

A EXPERIMENT SETUP IN DETAIL

We describe the experimental setup used to evaluate our input-adaptive inference mechanism in detail. We implemented our strategy on top of the codebases provided by the authors of the HAMT (Chen et al., 2021) and DUET (Chen et al., 2022). During inference, instead of using cached image features, we integrate a ViT-B/16 (Dosovitskiy et al., 2021) model to process the images directly.

Hardware and software. We run our experiments on a machine equipped with an Intel Xeon processor with 48 cores, 64GB of DRAM, and 8 NVIDIA A40 GPUs, with all inference tasks performed on a single GPU using a batch size of 1. Following the original HAMT study, we use Python v3.8.5 and PyTorch v1.7.1, along with CUDA v10.1. For GFLOPs calculations, we use the Python library thop.

Datasets. We describe the benchmarking datasets we use in detail:

- **R2R** (Anderson et al., 2018) is based on Matterport3D (Chang et al., 2017), containing 10,567 panorama views taken from 90 photo-realistic houses. The dataset includes 7,189 shortest-path trajectories, and each of them is associated with 3 natural language instructions. The training, validation (seen), validation (unseen), and test (unseen) sets include 61, 56, 11, and 18 houses, respectively. The validation (seen) set consists of houses in the training set, typically used to check the generalization status of a model during training, while the sets marked as ‘unseen’ are the houses not in the training set.
- **R2R-Back** (Chen et al., 2021) requires the agent to return to its starting point after reaching the destination. To complete the task, the agent must remember its navigation history. A return command is appended to each R2R instruction, and the reversed path is provided as guidance for the return trip.
- **R2R-Last** (Chen et al., 2021) uses only the last sentence from the original R2R instructions to describe the destination.
- **REVERIE** (Qi et al., 2020) provides high-level instructions, closer to those given by humans, replacing the step-by-step instructions of R2R. Instead of navigating to a target location, the agent is required to identify and localize the target object upon arrival, making the task more complex and realistic. The dataset includes 4,140 target objects, which are categorized into 489 distinct groups.
- **CVDN** (Thomason et al., 2020) requires the agent to navigate based on long, potentially unclear instructions. The agent interacts with a navigator through question and answer dialog to clarify and complete the task. In total, it has 2,050 human-human navigation dialogues, consisting of over 7,000 navigation trajectories accompanied by question-answer interactions, covering 83 matterport3D houses.
- **SOON** (Zhu et al., 2021) is similar to REVERIE but contains longer and more detailed instructions. The average length of these instructions is 47 words, with path lengths varying from 2 to 21 steps. It requires the agent to navigate by understanding the relationship between objects in the environment to accurately locate the target object.

B OPTIMAL HYPERPARAMETER CHOICE FOR ADAPTING MUE TO OUR WORK

To best evaluate MuE on VLN tasks, we perform a hyperparameter sweep over the threshold used for early-exiting. Figure 7 presents the performance (in SR) and GFLOPs across different early exit thresholds applied to the MuE version of ViT used in the HAMT agent, tested on the R2R dataset. The lowest threshold we report is 0.99, as lower thresholds caused a dramatic drop in performance (more than 50%). As the threshold increases, the success rate of the MuE agent increases substantially but at the cost of computational savings. Even for thresholds close to 1, meaning that the ViT is using a majority of its layers for each input, we still see a large performance drop compared to the baseline agent. As we discuss in Sec 4.2, this is likely because MuE statically applies early-exits, causing it to under-process important components of the panorama such as navigable views.

Why does MuE underprocess important views? The intuition behind MuE (Tang et al., 2023) is that the activations of Transformer-based vision models *saturate*, where their similarity between layers peaks early on, and is maintained at future stages of computation, suggesting a lack of new/useful information. MuE then exploits this property to skip the later layers without a significant loss in performance. So, for MuE to be successful, the similarity of activations must sufficiently saturate and

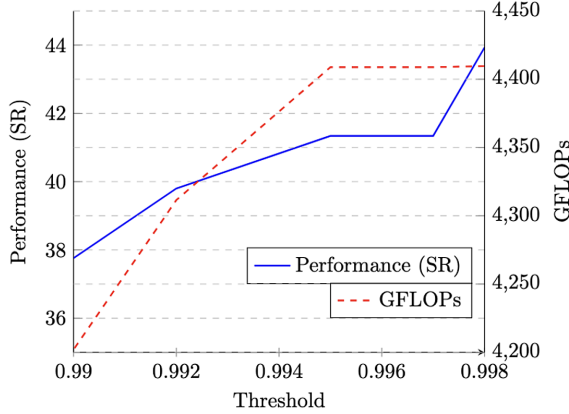


Figure 7: **Comparison of performance (in SR) and GFLOPs in MuE across different thresholds.**

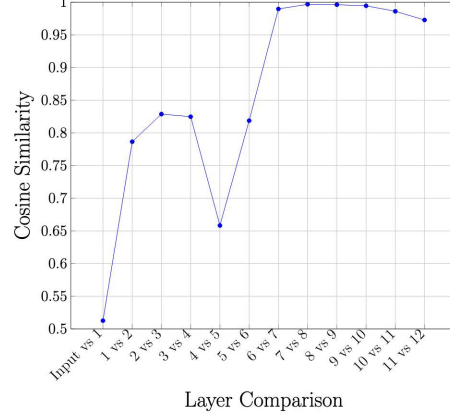


Figure 8: **Cosine similarity between adjacent layers of ViT used in HAMT.**

not decrease at later layers. However, as shown in Figure 8, the necessary saturation pattern is not observed in the VLN setting. The cosine similarity peaks between layers 7 and 8 but then decreases for all future layers. This explains the significant performance drop when MuE is directly applied to VLN agents, as it consistently early-exits despite saturation not being achieved.

C OUR LSH ALGORITHM IN DETAIL

A core mechanism we introduce in Sec 4.2.3 is our SimHash algorithm, used to avoid reprocessing previously seen or near-identical images. Algorithm 2 covers our implementation in detail.

(line 1-9) Hashing RGB vectors. Given an image, we first hash the raw RGB vector into a short binary encoding using random projection Charikar (2002); Andoni & Indyk (2008). The algorithm calculates the dot product between the image vector and each hyperplane. If the dot product is positive, it assigns a binary value of 1, otherwise it assigns 0. These binary values are sequentially appended to form a complete binary hash key. The length of the hash key is determined by the number of hyperplanes used in the projection.

(line 10-14) Adding embeddings to the hash table. This function is used to insert processed images and their corresponding embeddings into the hash table for future use.

(line 15-32) Retrieving a similar embedding. This function takes an image we have not yet processed and tries to find a suitable embedding candidate. We first obtain all embeddings with images similar to the current image by hashing it into its binary encoding and accessing the corresponding bucket in the hash table. We then loop through all images associated with the similar embeddings and find the one yielding the highest similarity score (in our main experiments, the score is computed using cosine similarity). If this score exceeds our threshold hyperparameter, we return the associated embedding; otherwise, we return nothing.

Running the algorithm. We employ the above three functions to run SimHash on an arbitrary panorama. For each extended navigable view (other views are omitted and explained in Algorithm 1), we attempt to use a high-similarity embedding from the hash table. If it exists, we reuse this embedding for the current view and continue to the next. If not, we need to process the view using the ViT adapted for MuE, and then add the image and its embedding to the hash table. After processing the entire panorama, we return the set of final embeddings to be used for agent navigation.

Storage overhead analysis. Here, we consider the storage overhead necessary to deploy our hashing algorithm on VLN agents. Our LSH technique stores pairs of images and embeddings. In the benchmarks we consider, these images are of size $3 \times 224 \times 224$. The embedding size depends on the model, which for HAMT and DUET is 197×768 (the number of ViT patches times the model’s hidden dimension). These are stored in full-precision floating-point format (32 bits per value), resulting in $(3 \times 224 \times 224 + 197 \times 768) \times 32$ bits of storage per cached pair, approximately 1.2 MB. In our experiments, the longest navigation route was roughly 12 steps (from R2R-Back), and if we assume

Algorithm 2 SimHash Algorithm

Input: a current view v_i
Output: a binary hash key

```

1: function HASH( $v_i$ )
2:    $key \leftarrow \emptyset$ 
3:   for each  $hp$  in Hyperplanes do
4:      $sign \leftarrow \text{DotProduct}(hp, v_i)$ 
5:      $hash\_val \leftarrow (sign > 0)$  ▷ converts to binary
6:      $key \leftarrow key + hash\_val$ 
7:   end for
8:   return  $key$ 
9: end function

```

Input: a hash table h , a current view v_i , an embedding e_i
Output: a hash table h

```

10: function ADDTOHASHTABLE( $h, v_i, e_i$ )
11:    $key \leftarrow \text{Hash}(v_i)$ 
12:    $h \leftarrow \text{InsertToHashTable}(key, v_i, e_i)$ 
13:   return  $h$ 
14: end function

```

Input: a hash table h , a current view v_i
Output: an embedding e_i

```

15: function FINDSIMILAR( $h, v_i$ )
16:    $s_{max} \leftarrow -1$ 
17:    $key \leftarrow \text{Hash}(v_i)$ 
18:    $bucket \leftarrow h.get(key)$ 
19:   for each ( $v_{candidate}, e_{candidate}$ ) in bucket do
20:      $s \leftarrow \text{CosineSimilarity}(v_i, v_{candidate})$ 
21:     if  $s > s_{max}$  then
22:        $s_{max} \leftarrow s$ 
23:        $e_{best} \leftarrow e_{candidate}$ 
24:     end if
25:   end for
26:   if  $s_{max} > threshold$  then
27:      $e_i \leftarrow e_{best}$ 
28:   else
29:      $e_i \leftarrow \emptyset$ 
30:   end if
31:   return  $e_i$ 
32: end function

```

all 36 images per panorama are cached, we obtain a worst-case overhead of 522.7 MB. In practice, however, we find that most tasks are 5–7 steps, and we cache at most 14 images per step, producing a more typical overhead of 84.7–118.6 MB. Considering that modern VLN agents Chen et al. (2021; 2022) are orders of magnitude larger, this is not a limiting factor to practical deployment.

D FULL EVALUATION RESULTS

Table 7 complements our main evaluation in Sec 5.1 with additional benchmarks: R2R-Back (Chen et al., 2021), REVERIE (Qi et al., 2020), R2R-Last (Chen et al., 2021), CVDN (Thomason et al., 2020), and SOON (Zhu et al., 2021). For CVDN, we report the additional evaluation metric Goal Progress (GP), which assigns a higher score as the agent moves closer to the goal, indicating better performance (Chen et al., 2021). For REVERIE and SOON, in addition to image features, object features are required during navigation. We were unable to find the original implementation for object feature extraction, so for these benchmarks we use cached object features and apply our strategy only to image feature extraction. To accommodate this in the performance calculations, we report the GFLOPs necessary for image feature processing and treat the computational cost of object feature extraction as a constant (the $+C$ in Table 7). Note that this prevents us from being able to report

Agent	Task	Method	Performance					GFLOPs
			TL	OSR	SR	SPL	GP	
HAMT	R2R-Back	Base	20.56	-	55.43	52.34	-	8181.55
		Ours (All)	20.53	-	49.21	46.47	-	3331.80
	REVERIE	Base	14.07	35.73	31.81	29.17	-	5434.71+C
		Ours (All)	13.70	26.75	24.96	23.13	-	2735.90+C
	R2R-Last	Base	12.28	54.24	47.85	42.27	-	4982.68
		Ours (All)	12.36	49.72	41.93	36.97	-	2589.44
	CVDN	Base	-	-	-	-	4.88	11022.03
		Ours (All)	-	-	-	-	4.45	4773.34
	REVERIE	Base	22.49	51.46	47.09	33.54	-	6185.15+C
		Ours (All)	21.59	46.44	41.32	28.90	-	3350.31+C
DUET	SOON	Base	35.87	50.38	36.19	22.67	-	9997.81+C
		Ours (All)	42.36	54.22	36.43	20.37	-	4533.83+C

Table 7: Comparison of the performance and efficiency of the baseline agents versus our improved-efficiency agents across multiple benchmarks. Here, we denote the cost of object feature extraction as C .

percentage-wise changes in total performance, so we consider the raw reduction in GFLOPs in these cases.

The upper section of the table compares the performance and efficiency of the baseline HAMT agent against our efficient HAMT agent. For R2R-Back, our strategy achieves a 60% reduction in computation with an 11% decrease in SR. For REVERIE, our efficient VLN model reduces computation by 2698.81 GFLOPs, with a 20% drop in SR. For R2R-Last, our method reduces computation by 48%, with a 12% reduction in SR. Finally, for the CVDN evaluation, our efficient model reduces computation by 57%, with only a 9% decrease in GP. The lower section of the table presents a comparison of the performance and efficiencies of the DUET agents. For REVERIE, our strategy saved 2834.84 GFLOPs with a 12% decrease in SR. For SOON, we observed a marginal increase in SR accompanied by a 10% drop in SPL, while saving 5463.98 GFLOPs. Despite the more significant performance drop in the REVERIE task using the HAMT agent, these results demonstrate that our efficiency strategies are applicable across different benchmarks, achieving substantial computational savings while maintaining an acceptable trade-off in performance.

Robustness to navigation length. It is possible that the errors introduced by our method *propagate*, resulting in worse agent navigation for longer trajectories. We study if this is the case by considering the *navigation error* (NE)—the distance of an agent’s final position to the target position (in meters)—on benchmarks with varying path lengths. We deploy all of our proposed methods (simultaneously) on the HAMT agent and report the changes in NE and GFLOPs compared to the baseline in Table 8.

Agent	Task	Average Path Length	Δ NE(\downarrow)	Δ GFLOPs(\downarrow)
HAMT	R2R	6.0	+0.53	-2845.63
	R2R-Last	6.0	+0.45	-2393.24
	R2R-Back	12.0	+0.54	-5463.98
DUET	R2R	6.0	+0.68	-2971.70
	SOON	9.6	-0.44	-5463.98

Table 8: Performance of our efficient HAMT agent on benchmarks with different path lengths. Δ NE and Δ GFLOPs are the changes in navigation error (NE) and GFLOPs compared to the baseline agent. The path length is the minimum number of navigation actions needed to reach the target destination.

We find our method is largely robust to longer path lengths. The NE does not increase for longer trajectories, and we even see a decrease for the SOON benchmark, which has an average path length 3.6 more steps than R2R. The results also show that our efficient VLN agent sees roughly proportional computational savings for longer paths. For example, the average path length in R2R-Back is double R2R, and we achieve a 1.92x larger reduction in GFLOPs for the HAMT agent.

Method	TL(\downarrow)	OSR(\uparrow)	SR(\uparrow)	SPL(\uparrow)	GFLOPs(\downarrow)
None (Base)	11.53	74.29	66.16	61.49	4763.24
k -extension	12.52	71.86	61.30	55.79	2,408.99
thresholds	12.33	72.46	62.62	57.39	3,867.46
LSH	11.53	74.20	66.11	61.47	3,894.76
k -extension+LSH	12.52	71.90	61.17	55.63	2,013.48
k -extension+thresholds	12.89	71.95	60.41	54.57	2,294.23
thresholds+LSH	12.33	72.41	62.49	57.33	3,190.66
All	12.87	71.95	60.41	54.50	1,917.61

Table 10: Performance of all combinations of our speed-up techniques (k -extensions, early-exiting, and LSH) with the HAMT agent on the R2R benchmark.

Runtime comparison. To validate that our approach improves efficiency in the real world, we report the wall-time comparison between our efficient VLN model and the baseline VLN for both HAMT and DUET agents, tested on the R2R validation unseen split, in Table 9. Evidently, our efficient strategy applied to the VLN agents results in significant runtime savings, with an approximate 40% reduction. It is important to note that the disparity between the 60% GFLOPs savings and the 40% runtime reduction can be attributed to various hardware and software related factors.

Task	Agent	Method	Wall-time (s)
R2R	HAMT	Base	200811
		Ours	119514
	DUET	Base	268962
		Ours	170464

Table 9: Wall-time comparison between the baseline agent and our efficient agent on the R2R task.

E PER-MECHANISM ANALYSIS

In most experiments, we treat our proposed mechanisms as a single unit by applying all three simultaneously. While this is the most flexible and offers the best trade-off between performance and efficiency, analyzing each mechanism independently can provide valuable insights into its effectiveness and robustness. Here, we present results on a per-mechanism basis.

Effectiveness. In Sec 5.1, we apply our k -extension technique and then add adaptive thresholding early-exiting (denoted thresholds in Table 2) and locality-sensitive hashing (LSH) as we found those combinations of techniques offer the most computational savings. Here, we study all combinations of three efficiency mechanisms. To use early-exiting and LSH without k -extension, we treat every non-navigable view as one that can be early-exited or hashed. Navigable views are still fully processed. We report results for the HAMT agent on the R2R benchmark in Table 10.

The results show that between individual techniques, k -extension offers the best computational savings with a 49% reduction compared to the baseline agent. Early-exiting and LSH only reduce GFLOPs by $\sim 18\%$ because early-exiting still requires processing every view, and LSH reuses only a minority of cached image embeddings. We find that LSH provides better performance than the other two individual mechanisms, with an SR only 0.05 lower than the baseline. This is likely because the cached embeddings reused by LSH are near-identical, having a negligible impact on performance when interchanged. However, it is far less efficient than when combined with our other techniques.

The combination we do not present in Table 2, early-exiting and LSH (**thresholds+LSH**), provides slightly better performance than combinations using k -extension but at the cost of 39–66% more GFLOPs. Like the individual mechanisms, this suggests that retaining and partially processing/reusing the non-navigable views mitigates performance drop but is not nearly as efficient as k -extension. Overall, we find that all combinations of our techniques fare well, offering different trade-offs between performance and efficiency.

Robustness to natural corruptions. Now, we complement Sec 5.3 and study the robustness of each of our proposed mechanisms to visual corruption. We select the Low Lighting and Motion Blur corruptions based on their varying impact on performance and being more likely to occur in

Corruption	Method	TL(↓)	OSR(↑)	SR(↑)	SPL(↑)	GFLOPs(↓)
Low Lighting	None (Base)	12.15	71.31	62.58	57.23	4903.06
	k -extension	13.86	71.14	57.34	50.78	2571.06
	thresholds	13.63	70.29	58.79	52.16	4099.21
	LSH	12.95	71.43	61.47	55.19	2444.05
Motion Blur	None (Base)	12.41	68.20	59.13	54.01	4996.64
	k -extension	14.03	65.13	53.77	48.01	2588.06
	thresholds	13.81	68.20	57.51	51.05	4073.04
	LSH	12.39	68.03	59.30	54.04	4030.52

Table 11: Performance under visual corruption of our methods applied *independently* to the HAMT agent on the R2R benchmark.

real-world VLN systems. We apply our methods to the HAMT agent and report results on the R2R benchmark in Table 11.

Our methods appear more robust to Low Lighting than Motion Blur, which corroborates our findings in Sec 5.3. Across both corruptions, k -extension and early-exiting see a slight increase of 150–200 GFLOPs compared to the results in Table 10. This can likely be attributed to the increased trajectory length, and for early-exiting, we also find that the OOD samples require more ViT layers before sufficiently saturating. Both mechanisms result in significant drops in performance, though less than when we apply all simultaneously (results shown in Table 6). Early-exiting is slightly more robust, achieving a 2–7% higher SR, which makes sense as it processes strictly more images than k -extension.

Interestingly, LSH functions extremely well when Low Lighting is applied. It offers a ~49% reduction in GFLOPs, compared to just 18% when no corruption is present. We hypothesize that the reduced lighting makes more images similar, causing our algorithm to find more matches and reuse more embeddings. It also offers significant robustness, only incurring a 1% point drop in SR. It seems like our caching mechanism is better suited for this environment, a finding we hope to explore in future work. For Motion Blur, LSH is less successful, being more robust than our other mechanisms but with minimal computational savings.

F RELATED WORK ON MODEL COMPRESSION

Research has proposed an orthogonal approach to reduce the computational demands and memory footprint of deep-learning models: *model compression*. Quantization and pruning are the leading practice in model compression. Quantization (Jacob et al., 2018; Choi et al., 2018; Louizos et al., 2018; Bhalgat et al., 2020; Uhlich et al., 2019; Banner et al., 2019; Choukroun et al., 2019; Li et al., 2021; Nagel et al., 2020) transforms the memory representation of model parameters from 32-bit floating point numbers to a lower-bit integers (e.g., 4-bit integers), thereby making it more storage efficient and lowering memory usage. Pruning (Molchanov et al., 2016; Fan et al., 2019; Fang et al., 2023; Nova et al., 2023; Han et al., 2015b;a; Hoang & Liu, 2023) aims to create sparse models by removing parameters that are less important for maintaining performance, effectively reducing model size and computation.

While quantization and pruning have been demonstrated in simpler unimodal encoder settings for image and text, they are much more challenging in vision-language model(VLM) settings (Wang et al., 2022; Sun et al., 2024) and largely unexplored in VLN. (Wang et al., 2022) highlighted the challenges of pruning VLMs due to the unequal weighting of visual and linguistic modalities. They mitigated this by using a modal-adaptive approach, adjusting pruning ratios across different model components based on downstream task sensitivity. Similarly, (Sun et al., 2024) demonstrated that naively applying post-training quantization to CLIP caused significant performance degradation, which they addressed by introducing prompt tuning and alignment modules.

We expect similar challenges to be exhibited by VLN agents, if not exacerbated. VLN models, in addition to processing language and visual modalities, involve sequential decision-making dependent

on actions taken at each time step. We anticipate the complex interactions between these information sources to require careful consideration while adapting model compression techniques. Future research on such techniques can be superposed along with our input-adaptive inference method to develop highly efficient models with an acceptable performance trade-off.

G COMPARISON OF ADDITIONAL SIMILARITY METRICS

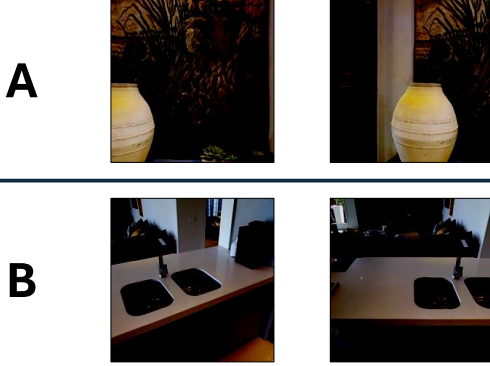


Figure 9: **Two sets of example views (A and B)** demonstrating non-identical but similar views that have been slightly shifted during navigation.

Simiarlity Metrics	Set A	Set B
SSIM (Wang et al., 2004)	0.24	0.32
FSIM (Zhang et al., 2011)	0.26	0.27
LPIPS (Zhang et al., 2018)	0.55	0.62
SURF (Bay et al., 2006)	0.31	0.32
SIFT (Lowe, 2004)	0.45	0.37
ORB (Rublee et al., 2011)	0.07	0.19

Figure 10: **Similarity scores measured on Set A and B.** We test 6 different similarity metrics.

Other than the three similarity metrics we use, we test three additional metrics for comparison: SURF (Bay et al., 2006), SIFT (Lowe, 2004), and ORB (Rublee et al., 2011). These are feature detection and description algorithms designed to identify and match keypoints in images. The similarity scores are computed by dividing the number of matching keypoints by the minimum number of keypoints detected in the two images. We test all six algorithms on two sets of scenes, reflecting shifts caused by an agent’s changing perspectives during navigation.

Figure 9 illustrates the two sets of scenes, and Table 10 summarizes the quantitative comparison. Among the three metrics we employ for our main evaluation, LPIPS demonstrates a higher similarity measure of approximately 60% for both sets. In contrast, SSIM and FSIM are less effective at capturing the similarity between views in Sets A and B. The three additional metrics (SURF, SIFT, and ORB) are also ineffective in providing reliable similarity scores for both image sets A and B. Our qualitative comparison of different similarity metrics applied to sets of similar scenes highlight the challenges these metrics face in accurately identifying true visual similarity. We believe that an accurate measure of scene similarity is crucial for further reducing the computational demands of a VLN agent, and we leave this for future work.

H ANALYZING PERFORMANCE-EFFICIENCY TRADE-OFF IN OUR METHOD

In order to illustrate our tunable performance-efficiency trade-off, we show that even when limiting the performance drop to under 5%, our input adaptive inference method applied to the HAMT agent achieves significant computational savings. For reference, the baseline HAMT model achieves a SR of 66.16 with a computational cost of 4763.24 GFLOPs. Figure 11 shows that with a 3–5% drop in SR, we still manage to achieve 43–50% savings in GFLOPs. These results were tested on the R2R validation unseen dataset.

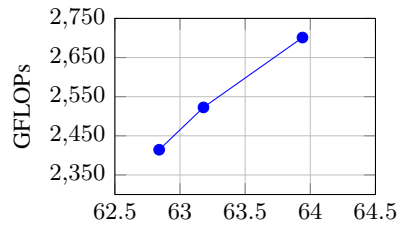


Figure 11: Trade-off between Performance (SR) and GFLOPs