



Scan for project page

## 1 Continual Learning

**Setting.** Large pre-trained models are fine-tuned on a sequence of new tasks via lightweight parameter updates (e.g. LoRA)

### CATASTROPHIC FORGETTING

**Pre-trained knowledge** — base model's general capabilities degrade as incrementally learning on new tasks

**Fine-tuned knowledge** — what earlier task learned, got overwritten by learning on new incoming tasks

## 2 Current Pitfalls – MoE-LoRA

The expert is an **indivisible block** — three coupled pitfalls follow:

① **Interference:** Coarse experts *bundle* many *low-specialty subspaces*. Activating one for a specific input inevitably triggers irrelevant subspaces — *interfered with knowledge of others*.

② **Redundancy:** New experts cannot *reuse* “atoms” of old coarse experts. If the *existing combination* doesn't cover the new task, the whole block must be relearned — *wasted capacity*.

③ **Routing/Retrieval collapse:** *Under-specialized* experts *confuse* the router. As experts accumulate, router struggles to index them — *routing drifts, old experts misrouted, forgetting*.

The router is where failure shows up — but the root cause is **coarse-grained, indivisible fixed-rank experts**.

## 3 Weight Matrix as a Memory

**Linear Associative Memory (LAM).** A weight matrix  $W \in \mathbb{R}^{d_{out} \times d_{in}}$  stores pattern pairs  $\{(k_i, v_i)\}_{i=1}^m$  — keys  $k_i \in \mathbb{R}^{d_{in}}$  and values  $v_i \in \mathbb{R}^{d_{out}}$  — as single matrix of outer products

$$y = Wx \approx \sum_{i=1}^m v_i (k_i^T x)$$

where the inner product  $(k_i^T x)$  computes the *relevance* (or activation strength) of the  $i$ -th memory slot to the input, which weights the retrieval of the *value* vector  $v_i$ .

Each **weight matrix** in the Pre-trained model is already a storage of *implicit (key, value) memory atoms*.

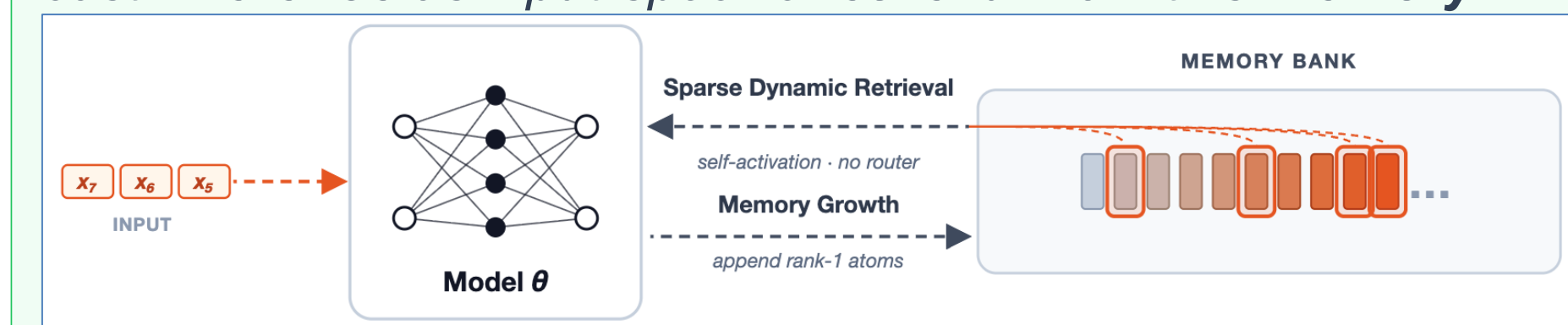
## 4 Fine-tuning as Memory Augmentation

**Low-Rank Updates (LoRA) as Rank-1 Memory Augmentation.** LoRA's update  $\Delta W = BA$  is itself a key-value memory: column vectors  $B_{:,i}$  are values, row vectors  $A_{i,:}$  are keys.

$$\Delta Wx = \sum_{i=1}^r \underbrace{B_{:,i}}_{\text{Value } v_i} (\underbrace{A_{i,:}}_{\text{Key } k_i^T} x)$$

A rank- $r$  update inserts  $r$  new key-value memory atoms into the pre-trained model's intrinsic linear associative memory.

*Continue Learning as maintaining and expanding a parametric memory — a growing collection of fine-grained, atomic units stored outside the base model's weights — and cast inference as input-specific retrieval from this memory.*



Each **memory unit** encodes an incremental **learning snapshot**, and relevant units are retrieved at test time to **specialize** the model for the input — **recall** and **reuse** of memorized snapshots.

## 5 Mixture of Rank-1 Memory (MoRAM)

**MoRAM Formulation.** Adaptation parameters become a dynamic memory bank of rank-1 atoms:

$$\mathcal{M}_t = \{(B_{:,i}, A_{i,:})\}_{i=1}^{r_t}$$

where  $r_t$  is the number of atoms accumulated by task  $t$ . For input  $x$ , the update is a sparse, input-dependent mixture over  $\mathcal{M}_t$ .

$$\Delta W^t = \sum_{i=1}^{r_t} w_i B_{:,i} A_{i,:}$$

### Self-Activation of Each Atom via Memory Keys

Each atom's own key  $A_{i,:}$  scores its relevance to the input:

$$s_i = \frac{A_{i,:} \cdot x}{\sum_{j=1}^{r_t} (A_{j,:} \cdot x)^2}$$

- ❌ Auxiliary learned router module → Router induced forgetting
- ✅ Content-addressable retrieval → Atom knows its relevance

### Sparse Memory Retrieval Drives Specialization

1. Top- $k$  masking — keep only the  $k$  most relevant atoms; mask the rest:

$$[\text{TopK}(s, k)]_i = \begin{cases} s_i, & \text{if } s_i \in \text{top-}k(s) \\ -\infty, & \text{otherwise.} \end{cases}$$

2. Sharpening — temperature softmax concentrates mass on specialists:

$$w_i = \text{softmax}\left(\frac{\text{TopK}(s, k)}{\tau_{\text{MoRAM}}}\right)_i$$

3. Test-time thresholding — prune weak signals at inference:

$$w_i := \mathbb{1}\{s_i \geq \delta\} \odot w_i$$

Forward: focused activation. Backward: gradients flow only to winning atoms — faster specialization, frozen atoms stay frozen.

### Learning “Little by Little”

For each new task: add  $r$  new atoms (Keys  $A$ , Values  $B$ ), freeze all prior atoms, let self-activation route over the union.

Each task becomes an atomic learning snapshot in memory — recalled and reused at test time, never overwritten.

Training: standard loss only. No regularization. No load-balancing. Fine-grained atoms specialize naturally; the self-activated memory retrieval stays stable as memory grows.

## 5 Experiments

### X-TAIL — CLIP (Vision-Language)

**Backbone:** CLIP ViT-B/16  
**Tasks:** Aircraft, Caltech101, DTD, EuroSAT, Flowers, Food, MNIST, OxfordPet, StanfordCars, SUN397.

**Metrics:**  
**Transfer** ↑ — zero-shot accuracy on unseen domains  
**Average** ↑ — running mean across all stages  
**Last** ↑ — accuracy after the full stream

### TRACE — LLMs (language)

**Backbones:** LLaMA-2-7B, Gemma-2B, LLaMA-3.2-1B  
**Tasks:** stance detection, finance QA, scientific summarisation, code generation, science QA, math word problems, German news, dialogue.  
**Metrics:**  
**OP** ↑ — overall performance (mean accuracy after the stream)  
**BWT** ↓ — backward transfer (accuracy drop on prior tasks; 0 = no forgetting)

### Results — X-TAIL (CLIP)

Method	Transfer ↑	Average ↑	Last ↑
CLIP zero-shot	62.4	—	—
LwF	47.7	53.2	64.0
WiSE-FT	52.3	54.2	58.0
iCaRL	61.7	61.0	64.0
ZSCL	59.0	60.0	63.4
MoE-Adapter	56.0	63.0	70.5
RAIL-Primal	62.4	70.7	79.1
CoDyRA	63.2	71.3	79.2
<b>MoRAM (ours)</b>	<b>63.3</b>	<b>72.7</b>	<b>80.9</b>

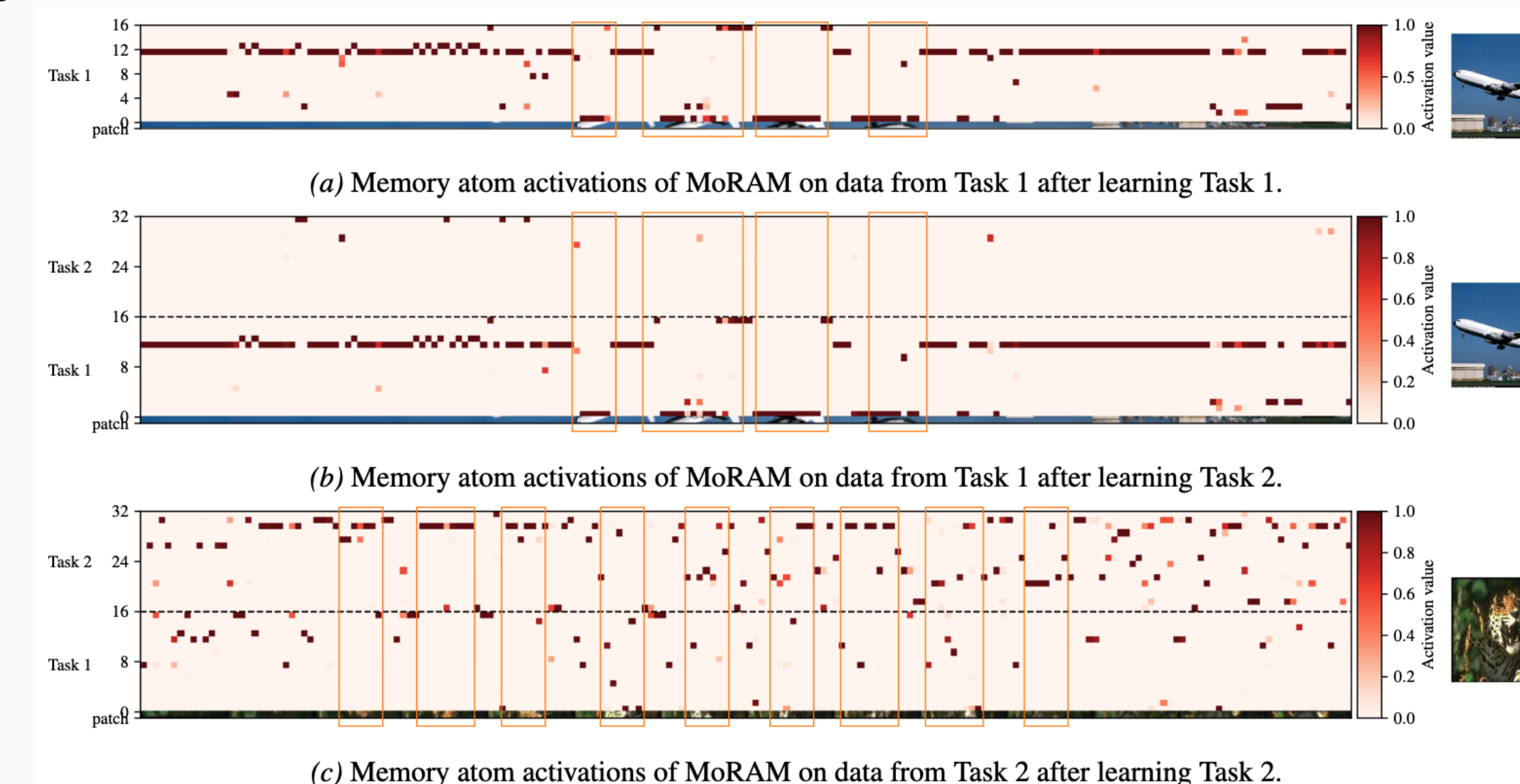
- **Transfer** — learns new tasks without losing pre-trained general capability.
- **Average & Last** — fine-grained atoms specialize without overwriting prior knowledge.

### Ablation of Memory Retrieval Strategies

Routing	Transfer	Average	Last
<i>Coarse-Grained (Rank-r Experts)</i>			
MoE-LoRA (Baseline)	62.56	69.45	74.53
w/ Temperature Scaling ( $\tau^*$ )	62.48	69.40	74.57
<i>Fine-Grained (Rank-1 Memory Experts)</i>			
External Router (Learned $W_r$ )	60.09	65.97	69.76
Self-Activated Retrieval	60.26	65.94	69.85
w/ Sparsity Constraint (Top- $k$ )	60.69	66.52	70.62
w/ Temperature Scaling ( $\tau_{\text{MoRAM}}$ )	62.07	71.15	79.62
w/ Threshold-based Selection ( $\delta$ )	60.78	66.83	71.08
<b>MoRAM (Full)</b>	<b>63.30</b>	<b>72.70</b>	<b>80.90</b>

- **Sparsity retrieval drives specialization on fine-grained atoms** — top- $k$ , temperature, and threshold each reduce interference and redundancy.
- **No gain on coarse MoE-LoRA** — the structural win is fine-grained memory, not the sparsity machinery on top.

### Analysis — What Do Atoms Learn?



**Visualization.** Per-atom mixture weights in an activation map: rows = atom indices (by task), columns = input patches, colour = activation value. (More detailed visualizations in the paper Fig. 6 – Forgetting mitigation and Fig. 7 – Knowledge reuse.)

**Specialization.** A *concentrated* subset of specialized atoms fire strongly on specific patches; the rest stay near zero. **Forgetting mitigation.** Activation patterns for Task-1 inputs stay *invariant* after learning new tasks. **Knowledge reuse.** Prior-task atoms *re-fire* on later-task inputs when similar semantics recurs (e.g., sky in task 1 re-fired for blue-sky task 9, see Fig. 7 in paper for more details) — no redundant expert/memory needed.

### Results — TRACE (LLMs)

	FIX(ICL)	SeqLoRA	OGD	GEM	EWC	L2P	DualPrompt	HiDeLoRA	O-LoRA	TreeLoRA	MoRAM
<i>meta-llama / LLaMA-2-7B-Chat</i>											
OP	38.94 ± 0.3	34.3 ± 1.2	42.09 ± 1.6	40.08 ± 1.6	42.36 ± 1.2	36.23 ± 0.8	37.69 ± 1.2	41.60 ± 0.8	42.78 ± 0.8	43.52 ± 1.0	<b>44.54 ± 0.9</b>
BWT	—	18.5 ± 0.8	8.06 ± 1.2	6.77 ± 1.2	5.97 ± 0.8	8.25 ± 0.8	8.03 ± 0.8	7.12 ± 0.4	7.16 ± 0.4	3.46 ± 0.4	<b>1.37 ± 0.3</b>
<i>google / Gemma-2B-it</i>											
OP	32.3 ± 0.2	31.89 ± 0.8	32.85 ± 1.4	26.48 ± 1.5	28.35 ± 1.6	31.14 ± 1.2	32.42 ± 1.0	33.25 ± 0.9	33.73 ± 0.8	33.41 ± 0.9	<b>36.27 ± 0.7</b>
BWT	—	15.28 ± 0.4	12.27 ± 0.9	18.25 ± 0.9	16.96 ± 1.2	15.77 ± 0.7	14.25 ± 0.5	13.66 ± 0.5	12.36 ± 0.4	8.50 ± 0.5	<b>2.74 ± 0.4</b>
<i>meta-llama / LLaMA-3-1B-Instruct</i>											
OP	31.16 ± 0.4	29.73 ± 1.6	30.12 ± 2.0	32.19 ± 2.0	31.96 ± 1.6	29.38 ± 1.2	30.76 ± 1.2	33.73 ± 1.2	32.94 ± 0.8	36.14 ± 0.7	<b>37.77 ± 0.8</b>
BWT	—	17.03 ± 1.2	15.2 ± 1.6	10.74 ± 1.6	11.62 ± 1.2	13.57 ± 0.8	11.34 ± 0.8	12.36 ± 0.8	12.89 ± 1.2	7.36 ± 0.8	<b>3.12 ± 0.8</b>

**Best OP and lowest BWT on every backbone.** Frozen keys + content-addressable retrieval over specialized memory atoms with the memory bank — old atoms don't drift as new tasks arrive.

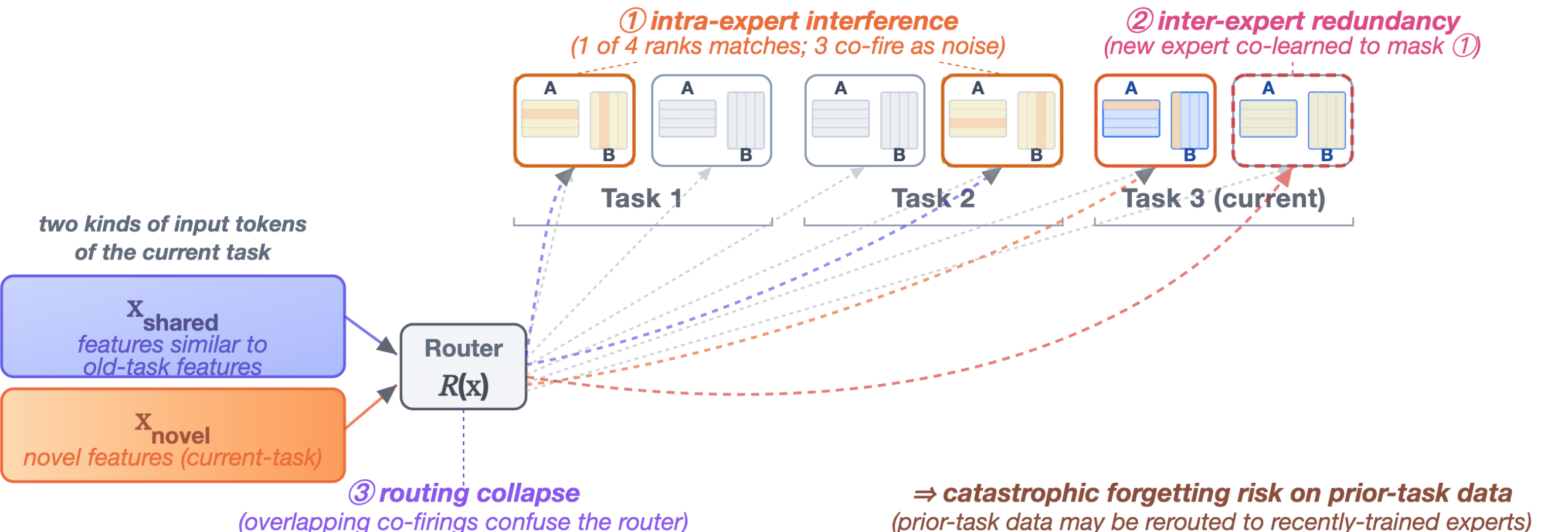
### Results — Supervised Fine-Tuning (SFT)

Method	HumanEval (Pass@1)	Out-of-Domain (Acc.)								Params (M)	%Params	
		Logic	Phil.	Reli.	Econ.	Pub. Rel.	STEM	Phys.	ML			MMLU
LLaMA-3.1-8B	38.40	42.06	<b>71.06</b>	<b>83.63</b>	70.17	<b>68.18</b>	54.84	39.22	<b>40.18</b>	63.45	—	—
LoRA ( $r=4$ )	41.46	39.68	70.09	81.87	71.43	65.45	54.77	<b>45.10</b>	40.17	63.28	10.5	0.13%
LoRA ( $r=8$ )	44.51	39.68	<b>70.74</b>	<b>81.87</b>	<b>71.85</b>	64.54	54.17	42.16	39.29	63.03	21.0	0.26%
LoRA ( $r=16$ )	<b>45.73</b>	41.27	68.49	80.70	<b>72.69</b>	<b>66.36</b>	54.96	44.11	38.39	63.35	41.9	0.52%
LoRA ( $r=32$ )	<b>47.56</b>	<b>42.85</b>	69.45	81.87	72.27	<b>66.36</b>	<b>55.44</b>	<b>45.10</b>	39.29	<b>63.59</b>	83.9	1.03%
<b>MoRAM</b>	<b>47.56</b>	<b>48.41</b>	70.09	<b>82.46</b>	<b>73.53</b>	<b>68.18</b>	<b>55.53</b>	<b>46.08</b>	<b>41.96</b>	<b>63.70</b>	41.9/26.2	0.52%/0.32%

Retained (and possibly improved) Out-of-Domain Generalization After SFT on Code-Alpaca.

### MoE-LoRA · adds new sets of rank- $r$ experts for each task

Each new task contributes 2 rank- $r$  experts ( $A \in \mathbb{R}^{r \times d}$ ,  $B \in \mathbb{R}^{d \times r}$ ,  $r=4$ ); router  $R(x)$  assigns experts — granularity is whole-expert



### MoRAM · linear associative memory of rank-1 key-value atoms

Each new task appends 4 rank-1 memory entries ( $A_i \in \mathbb{R}^{1 \times d}$  key,  $B_i \in \mathbb{R}^{d \times 1}$  value); content-addressable retrieval via  $A_i \cdot x$  — no router

