
Many-Shot In-Context Learning in Multimodal Foundation Models

Anonymous Authors¹

Abstract

Large language models are well-known to be effective at few-shot in-context learning (ICL). Recent advancements in multimodal foundation models have enabled unprecedentedly long context windows, presenting an opportunity to explore their capability to perform ICL with many more demonstrating examples. In this work, we evaluate the performance of multimodal foundation models scaling from few-shot to many-shot ICL. We benchmark GPT-4o and Gemini 1.5 Pro across 10 datasets spanning multiple domains (natural imagery, medical imagery, remote sensing, and molecular imagery) and tasks (multi-class, multi-label, and fine-grained classification). We observe that many-shot ICL, including up to almost 2,000 multimodal demonstrating examples, leads to substantial improvements compared to few-shot (<100 examples) ICL across all of the datasets. Further, Gemini 1.5 Pro performance continues to improve log-linearly up to the maximum number of tested examples on many datasets. Given the high inference costs associated with the long prompts required for many-shot ICL, we also explore the impact of batching multiple queries in a single API call. We show that batching up to 50 queries can lead to performance improvements under zero-shot and many-shot ICL, with substantial gains in the zero-shot setting on multiple datasets, while drastically reducing per-query cost and latency. Finally, we measure ICL data efficiency of the models, or the rate at which the models learn from more demonstrating examples. We find that while GPT-4o and Gemini 1.5 Pro achieve similar zero-shot performance across the datasets, Gemini 1.5 Pro exhibits higher ICL data efficiency than GPT-4o on most datasets. Our

results suggest that many-shot ICL could enable users to efficiently adapt multimodal foundation models to new applications and domains.

1. Introduction

Large language models (LLMs) have been shown to substantially benefit from the inclusion of a few demonstrating examples (*shots*) in the LLM context before the test query (Brown et al., 2020; Parnami & Lee, 2022; Wang et al., 2020). This phenomenon, commonly referred to as in-context learning (ICL), enables LLMs to learn from few shots without any updates to model parameters, and therefore improves specialization to new tasks without any further model training. More recently, large multimodal models (LMMs) have also demonstrated the capability of learning from in-context examples (Achiam et al., 2023; Han et al., 2023; Zhang et al., 2024). Han et al. (2023) and Zhang et al. (2024) both show that few-shot multimodal ICL specifically helps to improve LMM performance on out-domain or out-of-distribution tasks.

While few-shot ICL has enabled promising performance improvements for both LLMs and LMMs, limited model context windows have constrained research on the impact of increasing the number of demonstrating examples on performance. This is especially true for LMMs as most use a large number of visual tokens to represent images. However, due to recent advancements enabling substantially longer context windows – for example, 128,000 tokens for GPT-4o and up to one million tokens for Gemini 1.5 Pro – it is now possible to explore the effect of drastically increasing the number of demonstrating examples.

To investigate the capability of state-of-the-art multimodal foundation models to perform many-shot ICL, we conduct a large suite of experiments benchmarking model performance on 10 datasets spanning several domains and image classification tasks after scaling up the number of demonstrating examples by multiple orders of magnitude. Specifically, our contributions are as follows:

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the 1st In-context Learning Workshop at the International Conference on Machine Learning (ICML). Do not distribute.

055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107
108
109

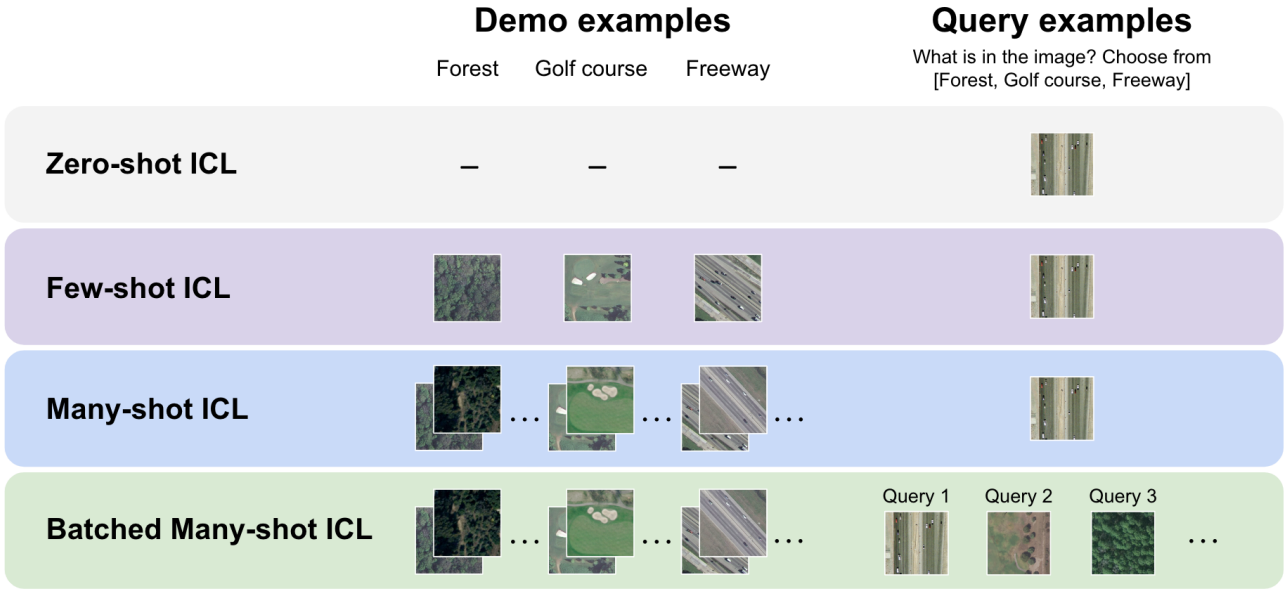


Figure 1. Many-shot multimodal in-context learning compared to zero-shot and few-shot multimodal ICL. In zero-shot and few-shot settings, respectively, no demonstrating examples or only a small number of demonstrating examples are provided in the context before the test query. In a many-shot ICL setting, we include a large number of demonstrating examples in the prompt, whereas in batched many-shot ICL, we perform multiple queries at once using query references.

1. We show that providing multimodal foundation models with many demonstrating examples leads to substantial performance improvements compared to providing only a few demonstrating examples. We observe that the performance of Gemini 1.5 Pro generally improves log-linearly as the number of demonstrating examples increases, whereas GPT-4o exhibits less stable improvements as the number of in-context examples increases.
2. We measure the data efficiency of the models under ICL as the number of demonstrating examples is increased, and find that Gemini 1.5 Pro exhibits higher ICL data efficiency than GPT-4o on most datasets.
3. We demonstrate that batching multiple queries into a single request can achieve similar or better performance than single query requests in a many-shot setting, while enabling substantially lower per-example latency and much cheaper per-example inference cost.
4. We find that batching multiple questions can lead to substantial performance improvements in a zero-shot setting. We design experiments to explain this phenomenon, and find that the improvements are due to a combination of domain calibration, class calibration, and self-generated demonstrating examples due to autoregressive decoding.

2. Related Work

Scaling ICL. The seminal work of Brown et al. (2020) discovered performance improvements for LLMs from increasing the number of in-context examples, but the tested number of demonstrating examples was low (10 to 100), likely due to the restrictive context size (2048 tokens for GPT3). Increasing the number of in-context examples has only been explored recently by a few works (Li et al., 2023; Agarwal et al., 2024; Bertsch et al., 2024). Both Li et al. (2023) and Agarwal et al. (2024) explore scaling in-context learning to more than 1,000 demonstrating examples and find performance improvements across multiple tasks. However, their experiments are limited to text-only benchmarks and do not compare performance across different models.

Multimodal ICL. Due to the recent emergence of LMMs, research on multimodal ICL is still nascent. One prior work developed a new model to leverage complex prompts composed of multimodal inputs in order to allow models to compare images (Zhao et al., 2023), while other recent works explored the generalizability of GPT-4V and Gemini to multimodal out-domain and out-of-distribution tasks, and found that ICL leads to performance benefits for both models across many tasks (Zhang et al., 2024; Han et al., 2023). However, none of these works have leveraged the new largely expanded context windows to investigate the effects of increasing the number of demonstrating examples.

110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164

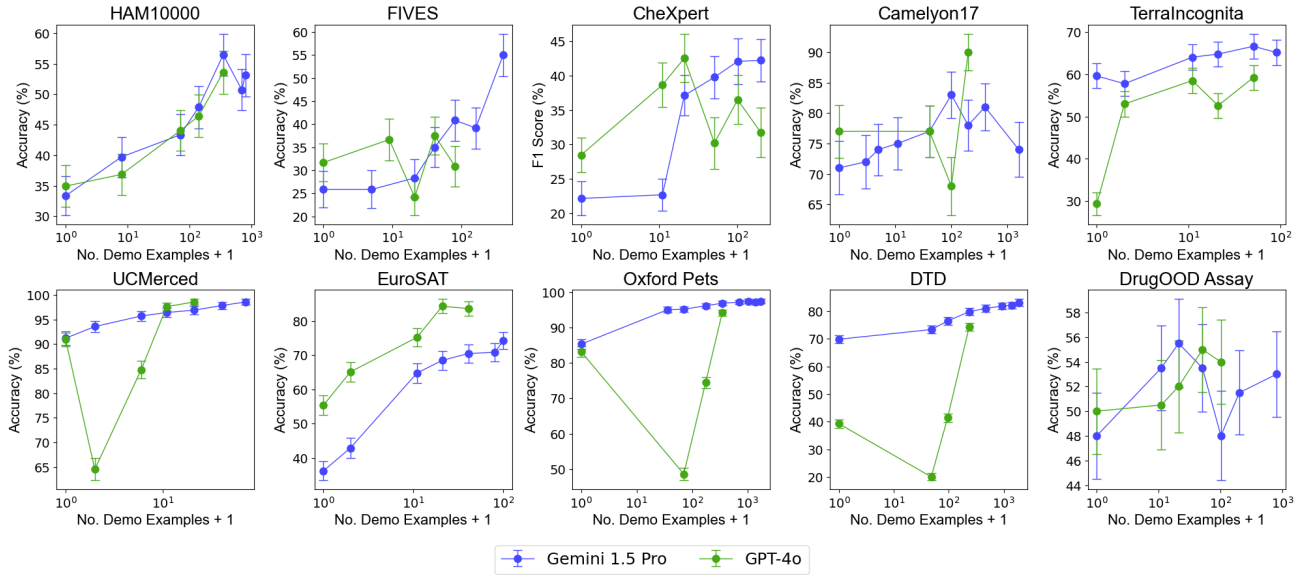


Figure 2. Gemini 1.5 Pro and GPT-4o performance from zero-shot to many-shot ICL. X-axis is in log scale. For Gemini 1.5 Pro, we observe log-linear improvement on 9 out of the 10 datasets and for GPT-4o we observe improvement from more demonstrating examples on most datasets, albeit substantially less stable than Gemini 1.5 Pro.

Batch Querying. Multiple prior works have explored batching queries (also commonly referred to as batch prompting) for more efficient and cheaper inference. Batch prompting was first introduced in Cheng et al. (2023), leading to comparable or better performance than single prompting, while achieving substantially reduced inference token cost and latency. Lin et al. (2023) observe performance degradation with batched prompts in longer contexts, and propose a variety of techniques to mitigate the performance loss. More recently, additional variations of batch prompting have been proposed, including grouping similar questions together (Liu et al., 2024), batching prompts of different tasks (Son et al., 2024), and concatenating multiple images into a single image collage (Xu et al., 2024). We again note that batch prompting with high numbers of demonstrating examples and high numbers of queries has only become feasible due to larger context windows of recent models.

3. Methods

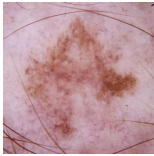


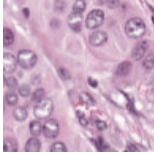




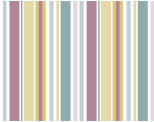
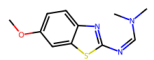
We conduct several experiments to test the effect of increasing the number of demonstrating examples on the performance of two state-of-the-art multimodal foundation models: GPT-4o and Gemini 1.5 Pro (Section 3.1). We benchmark their performance using standard performance metrics as well as an ICL data efficiency metric (Section 3.3) on 10 datasets spanning several vision domains and image classification tasks (Section 3.2). We conduct ablation studies to test the impact of batching queries on model performance

and explain the substantial improvement in zero-shot settings (Section 4.2). We refer to the many-shot in-context learning framework as many-shot ICL. Figure 1 provides an illustrative summary of many-shot ICL and batched many-shot ICL compared to zero-shot and few-shot ICL.

3.1. Models

We use three state-of-the-art multimodal foundation models with public API access, namely GPT-4o, GPT4(V)-Turbo (Achiam et al., 2023), and Gemini 1.5 Pro (Reid et al., 2024). Because GPT-4o performs substantially better than GPT4(V)-Turbo, we focus on the results of GPT-4o and Gemini 1.5 Pro in the main text, and include GPT4(V)-Turbo results in the Appendix. We do not utilize Claude3-Opus in our experiments, as it only accepts up to 20 images in one request at the time of writing. The specific endpoint for GPT-4o is “gpt-4o-2024-05-13”, for GPT-4(V)-Turbo is “gpt-4-turbo-2024-04-09”, and for Gemini 1.5 Pro is “gemini-1.5-pro-preview-0409”. We use the API service provided by OpenAI for GPT-4o and GPT-4(V)-Turbo, and the API service provided by Google Cloud on Vertex AI for Gemini 1.5 Pro. We set the temperature to zero for all models and a random seed for GPT-4(V)-Turbo and GPT-4o to obtain more deterministic responses. To prevent models from abstaining (which happens rarely), we rerun the query until an answer is provided.

Table 1. **Summary of benchmark datasets.** We use 10 datasets spanning multiple domains (natural imagery, medical imagery, remote sensing, molecular imagery) and tasks (multi-class, multi-label, and fine-grained classification).

Dataset	Task and image type	# Classes	Demo / test set size	Example image
HAM10000(Tschandl et al., 2018)	Skin disease classification on clinical photos	7	805 / 210	
FIVES (Jin et al., 2022)	Eye disease classification on fundus images	4	400 / 120	
CheXpert (Irvin et al., 2019)	Multi-label lung disease detection on chest X-rays	5	200 / 150	
Camelyon17 (Bandi et al., 2018)	Tumor detection on pathology images	2	2000 / 100	
TerraIncognita (Beery et al., 2018)	Animal species recognition on camera images	9	1035 / 270	
UCMerced(Yang & Newsam, 2010)	Land use classification on satellite images	21	1470 / 420	
EuroSAT (Helber et al., 2019)	Land use / land cover classification on satellite images	10	1000 / 300	
Oxford Pets (Parkhi et al., 2012)	Pet classification on camera images	35	1750 / 700	
DTD (Cimpoi et al., 2014)	Texture classification on synthetic images	47	2350 / 940	
DrugOOD Assay (Ji et al., 2022)	Drug binding prediction on molecular images	2	1600 / 200	

3.2. Datasets

We benchmark the model performance on 10 datasets spanning multiple domains (natural imagery, medical imagery,

Table 2. **Many-shot ICL performance and efficiency comparison.** We report the performance under a zero-shot regime and performance at the optimal demo set size as well as the many-shot ICL data efficiency of GPT-4o and Gemini 1.5 Pro. We measure performance using accuracy on all datasets except CheXpert, for which we use macro-average F1. We bold the highest ICL data efficiency between the two models on each dataset.

Dataset	GPT-4o			Gemini 1.5 Pro		
	Zero-shot	Best	Efficiency	Zero-shot	Best	Efficiency
HAM10000	34.93	53.59 (+18.66)	5.91	33.33	56.46 (+23.13)	6.94
FIVES	31.67	37.50 (+5.83)	0.30	25.83	55.00 (+29.17)	7.56
CheXpert	28.47	42.54 (+14.08)	3.70	22.16	42.23 (+20.08)	9.06
Camelyon17	77.00	90.00 (+13.00)	1.00	71.00	83.00 (+12.00)	3.00
TerraIncognita	29.26	59.26 (+30.00)	20.50	59.63	66.67 (+7.04)	3.50
UCMerced	90.95	98.57 (+7.62)	1.20	91.19	98.57 (+7.38)	4.36
EuroSAT	55.37	84.23 (+28.86)	19.40	36.24	74.16 (+37.92)	20.61
Oxford Pets	83.14	94.14 (+11.00)	-3.72	85.29	97.43 (+12.14)	4.26
DTD	39.26	74.47 (+35.21)	4.48	69.89	83.19 (+13.30)	3.89
DrugOOD Assay	50.00	55.00 (+5.00)	2.02	48.00	55.50 (+7.50)	2.03

remote sensing, and molecular imagery) and tasks (multi-class, multi-label, and fine-grained classification). We choose to focus on image classification tasks as other tasks such as region captioning would require substantially more tokens thereby limiting the total number of demonstrating examples, and most LMMs are not yet capable of accurately producing localizations required for other tasks like bounding boxes and segmentation masks (Wu et al., 2024; Zang et al., 2023). Table 1 provides a summary of the datasets used in this study.

For all datasets, we construct a set of demonstration (demo) examples from the original training and validation splits used for in-context learning and a test set from the original test split (if one exists) to evaluate the performance of the models. We randomly sample the demo and test sets from the original dataset without replacement. For the multi-class and fine-grained classification datasets, we perform a class-stratified sampling, ensuring an equal number of examples per class in both the demo and test sets. For the multi-label classification dataset (CheXpert), we sample an equal number of positive and negative samples per class in both the demo and test sets. We note that, since the task is multi-label, this sampling procedure does not result in an exactly equal number of examples per class. The per-dataset sizes of the full demo and test sets are shown in Table 1, and we increase the number of demonstration examples up to the numbers shown in the table while ensuring class balance for the scaling experiments.

3.3. Evaluation Metrics

We use standard metrics to evaluate model performance on each dataset. Specifically, we measure performance using accuracy for all multi-class classification datasets as they

are sampled to have a balanced class distribution. For multi-label classification on CheXpert, we use the macro-averaged F1 metric. In the rare case of parsing errors, we consider the response as incorrect. To estimate the variability around the evaluation metrics, we compute standard deviation using bootstrapping with 1,000 bootstrap replicates.

In addition to standard performance metrics, we measure the data efficiency of each model. Specifically, we compute a linear regression between $\log_{10}(N+1)$ (with N the number of examples) and model performance, enforcing that the line passes through the zero-shot performance point. This value approximates the amount of performance improvement from zero-shot expected from including an order of magnitude more demonstrating examples.

4. Results

We present many-shot ICL performance using batched queries in Section 4.1, investigate the impact of batching queries on performance in Section 4.2, and provide an analysis on cost and latency in Section 4.3. Results using GPT4(V)-Turbo are in Appendix C.

4.1. Increasing number of demonstrating examples

Main Results. Gemini 1.5 Pro exhibits consistent and substantial improvements as the number of demonstrating examples increases across all datasets except for DrugOOD Assay (Figure 2). Gemini 1.5 Pro shows particularly large improvements from many-shot ICL on HAM10000 (+23% accuracy compared to zero-shot, +16% compared to 7 examples), FIVES (+29% compared to zero-shot, +27% compared to 20 examples), and EuroSAT (+38% compared to zero-shot, +31% compared to 10 examples). Notably, for 5 out of the

275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329

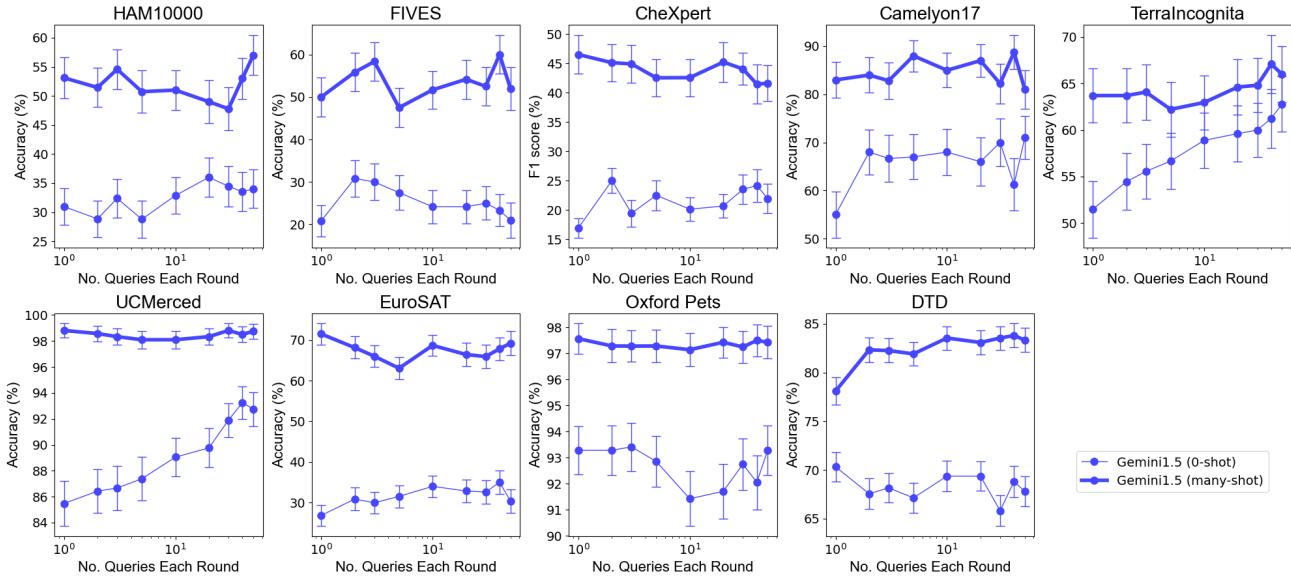


Figure 3. Gemini 1.5 Pro performance under many-shot and zero-shot ICL when varying the amount of queries included in every request. We show performance per batch size with the optimal number of demo examples (many-shot) and no demo examples (zero-shot). The x -axis is in log scale. Under the many-shot regime, batching queries leads to no substantial drop in performance compared to individual queries when we choose a suitable batch size. For zero-shot, including only one query is suboptimal for many datasets.

10 datasets (FIVES, UCMerced, EuroSAT, Oxford Pets, and DTD), Gemini 1.5 Pro performance continues to improve up to the highest number of demonstrating examples considered (~1,000 examples). On the other 5 datasets, the optimal performance occurs prior to the highest number of demo examples, with the maximum number of demo examples leading to similar or slightly worse performance than the optimal demo set size. On the other hand, Gemini 1.5 Pro performance on DrugOOD Assay does not substantially benefit from many-shot ICL, with high variance in performance across demo sizes and the peak performance at 40 demo examples.

Similarly, GPT-4o shows substantial performance improvements on all datasets except FIVES and DrugOOD Assay using many-shot ICL, but the improvement is not consistent. For many datasets, performance drops sharply at first and then improves significantly as the number of demonstrating examples increases further, resulting in V-shaped scaling curves (Figure 2). We also note that we were unable to increase the number of demo examples to the same level as considered for Gemini 1.5 Pro because GPT-4o has a shorter context window and is more prone to timeout errors with longer inputs. GPT-4o performance on DrugOOD Assay shows high variance, similar to Gemini 1.5 Pro, with the peak performance observed at 50 demo examples.

Sensitivity to prompt selection. We also explore a different set of prompts to test the robustness of many-shot ICL to dif-

ferences in prompt wording on two datasets. While there is a small deviation in performance between different prompts, the overall log-linear improvement trend is consistent across the prompts. Details can be found in Appendix B.

ICL data efficiency. We find Gemini 1.5 Pro demonstrates higher ICL data efficiency than GPT-4o across all datasets except TerraIncognita and DTD (Table 2). Gemini 1.5 Pro ICL efficiency is especially high on EuroSAT, with 20.61% improvement in accuracy for every 10x more demo examples, and lowest on DrugOOD Assay (2.03), Camelyon17 (3.00), and TerraIncognita (3.50). GPT-4o ICL data efficiency is especially high on TerraIncognita (20.50%) and EuroSat (19.40). Gemini 1.5 Pro has a positive efficiency on all datasets and GPT-4o has a positive data efficiency on 9 of the 10 datasets (excluding Oxford Pets). Importantly, both models benefit substantially from many-shot ICL at the optimal demo set size, with an average improvement of +17% for both Gemini 1.5 Pro and GPT-4o.

4.2. Impact of batching queries

As including a large set of demo examples in the prompt leads to much longer sequence lengths and therefore higher inference time and cost, we consider batching queries in a single prompt to reduce per-query cost, and examine the impact of different batch sizes on model performance. Due to its superior performance and free preview access, we use Gemini 1.5 Pro for these experiments.

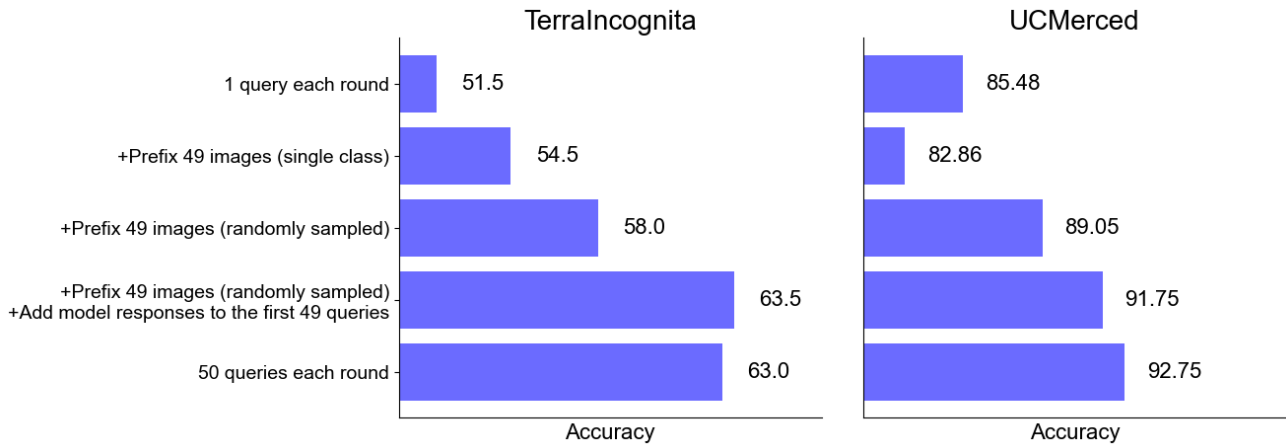


Figure 4. Ablation study to investigate why batching queries leads to performance improvements when using Gemini 1.5 Pro in a zero-shot setting. The first bar shows performance when including a single query, the second adds 49 unlabeled images from a single class, the third adds 49 unlabeled images in total from all classes, the fourth adds model responses to include self-generated demonstrations, and the last includes 50 queries in one request.

Main Results. We find minimal performance degradations, and sometimes performance improvements, as we increase the number of queries included in each batch across under both zero-shot and many-shot (at the optimal demo set size) regimes (Figure 3). Notably, using a single query each time with many-shot ICL is suboptimal across many of the datasets. We find that the optimal batch size is among the three largest sizes on every dataset except CheXpert and EuroSAT, which both see optimal performance with a single query at a time.

We additionally observe that including a single query at a time is suboptimal on most datasets in the zero-shot regime. Surprisingly, performance with the highest batch size is substantially higher across three datasets under the zero-shot regime, with a consistent performance improvement as the batch size is increased on both UCMerced and TerraIncognita.

Zero-shot performance improvements from batching queries. We conduct several additional experiments to investigate why batch querying can lead to large performance improvements under the zero-shot regime on TerraIncognita and UCMerced. We hypothesize that this improvement may be due to three potential benefits from ICL: (1) domain calibration, where the model benefits from seeing more images in the domain in order to adapt to it, (2) class calibration, where seeing images of different classes enables the model to better calibrate its outputs, and (3) self-ICL (shown to be effective in prior work (Chen et al., 2023)), where the model can learn from self-generated demonstrations due to autoregressive decoding. We design experiments to isolate the potential benefits from each of these types of ICL between

asking a single query to batching 50 queries together.

First, to measure potential improvement from domain calibration, we include 49 images from the same class in the prompt without including any label. We find a 3.0% improvement on TerraIncognita and 2.6% degradation on UCMerced, suggesting domain calibration is helpful for the former but not the latter. Second, to capture performance gains from class calibration, we include a random sample of 49 images in the prompt, again without including the label. We see a further 3.5% improvement on TerraIncognita (6.5% improvement from a single query) and a 4.5% improvement from a single query on UCMerced, suggesting including the context of class-balanced images is helpful even without labels. Third, to capture additional performance improvements from the self-generated labels, we obtain predicted labels from the zero-shot model using a single query for each of the 49 randomly sampled images and add them to the prompt. We observe further performance increase on both datasets, with 5.5% on TerraIncognita and 2.7% on UCMerced. The final total accuracy is similar to asking the 50 questions each round, which suggests these three components mostly explain the reason for improved zero-shot performance under a larger query batch size.

4.3. Cost and latency analysis

Many-shot ICL incurs zero additional training cost, but per-query inference can be costly and slow due to long input contexts. To quantitatively measure this, we compute the latency and cost associated with the zero-shot and many-shot requests with and without batching when using Gemini 1.5 Pro on HAM10000 and TerraIncognita. We calculate

Table 3. **Inference latency and cost using Gemini 1.5 Pro with and without query batching.** We use 50 queries per batch. In the zero-shot setting, we can achieve lower per-example latency with batching, but the per-example cost remains identical. In the many-shot setting, the per-example cost and per-example latency both drop substantially with query batching.

Dataset	No Query Batching			Query Batching		
	Per-batch latency	Per-example latency	Per-example cost	Per-batch latency	Per-example latency	Per-example cost
HAM10000 (zero-shot)	2.2s	2.2s	\$0.0038	11.4s	0.23s	\$0.0038
TerraIncognita (zero-shot)	2.0s	2.0s	\$0.0037	51.6s	1.0s	\$0.0038
HAM10000 (350-shot)	17.3s	17.3s	\$0.8420	26.9s	0.54s	\$0.0877
TerraIncognita (810-shot)	34.9s	34.9s	\$1.8420	85.9s	1.7s	\$0.0406

the costs using the Gemini 1.5 Pro preview pricing (\$7 per 1 million input tokens and \$21 per 1 million output tokens). We run the query three times under each setting and report the average.

In the zero-shot regime, we see substantial per-example latency reductions due to query batching, close to a 10x reduction on HAM10000 and 2x on TerraIncognita (Table 3). The per-example cost is similar between the two as there is no additional context needed for including demonstrating examples. In the many-shot regime, we observe substantial reductions in both per-example latency and cost. Specifically, for HAM10000, we find a near 35x reduction in latency and 10x reduction in cost, and 20x reduction in latency and 45x reduction in cost for TerraIncognita.

5. Discussion

In this study, we evaluate many-shot ICL of state-of-the-art multimodal foundation models across 10 datasets and find consistent performance improvements across most of the datasets. Batching queries with many-shot ICL further exhibits substantially reduced per-example latency and inference costs without compromising performance.

Our findings suggest that these multimodal foundation models have the capability of performing ICL with large numbers of demonstrating examples, which may have significant implications on their practical use. For example, it was previously impossible to adapt these large, private models to new tasks and domains, but many-shot ICL would enable users to leverage demonstrating examples to adapt the models. One significant advantage of many-shot ICL is its ability to get quick results even on the same day of model release, and that’s why we can finish our evaluation using GPT-4o within days. Furthermore, fine-tuning open-source models is the standard practice when practitioners have access to moderately sized datasets, but many-shot ICL may remove the need for fine-tuning, making it much easier to develop

customized approaches. We note that it remains to be seen how traditional fine-tuning of these models compares to many-shot ICL with foundation models in terms of absolute performance and data efficiency, so future work should explore this. In addition, it is important to study general issues which plague those foundation models, such as hallucinations and biases, under the context of many-shot ICL and batching queries. For example, it would be interesting to explore if carefully curated and large sets of demonstrating examples can reduce biases across different sub-groups. We leave this to future work.

Our study has limitations. First, we only explore performance under many-shot ICL on image classification tasks and with private foundation models. We believe these are the most practically relevant and common multimodal settings, but it is worthwhile for future work to explore potential benefits from many-shot ICL on other tasks and with upcoming open-source multimodal foundation models like LLaMA-3 (11a). Second, even after recent developments to increase context size, the size prohibits many-shot ICL from being used on datasets with a large number (several hundred or more) of classes. We anticipate that context window sizes will continue to increase in size over time which will mitigate this issue. Third, the datasets which were used to train these private models have not been disclosed, so it is difficult to tell whether the models have been trained on the datasets we selected. We argue that zero-shot performance across the datasets is far from perfect which provides evidence that the datasets have not been used for training, but we cannot determine that with certainty.

6. Conclusion

In summary, we show that multimodal foundation models are capable of many-shot ICL. We believe that these results pave a promising path forward to improve the adaptability and accessibility of large multimodal foundation models.

References

440 Introducing meta llama 3: The most capable openly avail-
441 able llm to date. URL [https://ai.meta.com/
442 blog/meta-llama-3/](https://ai.meta.com/blog/meta-llama-3/).
443
444

445 Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I.,
446 Aleman, F. L., Almeida, D., Altenschmidt, J., Altman, S.,
447 Anadkat, S., et al. Gpt-4 technical report. *arXiv preprint
448 arXiv:2303.08774*, 2023.
449

450 Agarwal, R., Singh, A., Zhang, L. M., Bohnet, B., Chan, S.,
451 Anand, A., Abbas, Z., Nova, A., Co-Reyes, J. D., Chu,
452 E., et al. Many-shot in-context learning. *arXiv preprint
453 arXiv:2404.11018*, 2024.
454

455 Bandi, P., Geessink, O., Manson, Q., Van Dijk, M., Balken-
456 hol, M., Hermsen, M., Bejnordi, B. E., Lee, B., Paeng, K.,
457 Zhong, A., et al. From detection of individual metastases
458 to classification of lymph node status at the patient level:
459 the camelyon17 challenge. *IEEE transactions on medical
460 imaging*, 38(2):550–560, 2018.

461 Beery, S., Van Horn, G., and Perona, P. Recognition in terra
462 incognita. In *Proceedings of the European conference on
463 computer vision (ECCV)*, pp. 456–473, 2018.
464

465 Bertsch, A., Ivgi, M., Alon, U., Berant, J., Gormley,
466 M. R., and Neubig, G. In-context learning with long-
467 context models: An in-depth exploration. *arXiv preprint
468 arXiv:2405.00200*, 2024.
469

470 Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan,
471 J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G.,
472 Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G.,
473 Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu,
474 J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M.,
475 Gray, S., Chess, B., Clark, J., Berner, C., McCandlish,
476 S., Radford, A., Sutskever, I., and Amodei, D. Language
477 models are few-shot learners, 2020.

478 Chen, W.-L., Wu, C.-K., and Chen, H.-H. Self-icl: Zero-shot
479 in-context learning with self-generated demonstrations.
480 *arXiv preprint arXiv:2305.15035*, 2023.
481

482 Cheng, Z., Kasai, J., and Yu, T. Batch prompting: Efficient
483 inference with large language model apis. *arXiv preprint
484 arXiv:2301.08721*, 2023.
485

486 Cimpoi, M., Maji, S., Kokkinos, I., Mohamed, S., and
487 Vedaldi, A. Describing textures in the wild. In *Pro-
488 ceedings of the IEEE conference on computer vision and
489 pattern recognition*, pp. 3606–3613, 2014.

490 Han, Z., Zhou, G., He, R., Wang, J., Xie, X., Wu, T., Yin,
491 Y., Khan, S., Yao, L., Liu, T., et al. How well does
492 gpt-4v (ision) adapt to distribution shifts? a preliminary
493 investigation. *arXiv preprint arXiv:2312.07424*, 2023.
494

Helber, P., Bischke, B., Dengel, A., and Borth, D. Eurosat:
A novel dataset and deep learning benchmark for land
use and land cover classification. *IEEE Journal of Se-
lected Topics in Applied Earth Observations and Remote
Sensing*, 12(7):2217–2226, 2019.

Irvin, J., Rajpurkar, P., Ko, M., Yu, Y., Ciurea-Ilcus, S.,
Chute, C., Marklund, H., Haghgoo, B., Ball, R., Shpan-
skaya, K., et al. Chexpert: A large chest radiograph
dataset with uncertainty labels and expert comparison. In
*Proceedings of the AAAI conference on artificial intelli-
gence*, volume 33, pp. 590–597, 2019.

Ji, Y., Zhang, L., Wu, J., Wu, B., Huang, L.-K., Xu, T.,
Rong, Y., Li, L., Ren, J., Xue, D., et al. Drugood:
Out-of-distribution (ood) dataset curator and benchmark
for ai-aided drug discovery—a focus on affinity predic-
tion problems with noise annotations. *arXiv preprint
arXiv:2201.09637*, 2022.

Jin, K., Huang, X., Zhou, J., Li, Y., Yan, Y., Sun, Y., Zhang,
Q., Wang, Y., and Ye, J. Fives: A fundus image dataset
for artificial intelligence based vessel segmentation. *Sci-
entific Data*, 9(1):475, 2022.

Li, M., Gong, S., Feng, J., Xu, Y., Zhang, J., Wu, Z., and
Kong, L. In-context learning with many demonstration
examples. *arXiv preprint arXiv:2302.04931*, 2023.

Lin, J., Diesendruck, M., Du, L., and Abraham, R. Batch-
prompt: Accomplish more with less. *arXiv preprint
arXiv:2309.00384*, 2023.

Liu, J., Yang, T., and Neville, J. Cliqueparcel: An approach
for batching llm prompts that jointly optimizes efficiency
and faithfulness. *arXiv preprint arXiv:2402.14833*, 2024.

Parkhi, O. M., Vedaldi, A., Zisserman, A., and Jawahar, C.
Cats and dogs. In *2012 IEEE conference on computer
vision and pattern recognition*, pp. 3498–3505. IEEE,
2012.

Parnami, A. and Lee, M. Learning from few examples:
A summary of approaches to few-shot learning. *arXiv
preprint arXiv:2203.04291*, 2022.

Reid, M., Savinov, N., Teplyashin, D., Lepikhin, D., Lilli-
crap, T., Alayrac, J.-b., Soricut, R., Lazaridou, A., Firat,
O., Schrittwieser, J., et al. Gemini 1.5: Unlocking multi-
modal understanding across millions of tokens of context.
arXiv preprint arXiv:2403.05530, 2024.

Son, G., Baek, S., Nam, S., Jeong, I., and Kim, S. Multi-task
inference: Can large language models follow multiple
instructions at once? *arXiv preprint arXiv:2402.11597*,
2024.

495 Tschandl, P., Rosendahl, C., and Kittler, H. The ham10000
496 dataset, a large collection of multi-source dermatoscopic
497 images of common pigmented skin lesions. *Scientific*
498 *data*, 5(1):1–9, 2018.

499 Wang, Y., Yao, Q., Kwok, J. T., and Ni, L. M. Generalizing
500 from a few examples: A survey on few-shot learning.
501 *ACM computing surveys (csur)*, 53(3):1–34, 2020.

503 Wu, Y., Wang, Y., Tang, S., Wu, W., He, T., Ouyang, W.,
504 Wu, J., and Torr, P. Dettoolchain: A new prompting
505 paradigm to unleash detection ability of mllm. *arXiv*
506 *preprint arXiv:2403.12488*, 2024.

508 Xu, S., Wang, Y., Liu, D., and Xu, C. Collage prompting:
509 Budget-friendly visual recognition with gpt-4v. *arXiv*
510 *preprint arXiv:2403.11468*, 2024.

511 Yang, Y. and Newsam, S. Bag-of-visual-words and spatial
512 extensions for land-use classification. In *Proceedings*
513 *of the 18th SIGSPATIAL international conference on ad-*
514 *vances in geographic information systems*, pp. 270–279,
515 2010.

517 Zang, Y., Li, W., Han, J., Zhou, K., and Loy, C. C. Con-
518 textual object detection with multimodal large language
519 models. *arXiv preprint arXiv:2305.18279*, 2023.

521 Zhang, X., Li, J., Chu, W., Hai, J., Xu, R., Yang, Y., Guan,
522 S., Xu, J., and Cui, P. On the out-of-distribution gener-
523 alization of multimodal large language models. *arXiv*
524 *preprint arXiv:2402.06599*, 2024.

525 Zhao, H., Cai, Z., Si, S., Ma, X., An, K., Chen, L., Liu, Z.,
526 Wang, S., Han, W., and Chang, B. Mmicl: Empower-
527 ing vision-language model with multi-modal in-context
528 learning. *arXiv preprint arXiv:2309.07915*, 2023.

530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549

550 A. Prompts used for ICL experiments

551 A.1. Prompt used for image classification experiments

```
552 prompt = ""
553 for demo in demo_examples:
554     prompt += f"""<IMG>>Given the image above, answer the following question-
555 using the specified format.
556 Question: What is in the image above?
557 Choices: {str(class_desp)}
558 Answer Choice: {demo.answer}
559 """
560
561 prompt += f"""<IMG>>Given the image above, answer the following question-
562 using the specified format.
563 Question: What is in the image above?
564 Choices: {str(class_desp)}
565
566 Please respond with the following format:
567 --BEGIN FORMAT TEMPLATE---
568 Answer Choice: [Your Answer Choice Here]
569 Confidence Score: [Your Numerical Prediction Confidence Score Here From 0 To 1]
570 --END FORMAT TEMPLATE---
571
572 Do not deviate from the above format. Repeat the format template for the answer."""
```

575 A.2. Prompts used for image classification experiments with batching

```
576 prompt = ""
577 for demo in demo_examples:
578     prompt += f"""<IMG>>Given the image above, answer the following question-
579 using the specified format.
580 Question: What is in the image above?
581 Choices: {str(class_desp)}
582 Answer Choice: {demo[1]}
583 """
584
585 for idx, i in enumerate(test_df.iloc[start_idx:end_idx].itertuples()):
586     prompt += f"""<IMG>>Given the image above, answer the following question-
587 using the specified format.
588 Question {qn_idx}: What is in the image above?
589 Choices {qn_idx}: {str(class_desp)}
590
591 """
592
593 for i in range(start_idx, end_idx):
594     qn_idx = i-start_idx+1
595     prompt += f"""
596 Please respond with the following format for each question:
597 --BEGIN FORMAT TEMPLATE FOR QUESTION {qn_idx}---
598 Answer Choice {qn_idx}: [Your Answer Choice Here for Question {qn_idx}]
599 Confidence Score {qn_idx}: [Your Numerical Prediction Confidence Score Here-
600 From 0 To 1 for Question {qn_idx}]
601 --END FORMAT TEMPLATE FOR QUESTION {qn_idx}---
602
603
604
```

605 Do not deviate from the above format. Repeat the format template for the answer. """
606

607 **A.3. Prompts used for batching ablation experiments**

608 609 **A.3.1. PREFIXING IMAGES**

```
610 prompt = ""
611 for demo in prefix_image_paths:
612     prompt += f"""<IMG>>
613
614 """
615 prompt += "Above are some images from the same dataset. "
616 qns_idx = []
617 for idx, i in enumerate(test_df.iloc[start_idx:end_idx].itertuples()):
618     qn_idx = idx+1
619     prompt += f"""<IMG>> Given the image above, answer the following question-
620 using the specified format.
621 Question {qn_idx}: What is in the image above?
622 Choices {qn_idx}: {str(class_desp)}
623
624 """
625 for i in range(start_idx, end_idx):
626     qn_idx = i-start_idx+1
627     prompt += f"""
628 Please respond with the following format for each question:
629 --BEGIN FORMAT TEMPLATE FOR QUESTION {qn_idx}---
630 Answer Choice {qn_idx}: [Your Answer Choice Here for Question {qn_idx}]
631 Confidence Score {qn_idx}: [Your Numerical Prediction Confidence Score Here-
632 From 0 To 1 for Question {qn_idx}]
633 --END FORMAT TEMPLATE FOR QUESTION {qn_idx}---
634
635 Do not deviate from the above format. Repeat the format template for the answer. """
```

637 **B. Prompt selection**

638
639 We utilize a different set of prompts to test the robustness of ManyICL to differences in prompt wording. We randomly
640 sample two datasets (HAM10000 and EuroSAT) for this experiment due to budget limit.

642 **B.1. Prompts used for prompt selection experiments**

643
644 Note that only the question section is shown here, and prompt 1 is used for all other image classification experiments.

646 **B.1.1. PROMPT 1**

```
647 <<IMG>>Given the image above, answer the following question using the specified format.
648 Question {qn_idx}: What is in the image above?
649 Choices {qn_idx}: {str(class_desp)}
```

651 **B.1.2. PROMPT 2**

```
652 <<IMG>>Given the image above, answer the following question using the specified format.
653 Question {qn_idx}: Which class does this image belong to?
654 Choices {qn_idx}: {str(class_desp)}
```

656
657
658
659

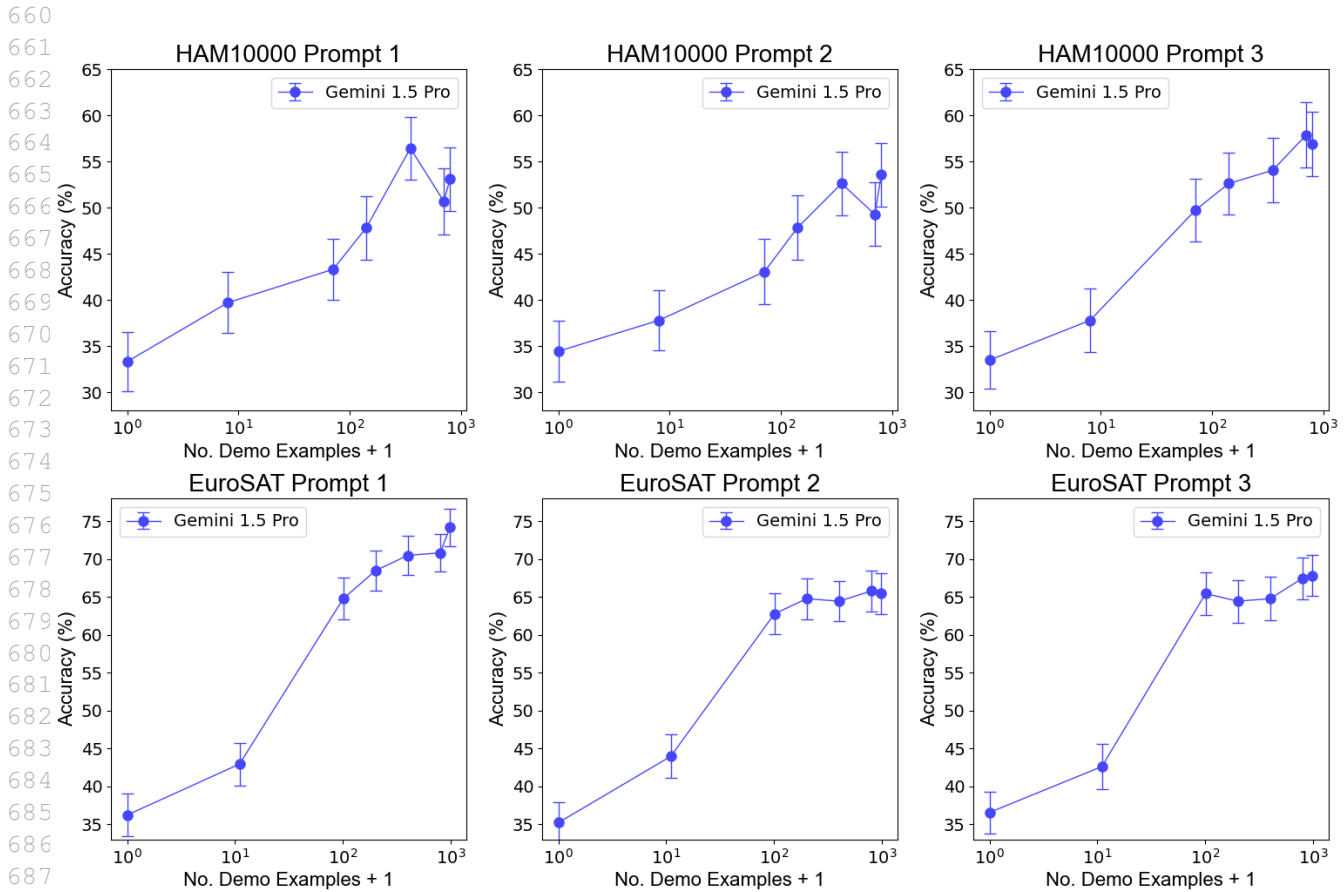


Figure 5. Sensitivity analysis of many-shot ICL. These plots show the change in task performance on two datasets as the number of demonstrating examples increases, using three different prompts. For all experiments on sensitivity analysis, the Gemini 1.5 Pro model is used. The x -axis is in the logarithmic scale, representing the number of demonstrating examples plus one. The log-linear improvement until the optimal performance is consistent across all prompts selected.

B.1.3. PROMPT 3

Question {qn_idx}: <>Classify the image above, choose from {str(class_desp)}

B.2. Prompt selection results

Figure 5 shows the sensitivity of performance to prompt selection on two datasets with three prompts. While there exists a small deviation in performance, but the overall log-linear improvement trend is consistent.

C. GPT4(V)-Turbo performance under many-shot ICL

GPT4(V)-Turbo shows mixed results for many-shot ICL, with substantial performance improvements on HAM1000, UCMerced, EuroSAT, and DTD, but minimal improvements or no improvement across the other six datasets (Figure 6). However, we note that we were unable to increase the number of demo examples to the same level as Gemini 1.5 Pro because GPT4(V)-Turbo has a shorter context window and is more prone to timeout errors when scaling. Additionally, GPT4(V)-Turbo seems to generally underperform Gemini 1.5 Pro across the datasets excluding FIVES and EuroSAT for which it seems to mostly match the Gemini 1.5 Pro performance. GPT4(V)-Turbo performance on DrugOOD Assay shows high variance, resembling that of Gemini 1.5 Pro with the peak performance at 40 demo examples.

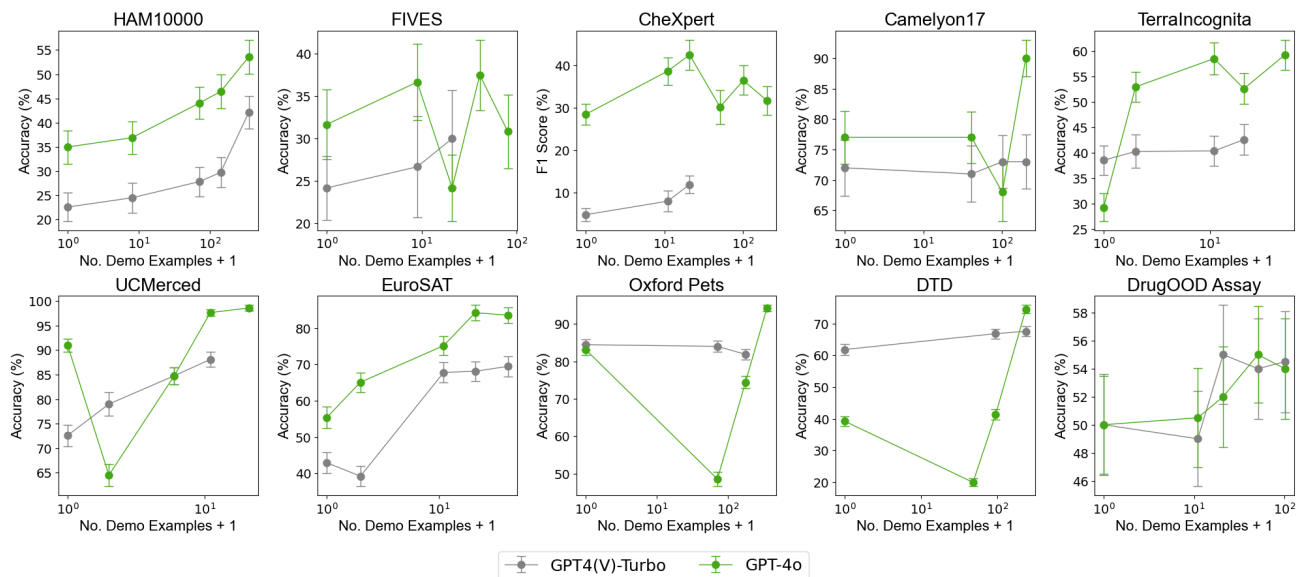


Figure 6. GPT4(V)-Turbo and GPT-4o performance from zero-shot to many-shot ICL. X-axis is in log scale.

D. Performance of many-shot ICL on medical QA tasks

D.1. Prompt used for medical QA experiments (MedQA, MedMCQA)

```
prompt = "You are an expert in answering medical exam questions. "
for demo in demo_examples:
    prompt += f""Question: {demo.question}
    Choices: {demo.options}
    Answer: {demo.answer}
    "" "
```

```
prompt += f""Question: {actual.question}
    Choices: {actual.options}
```

Please respond with the following format:

```
---BEGIN FORMAT TEMPLATE---
```

```
Answer: [Your Answer Choice Here]
```

```
Confidence Score: [Your Numerical Prediction Confidence Score Here From 0 To 1]
```

```
---END FORMAT TEMPLATE---
```

Do not deviate from the above format. Repeat the format template for the answer. "" "

D.2. Results

Figure 7 shows the results on medical QA tasks.

770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824

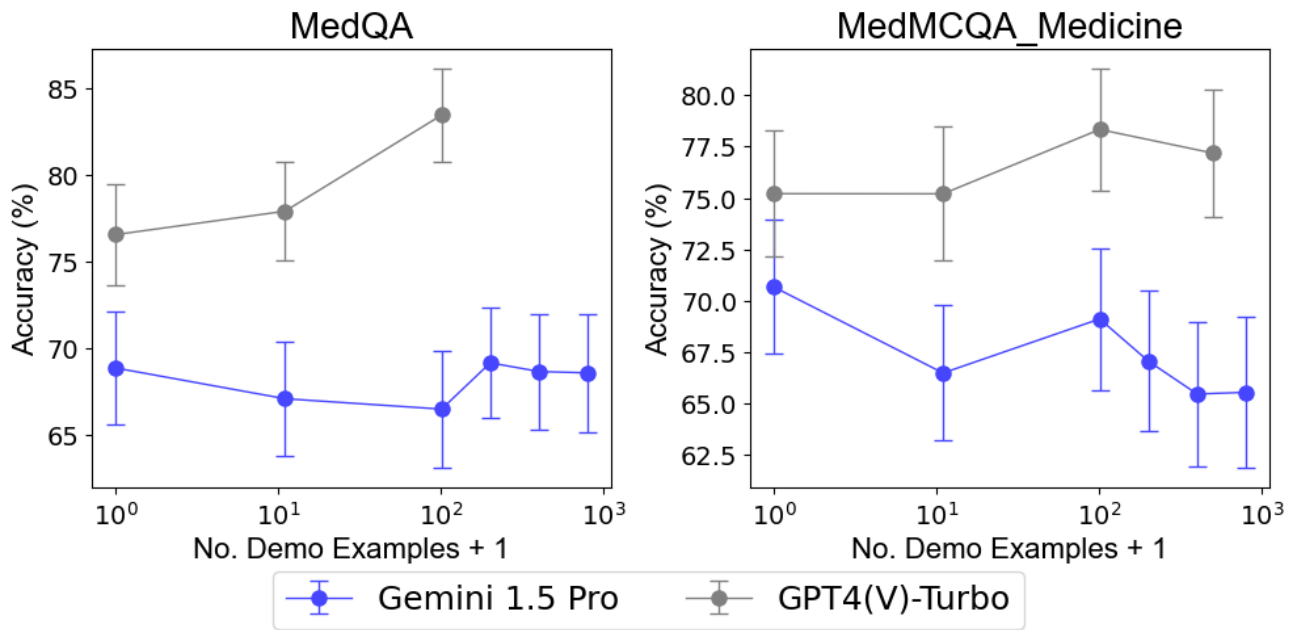


Figure 7. Many-shot ICL performances of medical QA tasks.