

AdaHAT: Adaptive Hard Attention to the Task in Task-Incremental Learning

Pengxiang Wang¹ Hongbo Bo^{2,3} Jun Hong⁴ Weiru Liu³ Kedian Mu¹

¹ Peking University, School of Mathematical Sciences, Beijing, China

² Newcastle University, Population Health Sciences Institute, Newcastle, UK

³ University of Bristol, School of Engineering Mathematics and Technology, Bristol, UK

⁴ University of the West of England, School of Computing and Creative Technologies, Bristol, UK



Introduction

Continual Learning

Continual Learning

- ▶ A machine learning paradigm
- ▶ Learn continual tasks and adapt over time
- ▶ One of the key features of human intelligence

Catastrophic Forgetting

- ▶ Drastic performance drops on previous tasks after learning new tasks
- ▶ A major issue for continual learning algorithm to address

Problem Definition

Continual Learning (CL): learning a sequence of tasks $t = 1, \dots, N$ in order, with datasets $D^t = \{x^t, y^t\}$

Task-Incremental Learning (TIL): continual learning scenario, aim to train a model that performs well on all learned tasks

$$\max_f \sum_{t=1}^N \text{metric}(f(x^t), y^t), \{x^t, y^t\} \in D^t$$

Key assumptions when training and testing task t :

- ▶ No access to the whole data from previous tasks $1, \dots, t - 1$
- ▶ Testing on all seen tasks $1, \dots, t$
- ▶ For TIL testing, task ID t of each test sample is known by the model

Existing Approaches for TIL

Replay-based Approaches

- ▶ Prevent forgetting by storing parts of the data from previous tasks
- ▶ Replay algorithms use them to consolidate previous knowledge
- ▶ E.g. iCaRL, GEM, DER, ...

Regularization-based Approaches

- ▶ Add regularization terms constructed using information about previous tasks to the loss function when training new tasks
- ▶ E.g. LwF, EWC, SI, IMM, VCL, ...

Architecture-based Approaches

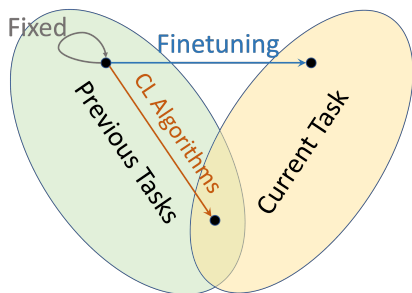
- ▶ Dedicate network parameters in different parts of the network to different tasks (inherent nature of parameter separability)
- ▶ Keep the parameters learned in previous tasks from being significantly changed
- ▶ Focus on reducing representational overlap between tasks
- ▶ E.g. Progressive Networks, PackNet, UCL, Piggyback, HAT, CPG, SupSup, ...

Stability-Plasticity Dilemma

Continual learning is a trade-off between stability and plasticity.

- ▶ **Stability:** preserve knowledge for previous tasks
- ▶ **Plasticity:** reserve representational space for new tasks

We must trade them off to get higher performance averaged on all tasks.



For replay, regularization approaches:

- ▶ Emphasis on stability in their forgetting prevention mechanisms
- ▶ But generally still lean towards plasticity

For architecture approaches:

- ▶ Distinctly different strategies that overly prioritize stability
- ▶ Tilting the trade-off towards stability instead

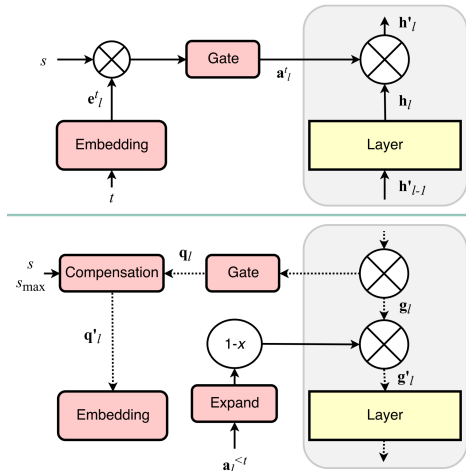
Related Work

HAT: Hard Attention to the Task

HAT (Hard Attention to the Task) is one of the most representative architecture-based approaches. Our work AdaHAT provides an extension to HAT.

Key features:

- ▶ Hard (binary) attention vectors (masks) on layers, dedicating the part of each task
- ▶ Treat the masks as model parameters, which means masks are learned
- ▶ Masks condition on gradients directly. Masked parameters won't be updated



Mechanism Details of HAT

Layer-wise attention vectors (masks) are learned to pay hard (binary) attention on units in each layer $l = 1, \dots, L - 1$ to a new task t :

$$\mathbf{m}_l^{\leq t} = \max(\mathbf{m}_l^t, \mathbf{m}_l^{\leq t-1})$$

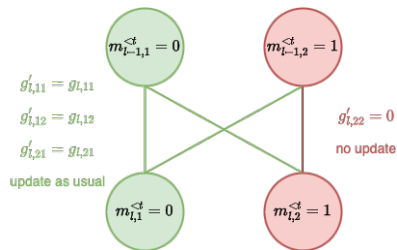
Binary values are gated from real-value task embeddings which is learnable:

$$\mathbf{m}_l^t = \sigma(\mathbf{se}_l^t)$$

Masks **hard-clip** gradients of parameters:

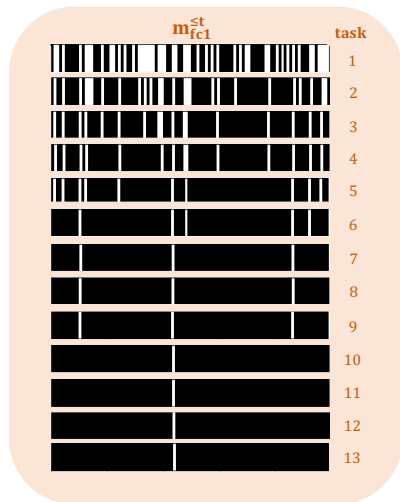
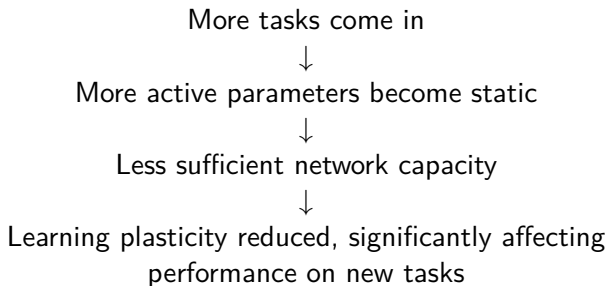
$$g'_{l,ij} = a_{l,ij} \cdot g_{l,ij}, \quad a_{l,ij} \in \{0, 1\}$$

$$a_{l,ij} = 1 - \min(m_{l,i}^{\leq t}, m_{l-1,j}^{\leq t})$$



Problem 1: Insufficient Network Capacity

Architecture-based approaches all suffer from **network capacity problem** especially in long sequence of tasks, sacrificing plasticity for stability. HAT's hard-clipping mechanism allows no update for parameters masked by previous tasks.



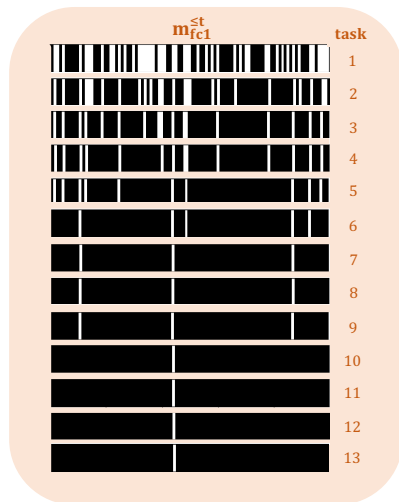
Problem 1: Insufficient Network Capacity

HAT tries to solve it by sparsity regularization on learnable masks:

$$\mathcal{L}' = \mathcal{L}(f(x_t), y_t) + cR(M^t, M^{<t})$$

$$R(M^t, M^{<t}) = \frac{\sum_{l=1}^{L-1} \sum_{i=1}^{N_l} m_{l,i}^t (1 - m_{l,i}^{<t})}{\sum_{l=1}^{L-1} \sum_{i=1}^{N_l} (1 - m_{l,i}^{<t})}$$

- ▶ Meant to promote low network capacity usage and high compactness of the masks
- ▶ Helps alleviate the issue on network capacity to a certain extent
- ▶ However, the network capacity will eventually run out



Problem 2: Non-adaptive Hyperparameters

Most architecture-based approaches:

- ▶ Use several hyperparameters to manually allocate network capacity usage
- ▶ Without leveraging any information about previous tasks
- ▶ E.g. PackNet use pruning ratios, HAT uses s_{\max}

In continual learning:

- ▶ We never know how many tasks in future, maybe infinite ...
- ▶ We're not able to decide the capacity allocation beforehand
- ▶ Manual hyperparameter tuning is infeasible!

In our work, an adaptive strategy smartly allocates the network capacity with taking into account the information about previous tasks.

Methodology

AdaHAT: Adaptive Hard Attention to the Task

Our Proposed AdaHAT **soft-clips gradients**, which allows minor updates for parameters masked by previous tasks:

$$g'_{l,ij} = a_{l,ij}^* \cdot g_{l,ij}, \quad a_{l,ij}^* \in [0, 1]$$

The adjustment rate $a_{l,ij}^*$ now is an adaptive controller, guided by two pieces of information about previous tasks directly from HAT architecture:

- ▶ **Parameter Importance**
- ▶ **Network Sparsity**

AdaHAT: Parameter Importance

The attention vectors (masks) indicate **the importance of parameter**.

Cumulative Attention Vectors (HAT)

$$\mathbf{m}_l^{\leq t} = \max(\mathbf{m}_l^t, \mathbf{m}_l^{\leq t-1})$$

- ▶ Binary $\{0, 1\}$, represents if it's masked by previous tasks

Summative Attention Vectors (AdaHAT)

$$\mathbf{m}_l^{\leq t, \text{sum}} = \mathbf{m}_l^t + \mathbf{m}_l^{\leq t-1, \text{sum}}$$

- ▶ Range from 0 to $t - 1$, represents how many previous tasks it's masked
- ▶ Encapsulates more information about previous tasks

Adaptive process: Higher summative vectors \rightarrow More important to previous tasks
 \rightarrow Smaller adjustment rate $a_{l,ij}^*$ \rightarrow Smaller updates for the parameter

AdaHAT: Network Sparsity

The sparsity regularization term $R(M^t, M^{<t})$ measures the compactness of masks.

It is closely related to the current network capacity:

Generally, when a smaller proportion of parameters in the network are static (i.e., sufficient network capacity), the regularization value tends to be larger, as there is a great possibility for the hard attention to be paid to active parameters.

Adaptive process: Higher sparsity regularization \rightarrow (Suggesting) more unmasked space available for new tasks \rightarrow Less need to adjust the static space for previous tasks, should go for active parameters \rightarrow Smaller adjustment rate $a_{l,ij}^*$ in general

Note: in this way, AdaHAT tries its best to mimic HAT before the network capacity limit, retaining maximum stability before affecting plasticity for new tasks.

AdaHAT: The Adjustment Rate

Adaptive Adjustment Rate (AdaHAT)

$$a_{l,ij}^* = \frac{r_l}{\min(m_{l,i}^{<t,\text{sum}}, m_{l-1,j}^{<t,\text{sum}}) + r_l}, \quad r_l = \frac{\alpha}{R(M^t, M^{<t}) + \epsilon}$$

Adjustment Rate (HAT)

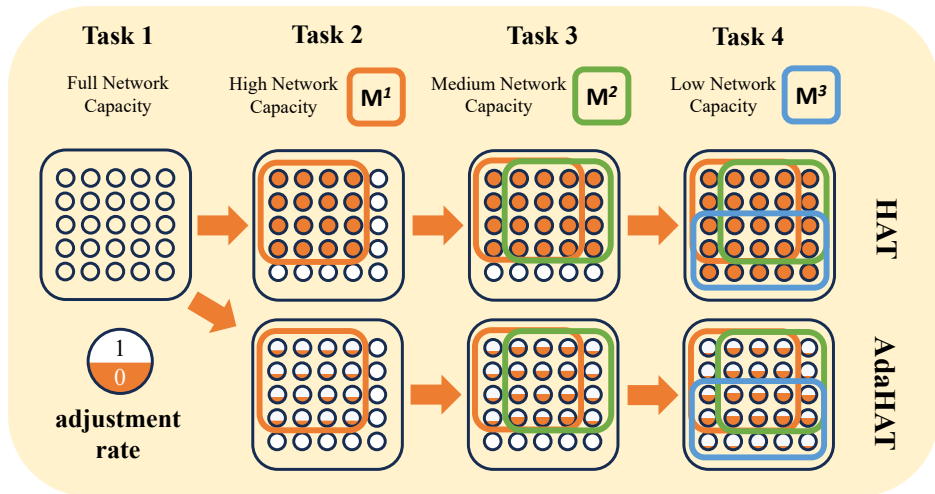
$$a_{l,ij} = 1 - \min(m_{l,i}^{<t}, m_{l-1,j}^{<t})$$

The adjustment rate in AdaHAT adaptively incorporates both information about previous tasks:

- ▶ The higher parameter importance $\min(m_{l,i}^{<t,\text{sum}}, m_{l-1,j}^{<t,\text{sum}})$, the lower adjustment rate
- ▶ The higher network sparsity $R(M^t, M^{<t})$, the higher adjustment rate

While in HAT only the accumulation of masks.

AdaHAT: Adaptive Hard Attention to the Task



Experiments

Main Results

Table 1. Results on performance and stability-plasticity trade-off metrics (mean \pm std) of different approaches on the two datasets (20 tasks).

Dataset	Approach	AA(%)	FR (%)	BWT(%)	FWT (%)
Permuted MNIST	Finetuning	32.62 \pm 1.60	-73.78 \pm 1.84	-68.10 \pm 1.68	63.51 \pm 0.03
	LwF	26.95 \pm 1.80	-80.35 \pm 2.08	-72.59 \pm 1.91	62.04 \pm 0.09
	EWC	52.25 \pm 2.46	-51.38 \pm 2.83	-42.04 \pm 2.67	58.12 \pm 0.15
	HAT	67.64 \pm 1.27	-33.70 \pm 1.46	-0.11 \pm 0.18	32.49 \pm 1.12
	HAT-random	66.43 \pm 1.21	-35.10 \pm 1.39	-0.27 \pm 0.49	31.40 \pm 1.22
	HAT-const-alpha	68.08 \pm 1.18	-33.20 \pm 1.36	-1 * e ⁻³ \pm 0.00	32.92 \pm 1.23
	HAT-const-1	48.83 \pm 4.35	-55.14 \pm 5.02	-49.68 \pm 4.40	62.26 \pm 0.21
	AdaHAT	79.90 \pm 2.40	-19.43 \pm 2.76	-14.68 \pm 2.48	59.96 \pm 0.09
Split CIFAR- 100	Finetuning	24.34 \pm 0.73	-91.66 \pm 1.32	-54.00 \pm 1.00	53.10 \pm 0.55
	LwF	34.56 \pm 0.94	-70.91 \pm 2.05	-48.03 \pm 1.01	57.61 \pm 0.40
	EWC	30.23 \pm 1.61	-79.84 \pm 3.13	-54.05 \pm 1.28	59.20 \pm 0.50
	HAT	32.44 \pm 1.58	-74.71 \pm 3.37	-45.59 \pm 1.49	53.11 \pm 0.34
	HAT-random	31.41 \pm 1.29	-76.98 \pm 2.45	-48.80 \pm 1.33	52.76 \pm 0.57
	HAT-const-alpha	32.16 \pm 2.48	-75.04 \pm 5.16	-44.49 \pm 2.57	51.86 \pm 0.82
	HAT-const-1	32.40 \pm 1.40	-75.58 \pm 3.08	-48.80 \pm 1.72	56.30 \pm 0.36
	AdaHAT	38.74 \pm 2.24	-62.37 \pm 4.64	-42.11 \pm 2.02	56.33 \pm 0.82

Main Results

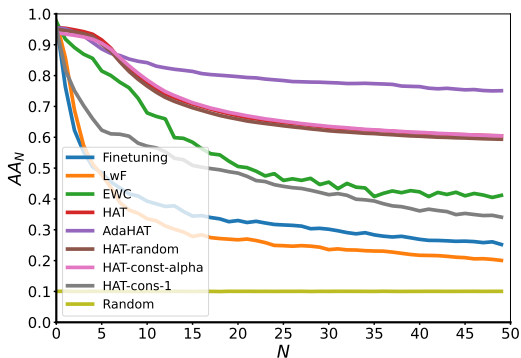
- ▶ Datasets: Permuted MNIST, Split CIFAR-100, 20 tasks
- ▶ Main metrics:
 - ▶ Average Accuracy (AA) over all tasks
 - ▶ Forgetting Rate (FR)
- ▶ Metrics for stability-plasticity trade-off:
 - ▶ Backward Transfer (BMT) for stability
 - ▶ Forward Transfer (FWT) for plasticity

Results:

- ▶ AdaHAT outperforms all baselines
- ▶ AdaHAT balances stability-plasticity better, while
 - ▶ HAT: high BWT, low FWT
 - ▶ Finetuning: low BWT, high FWT
 - ▶ HAT-const-1: low BWT, high FWT

Conclusions: it is important to maintain a balanced stability-plasticity trade-off for optimal performance.

Results on Longer Task Sequences



Dataset: Permutated MNIST, **50 tasks**
(longer)

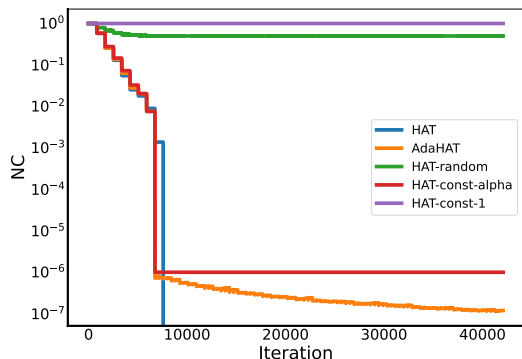
Results:

- ▶ HAT slightly outperforms before task 8 then drastically drops
- ▶ AdaHAT keeps significant superiority after the turning point
- ▶ AdaHAT is still close to HAT before task 8

Conclusions:

- ▶ There is a turning point for HAT when it exhausts network capacity
- ▶ AdaHAT mimics HAT well before the network capacity limit, and shows much more capability for long task sequence settings

Network Capacity Usage



Network Capacity Measurement

$$NC = \frac{1}{\sum_l N_l} \sum_{l,i,j} a_{l,ij}$$

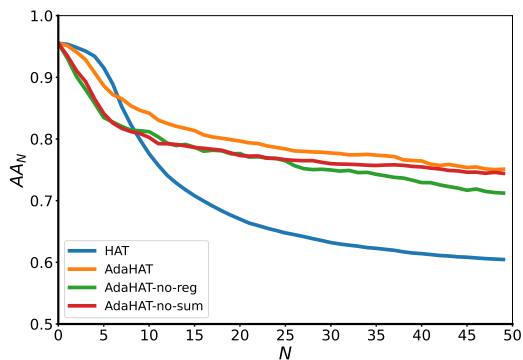
0 = all parameters can be updated freely

1 = no parameter can be updated

Results and Conclusions:

- ▶ HAT runs out of network capacity very soon at a fixed turning point (task 8)
- ▶ AdaHAT again behaves very similarly to HAT at first
- ▶ After the turning point, it manages it **adaptively** over time (through an adaptive adjustment rate), make it converge to 0 but never reach 0

Ablation Study



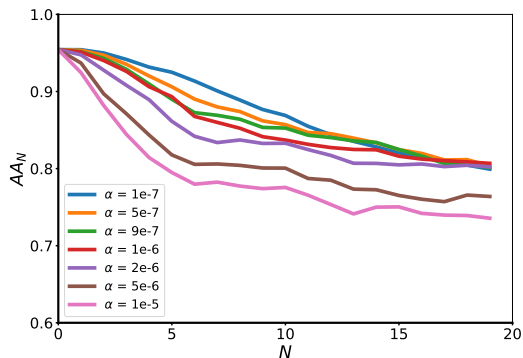
Ablation of two pieces of information:

- ▶ **AdaHAT-no-sum:** fix summative $\min(m_{l,i}^{<t,\text{sum}}, m_{l-1,j}^{<t,\text{sum}})$ at constant t
- ▶ **AdaHAT-no-reg:** fix regularization term $R(M^t, M^{<t})$ at constant 0

Results: both underperform AdaHAT but outperform HAT

Conclusions: both information (parameter importance, network sparsity) play crucial roles for an adaptive extension of HAT.

Hyperparameters



AdaHAT introduces only one additional hyperparameter:

α – overall intensity of gradient adjustment

Results: $\alpha = 10^{-6}$ is optimal.

Conclusions:

- ▶ Neither small nor large gradient adjustment balances the stability-plasticity trade-off, thus underperforms
- ▶ The optimal is still a relatively small value, indicating the importance to design a proper and well-guided adjustment rate

Conclusions

Conclusions

Existing architecture-based approaches (like HAT):

- ▶ Tend to tilt the stability-plasticity trade-off towards stability
- ▶ Suffer from insufficient network capacity problem in long sequence of tasks

Our proposed AdaHAT:

- ▶ Balances the trade-off in an adaptive adjustment mechanism
- ▶ Also retains maximum stability benefits before the network capacity limit
- ▶ Effectively leverages information about previous tasks which was seldom used in architecture-based approaches
- ▶ All of them leads to better performance than HAT

Future work:

- ▶ Explore and exploit more subtle information about previous tasks

Thank You

Thank you for your attention!

Please feel free to ask any questions or reach out to us at:

wangpengxiang@stu.pku.edu.cn

Project page:

<https://pengxiang-wang.com/projects/continual-learning-arena>

GitHub:

<https://github.com/pengxiang-wang/continual-learning-arena>

