









































Evaluation Datasets and Their Difficulties

• Each Dataset: contain one task like classify movie tags or summarize conversation content

- Level of difficulty
 - Classification: LaMP-2 < LaMP-3 < LaMP-1
 - Generation: LaMP-6 < LaMP-5 < LaMP-7 < LaMP-4
 - · Classification in general is easier than generation

Dataset	GPT-4	Claude 3 Opus	Gemini 1.0 Pro	Llama 3 70B	Average	Normalized Performance	Task Type
LaMP-1 (Accuracy)	0.539	0.539	0.490	0.422	0.4975	0.995	Classification
LaMP-2 (Accuracy)	0.355	0.320	0.300	0.400	0.3438	5.157	Classification
LaMP-3 (Accuracy)	0.667	0.657	0.510	0.755	0.6472	3.236	Classification
LaMP-4 (ROUGE-1)	0.143	0.171	0.139	0.131	0.1460	0.1460	Generation
LaMP-5 (ROUGE-1)	0.386	0.374	0.405	0.084	0.3123	0.3123	Generation
LaMP-6 (ROUGE-1)	0.351	0.356	0.405	0.278	0.3475	0.3475	Generation
LaMP-7 (ROUGE-1)	0.326	0.136	0.237	0.255	0.2385	0.2385	Generation





























Estimate the information volume, Keep domain-specific	c data, Dr	op correlated data
Metric 1: Entropy of Embedding		
• Get embedding from pretrained LLM (PLM)	Domain Admin	Example Lexicons dose vial inhale inject ml pills ingredient
• Use Shannon Entropy, normalized by logistic sequence length	Anatomy Drug	Pelvis arm sinus breast chest lymph tonsil ACOVA ACTONEL CARTIA EMGEL
to estimate the amount of information in data	Fear Surprise Trust	bunker cartridge cautionary chasm cleave amazingly hilarious lucky merriment advocate alliance canons cohesion
Keep high information data	GloVeTW GloVeCC	extreme potential activity impact movementsymptomatic thrombosis fibrillation
Metric 2: Domain Specific Score	GloveTW	775 nyquil benadryl midol pepto midol ritalin
• Get the input domain		
• Count the domain-related tokens	Metric	1 EOE $(\vec{\mathbf{E}}_i) = \frac{-\sum_{e_i \in \vec{\mathbf{E}}} p(e_i) \log p(e_i)}{\log(n)}$
Keep the most domain-related input data		m m in a li
Metric 3: In-Domain Dissimilarity	Metric	2 DSS(<i>T</i> , <i>L</i>) = $\frac{1}{m} \sum_{i=1}^{m} \frac{ T \cap l_i }{n}$
• Within the same domain, we keep the most distinct data		$1 \frac{R}{R}$
Reuse embeddings from PLM	Metric	3 IDD($\vec{\mathbf{E}}, B$) = $\frac{1}{R} \sum_{i=1}^{N} (1 - \cos(\vec{\mathbf{E}}, \vec{\mathbf{E}}_{Dom_d}^i))$











MIPS in RAG

- What to store:
 - One user-generated text → Sentence Embedding Model → One stored vector
 - Many user-generated texts → Vectors → "Can formalize a matrix"
- Query:
 - Original prompt → Sentence Embedding Model → One query vector
- Find the appropriate user-generated text:
 - Inner products between query vector and every store vector
 - Rank the stored vectors based on the products
 - Max Inner Product Search (MIPS)

- User-generated data:
 - Save all the user-generated data →
 Resource on edge is limited
- Manage data:
 - Save on RAM: Easy for compute, but take up resources for other applications
 - Save on Disk: Data movement can lead to latency (Longer than LLM Inference)





















RoCR: Noise-ware Trainin	g								
 Noise injection: [1, 0.75], [0.75, 0.5], [0.5, 0.25], [0.25, 0] (4 states in NVM), each range corresponding a variance level, shown as Figure 1 Concatenating with gaussian distribution (default 		Name RRAM ₁ (Device-1) FeFET ₂ (Device-2) FeFET ₃ (Device-3) RRAM ₄ (Device-3) Figure 1: D different	# of Levels	L ₀ 0.0100 0.0067 0.0049 0.0038 0.0026 non-id	$ \begin{array}{r} \hline Device Va \\ L_1 \\ \hline 0.0100 \\ 0.0135 \\ 0.0146 \\ 0.0155 \\ \hline 0.0155 \\ \hline eality 1 \\ the size \\ \end{array} $	$\frac{L_2}{\frac{L_2}{0.0100}}$	L ₃ 0.0100 0.0067 0.0049 0.0038 0.0026 ng for ces		
 to 0.1) During training: Noise are added to embedding, shown as Figure 2 Rationale: 	1 def add_noise(embe 2 w_n = embeddin 3 w_n(w_n > 0.75 4 mask = (w_n < 5 w_n(mask) = w_ 6 mask = (w_n < 7 w_n(mask) = w_ 8 mask = (w_n <	ddings, noise_factor gs / embeddings.abs(] = $w_n f w_n > 0.75$] 0.75) * ($w_n > 0.5$ (mask) + torch.rand 0.5) * ($w_n > 0.25$ n[mask] + torch.rand 0.25)	_1, noise_f).max().ite + torch.ran) n_like(w_n[) n_like(w_n[actor_2, n() dn_like(w nask]) *	noise_facto _n[w_n > 0. noise_facto noise_facto	or_3, noise 75]) * noi or_2 * gaus or_3 * gaus	_factor_4, g se_factor_1 sian_noise_s sian_noise_s	aussian_noise_ * gaussian_noi igma igma	_sigma) ise_sig
 When injected noise will not lead to undesirable LLM generating, we stop training 	<pre>9 w_n[mask] = w_n[mask] + torch.randn_like(w_n[mask]) * noise_factor_4 * gaussian_noise_sigma 10 11 return w_n Figure 2: Noise injection</pre>								
									4























Co-design: Retrieval algorithm

- High-level Description: Adapter-level search (More complicated than MIPS)
- Concept:
 - Virtual Tokens (adapter), different from that in RAG, are integrated into a matrix.
 - Instead of vector (input) and matrix (stored data) multiplication, **matrix** (input) and **matrices** (stored adapters) **multiplications** are need
- Propose: Weighted Multi-Scale Dot Product search (WMSDP)
 - Scale: average pooling adapters
 - Weighted: on designed factors 1, 2, 4
 - Dot Product: Between input matrix and every stored matrix
 - Rationale: Information stored in adapter is more hidden, compared to sentence embedding data in RAG

61

























Performance and Conclusion

- Our framework (Tiny-Align) demonstrates decent performance and efficiency on different audio datasets
- It can benefit people with Dementia, Aphasia, and Specific Language Impairment. Why? On-device learning for personalized audio input, so the LLM can understand such audio input, and process with its strong reasoning capability
- Echo:
 - RAG-CiM and NVCiM-PT: Can we use novel circuits and energy efficient hardware to further optimize Tiny-Align?
 - Empirical Study and Data Selection: Can we use traditional devices and algorithms to optimize Tiny-Align?

Dataset	Mathod	ASR + Llama-3.2-1B			ASR + Llama-3.2-3B			ASR + Gemma-2-2B			ASR + Phi-3.5-mini		
	wiedlou	R-1	R-L	C-T(10s)	R-1	R-L	C-T(10s)	R-1	R-L	C-T(s)	R-1	R-L	C-T(10s)
eSS	A1	0.021	0.025	293	0.032	0.027	293	0.051	0.049	2249	0.103	0.072	2520
	A2	0.152	0.138	3107	0.134	0.128	3107	0.196	0.112	7338	0.213	0.119	11875
DR	A3	0.104	0.096	1128	0.108	0.103	1128	0.185	0.119	1841	0.029	0.026	2093
A	Ours	0.192	0.164	111	0.205	0.193	111	0.268	0.154	97	0.225	0.121	104
Baycrest	A1	0.024	0.024	792	0.024	0.024	792	0.017	0.016	11517	0.050	0.047	7375
	A2	0.141	0.104	949	0.152	0.109	949	0.183	0.114	4614	0.202	0.110	3847
	A3	0.065	0.057	3015	0.071	0.067	3015	0.088	0.060	9951	0.012	0.011	6903
	Ours	0.184	0.167	167	0.197	0.173	167	0.233	0.141	133	0.205	0.112	183
Perf	ormance	compar	ison be	tween our	framew	vork wi	th existing	metho	ds. Met	rics incl	uding R	OUGE	-1

(R-1), ROUGE-L (R-L), and Convergence Time (C-T)

Acknowledgement: Ruiyang Qin

Ph.D. candidate at Notre Dame (Aug 2022 - Present)

Research Area:

- Edge AI
- AI Healthcare
- Emerging Semiconductor Devices

Publications:

- EDA (DAC, ICCAD, DATE, ASP-DAC)
- Information System (TOIS, TMIS)
- HCI (RO-MAN, CHI-EA)

Experience:

- BS/MS in CS at Georgia Tech (2017-2021)
- Visiting Scholar at SUNY-Buffalo (2024)

Services:

NeurIPS, ICLR, ICML, AAAI, AISTATS, RO-MAN, JMIR,

ruiyangqin2016.github.io/

rqin@nd.edu

TNSRE (Trans on Neural Systems & Rehabilitation Engineering)

The College of Engineering at the University of Notre Dame