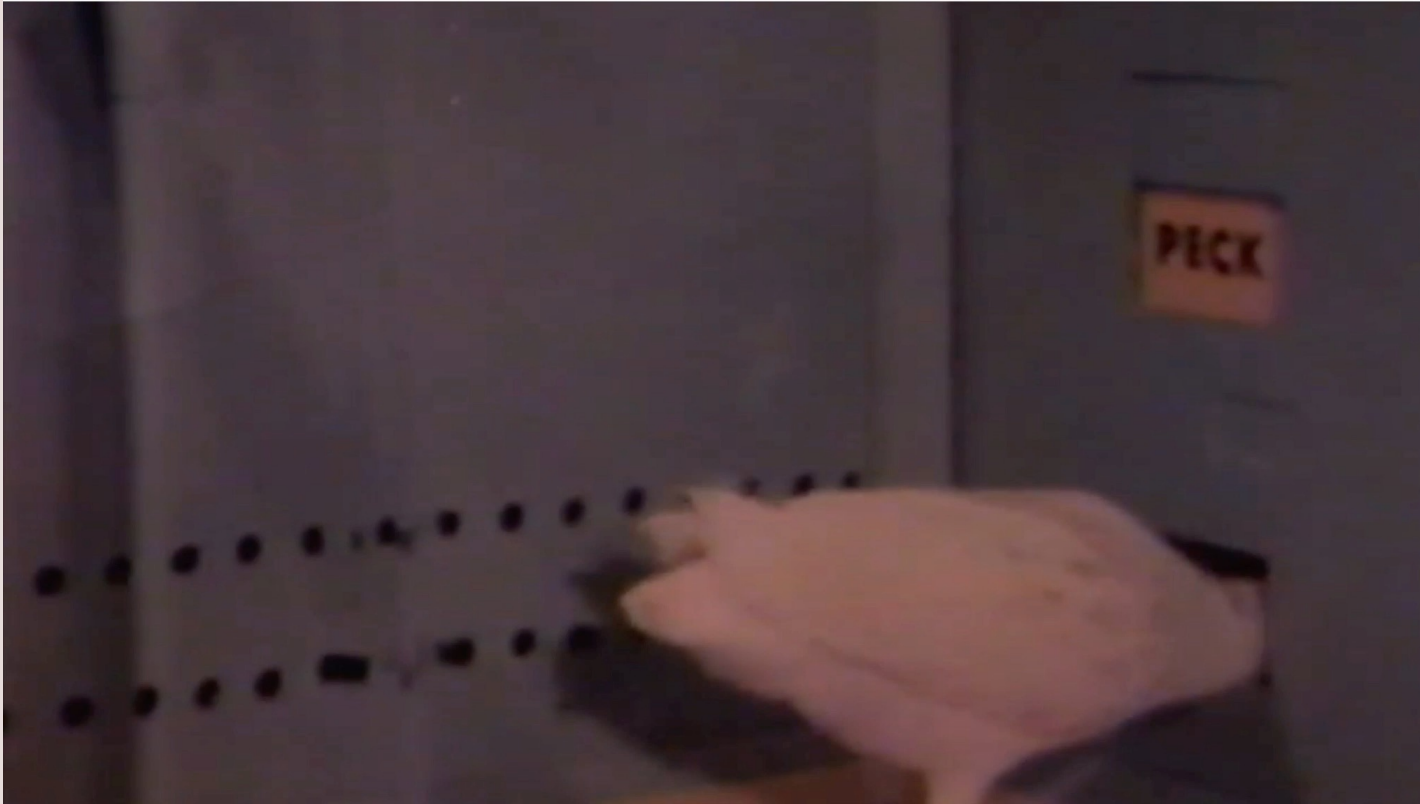
- Procedural programming requires an explicit understanding of the steps that map input to output.

Good Homemade Biscuits

3 cups self-rising flour
1 Tablespoon baking powder
1 Teaspoon Cream of Tartar
3/4 cup shortening (Crisco)

Butter milk to make soft dough
mix well. Put on to a floured
surface and knead until firm
enough to roll out. Cut with
biscuits cutter, brush tops with
melted butter. Bake in hot oven
400° until golden brown.
Yummy!

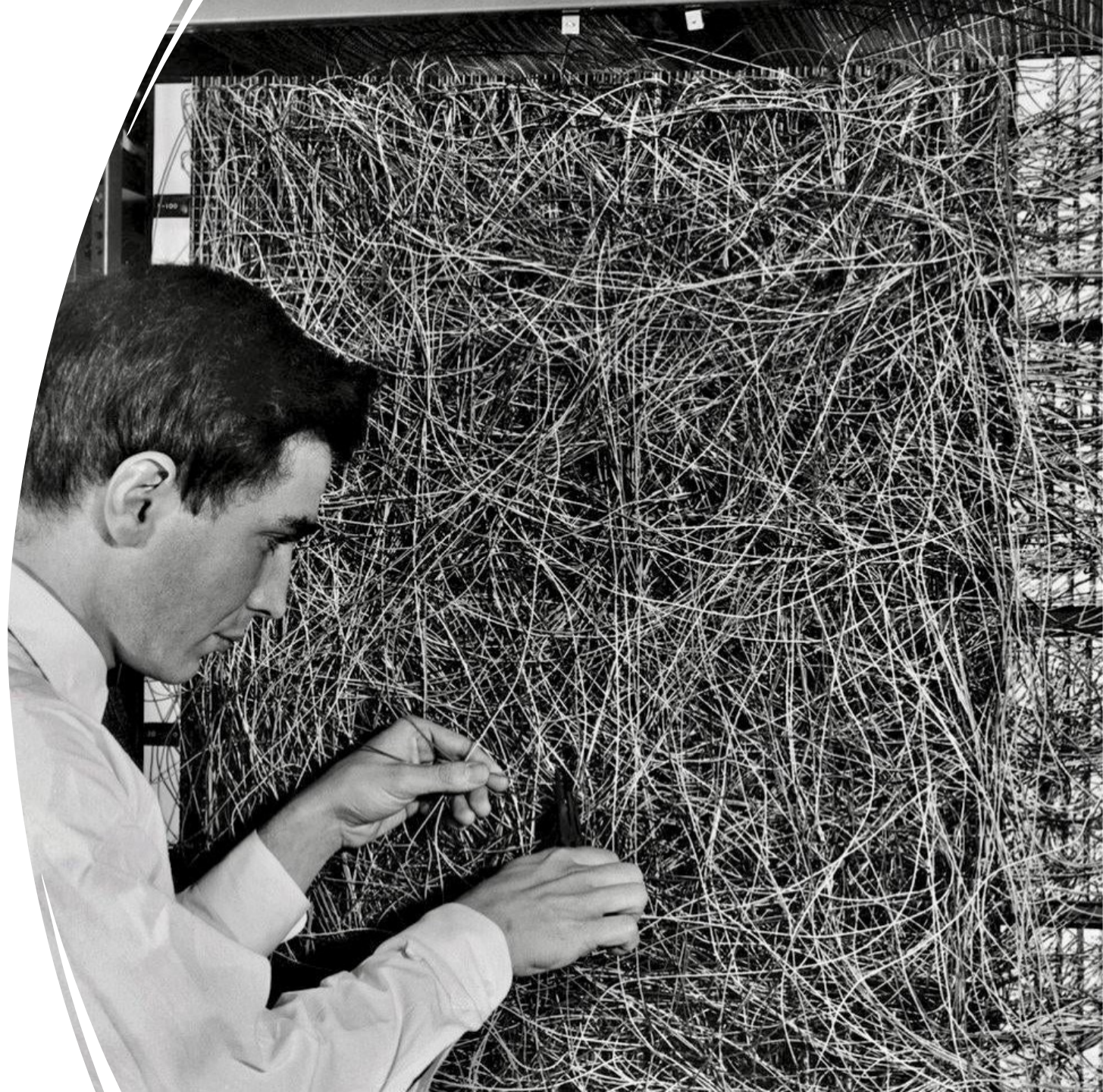
Procedure Driven vs Data Driven



- A data driven program maps input to output based on lots of example inputs and outputs.
- B.F. Skinner 1940s, Psychologist “Operant Conditioning”
- Pigeon trained to classify visual input.
- Food as reward

Neural Networks

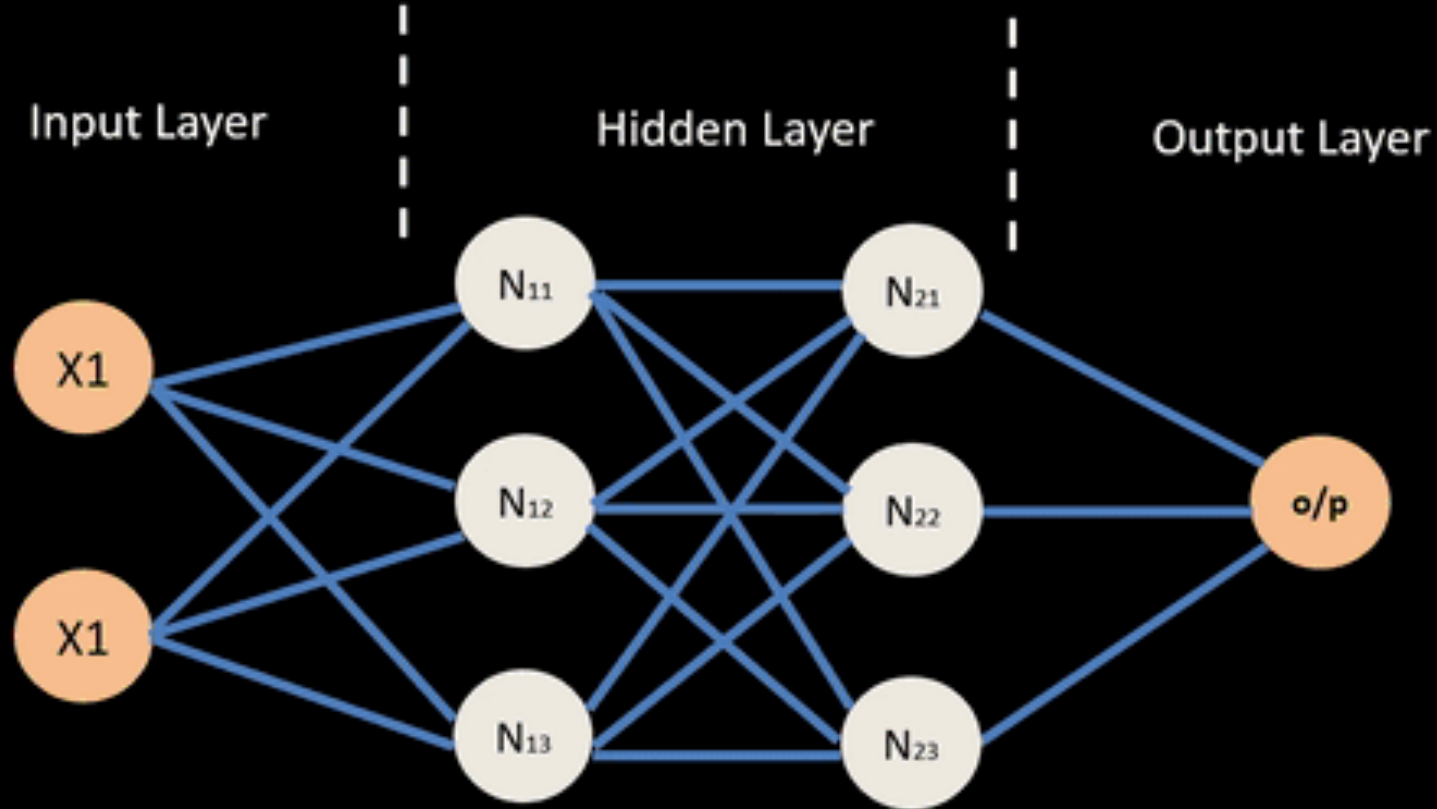
- First Perceptron created in 1958 by psychologist **Frank Rosenblatt**
- Seen here with a Mark I neural network.
- Rosenblatt referenced Skinner's work in several of his papers.



Neural Network – Backpropagation



INPUT

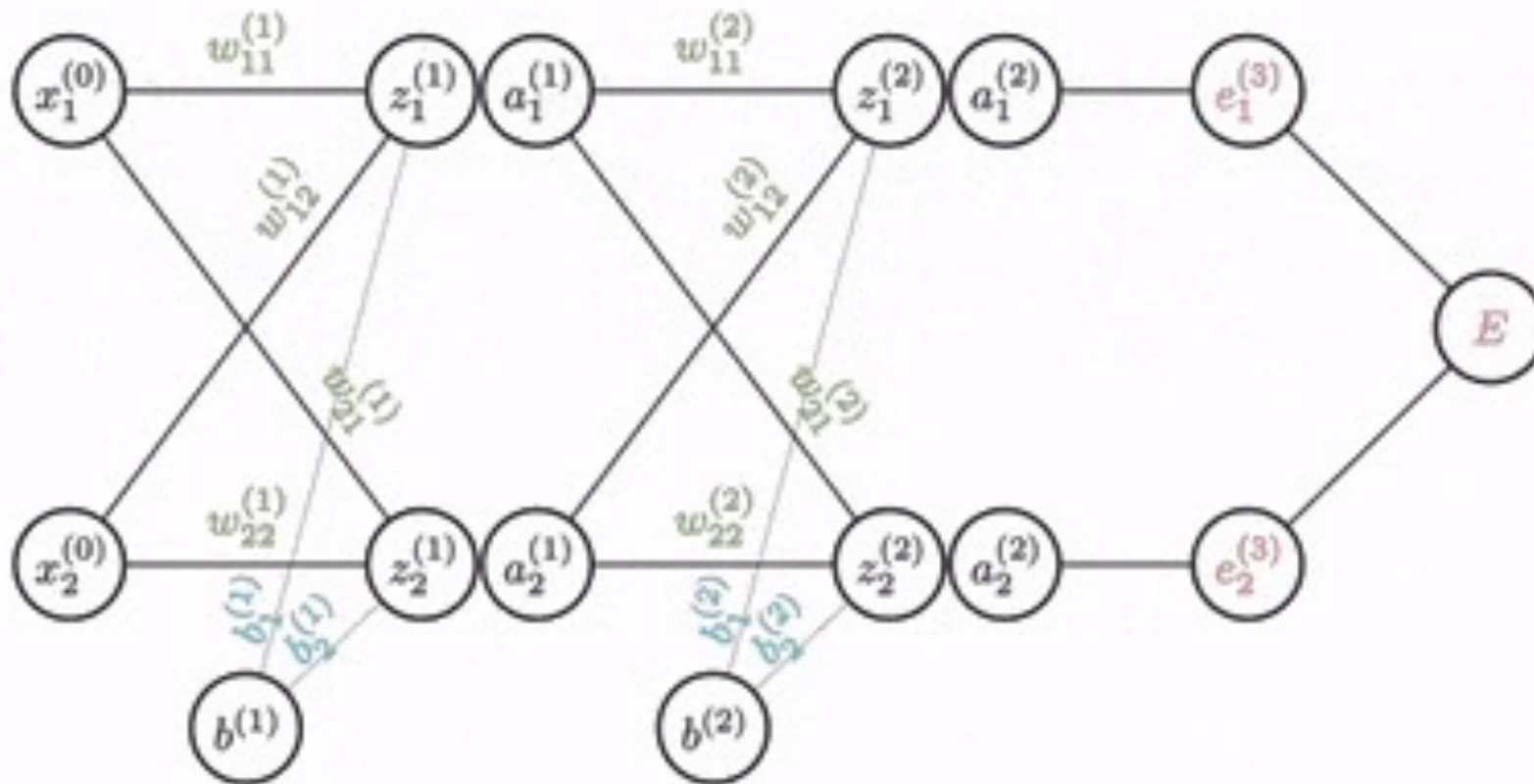
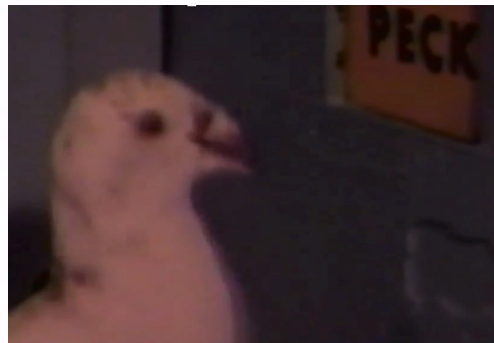


Turn
or
Peck

Multi-Layer Perceptron

$$\begin{bmatrix} \frac{\partial E}{\partial w_{11}^{(1)}} & \frac{\partial E}{\partial w_{12}^{(1)}} \\ \frac{\partial E}{\partial w_{21}^{(1)}} & \frac{\partial E}{\partial w_{22}^{(1)}} \end{bmatrix} = \begin{bmatrix} \delta_1^L \\ \delta_2^L \end{bmatrix}^\top \cdot \begin{bmatrix} \frac{\partial z_1^{(2)}}{\partial a_1^{(1)}} & \frac{\partial z_1^{(2)}}{\partial a_2^{(1)}} \\ \frac{\partial z_2^{(2)}}{\partial a_1^{(1)}} & \frac{\partial z_2^{(2)}}{\partial a_2^{(1)}} \end{bmatrix} \odot \begin{bmatrix} \frac{\partial a_1^{(1)}}{\partial z_1^{(1)}} \\ \frac{\partial a_2^{(1)}}{\partial z_2^{(1)}} \end{bmatrix} \odot \begin{bmatrix} \frac{\partial z_1^{(1)}}{\partial w_{11}^{(1)}} & \frac{\partial z_1^{(1)}}{\partial w_{12}^{(1)}} \\ \frac{\partial z_2^{(1)}}{\partial w_{21}^{(1)}} & \frac{\partial z_2^{(1)}}{\partial w_{22}^{(1)}} \end{bmatrix}$$

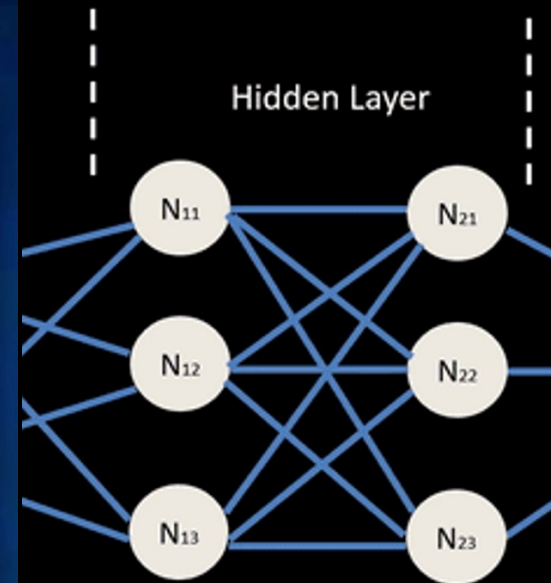
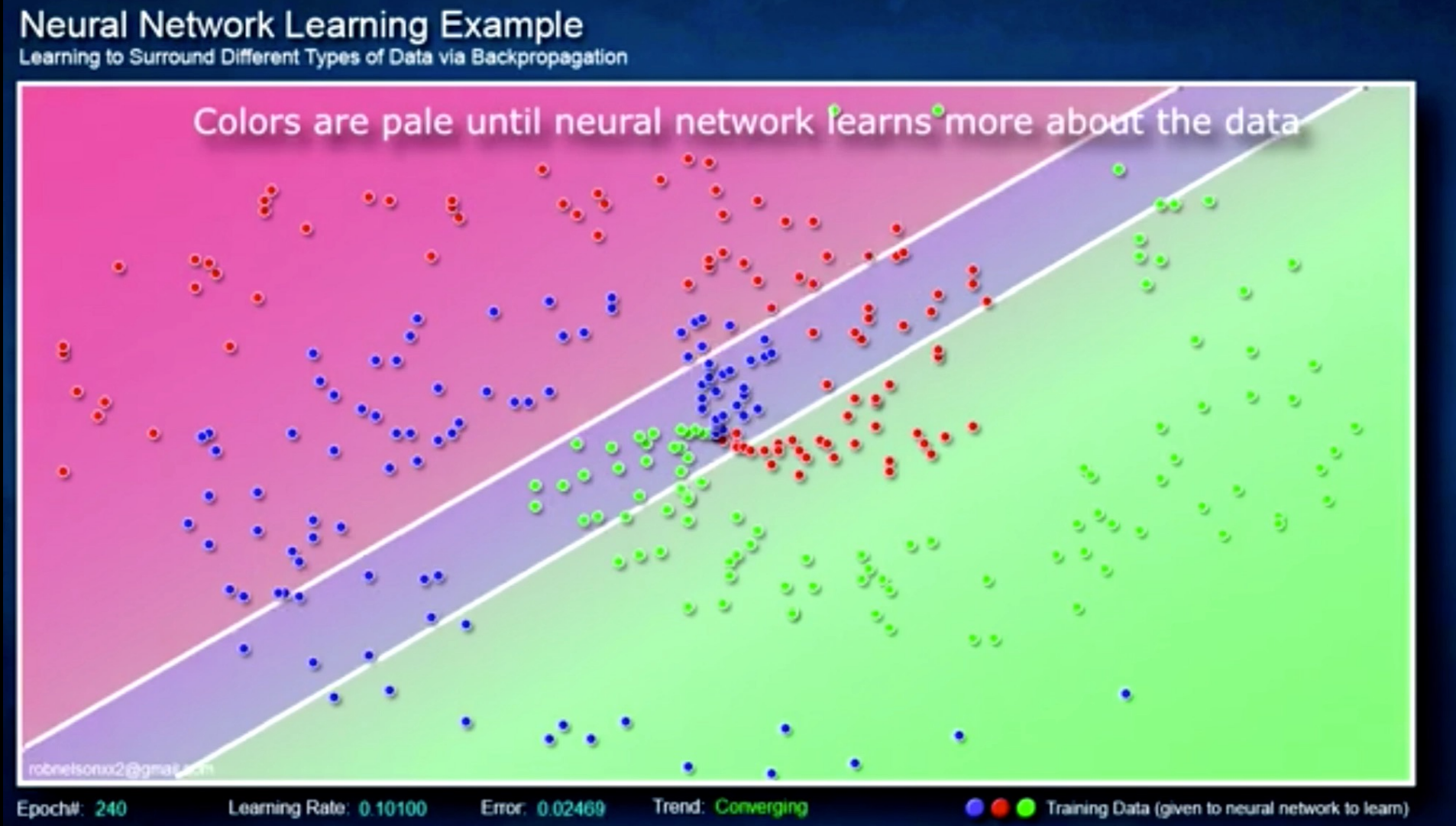
INPUT



Turn
or
Peck

Multi-Layer Perceptron

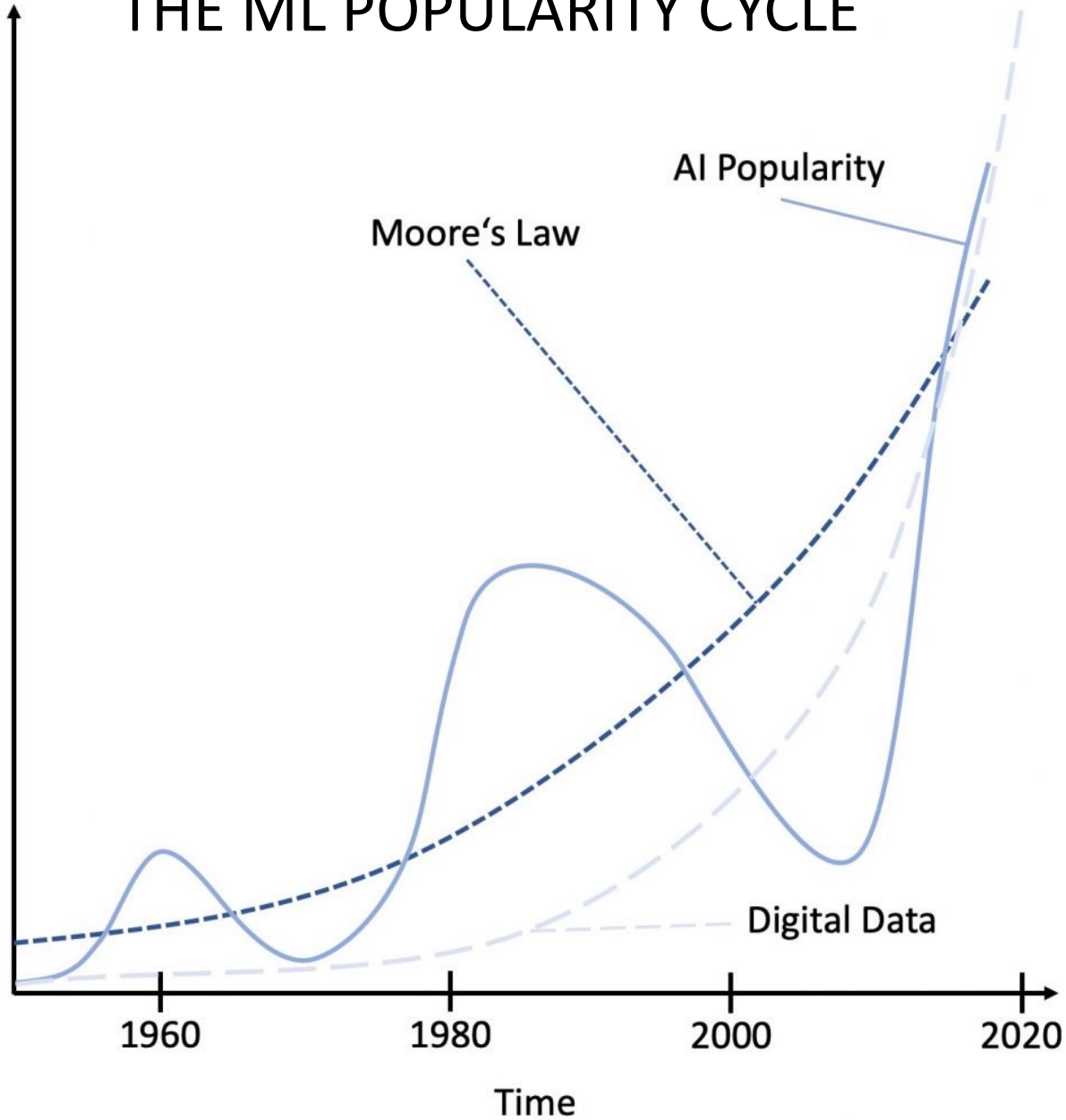
3 Category (red green blue) Classifier of data with 2 features (x,y position)



Why Does it Matter

- Machine learning allows us to “solve” (approximate) problems that we cannot solve analytically.

THE ML POPULARITY CYCLE



Current Upswing



TensorFlow
Google

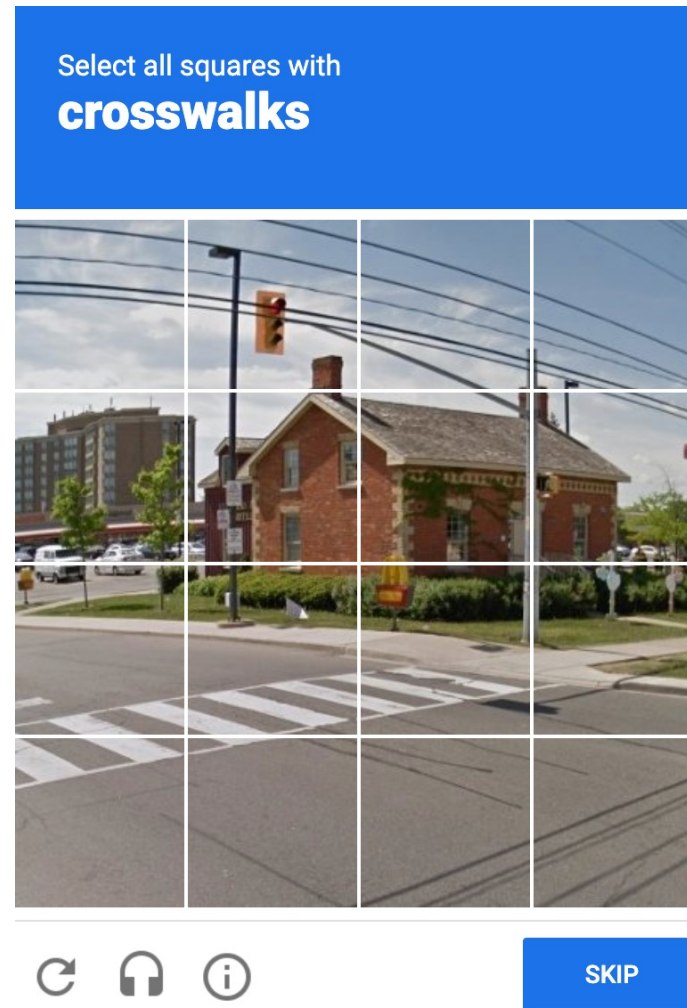
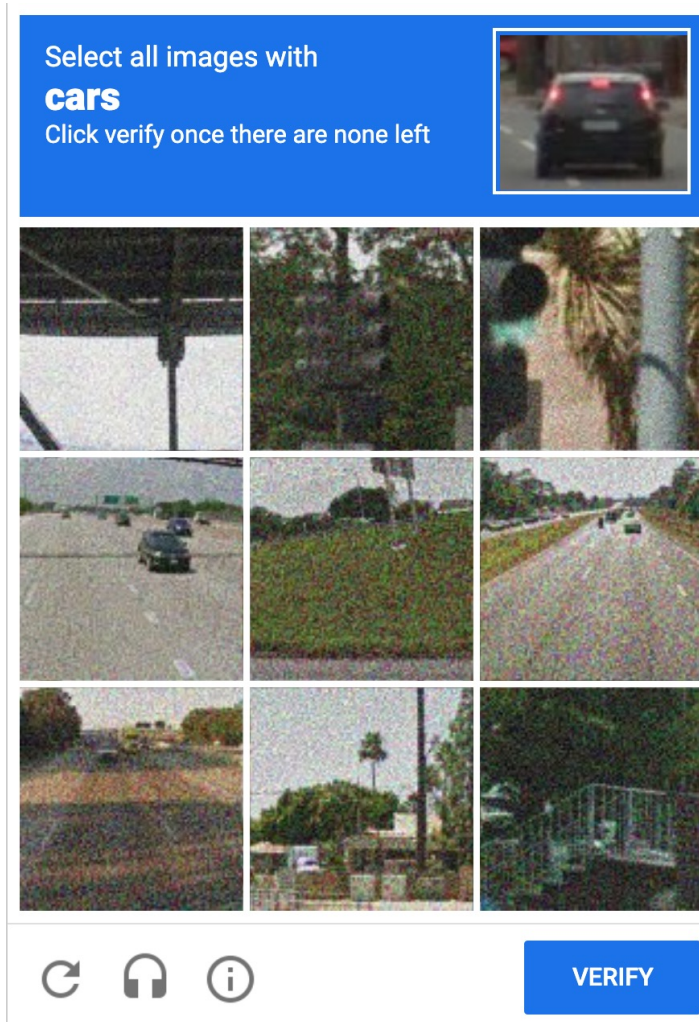
\$1.809 Trillion

facebook

PYTORCH

\$650 Billion

Current Upswing BIG DATA



reCAPTCHA is a **free** Google service



TensorFlow
Google

\$1.809 Trillion

facebook

PYTORCH

\$650 Billion

Convolutional Neural Networks (CNNs) for Image Classification.

Input image



Convolution
Kernel

$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$

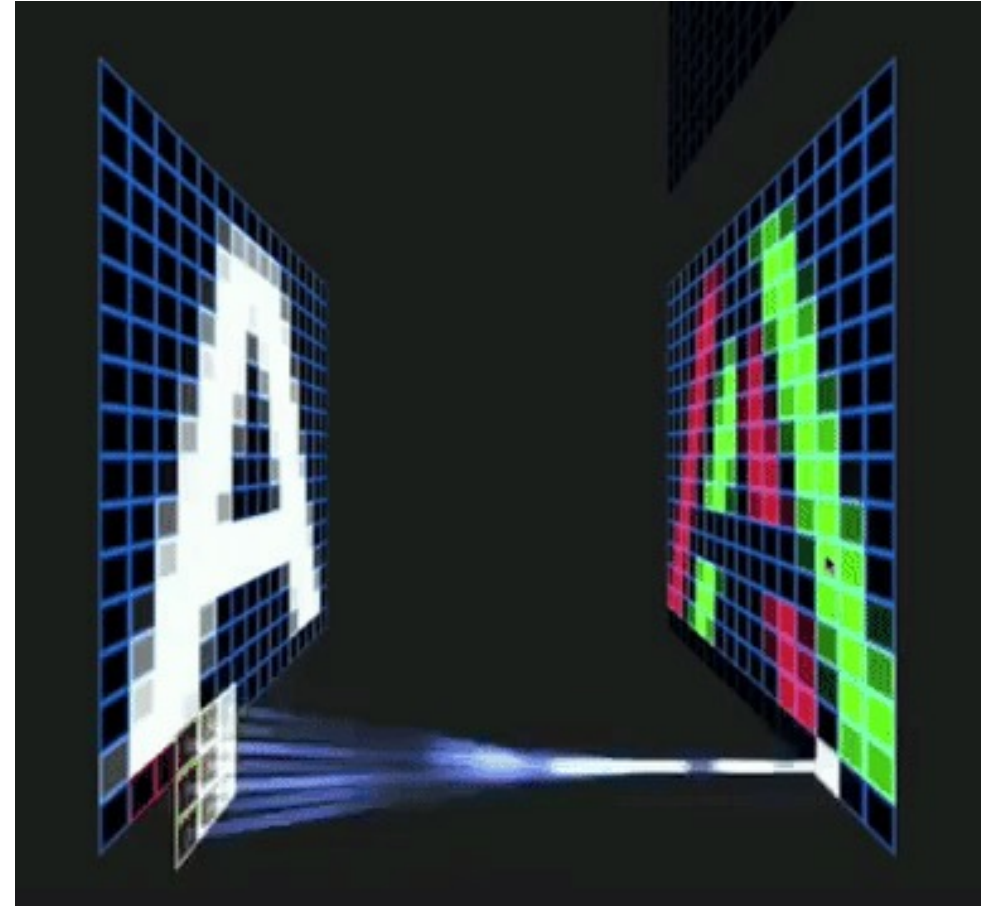
Feature map

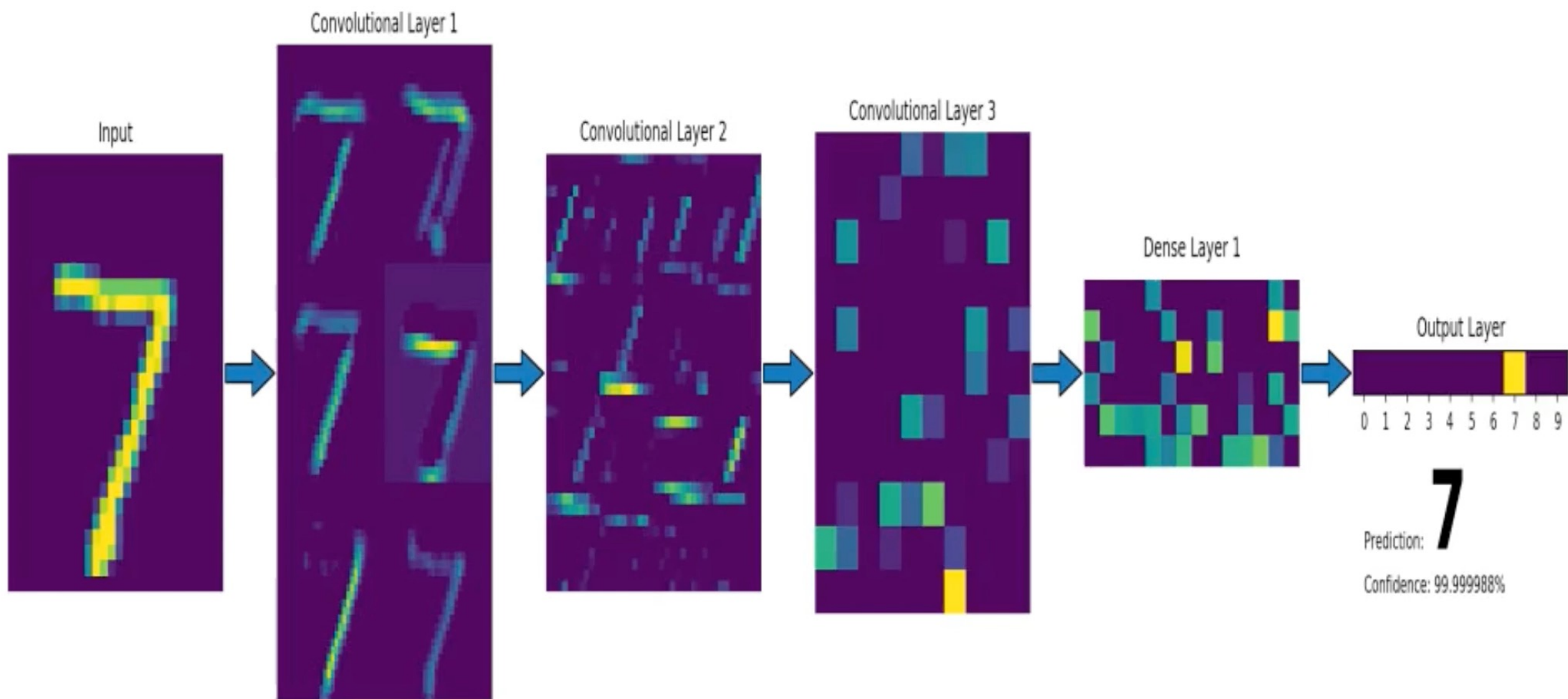


Allows the network itself to

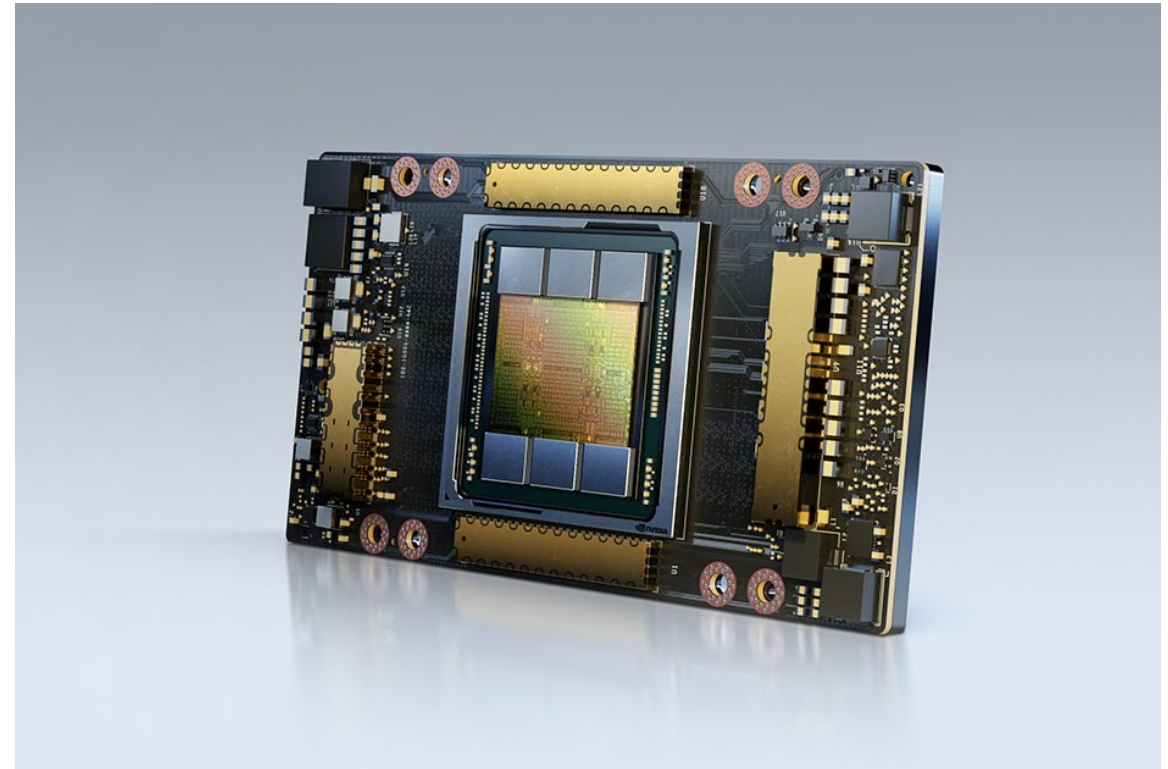
- Detect Edges
- Sharpen
- Blur
- Rotate

And generally, remove extraneous information and focus on the things that help it classify.

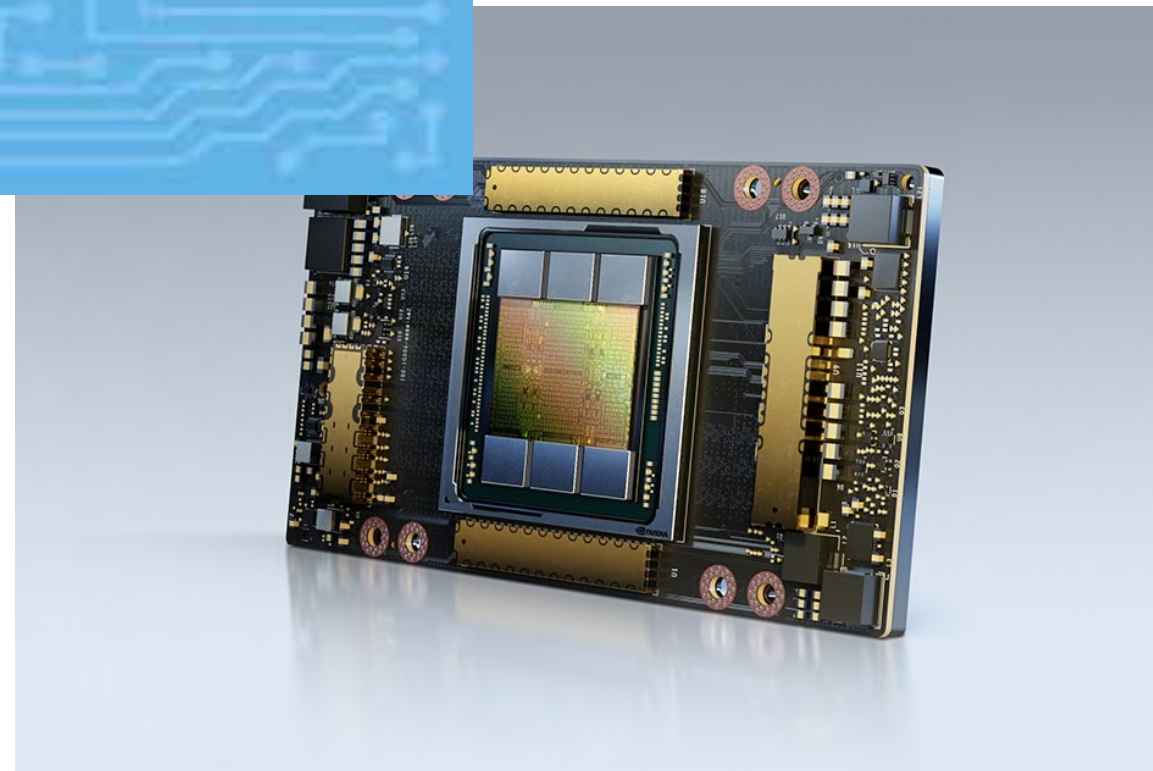
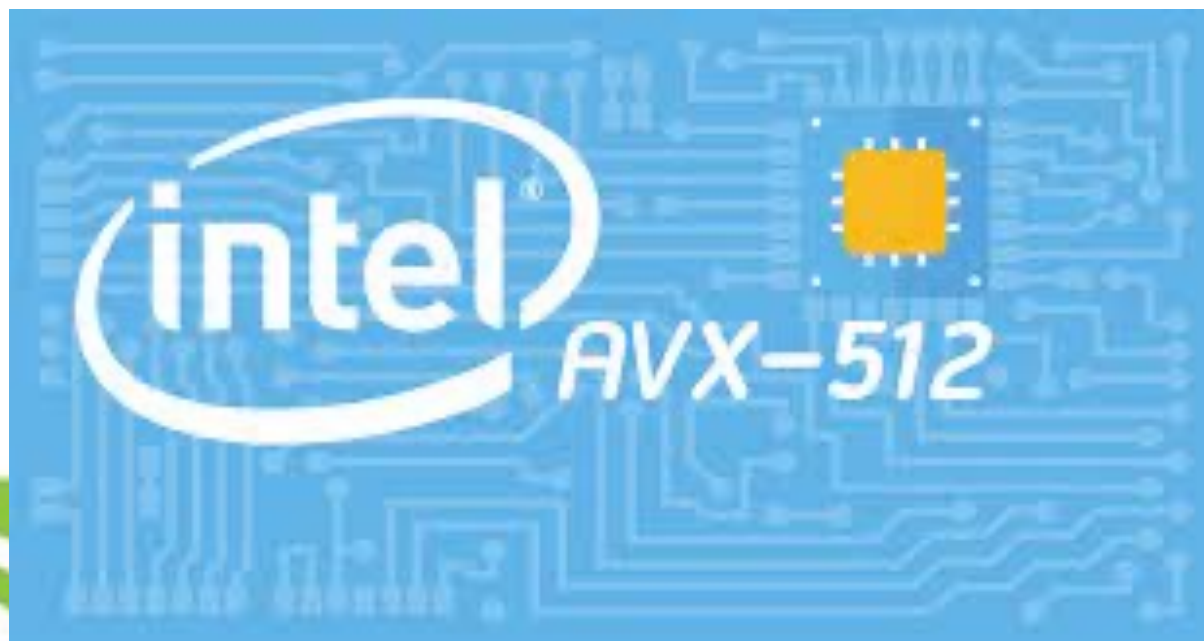




Hardware that's Good at Matrix Multiplication



ML has become so important that it is driving GPU and CPU architectures.

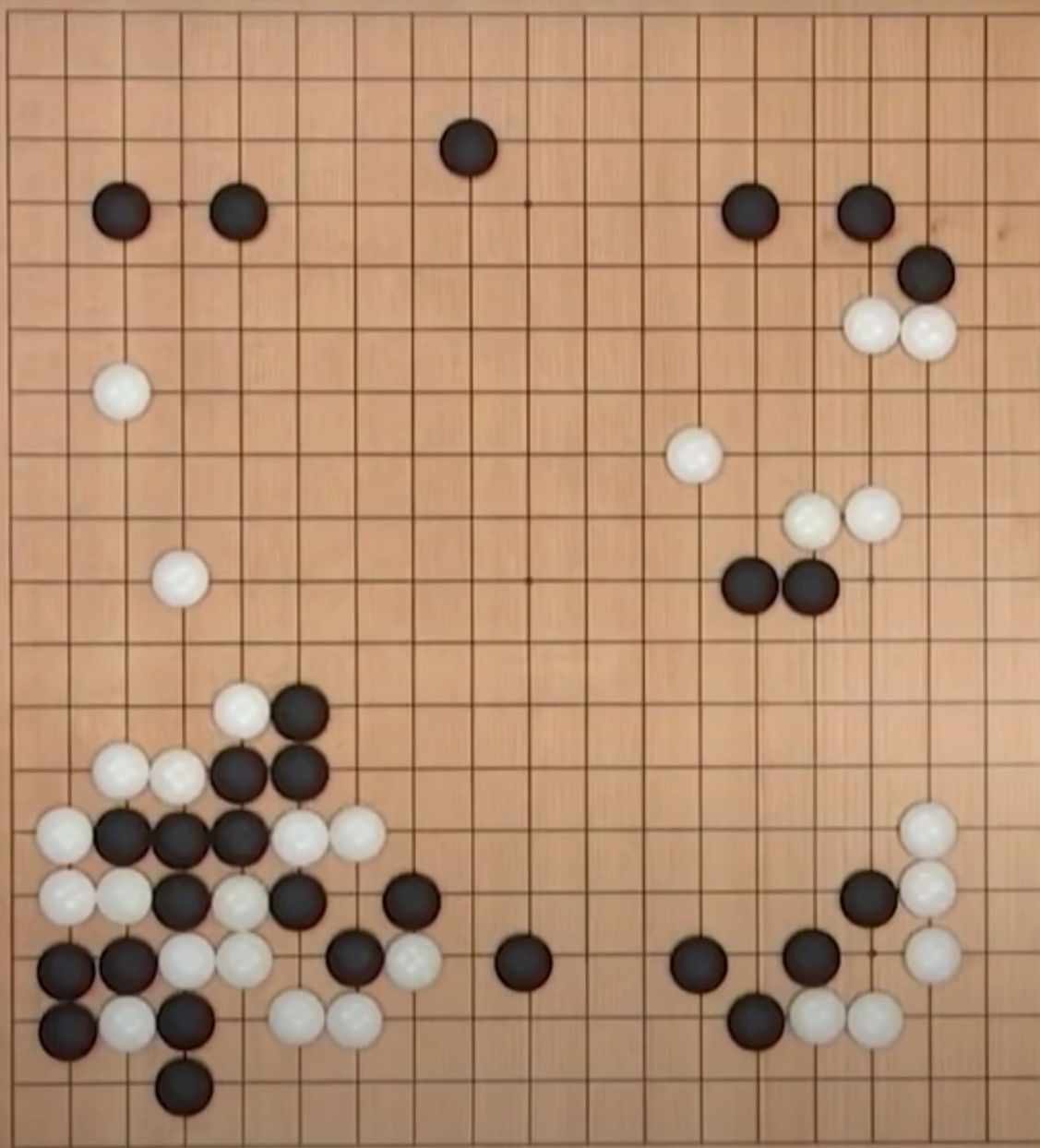


ML has become so important that it is driving GPU and CPU architectures.

The Good and the Bad and the Surprising

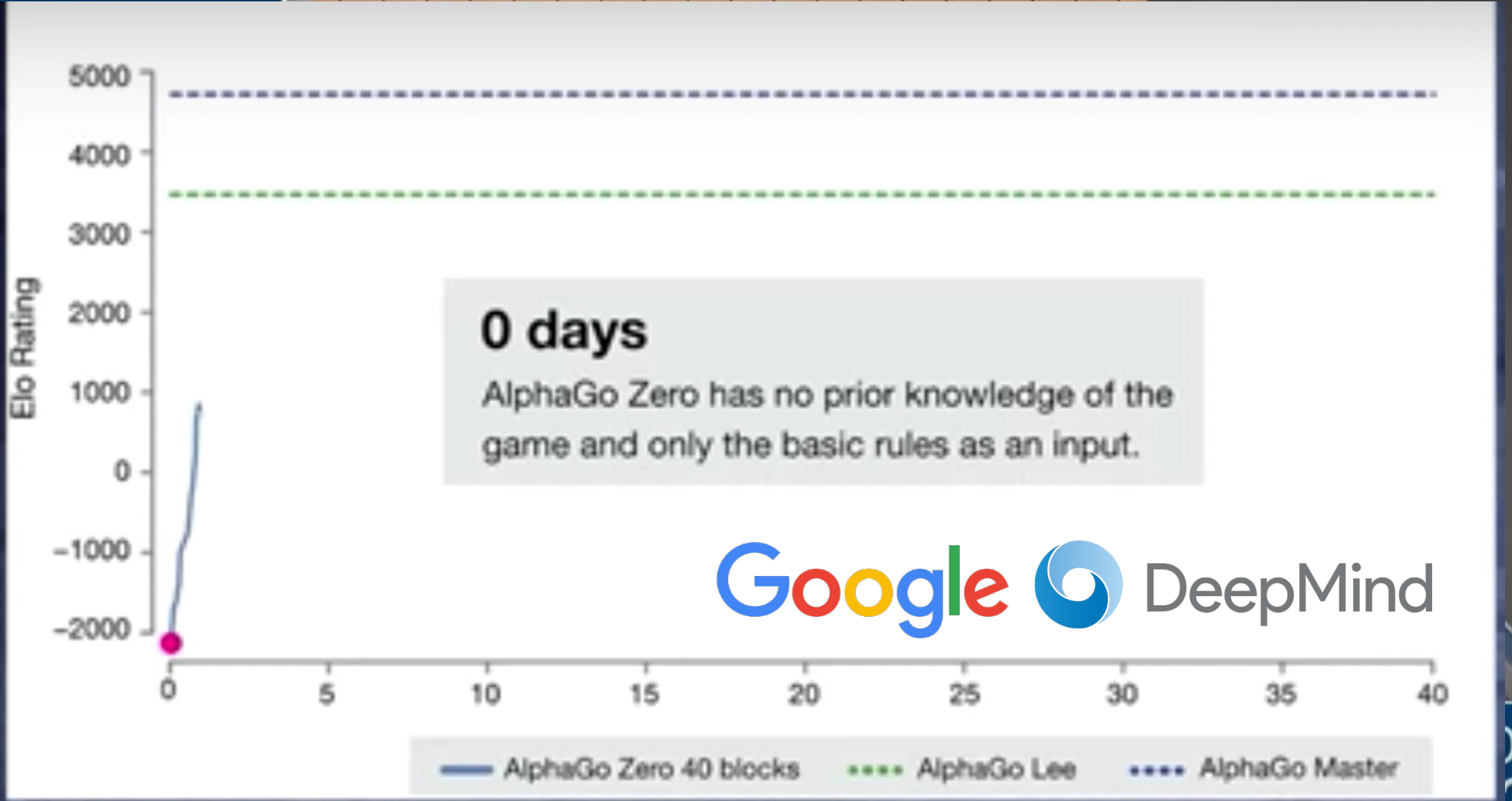
- Being able to train machine learning models from vast datasets has some great applications:
- But the model trained depends entirely on the data used to train it.

LEE SEDOL
01:07:44

The logos for Google and DeepMind are displayed in the top right corner. The Google logo is in its multi-colored font, and the DeepMind logo is in a grey, sans-serif font.

The surprising: 2016

ALPHAGO
01:28:32



BETA

[Search](#)Examples: [Free fatty acid receptor 2](#) [At1g58602](#) [Q5VSL9](#) [E. coli](#)Help: [AlphaFold DB search help](#)

T-cell immunomodulatory protein homolog

AlphaFold structure prediction

Download

[PDB file](#)[mmCIF file](#)[Predicted aligned error](#)[NEW Feedback on structure](#)[Looks great](#)[Could be improved](#)[Information](#)

'The game has changed.' AI triumphs at solving protein structures

In milestone, software predictions finally match structures calculated from experimental data – Science Magazine

CARC Alphafold users: - Cristian Bologa - Melanie Moses



The Good

AlphaFold

Search for

Examples:

T-cell

AlphaFold

Download

NEW Feature

Information

'The game
structures

In milestone
calculated from experimental data – Science Magazine



The Bad

- Machine learning has been applied to
 - Mortgage Approval
 - Parole Hearings in New Mexico
 - Google hiring and promotion
- The input data is the history of parole, mortgage decisions, and hiring.
- See the problem?



(a) Three samples in criminal ID photo set S_c .



(b) Three samples in non-criminal ID photo set S_n

Figure 1. Sample ID photos in our data set.

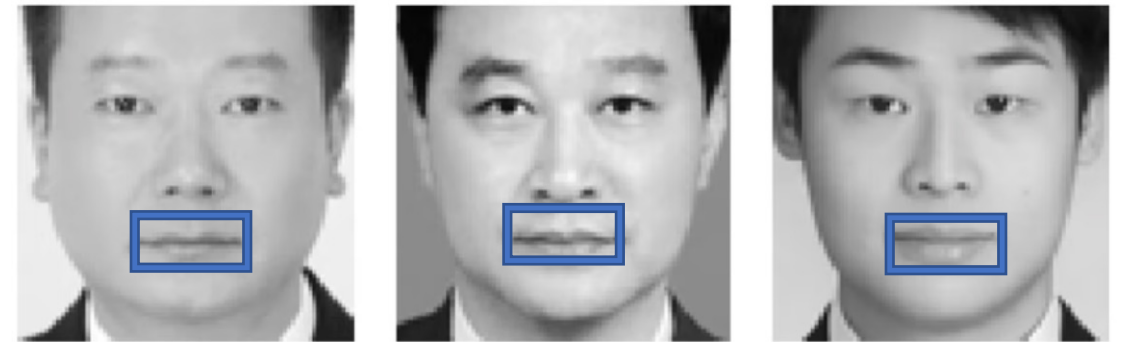
Xiaolin, Wu, and X. Zhang. "Automated Inference on Criminality using Face Images." *CoRR.*–2016 (2016).

The Bad

- Machine learning has been applied to
 - Mortgage Approval
 - Parole Hearings in New Mexico
 - Google hiring and promotion
- The input data is the history of parole, mortgage decisions, and hiring.
- See the problem?
- These models all learned to associate gender and race with the outcome.



(a) Three samples in criminal ID photo set S_c .



(b) Three samples in non-criminal ID photo set S_n

Figure 1. Sample ID photos in our data set.

Xiaolin, Wu, and X. Zhang. "Automated Inference on Criminality using Face Images." *CoRR.*–2016 (2016).

Office of the General Counsel
Rules Docket Clerk
Department of Housing and Urban Development
451 7th Street, SW, Room 10276
Washington, DC 20410-0001.



SANTA FE
INSTITUTE



Cris Moore

Melanie Moses

Regarding Docket No. FR-6111-P-02, HUD's Implementation of the Fair Housing Act's Disparate Impact Standard

We are a group of computer scientists, social scientists and legal scholars who are writing to express our concern about the proposed amendments to HUD's implementation of the Fair Housing Act, and in particular those amendments related to the use of algorithms.

We believe that the proposed amendments are based on a failure to recognize how modern algorithms can result in disparate impact, even in the absence of discriminatory intent, and how subtle the process of auditing algorithms

ML for CoVID Spike Modelling – GOOD!

Inferring Criminality from Appearance – BAD!



(a) Three samples in criminal ID photo set S_c .



(b) Three samples in non-criminal ID photo set S_n

Figure 1. Sample ID photos in our data set.

Why CARC?

- As we have seen, machine learning depends on data. To process all that data one needs significant computational resources.
- CARC has those resources.
- ML has become so important that it is driving GPU and CPU architectures.
- K40s, v100s, a100 Nvidia GPUs.
- Avx512 CPUs on Hopper.
- Large memory nodes – 1 TB and 3 TB RAM nodes



Xena Cluster

*Susan Atlas



x32

```
mfricke@xena:~ $ qgrok
```

queues	free	busy	offline	jobs	nodes	CPUs
-----	-----	-----	-----	-----	-----	-----
singleGPU	11	12	0	12	23	368
dualGPU	3	1	0	1	4	64
bigmem-1TB	1	1	0	2	2	128
bigmem-3TB	1	1	0	1	2	128
debug	1	0	0	0	1	16
systems	1	0	0	0	1	16
totals:	18	15	0	16	33	720

Machine Learning QuickByte Tutorials

[Http://carc.unm.edu](http://carc.unm.edu)

- R
 - [R Programming in HPC](#)
 - [R at CARC](#)
 - [Running R in Parallel with Future](#)
 - [Gurobi optimizer with R](#)
- Machine Learning
 - [Tensorflow](#)
 - [Machine Learning Conda Environments \(including PyTorch K40 GPU libraries\)](#)
 - [Tensorflow with multiple GPUs](#)
 - [Parallelization with JupyterHub using Dask and SciKit-learn](#)
- Bioinformatics
 - [Genomic variant calling with GATK](#)
 - [Genome evaluation with QUAST and BUSCO](#)
 - [Single genome demographic history with PSMC](#)

Files Running IPython Clusters

Select items to perform actions on them.

Upload New ↕

<input type="checkbox"/>	0	▼	📁 / workshops / ML
<input type="checkbox"/>	📁 ..		
<input type="checkbox"/>	📁 checkpoints_best_cnn		
<input type="checkbox"/>	📁 checkpoints_best_mlp		
<input type="checkbox"/>	📁 data		
<input type="checkbox"/>	📄 GPUvsCPUMatMul.ipynb		
<input type="checkbox"/>	📄 GPUvsCPUMatMulSaved.ipynb		
<input type="checkbox"/>	📄 SVHN_classifier.ipynb		
<input type="checkbox"/>	📄 SVHN_classifierSaved.ipynb		
<input type="checkbox"/>	📄 TensorBoard.ipynb		

Notebook:	
Julia 1.5.2	
Python 3	
Python [conda env:.conda-MSML]	
Python [conda env:.conda-MSML_TF_GPU]	
Python [conda env:.conda-MSML_tensorflow_decision_forest]	
Python [conda env:.conda-R]	
Python [conda env:.conda-SVM]	kB
Python [conda env:.conda-alphafold-env]	kB
Python [conda env:.conda-chem501]	kB
Python [conda env:.conda-cjs14]	kB
Python [conda env:.conda-comp_immunology]	kB
Python [conda env:.conda-curvefit]	kB

CPUs Total:
Hosts up:
Hosts down:

384
24
0

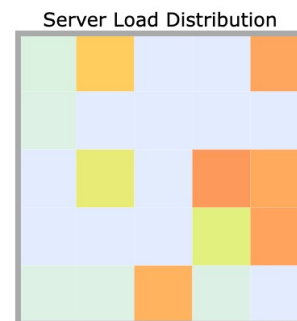
Current Load Avg (15, 5, 1m):

22%, 22%, 22%

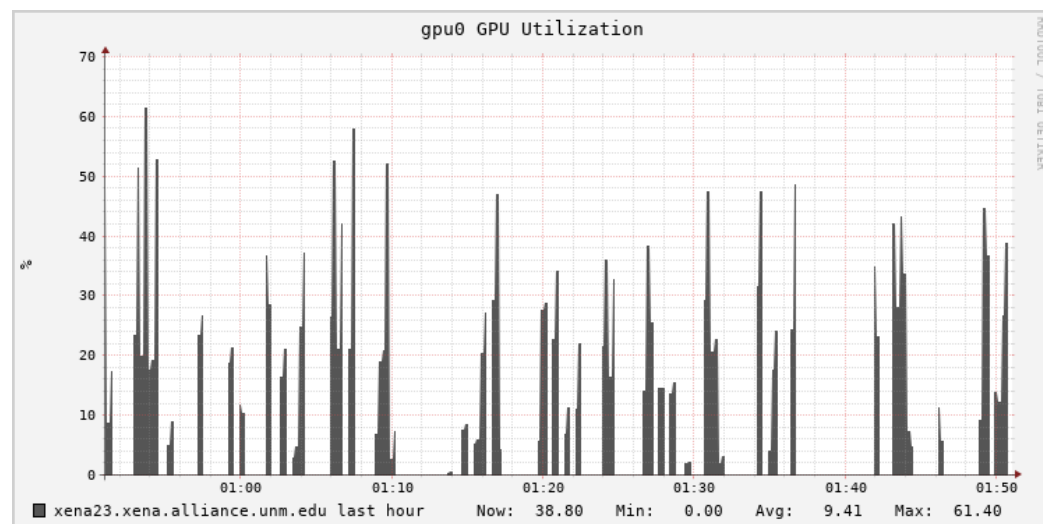
Avg Utilization (last hour):

23%

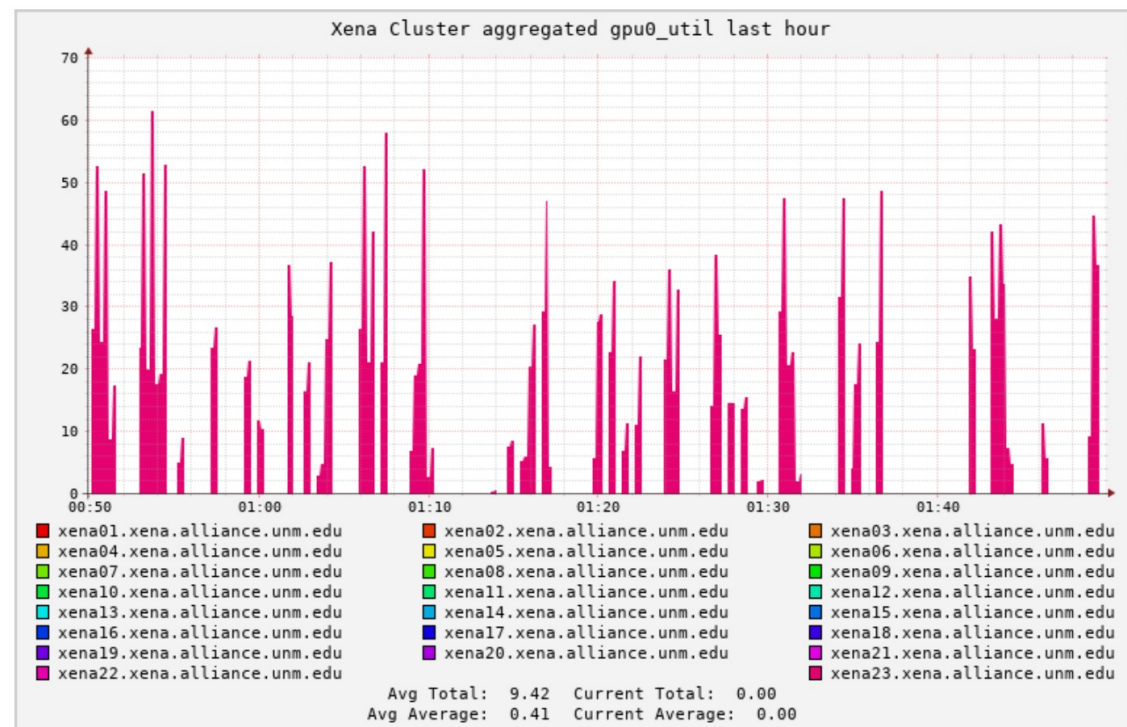
<http://xena.alliance.unm.edu/ganglia>



Stacked Graph - gpu0_util (%)



Xena23 GPU0 Utilization



Happy Hour!

Please come to my open office hours!

Send email to help@carc.unm.edu