



**ICML**  
International Conference  
On Machine Learning



# Continual Model Routing in Evolving Model Hubs

---

Jack Bell<sup>1</sup> Giacomo Carfi<sup>1</sup> Gerlando Gramaglia<sup>1</sup> Vincenzo Lomonaco<sup>2</sup>

<sup>1</sup>University of Pisa <sup>2</sup>LUISS University, Rome

jack.bell@phd.unipi.it

# Motivation

How should we route a user request to the right model when model hubs contain thousands of rapidly evolving experts?

- Existing routers are inherently static
- Formalise routing as a continual learning problem
- CARvE learns a contrastive embedding that routes queries to models while anchoring past knowledge via checkpoint replay.

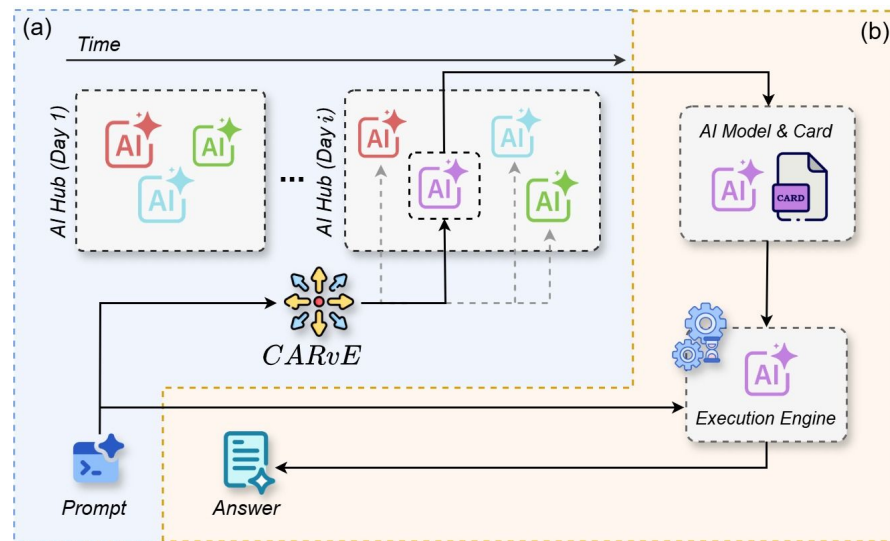


Figure 1: (a) Our adaptive router CARvE learns continually from selection samples to dynamically route a prompt to the most appropriate model. (b) Actual execution of the model (not the focus of this paper)

# CMR Formulation

## Formulation

Given a growing hub  $\mathcal{H}$  and a stream of experiences  $\mathcal{E}_1, \dots, \mathcal{E}_n$ , each providing query-model pairs  $D_i = \{(q, m^*)\}$ , learn a router  $f_\theta : \mathcal{Q} \rightarrow \mathcal{H}$  that maximises routing accuracy across all experiences, with access only to  $D_i$  at step  $i$  and without executing any  $m \in \mathcal{H}$ .



### Scale

Hubs contain millions of models - routing must remain efficient as the hub grows.



### Forgetting

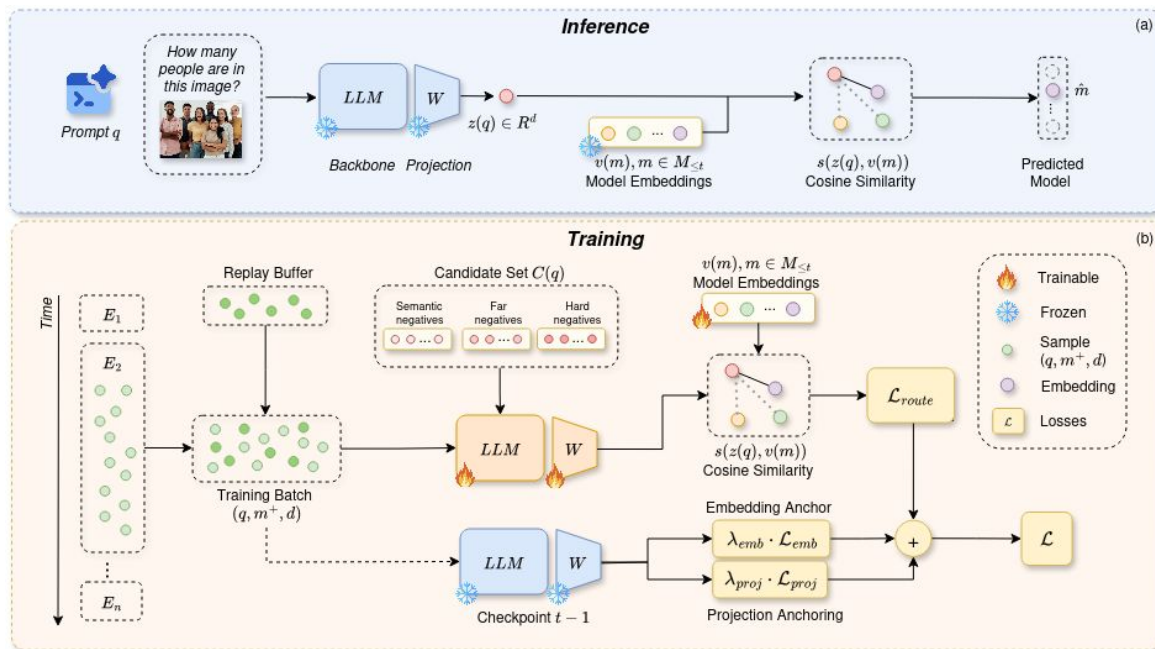
Sequential training on new domains causes catastrophic forgetting of prior routing decisions.



### Label Shift

Each new experience introduces new model families and query distributions, shifting the label space.

# CARvE: Method Overview



① Query Prompt

② CARvE Router

③ Selected Model

# CMRBench: Experimental Setup

Datasets	Baselines	Metrics
<ul style="list-style-type: none"><li>• APIBench</li></ul>	<ul style="list-style-type: none"><li>• Sequential Fine-tuning</li></ul>	<ul style="list-style-type: none"><li>• Model-ID Acc (M-Acc)</li></ul>
<ul style="list-style-type: none"><li>• ToolMMBench</li></ul>	<ul style="list-style-type: none"><li>• SFT with Random Replay (10%)</li></ul>	<ul style="list-style-type: none"><li>• Family Acc (F-Acc)</li></ul>
<ul style="list-style-type: none"><li>• HuggingBench</li></ul>	<ul style="list-style-type: none"><li>• EWC (regularisation)</li></ul>	<ul style="list-style-type: none"><li>• Domain Acc (D-Acc)</li></ul>
<ul style="list-style-type: none"><li>• 4 experiences · 2,000+ models</li></ul>	<ul style="list-style-type: none"><li>• Cumulative (upper bound)</li></ul>	<ul style="list-style-type: none"><li>• Forgetting</li></ul>



Implementation: PyTorch + TRL framework · GPU: H100 (×1) · Code available at: <https://github.com/collagelab/CMR>

# Main Results

**80.7%**

Domain Accuracy  
(CARvE, 10% replay)

**5.9%**

Domain Forgetting

**+4.8pp**

over Random Replay D-Acc

Setting	Accuracy (%) ( $\uparrow$ )			Forgetting (%) ( $\downarrow$ )		
	M-Acc	F-Acc	D-Acc	M-Fgt	F-Fgt	D-Fgt
<b>Retrieval / LLM-controller baselines</b>						
BGE-M3 (Chen et al., 2024a)	<b>13.6 <math>\pm</math> 0.0</b>	<b>16.2 <math>\pm</math> 0.0</b>	44.0 $\pm$ 0.0	2.5 $\pm$ 0.0	3.9 $\pm$ 0.0	3.3 $\pm$ 0.0
Gorilla RAG (Patil et al., 2024)	<u>6.7 <math>\pm</math> 0.2</u>	<u>10.4 <math>\pm</math> 0.2</u>	43.0 $\pm$ 0.5	<b>0.0 <math>\pm</math> 0.2</b>	<b>0.8 <math>\pm</math> 0.2</b>	<b>0.1 <math>\pm</math> 0.2</b>
HuggingGPT Qwen3-32B (Shen et al., 2023)	–	–	<b>51.7 <math>\pm</math> 0.2</b>	–	–	–
<b>Continual learning experiments</b>						
Random Replay (5%)	36.5 $\pm$ 0.4	44.2 $\pm$ 0.4	70.2 $\pm$ 4.2	19.5 $\pm$ 0.5	23.2 $\pm$ 0.5	24.0 $\pm$ 4.4
Random Replay (10%)	39.1 $\pm$ 0.5	47.3 $\pm$ 0.7	75.9 $\pm$ 0.2	14.0 $\pm$ 0.1	16.8 $\pm$ 0.3	13.1 $\pm$ 0.2
Random Replay (20%)	41.3 $\pm$ 0.2	<u>49.8 <math>\pm</math> 0.2</u>	78.1 $\pm$ 0.1	<u>8.6 <math>\pm</math> 0.4</u>	<u>11.1 <math>\pm</math> 0.3</u>	7.8 $\pm$ 0.1
Random Replay Qwen2.5-7B (10%)	40.7 $\pm$ 0.5	49.2 $\pm$ 0.4	78.2 $\pm$ 0.2	14.2 $\pm$ 0.6	16.5 $\pm$ 0.6	10.3 $\pm$ 0.5
Random Replay Qwen3-4B (10%)	39.6 $\pm$ 0.2	47.9 $\pm$ 0.1	77.2 $\pm$ 0.2	14.9 $\pm$ 0.2	18.2 $\pm$ 0.4	11.2 $\pm$ 0.4
Sequential Finetuning	28.0 $\pm$ 0.1	34.8 $\pm$ 0.1	64.3 $\pm$ 0.2	34.6 $\pm$ 0.2	41.8 $\pm$ 0.2	37.2 $\pm$ 0.2
Model Merging (TIES) (Yadav et al., 2023)	7.6 $\pm$ 0.3	10.9 $\pm$ 0.2	28.6 $\pm$ 0.5	15.0 $\pm$ 0.2	19.7 $\pm$ 0.1	32.6 $\pm$ 0.2
LwF (Li & Hoiem, 2018)	28.8 $\pm$ 0.3	35.9 $\pm$ 0.5	56.4 $\pm$ 0.2	20.8 $\pm$ 0.3	25.7 $\pm$ 0.2	39.5 $\pm$ 0.2
EWC (Kirkpatrick et al., 2017)	31.3 $\pm$ 0.4	38.4 $\pm$ 0.4	66.2 $\pm$ 0.1	26.6 $\pm$ 0.3	32.7 $\pm$ 0.2	31.4 $\pm$ 0.5
CARvE EWC (10% replay)	42.2 $\pm$ 2.4	47.7 $\pm$ 2.3	80.5 $\pm$ 2.2	9.1 $\pm$ 1.9	9.8 $\pm$ 3.3	6.1 $\pm$ 4.8
CARvE (5% replay)	40.5 $\pm$ 0.7	45.4 $\pm$ 0.6	78.5 $\pm$ 0.4	15.9 $\pm$ 1.7	17.4 $\pm$ 1.1	10.1 $\pm$ 0.7
CARvE (10% replay)	<u>43.1 <math>\pm</math> 0.2</u>	48.5 $\pm$ 0.3	80.7 $\pm$ 0.0	11.4 $\pm$ 0.4	12.5 $\pm$ 0.3	<u>5.9 <math>\pm</math> 0.3</u>
CARvE Qwen3-4B (10% replay)	36.7 $\pm$ 3.4	42.1 $\pm$ 3.9	78.6 $\pm$ 3.3	13.5 $\pm$ 3.3	14.3 $\pm$ 4.5	9.1 $\pm$ 5.0
CARvE Qwen2.5-7B (10% replay)	42.4 $\pm$ 0.1	48.1 $\pm$ 0.3	<u>81.5 <math>\pm</math> 0.3</u>	11.2 $\pm$ 0.2	12.3 $\pm$ 0.1	6.2 $\pm$ 0.6
CARvE (20% replay)	<b>46.4 <math>\pm</math> 0.1</b>	<b>51.9 <math>\pm</math> 0.2</b>	<b>82.9 <math>\pm</math> 0.2</b>	<b>6.4 <math>\pm</math> 0.8</b>	<b>7.3 <math>\pm</math> 0.8</b>	<b>3.0 <math>\pm</math> 0.3</b>

# Ablation Study

Method	D-Acc (%)	D-Fgt (%)	$\Delta$ D-Fgt
CARvE (full)	80.7	5.9	—
No emb. anchor	78.4	7.4	+1.5
No proj. anchor ★	74.2	9.4	+3.5
Domain emb. initialisation	78.8	10.2	+4.3
Sequential FT	64.3	35.7	+29.8

► Projection anchor is the most critical component, removing it raises D-Fgt from 5.9 to 9.4 (+3.5 pp forgetting). Random embedding initialisation is also an important design choice to avoid constraining learning.

# Conclusion

---

## Key Takeaways

- We formalised Continual Model Routing (CMR) - selecting the right model from a growing hub without forgetting past domains.
- CARvE achieves 80.7% domain accuracy on CMRBench with only 5.9% forgetting
- CMRBench is released as a new evaluation suite for the growing model-routing research community.

## Future Work

Hybrid CARvE with retrieval

Post inference routing

Theoretical analysis of the contrastive anchoring objective

# Thank You!

---

 Paper: <https://arxiv.org/abs/2605.28577>

 Code: <https://github.com/collagelab/CMR>

 Email: [jack.bell@phd.unipi.it](mailto:jack.bell@phd.unipi.it)

 Lab: [vincenzolomonaco.com/collagelab](http://vincenzolomonaco.com/collagelab)

