### Introducing My PhD Work in 20 Minutes

#### Pengxiang Wang

Peking University, School of Mathematical Sciences

2025-10-21



# My Work in Continual Learning

My PhD research mainly focuses on continual learning with:

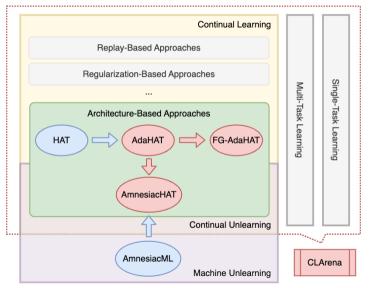
#### 3 papers:

- 1. AdaHAT: Adaptive Hard Attention to the Task in Task-Incremental Learning
- FG-AdaHAT: Fine-Grained Neuron Importance Guided Architecture-Based Continual Learning
- 3. **AmnesiacHAT**: Continual Unlearning via Capacity Recycling in Architecture-Based Continual Learning

#### 1 open-source software:

▶ **CLArena**: A Python package for continual learning research

#### Connections Among My Work



## My Work in Continual Learning

Work	Scope	Approach		
AdaHAT	AdaHAT Task-Incremental Learning			
FG-AdaHAT	Task-Incremental Learning	Architecture-Based		
AmnesiacHAT	Continual Unlearning	Architecture-Based		
CLArena	Continual Learning $/$ Unlearning	All Kinds		

Task-Incremental Learning (TIL): A typical setting in continual learning

**Continual Unlearning (CUL)**: A new area combining machine unlearning and continual learning, hardly explored before

**Architecture-Based Approach**: One of the main branches of continual learning approaches

# Continual Learning Background

# Definition of Continual Learning

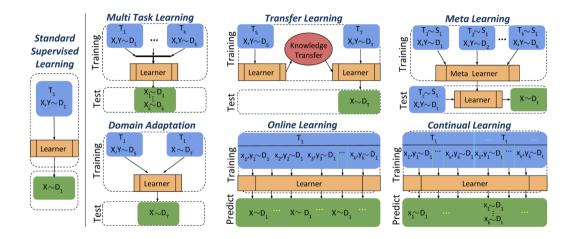
**Continual Learning (CL)**, which becomes popular after 2017, is a machine learning paradigm where an algorithm receives the data from tasks sequentially without the access to previous ones to learn a model that performs the best for all tasks.

(a.k.a. lifelong learning, incremental learning, sequential learning)

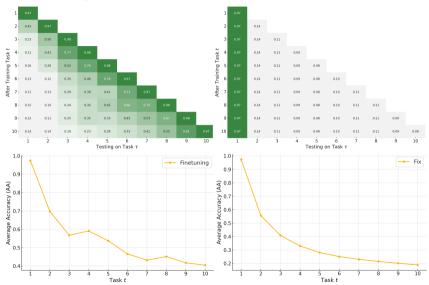
#### Key assumptions:

- Sequential tasks from different distributions
- Test for all tasks, but no access to previous task's data

#### Differences from Other Paradigms



## Baselines: Finetuning and Fix



## Challenges

#### Catastrophic Forgetting

- Previous knowledge can hardly be preserved within neural networks after training different tasks
- ▶ The main problem that most CL algorithms seek to address

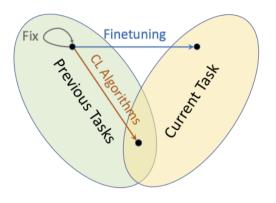
#### Stability-Plasticity Dilemma

- The model cannot achieve both stability and plasticity at the same time
- Need to trade off the balance of stability-plasticity

#### Network Capacity Problem

- Any fixed model will eventually get full as infinite tasks arrive
- Cannot select a proper-sized network beforehand under the infinite task assumption

## Challenges



For more about continual learning, check out my article: pengxiang-wang.com/posts/continual-learning-beginners-guide

## Solved Problems of My Work

#### AdaHAT / FG-AdaHAT:

- ▶ Balances the stability-plasticity trade-off in continual learning
- ▶ Alleviates the network capacity problem that architecture-based approaches generally suffer from
- Addresses the challenges of continual learning on long task sequences

#### AmnesiacHAT:

- Provides a solution to continual unlearning problem, which is: how to erase knowledge learned from specific tasks without affecting other learned tasks
- Addresses the network capacity problem from a new perspective: capacity recycling through unlearning

## Significance of My Work

#### AdaHAT / FG-AdaHAT:

- Provides effective architecture-based approaches for continual learning
- Proposes a unified gradient adjustment algorithm framework
- Paves the way for leveraging task information (especially fine-grained) in architecture-based approaches

#### AmnesiacHAT:

- Could be one of the foundational works in the new continual unlearning area
- Could be one of the first to explore the potential benefits from unlearning

#### CLArena:

- Another open-source package in the research community of continual learning
- Solid and robust codebase supporting my 3 papers above

I will introduce them as follows. You can ask questions at any time!

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

## Background: Architecture-Based Approaches

**Architecture-based approaches** (a.k.a. parameter isolation) are one of the main branches of continual learning algorithms:

- ▶ A distinctly different strategy that decomposes the network architecture
- Dedicate different parts of a neural network to different tasks

**HAT (Hard Attention to the Task)** is one of the most representative architecture-based approaches:

- ► Hard (binary) masks on layers
- Treat the masks as model parameters
- Masks condition on gradients directly

## Network Capacity Problem in HAT

HAT's hard gradient clipping mechanism allows no update for parameters masked by previous tasks:

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

More tasks come in

More active parameters become static

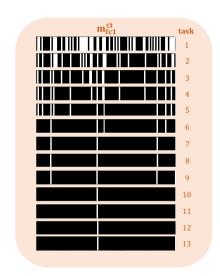
Insufficient network capacity



Significantly affect performance on new tasks



Learning plasticity reduced, inbalanced stability-plasticity trade-off



## AdaHAT: Adaptive Hard Attention to the Task

AdaHAT **soft-clips gradients**, allowing minor updates for parameters masked by previous tasks:

$$g'_{l,ij} = a^{\star}_{l,ij} \cdot g_{l,ij}, \ a^{\star}_{l,ij} \in [0,1]$$

The adjustment rate  $a_{l,ij}^{\star}$  now is an adaptive controller, guided by two pieces of information about previous tasks:

- ▶ Parameter Importance: how many tasks the target parameter has been used for
- ▶ Network Sparsity: how many parameters are currently used for all previous tasks

#### AdaHAT: Adaptive Hard Attention to the Task

Adaptive Adjustment Rate (AdaHAT)

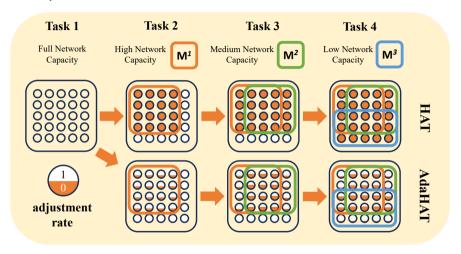
$$a_{l,ij}^{\star} = \frac{r_l}{\min\left(m_{l,i}^{< t, \mathsf{sum}}, m_{l-1,j}^{< t, \mathsf{sum}}\right) + r_l}, \ r_l = \frac{\alpha}{R\left(\mathsf{M}^t, \mathsf{M}^{< t}\right) + \epsilon}$$

Adjustment Rate (HAT)

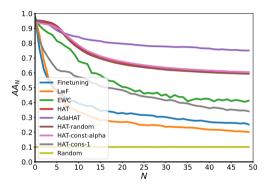
$$a_{l,ij} = 1 - \min\left(m_{l,i}^{< t}, m_{l-1,j}^{< t}\right)$$

For more about this work, check out my website: pengxiang-wang.com/papers/adahat

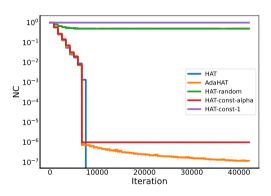
#### AdaHAT: Adaptive Hard Attention to the Task



 $HAT \Rightarrow AdaHAT$  (Adaptive HAT)



Results on long task sequences



Network capacity usage

# FG-AdaHAT: Fine-Grained Neuron Importance Guided Architecture-Based Continual Learning

#### Limitations of AdaHAT

The adaptive gradient clipping mechanism in AdaHAT is coarse:

- Parameter Importance is an integer value from 0 to t-1
- Network Sparsity is a single scalar value for the whole network

Fine-grained task information plays a crucial role among other continual learning approaches, meaningfully contributing to guiding the training of new tasks.

AdaHAT ⇒ **FG-AdaHAT** (Fine-Grained AdaHAT)

$$a_{l,ij}^{\mathsf{FG-AdaHAT}} = \alpha \left[ c^t \cdot \mathsf{Agg} \left( I_{l,i}^{< t}, I_{l-1,j}^{< t} \right) + \alpha \right]^{-1}, c^t = (t+b_L) \cdot \left[ R \left( \mathsf{M}^t, \mathsf{M}^{< t} \right) + b_R \right]$$

#### Fine-Grained Neuron Importance in FG-AdaHAT

**Neuron Importance** measures how important a neuron is to previous tasks, which is finer-grained than Parameter Importance and Network Sparsity.

$$I_{l,i}^{< t} = \sum_{\tau < t} I_{l,i}^{\tau}, \ I_{l,i}^{\tau} = \frac{1}{|D^{\tau}|} \sum_{(\mathbf{x}, y) \in D^{\tau}} I_{l,i}^{\tau}(\mathbf{x}, y)$$

It can be constructed from training information:

- ▶ Input Gradients (IG):  $I_{l,i}^{\tau}(\mathbf{x},y) = \sum_{i} |g_{l,ij}'|$
- lacksquare Output Gradients (OG):  $I_{l,i}^{ au}(\mathbf{x},y) = \sum_i |g_{l+1,ij}'|$
- Contribution Utility (CU):  $I_{l,i}^{\tau}(\mathbf{x},y) = |h_{l,j}(\mathbf{x})| \sum_{i} |w_{l+1,ij}|$
- ▶ Input Contribution Utility (ICU):  $I_{l,i}^{\tau}(\mathbf{x},y) = |h_{l,i}(\mathbf{x})| \sum_{i} |w_{l,ij}|$
- Activation Fisher Information (AFI):  $I_{l,i}^{\tau}(\mathbf{x},y) = |h_{l,i}(\mathbf{x})| \sum_{i} |g_{l,ij}'|^2$

#### Fine-Grained Neuron Importance in FG-AdaHAT

Layer attribution methods from **Explainable AI (XAI)** provides more fine-grained neuron importance measures:

- Feature Ablation (FA): perturbation-based attribution
- Internal Influence (II): gradient-based attribution
- ▶ Gradient SHAP (GS): gradient-based attribution, using Shapley values
- ▶ DeepLIFT (DL): backpropagation-based attribution

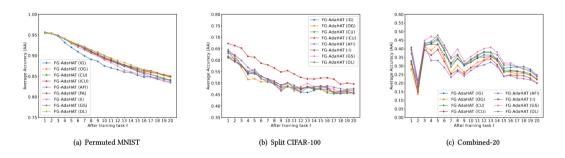
For more about this work, check out my website: pengxiang-wang.com/papers/fg-adahat

Approach	Permuted MNIST		Split CIFAR-100		Combined-20	
	AA (↑)	FR (↑)	AA (↑)	FR (↑)	AA (†)	FR (†)
Finetuning	$32.06 \pm 0.96$	$-73.19 \pm 1.11$	$34.19 \pm 1.84$	$-72.47 \pm 3.34$	$10.88 \pm 0.44$	$-74.45 \pm 1.72$
LwF	$32.95 \pm 1.35$	$-72.12 \pm 1.55$	$31.78 \pm 2.36$	$-77.20 \pm 5.34$	$10.92 \pm 0.91$	$-74.36 \pm 1.04$
HAT	$67.10 \pm 0.56$	$-30.68 \pm 0.65$	$40.12 \pm 2.17$	$-59.12 \pm 4.01$	$17.05 \pm 2.57$	$-74.63 \pm 4.16$
AdaHAT	$79.91 \pm 1.47$	$-12.43 \pm 6.46$	$44.93 \pm 1.29$	$-50.28 \pm 2.51$	$16.46 \pm 2.49$	$-68.31 \pm 4.74$
FG-AdaHAT (IG)	$84.05 \pm 1.71$	$-10.12 \pm 1.86$	$45.41 \pm 3.30$	$-49.36 \pm 5.87$	$24.67 \pm 4.15$	$-56.65 \pm 4.74$
FG-AdaHAT (OG)	$85.05 \pm 0.71$	$-8.92 \pm 0.82$	$45.54 \pm 2.19$	$-49.55 \pm 3.69$	$22.05 \pm 2.77$	$-59.11 \pm 4.17$
FG-AdaHAT (CU)	$84.11 \pm 0.91$	$-10.05 \pm 1.05$	$45.71 \pm 3.47$	$-49.12 \pm 6.42$	$22.26 \pm 1.69$	$-59.98 \pm 3.76$
FG-AdaHAT (ICU)	$84.97 \pm 0.63$	$-9.01 \pm 0.72$	$49.66 \pm 1.88$	$-40.93 \pm 3.75$	$19.95 \pm 3.93$	$-62.37 \pm 4.80$
FG-AdaHAT (AFI)	$83.43 \pm 1.28$	$-10.88 \pm 1.47$	$47.00 \pm 1.79$	$-46.66 \pm 3.31$	$20.09 \pm 5.16$	$-61.29 \pm 5.90$
FG-AdaHAT (FA)	$84.86 \pm 0.89$	$-9.14 \pm 1.03$	_	_	_	_
FG-AdaHAT (II)	$83.66 \pm 0.74$	$-10.60 \pm 0.85$	$45.70 \pm 1.66$	$-48.85 \pm 2.85$	$22.32 \pm 1.01$	$-59.10 \pm 4.24$
FG-AdaHAT (GS)	$83.91 \pm 0.70$	$-10.30 \pm 0.81$	$46.43 \pm 3.83$	$-47.40 \pm 7.32$	$24.25 \pm 6.52$	$-56.05 \pm 9.23$
FG-AdaHAT (DL)	$84.70 \pm 0.57$	$-9.34\pm0.66$	$47.70 \pm 4.64$	$-45.02 \pm 7.92$	$23.71 \pm 4.49$	$-57.79 \pm 5.34$

Performance on 3 continual learning benchmarks

Approach	Permuted MNIST BWT (↑) FWT (↑)		Split CIFAR-100 BWT (↑) FWT (↑)		Combined-20 BWT (↑) FWT (↑)	
	D ( ( )	1 ** 1 ( )	DW1()	1 ** 1 ( )	D W I ( )	1 ** 1 ( )
Finetuning	$-68.52 \pm 0.96$	$0.42 \pm 0.02$	$-51.17 \pm 1.87$	$10.70 \pm 0.40$	$-64.88 \pm 0.57$	$9.87 \pm 0.22$
LwF	$-67.42 \pm 1.34$	$0.25 \pm 0.05$	$-51.78 \pm 3.40$	$8.88 \pm 1.06$	$-64.75 \pm 0.81$	$9.70 \pm 0.21$
HAT	$-0.00 \pm 0.00$	$-31.12 \pm 0.57$	$-21.40 \pm 2.40$	$-12.43 \pm 0.72$	$-19.18 \pm 3.98$	$-27.31 \pm 4.40$
AdaHAT	$-12.46 \pm 1.42$	$-5.19 \pm 0.09$	$-27.23 \pm 1.06$	$-1.60 \pm 0.97$	$-40.91 \pm 2.29$	$-6.16 \pm 1.29$
FG-AdaHAT (IG)	$-7.73 \pm 1.67$	$-5.55 \pm 0.17$	$-19.90 \pm 3.13$	$-8.13 \pm 0.82$	$-33.90 \pm 2.45$	$-5.11 \pm 2.62$
FG-AdaHAT (OG)	$-5.33 \pm 0.70$	$-6.90 \pm 0.19$	$-19.66 \pm 2.08$	$-8.12 \pm 1.01$	$-34.55 \pm 2.13$	$-6.56 \pm 1.80$
FG-AdaHAT (CU)	$-6.09 \pm 0.90$	$-7.14 \pm 0.11$	$-22.61 \pm 2.62$	$-5.00 \pm 0.77$	$-35.31 \pm 1.57$	$-5.79 \pm 0.71$
FG-AdaHAT (ICU)	$-5.57 \pm 0.67$	$-6.75 \pm 0.11$	$-23.60 \pm 1.70$	$-0.41 \pm 0.69$	$-38.10 \pm 2.72$	$-5.37 \pm 1.55$
FG-AdaHAT (AFI)	$-7.32 \pm 1.24$	$-6.62 \pm 0.14$	$-21.92 \pm 2.23$	$-4.40 \pm 1.18$	$-37.06 \pm 2.86$	$-6.58 \pm 2.71$
FG-AdaHAT (FA)	$-6.71 \pm 0.87$	$-5.73 \pm 0.09$	_	_	_	_
FG-AdaHAT (II)	$-6.02 \pm 0.70$	$-7.69 \pm 0.12$	$-20.57 \pm 1.54$	$-7.01 \pm 1.58$	$-34.48 \pm 4.14$	$-6.59 \pm 3.23$
FG-AdaHAT (GS)	$-4.75 \pm 0.69$	$-8.68 \pm 0.25$	$-20.92 \pm 2.61$	$-6.06 \pm 1.83$	$-34.58 \pm 4.19$	$-4.83 \pm 1.80$
FG-AdaHAT (DL)	$-3.89 \pm 0.57$	$-8.70\pm0.25$	$-20.03 \pm 3.92$	$-5.63 \pm 1.11$	$-34.65 \pm 2.26$	$-5.29\pm1.70$

Stability-plasticity trade-off on 3 continual learning benchmarks



Comparison of 9 fine-grained neuron importance on 3 benchmarks

# AmnesiacHAT: Continual Unlearning via Capacity Recycling in Architecture-Based Continual Learning

## Background: Continual Unlearning

**Machine unlearning**, which becomes popular after 2020, deliberately removes the influence of specific data from a trained Al model.

Applications: data privacy, security, data quality, fairness, etc.



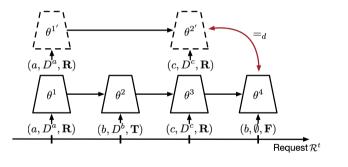
- ▶ AI models like neural networks are usually holistic, where unlearning causes side effects
- Retraining from scratch without the unlearning data is:
  - 1. Computationally expensive
  - 2. Not always feasible when the original data is unavailable



# Background: Continual Unlearning

Continual unlearning: Unlearning specific tasks from a continual learning model

(A field that has been hardly explored before)

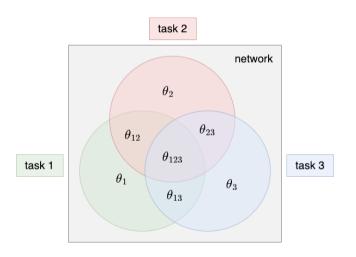


#### Objectives:

- 1. **Learning** objective: Maximal average performance on the remaining tasks
- 2. **Unlearning** objective: Minimal difference from retraining without the unlearning tasks

For more about unlearning, check out my slides: pengxiang-wang.com/slides/slides-continual-learning-and-machine-unlearning.pdf

#### Motivation



Architecture-based approaches are naturally suitable for continual unlearning:

- Less overlapping, less task interference
- We choose our state-of-the-art fix-sized architecture-based approach AdaHAT

AdaHAT ⇒ AmnesiacHAT

#### Unlearning Algorithm in AmnesiacHAT

The AdaHAT architecture needs to be modified to enable unlearning:

1. Store task-wise update record (inspired by AmnesiacML):

$$\theta_{l,ij}^{(t)} = \theta_{l,ij}^{(0)} + \Delta \theta_{l,ij}^{(1)} + \dots + \Delta \theta_{l,ij}^{(t)}$$

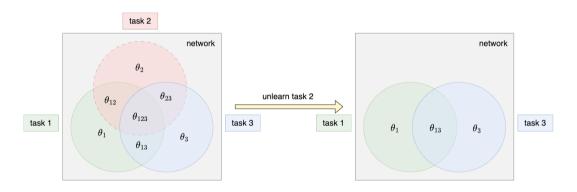
The update in HAT is sparse which can be stored efficiently by sparse storage format and pruning.

- 2. Unlearn task au= Remove its parameter update  $\Delta \theta_{l,ij}^{( au)}$
- 3. Deal with the side effects on overlapping parameters

If the parameter is first used for task  $\tau$ , compensate the adjustment rate back in its next task's update. Small-scale retraining for fixing is applied after.

# Unintended Benefit: Capacity Recycling

When a task is unlearned, the parameters that were dedicated to that task become free and can be reused for future tasks.

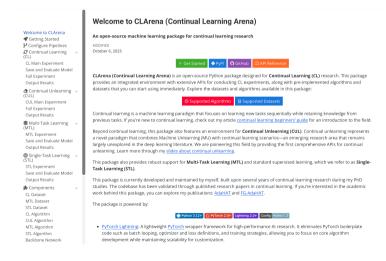


How to measure it? Compare the performance without unlearning

# Continual Learning Arena (CLArena)

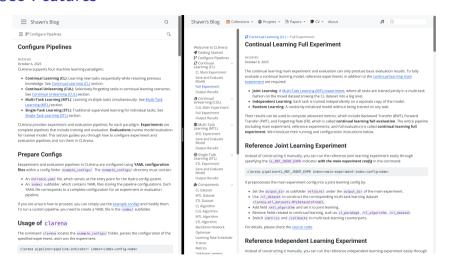
## Continual Learning Arena (CLArena)

**Continual Learning Arena (CLArena)** is an open-source Python package for continual learning research, developed and documented entirely by myself.

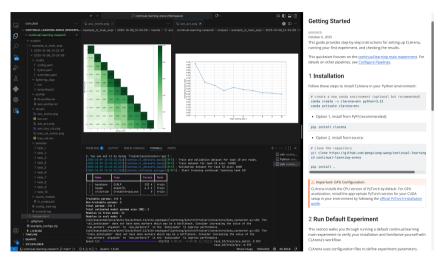


#### Key Features of CLArena

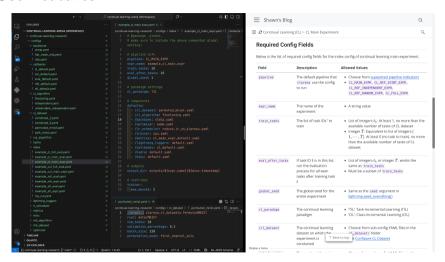
- 1. Provide pipelines for various machine learning paradigms:
- Continual Learning (CL)
- Continual Unlearning (CUL)
- Multi-Task Learning (MTL)
- ► Single-Task Learning (STL)
- 2. Provide implementations of:
- Mainstream CL algorithms: LwF, EWC, HAT, WSN, CBP, etc.
- Common CL / MTL / STL benchmarks: permuted, split, combined, etc.
- Neural network architectures: MLP, ResNet, HAT-based models, etc.
- 3. Provide framework for building custom datasets, models, algorithms, metrics, callbacks, etc.
- 4. Full code support for my 3 papers above, ensuring reproducibility



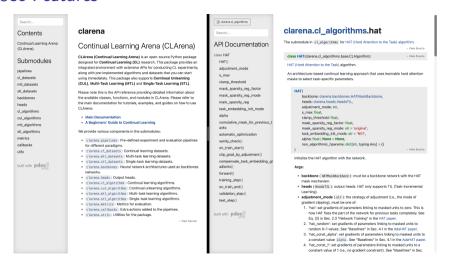
Well-structured and detailed documentation



Get started quickly with tutorials and examples



An easy configuration system based on YAML, each config field clearly explained



Complete API reference for custom implementation in the framework

## Package Information

```
Main Page / Documentation
pengxiang-wang.com/projects/continual-learning-arena
GitHub (feel free to give it a star! :D)
github.com/pengxiang-wang/continual-learning-arena
PyPI ('pip install clarena')
pypi.org/project/clarena
```

**API** Reference

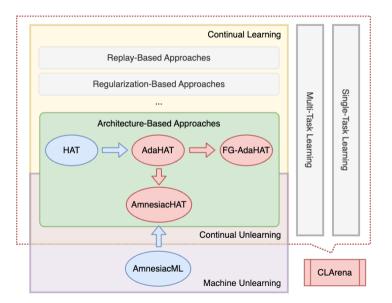
pengxiang-wang.com/projects/continual-learning-arena/docs/api-reference

#### Thank You

Overall, my PhD work dives deep into continual learning through architecture-based approaches, across comprehensive contributions to this area in both width and depth:

- New well-performed continual learning approaches: AdaHAT, FG-AdaHAT variants
- New algorithm framework: **FG-AdaHAT**
- New paradigm problem hardly explored before: continual unlearning, and its solution: AmnesiacHAT
- New software platform for continual learning research: **CLArena**

#### Thank You



#### Thank You

Thank you for your attention!

My website: pengxiang-wang.com

My email: wangpengxiang@stu.pku.edu.cn

