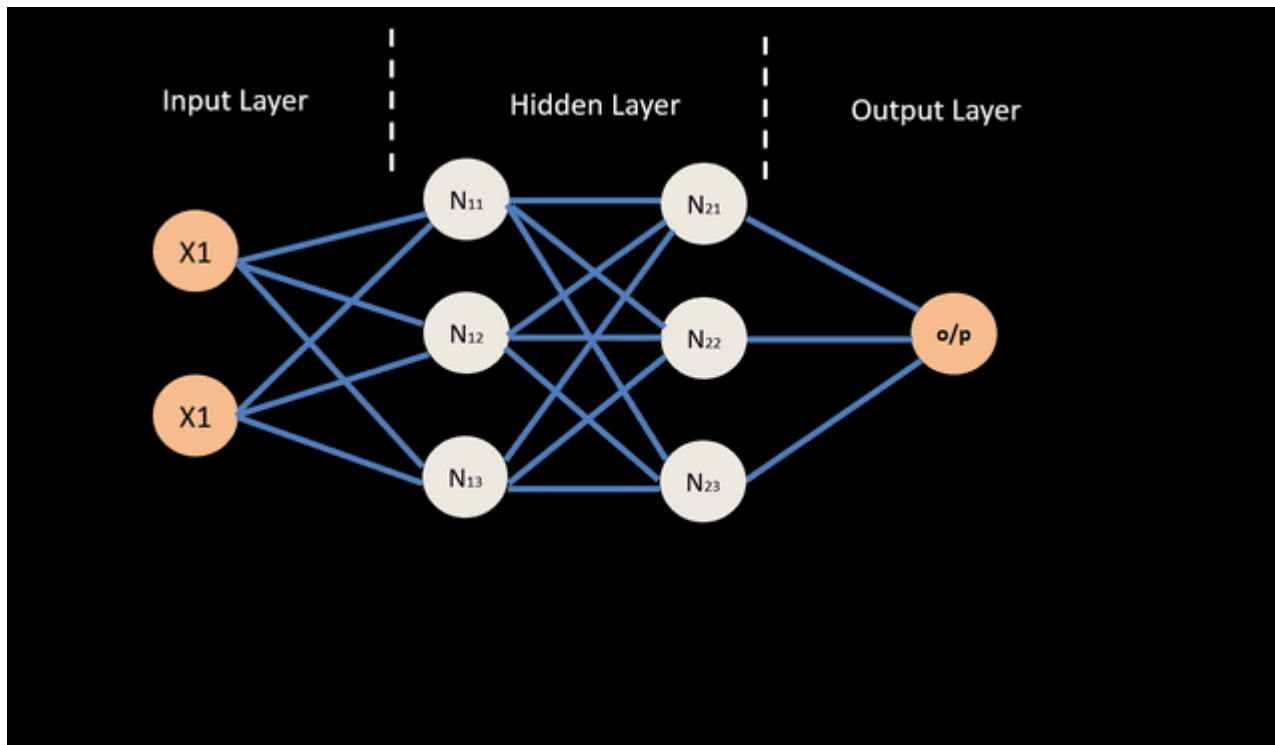# Lecture 11-a

Unconventional Algorithms and Hardware for Neural Networks
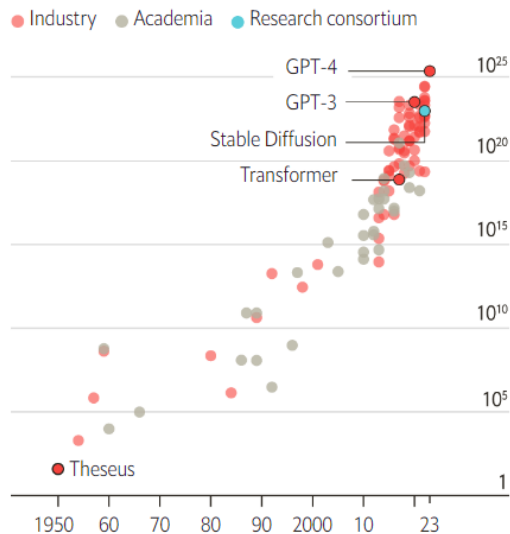
# Artificial Neural Networks

# Artificial Neural Networks
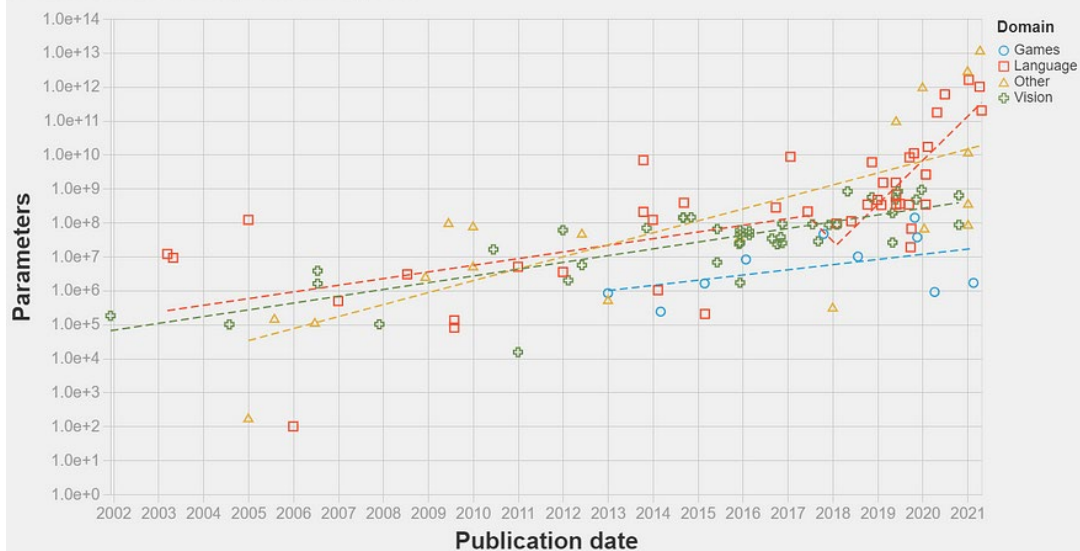
Computing power used in training AI systems
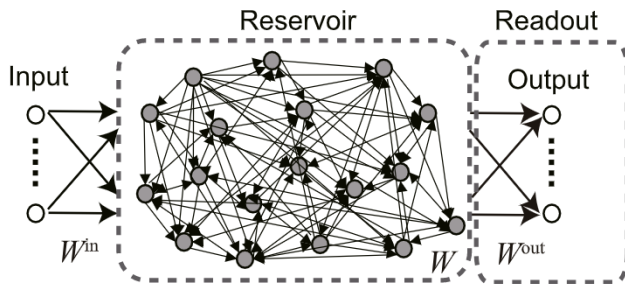Selected systems, floating-point operations, log scale

● Industry  ● Academia  ● Research consortium

Sources: Sevilla et al., 2023; Our World in Data



Parameter count of ML systems through time
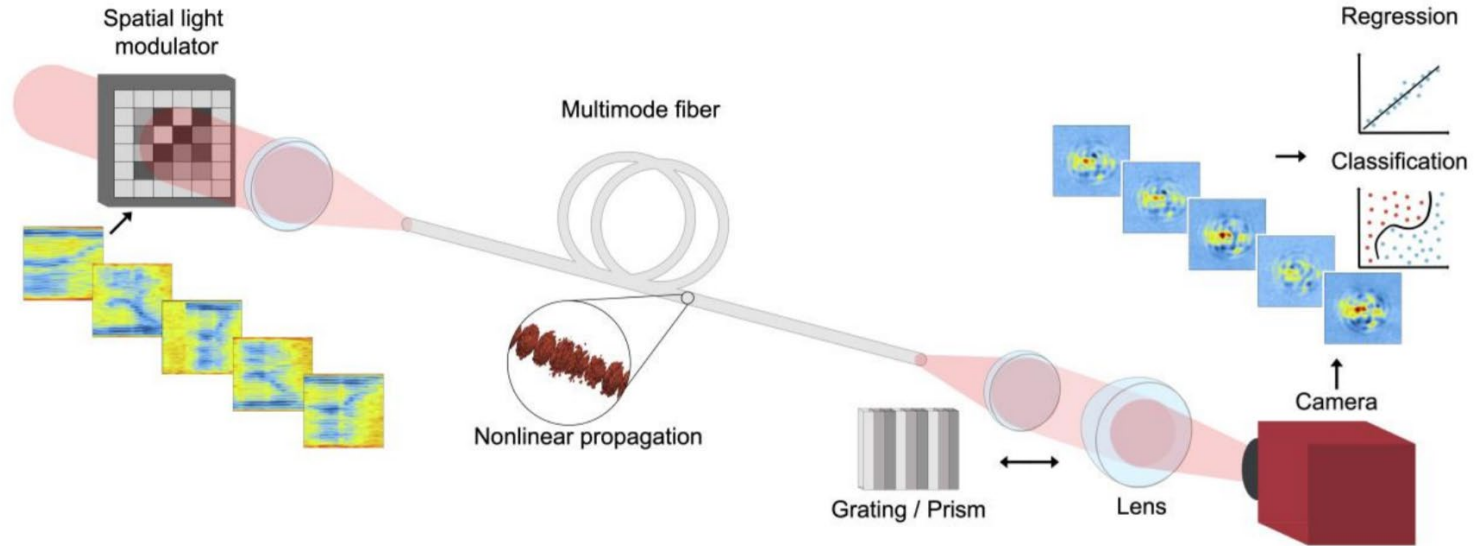
# Reservoir Computing

- Consists of high-dimensional, fixed, non-linear connections for transforming data, only output weights are optimized.

- Especially useful for fast operation and low training cost

- Can be realized by **high dimensional** and **nonlinear** physical systems.



(a) Conventional RC

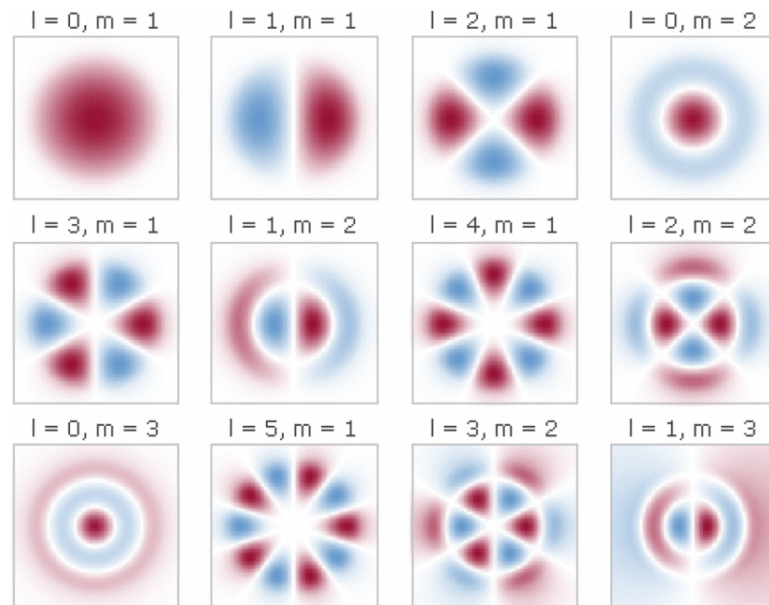# Computing with Spatiotemporal Nonlinearities of MMFs

Spatial light modulator

Multimode fiber

Nonlinear propagation

Grating / Prism

Lens

Camera

Regression

Classification

Teğin, U., Yıldırım, M., Oğuz, İ., Moser, C. & Psaltis, D. Scalable optical learning operator. *Nat. Comput. Sci. 2021 18* **1**, 542–549 (2021).

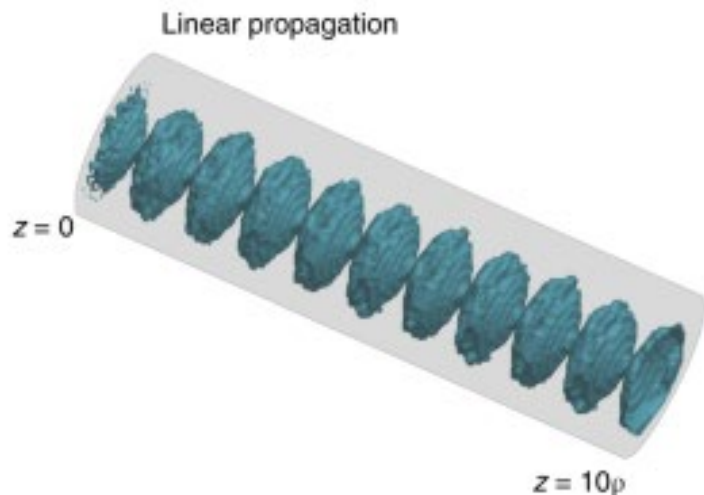# Computing with Spatiotemporal Nonlinearities of MMFs

- In multimode fibers light propagate in discrete channels, depending on their properties MMFs can support up to millions of channels:

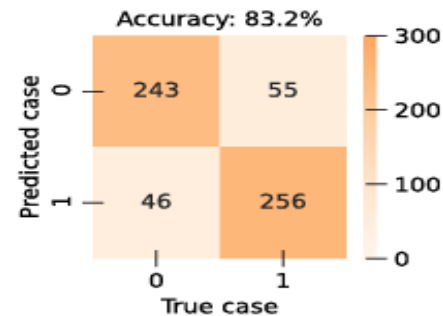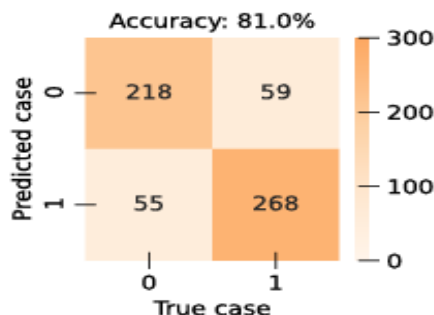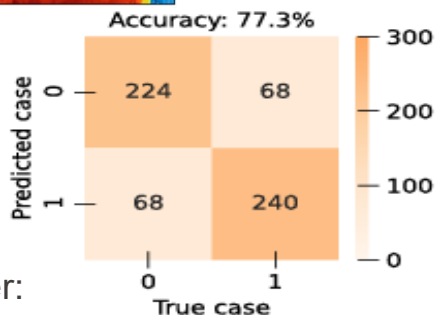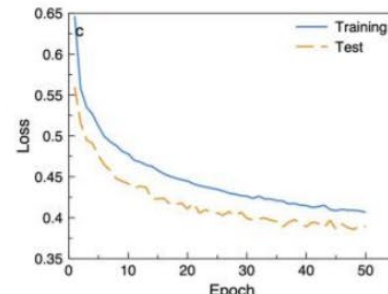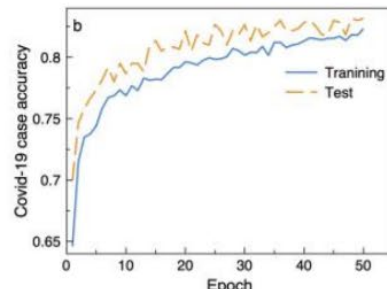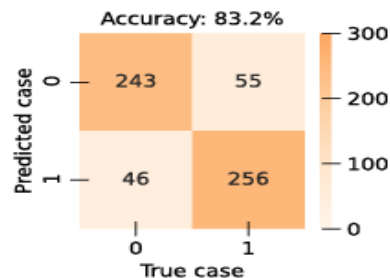$$E(x, y, \omega) = \sum_n^N A_n F_n(x, y, \omega)$$

- At high intensities modes start to couple each other due to light-matter interactions.

# Computing with Spatiotemporal Nonlinearities of MMFs

Linear propagation

$z = 0$

$z = 10p$

# Computing with Spatiotemporal Nonlinearities of MMFs

W/out fiber: 77%

# Programming
# Propagation inside MMFs

Programming Algorithm

Higher Performance

?

Reservoir

Readout

Input

Output

$W^{out}$

# Programming Propagation inside MMFs

# Programming Propagation inside MMFs

Programming Parameters

Mode Shape Coefficients

Experimental Parameters

$$WF(x,y) = \sum_{n=1}^{N} (a_{2n-1} + i * a_{2n}) * F_n(x,y)$$

**Programming Wavefront**

Intensity

Diffraction Angle of Blazed Grating

Pattern Displacement

Defocus

# Optimization of Programming Parameters

$$f^\sim(x, w') \to y$$

Surrogate Model: Simple function(RBFs) to fit to experimental results

New Sampling Points for Exploring/ Exploiting a Solution

Outputs from Experiment for Refining the Surrogate Model

Experiment: Unknown function, costly to evaluate, tens of parameters

$$f(x, w) \to y$$

STEP 7

- - - TARGET FUNCTION
• • • DATA
—— PREDICTION
    $2\sigma$ CREDIBLE REGION

EFFECTIVE LIKELIHOOD

2.0

1.5

1.0

0.5

0.0

**EPFL**

---

# The Forward-Forward Algorithm: Some Preliminary Investigations

---

**Geoffrey Hinton**
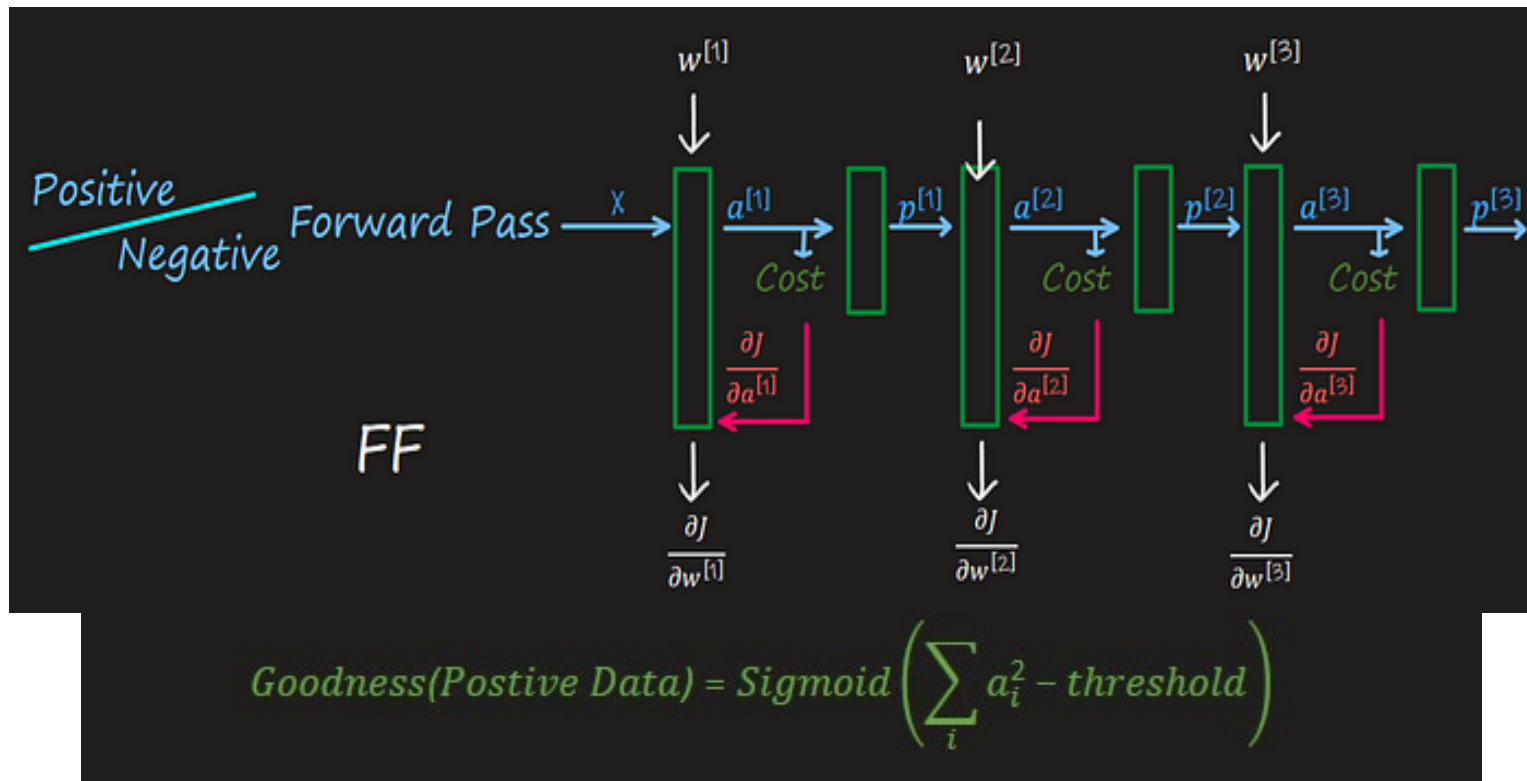Google Brain
geoffhinton@google.com

Training with FF:
for i in layerCount:

    1- optimize $W_{ij}$ such that goodness is high for positive samples and low for negative samples:

$$p(\,positive\,) = \sigma\big(\textstyle\sum_j \, y_j^2 - \theta\big), \text{ where } y_j = f(\textstyle\sum_i^N W_{ij} x_i)$$

2- Normalize $\sum_j \, y_j^2$'s for each sample before passing it to the next layer

# Forward-Forward Algorithm

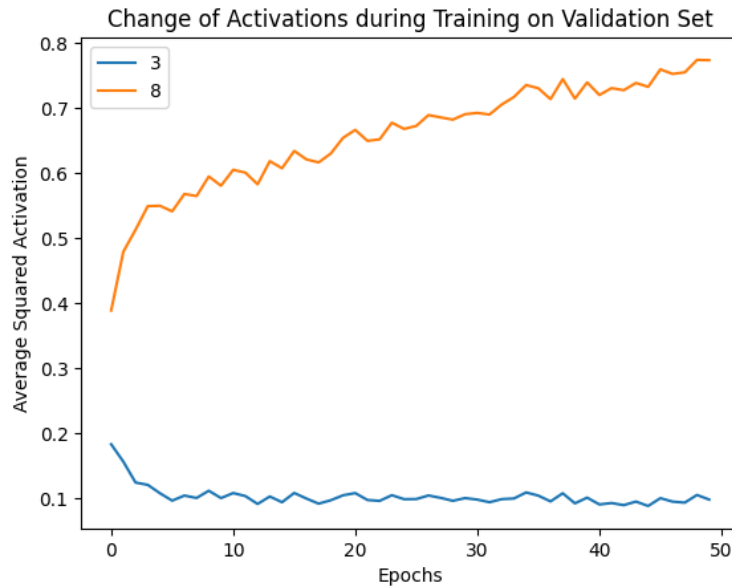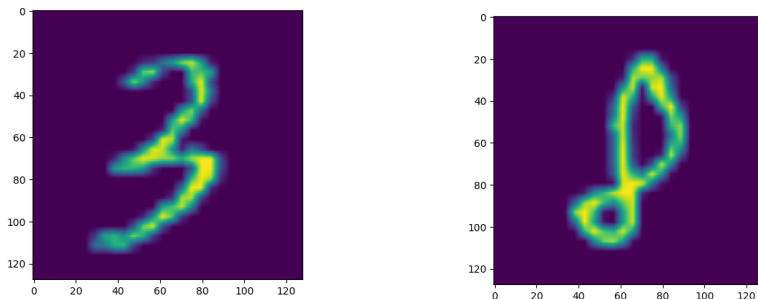$$Goodness(Postive\ Data) = Sigmoid\left(\sum_i a_i^2 - threshold\right)$$
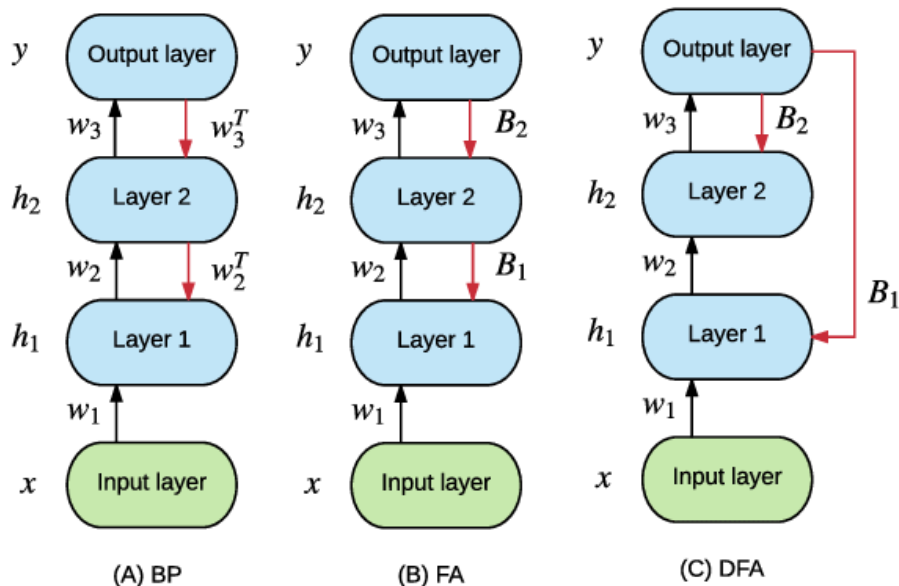
# Forward-Forward Algorithm

- The connections between layers can be unknown transforms ➔ Can be used by very-low power analog computing devices

- Each layer is updated at a time ➔ No need to save all activations ➔ Less memory consumption

- Biologically plausible

- Can be used by very-low power analog computing devices

# Forward-Forward Algorithm on MNIST

- 8's are labelled as positive and 3's are labelled as negative data

- Fully connected layer with 1024 neurons are trained FFA

- Thresholding activations: 92.7 %

- Backprop training: 96.5 %







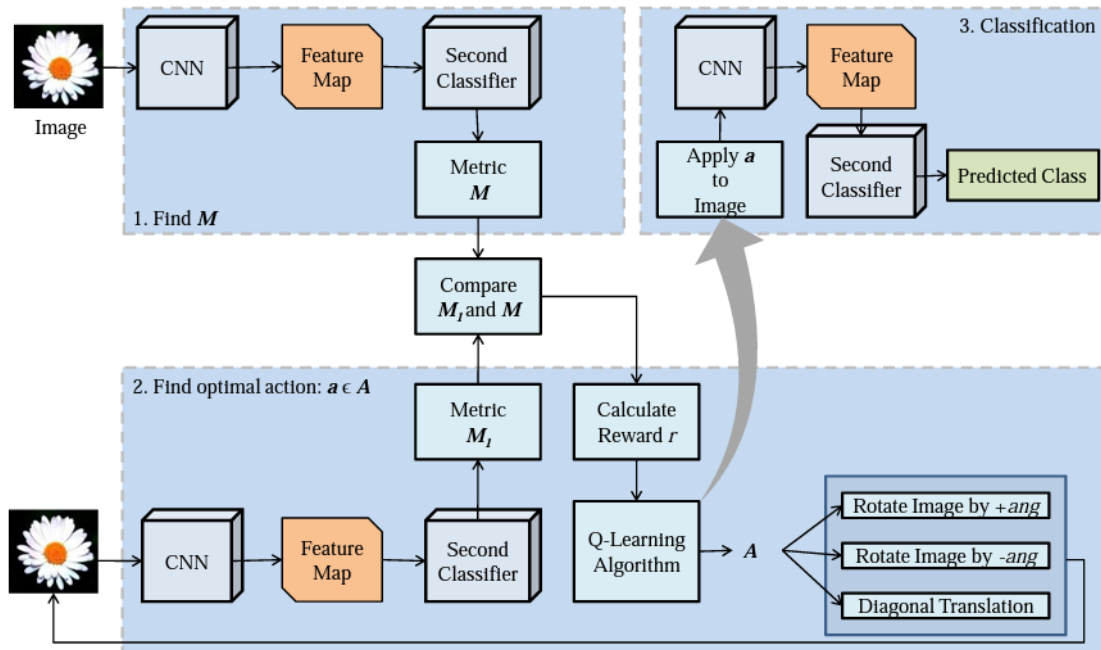Change of Activations during Training on Validation Set

# Feedback Alignment

- Instead of backpropagating errors with same weights as the forward pass, feedback alignment uses random backward connection.

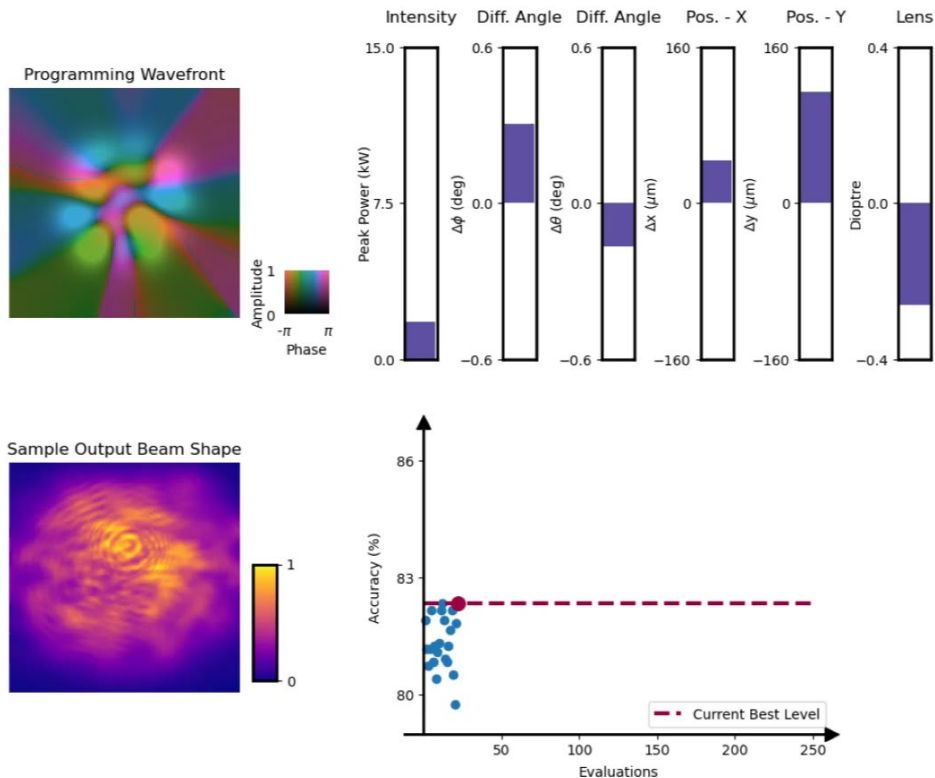- When w and B matrices are from similar distribution, convergence is observed.



(A) BP    (B) FA    (C) DFA

Error Rate

| | Dataset | |
|---|---|---|
| | MNIST | Fashion MNIST |
| BP | 0.91 | 9.20 |
| FA | 1.7 | 13.06 |
| DFA | 1.61 | 12.81 |

# An Alternative: Reinforcement Learning

Hafiz, Abdul Mueed. "Image Classification by Reinforcement Learning With Two-State Q-Learning." *Handbook of Intelligent Computing and Optimization for Sustainable Development* (2022): 171-181.

# Programming Propagation for MNIST-Fashion

Programming Wavefront

Sample Output Beam Shape

# Programming Propagation for MNIST-Fashion

1200 training and 300 test samples

Bag    Boot

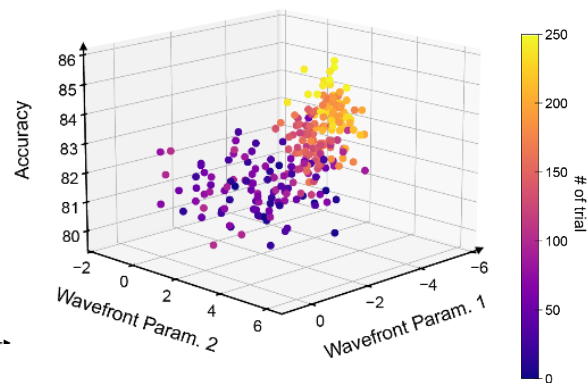$\Delta Acc = +5.0\%$

# Programming Propagation for MNIST-Fashion

Complex Field Control



Convolution with Kernel

# Programming Propagation for All-Optical Classification

# Programming Propagation for All-Optical Classification
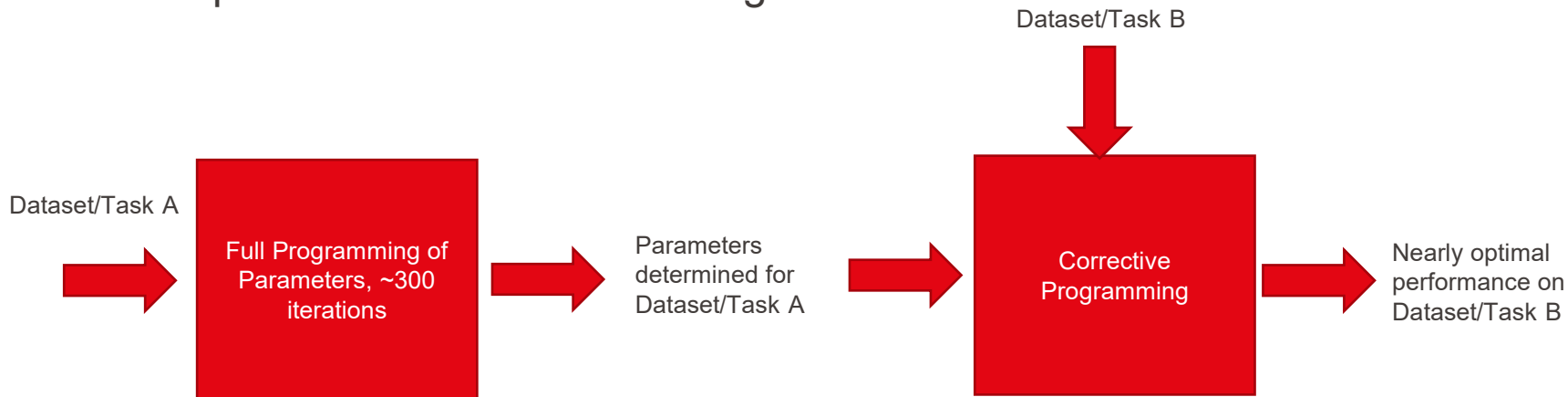
# Programming Propagation for All-Optical Classification

| Network Structure | Total Number of Parameters | Operations per Sample on Digital Computer (FLOP) | Accuracy on Melanoma dataset (%) | Accuracy on COVID-19 dataset (%) |
|---|---|---|---|---|
| LeNet-5 | 82826 | 1175640 | 64.9 | 74.6 |
| MMF + classification with output location (with programming) | 55 | 2029 | 61.3 | 77.0 |

# Transferring Programming Parameters between Different Tasks

- Current method requires ~300 iterations over the dataset for programming to converge.

- For 1500 samples at 50 images per second, full programming corresponds to ~3 hours of training.

Dataset/Task B

Dataset/Task A

Full Programming of Parameters, ~300 iterations

Parameters determined for Dataset/Task A

Corrective Programming

Nearly optimal performance on Dataset/Task B

# Transferring Programming Parameters between Different Tasks

CelebA - Gender

CelebA – Young/ Not Young

Full P... Para...

Eyeglasses

Bangs

Pointy Nose

Oval Face

Wearing Hat

Wavy Hair

Mustache

Smiling

*Test Acc...*

readout ngle l pass, version

*Test Acc*: 69.0% ⇒ 76.0%

- Training from scratch : 76.3%

# Comparison with GPU-based NNs

| Network Structure | Total Number of Parameters | Operations per Sample on Digital Computer (FLOP) | Test Accuracy on Age Task | Test Accuracy on Gender Task |
|---|---|---|---|---|
| LeNet-5 | ~82k | ~1.2M | 63.0 | 75.2 |
| 7-layer Convolutional NN | ~410k | ~65M | 65.3 | 80.1 |
| MMF + linear output layer | 2026 | 4050 | 59.0 | 69.0 |
| Programmed MMF for Age Task + linear output layer | 2078 | 6075 | 67.0 | 76.0 |
| Programmed MMF for Gender Task + linear output | 2078 | 6075 | 64.7 | 76.3 |

**Digital** (first two rows)

**Optical + Digital** (last three rows)

# THE LOTTERY TICKET HYPOTHESIS: FINDING SPARSE, TRAINABLE NEURAL NETWORKS

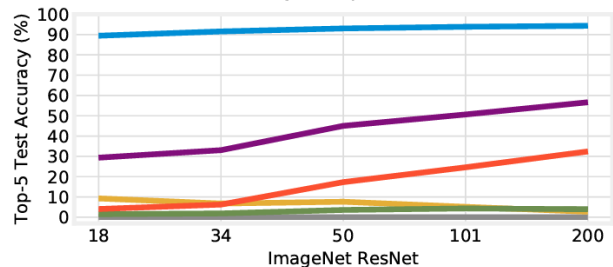**Jonathan Frankle**
MIT CSAIL
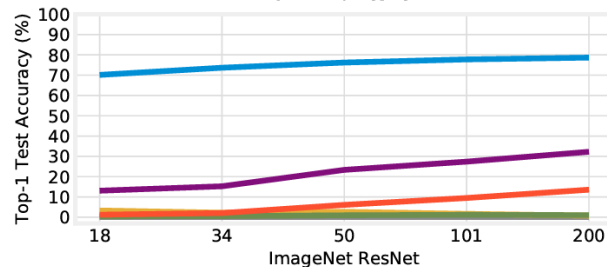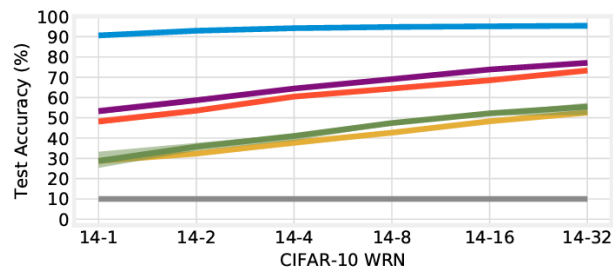
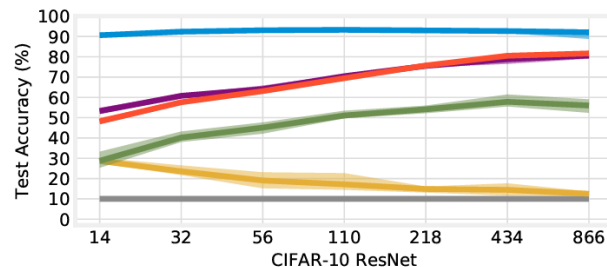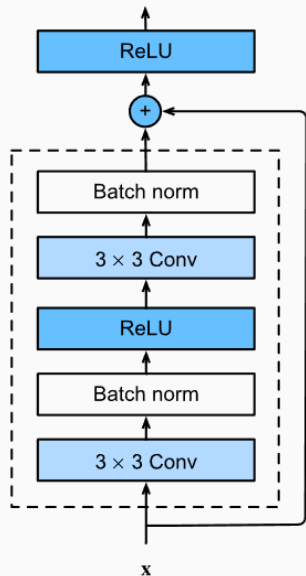**Michael Carbin**
MIT CSAIL

# TRAINING BATCHNORM AND ONLY BATCHNORM: ON THE EXPRESSIVE POWER OF RANDOM FEATURES IN CNNS

**Jonathan Frankle**[*]
MIT CSAIL
jfrankle@mit.edu

**David J. Schwab**
CUNY Graduate Center, ITS
Facebook AI Research
dschwab@fb.com

**Ari S. Morcos**
Facebook AI Research
arimorcos@fb.com