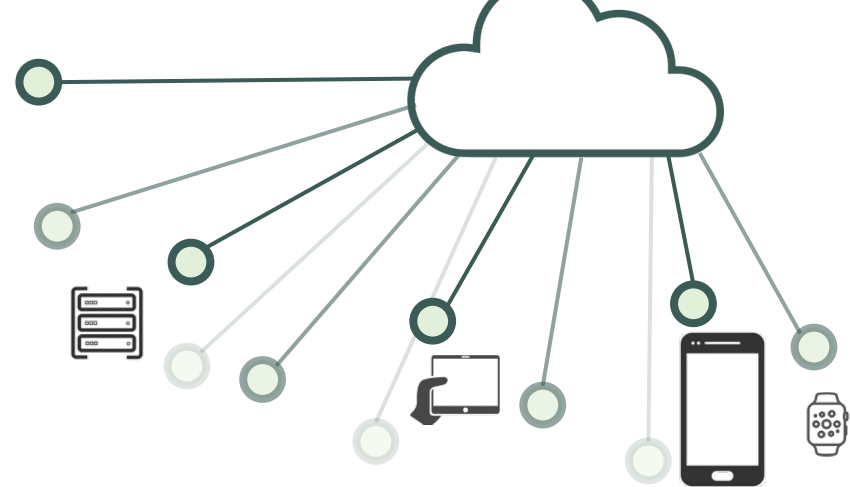




COMPUTER
VISION
TALKS



MECTA: Memory-Economic Continual Test-Time Model Adaptation

Junyuan Hong^{1*}, Lingjuan Lyu², Jiayu Zhou¹, Michael Spranger²

¹Michigan State University, ²Sony AI

*Work done during internship at Sony AI



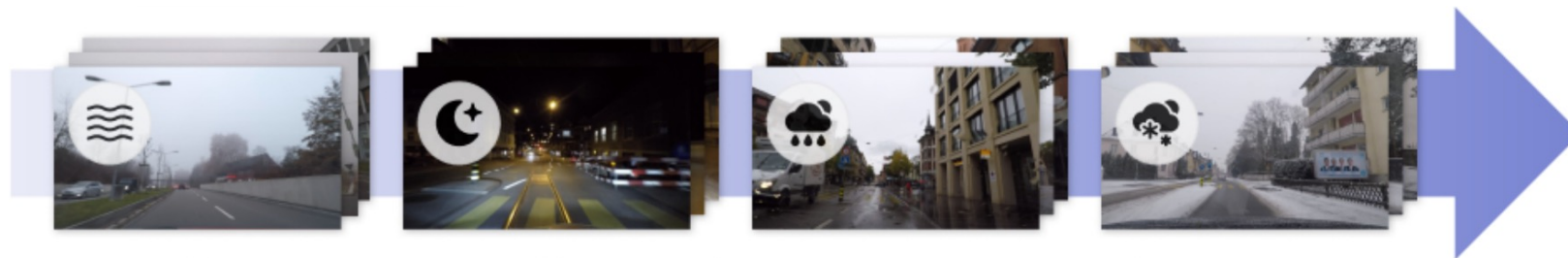
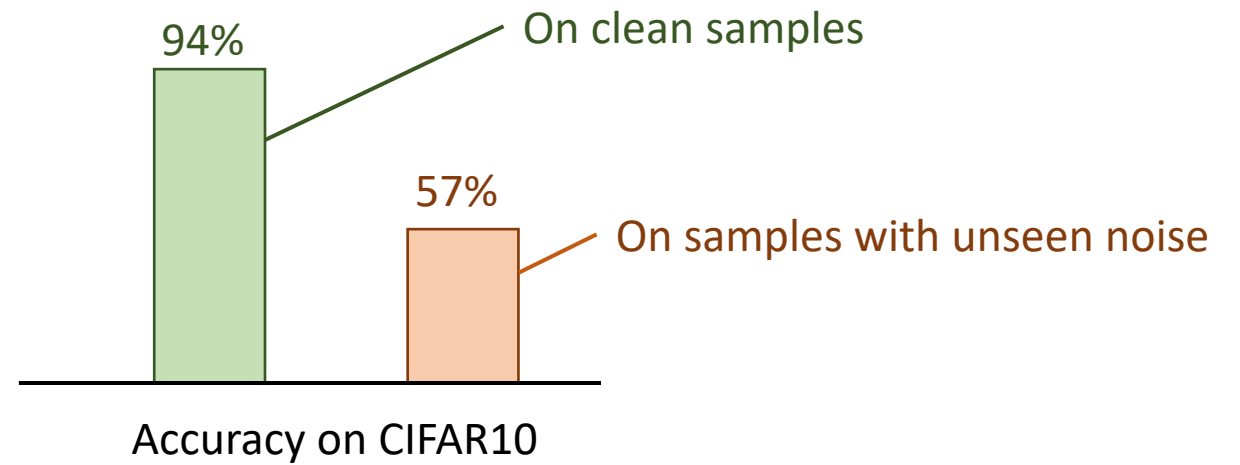
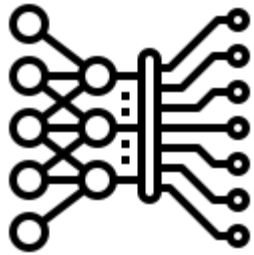
ICLR
International Conference On
Learning Representations

Sony AI

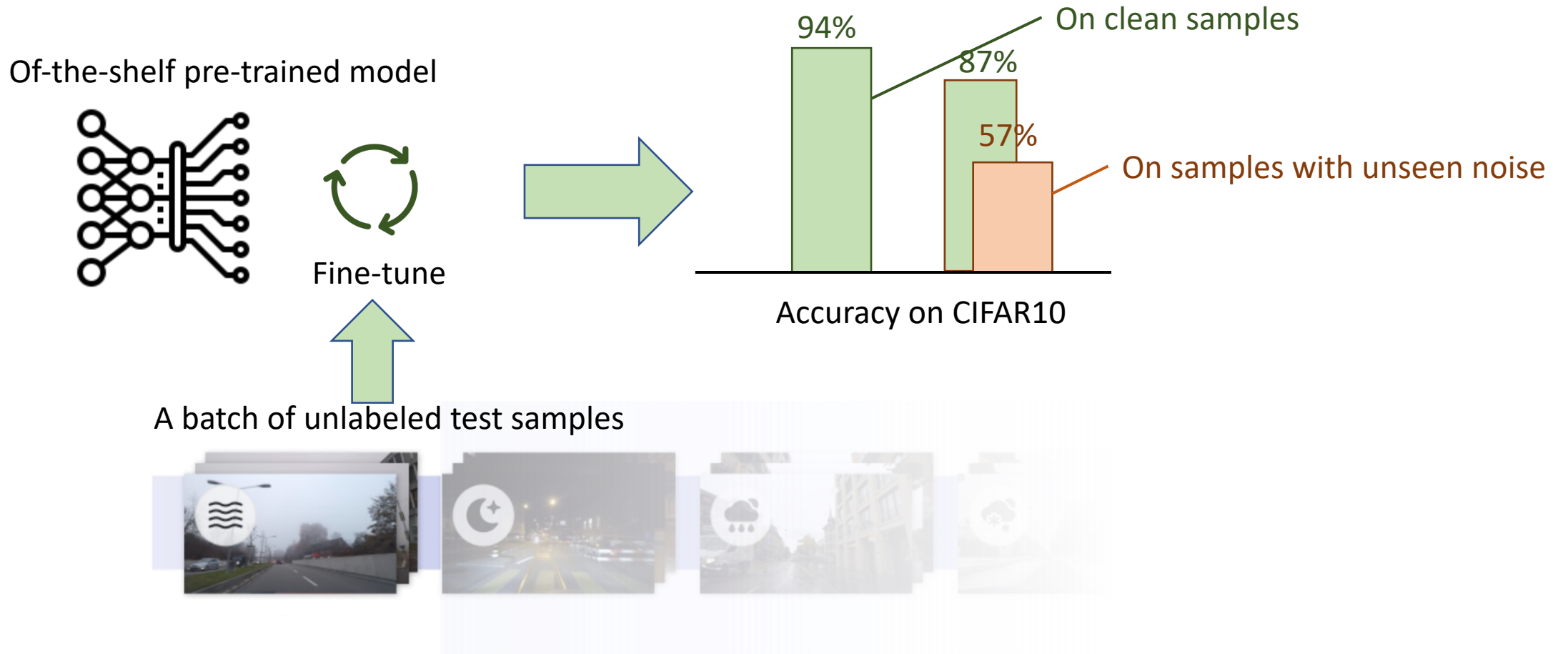


Continually Changing Environments

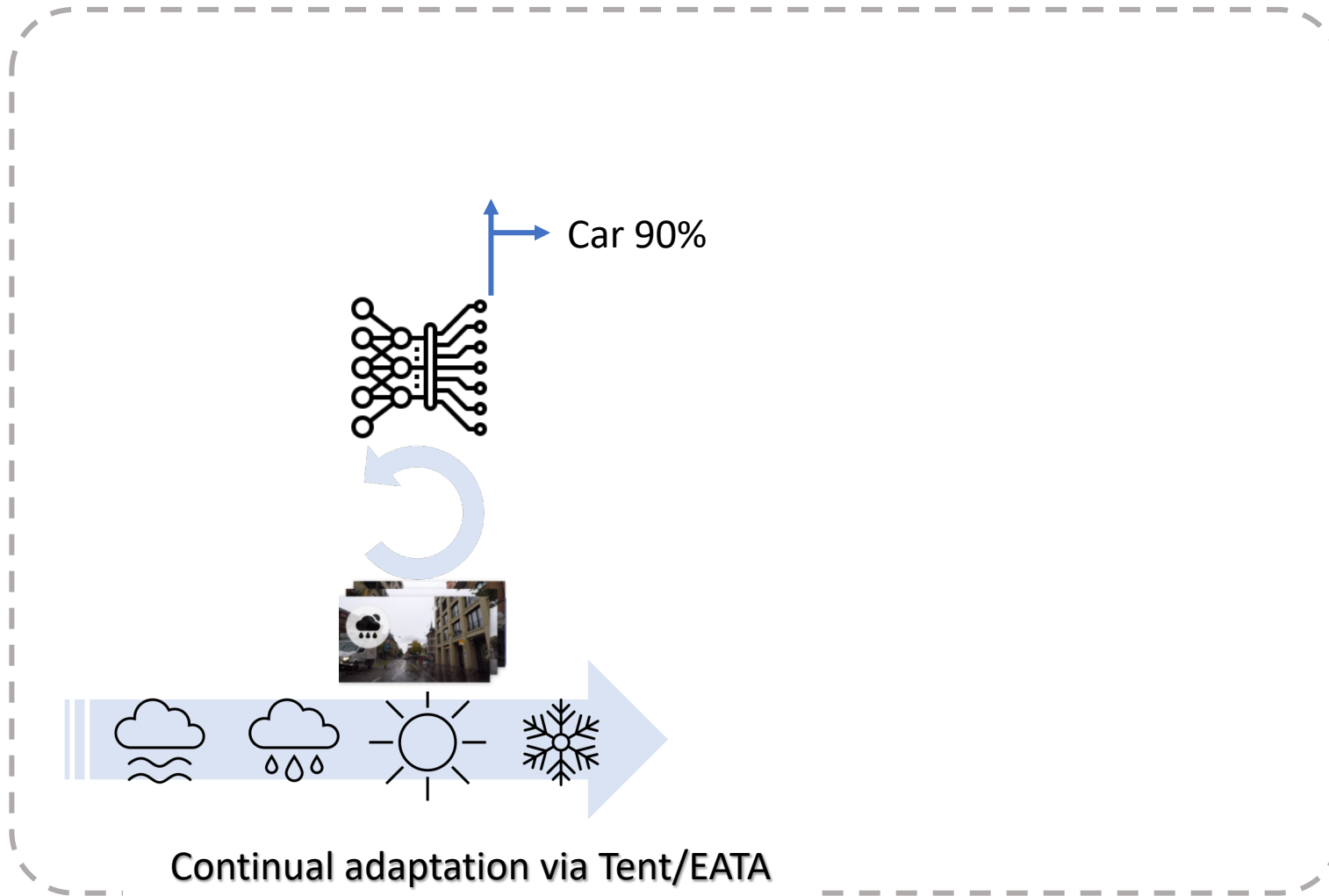
Of-the-shelf pre-trained model



Continual Test-time Adaptation (CTA)



Continual Test-time Adaptation (CTA)

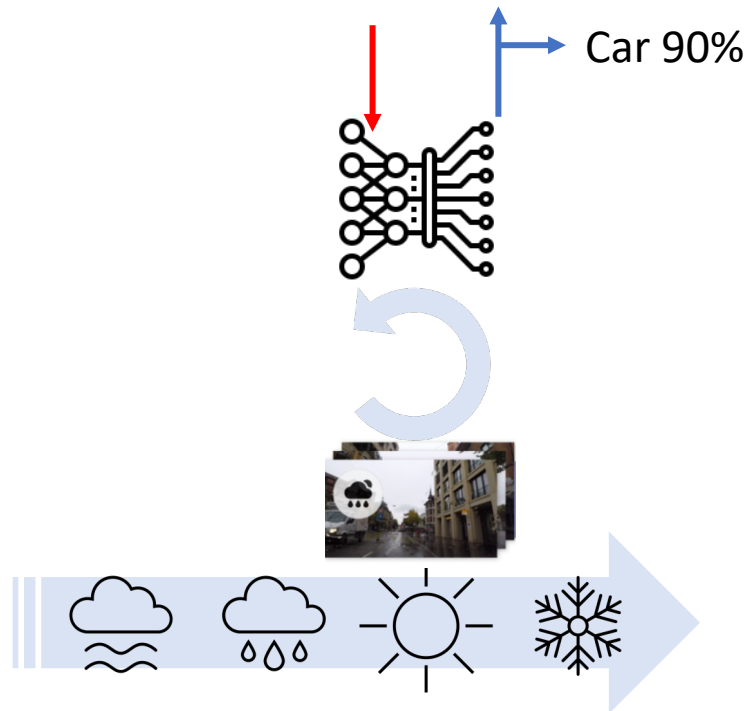


Wang, D., Shelhamer, E., Liu, S., Olshausen, B., & Darrell, T. (2021). Tent: Fully test-time adaptation by entropy minimization. ICLR.

Continual Test-time Adaptation (CTA)

Entropy Minimization (Tent)

$$\theta_t = \text{Optimize}_{\theta \in \Theta_t} (\mathbb{E}_{x \sim P_t(x)} [H(f_\theta(x))], \theta_{t-1}),$$



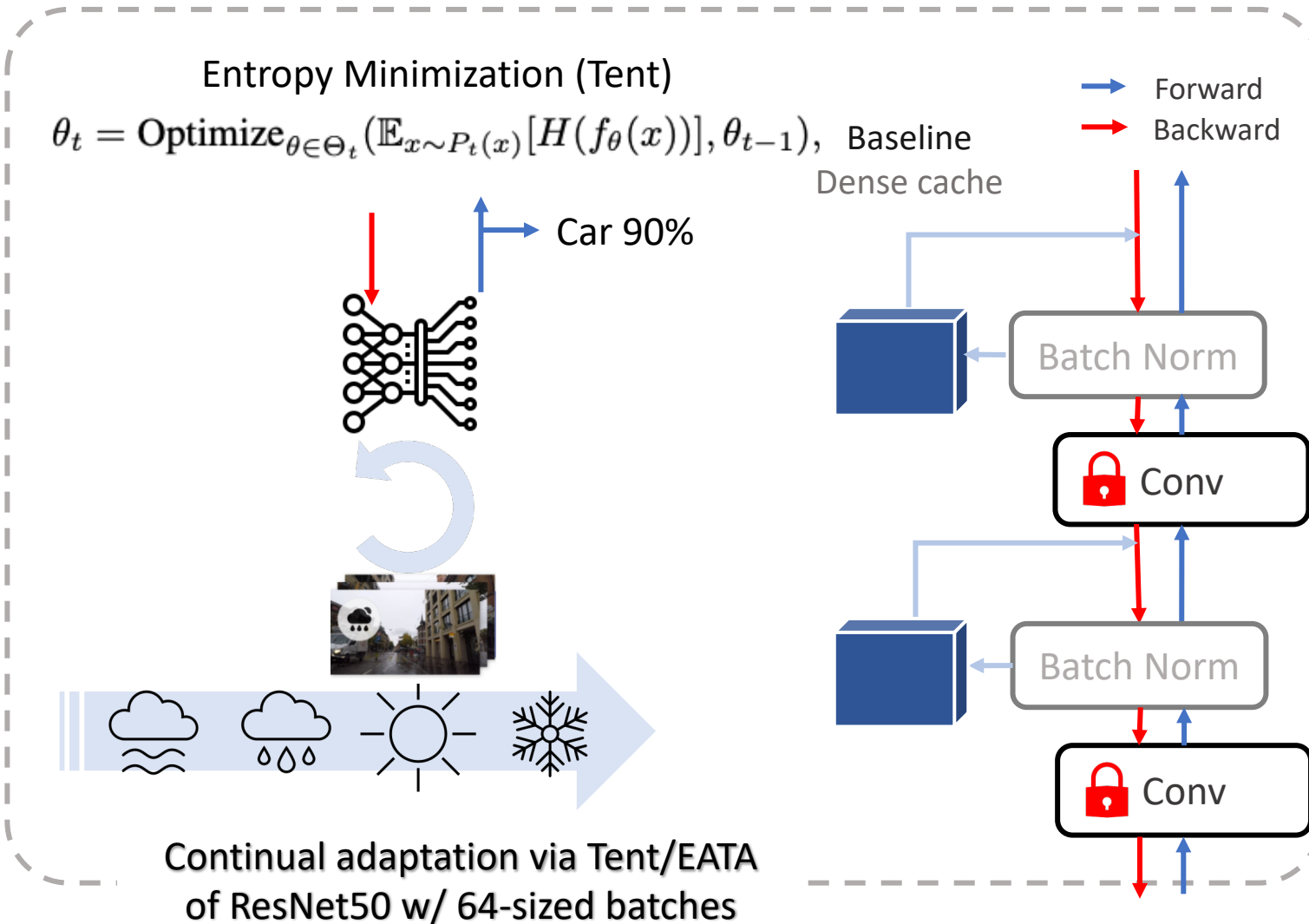
Continual adaptation via Tent/EATA
of ResNet50 w/ 64-sized batches

CTA:

- Unsupervised finetuning.

Wang, D., Shelhamer, E., Liu, S., Olshausen, B., & Darrell, T. (2021). Tent: Fully test-time adaptation by entropy minimization. ICLR.

Continual Test-time Adaptation (CTA)



CTA:

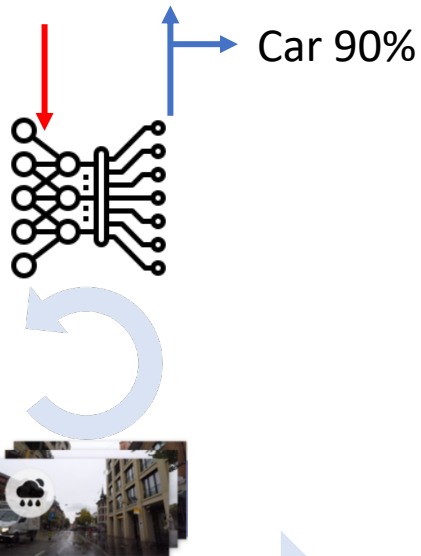
- Unsupervised finetuning.
- Parameter efficient: Update BN only.

Wang, D., Shelhamer, E., Liu, S., Olshausen, B., & Darrell, T. (2021). Tent: Fully test-time adaptation by entropy minimization. ICLR.

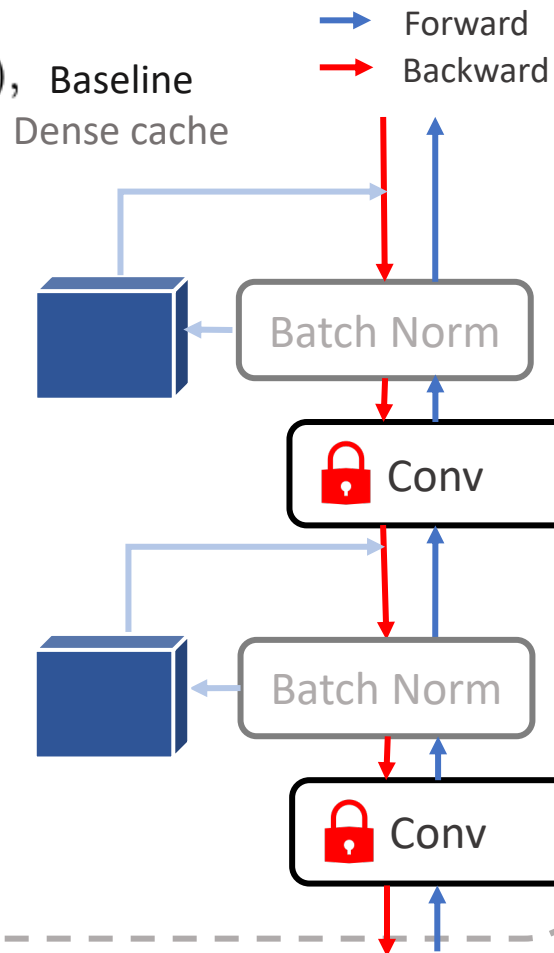
Continual Test-time Adaptation (CTA)

Entropy Minimization (Tent)

$$\theta_t = \text{Optimize}_{\theta \in \Theta_t} (\mathbb{E}_{x \sim P_t(x)} [H(f_{\theta}(x))], \theta_{t-1}), \text{ Baseline}$$



Continual adaptation via Tent/EATA of ResNet50 w/ 64-sized batches



Batch Norm (BN)

$$z_{n,i,j,k}^l = \frac{x_{n,i,j,k}^l - \mu_i^l}{\sqrt{\sigma_i^2 + \epsilon_0}} \quad a_{n,i,j,k}^l = \gamma_i^l z_{n,i,j,k}^l + b_i^l$$

CTA:

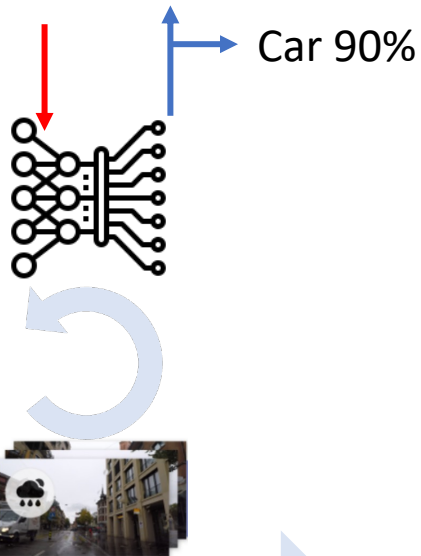
- Unsupervised finetuning.
- Parameter efficient: Update BN only.
- Batch-estimated BN statistics (μ, σ) for capturing new environment.

Wang, D., Shelhamer, E., Liu, S., Olshausen, B., & Darrell, T. (2021). Tent: Fully test-time adaptation by entropy minimization. ICLR.

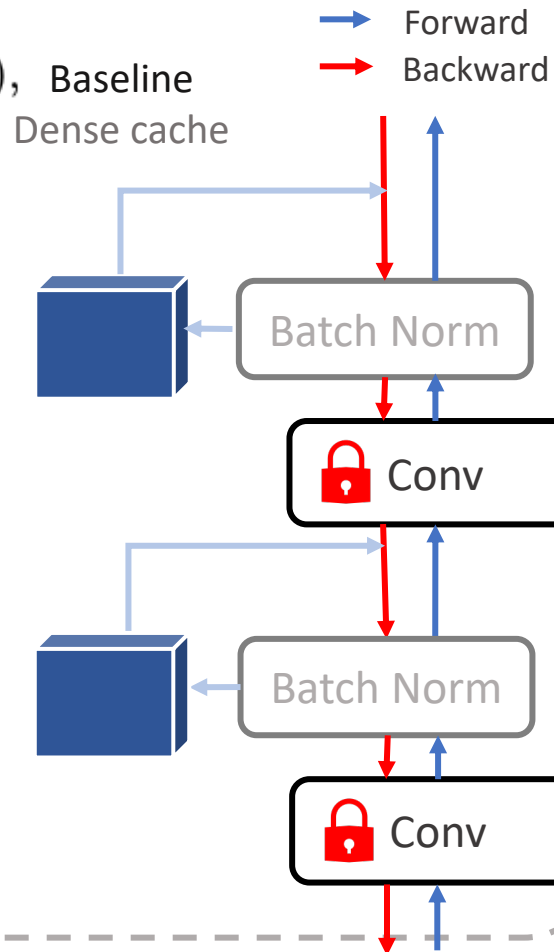
Continual Test-time Adaptation (CTA)

Entropy Minimization (Tent)

$$\theta_t = \text{Optimize}_{\theta \in \Theta_t} (\mathbb{E}_{x \sim P_t(x)} [H(f_{\theta}(x))], \theta_{t-1}), \text{ Baseline}$$



Continual adaptation via Tent/EATA of ResNet50 w/ 64-sized batches



Batch Norm (BN)

$$z_{n,i,j,k}^l = \frac{x_{n,i,j,k}^l - \mu_i^l}{\sqrt{\sigma_i^2 + \epsilon_0}} \quad a_{n,i,j,k}^l = \gamma_i^l z_{n,i,j,k}^l + b_i^l$$

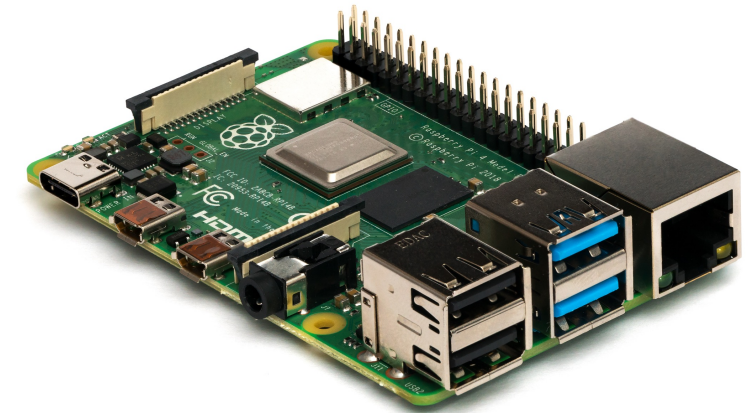
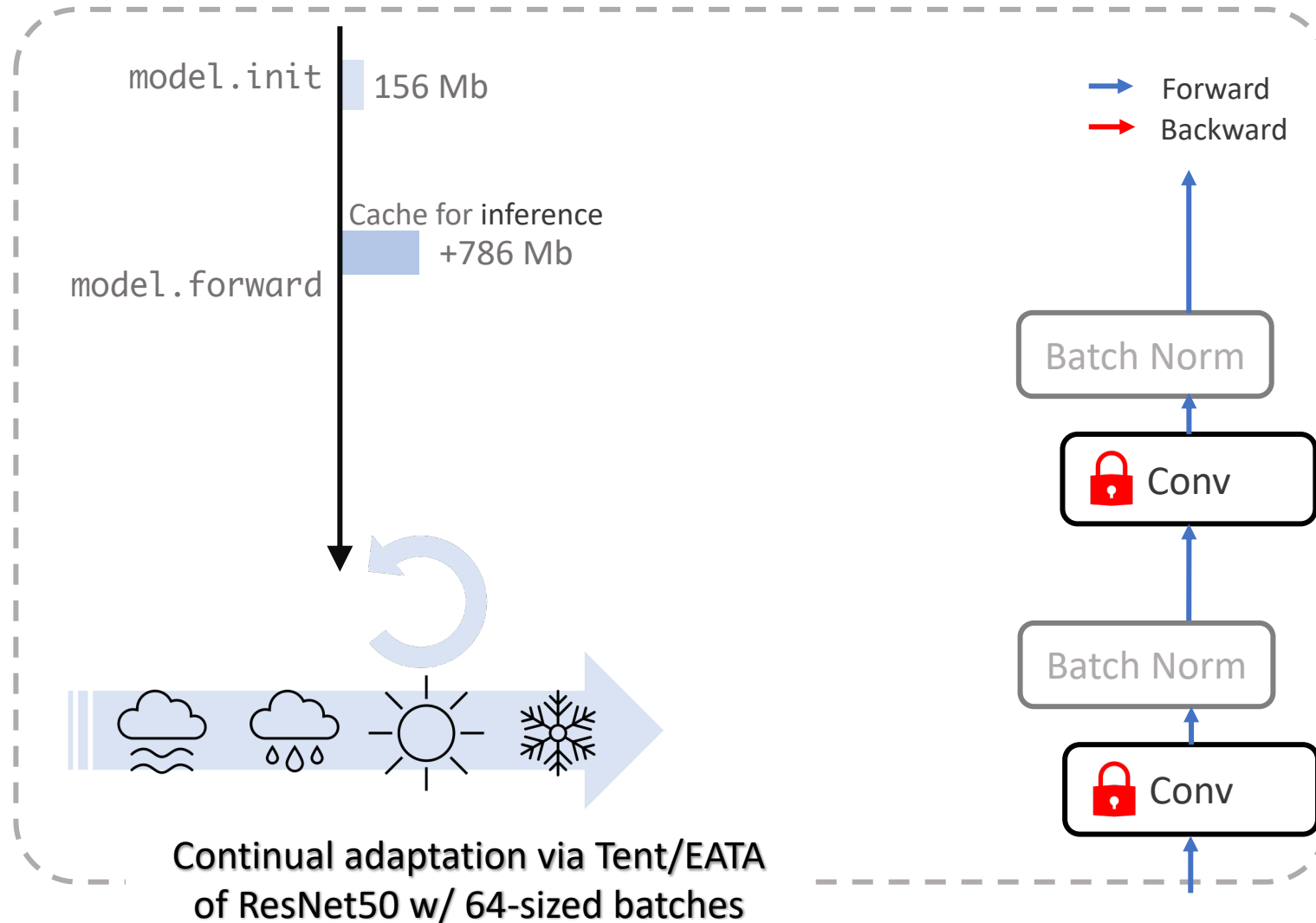
CTA:

- Unsupervised finetuning.
- Parameter efficient: Update BN only.
- Batch-estimated BN statistics (μ, σ) for capturing new environment.
- (EATA by Niu et al., 2022) Computation efficient.

Wang, D., Shelhamer, E., Liu, S., Olshausen, B., & Darrell, T. (2021). Tent: Fully test-time adaptation by entropy minimization. ICLR.

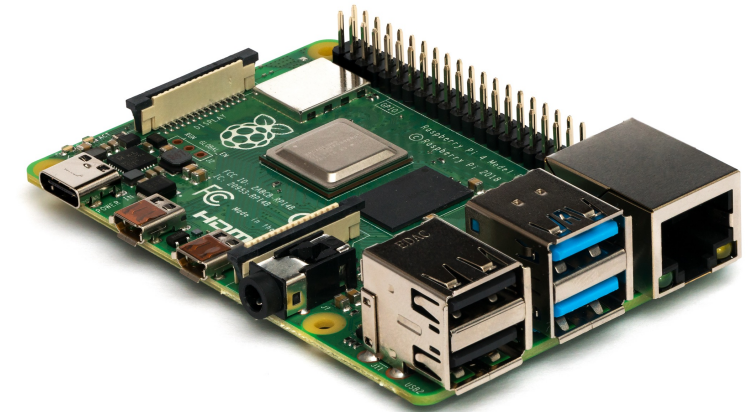
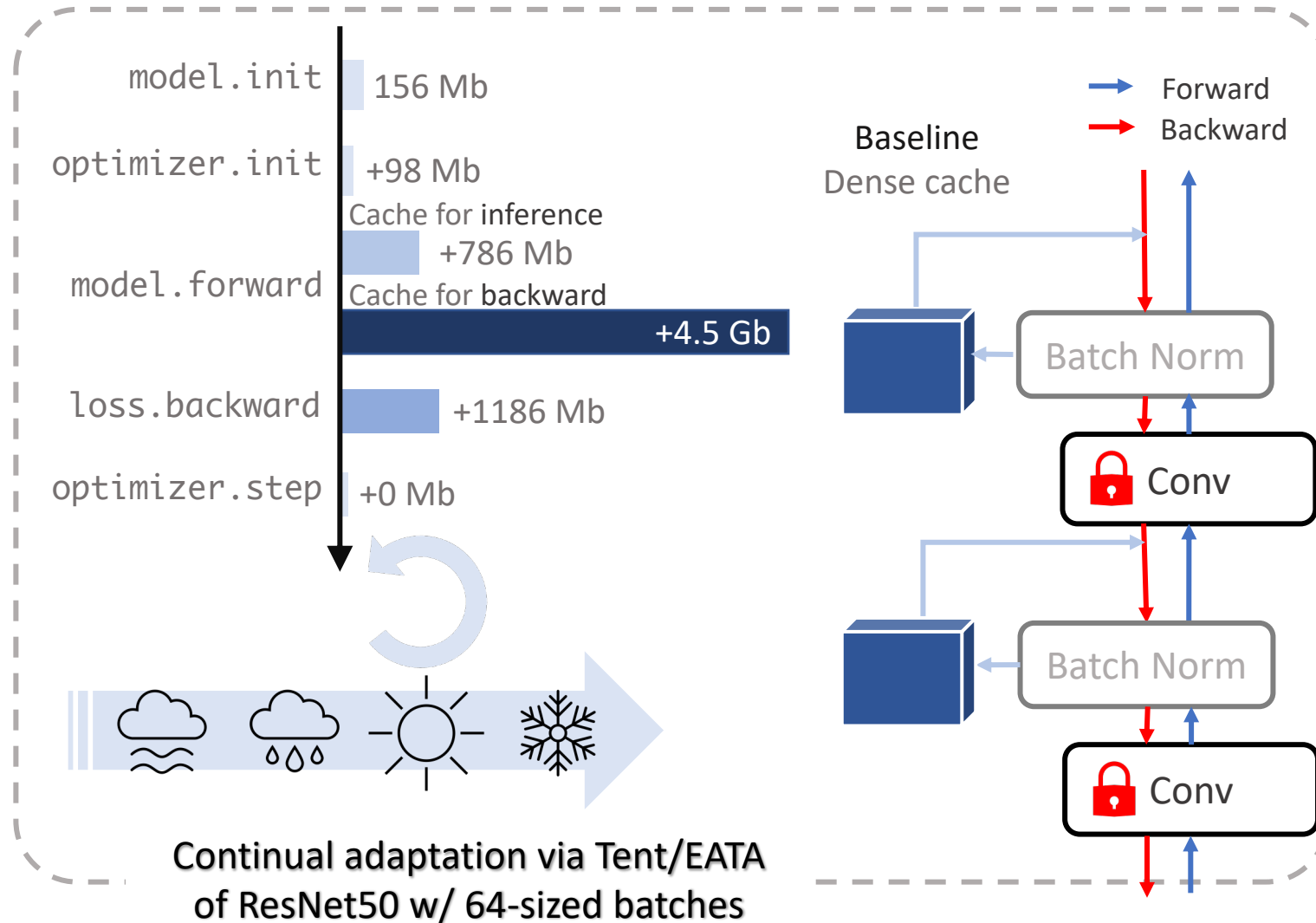
Niu, S., Wu, J., Zhang, Y., Chen, Y., Zheng, S., Zhao, P., & Tan, M. (2022). Efficient test-time model adaptation without forgetting. ICML.

High memory load for CTA on edge



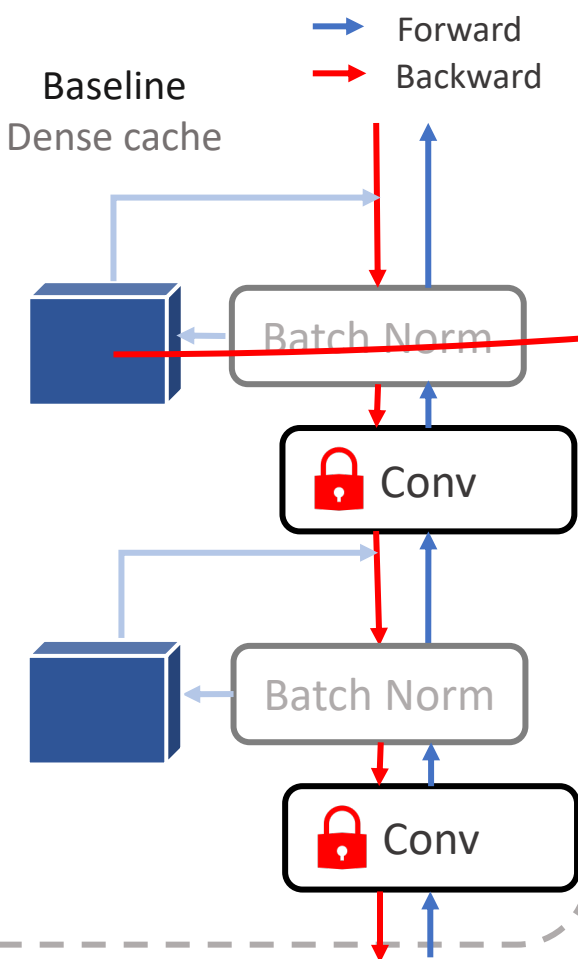
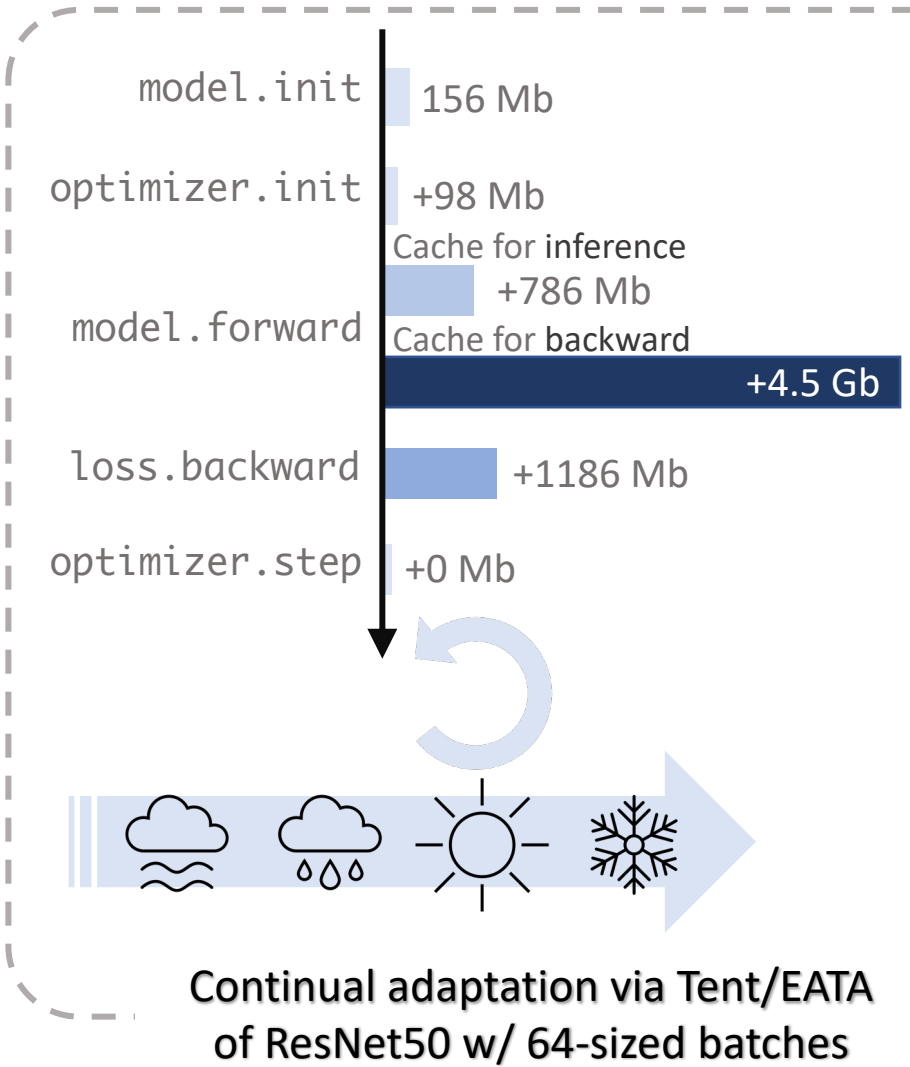
Raspberry Pi: 1-4Gb RAM
* Enough for inference

High memory load for CTA on edge



Raspberry Pi: 1-4Gb RAM
* Enough for inference **but not for adaptation**

High cache memory for back-propagation



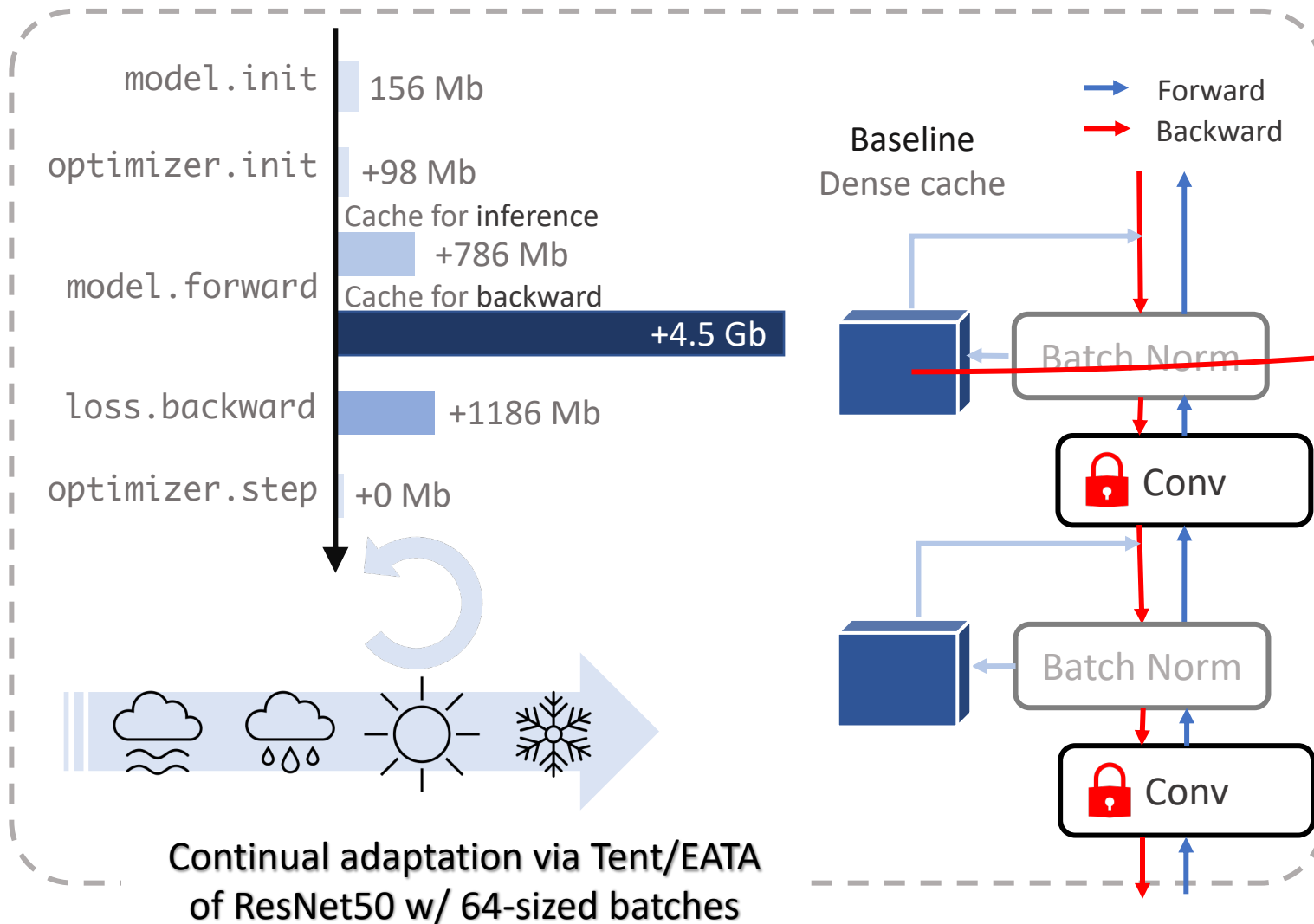
Batch-norm

$$z_{n,i,j,k}^l = \frac{x_{n,i,j,k}^l - \mu_i^l}{\sqrt{\sigma_i^2 + \epsilon_0}} \quad a_{n,i,j,k}^l = \gamma_i^l z_{n,i,j,k}^l + b_i^l$$

$$\sum_{n=1}^B \frac{\partial l_n}{\partial \gamma_i^l} = \sum_{n=1}^B \sum_{j=1}^W \sum_{k=1}^H \frac{\partial l_n}{\partial a_{i,j,k}^l} z_{n,i,j,k}^l$$

cache

High cache memory for back-propagation



Batch-norm

$$z_{n,i,j,k}^l = \frac{x_{n,i,j,k}^l - \mu_i^l}{\sqrt{\sigma_i^2 + \epsilon_0}} \quad a_{n,i,j,k}^l = \gamma_i^l z_{n,i,j,k}^l + b_i^l$$

$$\sum_{n=1}^B \frac{\partial l_n}{\partial \gamma_i^l} = \sum_{n=1}^B \sum_{j=1}^W \sum_{k=1}^H \frac{\partial l_n}{\partial a_{i,j,k}^l} z_{n,i,j,k}^l$$

Cache size

$$R_{\text{fwd}} = \max_{l \in \{1, \dots, L\}} B \times C^l \times W^l \times H^l$$

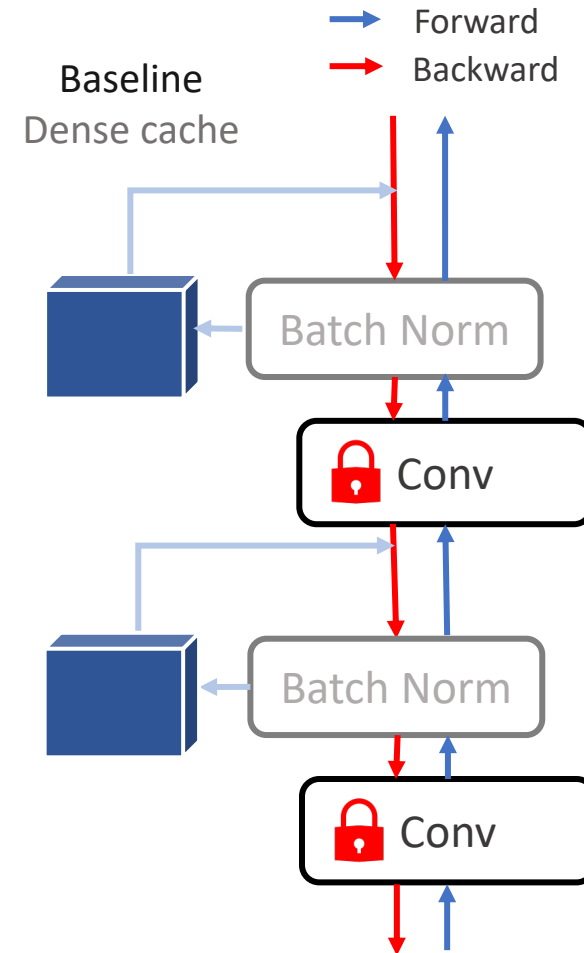
$$R_{\text{bwd}} = \sum_{l=1}^L B \times C^l \times W^l \times H^l \geq R_{\text{fwd}}$$

Traditional CTA cannot fit into low-memory devices.

- **(B)** Require large batch size for statistic estimation.
- **(L&C)** The cache tensor z scales by number of layers and channels.

Memory-Efficient Adaptation by MECTA

$$z_{n,i,j,k}^l = \frac{x_{n,i,j,k}^l - \mu_i^l}{\sqrt{\sigma_i^2 + \epsilon_0}}$$
$$R_{\text{bwd}} = \sum_{l=1}^L B \times C^l \times W^l \times H^l \geq R_{\text{fwd}}$$



Memory-Efficient Adaptation by MECTA

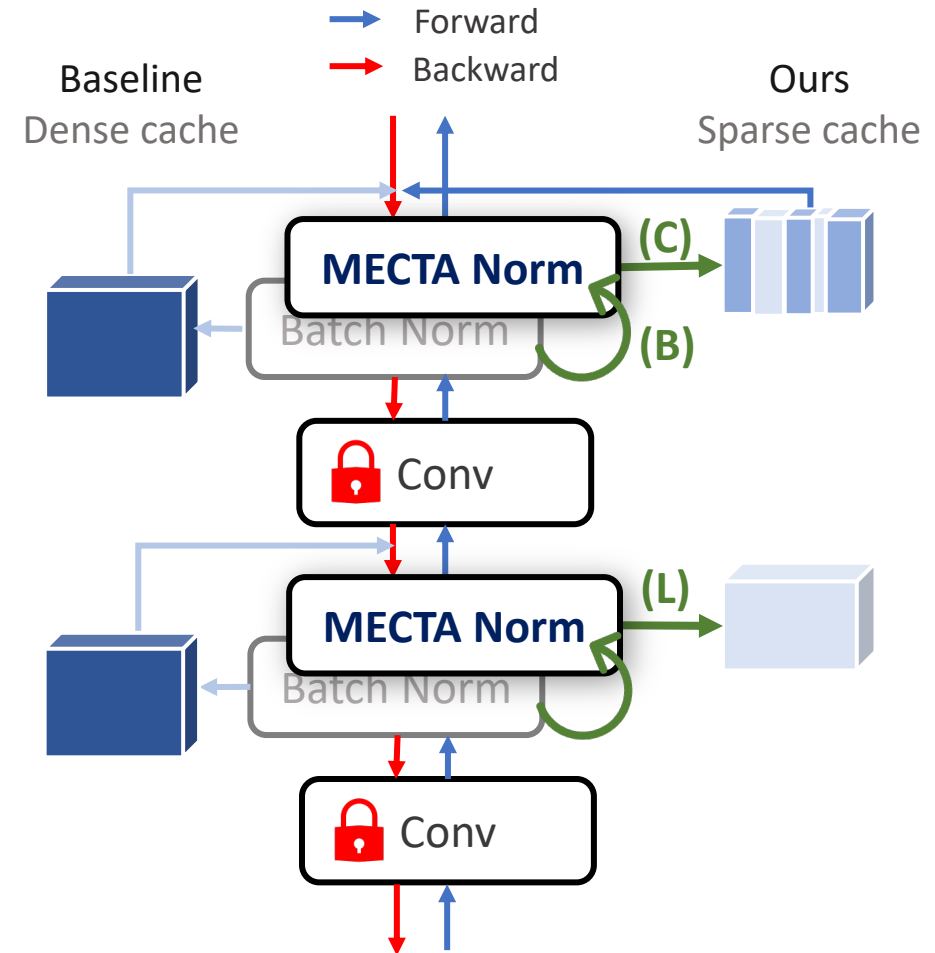
- (B) Reduce batch size and maintain accurate statistic estimation.

$$\mu = \frac{1}{B} \frac{1}{W} \frac{1}{H} \sum_{n,j,k}^{B,W,H} x_{n,i,j,k}$$

Accurate **statistics** $\phi = [\mu, \sigma]$ require more samples in a batch (larger B).

$$z_{n,i,j,k}^l = \frac{x_{n,i,j,k}^l - \mu_i^l}{\sqrt{\sigma_i^2 + \epsilon_0}}$$

$$R_{\text{bwd}} = \sum_{l=1}^L B \times C^l \times W^l \times H^l \geq R_{\text{fwd}}$$



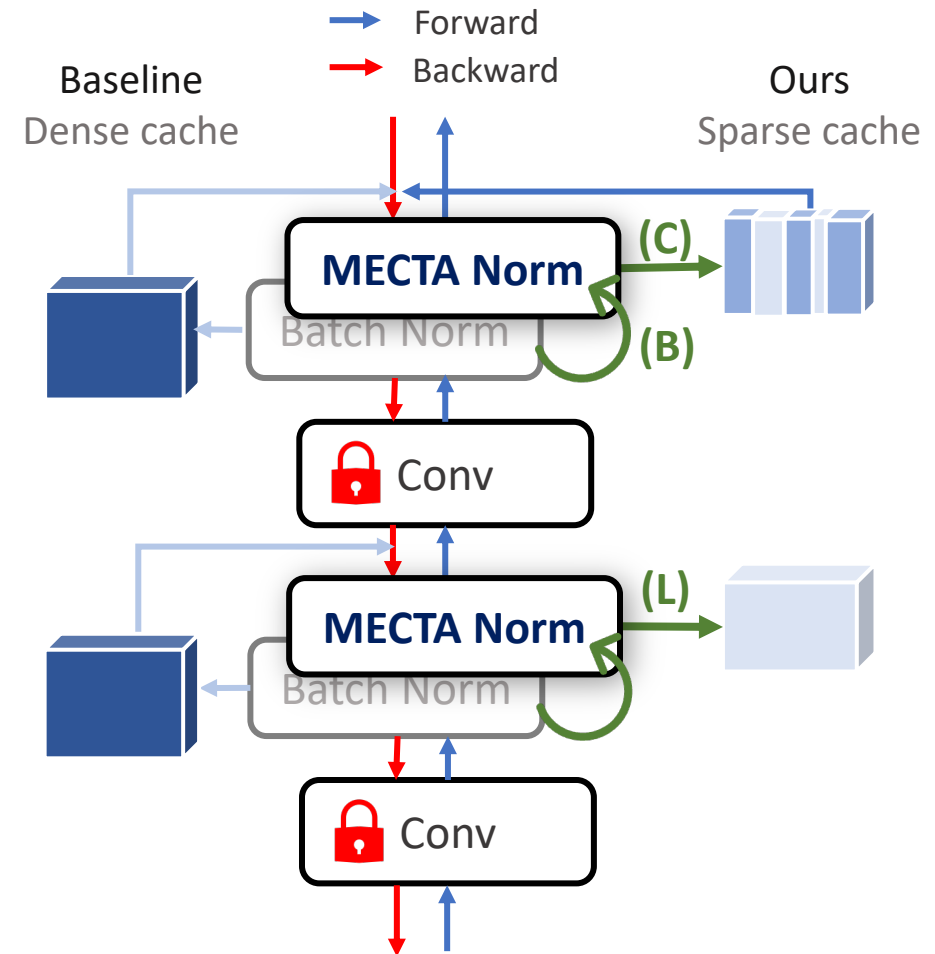
Memory-Efficient Adaptation by MECTA

- (B) Reduce batch size and maintain accurate statistic estimation.
- (C) Reduce channels to update.
- (L) Reduce layers to update.

$$z_{n,i,j,k}^l = \frac{x_{n,i,j,k}^l - \mu_i^l}{\sqrt{\sigma_i^2 + \epsilon_0}}$$

$$R_{\text{bwd}} = \sum_{l=1}^L B \times C^l \times W^l \times H^l \geq R_{\text{fwd}}$$

Drop cache \Leftrightarrow Drop gradient
 Maintain effective adaptation requires **enough gradient** information.



Memory-Efficient Adaptation by MECTA



- (B) Reduce batch size and **maintain** accurate statistic estimation **per BN layer**.

Exponential Moving Average (EMA)

ϕ represents $[\mu, \sigma]$

$$\hat{\phi}_t$$

Statistic of current batch

$$\mu = \frac{1}{B} \frac{1}{W} \frac{1}{H} \sum_{n,j,k}^{B,W,H} x_{n,i,j,k}$$

Inaccurate if $B \downarrow$

Memory-Efficient Adaptation by MECTA



- (B) Reduce batch size and **maintain** accurate statistic estimation **per BN layer**.
 - Moving average accommodate small batch sizes towards robust and accurate statistic estimations.

Exponential Moving Average (EMA)

ϕ represents $[\mu, \sigma]$

$$\phi_t = (1 - \beta)\phi_{t-1} + \beta\hat{\phi}_t,$$

Previous memory

Statistic of current batch

$$\mu = \frac{1}{B} \frac{1}{W} \frac{1}{H} \sum_{n,j,k}^{B,W,H} x_{n,i,j,k}$$

Memory-Efficient Adaptation by MECTA



- (B) Reduce batch size and **maintain** accurate statistic estimation.
 - Moving average accommodate small batch sizes towards robust and accurate statistic estimations.

Exponential Moving Average (EMA)

ϕ represents $[\mu, \sigma]$

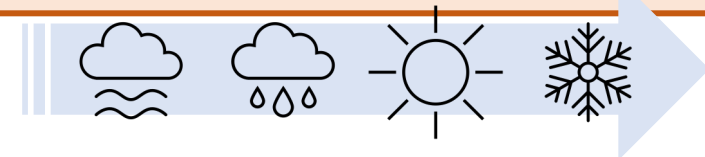
$$\phi_t = (1 - \beta)\phi_{t-1} + \beta\hat{\phi}_t,$$

Previous memory

Statistic of current batch

$\beta = 1$: Only using current **small** batch result in **inaccurate** estimation.

How to set β in changing env?



Memory-Efficient Adaptation by MECTA

- (B) Reduce batch size and **maintain** accurate statistic estimation.
 - Moving average accommodate small batch sizes towards robust and accurate statistic estimations.

Exponential Moving Average (EMA)

ϕ represents $[\mu, \sigma]$

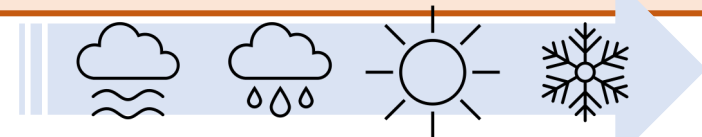
$$\phi_t = (1 - \beta)\phi_{t-1} + \beta\hat{\phi}_t,$$

Previous memory

Statistic of current batch

$\beta = 1$: Only using current **small** batch result in **inaccurate** estimation.

How to set β in changing env?



$\beta = 0$: Only using **training/past** stat (without updates) result in **non-adaptive/non-robust** estimation.

Memory-Efficient Adaptation by MECTA

- (B) Reduce batch size and maintain accurate statistic estimation.
 - Moving average accommodate small batch sizes towards robust and accurate statistic estimations.

Exponential Moving Average (EMA)

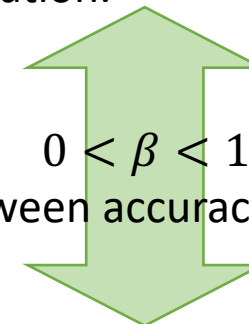
ϕ represents $[\mu, \sigma]$

$$\phi_t = (1 - \beta)\phi_{t-1} + \beta\hat{\phi}_t,$$

Previous memory

Statistic of current batch

$\beta = 1$: Only using current **small** batch result in **inaccurate** estimation.



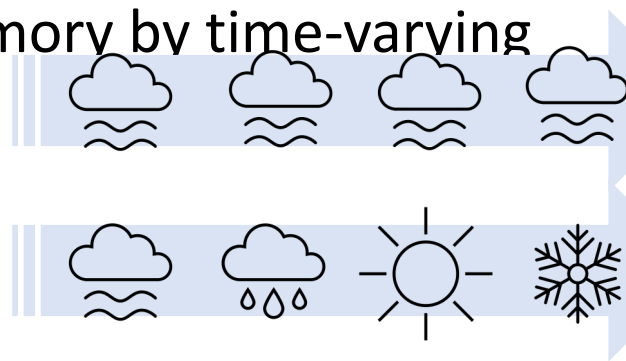
$$0 < \beta < 1$$

Trade-off between accuracy and adaptivity

$\beta = 0$: Only using **training** batch (without updates) result in **non-adaptive** estimation.

Memory-Efficient Adaptation by MECTA

- (B) Reduce batch size and maintain accurate statistic estimation.
 - Moving average accommodate small batch sizes towards robust and accurate statistic estimations.
 - Adaptive memory by time-varying β .



Exponential Moving Average (EMA)

ϕ represents $[\mu, \sigma]$

$$\phi_t = (1 - \beta)\phi_{t-1} + \beta\hat{\phi}_t,$$

Previous memory

Statistic of current batch

$\beta = 1$: Only using current **small** batch result in **inaccurate** estimation.

$P_{t-1} = P_t$: Stable $\rightarrow \beta \rightarrow 0$ (accurate)
 $P_{t-1} \neq P_t$: On change. $\rightarrow \beta \rightarrow 1$ (fast adapt)

$\beta = 0$: Only using **training** batch (without updates) result in **non-adaptive** estimation.

Memory-Efficient Adaptation by MECTA

- **(B)** Reduce batch size and maintain accurate statistic estimation.
 - Moving average accommodate small batch sizes towards robust and accurate statistic estimations.
 - Adaptive memory by time-varying β .

$$\beta_t = 1 - e^{-D(\phi_{t-1}, \hat{\phi}_t)},$$

$$D(\phi_{t-1}, \hat{\phi}_t) = \frac{1}{C} \sum_{i=1}^C KL(\phi_{t-1,i} \| \hat{\phi}_{t,i}) + KL(\hat{\phi}_{t,i} \| \phi_{t-1,i}).$$

Exponential Moving Average (EMA)

ϕ represents $[\mu, \sigma]$

$$\phi_t = (1 - \beta)\phi_{t-1} + \beta\hat{\phi}_t,$$

Previous memory

Statistic of current batch

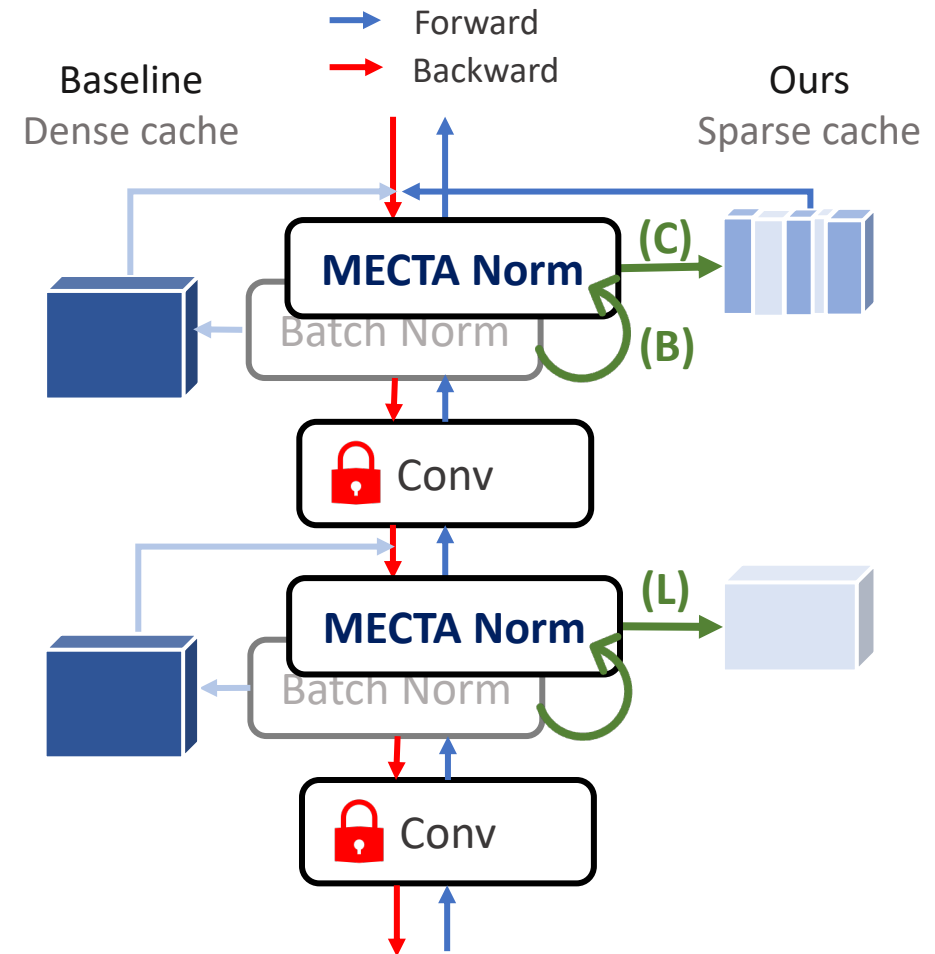
$\beta = 1$: Only using current **small** batch result in **inaccurate** estimation.

$\phi_{t-1} = \hat{\phi}_t$: Stable $\rightarrow \beta \rightarrow 0$ (accurate)
 $\phi_{t-1} \neq \hat{\phi}_t$: On change. $\rightarrow \beta \rightarrow 1$ (fast adapt)

$\beta = 0$: Only using **training** batch (without updates) result in **non-adaptive** estimation.

Memory-Efficient Adaptation by MECTA

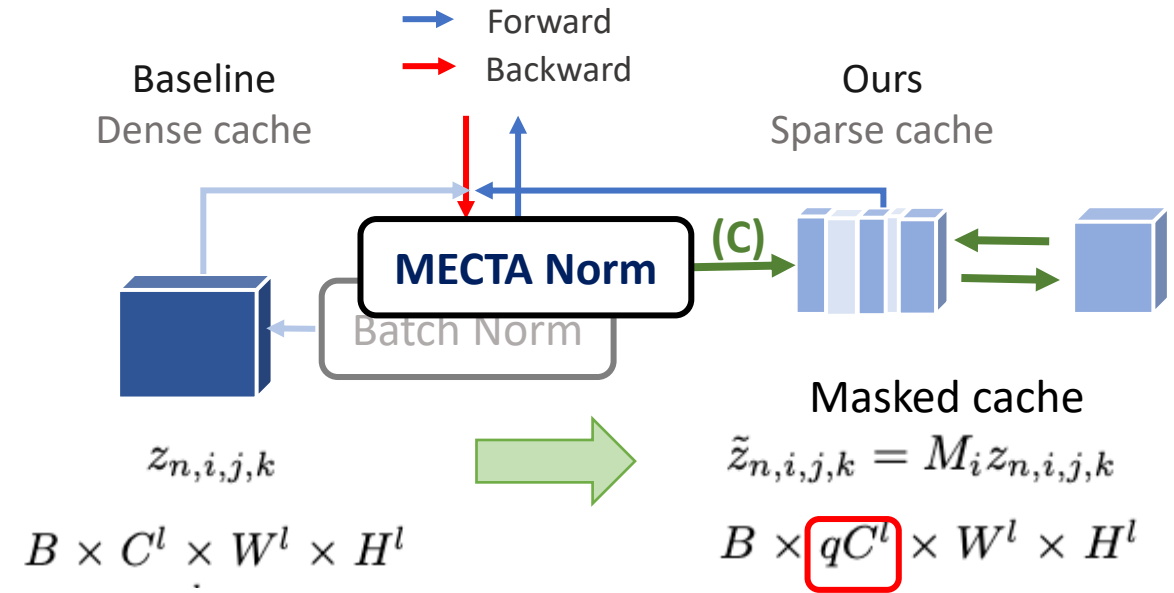
- **(B)** Reduce batch size and maintain accurate statistic estimation.
- **(C)** Reduce channels to update.
- **(L)** Reduce layers to update.



Memory-Efficient Adaptation by MECTA

- (C) Reduce channels to update.
 - Drop $q \times 100\%$ channels in cache.

Drop cache \Leftrightarrow Drop gradient
Maintain effective adaptation requires **enough gradient** information.



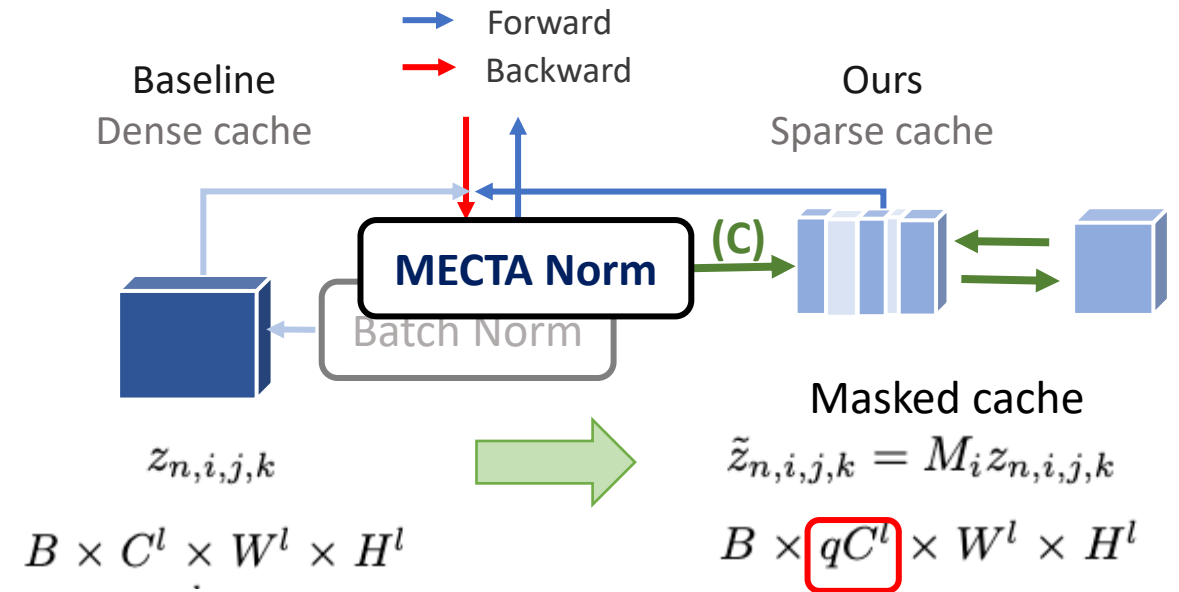
Memory-Efficient Adaptation by MECTA

- (C) Reduce channels to update.
 - Stochastically drop $q \times 100\%$ channels in cache.

Stoc. drop cache \Leftrightarrow **Not** drop gradient.
Maintain effective adaptation requires **enough gradient** information.

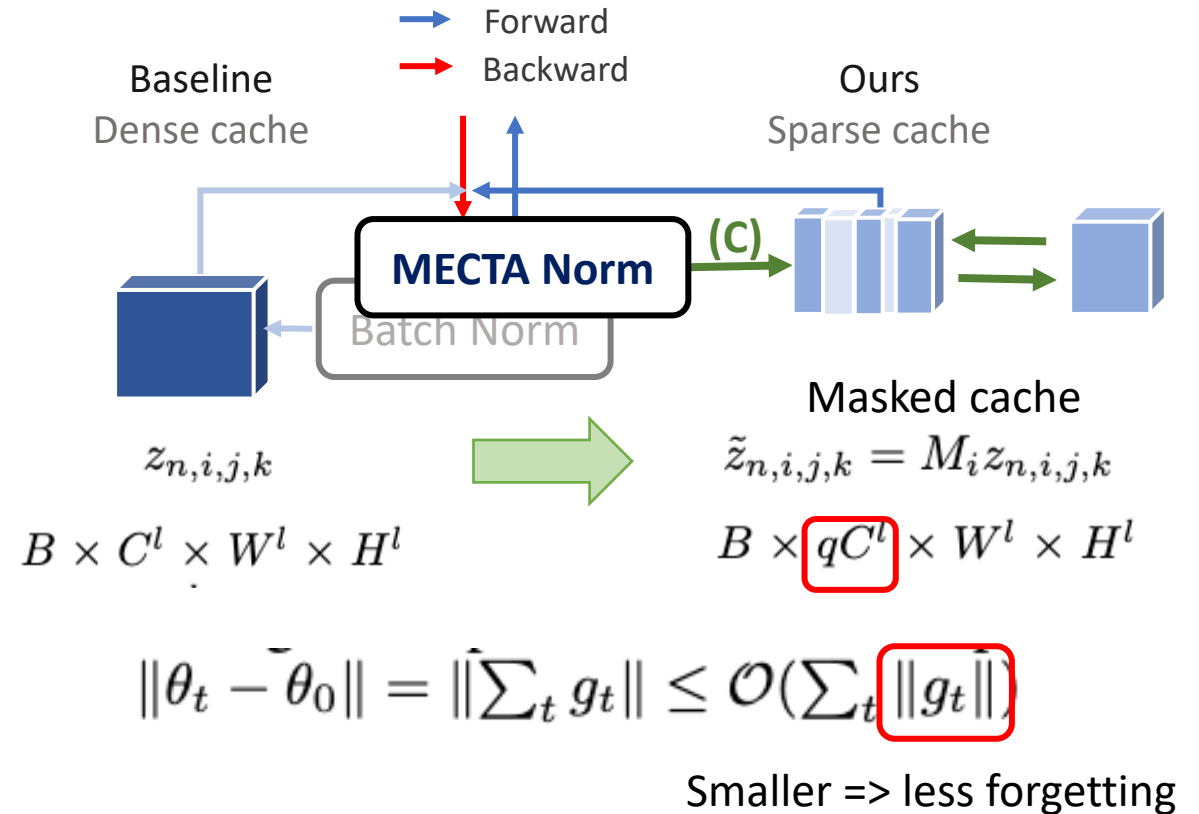
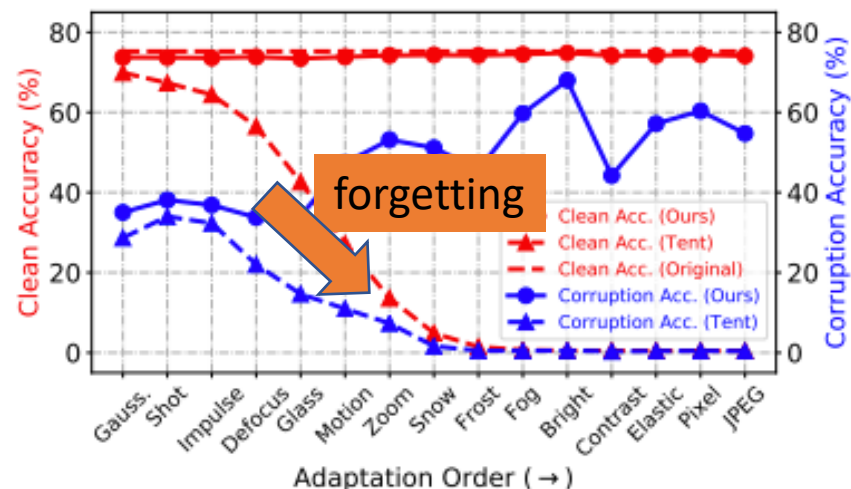
$$E[M_i z_{n,i,k}] = qE[z_{n,i,j,k}]$$

Expectation over continual batches can **restore** the dropped channels



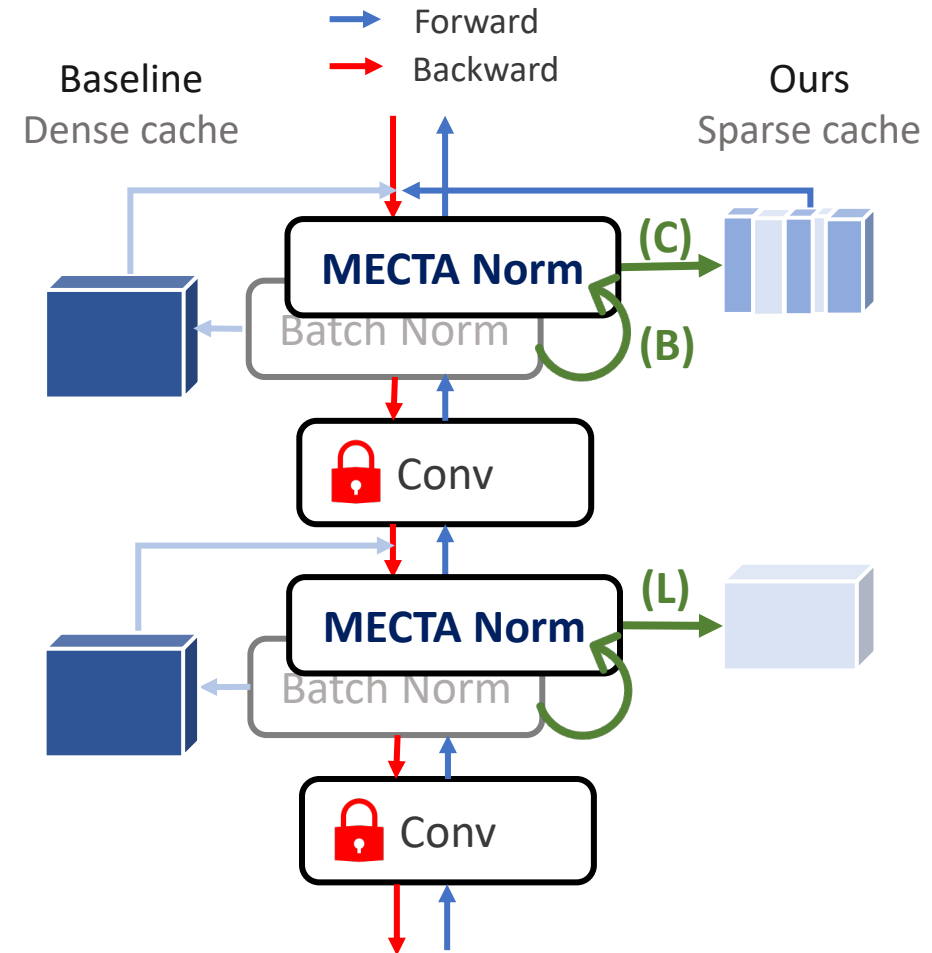
Memory-Efficient Adaptation by MECTA

- (C) Reduce channels to update.
 - Stochastically drop $q \times 100\%$ channels in cache.
 - Implicit gradient regularization which **mitigates forgetting**.



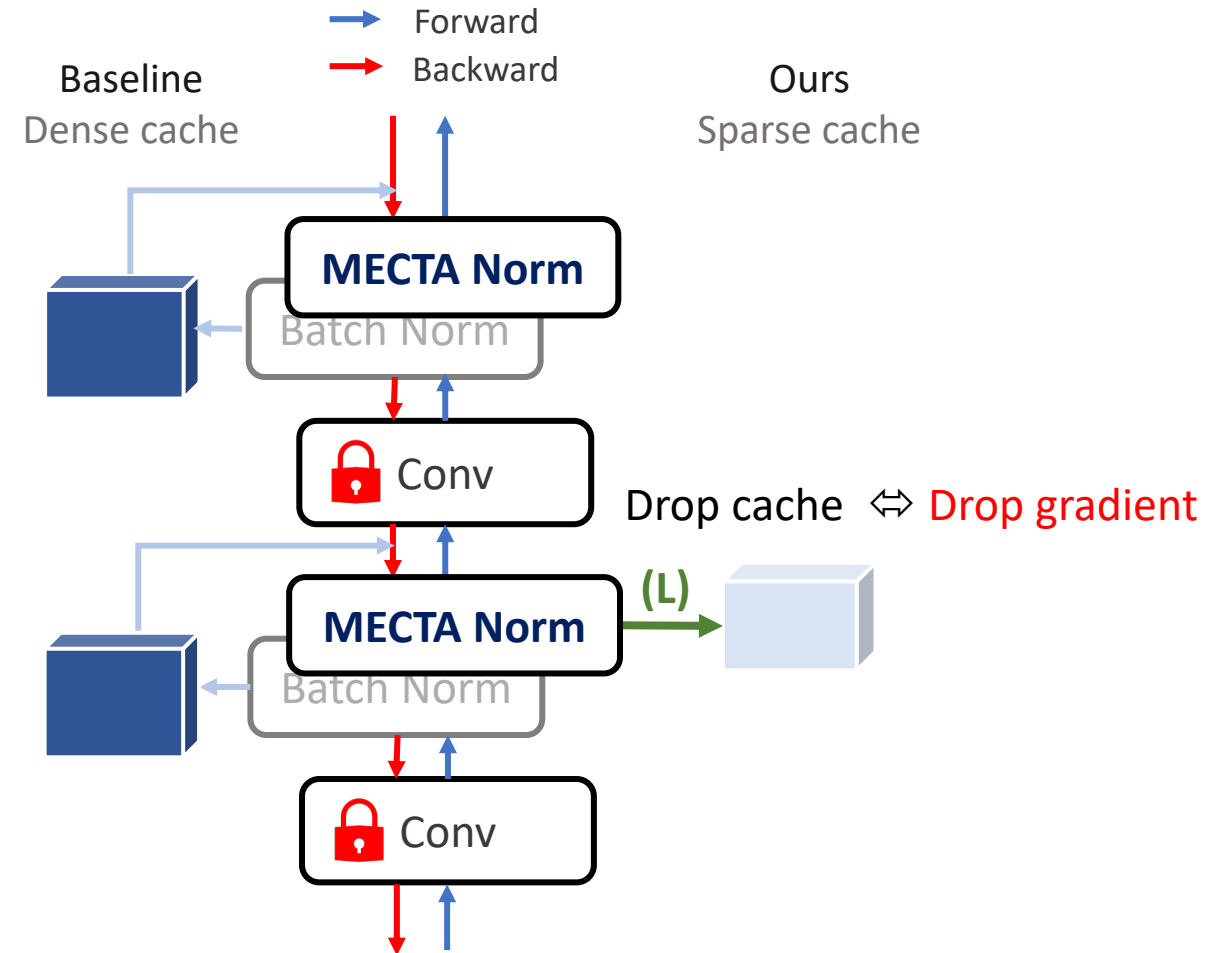
Memory-Efficient Adaptation by MECTA

- **(Reduce B)** Adaptive and online statistic estimation on dynamic distributions for accurate statistics on small batch sizes.
- **(Reduce C)** Channel-sparse gradients via stochastically-pruned caches.
- **(Dynamic L)** Cache and train layers on demand.



Memory-Efficient Adaptation by MECTA

- (Dynamic L) Cache and train layers on demand.
 - If environment is **stable**, there is no need to continually adapt.
 - Stop gradient to save caches.



Memory-Efficient Adaptation by MECTA

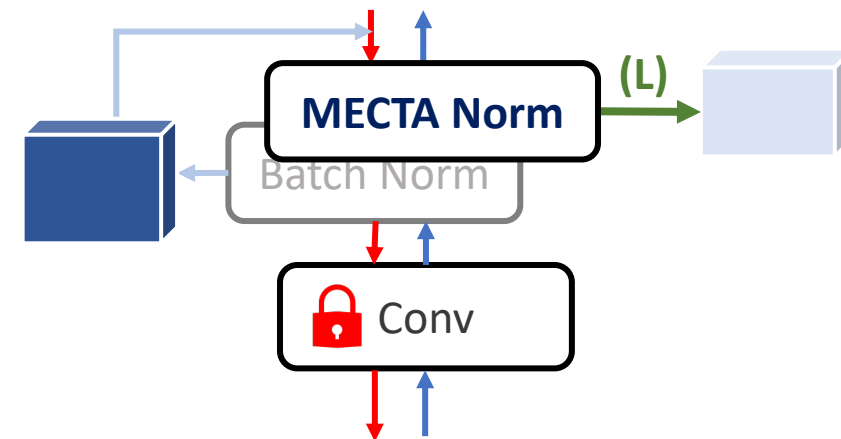
- (Dynamic L) Cache and train layers on demand.
 - If environment is **stable**, there is no need to continually adapt.
 - Stop gradient to save caches.

$$\phi_t = (1 - \beta)\phi_{t-1} + \beta\hat{\phi}_t,$$

Previous memory

Statistic of current batch

- When to stop? Environment is stable.
- When to restart? Environment changes.



Memory-Efficient Adaptation by MECTA

- (Dynamic L) Cache and train layers on demand.
 - If environment is **stable**, there is no need to continually adapt.
 - Stop gradient to save caches.

$$\phi_t = (1 - \beta)\phi_{t-1} + \beta\hat{\phi}_t,$$

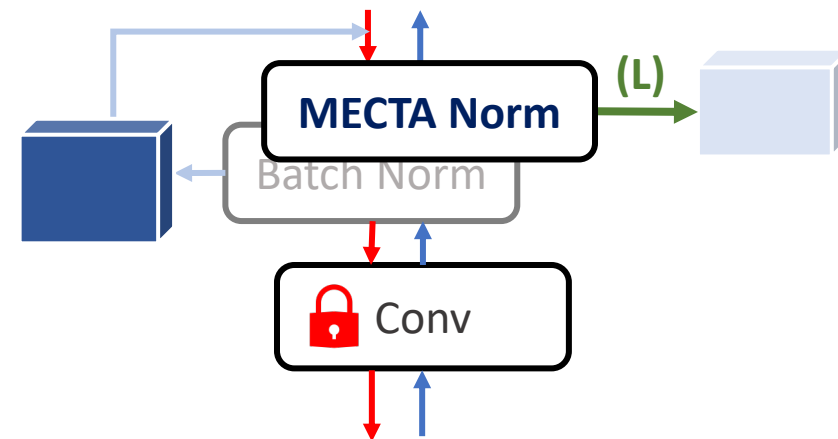
Previous memory

Statistic of current batch

$P_{t-1} = P_t$: Stable.

$P_{t-1} \neq P_t$: On change.

- When to stop? Environment is stable.
- When to restart? Environment changes.



Memory-Efficient Adaptation by MECTA

- (Dynamic L) Cache and train layers on demand.
 - If environment is **stable**, there is no need to continually adapt.
 - Stop gradient to save caches.

$$\phi_t = (1 - \beta)\phi_{t-1} + \beta\hat{\phi}_t,$$

Previous memory

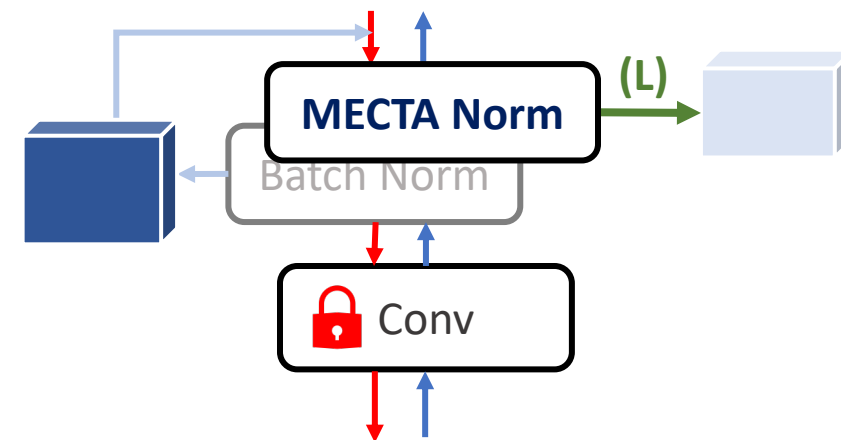
Statistic of current batch

$$P_{t-1} = P_t : \text{Stable} \quad \leftarrow \beta \rightarrow 0$$

$$P_{t-1} \neq P_t : \text{On change.} \quad \leftarrow \beta \rightarrow 1$$

$$\beta_t = 1 - e^{-D(\phi_{t-1}, \hat{\phi}_t)},$$

- When to stop? Environment is stable.
- When to restart? Environment changes.



Memory-Efficient Adaptation by MECTA

- (Dynamic L) Cache and train layers on demand.
 - If environment is **stable**, there is no need to continually adapt.
 - Stop gradient to save caches.

$$\phi_t = (1 - \beta)\phi_{t-1} + \beta\hat{\phi}_t,$$

Previous memory

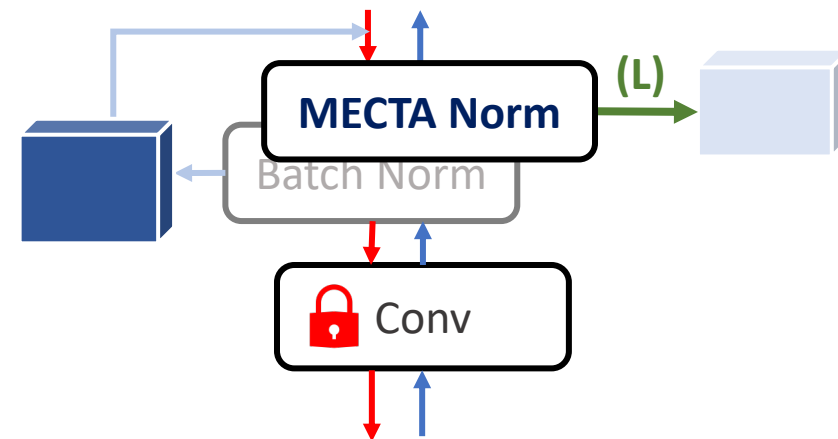
Statistic of current batch

$$P_{t-1} = P_t : \text{Stable} \quad \leftarrow \beta \rightarrow 0$$

$$P_{t-1} \neq P_t : \text{On change.} \quad \leftarrow \beta \rightarrow 1$$

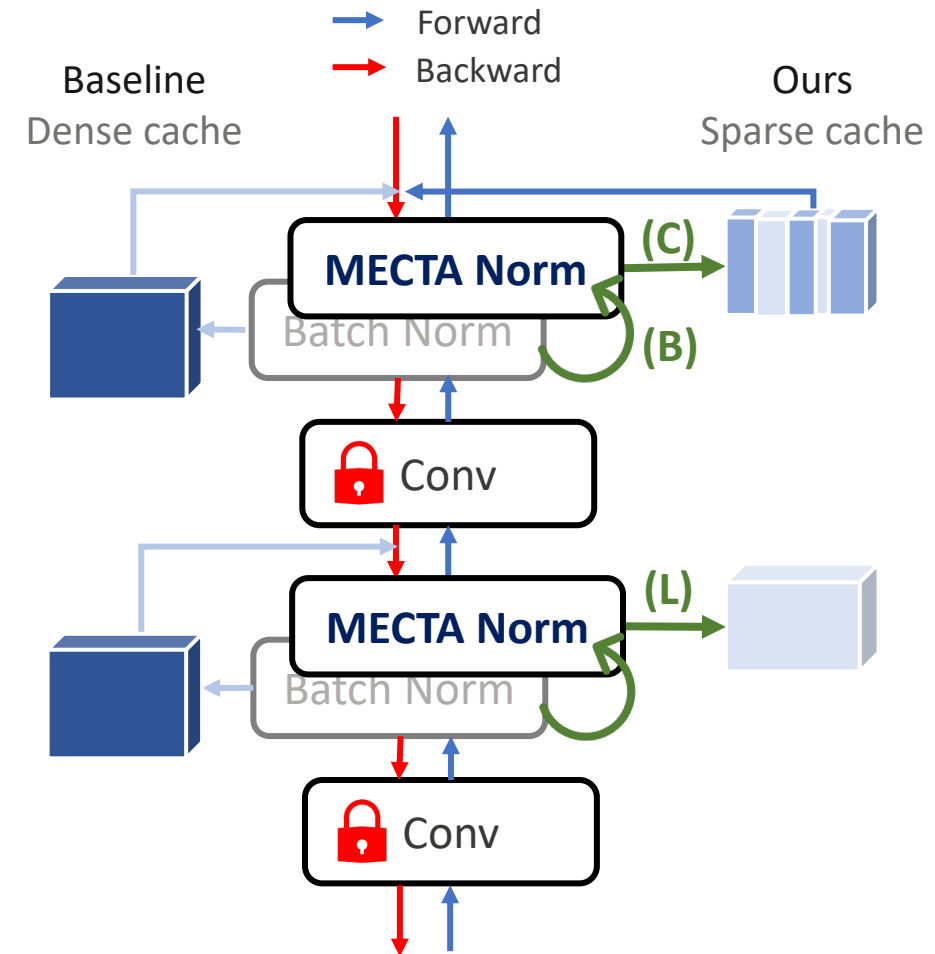
$$\beta_t = 1 - e^{-D(\phi_{t-1}, \hat{\phi}_t)}$$

- When to stop? Environment is stable: $\beta_t < \beta_{th}$
- When to restart? Environment changes: $\beta_t > \beta_{th}$

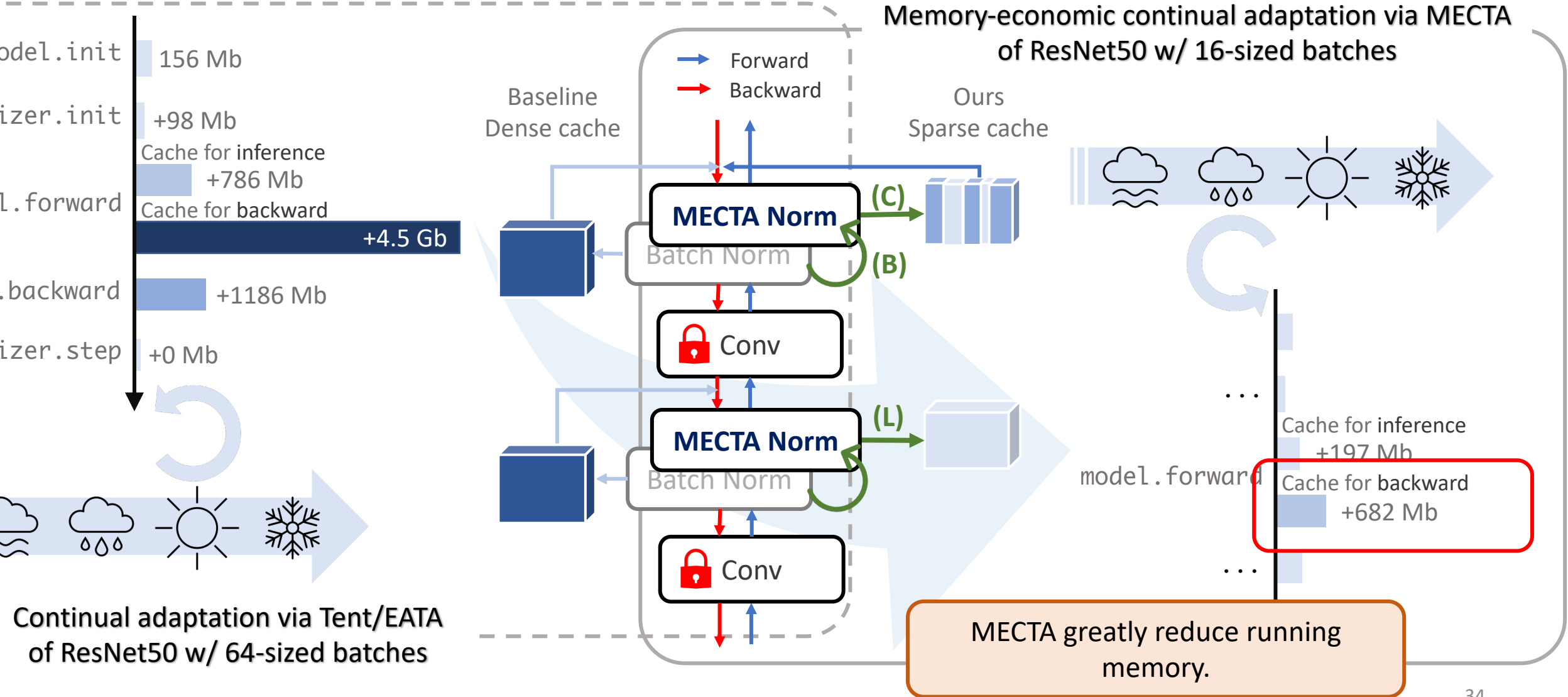


Memory-Efficient Adaptation by MECTA

- **(Reduce B)** Adaptive and online statistic estimation on dynamic distributions for accurate statistics on small batch sizes.
- **(Reduce C)** Channel-sparse gradients via stochastically-pruned caches.
- **(Dynamic L)** Cache and train layers on demand.



Memory-Efficient Adaptation by MECTA



Benchmark: Accuracy & Memory Efficiency

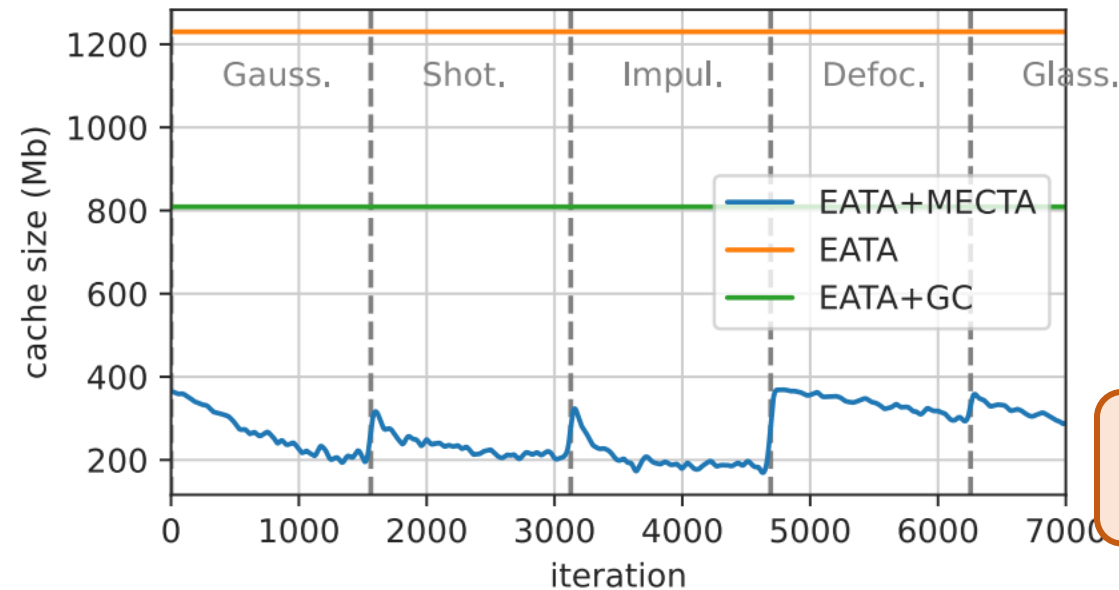
Changing test-time noise



Alg.	Noise			Blur			Weather			Digital			Acc.	Cache	#fwd	#bwd			
	Gauss.	Shot.	Impul.	Defoc.	Glass.	Motion	Zoom.	Snow	Frost	Fog	Bright.	Contr.	Elast.	Pixel.	JPEG	Avg	Max (Mb)		
ResNet50; Reset model per perturbation																			
TTT (GN+JT)	31.0	33.6	33.4	28.1	7.8	33.2	36.8	40.9	19.0	51.0	61.8	38.9	49.4	51.7	48.0	37.6	2460×20	21	20
BN	15.5	16.1	16.3	20.0	20.0	28.5	40.0	34.8	35.0	48.5	65.9	24.1	45.8	50.7	41.1	33.5	206	1	0
TTA	4.1	4.9	4.5	12.5	8.2	12.9	25.8	14.0	19.1	21.3	53.0	12.4	14.6	24.6	33.6	17.7	206	64	0
MEMO	7.5	8.7	9.0	19.7	13.0	20.7	27.6	25.3	28.8	32.1	61.0	11.0	23.8	33.0	37.5	23.9	2460×65	65	65
ResNet50; Lifelong adaptation																			
CoTTA(+GC)	16.9	20.3	22.8	20.6	22.0	31.7	42.4	34.5	34.0	47.2	58.9	24.1	44.5	48.6	42.4	34.1	2845 (1618)	33	1
Tent(+GC)	28.4	34.1	31.7	19.3	12.2	6.9	4.0	1.4	0.8	0.7	0.9	0.4	0.6	0.7	0.6	9.5	2845 (1618)	1	1
Tent+MECTA	24.5	29.5	28.3	22.0	23.5	27.4	37.2	28.2	27.1	36.8	50.7	15.5	38.0	40.2	34.7	30.9	847	1	1
EATA(+GC)	35.0	38.1	36.8	33.8	34.2	47.3	53.2	51.1	45.6	59.7	68.0	44.2	57.2	60.4	54.7	48.0	2845 (1618)	1	0.56
EATA+MECTA	<u>33.7</u>	<u>39.1</u>	<u>37.8</u>	<u>31.7</u>	<u>33.1</u>	<u>42.2</u>	<u>50.3</u>	<u>46.3</u>	<u>43.0</u>	<u>56.9</u>	<u>65.4</u>	<u>41.2</u>	<u>55.2</u>	<u>58.2</u>	<u>53.7</u>	<u>45.9</u>	<u>847</u>	1	0.56

Best trade-off between memory and accuracy.

MECTA dynamically cache data on demand



Cache on demand of environment change

Benchmark with Constrained Cache

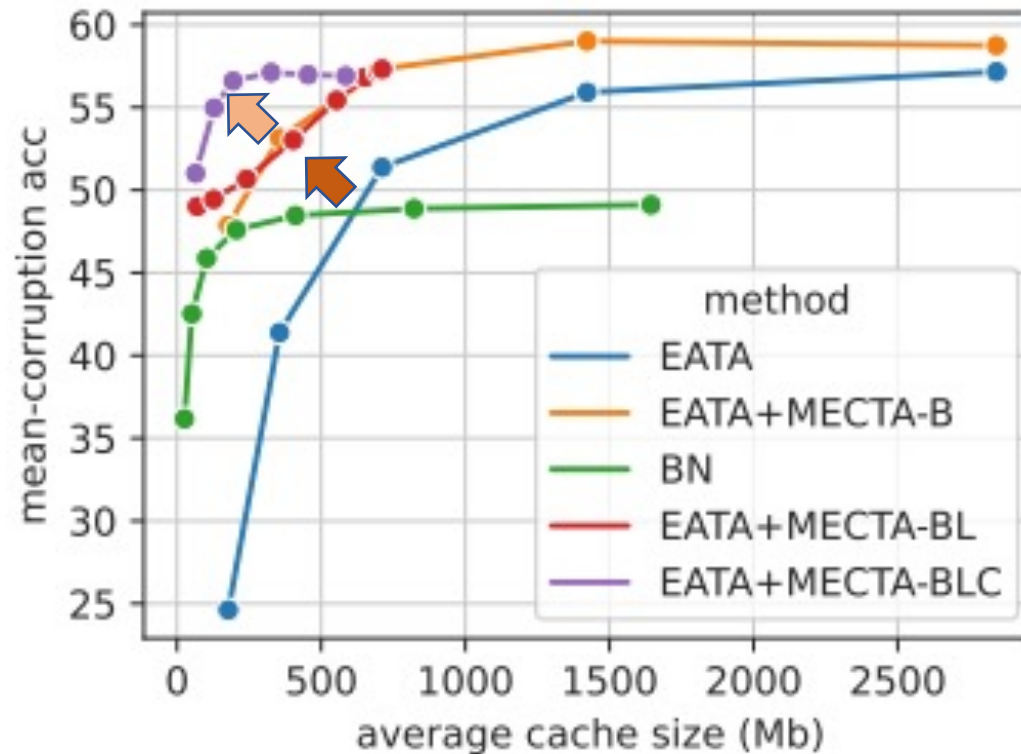
Alg.	BS	Noise			Blur			Weather				Digital				Orig.	Acc.	Cache		GFLOPs	
		Gauss.	Shot.	Impul.	Defoc.	Glass.	Motion	Zoom.	Snow	Frost	Fog	Bright.	Contr.	Elast.	Pixel.		JPEG	Avg	Avg		Max
<i>ImageNet-C</i>																					
BN	128	39.2	42.6	39.6	29.9	32.9	40.8	47.4	45.0	47.7	55.8	68.5	36.0	54.8	65.4	55.7	74.2	48.5	411	411	4.1
TENT	8	33.8	16.5	0.8	0.4	0.3	0.4	0.4	0.4	0.4	0.3	0.4	0.2	0.3	0.4	0.3	0.4	3.5	355	355	8.183
+GC	16	43.3	46.1	42.8	25.8	14.8	5.0	1.3	0.7	0.7	0.7	0.8	0.6	0.7	0.7	0.7	0.7	11.6	404	404	10.3
+MECTA	30	48.6	50.9	48.5	35.7	38.3	39.6	44.2	37.0	37.4	42.1	51.9	31.7	42.9	47.6	42.5	53.6	43.3	338	397	8.190
EATA	8	34.1	37.0	35.0	27.5	28.1	35.5	38.6	39.6	39.7	47.8	56.6	35.5	44.1	53.3	46.7	63.2	41.4	355	355	8.183
+GC	16	44.4	47.1	45.4	39.0	39.4	47.4	49.7	49.7	48.4	57.6	64.3	47.8	54.5	61.7	56.3	69.5	51.4	404	404	10.3
+MECTA	30	50.6	53.3	51.7	44.7	46.1	52.2	56.1	53.4	53.0	62.0	68.9	52.9	60.4	67.1	61.7	73.6	56.7	342	397	8.190

MECTA norm avoid forgetting of TENT.

Better accuracy on all noise, reasonable computation load.

Which component matters more?

Trade-off between robust accuracy and memory (cache size)



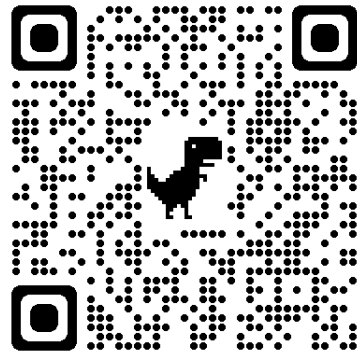
B, C are more effective while L is useful in low memory region.

MECTA: Memory-Economic Continual Test-Time Model Adaptation

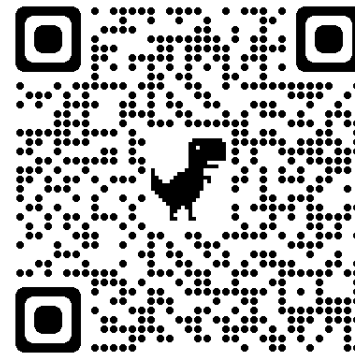
Improve **on-device** machine learning memory efficiency on **changing** environments.

- **New Problem:** We initiate the study on the **memory efficiency** of continual test-time adaptation (CTA), revealing the substantial obstacle in practice.
- **New Method:** We propose a novel method with a simple **plug-in MECTA Norm** layer that improves the memory efficiency of different CTA methods.
- **Better Memory-Robustness Trade-off:** Our method maintains comparable performance to full back-propagation methods while significantly reducing the dynamic and maximal cache overheads.

Thank you!



paper



code