

Associated learning

Hung-Hsuan Chen

1/11/21

1

1

Review: backpropagation

- A method to efficiently compute the gradient of the loss with respect to the parameters
 - Chain rule
 - Layer-wise update

1/11/21

2

2

Backward locking

- Since the parameters in a DNN are updated in a layer-wise fashion:
 - Updating the weights of one layer *locks* the weight updates in the other layers
- Challenging to apply parallel computing or a pipeline structure to update the weights in different layers simultaneously

1/11/21

3

3

Time complexity of training a DNN

- Training accuracy:
 - n : number of training instances
 - ℓ : number of hidden layers.
 - Time complexity: $O(n\ell)$
- As a network becomes deeper, more training time are needed

1/11/21

4

4

Can we “parallelize” the training?

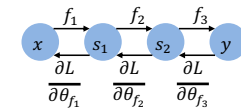
- Two types of parallelizations
- Type 1
 - Compute SGD at different cores using different subsamples
 - Average the gradients of different cores, and apply parameter updates
- Type 2
 - Update parameters of different layers simultaneously
 - Seems more challenging
 - We focus on this type

1/11/21

5

5

Naively using SGD and backpropagation



	t = 1	t = 2	t = 3	t = 4	t = 5	t = 6	t = 7	t = 8	t = 9	t = 10	t = 11	t = 12
$(x^{(1)}, y^{(1)})$	$f_1^{(0)}$	$f_2^{(0)}$	$f_3^{(0)}$	$\frac{\partial L}{\partial \theta_{f_3}}$	$\frac{\partial L}{\partial \theta_{f_2}}$	$\frac{\partial L}{\partial \theta_{f_1}}$						
$(x^{(2)}, y^{(2)})$							$f_1^{(1)}$	$f_2^{(1)}$	$f_3^{(1)}$	$\frac{\partial L}{\partial \theta_{f_3}}$	$\frac{\partial L}{\partial \theta_{f_2}}$	$\frac{\partial L}{\partial \theta_{f_1}}$

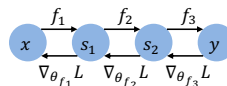
- Notation
 - f_i : function at layer i
 - $f_i^{(j)}$: parameter values for layer i after j updates
 - θ_{f_i} : parameters for layer i
- Two hidden layers
 - Updating all parameters once requires 6 time units
- If ℓ hidden layers
 - Updating all parameters once requires $2(\ell + 1)$ time units
- Training time: $O(\ell)$
 - Training time is long for “deep” network

1/11/21

6

6

Add pipeline directly



	t = 1	t = 2	t = 3	t = 4	t = 5	t = 6	t = 7	t = 8	t = 9	t = 10	t = 11	t = 12
$(x^{(1)}, y^{(1)})$	$f_1^{(0)}$	$f_2^{(0)}$	$f_3^{(0)}$	$\frac{\partial L}{\partial \theta_{f_3}}$	$\frac{\partial L}{\partial \theta_{f_2}}$	$\frac{\partial L}{\partial \theta_{f_1}}$						
$(x^{(2)}, y^{(2)})$		$f_1^{(0)}$	$f_2^{(0)}$	$f_3^{(0)}$	$\frac{\partial L}{\partial \theta_{f_3}}$	$\frac{\partial L}{\partial \theta_{f_2}}$	$\frac{\partial L}{\partial \theta_{f_1}}$					
$(x^{(3)}, y^{(3)})$			$f_1^{(0)}$	$f_2^{(0)}$	$f_3^{(1)}$	$\frac{\partial L}{\partial \theta_{f_3}}$	$\frac{\partial L}{\partial \theta_{f_2}}$	$\frac{\partial L}{\partial \theta_{f_1}}$				
$(x^{(4)}, y^{(4)})$				$f_1^{(0)}$	$f_2^{(2)}$	$f_3^{(2)}$	$\frac{\partial L}{\partial \theta_{f_3}}$	$\frac{\partial L}{\partial \theta_{f_2}}$	$\frac{\partial L}{\partial \theta_{f_1}}$			
$(x^{(5)}, y^{(5)})$					$f_1^{(0)}$	$f_2^{(1)}$	$f_3^{(3)}$	$\frac{\partial L}{\partial \theta_{f_3}}$	$\frac{\partial L}{\partial \theta_{f_2}}$	$\frac{\partial L}{\partial \theta_{f_1}}$		

- If ℓ hidden layers and fully pipelined
 - Starting from $2(\ell + 1)$ time unit, all parameters are updated once for every time unit
- However, the updating is wrong

– E.g., at $t = 6$, $\nabla_{\theta_{f_1}} L = \frac{\partial L}{\partial \theta_{f_3}} \Big|_{f_3=f_3^{(2)}} \cdot \frac{\partial f_3}{\partial \theta_{f_2}} \Big|_{f_2=f_2^{(1)}} \cdot \frac{\partial f_2}{\partial \theta_{f_1}} \Big|_{f_1=f_1^{(0)}}$

1/11/21

7

7

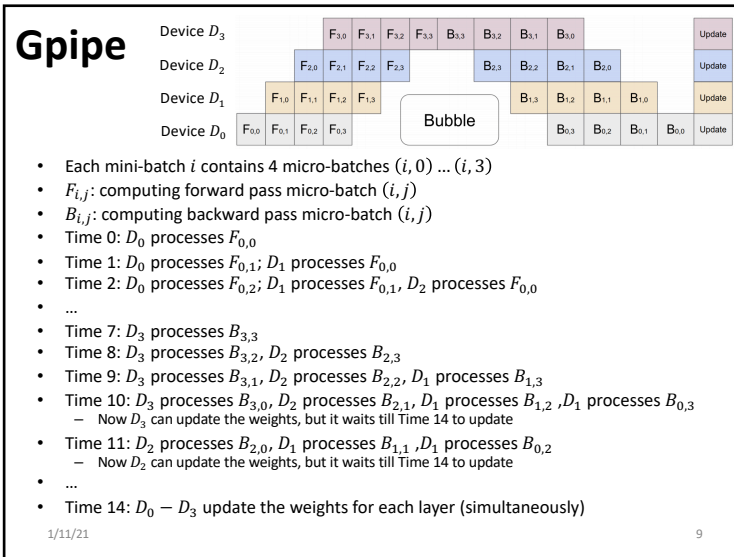
How to pipeline?

- Gpipe
 - Yanping Huang et al.
 - NeurIPS 2019
- Associated learning
 - Yu-Wei Kao and Hung-Hsuan Chen
 - Neural Computation 33(1) 2021.

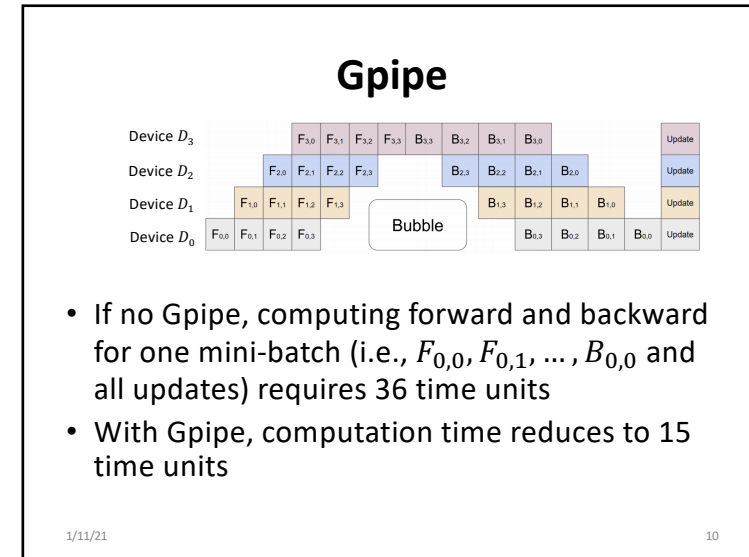
1/11/21

8

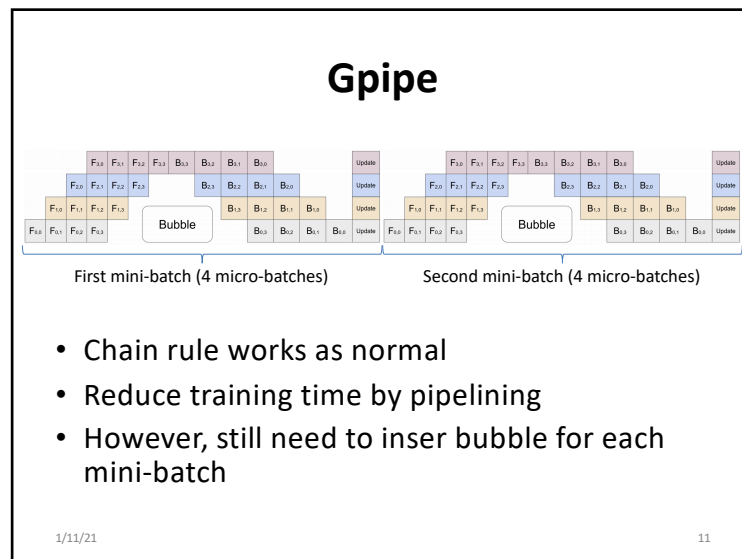
8



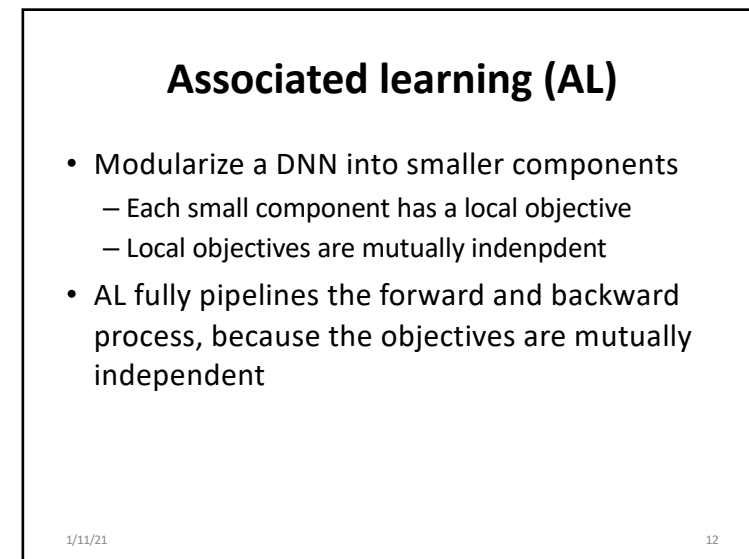
9



10

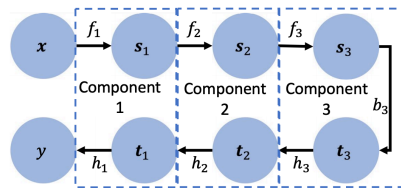
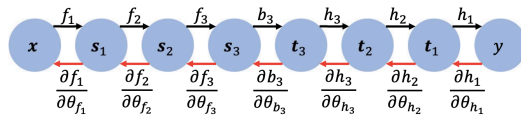


11



12

Modularize DNN (e.g., 6 hidden layers)



1/11/21

13

13

How to pipeline

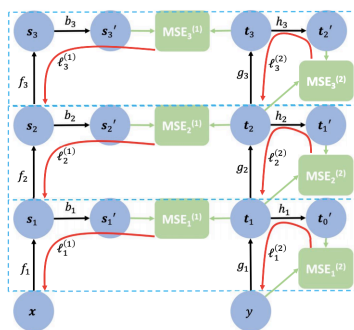
Time unit	1	2	3	4	5	6	7	...
1 st mini-batch	Task 1	Task 2	Task 3					
2 nd mini-batch		Task 1	Task 2	Task 3				
3 rd mini-batch			Task 1	Task 2	Task 3			
4 th mini-batch				Task 1	Task 2	Task 3		
5 th mini-batch					Task 1	Task 2	Task 3	
...								

1/11/21

14

14

Associated learning – procedure



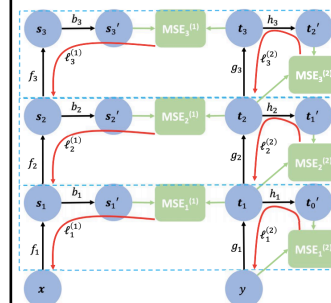
- Black arrow: forward
- Red arrow: backward gradient flow
- Green arrow: connect the variables that should be compared to minimize their associated distance

1/11/21

15

15

Associated learning – procedure



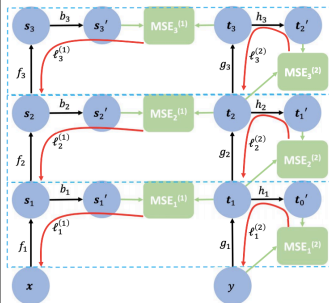
- Each component contains 2 local forward functions f_i and g_i
 - $s_i = f_i(s_{i-1})$
 - s_0 is x
 - $t_i = g_i(t_{i-1})$
 - t_0 is y
- Initially, we want to find f_i and g_i such that $s_i \approx t_i$ (local objective)
 - Turn out challenging
- Add “bridge function” b_i : find f_i, g_i, b_i such that
 - $s'_i = b_i(s_i) \approx t_i$

1/11/21

16

16

Associated learning – procedure

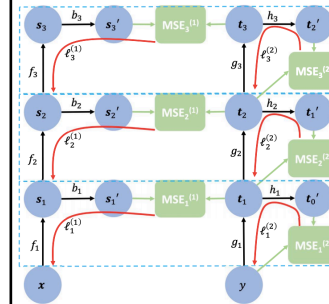


1/11/21

17

- The above training allows finding f_1, f_2, f_3, b_3 , so we can perform $b_3(f_3(f_2(f_1(x))))$
- We assume $s'_3 \approx t_3$, so $t_3 \approx b_3(f_3(f_2(f_1(x))))$
- How to do the following transform?
 $t_3 \rightarrow t_2 \rightarrow t_1 \rightarrow t_0 (= y)$

Associated learning – procedure



1/11/21

18

- Look for h_i such that $h_i(g_i(t_{i-1})) \approx t'_{i-1}$
- Regarded as an autoencoder
 $t_{i-1} \xrightarrow{g_i} t_i \xrightarrow{h_i} t'_{i-1}$
- If we can find such h_i s, then
 $t_{i-1} \approx h_i(t_i)$
- So we know how to
 $t_3 \rightarrow t_2 \rightarrow t_1 \rightarrow t_0$

- Prediction function

$$y = h_1 \left(h_2 \left(h_3 \left(b_3 \left(f_3 \left(f_2 \left(f_1(x) \right) \right) \right) \right) \right) \right)$$

Associated learning

- Prediction function

$$y = h_1 \left(h_2 \left(h_3 \left(b_3 \left(f_3 \left(f_2 \left(f_1(x) \right) \right) \right) \right) \right) \right)$$

backward function
Forward function

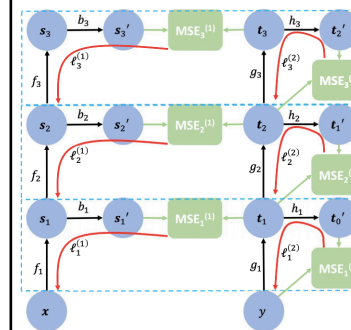
The last bridge function

1/11/21

19

19

Effective and affiliated parameters



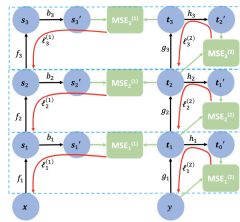
1/11/21

20

20

- During training, the model learns
– f_i s, b_i s, g_i s, h_i s
- At prediction, only
– f_i s, h_i s and last b_i are used
- We call f_i s, h_i s and last b_i effective parameters
- The others are affiliated parameters

Associated learning and pipeline



	$t = 1$	$t = 2$	$t = 3$	$t = 4$	$t = 5$	$t = 6$
$(x^{(1)}, y^{(1)})$	$f_1^{(0)}, b_1^{(0)}, g_1^{(0)}, h_1^{(0)}$	$f_2^{(0)}, b_2^{(0)}, g_2^{(0)}, h_2^{(0)}$	$f_3^{(0)}, b_3^{(0)}, g_3^{(0)}, h_3^{(0)}$			
$(x^{(2)}, y^{(2)})$		$f_1^{(1)}, b_1^{(1)}, g_1^{(1)}, h_1^{(1)}$	$f_2^{(1)}, b_2^{(1)}, g_2^{(1)}, h_2^{(1)}$	$f_3^{(1)}, b_3^{(1)}, g_3^{(1)}, h_3^{(1)}$		
$(x^{(3)}, y^{(3)})$			$f_1^{(2)}, b_1^{(2)}, g_1^{(2)}, h_1^{(2)}$	$f_2^{(2)}, b_2^{(2)}, g_2^{(2)}, h_2^{(2)}$	$f_3^{(2)}, b_3^{(2)}, g_3^{(2)}, h_3^{(2)}$	
$(x^{(4)}, y^{(4)})$				$f_1^{(3)}, b_1^{(3)}, g_1^{(3)}, h_1^{(3)}$	$f_2^{(3)}, b_2^{(3)}, g_2^{(3)}, h_2^{(3)}$	$f_3^{(3)}, b_3^{(3)}, g_3^{(3)}, h_3^{(3)}$
$(x^{(5)}, y^{(5)})$					$f_1^{(4)}, b_1^{(4)}, g_1^{(4)}, h_1^{(4)}$	$f_2^{(4)}, b_2^{(4)}, g_2^{(4)}, h_2^{(4)}$
$(x^{(5)}, y^{(5)})$						$f_1^{(5)}, b_1^{(5)}, g_1^{(5)}, h_1^{(5)}$

1/11/21

21

21

Experimental settings

- Dataset
 - MNIST
 - CIFAR-10
 - CIFAR-100
- NN structure
 - MLP
 - Vanilla CNN
 - VGG Net
 - ResNet-20
 - ResNet-32

1/11/21

22

22

Experimental results on MNIST

	BP	AL
MLP	98.5 ± 0.0%	98.6 ± 0.0%
Vanilla CNN	99.4 ± 0.0%	99.5 ± 0.0%

- MLP: 5 hidden layers
 - 1024 x 1024 x 5120 x 1024 x 1024
- Vanilla CNN: 13 hidden layers
 - First four layers: (3 x 3 x 32) kernels
 - Next four layers: (3 x 3 x 64) kernels
 - Last five layers: fully connected 1280 x 256 x 256 x 256 x 256
- We didn't try more complex networks, as MLP and Vanilla CNN are extremely good

1/11/21

23

23

Experimental results on CIFAR-10

	BP	AL
MLP	60.6 ± 0.3%	62.8 ± 0.2%
Vanilla CNN	85.2 ± 0.4%	85.8 ± 0.1%
ResNet-20	91.2 ± 0.4%	89.1 ± 0.5%
ResNet-32	92.0 ± 0.2%	88.7 ± 0.4%
VGG	92.3 ± 0.2%	92.6 ± 0.1%

1/11/21

24

24

Experimental results on CIFAR-100

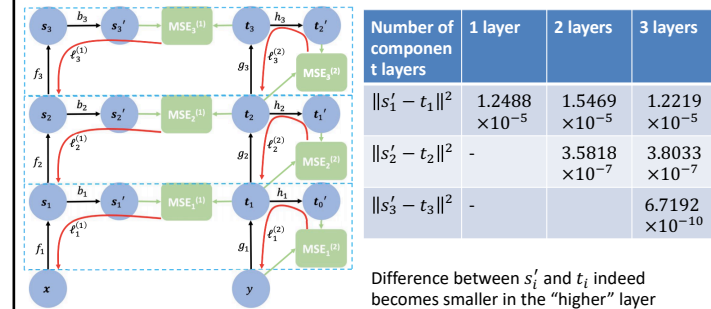
	BP	AL
MLP	$26.5 \pm 0.4\%$	$29.7 \pm 0.2\%$
Vanilla CNN	$51.1 \pm 0.2\%$	$52.2 \pm 0.5\%$
ResNet-20	$63.7 \pm 0.2\%$	$61.0 \pm 0.6\%$
ResNet-32	$63.7 \pm 0.3\%$	$59.0 \pm 1.6\%$
VGG	$65.8 \pm 0.3\%$	$67.1 \pm 0.3\%$

1/11/21

25

25

Associated loss at different layers on MNIST after 200 epochs



1/11/21

26

26

Number of layers vs training/test accuracies

Number of component layers	1 layer	2 layers	3 layers
Training accuracy	1.0	1.0	1.0
Test accuracy	0.9849	0.9860	0.9871

Increasing layers increases test accuracies

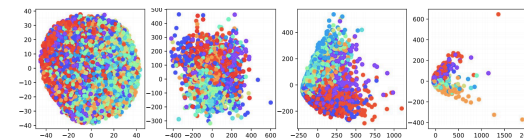
1/11/21

27

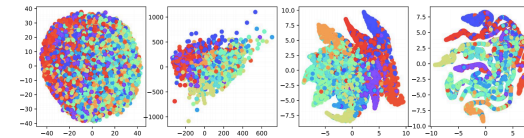
27

t-SNE visualization of MLP on CIFAR-10

- Raw data, 2nd layer, 4th layer, and output layer when training by BP



- Raw data, 2nd layer, 4th layer, and output layer when training by AL



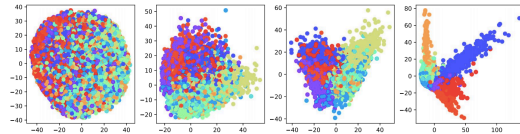
1/11/21

28

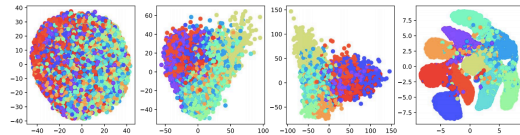
28

t-SNE visualization of CNN on CIFAR-10

- Raw data, 4th layer, 8th layer, and 12th layer when training by BP



- Raw data, 4th layer, 8th layer, and 12th layer when training by AL



1/11/21

29

29

Inter- and intra-distances and ratios of the two

Dataset	Network	Method	Interclass distance	Intraclass distance	Inter:Intra ratio
CIFAR-10	MLP	BP	39.36	67.97	0.58
		AL	0.73	0.66	1.11
	Vanilla CNN	BP	41.82	26.87	1.56
		AL	1.17	0.36	3.25
CIFAR-100	MLP	BP	114.42	342.65	0.33
		AL	0.23	0.28	0.82
	Vanilla CNN	BP	114.71	163.43	0.70
		AL	0.55	0.51	1.08

1/11/21

30

30

Summary

- Propose a method, AL, to decompose end-to-end backpropagation into many short gradient flows
- Fully pipelined the learning
- Experiments of multiple network structures on multiple datasets show that AL is comparable (and sometimes better) than BP

1/11/21

31

31

Discussion

- We implemented AL, but we didn't implement a "pipelined" AL
 - How to implement pipeline in Tensorflow or PyTorch?
- Most components in AL do not receive messages directly from the target, but AL performs comparable to BP.
 - This is a nice surprise, but we don't know why
 - If we know why AL works, perhaps we can design better network structure?
- Does AL converge? Does the loss function converge to minima?
- Experiments on other (larger) datasets?
- If you are interested in any of the above topics, please let me know

1/11/21

32

32

More information

- Associated Learning: Decomposing End-to-End Backpropagation Based on Autoencoders and Target Propagation
 - Yu-Wei Kao (高聿緯), Hung-Hsuan Chen (陳弘軒)
 - MIT Neural Computation 33(1), 2021
 - <https://arxiv.org/pdf/1906.05560.pdf>

1/11/21

33

33

Quiz

- What is backward locking?
- Compare pipeline and parallel computing
- What is the key concept of associated learning?

1/11/21

34

34