# Enhancing Large Language Models for a Dynamic World

**Presenter:** Zixuan Ke
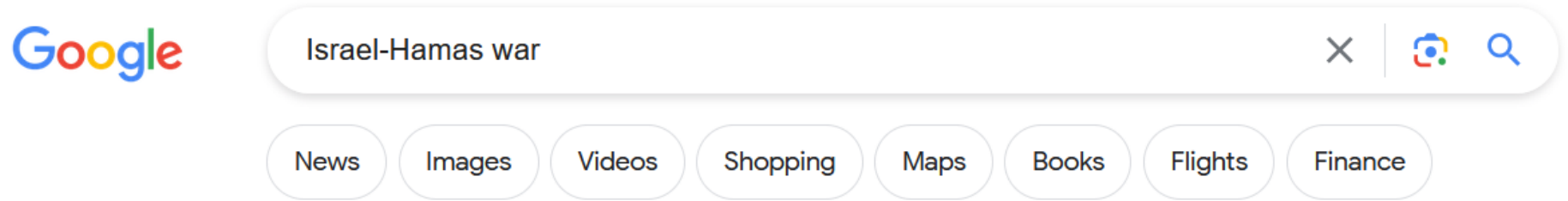https://vincent950129.github.io/

# LLM



Packed with
knowledge and excels
in many tasks

# LLM in A Fixed World



LLMs — Packed with knowledge and excels in many tasks

ASSUMPTION

The world is **fixed** (i.i.d)

# The World Changes Quickly

# LLM in A Dynamic World



Packed with knowledge and excels in many tasks



The world is **ever-changing**

# LLM in A Dynamic World

How to make knowledge in LLM more **reusable** and **updatable** in the **dynamic** world?

Packed with knowledge and excels in many tasks

The world is **ever-changed**

# LLM in A Dynamic World

How to make knowledge
in LLM more **reusable**
and **updatable** in the
**dynamic** world?

External
Memory

LLMs

Tool

# Retrieval-augmented LLM



**Who is the CEO of Twitter?**

As of my **knowledge cutoff in September 2021**, the CEO of Twitter is **Jack Dorsey**....

Google — Who is the CEO of Twitter?

All    News    Images    Shopping    Videos    More    Tools

About 1,090,000,000 results (0.45 seconds)

Twitter / CEO

**Linda Yaccarino**

Jun 5, 2023–

- The datastore can be easily **updated** and **expanded** - even without retraining!

Datastore
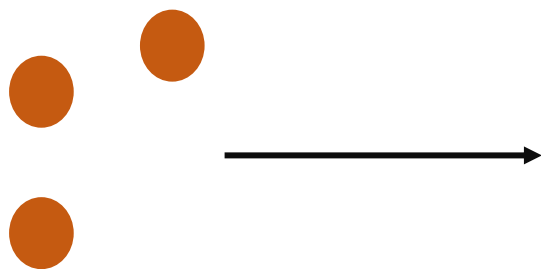
New!

Input

Query → Index → LM +

# LLM in A Dynamic World

How to make knowledge in LLM more **reusable** and **updatable** in the **dynamic** world?

# Continual Learning

What will happen if we **update the LLM** in a **changing world**?

$X$ $\qquad$ $f_\theta: X \rightarrow Y$ $\qquad$ $Y$

$(t)$

$X^{(1)}$

LLMs

$Y^{(1)}$

$f_\theta \colon X^{(1)} \to Y^{(1)}$

$X^{(2)}$

$Y^{(2)}$

Another dimension, can be
**Task** (sentiment classification, news classification...)
**Domain** (restaurant, phone, camera...)
**Class** (dog, cat, car, horse...)
**Time** (2020, 2023...)
etc.

$(t)$

$X^{(1)}$

$Y^{(1)}$

$f_\theta: X^{(1)} \rightarrow Y^{(1)}$

$X^{(2)}$

$Y^{(2)}$

The dimension *t* can be called "task" in continual learning

$(t)$

$X^{(1)}$

LLMs

$Y^{(1)}$

$X^{(2)}$

$f_\theta: X^{(1)} \cup X^{(2)} \rightarrow Y^{(1)} \cup Y^{(2)}$

$Y^{(2)}$

$(t)$

$\mathrm{X}^{(1)}$

$Y^{(1)}$

$f_\theta: \mathrm{X}^{(1)} \cup \mathrm{X}^{(2)} \rightarrow Y^{(1)} \cup Y^{(2)}$

$\mathrm{X}^{(2)}$

$Y^{(2)}$

Cons
- Train from scratch whenever update is needed
- Save all the past data
- Privacy concern (user may not want to share their personal data)
- etc.

$(t)$

$X^{(1)}$

$f_\theta^{(1)}: X^{(1)} \rightarrow Y^{(1)}$

$Y^{(1)}$

$X^{(2)}$

$f_\theta^{(2)}: X^{(2)} \rightarrow Y^{(2)}$

$Y^{(2)}$

$f_\theta^{(2)} : X^{(2)} \rightarrow Y^{(2)}$

Cons {
Models are isolated. They are hard to help each other
Need to know the task belonging in testing
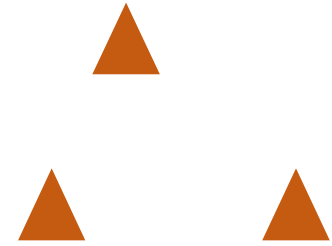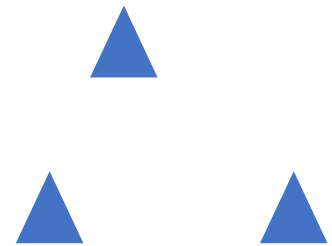Need to train a new model whenever update is needed
etc.
}

17

$(t)$

$X^{(1)}$

$X^{(2)}$

LLMs

$f_\theta: X^{(2)} \rightarrow Y^{(2)}$

$Y^{(1)}$

$Y^{(2)}$

18

$(t)$

$X^{(1)}$

$X^{(2)}$

LLMs

$Y^{(1)}$

$Y^{(2)}$

$f_\theta : X^{(2)} \to Y^{(2)}$

Challenges:
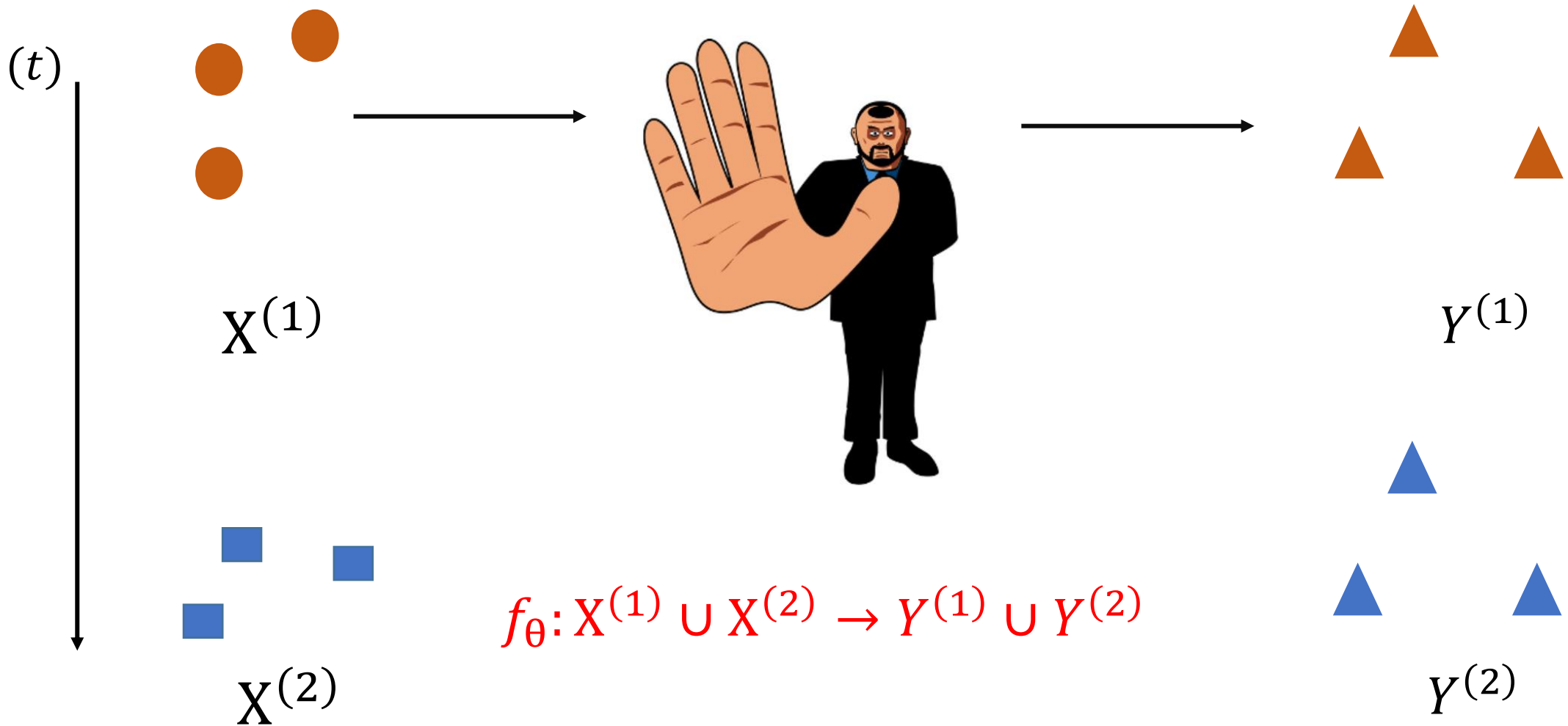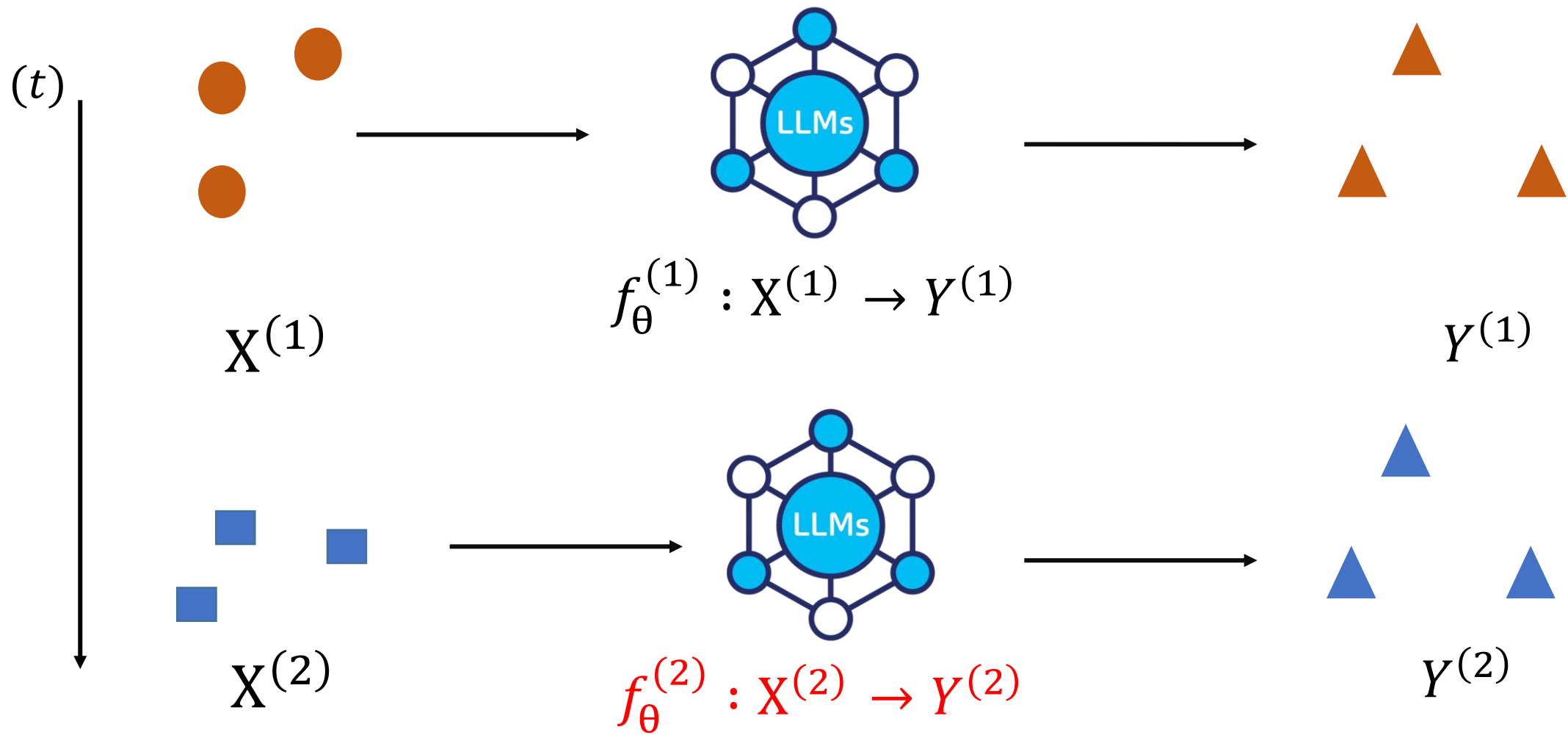**Catastrophic Forgetting (CF)**
**Knowledge Transfer (KT)**

# Catastrophic Forgetting

- Simple case:
  - 2 features
  - Learn a line in 2D plane
- You learn a perfect line for task 1

$(X^{(1)}, Y^{(1)})$

Dog

Cat

# Catastrophic Forgetting

- Simple case:
  - 2 features
  - Learn a line in 2D plane

- **Update** the learned parameters, and learn **another** perfect separate line for task 2

$(X^{(2)}, Y^{(2)})$

Car

Horse

# Catastrophic Forgetting

- Simple case:
  - Evaluate the final model on **both** learned tasks (assuming task id unknown)

- After learning a second task, you **forget** how to deal with the first task!



Dog
Cat
Car
Horse

**Catastrophic Forgetting (CF)**

# Knowledge Transfer

- Simple case:
  - 2 features
  - Learn a line in 2D plane
- This time, you learn an **imperfect** line for task 1

$(X^{(1)}, Y^{(1)})$

Dog
Cat

# Knowledge Transfer

- Simple case:
  - 2 features
  - Learn a line in 2D plane

- Task 2 has **similar label** as task 1, but the input images becomes **binary**

- You learn a **perfect** line for task 2

$(\mathrm{X}^{(2)}, Y^{(2)})$

Dog

Cat

# Knowledge Transfer

- Simple case:
  - Evaluate the final model on **both** learned tasks

- After Learning a second task, the **old task improved**

- Because the knowledge from task 2 is **helpful** to task 1

Dog
Cat
Dog
Cat

**Knowledge Transfer (KT)**

# Continual Learning

How to
(1) **Mitigate forgetting**, i.e., perform reasonably well on what has been learned
(2) **Knowledge transfer**, i.e., relevant tasks can help each other

# Enhancing LLM for A Dynamic World

How to make knowledge in LLM more **reusable** and **updatable**?



Continual Post-training of Language Models, Ke et al., ICLR 2023

# Plan

- Motivation
- Introduction
  - Continual Learning
- Continual Post-training of Language model
- Conclusion and future work

# Continual Post-training of Language Model

Pre-training



$t$       $(t = 0)$

Huge amount
of general data

# Continual Post-training of Language Model

Pre-training

Continual (Domain-adaptive) Pre-training
Post-training / Pre-finetuning
(Our focus)

Restaurant

Phone

Camera

......

$t$      $(t = 0)$      $(t = 1)$      $(t = 2)$      $(t = 3)$

Domain-specific data

# Continual Post-training of Language Model

Pre-training

Continual (Domain-adaptive) Pre-training
Post-training / Pre-finetuning
(Our focus)

Restaurant          Phone          Camera          ......

$t$     $(t = 0)$     $(t = 1)$     $(t = 2)$     $(t = 3)$

Current task

Accessibility

MLM Head

## (A) <span style="color:red">Continual</span> Post-training

Add & Layer Norm ×L

+

FFN

+

Add & Layer Norm

+

Attention

Hidden States

(We use RoBERTa in this work)

Restaurant  Phone  Camera

**First,** we continually post-trains **a sequence of domains**

## (A) **Continual** Post-training

MLM Head

| Add & Layer Norm | × L |

FFN

⊕

Add & Layer Norm

⊕

Attention

Hidden States

**First,** we continually pre-trains **a sequence of domains**

Restaurant  Phone  Camera

---

## (B) **Individual** Fine-tuning

Classification Head

| Add & Layer Norm | × L |

⊕

FFN

⊕

Add & Layer Norm

⊕

Attention

Hidden States

End-tasks

ASC-Restaurant

ASC-Phone

ASC-Camera

ASC: Aspect Sentiment Classification

**After (A),** the performance is **evaluated** by end-tasks

Each end-task **corresponding** to one domain and has its **own** training and testing set. It is trained individually and **will not** affect the continual learning

33

# Continual Post-training of Language Model

**6 domains**

| Unlabelde Domain Datasets | | | End-Task Classification Datasets | | | | |
|---|---|---|---|---|---|---|---|
| Source | Dataset/Domain | Size | Dataset/Domain | Task | #Training | #Testing | #Classes |
| | Yelp Restaurant | 758MB | Restaurant | Aspect Sentiment Classification (ASC) | 3,452 | 1,120 | 3 |
| Reviews | Amazon Phone | 724MB | Phone | Aspect Sentiment Classification (ASC) | 239 | 553 | 2 |
| | Amazon Camera | 319MB | Camera | Aspect Sentiment Classification (ASC) | 230 | 626 | 2 |
| | ACL Papers | 867MB | ACL | Citation Intent Classification | 1,520 | 421 | 6 |
| Academic Papers | AI Papers | 507MB | AI | Relation Classification | 2,260 | 2,388 | 7 |
| | PubMed Papers | 989MB | PubMed | Chemical-protein Interaction Prediction | 2,667 | 7,398 | 13 |

**Continual post-training**

**Individual Fine-tuning**

# Continual Post-training of Language Model

- Setting
  - Post-train a sequence of domains **without** access to the data that used in **pre-training** and **previously learned domains**
  - End-task doesn't know its domain belonging
- Goals
  - CF prevention
  - KT (backward and forward)
- Related Work
  - There are CL and Post-training work **but** none directly on continual post-training.
- Approach
  - **C**ontinual **P**ost-training with **S**oft-masking (**CPS**)

# Continual Post-training of Language Model

$L_{\mathrm{MLM}}$

Transformer Layer $l$

× L

**Forward**

$\nabla_l$

Transformer Layer $l$

× L

**Backward**

Sequence of domains

Restaurant    Phone    Camera    ● ● ●

# Continual Post-training of Language Model

$L_{\mathrm{MLM}}$

⋮

Transformer Layer $l$   × L

**Forward**

$\nabla_l$   × L

Transformer Layer $l$

**Backward**

**1st Issue:** CF on the general knowledge

General knowledge means the knowledge in the **original pre-trained** LM

The knowledge learned from each domain alone **will not be sufficient** to recover it and give good end-task performances

# Continual Post-training of Language Model

$L_{\text{MLM}}$

⋮



**Forward**



**Backward**

**2nd Issue:** CF on the previously learned domain knowledge

Because we post-train a sequence of domains

# Continual Post-training of Language Model



Pre-trained
LM

# Continual Post-training of Language Model



$t = 1$

No training

**Importance Computation**

**Key Idea**

1) Detect importance of units (attention heads and neurons) for general and domain knowledge

**KEY TECHNOLOGY**

1) How to detect importance for the two types of knowledge

# Continual Post-training of Language Model



LLMs

$t = 1$

No training

Post-training

**Importance Computation**

**Soft-masking**

Backward

Key Idea

1) Detect importance of units for general and domain knowledge

2) Soft-masking the important units when training new tasks

KEY TECHNOLOGY

1) How to detect importance for the two types of knowledge

2) How to soft-mask

# Continual Post-training of Language Model



**Goal:** Compute the importance of units for **general** (and domain) knowledge

**Why?**
1) Not all units are important
2) Given the important units, we can protect them afterward

No training involved. We only need the importance

# Importance Computation



$L_{\mathrm{MLM}}$

⋮

Transformer Layer $l$

$\times$ L

**Forward**

To compute the importance

Element-wise multiplication

⋮

$g_l$ ⊗

Transformer Layer $l$

$\times$ L

**Forward**

**First,** we added **virtual parameters** $\boldsymbol{g}_l$.

$\boldsymbol{g}_l$ is the **virtual parameters.** Each virtual parameter $g_{l,i}$ in $\boldsymbol{g}_l$ corresponding to an attention head or neurons (units)

It is **initialized as all 1's**, and has its gradient but will **never change**.

**Why?** We only use its gradient to compute importance

# Importance Computation



(the loss for importance computation)

$L_{\mathrm{MLM}}$

$L_{\mathrm{impt}}$

× L

To compute the importance

Transformer Layer $l$

**Forward**

$g_l \rightarrow \otimes$    × L

Transformer Layer $l$

**Forward**

The gradient of $\boldsymbol{L}_{impt}$ w.r.t $\boldsymbol{g}_l$ will be used to compute importance.

# Importance Computation

$L_{\mathrm{MLM}}$

$\vdots$

Transformer Layer $l$ $\times$ L

**Forward**

To compute the importance

$\Longrightarrow$

$L_{\mathrm{impt}}$

$\vdots$

$g_l \rightarrow \otimes$ $\times$ L

Transformer Layer $l$

**Forward**

For **domain knowledge,**

$L_{\mathrm{impt}} = L_{\mathrm{MLM}}$

$$\boldsymbol{\nabla}^m_{\boldsymbol{g}_l} = \frac{\partial L_{\mathrm{impt}}(\boldsymbol{x}^{(t)}_m, \boldsymbol{y}^{(t)}_m)}{\partial_{\boldsymbol{g}_l}}$$

$$\boldsymbol{I}^{(t)}_l = \frac{1}{M}\sum_M |\boldsymbol{\nabla}^m_{\boldsymbol{g}_l}|$$

Use **absolute gradient** to indicate importance[1]

[1] Michel et al. Are sixteen heads really better than one? NeurIPS, 2019.

# Importance Computation

**However,** for **general knowledge**, we cannot do $L_{\mathrm{impt}} = L_{\mathrm{MLM}}$ as we do not have the pre-training data.

We need another $L_{\mathrm{impt}}$

# Importance Computation



Random dropout

# Importance Computation



Random dropout

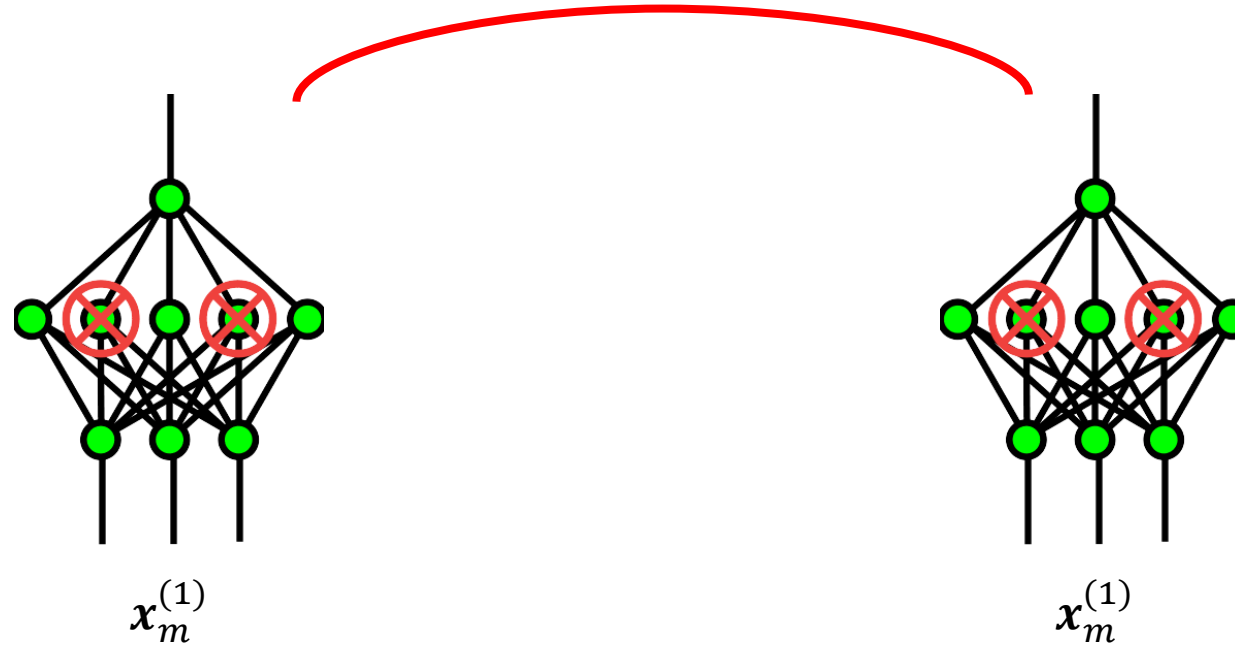Random noise

# Importance Computation

Same input but different output representation
The distance indicate the **robustness**



$$x_m^{(1)}$$

$$x_m^{(1)}$$

# Importance Computation

their changes will cause the pre-trained LM to change significantly

Units that are important to the robustness

Units that are important to the pre-trained/general knowledge

# Importance Computation

$$L_{\text{impt}} = \text{KL}(f_{\text{LM}}^1(x_m^{(1)}), f_{\text{LM}}^2(x_m^{(1)}))$$



$f_{\text{LM}}^1$ Pre-trained LM    × L

Transformer Layer $l$

$x_m^{(1)}$

$f_{\text{LM}}^2$ Pre-trained LM    × L

Transformer Layer $l$

$x_m^{(1)}$

Based on the intuition, we propose another $L_{\text{impt}}$, which do not need pre-training data

**KL**: how different given two representations

$f_{LM}^1 / f_{LM}^2$: Transformer with different dropouts

$x_m^{(1)}$ : We only use first domain data because we want the importance of units for the pre-trained knowledge

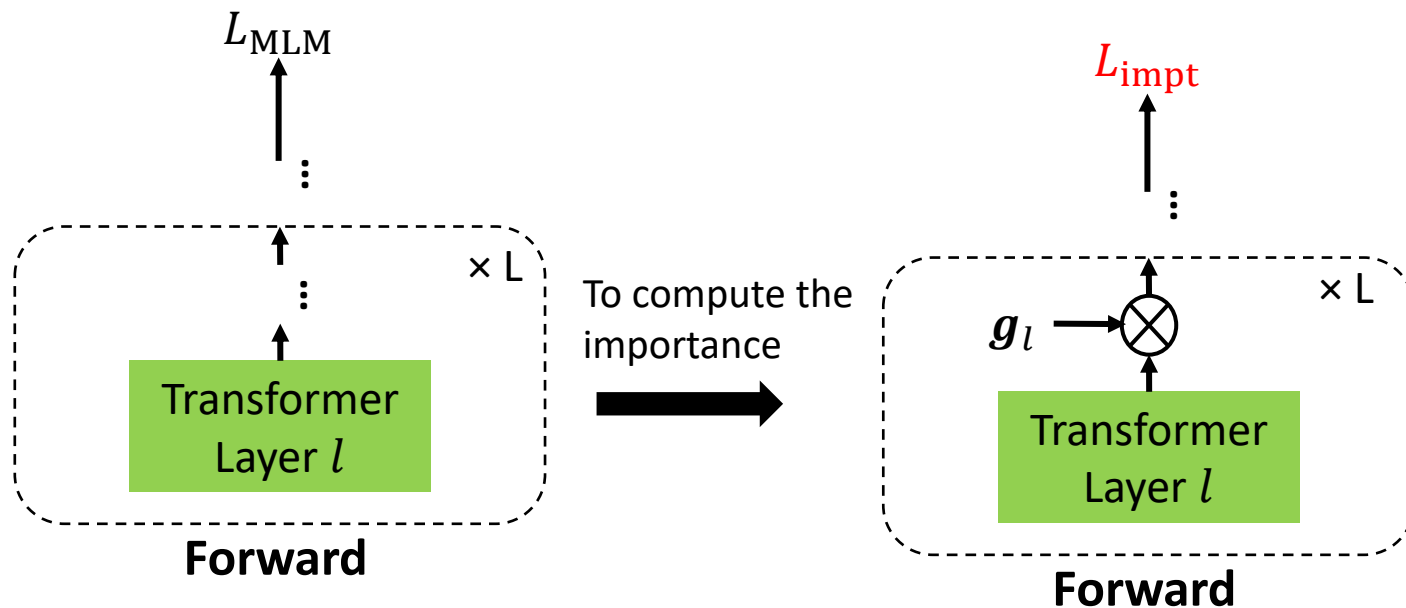# Importance Computation



$L_{\mathrm{MLM}}$

$L_{\mathrm{impt}}$

× L

Transformer Layer $l$

**Forward**

To compute the importance

$g_l$

× L

Transformer Layer $l$

**Forward**

For **general knowledge,**

$$L_{\mathrm{impt}} = \mathsf{KL}(f_{\mathrm{LM}}^1(\boldsymbol{x}_m^{(1)}), f_{\mathrm{LM}}^2(\boldsymbol{x}_m^{(1)}))$$

$$\boldsymbol{\nabla}_{\boldsymbol{g}_l}^m = \frac{\partial L_{\mathrm{impt}}(\boldsymbol{x}_m^{(1)})}{\partial_{\boldsymbol{g}_l}}$$

$$\boldsymbol{I}_l^{(0)} = \frac{1}{M}\sum_M |\boldsymbol{\nabla}_{\boldsymbol{g}_l}^m|$$

Importance of units for general knowledge

# Continual Post-training of Language Model

No training

Importance Computation

$\{I_l^{(k)}\}_{k=1}^{t-1}$

Post-training

Soft-masking

Backward
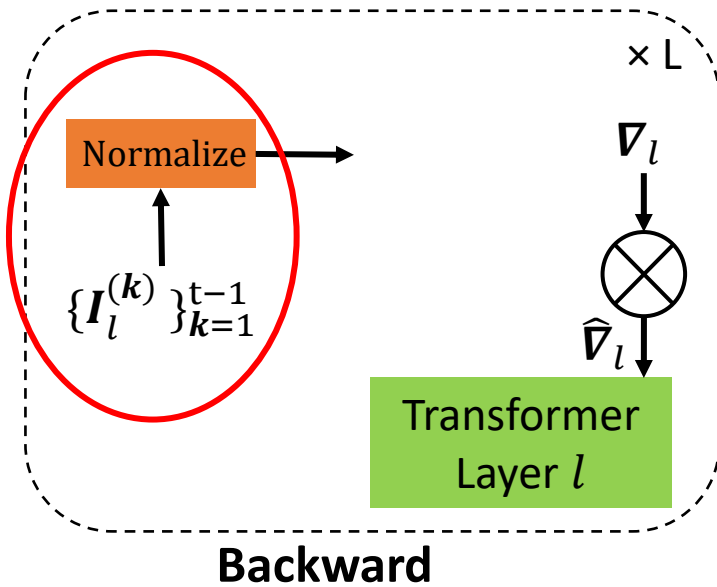
**Goal:** Soft-mask the **gradient** based on the importance

**Why?**
1) We need to protect them when training new domain
2) We want to allow knowledge transfer

# Soft-masking



**Backward**

First, we normalized the importance so that they are comparable

$$I_l^{(k)} = |\text{Tanh}(\text{Norm}(I_l^{(k)}))|$$

# Soft-masking



First, we normalized the importance so that they are comparable

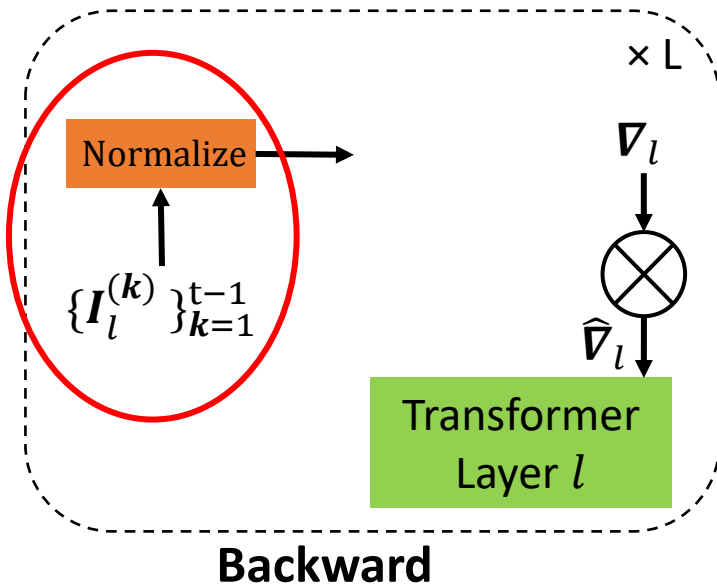$$I_l^{(k)} = |\text{Tanh}(\text{Norm}(I_l^{(k)}))|$$   make sure the importance is [0,1]

# Soft-masking



**Backward**

First, we normalized the importance so that they are comparable
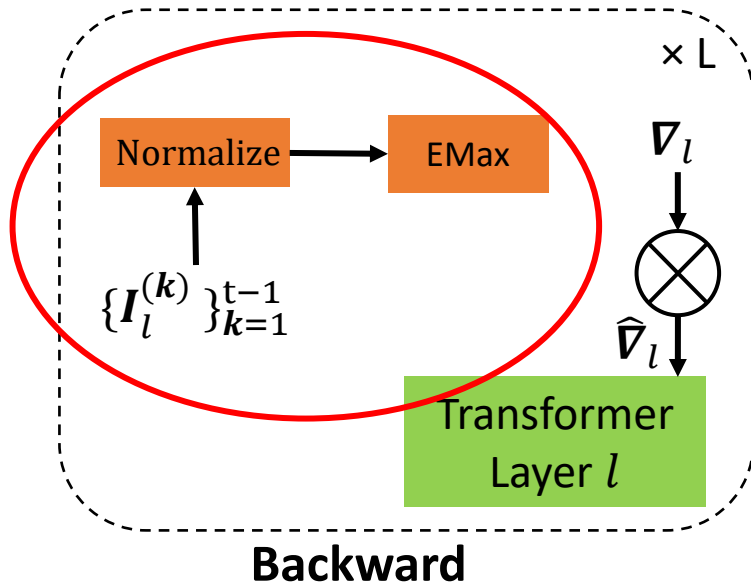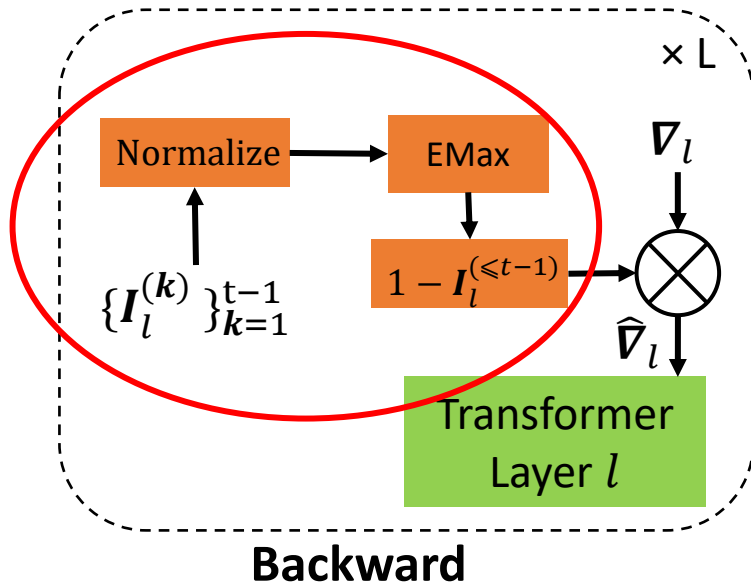
$$I_l^{(k)} = |\text{Tanh}(\text{Norm}(I_l^{(k)}))|$$

Second, we accumulate all importance before current domain $t$

$$I_l^{(\leqslant t-1)} = \text{EMax}(\{I_l^{(t-1)}, I_l^{(t-2)}\})$$

# Soft-masking



**Backward**

First, we normalized the importance so that they are comparable

$$\boldsymbol{I}_l^{(k)} = |\mathrm{Tanh}(\mathrm{Norm}(\boldsymbol{I}_l^{(k)}))|$$

Second, we accumulate the importance

$$\boldsymbol{I}_l^{(\leqslant t-1)} = \mathrm{EMax}(\{\boldsymbol{I}_l^{(t-1)}, \boldsymbol{I}_l^{(t-2)}\})$$

Third, we soft-mask the gradient (in backward pass)

$$\boldsymbol{\nabla'}_l = \left(1 - \boldsymbol{I}_l^{(\leqslant t-1)}\right) \otimes \boldsymbol{\nabla}_l$$

# Continual Post-training of Language Model

No training

<span style="color:red">No training</span>

Importance Computation

$t \mathrel{+}= 1$

$\{I_l^{(k)}\}_{k=1}^{t-1}$

Post-training

Soft-masking

$$\boldsymbol{\nabla'}_l = \left(1 - \boldsymbol{I}_l^{(\leqslant t-1)}\right) \otimes \boldsymbol{\nabla}_l$$

Backward

**Initialization**

**(A)**

$\mathrm{KL}(f_{\mathrm{LM}}^1(\boldsymbol{x}_m^{(1)}), f_{\mathrm{LM}}^2(\boldsymbol{x}_m^{(1)}))$

KL loss as $L_{\mathrm{impt}}$

$\widehat{\boldsymbol{o}}_l$

$\boldsymbol{g}_l \longrightarrow \bigotimes$

$\boldsymbol{o}_l$

Transformer Layer $l$

**Forward**

$\frac{1}{M}\sum_{\boldsymbol{M}}|\nabla_{\boldsymbol{g}_l}^m| \longrightarrow \boldsymbol{I}_l^{(0)}$

$\nabla_{\boldsymbol{g}_l}$

Transformer Layer $l$

Use gradient to indicate importance, but the gradient does not optimize the layer

$\boldsymbol{I}_l^{(0)}$ indicates the importance for general knowledge

**Backward**

59

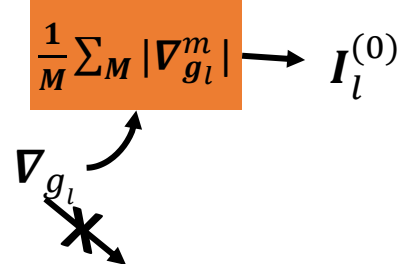**Initialization**

**(A)**

$$\text{KL}(f_{\text{LM}}^1(\boldsymbol{x}_m^{(1)}), f_{\text{LM}}^2(\boldsymbol{x}_m^{(1)}))$$

⋮

$\widehat{\boldsymbol{o}}_l$

$\boldsymbol{g}_l \longrightarrow \otimes$

$\boldsymbol{o}_l$

Transformer Layer $l$

**Forward**

$\frac{1}{M}\sum_M |\nabla_{g_l}^m| \longrightarrow \boldsymbol{I}_l^{(0)}$

$\nabla_{g_l}$

Transformer Layer $l$

**Backward**

**Continual Learning**

Next, we start **continual learning**

60

**Initialization**

**(A)**

$$\text{KL}(f_{\text{LM}}^1(\boldsymbol{x}_m^{(1)}), f_{\text{LM}}^2(\boldsymbol{x}_m^{(1)}))$$

$\hat{\boldsymbol{o}}_l$

$\boldsymbol{g}_l \longrightarrow \otimes$

$\boldsymbol{o}_l$

Transformer Layer $l$

**Forward**

$\frac{1}{M}\sum_{\boldsymbol{M}}|\boldsymbol{\nabla}_{g_l}^m| \longrightarrow \boldsymbol{I}_l^{(0)}$

$\boldsymbol{\nabla}_{g_l}$

Transformer Layer $l$

**Backward**

**Continual Learning**

**(B)**　　　　　　　　　　**(C)**

$L_{\text{MLM}}$

Transformer Layer $l$

**Forward**

Nothing changed in forward pass

**Initialization**

**Continual Learning**

**(A)**

$\mathrm{KL}(f_{\mathrm{LM}}^1(\boldsymbol{x}_m^{(1)}), f_{\mathrm{LM}}^2(\boldsymbol{x}_m^{(1)}))$

$\hat{\boldsymbol{o}}_l$

$\boldsymbol{g}_l \longrightarrow \bigotimes$

$\boldsymbol{o}_l$

Transformer Layer $l$

**Forward**

$\frac{1}{M}\sum_{\boldsymbol{M}}|\boldsymbol{\nabla}_{\boldsymbol{g}_l}^m| \longrightarrow \boldsymbol{I}_l^{(0)}$

$\boldsymbol{\nabla}_{\boldsymbol{g}_l}$

Transformer Layer $l$

**Backward**

**(B)**

$L_{\mathrm{MLM}}$

Transformer Layer $l$

**Forward**

Normalize $\longrightarrow$ EMax

$\{\boldsymbol{I}_l^{(k)}\}_{\boldsymbol{k}=1}^{\mathrm{t}-1}$

$1 - \boldsymbol{I}_l^{(\leqslant t-1)} \longrightarrow \bigotimes$

$\boldsymbol{\nabla}_l$

$\hat{\boldsymbol{\nabla}}_l$

Transformer Layer $l$

**Backward**

**(C)**

Accumulate all importance of units
Use it to soft-mask the gradient

# Evaluation

Goals

CF Prevention

Knowledge Transfer

Metrics

Forgetting Rate

Final Performance

# Metrics

Restaurant

Restaurant

$A_{1,1}$

# Metrics

Restaurant    ACL

Restaurant

ACL

$A_{1,1}$

$A_{2,1}$    $A_{2,2}$

# Metrics

End-tasks of the domains

Restaurant    ACL    AI

Domains that have post-trained

Restaurant     $A_{1,1}$

ACL     $A_{2,1}$     $A_{2,2}$

AI     $A_{3,1}$     $A_{3,2}$     $A_{3,3}$

# Metrics

End-tasks of the domains

| | Restaurant | ACL | AI | Phone | PubMed | ...... |
|---|---|---|---|---|---|---|

Domains that have post-trained

| | Restaurant | ACL | AI | Phone | PubMed |
|---|---|---|---|---|---|
| Restaurant | $A_{1,1}$ | | | | |
| ACL | $A_{2,1}$ | $A_{2,2}$ | | | |
| AI | $A_{3,1}$ | $A_{3,2}$ | $A_{3,3}$ | | |
| Phone | $A_{4,1}$ | $A_{4,2}$ | $A_{4,3}$ | $A_{4,4}$ | |
| PubMed | $A_{5,1}$ | $A_{5,2}$ | $A_{5,3}$ | $A_{5,4}$ | $A_{5,5}$ |
| ⋮ | | | | | |
| | $A_{t,1}$ | ...... | | | $A_{t,t}$ |

# Final Performance

End-tasks of the domains

Restaurant    ACL    AI    Phone    PubMed    ……

Domains that have post-trained

Restaurant  $A_{1,1}$

ACL  $A_{2,1}$    $A_{2,2}$

AI  $A_{3,1}$    $A_{3,2}$    $A_{3,3}$

Phone  $A_{4,1}$    $A_{4,2}$    $A_{4,3}$    $A_{4,4}$

PubMed  $A_{5,1}$    $A_{5,2}$    $A_{5,3}$    $A_{5,4}$    $A_{5,5}$

$A_{t,1}$    ……    $A_{t,t}$

**Final Performance:** $\frac{1}{T}\sum_{i=1}^{T} R_{t,i}$

**The higher, the better.** The most popular metric

# Forgetting Rate

End-tasks of the domains

Restaurant    ACL    AI    Phone    PubMed    ......

goes down=forgetting

goes up = positive backward transfer

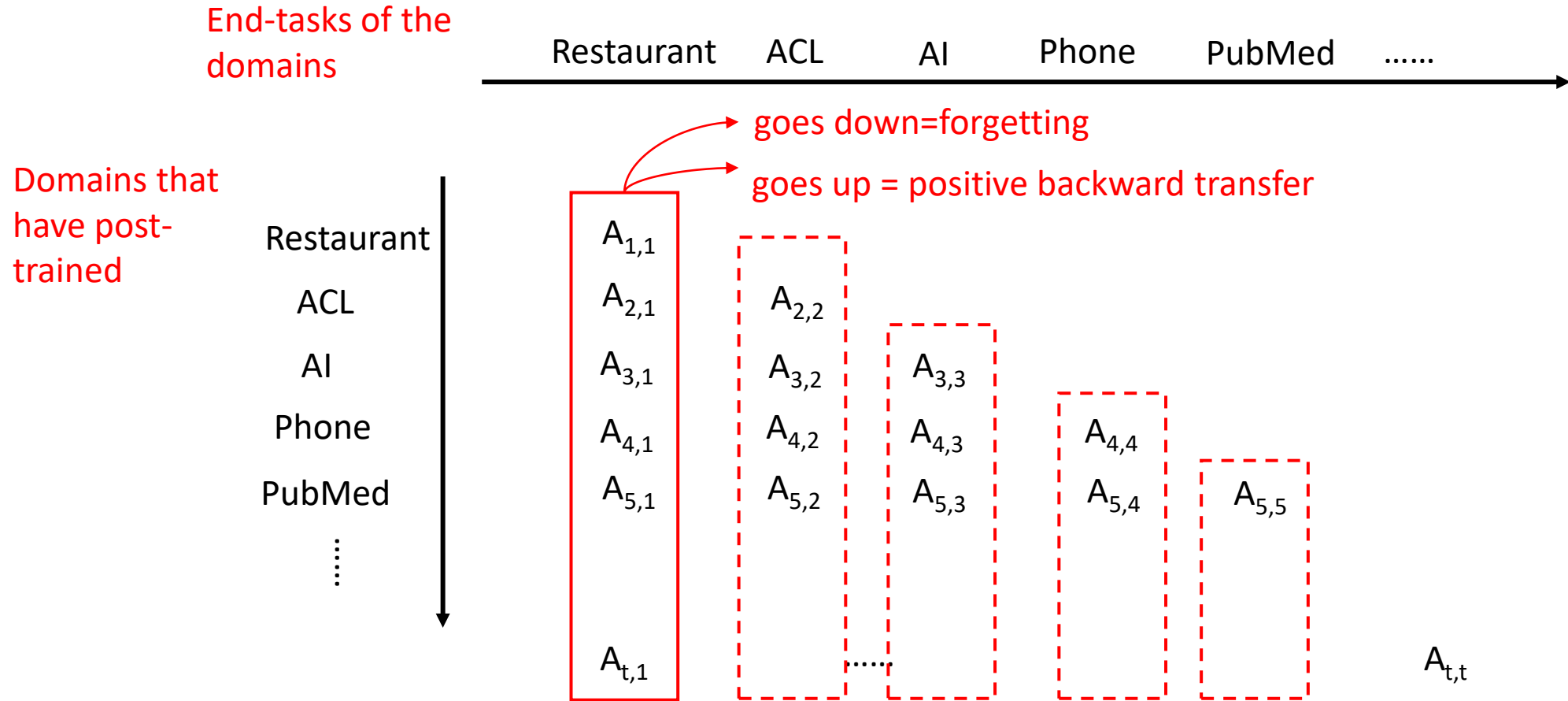Domains that have post-trained

| | Restaurant | ACL | AI | Phone | PubMed |
|---|---|---|---|---|---|
| Restaurant | $A_{1,1}$ | | | | |
| ACL | $A_{2,1}$ | $A_{2,2}$ | | | |
| AI | $A_{3,1}$ | $A_{3,2}$ | $A_{3,3}$ | | |
| Phone | $A_{4,1}$ | $A_{4,2}$ | $A_{4,3}$ | $A_{4,4}$ | |
| PubMed | $A_{5,1}$ | $A_{5,2}$ | $A_{5,3}$ | $A_{5,4}$ | $A_{5,5}$ |
| ⋮ | $A_{t,1}$ | ..... | | | $A_{t,t}$ |

**Forgetting Rate:** $\frac{1}{T-1}\sum_{k=1}^{t-1} A_{k,k} - A_{t,k}$

The difference between a task **first learned** performance and its **final** performance
**Positive**=forgetting; **Negative**=positive backward transfer

# Overall Performance

Non-Continual Learning

Without post-train (directly fine-tune the LM)

Individual post-training

| Restaurant | ACL | AI | Phone | PubMed | Camera | Average |
|---|---|---|---|---|---|---|
| 79.81 | 66.11 | 60.98 | **83.75** | 72.38 | 78.82 | 73.64 |
| **80.84** | **68.75** | **68.97** | 82.59 | **72.84** | **84.39** | **76.4** |

w/o Pre-trained < Individual Post-trained

✓ This is not surprising, as post-training has been demonstrated to improve performance in the literature.

# Overall Performance

Now we can look at continual learning

Without post-train (directly fine-tune the LM)

Individual post-training

Our continual post-training method (**CPS**)

| Restaurant | ACL | AI | Phone | PubMed | Camera | Average |
|---|---|---|---|---|---|---|
| 79.81 | 66.11 | 60.98 | 83.75 | 72.38 | 78.82 | 73.64 |
| **80.84** | 68.75 | 68.97 | 82.59 | **72.84** | 84.39 | 76.4 |
| 80.34 | **69.36** | **70.93** | **85.99** | 72.8 | **88.16** | **77.93** |

w/o Pre-trained < Individual Post-trained < CPS

✓ CPS is better than individual post-training
CPS can not only mitigate forgetting but also encourage knowledge transfer

# Overall Performance

Continual Learning **v.s. CPS**

Forgetting Rate: $\frac{1}{T-1}\sum_{k=1}^{t-1} A_{k,k} - A_{t,k}$

| | Restaurant | ACL | AI | Phone | PubMed | Camera | Average | Forgetting Rate |
|---|---|---|---|---|---|---|---|---|
| No post-train | 79.81 | 66.11 | 60.98 | 83.75 | 72.38 | 78.82 | 73.64 | --- |
| Individual post-train | 80.84 | 68.75 | 68.97 | 82.59 | 72.84 | 84.39 | 76.4 | --- |
| | 79.52 | 68.39 | 67.94 | 84.1 | 72.49 | 85.71 | 76.36 | 1.14 |
| | 80.34 | **69.36** | **70.93** | **85.99** | **72.8** | **88.16** | **77.93** | **-1.09** |

Non-Continual-learning

Naïve continual learning (**NCL**): continual learning without any specific technique

**CPS**

❌ + forgetting rate in NCL, indicates it does suffer from forgetting

✅ - forgetting rate in CPS, indicating it has positive transfer

|  | Restaurant | ACL | AI | Phone | PubMed | Camera | Average | Forgetting Rate |
|---|---|---|---|---|---|---|---|---|
| No post-train | 79.81 | 66.11 | 60.98 | 83.75 | 72.38 | 78.82 | 73.64 | --- |
| Individual post-train | 80.84 | 68.75 | 68.97 | 82.59 | 72.84 | 84.39 | 76.4 | --- |
| Naïve continual post-training | 79.52 | 68.39 | 67.94 | 84.1 | 72.49 | 85.71 | 76.36 | 1.14 |
| EWC | 80.98 | 65.94 | 65.04 | 82.32 | 71.43 | 83.35 | 74.84 | 0.02 |
| DER++ | 79 | 67.2 | 63.96 | 83.22 | 72.58 | 87.1 | 75.51 | 2.36 |
| HAT | 79.29 | 68.25 | 64.84 | 81.44 | 71.61 | 82.37 | 74.63 | -0.23 |
| BCL | 78.97 | 70.71 | 66.26 | 81.7 | 71.99 | 85.06 | 75.78 | -0.06 |
| **CPS** | 80.34 | **69.36** | **70.93** | **85.99** | **72.8** | **88.16** | **77.93** | **-1.09** |

Non-Continual-learning → No post-train, Individual post-train

Naïve continual post-training

SoTA continual learning baselines → EWC, DER++, HAT, BCL

✓ CPS outperforms SoTA

❌ Most of the SoTA only focus on mitigating forgetting, which is not enough

❌ Even replay-based method (DER++) is not good as post-training need much more replay data
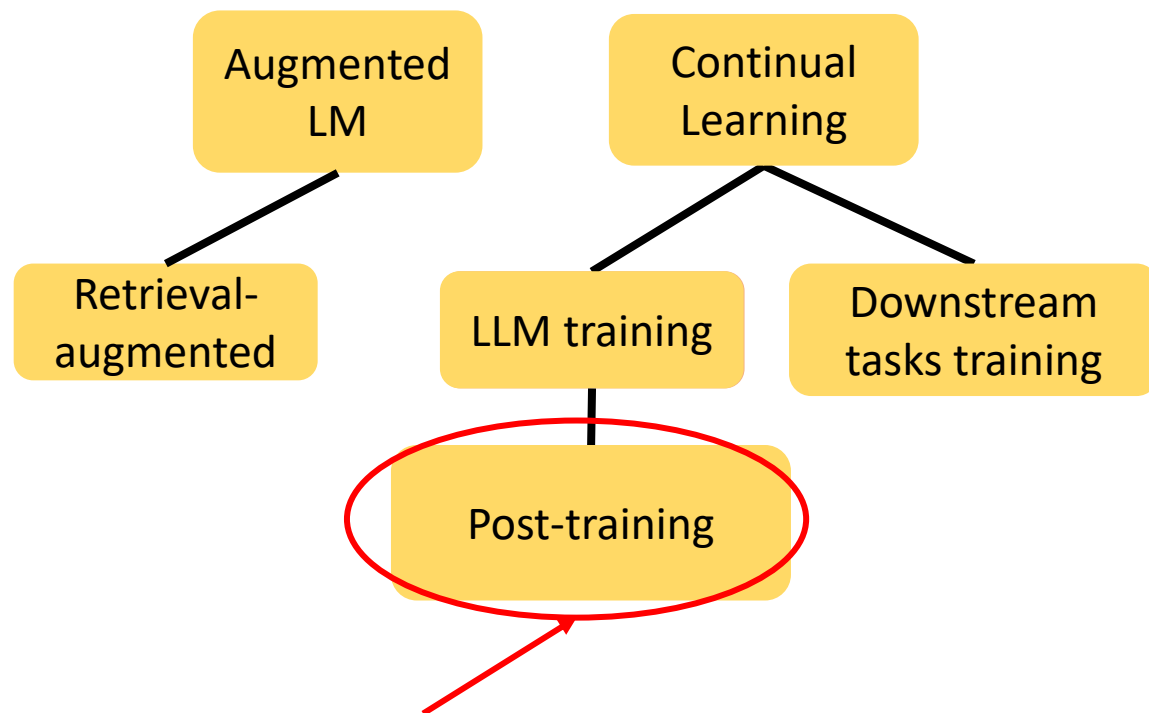
# Continual Post-training of Language Model

- Computing **importance** of units for general and domain knowledge, with **different** $L_{\mathrm{impt}}$

- **Soft-masking** the backward propagation based on importance (which help CF and KT)

# Enhancing LLM for A Dynamic World

How to make knowledge in LLM more **reusable** and **updatable**?

Augmented LM

Continual Learning

Retrieval-augmented

LLM training

Downstream tasks training

Post-training

Continual Post-training of Language Models, Ke et al., ICLR 2023

# Enhancing LLM for A Dynamic World

## Why it could be increasingly important?

# Enhancing LLM for A Dynamic World

## Why it could be increasingly important?

- The fixed world assumption is way too limited!

Over just a few months, ChatGPT went from correctly answering a simple math problem 98% of the time to just 2%, study finds

BY **PAOLO CONFINO**
July 19, 2023 at 6:29 PM CDT

How Is ChatGPT's Behavior Changing over Time?

Lingjiao Chen[†], Matei Zaharia[‡], James Zou[†]

[†]Stanford University  [‡]UC Berkeley

# Enhancing LLM for A Dynamic World

## Why it could be increasingly important?

- The fixed world assumption is way too limited!

- LLMs are increasingly replacing/eliminating building blocks and memorizing more and more knowledge, **yet** these still depends on human efforts. A more ambitious goal is to make this **fully autonomous**, which require LLMs to **self-initiate** and **adapt to new circumstances**.
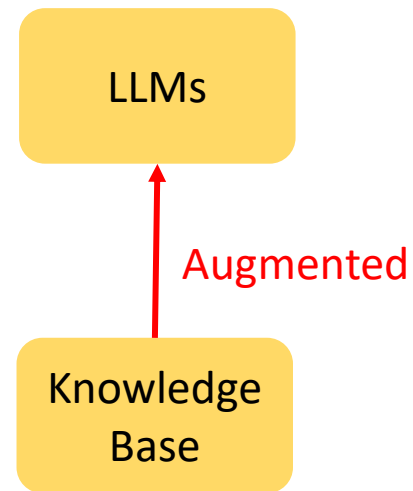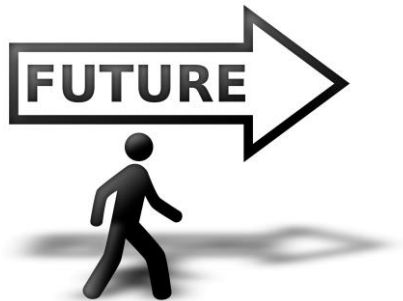
# Enhancing LLM for A Dynamic World

Why it could be increasingly important?

- The fixed world assumption is way too limited!

- LLMs are increasingly replacing/eliminating building blocks and memorizing more and more knowledge, these still depends on human efforts. A more ambitious vision is to make this **fully autonomous**, which require LLMs to **self-initiate** and **adapt to new circumstances**.

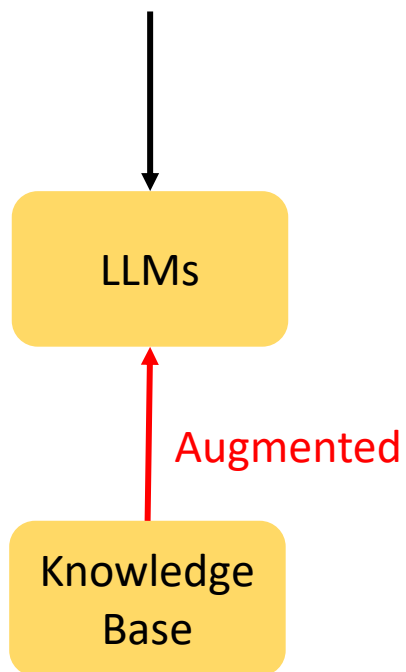- It is still cutting-edge, and still, plenty of room to improve (see next!)

**FUTURE**

What research questions can lead us toward a more autonomous LLMs?

LLMs

Augmented

Knowledge Base

**FUTURE**

What research questions can lead us toward a more autonomous LLMs?
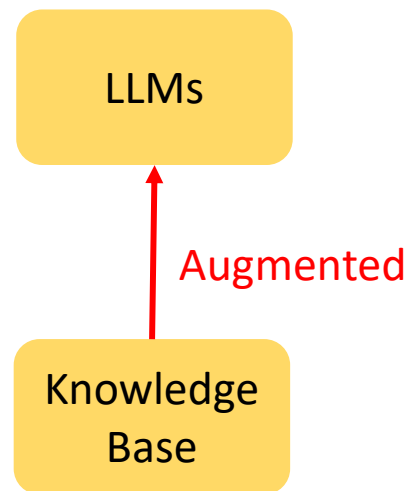
*Complete the sentence in Trump's tone: "Between a wall and an egg that breaks it, I will always stand on the side of "*

LLMs

Augmented

Knowledge Base

Retrieve Trump's speeches from the KB and augment the learner's working memory (context).

What research questions can lead us toward a more autonomous LLMs?

LLMs

Augmented

Knowledge Base

**Research questions:**

- What to retrieve (rerank/selection…)?
- How to better combine the retriever and LLM?
- When to use retrieval and when to update/use LLMs' parameters?

FUTURE

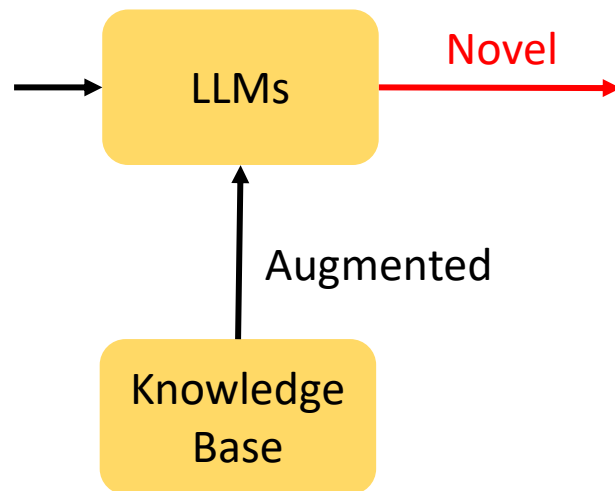What research questions can lead us toward a more autonomous LLMs?

Complete the sentence in Trump's tone: "Between a wall and an egg that breaks it, I will always stand on the side of "
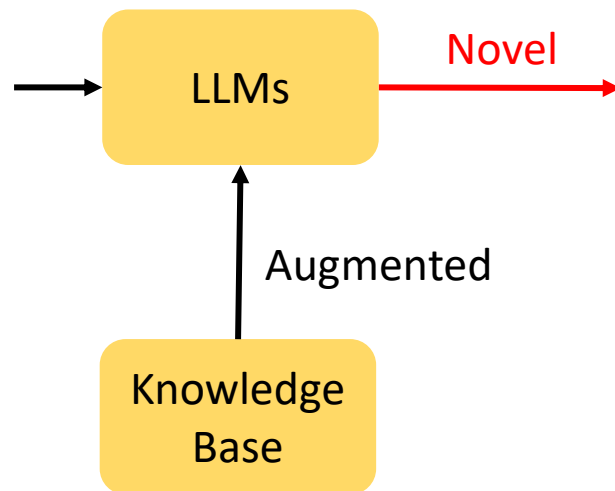
LLMs

Augmented

Knowledge Base

Application

# What research questions can lead us toward a more autonomous LLMs?

*Complete the sentence in **Zixuan's** tone: "Between a wall and an egg that breaks it, I will always stand on the side of "*

```
            ┌─────────┐   Novel
       ───► │  LLMs   │ ──────────►
            └─────────┘
                 ▲
              Augmented
                 │
            ┌─────────┐
            │Knowledge│
            │  Base   │
            └─────────┘
```

In some instances, there could be "novelty"

**Novelty/Unknown/Unexpected/Unclear:** anything that the LLM does not fully understand in order to accomplish the task

# What research questions can lead us toward a more autonomous LLMs?

**_Complete the sentence in Trump's tone_**:
"Between a wall and an egg that breaks it, I will always stand on the side of "

LLMs → **Novel**

↑ Augmented

Knowledge Base

It could also be unclear to the LLM what '_in one's tone_' means or what aspects should be focused on.

**Novelty/Unknown/Unexpected/Unclear:** anything that the LLM does not fully understand in order to accomplish the task

# What research questions can lead us toward a more autonomous LLMs?

*Complete the sentence in **Zixuan's** tone: "Between a wall and an egg that breaks it, I will always stand on the side of "*

LLMs → **Novel**

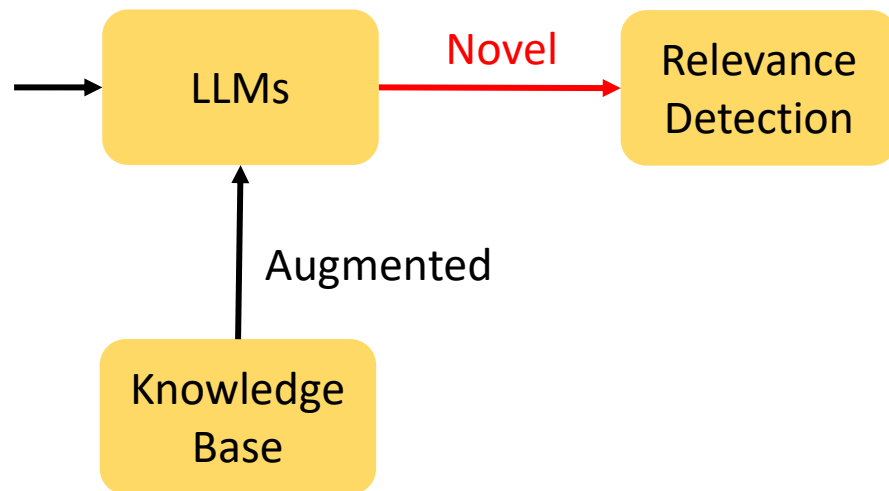Augmented

Knowledge Base

**Research question:**

How to detect novelty (knowledge that LLM does not already know)?

Application

**What research questions can lead us toward a more autonomous LLMs?**

Hello, my name is **_Vincent Bing._** _Please complete the sentence in_ **_Zixuan's_** _tone: "Between a wall and an egg that breaks it, I will always stand on the side of "_

LLMs

Novel

Relevance Detection

Augmented

Knowledge Base

It is possible that the novelty occurs but is not related to the application.

Application

# FUTURE
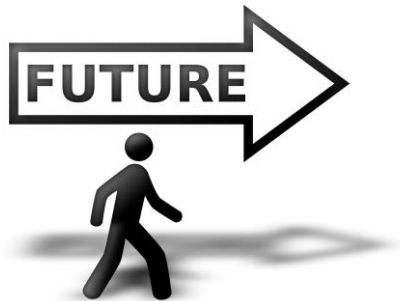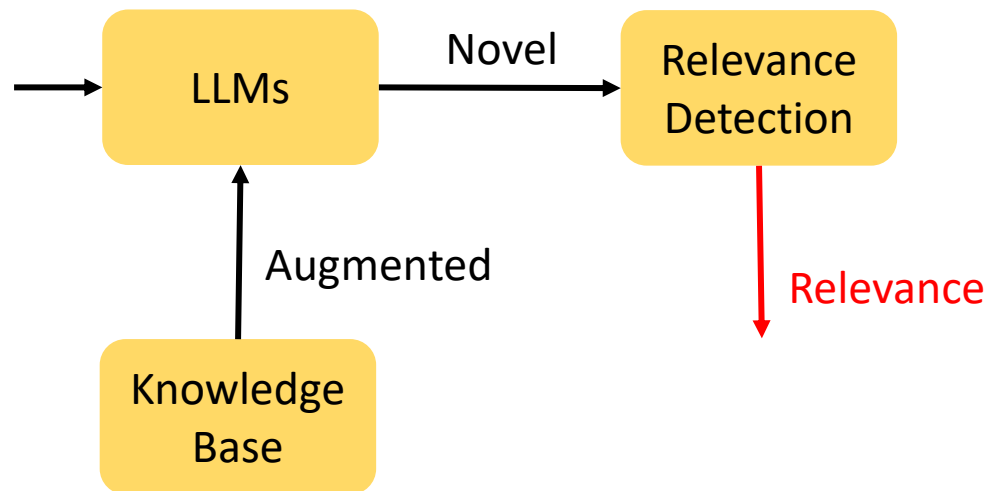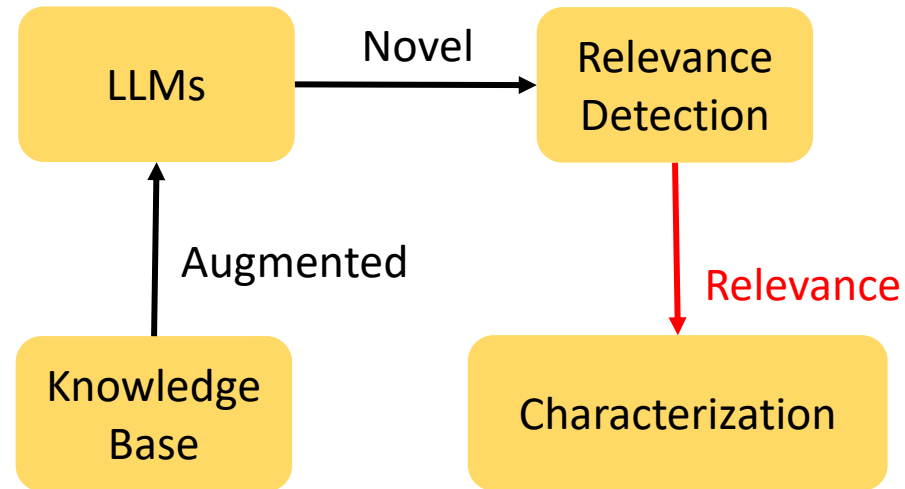
## What research questions can lead us toward a more autonomous LLMs?

**Hello, my name is Bing.**
*Please complete the sentence in **Zixuan's** tone: "Between a wall and an egg that breaks it, I will always stand on the side of "*



**Research question:**

How can we determine if the novelty is relevant to the final application? (there could be some noise!)

# What research questions can lead us toward a more autonomous LLMs?

*Complete the sentence in **Zixuan's** tone: "Between a wall and an egg that breaks it, I will always stand on the side of "*

LLMs

Novel →

Relevance Detection

Augmented ↑

Knowledge Base
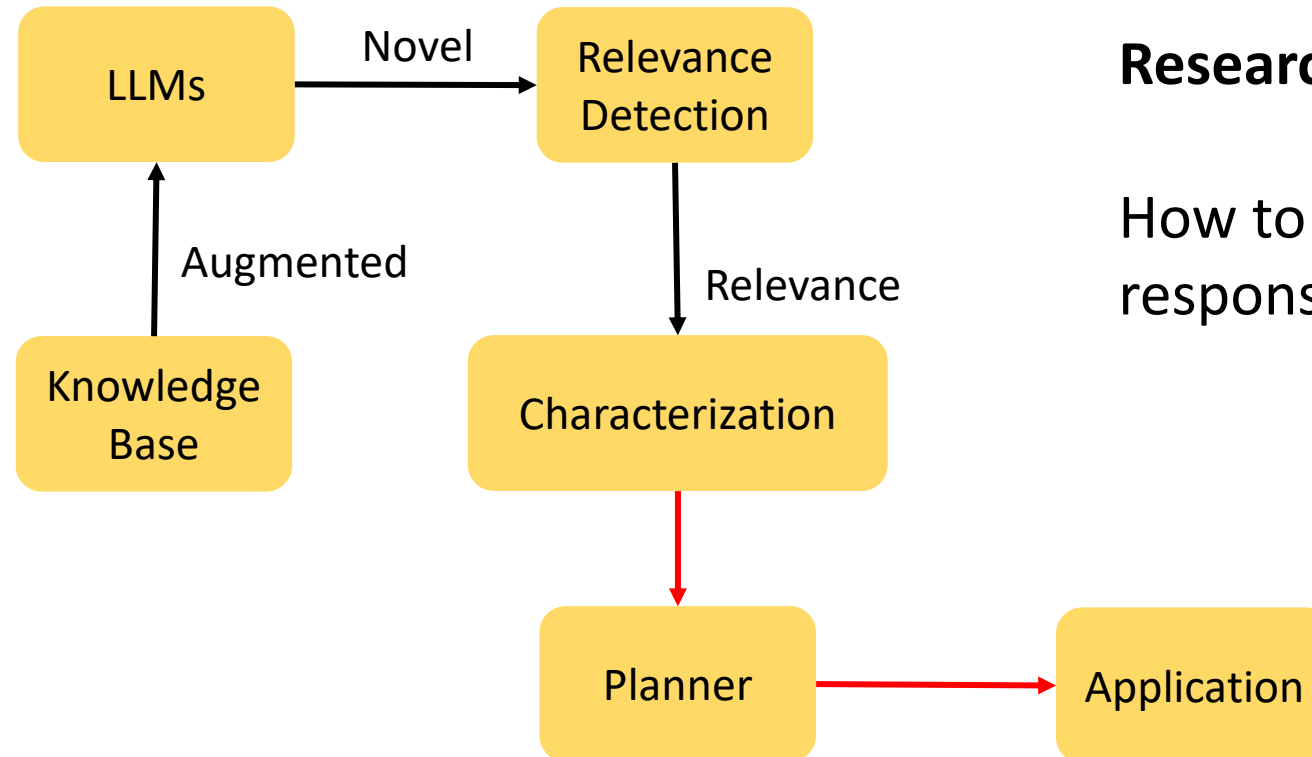
Relevance ↓

Characterization

Planner → Application

May want to response with:

*"I'm not familiar with "Zixuan". If you could provide more context or details about who Zixuan is or the style you have in mind, I will be happy to help*

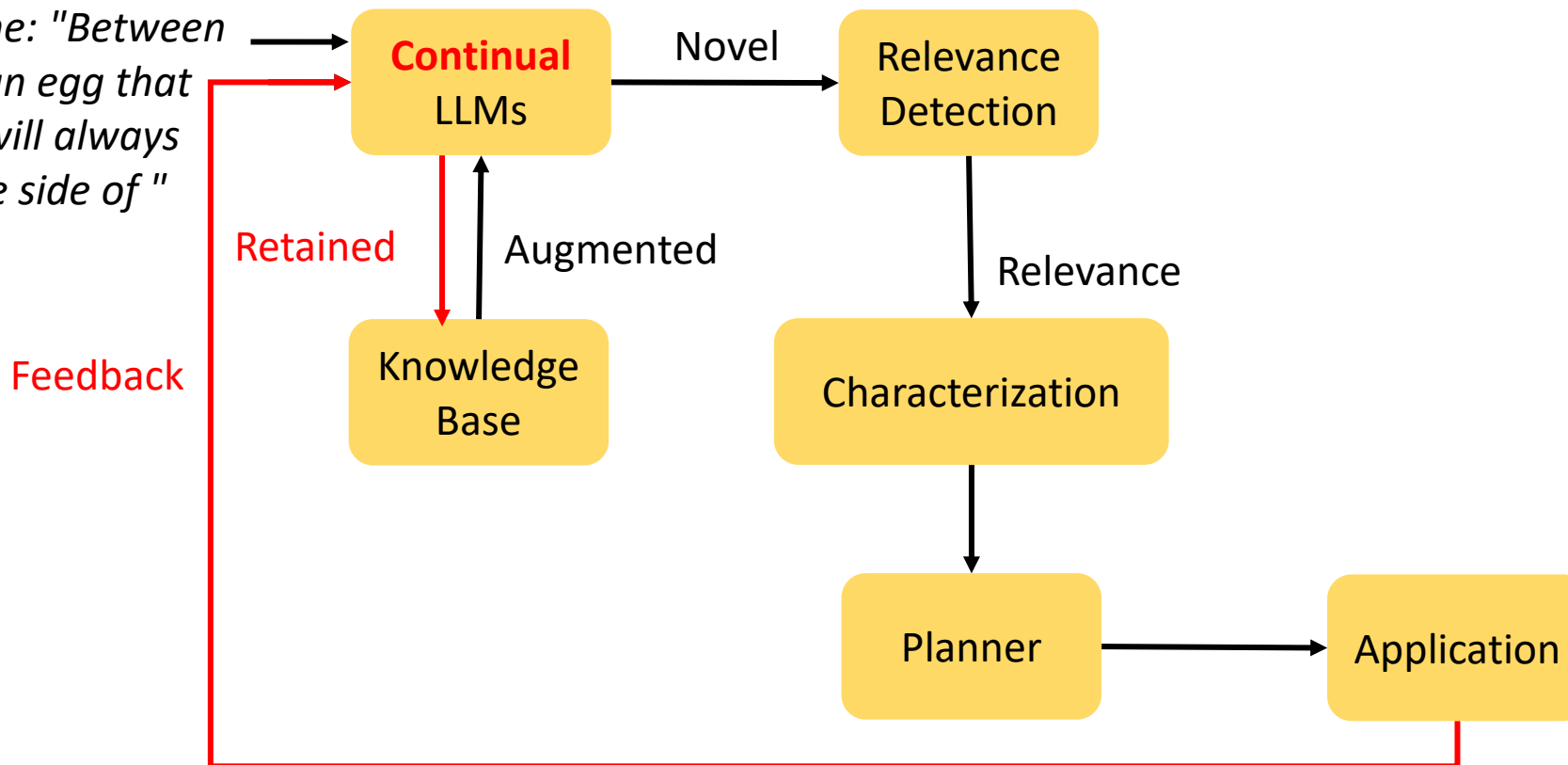What research questions can lead us toward a more autonomous LLMs?

Complete the sentence in *Zixuan's* tone: "Between a wall and an egg that breaks it, I will always stand on the side of "

Retained

Feedback

Continual LLMs

Novel

Relevance Detection

Augmented

Knowledge Base
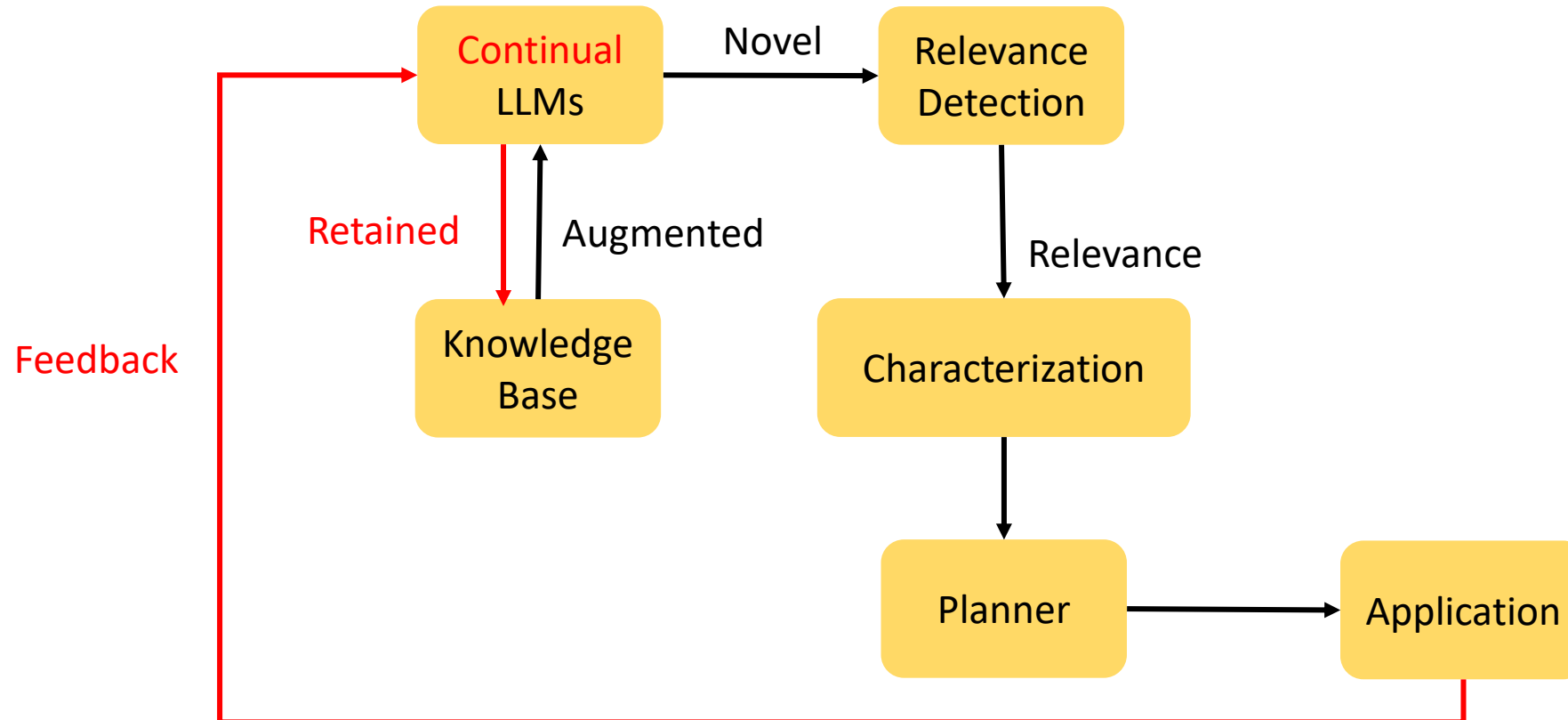
Relevance

Characterization

Planner

Application

User/another agent may give feedback:

"Zixuan is a final-year PhD student working on LLM (pretrain, post-train, frontiers like retrieval-augmented LLM) and continual learning…."

What research questions can lead us toward a more autonomous LLMs?
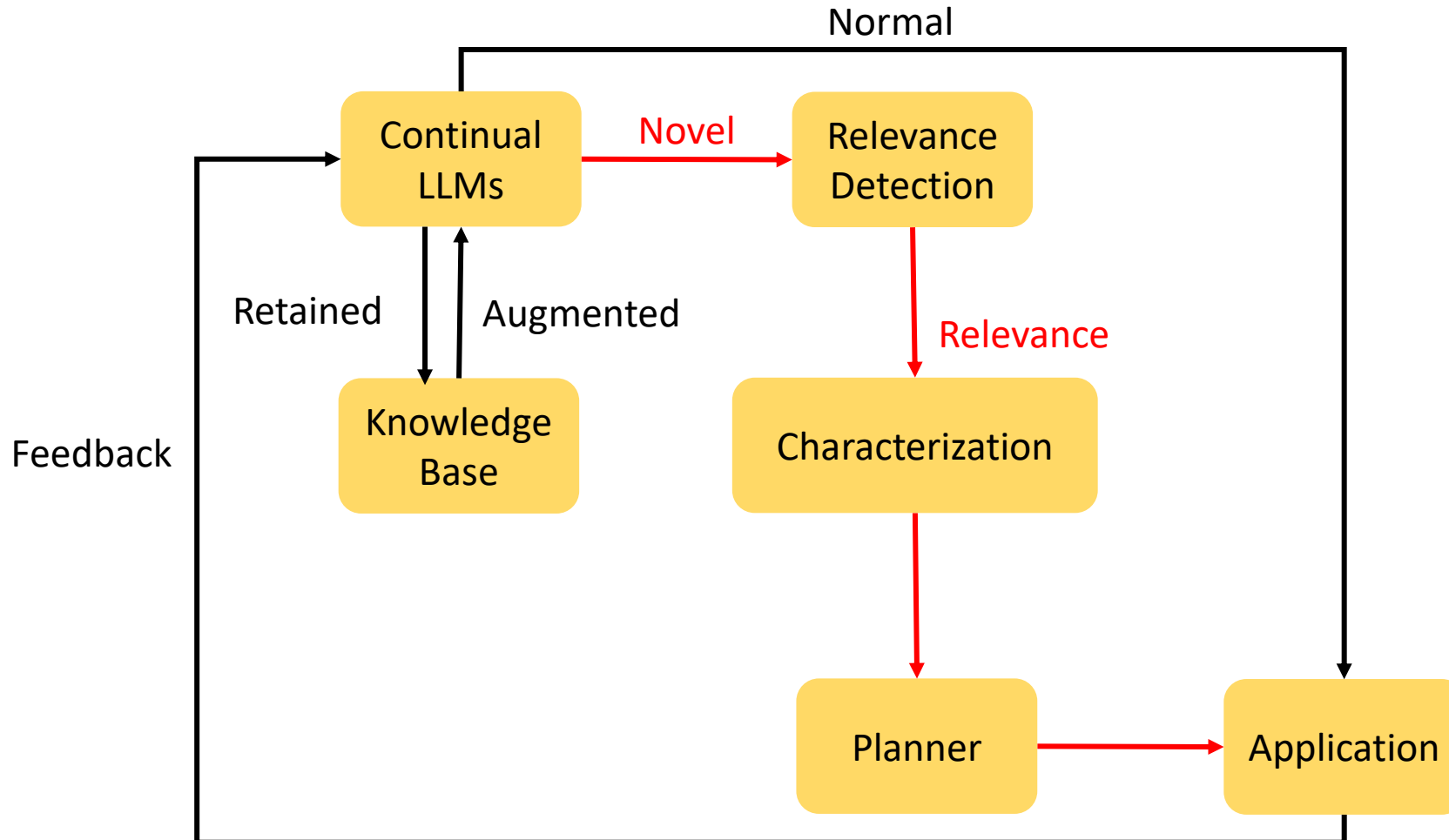
**Research question:**

How to use the feedback to continually update the LLM and retain the useful knowledge?

Using external memory, working memory (context), or updating the LLM?

# FUTURE

## What research questions can lead us toward a more autonomous LLMs?
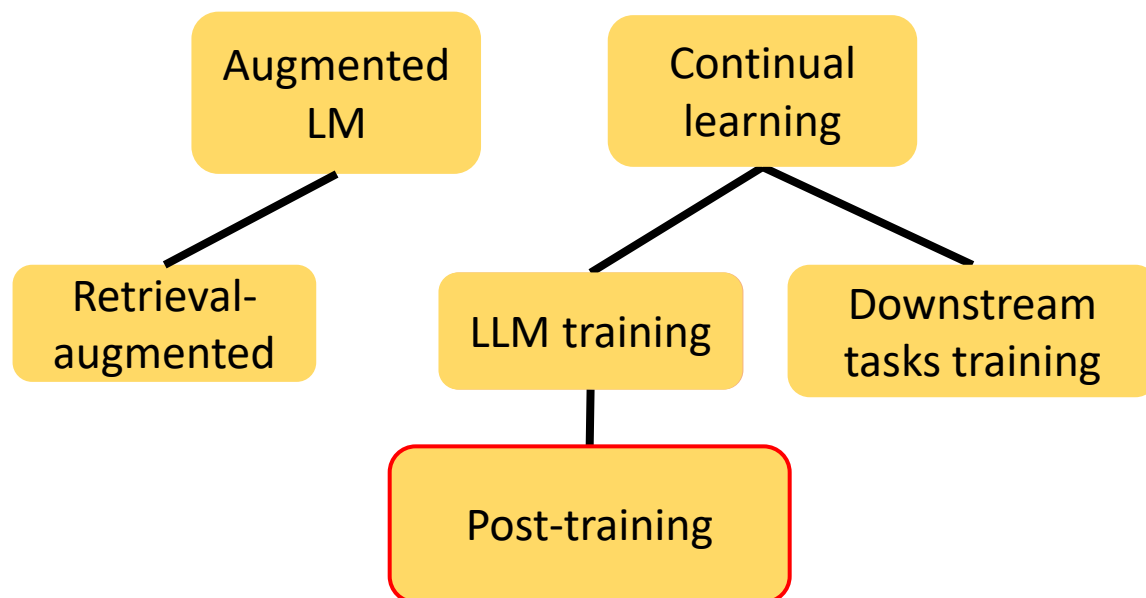
Most existing works are dedicated to the black part, which includes active research areas like retrieval-augmented LLM and continual learning.

The other component remains largely unexplored!

# Enhancing LLM for A Dynamic World

How to make knowledge in LLM more **reusable** and **updatable**?

Augmented LM

Continual learning

Retrieval-augmented

LLM training

Downstream tasks training

Post-training

**Ambitious goal:** Fully autonomous LLM

Thank you