

# Beyond Backprop:

## Different Approaches to Credit Assignment in Neural Nets

Irina Rish (Mila & UdeM)

NYU



Anna  
Choromanska



Benjamin  
Cowen

IBM  
Research



Sadhana  
Kumaravel



Ronny  
Luss



Mattia  
Rigotti



Irina  
Rish

MIT



Ravi  
Tejwani



Brian  
Kingsbury



Paolo  
Di Achillele



Viatcheslav  
Gurev



Djallel  
Bouneffouf

# RETHINKING DEEP LEARNING?

Let's "backtrack" in the decision-tree of AI history:  
welcome to the "garden of forking path"(a la Borges ☺ )

## More bio-plausible neurons:

E.g., segregated dendrites  
(Bengio NIPS 2018;  
Lillicrap, Richards; etc)

"Capsules"[Hinton et al]

## More bio-plausible learning algorithms:

Target Propagation variants [LeCun, 1986],  
[Bengio 2015], [Hinton 2018]

Krotov&Hopfield [2018], Chklovski [2018]

Auxiliary-variable methods

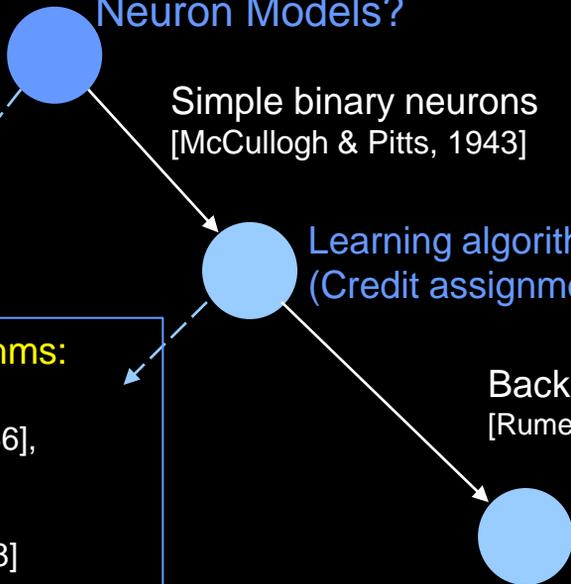
[Carreira-Perpinan 2014], [Taylor 2016],  
[Choromanska & AI Learning team 2018]

## Neuron Models?

Simple binary neurons  
[McCulloch & Pitts, 1943]

Learning algorithms?  
(Credit assignment)

Backpropagation  
[Rumelhart et al 1986]



# RETHINKING DEEP LEARNING?

“Artificial intelligence pioneer says we need to start over”  
interview with Geoff Hinton, 2017



Geoff Hinton is “**deeply suspicious of back-propagation**”, workhorse method that underlies most of the advances we are seeing in the AI field today...

“My view is **throw it all away and start again.**”



## Back Propagation

Is It The Achilles  
Heel Of Today's  
Artificial  
Intelligence



# WHAT'S WRONG WITH BACKPROP?

## Biologically implausibility:

- Error feedback does not influence neural activity, and hence does not conform to known biological feedback mechanisms underlying neural communication
- Weight transport problem: symmetric weight connectivity for feedforward and feedback directions
- Many other issues (precise clocking between feedforward and backprop phases, violation of Dale's law, etc)

## Computational Issues:

- Vanishing gradients (due to chain of derivatives)
- Difficulty handling non-differentiable nonlinearities (e.g., binary spikes)
- Lack of cross-layer parallelism

# ALTERNATIVES: PRIOR WORK

LeCun, Yann. Learning process in an asymmetric threshold network.

In Disordered systems and biological organization, 1986.



LeCun, Yann. Modeles connexionnistes de l'apprentissage. PhD thesis, Universite Paris 6, 1987.

Lee, Zhang, Fischer and Bengio. Difference target propagation. ECML-2015



Bengio, Y. How auto-encoders could provide credit assignment in deep networks via target propagation. arXiv:1407.7906, 2014.

Bartunov, S.; Santoro, A.; Richards, B. A.; Hinton, G. E.; and Lillicrap, T.

Assessing the scalability of biologically-motivated deep learning algorithms and architectures, NeurIPS 2018.



# ALTERNATIVES: PRIOR WORK

- **Offline Auxiliary-variable methods**

- MAC (Carreira-Perpiñán & Wang, 2014) and other BCD methods (Zhang & Brand, 2017; Zhang & Kleijn, 2017; Askari et al., 2018; Zeng et al., 2018; Lau et al., 2018; Gotmare et al., 2018)
- ADMM (Taylor et al., 2016; Zhang et al., 2016)
- offline (batch) is not scalable to large data and continual learning



- **Target propagation methods**

- [LeCun 1986] [Lee, Fisher, Bengio 2015] [Bartunov et al, 2018]
- Below backprop-SGD performance levels on standard benchmarks

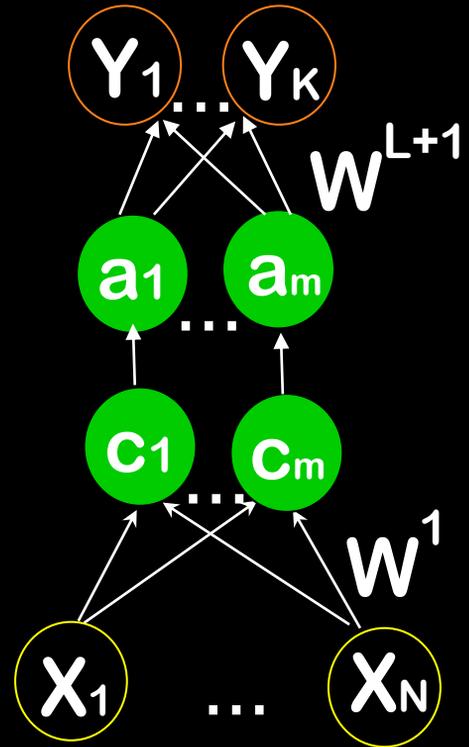
- **Proposed method:**

- Online (mini-batch, stochastic) auxiliary-variable alternating-minimization

# OUR APPROACH

## Breaking gradient chains with auxiliary activation variables:

- Relaxing nonlinear activations to noisy (Gaussian) linear activations followed by nonlinearity (e.g., ReLU)
- Alternating minimization over activations and weights: explicit activation propagation
- Weight updates are layer-local, and thus can be parallel (distributed, asynchronous)



# NEURAL NETWORK FORMULATIONS

Standard neural network objective function:

Nested

$$\min_W \mathcal{L}(y, f(\mathbf{W}, \mathbf{x}_L))$$

where  $f(\mathbf{W}, \mathbf{x}_L) = f_{L+1}(\mathbf{W}_{L+1}, f_L(\mathbf{W}_L, f_{L-1}(\mathbf{W}_{L-1}, \dots, f_1(\mathbf{W}_1, \mathbf{x}) \dots))$

Add auxiliary activation variables (hard constrained problem)

Constrained

$$\min_{\mathbf{W}, \mathbf{C}} \sum_{t=1}^n \mathcal{L}(\mathbf{y}_t, \mathbf{a}_t^L, \mathbf{W}^{L+1}), \text{ where } \mathbf{a}_t^l = \sigma_l(\mathbf{c}_t^l),$$

s.t.  $\mathbf{c}_t^l = \mathbf{W}^l \mathbf{a}_t^{l-1}, l = 1, \dots, L, \text{ and } \mathbf{a}_t^0 = \mathbf{x}_t$

Relax constraints and now amenable to alternating minimization

Relaxed

$$\min_{\mathbf{W}, \mathbf{C}} \sum_{t=1}^n \mathcal{L}(\mathbf{y}_t, \sigma_L(\mathbf{c}_t^L), \mathbf{W}^{L+1}) + \mu \sum_{t=1}^n \sum_{l=1}^L \|\mathbf{c}_t^l - \mathbf{W}^l \sigma_{l-1}(\mathbf{c}_t^{l-1})\|_2^2$$

# ONLINE ALTERNATING MINIMIZATION

Offline algorithms of prior works are not scalable to extremely large datasets and not suitable for incremental, continual/lifelong learning, hence ...

```
1: while more samples do
2:   Input  $(x_t, y_t)$ 
3:    $C \leftarrow \text{encodeInput}(x_t, W_{t-1})$  Forward: compute linear activations at layers 1,...,L
4:    $C \leftarrow \text{updateCodes}(C, y_t, W_{t-1}, \mu)$  Backward: error propagation by code changes
5:    $W_t \leftarrow \text{updateWeights}(W_{t-1}, x_t, y_t, C, \mu, \eta, Mem)$  Parallelizable
6: end while
7: return  $W_t$ 
```

Note: **updateWeights** has two options: Apply SGD to the current mini-batch or apply BCD to version that includes memory of previous samples using the following (via Mairal et al., 2009):

$$\sum_{i=1}^t \|c_i^l - W a_i^l\|_2^2 = Tr(W^T W A_t^l) - 2Tr(W^T B_t^l)$$

---

## Algorithm 2.1 Online Alternating Minimization (AM)

---

**Require:**  $(\mathbf{x}, \mathbf{y}) \sim p(\mathbf{x}, \mathbf{y})$  (data stream sampled from distribution  $p(\mathbf{x}, \mathbf{y})$ ); initial weights  $\mathbf{W}_0$ ;  $\mu \in \mathbb{R}^+$  (quadratic penalty weight);  $\eta \in \mathbb{R}^+$  (top-layer weight update step size); *Mem* (indicates the type of optimization method for **updateWeights**; if "yes", input initial memory matrices  $\mathbf{A}_0$  and  $\mathbf{B}_0$ ).

- 1: **while** more samples **do**
  - 2:   Input  $(\mathbf{x}_t, \mathbf{y}_t)$
  - 3:    $\mathbf{C} \leftarrow \mathbf{encodeInput}(\mathbf{x}_t, \mathbf{W}_{t-1})$  % forward: compute linear activations at layers 1, ...,  $L$
  - 4:    $\mathbf{C} \leftarrow \mathbf{updateCodes}(\mathbf{C}, \mathbf{y}_t, \mathbf{W}_{t-1}, \mu)$  % backward: error propagation by activation (code) changes
  - 5:    $\mathbf{W}_t \leftarrow \mathbf{updateWeights}(\mathbf{W}_{t-1}, \mathbf{x}_t, \mathbf{y}_t, \mathbf{C}, \mu, \eta, \mathit{Mem})$
  - 6: **end while**
  - 7: **return**  $\mathbf{W}_t$
-

---

**Algorithm 2.2** Activation Propagation (Code Update) Steps

---

**encodeInput**(  $\mathbf{x}$ ,  $\mathbf{W}$  )

- 1:  $\mathbf{c}^0 = \mathbf{x}$
- 2: **for**  $l = 1$  to  $L$  **do**
- 3:    $\mathbf{c}^l = \mathbf{W}^l \sigma_{l-1}(\mathbf{c}^{l-1})$   
    %  $\sigma_0(\mathbf{x}) = \mathbf{x}$ ,  $\sigma_l(\mathbf{x}) = \text{ReLU}(\mathbf{x})$  for  $l = 1, \dots, L$
- 4: **end for**
- 5: **return**  $\mathbf{C}$

**updateCodes**(  $\mathbf{C}$ ,  $\mathbf{y}$ ,  $\mathbf{W}$ ,  $\lambda_C$ ,  $\mu$  )

- 1:  $\mathbf{c}^L \leftarrow$  Solve Problem (4),  $\mathbf{c}^0 = \mathbf{x}$
  - 2: **for**  $l = L - 1$  to 1 **do**
  - 3:    $\mathbf{c}^l \leftarrow$  Solve Problem (5)
  - 4: **end for**
  - 5: **return**  $\mathbf{C}$
-

---

**Algorithm 2.3** Weight and Memory Update Steps

---

**updateWeights**(  $W$ ,  $x$ ,  $y$ ,  $C$ ,  $\mu$ ,  $\eta$ ,  $Mem$ )

```
1:  $W^{L+1} = W^{L+1} - \eta \nabla_W \mathcal{L}(y, \sigma_L(c^L), W^{L+1})$ 
2: for  $l = 1$  to  $L$  do
3:   if  $Mem$  then
4:      $(A^l_t, B^l_t) \leftarrow \text{updateMemory}(A^l_{t-1}, B^l_{t-1}, C^l)$ 
5:     %  $\hat{f}^l(W^l) \equiv Tr(W^T W A^l) - 2Tr(W^T B^l)$ 
6:      $W^l = \arg \min_W \hat{f}^l(W)$ 
7:   else
8:     % (parallel) local update of each layer weights,
9:     % (independently of other layers (unlike backprop))
10:     $W^l \leftarrow \text{SGD}(W^l, x, y, C^l, \mu, \eta)$ 
11:   end if
12: end for
13: return  $W$ 
```

**updateMemory**(  $A$ ,  $B$ ,  $C$ )

```
1: for  $l = 1$  to  $L$  do
2:    $a = \sigma_{l-1}(c^{l-1})$ ,  $A^l \leftarrow A^l + a a^T$ ,  $B^l \leftarrow B^l + c^l a^T$ 
3: end for
4: return  $A, B$ 
```

---

# NOVEL THEORETIC RESULT FOR STOCHASTIC ALTMIN

**Theorem 3.2.** *Given the stochastic AM gradient iterates of the version of Algorithm 2.1 given in eq. 8 with decaying step size  $\eta^t = \frac{3/2}{[2\xi - 3\gamma(K-1)](t+2) + \frac{3}{2}(K-1)\gamma}$  and assuming that  $\gamma < \frac{2\xi}{3(K-1)}$ , the error at iteration  $t + 1$  satisfies*

$$\mathbb{E} \left[ \sum_{d=1}^K \|\Delta_d^{t+1}\|_2^2 \right] \leq \mathbb{E} \left[ \sum_{d=1}^K \|\Delta_d^0\|_2^2 \right] \left( \frac{2}{t+3} \right)^{\frac{3}{2}} + \sigma^2 \frac{9}{[2\xi - 3\gamma(K-1)]^2(t+3)}, \quad (10)$$

where  $\Delta_d^{t+1} := \theta_d^{t+1} - \theta_d^*$  for  $d = 1, 2, \dots, K$ ,  $\gamma := \max_{i=1,2,\dots,K} \gamma_i$ , and  $\xi := \min_{i=1,2,\dots,K} \frac{2\mu_i \lambda_i}{\mu_i + \lambda_i}$ .

# FULLY-CONNECTED NETS

AM greatly **outperforms all off-line methods** (ADMM of Taylor et al, and offline AM), and often matches Adam and SGD (50 epochs)

## MNIST

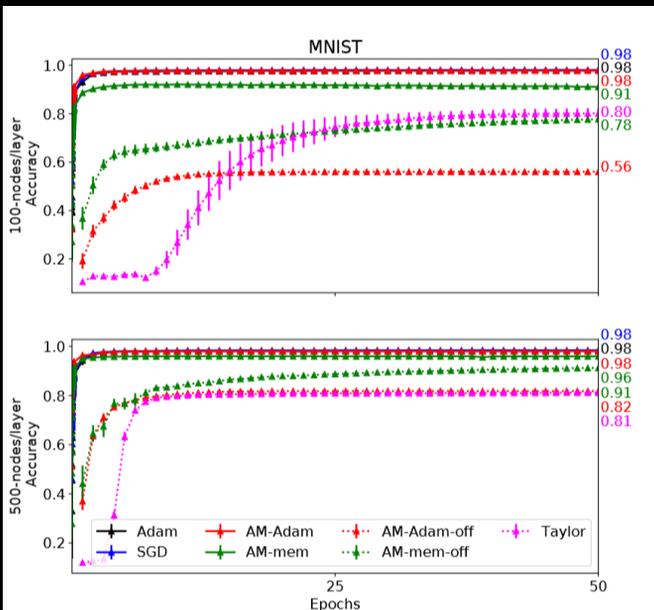


Figure 2. MNIST (fully-connected nets, 2 layers): online vs. of-line methods vs. Taylor's ADMM, 50 epochs.

## CIFAR-10

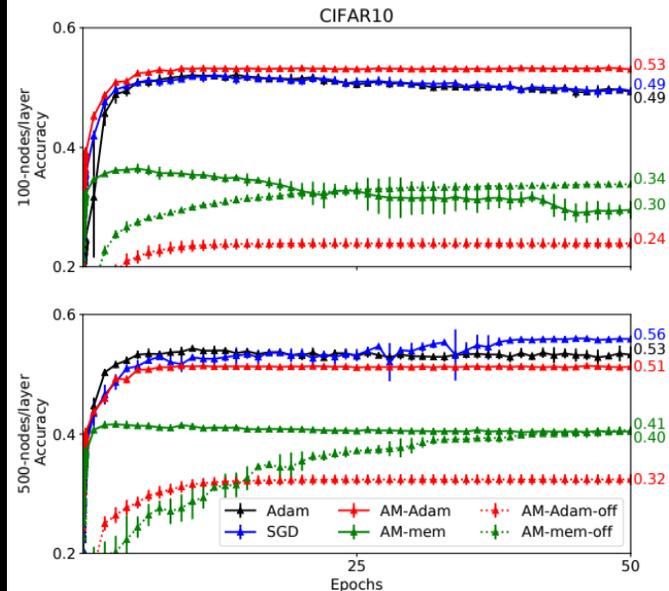
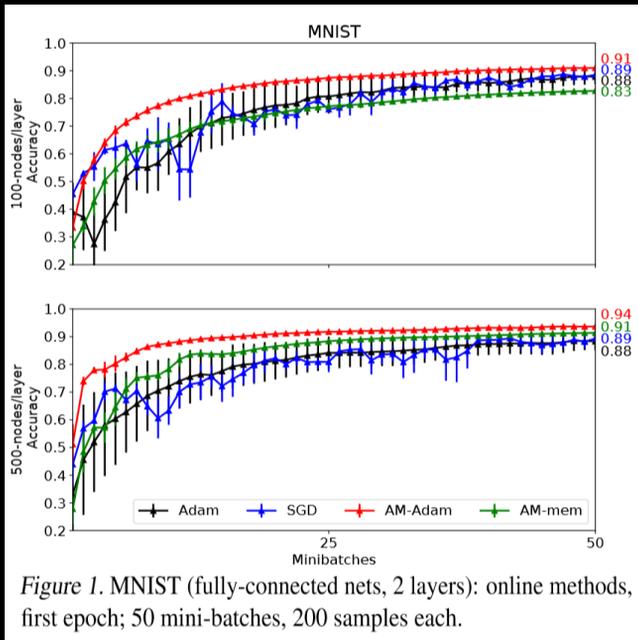


Figure 4. CIFAR10 (fully-connected networks): online vs. offline, 50 epochs. Similar experiments to Figure 2.

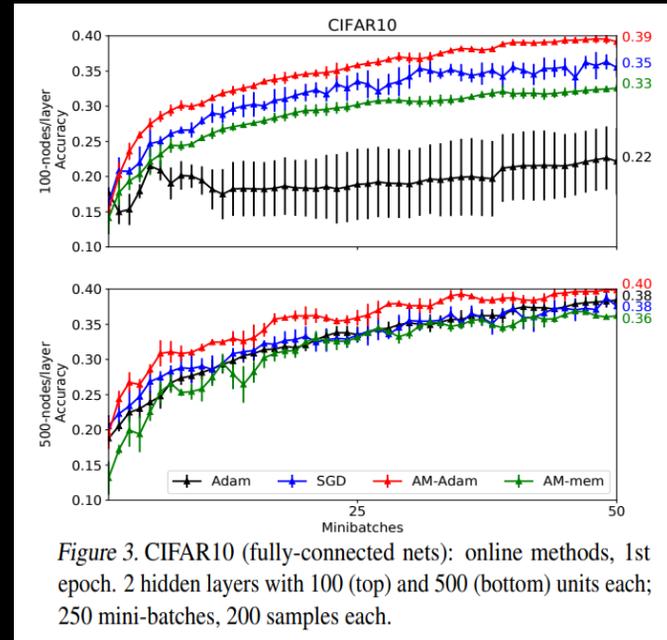
# FASTER INITIAL LEARNING: POTENTIAL USE AS A GOOD INIT?

- AM often **learns faster than SGD & Adam (backprop-based)** in the 1<sup>st</sup> epoch, then **matches** their performance

## MNIST



## CIFAR-10



# CONVNETS: LENET5, MNIST

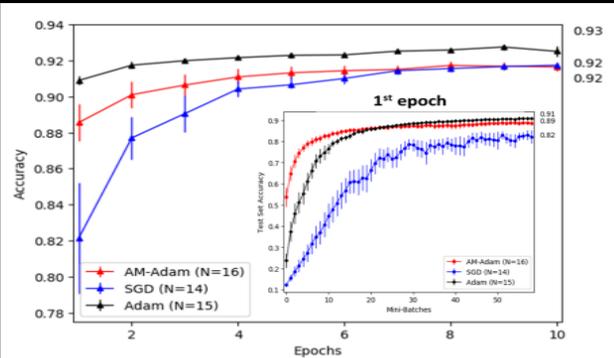


Figure 6. RNN-15, Sequential MNIST.

# RNN: SEQUENTIAL MNIST

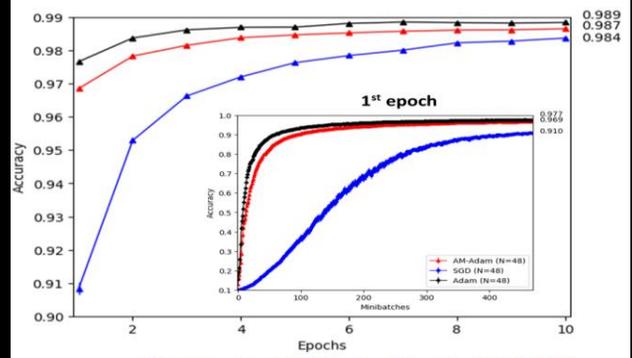


Figure 7. CNN: LeNet5, MNIST.

# HIGGS DATASET, FULLY-CONNECTED

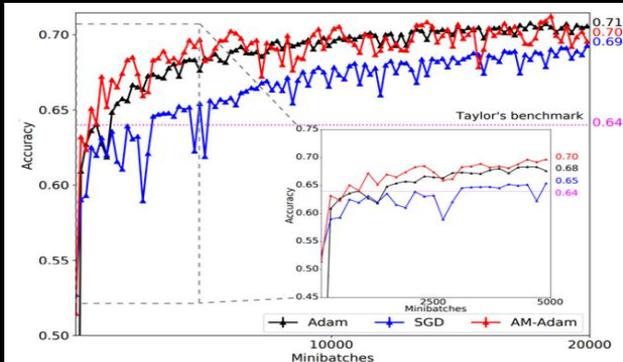


Figure 5. HIGGS dataset.

- AM performs similarly to Adam, outperforms SGD
- All methods greatly outperform offline ADMM (Taylor's 0.64 benchmark) using **less than 0.01%** of 10.5M-sample HIGGS data

# NONDIFFERENTIABLE (BINARY) NETS

- Backprop replaced by Straight-Through Estimator (STE)
- Comparing with Difference Target Propagation (DTP)
- DTP took about 200 epochs to reach 0.2 error, matching the STE performance (Lee et al., 2015)
- AM-Adam with binary activations reaches same error in < than 20 epochs

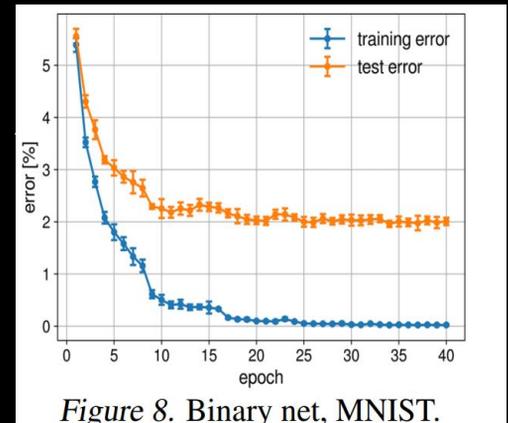


Figure 8. Binary net, MNIST.

# SUMMARY: CONTRIBUTIONS

- **Algorithm(s):** novel online (stochastic) auxiliary-variable approach for training neural networks (prior methods are offline/batch); two versions of the approach (memory-based and local-SGD-based)
- **Theory:** first general theoretical convergence guarantees for alternating minimization in the stochastic setting: the error decays at the sub-linear rate  $O((1/t)^{3/2} + 1/t)$  in  $t$  iterations
- **Extensive Evaluations:** variety of architectures and datasets demonstrating advantages of online vs offline approaches and performance similar to SGD (Adam), with faster initial convergence