

Iterative Context Vectors: Boost In-Context Learning within Activations

Anonymous ACL submission

Abstract

In-context learning has become a standard learning paradigm for language models. However, current prompt engineering methods, which function within the token space, may restrict their effectiveness. We propose to explore the potential of activation space through Iterative Context Vectors (ICVs), a technique aimed at improving task performance without backpropagation. ICVs are employed by first extracting and iteratively refining activations within a language model, then applying them during inference with minimal computational and memory overhead. We evaluate ICVs across a range of tasks using various models and observe significant improvements. Our findings suggest that activation steering can serve as a promising direction for in-context learning, thereby opening new avenues for future research.

1 Introduction

Few-shot learning has long been a prominent research focus. Recently, language models (LMs) have shown the capability to execute few-shot learning through in-context learning (ICL) (Brown et al., 2020). In this approach, learning a new task involves conditioning on a few support examples and predicting the most suitable tokens to complete a query input, all without the need for any parameter updates. This method is appealing because it relies solely on inference, allowing for quick adaptation to various downstream tasks.

However, it has been noted that despite its potential, the predictions of LMs can be highly volatile when conditioned on prompts. The outcomes depend significantly on the templates, demonstrations, their permutations, and can even ignore or violate the instructions of the prompt (Webson and Pavlick, 2022; Min et al., 2022b).

In this paper, we introduce Iterative Context Vectors (ICVs) to offer a new perspective. As illustrated in Figure 1, rather than remaining in the

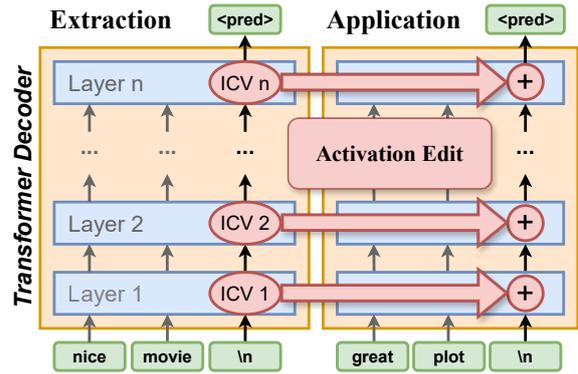


Figure 1: Iterative Context Vectors improve ICL performance by modifying model activations.

discrete prompt space, ICVs delve into the extensive activation space of the model. This exploration reveals a largely uncharted area for developing new methods, with our pioneering efforts to demonstrate how ICL can be enhanced from the representations within the model.

ICV contrasts with existing prompt tuning methods (Li and Liang, 2021; Lester et al., 2021), which operates in a continuous parameter space but still as part of the prompt and requires training via backpropagation. Again, unlike Parameter-Efficient Fine-Tuning (PEFT) methods, e.g. LoRA (Hu et al., 2021), ICV does not seek to tune the parameters of the model but rather modifies the activations during inference.

The essential traits of ICV, which will be elaborated upon in the rest of the paper, are highlighted as follows:

1. ICV has an intuitive theoretical support.
2. ICV is independent of instruction, prompt, label and permutation choices.
3. ICV does not require backpropagation.

To the best of our knowledge, we are the first to investigate the application of activation vectors on diverse real-world in-context learning tasks and to demonstrate their potential with in-context exam-

069 ples during inference.

070 2 Method

071 We begin by establishing the evaluation framework.

072 2.1 Activation Vector Evaluation

073 We adhere to standard few-shot benchmarking pro-
074 tocols (Vinyals et al., 2016; Finn et al., 2017; Snell
075 et al., 2017) to define the activation vector evalua-
076 tion setting. For a given split of an n -way k -shot
077 classification task $\mathcal{T} = \{\mathcal{T}_{\text{train}}, \mathcal{T}_{\text{val}}, \mathcal{T}_{\text{test}}\}$, which
078 includes textual query-answer pairs (x, y) , an ICL
079 *episode* is sampled as:

$$080 E = [(x_1, y_1), \dots, (x_{n \times k}, y_{n \times k}), (x_q, y_q)]. \quad (1)$$

081 Here, (x_q, y_q) represents the query and its label,
082 preceded by the $n \times k$ support examples. To avoid
083 the impact of unbalanced samples during extrac-
084 tion, we uniformly sample k examples from each
085 of the n classes and shuffle them to mitigate any
086 bias arising from sample permutation.

087 The episode must first be converted into a pure
088 text sequence before the language model $\text{LM}(\cdot)$
089 can process it. This conversion is handled by a *ver-*
090 *balizer*, which uses a predefined prompt template
091 to instantiate the samples. The template contains
092 two key components: the *input-output separator*
093 that links a question with its answer, and the *exam-*
094 *ple separator* that joins the given support set. To
095 preserve the simplicity of the template, we have
096 chosen to use one newline ($\backslash n$) for the input-output
097 separator and three newlines for the example sepa-
098 rator, as adopted in Min et al. (2022a). In the sub-
099 sequent text, the verbalizer will be considered an
100 integral component of the Language Model (LM)
101 and will not be explicitly referenced for the sake of
102 conciseness.

103 When the language model $\text{LM}(\cdot)$ is given an
104 episode E , it executes autoregressive inference on
105 each of its tokens. The input-output tokens are
106 particularly noteworthy because they are responsi-
107 ble for producing the answers. The prediction of
108 the LM can be obtained by applying the softmax
109 function to the logits of the possible labels.

$$110 \hat{y}_{\text{clean}} = \text{LM}(E). \quad (2)$$

111 In contrast, an “edited” run utilizes an *activation*
112 *vector editor* f_{edit} , represented as

$$113 \hat{y}_{\text{edit}} = \text{LM}(E; f_{\text{edit}}(\mathbf{v}, p)), \quad (3)$$

114 which relies on the extracted vectors \mathbf{v} from an *ac-*
115 *tivation vector extractor* f_{ext} with hyperparameters
116 p :

$$117 \mathbf{v} = f_{\text{ext}}(\mathcal{T}_{\text{train}}; p). \quad (4)$$

118 The extractor retrieves its target vectors \mathbf{v} from
119 $\mathcal{T}_{\text{train}}$ and identifies the optimal hyperparameters p^*
120 from \mathcal{T}_{val} by maximizing the metric M :

$$121 p^* = \arg \max_p M_{E \sim \mathcal{T}_{\text{val}}}(\hat{y}_{\text{edit}}, y_q) \quad (5)$$

$$122 \mathbf{v}^* = f_{\text{ext}}(\mathcal{T}_{\text{train}}; p^*).$$

123 For single-token classification tasks, macro-F1,
124 micro-F1, and weighted-F1 scores can serve as
125 the metrics. The vectors \mathbf{v}^* and the optimal hyper-
126 parameters p^* are then applied to the test set $\mathcal{T}_{\text{test}}$
127 to evaluate the final results $M_{E \sim \mathcal{T}_{\text{test}}}(\hat{y}_{\text{edit}}, y_q)$.

128 Having outlined the evaluation framework, we
129 will now move on to the theoretical grounds of our
130 method.

131 2.2 Theoretical Foundation

132 Given the significance of ICL, many theories have
133 been suggested to explain its mechanism, e.g. Xie
134 et al. (2022); Chan et al. (2022); Ye et al. (2023);
135 Oswald et al. (2023). Drawing inspiration from Irie
136 et al. (2022), we construct our ICV based on the
137 empirical evidence provided by Dai et al. (2023).

138 Irie et al. (2022) revisited the dual form of the
139 perceptron and applied it in the modern context
140 of deep NNs. They demonstrated that the forward
141 operation of any linear layer in neural networks
142 trained via gradient descent can be viewed as a key-
143 value-query attention mechanism (Vaswani et al.,
144 2017). In this framework, the training data points
145 act as the keys, the corresponding gradients serve
146 as the values, and the test input generates the query.
147 A more detailed introduction to the dual form is
148 provided in Appendix B.

149 With the help of the dual form, Dai et al. (2023)
150 showed that ICL can be interpreted as a meta-
151 optimization process. This was achieved by revers-
152 ing the direction of the equivalence and breaking
153 down the attention key and value terms for the ICL
154 query token into its zero-shot and demonstration
155 components. Under the relaxed normalization set-
156 ting, the pretrained LM acts as a meta-optimizer.
157 Through forward computation, the LM generates
158 meta-gradients from the demonstration examples,
159 which are then applied to the original language
160 model via attention, culminating in the formation
161 of the ICL inference capability. Their experiments

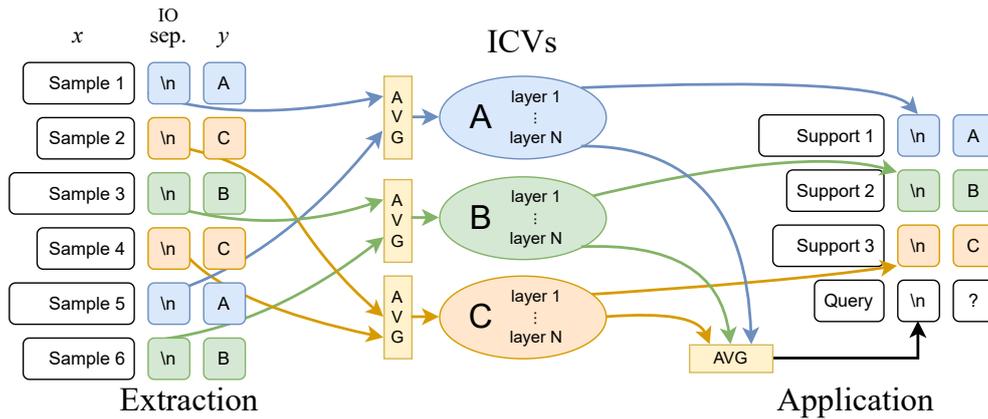


Figure 2: Illustration of the extraction and application phases of ICV. For clarity, contrastive subtraction and iterative updates have been omitted.

verified that ICL behaves similarly to explicit fine-tuning from multiple perspectives.

2.3 Iterative Context Vectors

We have determined that attention layers significantly influence ICL, with demonstrations acting as meta-gradients to help the model adapt to the task during inference. This explanation provides an intuitive understanding of how the LM uses in-context examples, but it also highlights why ICL performance can be unstable. Specifically, meta-gradients derived from limited in-context examples may not fully capture the task and may not fit well with the initial parameters. For this reason, we propose to extract the meta-gradients from the LM’s inference process to improve their accuracy and robustness. This would allow us to apply these meta-gradients directly in future inference tasks, eliminating the need to compute them afresh with ICL each time a query is evaluated.

To define ICV, we first specify the extractor f_{ext} . To simulate the gradients, we generate two versions of a given n -way k -shot episode E in a contrastive manner, where $k \geq 1$ is a hyperparameter. The positive sequence is the standard shuffled verbalization, serving as the target for the gradients. The negative sequence can have various design choices; we choose to use a zero-shot query, which provides no information about the task.

The extractor then identifies the activations for all $n \times k$ input-output tokens of the support set (if one exists) and the final input-output token for the query in each attention layer of the LM. When $k > 1$, we initially average the activations for each class. Subsequently, we subtract the negative activations from the positive activations, thereby obtaining the gradients for a single episode. Given

that there is no support set in the negative sequence, all activations from the positive support set share a common subtrahend of the negative query. By averaging over the training set, a preliminary version of the vectors can be calculated, as illustrated in Figure 2.

Next, to better utilize the forward pass computation, we propose to apply the vectors during the extraction phase, thus introducing the concept of *Iterative* Context Vectors. Specifically, we implement a batch-like update strategy, simulating standard batched gradient, which has been generally adopted to reduce the instability of single-step gradients. After every b extraction episodes, the vectors extracted from all previous episodes are averaged and used as the ICVs during subsequent extractions, leading us to the definition of the editor f_{edit} .

For the l -th attention layer $\text{Attn}_l(\cdot)$, we have the corresponding extracted ICV v_l . During inference, the editing is executed following Eq. 10

$$\text{EditAttn}_l(x) := \text{Attn}_l(x) + \alpha \times v_l, \quad (7)$$

where two additional hyperparameters are introduced: the extraction strength α_1 and the inference strength α_2 , adopted during the extraction and inference phrases, respectively. In summary, the hyperparameters for the ICVs are $p = \{k, b, \alpha_1, \alpha_2\}$.

3 Experiments

We apply our ICVs to three popular models across 12 tasks. The results are shown in Table 1. Details of the datasets can be found in Appendix C.

During the few-shot testing process, the model cannot ascertain the true class distribution of the test set, which is often imbalanced. Therefore, we adhere to the one-shot testing design, which sup-

Model	Task	agnews	emot.	hate	irony	offe.	sent.	abor.	athe.	clim.	femi.	hill.	trec	Avg.
gpt-j-6b	Clean	53.53	24.07	49.38	55.93	51.98	36.94	32.96	25.38	27.11	31.80	35.74	44.83	39.14
	FV	37.95	10.72	36.36	37.80	41.91	21.92	27.07	28.21	28.20	28.03	24.86	11.26	27.86
	TV	62.46	26.12	50.17	55.53	52.05	38.72	31.56	25.57	29.45	31.67	35.83	51.94	40.92
	ICV	62.63	20.26	51.19	63.27	52.54	34.55	37.74	33.09	36.46	38.66	38.65	40.66	42.48
llama-2-7b	Clean	61.94	54.45	53.27	58.65	51.86	38.96	27.52	22.13	28.60	29.27	29.42	56.56	42.72
	FV	23.68	16.67	50.76	38.70	21.76	12.60	23.26	10.33	27.05	27.30	20.08	5.92	23.18
	TV	70.93	59.68	52.44	50.48	54.05	43.67	27.90	21.83	32.04	29.31	32.99	56.61	44.33
	ICV	67.15	38.19	57.36	66.03	58.39	45.56	31.00	22.66	32.70	29.16	30.09	61.80	45.00
llama-2-13b	Clean	76.23	61.89	53.83	55.17	60.34	38.77	34.96	27.11	20.96	37.13	45.53	61.10	47.75
	FV	51.54	10.15	36.35	54.79	21.76	23.97	9.16	7.48	28.21	8.98	13.65	13.11	23.26
	TV	76.03	63.74	54.47	55.36	60.55	38.38	35.12	30.08	28.33	37.15	44.66	65.69	49.13
	ICV	83.48	65.51	54.43	53.66	62.01	44.03	34.98	26.11	36.63	47.08	54.69	72.67	52.94

Table 1: Main experiment results with macro-F1 as the metric. "Clean" denotes a standard one-shot ICL result. The models are GPT-J-6B (Wang and Komatsuzaki, 2021) and Llama 2 (Touvron et al., 2023).

plies the model with minimal yet sufficient information through a uniformly distributed support set. We evaluate over 200 episodes for both extraction and hyperparameter search, with a fixed iterative batch size $b = 10$ for all tasks. For other hyperparameters of ICVs, we search for the extraction shot $k \in \{1, 2, 3, 4\}$, the extraction strength and the inference strength $\alpha_1, \alpha_2 \in \{0.1, 0.3, 0.5, 0.7, 0.9\}$. The final testing results are averaged over 10,000 randomly sampled episodes.

As further proof of concept and baselines for comparison, we also include two recent activation vector proposals: Function Vectors (Todd et al., 2023) and Task Vectors (Hendel et al., 2023). Although these methods were not originally designed to operate under the ICL evaluation setting, we adapted them to utilize the training set by averaging the activations. We search over their respective hyperparameters as well as the extraction shot k to ensure a fair comparison. Please refer to Appendix A for a discussion of their designs.

The results demonstrate that ICVs generally enhance ICL performance, surpassing the baselines in most tasks as well as in the overall average. Despite its simpler design, Task Vectors prove to be surprisingly competitive and can serve as a robust baseline. Conversely, Function Vectors hardly contribute to performance enhancement. Due to FVs' high search time stemming from their design, they may necessitate substantially more effort for optimization in real ICL applications.

To achieve a more comprehensive understanding of ICVs, we compare their performance against the standard ICL using additional shots, as illustrated in Table 2. Firstly, it is observed that ICVs not only surpass the majority of 1-shot performances but also often match or exceed performances with

Task	agnews	hate	irony	offe.	sent.	abor.
0-shot	51.14	45.99	59.25	47.60	32.83	26.51
1-shot	62.89	49.71	44.65	52.85	41.27	27.25
2-shot	76.30	57.74	53.11	57.03	45.01	22.98
3-shot	80.05	60.10	58.31	59.45	45.83	18.90
4-shot	81.44	61.47	54.45	58.10	48.39	18.20
ICV	67.15	57.36	66.03	58.39	45.56	31.00

Table 2: Comparison between ICV and standard ICL on Llama-2-7b with macro-F1 as the metric. See Table 6 for complete results.

higher numbers of shots. Secondly, the data indicates that the performance of standard ICL does not always improve with an increased number of demonstration examples. This suggests potential limitations of solely relying on prompt engineering to enhance performance. Thirdly, it is crucial to note that the temporal cost of standard ICL theoretically scales with $O(n^2)$. Although it is feasible to cache keys and values in practice, these caches cannot be reused when a new set of examples is introduced, which is likely to occur due to the inherent difficulty in identifying and ensuring a single set of examples that are effective for all inputs.

Finally, an ablation study examining the iterative batch size b has been conducted and is presented in Appendix D.

4 Conclusion

In our study, we have derived the Iterative Context Vectors (ICVs) from an intuitive theoretical framework, defined the evaluation protocols and subsequently conducted a series of experiments. Despite ICVs' simplicity, the results obtained are highly encouraging, indicating that activation vectors show significant potential for further exploration.

293 Limitations

294 This study examines the application of Iterative
295 Context Vectors in the context of one-shot exam-
296 ples as a compromise between inference time and
297 in-context information. Although applying ICVs to
298 zero-shot inference would be more efficient, a com-
299 putational sequence of insufficient length might
300 hinder the model’s ability to effectively solve the
301 given task. (Feng et al., 2023)

302 We have opted for classification tasks wherein
303 a single output token is sufficient to distinguish
304 between the classes. The development and appli-
305 cation of activation vectors in more complex tasks,
306 as well as in generative tasks, represent areas for
307 future investigation. Nevertheless, it is worth not-
308 ing that the concept of ICVs and the associated
309 evaluation protocol can potentially be expanded to
310 encompass these more advanced applications.

311 References

312 Francesco Barbieri, Jose Camacho-Collados, Luis Es-
313 pinosa Anke, and Leonardo Neves. 2020. [TweetEval:
314 Unified Benchmark and Comparative Evaluation for
315 Tweet Classification](#). In *Findings of the Association
316 for Computational Linguistics: EMNLP 2020*, pages
317 1644–1650, Online. Association for Computational
318 Linguistics.

319 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
320 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
321 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
322 Askell, Sandhini Agarwal, Ariel Herbert-Voss,
323 Gretchen Krueger, Tom Henighan, Rewon Child,
324 Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens
325 Winter, Chris Hesse, Mark Chen, Eric Sigler, Ma-
326 teusz Litwin, Scott Gray, Benjamin Chess, Jack
327 Clark, Christopher Berner, Sam McCandlish, Alec
328 Radford, Ilya Sutskever, and Dario Amodei. 2020.
329 [Language Models are Few-Shot Learners](#). In *Ad-
330 vances in Neural Information Processing Systems*,
331 volume 33, pages 1877–1901. Curran Associates,
332 Inc.

333 Stephanie C. Y. Chan, Adam Santoro, Andrew K.
334 Lampinen, Jane X. Wang, Aaditya Singh, Pierre H.
335 Richemond, Jay McClelland, and Felix Hill. 2022.
336 [Data Distributional Properties Drive Emergent In-
337 Context Learning in Transformers](#). *Preprint*,
338 arxiv:2205.05055.

339 Damai Dai, Yutao Sun, Li Dong, Yaru Hao, Shuming
340 Ma, Zhifang Sui, and Furu Wei. 2023. [Why Can
341 GPT Learn In-Context? Language Models Implicitly
342 Perform Gradient Descent as Meta-Optimizers](#).
343 *Preprint*, arxiv:2212.10559.

344 Guhao Feng, Bohang Zhang, Yuntian Gu, Haotian Ye,
345 Di He, and Liwei Wang. 2023. [Towards Revealing](#)

the Mystery behind Chain of Thought: A Theoret-
ical Perspective. *Advances in Neural Information
Processing Systems*, 36:70757–70798.

Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017.
Model-agnostic meta-learning for fast adaptation of
deep networks. In *International Conference on Ma-
chine Learning*, pages 1126–1135. PMLR.

Roe Hendel, Mor Geva, and Amir Globerson. 2023.
[In-Context Learning Creates Task Vectors](#). In *Find-
ings of the Association for Computational Linguis-
tics: EMNLP 2023*, pages 9318–9333, Singapore.
Association for Computational Linguistics.

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan
Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and
Weizhu Chen. 2021. [LoRA: Low-Rank Adaptation
of Large Language Models](#). In *International Confer-
ence on Learning Representations*.

Kazuki Irie, Róbert Csordás, and Jürgen Schmidhuber.
2022. [The Dual Form of Neural Networks Revisited:
Connecting Test Time Predictions to Training Pat-
terns via Spotlights of Attention](#). In *Proceedings of
the 39th International Conference on Machine Learn-
ing*, pages 9639–9659. PMLR.

Brian Lester, Rami Al-Rfou, and Noah Constant. 2021.
[The Power of Scale for Parameter-Efficient Prompt
Tuning](#). In *Proceedings of the 2021 Conference on
Empirical Methods in Natural Language Processing*,
pages 3045–3059, Online and Punta Cana, Domini-
can Republic. Association for Computational Lin-
guistics.

Quentin Lhoest, Albert Villanova del Moral, Yacine
Jernite, Abhishek Thakur, Patrick von Platen, Suraj
Patil, Julien Chaumond, Mariama Drame, Julien Plu,
Lewis Tunstall, Joe Davison, Mario Šaško, Gun-
jan Chhablani, Bhavitvya Malik, Simon Brandeis,
Teven Le Scao, Victor Sanh, Canwen Xu, Nicolas
Patry, Angelina McMillan-Major, Philipp Schmid,
Sylvain Gugger, Clément Delangue, Théo Matus-
sière, Lysandre Debut, Stas Bekman, Pierric Cistac,
Thibault Goehringer, Victor Mustar, François Lagu-
nas, Alexander M. Rush, and Thomas Wolf. 2021.
[Datasets: A Community Library for Natural Lan-
guage Processing](#). *Preprint*, arxiv:2109.02846.

Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter
Pfister, and Martin Wattenberg. 2023. [Inference-
Time Intervention: Eliciting Truthful Answers from
a Language Model](#). *Advances in Neural Information
Processing Systems*, 36:41451–41530.

Xiang Lisa Li and Percy Liang. 2021. [Prefix-Tuning:
Optimizing Continuous Prompts for Generation](#). In
*Proceedings of the 59th Annual Meeting of the Asso-
ciation for Computational Linguistics and the 11th
International Joint Conference on Natural Language
Processing (Volume 1: Long Papers)*, pages 4582–
4597, Online. Association for Computational Lin-
guistics.

A Related Work

A.1 Activation Vectors

Some preliminary works have recently explored steering LMs in the representation space. Task Vectors (Hendel et al., 2023) are extracted from one layer of the model during ICL inference and then applied to a zero-shot query to determine whether they can preserve task-relevant information. Function Vectors (Todd et al., 2023), on the other hand, select activations from the top attention heads based on their causal effect in promoting the correct answer, average these activations, and add them to a specific layer.

Although these vectors are designed with intentions similar to ours, they are tested primarily on simple synthetic tasks like antonym, country-capital, and singular-plural pairs. In contrast, we target a practical setting by evaluating on real-world datasets, providing a more comprehensive assessment ground.

A.2 Generative Steering

Another research direction focuses on modifying LMs’ activations for generation and transfer purposes. Latent Steering Vectors (Subramani et al., 2022) aim at sentence recovery and sentiment transfer. Inference-Time Intervention (Li et al., 2023) involves probing each attention head and guiding the model with the probe vector to enhance the truthfulness of the generated text. Studies by Turner et al. (2023) and Liu et al. (2023) address style and sentiment transfer by employing positive and negative sentence pairs to extract contrastive guidance.

Despite their similarities, these methods either require training with backpropagation or are specifically tailored for generative or transfer tasks between sentence pairs. Consequently, they cannot be directly integrated into our approach.

B The Dual Form of Attention Layers

Formally, assume a linear layer trained via gradient descent utilizing T training inputs $(\mathbf{x}_1, \dots, \mathbf{x}_T)$ and their corresponding (backpropagated) error signals $(\mathbf{e}_1, \dots, \mathbf{e}_T)$, where $\mathbf{x}_t \in \mathbb{R}^{d_{in}}$ and $\mathbf{e}_t \in \mathbb{R}^{d_{out}}$. If standard gradient descent is applied, a loss function \mathcal{L} produces the error signal $\mathbf{e}_t = -\eta_t(\nabla_{\mathbf{y}}\mathcal{L})_t$, where $\eta_t \in \mathbb{R}$ is the learning rate, and $\mathbf{y}_t = \mathbf{W}\mathbf{x}_t$ is the output of the linear layer. Its weight matrix is given by

$$\mathbf{W} = \mathbf{W}_0 + \sum_{t=1}^T \mathbf{e}_t \otimes \mathbf{x}_t, \quad (8)$$

where $\mathbf{W}_0 \in \mathbb{R}^{d_{out} \times d_{in}}$ represents the initial value of the weights. This linear layer transforms an input $\mathbf{x} \in \mathbb{R}^{d_{in}}$ into an output $S_1(\mathbf{x}) \in \mathbb{R}^{d_{out}}$:

$$S_1(\mathbf{x}) = \mathbf{W}\mathbf{x}. \quad (9)$$

Next, consider a composite layer S_2 that stores T key-value pairs, represented by a key matrix $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_T) \in \mathbb{R}^{d_{in} \times T}$ and a value matrix $\mathbf{E} = (\mathbf{e}_1, \dots, \mathbf{e}_T) \in \mathbb{R}^{d_{out} \times T}$, along with a weight matrix $\mathbf{W}_0 \in \mathbb{R}^{d_{out} \times d_{in}}$. This layer transforms an input $\mathbf{x} \in \mathbb{R}^{d_{in}}$ into an output $S_2(\mathbf{x}) \in \mathbb{R}^{d_{out}}$ by

$$S_2(\mathbf{x}) = \mathbf{W}_0\mathbf{x} + \text{Attn}(\mathbf{X}, \mathbf{E}, \mathbf{x}), \quad (10)$$

where the parameters of the unnormalized attention operator $\text{Attn}(\cdot)$ are, in order, the key, value, and query.

It can be shown that S_1 and S_2 are equivalent by expanding the attention operation as

$$\text{Attn}(\mathbf{X}, \mathbf{E}, \mathbf{x}) = \mathbf{E}\mathbf{X}^\top \mathbf{x} = \left(\sum_{t=1}^T \mathbf{e}_t \otimes \mathbf{x}_t \right) \mathbf{x}. \quad (11)$$

C Dataset and Tasks

All datasets utilized in this research are obtained from Huggingface (Lhoest et al., 2021). A full list of these datasets, along with their corresponding access labels, is detailed in Table 3.

AG News (Zhang et al., 2015) is a subdataset of AG’s corpus of news articles constructed by assembling titles and description fields of articles from the 4 largest classes (“World”, “Sports”, “Business”, “Sci/Tech”) of AG’s Corpus.

TweetEval (Barbieri et al., 2020) introduces an evaluation framework consisting of a series of Twitter-specific classification tasks. We selected all single-token classification tasks from the dataset.

Text Retrieval Conference Question Answering (TrecQA) (Wang et al., 2007) is a dataset created from the TREC-8 (1999) to TREC-13 (2004) Question Answering tracks.

Our few-shot evaluation methodology employs episodic sampling to regulate the duration of both extraction and inference processes, rather than relying solely on the absolute number of samples.

Name	Abbr.	Huggingface Label
AG News	agnews	ag_news
Emotion	emot.	tweet_eval/emotion
Hate	hate	tweet_eval/hate
Irony	irony	tweet_eval/irony
Offensive	offe.	tweet_eval/offensive
Sentiment	sent.	tweet_eval/sentiment
Abortion	abor.	tweet_eval/stance_abortion
Atheism	athe.	tweet_eval/stance_atheism
Climate	clim.	tweet_eval/stance_climate
Feminist	femi.	tweet_eval/stance_feminist
Hillary	hill.	tweet_eval/stance_hillary
TREC	trec	trec

Table 3: The datasets and tasks employed, along with their corresponding abbreviations used in the result tables, and their respective labels as hosted on Hugging Face.

Consequently, not all available samples are utilized during the experimental procedures. This aspect underscores an additional dimension of efficiency inherent in activation vectors.

We utilize the default labels provided with the datasets to emphasize the independence of prompt formatting. It is noteworthy that, for the TREC dataset, the observed low zero-shot performance can be attributed to the default labels, which are capitalized abbreviations such as "ABBR", "ENTY", and "LOC". This particular dataset functions as a special case within our experiments, aimed at investigating whether ICVs can adapt to labels that are less semantically meaningful.

D Additional Results

More metrics of the main experiment. We present the results on the other two metrics, namely micro-F1 and weighted-F1, derived from our main experiment, in Table 4 and Table 5, respectively. Under these evaluation criteria, ICVs demonstrate performance on par with, and often surpassing, that of FV and TV across the majority of tasks.

All experiments were conducted utilizing a pre-determined random seed (42) to mitigate selection bias. To ensure a robust representation of result distributions, the tests from the main experiments were averaged over a substantial number of episodes, specifically 10,000.

Notably, the exception lies in the GPT-J-6B & micro-F1 setting, where FV exhibits superior performance. Given the pronounced underperformance of FVs across other cases, we hypothesize

that this anomalous result may indicate a strong bias of FVs towards the majority classes and the specific model. This bias results in an elevated micro-F1 score, while simultaneously failing to perform effectively under other evaluation settings.

Ablation on the extraction batch size. Additionally, we make an experiment that ablate on the extraction batch size b . The results are shown in Figure 3. It is evident that without our iterative batching strategy, whether when $b = 1$ or b is excessively large, the extracted ICVs demonstrate suboptimal performance.

As stated in Section 3, we consistently employ $b = 10$ across all models and tasks in other experiments, concentrating our search efforts on other hyperparameters. This approach is intended to limit the search time and to prioritize the identification and optimization of the more significant parameters.

Furthermore, the results suggest that there exists substantial potential for optimization within the design and hyperparameter space of ICVs. This potential will be reserved for future research endeavors aimed at refining and advancing more sophisticated methodologies.

E Code and Reproducibility

The core code defining the ICVs as well as the evaluation protocol is provided within the supplementary material. We will release the complete code repository necessary for reproducing all of our experiments to promote transparency and facilitate future research endeavors.

Model	Task	agnews	emot.	hate	irony	offe.	sent.	abor.	athe.	clim.	femi.	hill.	trec	Avg.
gpt-j-6b	Clean	57.97	31.91	49.39	59.86	63.22	38.73	39.17	30.49	30.92	37.70	40.33	54.01	44.48
	FV	41.45	27.30	57.13	60.76	72.16	48.98	68.37	73.35	73.31	64.94	59.47	18.04	55.44
	TV	66.85	33.23	50.20	59.69	60.83	40.01	37.87	31.27	37.15	37.52	40.80	61.51	46.41
	ICV	63.37	31.55	51.29	64.64	60.97	38.89	55.97	47.76	53.12	43.62	54.14	53.46	51.57
llama-2-7b	Clean	63.40	57.31	53.64	62.22	53.67	40.02	28.69	24.90	34.88	30.25	30.05	60.77	44.98
	FV	34.46	39.73	57.11	61.01	27.82	19.94	24.71	14.02	33.48	64.97	25.82	18.88	35.16
	TV	71.72	62.94	52.92	56.26	56.48	43.10	29.29	24.63	65.82	30.29	33.86	60.81	49.01
	ICV	69.72	36.41	57.38	66.58	64.78	48.22	37.44	26.12	54.66	33.60	37.49	59.02	49.28
llama-2-13b	Clean	77.96	65.42	54.00	55.19	63.56	41.41	52.57	42.78	20.36	55.94	56.83	67.02	54.42
	FV	58.42	25.40	57.10	54.85	27.82	48.83	15.92	12.64	73.37	15.57	25.74	28.43	37.01
	TV	77.89	67.90	54.59	55.40	63.70	40.31	53.25	44.48	42.68	55.95	56.61	70.52	56.94
	ICV	83.77	69.07	54.72	54.23	73.52	43.51	51.01	44.95	49.01	57.68	62.80	75.20	59.96

Table 4: Main experiment results with micro-F1 as the metric. "Clean" denotes a standard one-shot ICL result.

Model	Task	agnews	emot.	hate	irony	offe.	sent.	abor.	athe.	clim.	femi.	hill.	trec	Avg.
gpt-j-6b	Clean	53.69	22.48	49.46	58.64	62.47	33.50	42.61	34.82	34.83	40.34	42.14	51.52	43.88
	FV	37.89	11.71	41.54	45.93	60.49	32.21	55.53	62.07	62.02	51.96	44.36	11.08	43.07
	TV	62.63	23.97	50.34	58.33	61.26	35.46	41.46	35.53	42.27	40.15	42.56	57.36	45.94
	ICV	62.74	23.08	51.50	64.78	61.40	28.80	55.25	51.70	56.21	46.26	50.60	48.26	50.05
llama-2-7b	Clean	76.36	65.73	53.46	54.99	65.44	33.47	51.80	45.57	19.77	53.00	55.25	68.46	53.61
	FV	51.57	10.32	41.51	54.43	12.11	32.04	4.37	2.84	62.10	4.20	10.54	15.84	25.16
	TV	76.18	68.12	54.15	55.11	65.54	32.67	52.24	48.36	45.48	52.99	54.78	70.22	56.32
	ICV	83.56	69.51	53.91	52.56	71.28	42.13	51.52	47.12	52.56	59.37	61.99	75.31	60.07
llama-2-13b	Clean	62.03	57.45	53.83	61.15	56.07	35.33	30.58	27.50	38.72	31.75	27.79	64.49	45.56
	FV	23.85	23.29	52.94	46.74	12.11	8.91	24.87	4.18	34.90	51.57	13.49	6.60	25.29
	TV	71.05	63.17	53.09	53.30	58.79	41.54	31.37	27.49	61.51	31.87	32.48	64.41	49.17
	ICV	67.27	38.23	57.49	66.95	65.62	40.49	40.43	30.16	55.75	34.57	36.80	56.14	49.16

Table 5: Main experiment results with weighted-F1 as the metric. "Clean" denotes a standard one-shot ICL result.

Task	agnews	emot.	hate	irony	offe.	sent.	abor.	athe.	clim.	femi.	hill.	trec	Avg.
0-shot	51.14	30.02	45.99	59.25	47.60	32.83	26.51	23.58	13.22	24.75	31.18	0.59	32.22
1-shot	62.89	44.57	49.71	44.65	52.85	41.27	27.25	23.39	28.55	29.65	29.24	56.17	40.85
2-shot	76.30	53.86	57.74	53.11	57.03	45.01	22.98	21.43	30.88	22.57	31.67	68.62	45.10
3-shot	80.05	58.76	60.10	58.31	59.45	45.83	18.90	17.50	32.09	19.92	30.27	69.71	45.91
4-shot	81.44	59.04	61.47	54.45	58.10	48.39	18.20	14.12	27.67	21.03	27.35	71.33	45.22
ICV	67.15	38.19	57.36	66.03	58.39	45.56	31.00	22.66	32.70	29.16	30.09	61.80	45.01

Table 6: Full comparison between ICV and standard ICL on Llama-2-7b with macro-F1 as the metric. The clean results shown here are averaged over 1,000 episodes.

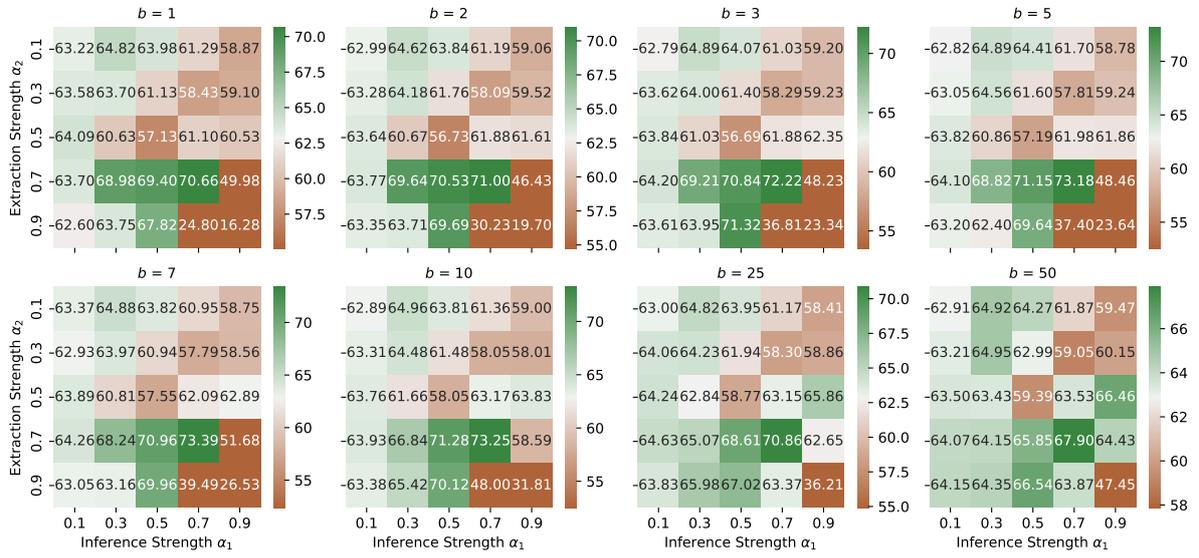


Figure 3: Ablation study on the extraction batch size b . Eight different values of b were evaluated using Llama-2-7b on the AG News dataset. The experiments were conducted under the conditions of 200 one-shot extraction episodes and 1,000 testing episodes. The white color in the heatmaps indicates the corresponding clean one-shot macro-F1 score (62.89, as presented in Table 2).