

SPARQ ATTENTION: BANDWIDTH-EFFICIENT LLM INFERENCE

Luka Ribar*, Ivan Chelombiev*, Luke Hudlass-Galley*

Charlie Blake, Carlo Luschi, Douglas Orr

Graphcore Research

London, UK

{lukar, ivanc, lukehg, charlieb, carlo, douglaso}@graphcore.ai

ABSTRACT

Through an analysis of the statistical properties of pre-trained large language models (LLMs), we highlight two opportunities for sparse memory access: first in the components of query and keys and second in the attention scores corresponding to key, value pairs. Based on this, we introduce **SparQ Attention**, a technique for increasing the inference throughput of LLMs by utilising memory bandwidth more efficiently within attention layers, through selective fetching of the cached history. Our proposed technique can be applied directly to off-the-shelf LLMs during inference, without requiring any modification to the pre-training setup or additional fine-tuning. We show that SparQ Attention brings up to $8\times$ savings in attention data-transfers without substantial drops in accuracy, by evaluating Llama 2, Mistral and Pythia models on a wide range of downstream tasks.

1 INTRODUCTION

Transformer-based large language models (LLMs) trained on large corpora of text have recently shown remarkable performance on complex natural language processing tasks (Achiam et al., 2023; Touvron et al., 2023). Many applications require LLMs to support long input sequences (e.g. long instructions, chat histories, and relevant documents). However, the standard optimisation used during batched sequence generation, *key-value (KV) caching* (Pope et al., 2023), is constrained by the need to fetch a large amount of data from memory. This in turn limits the speed at which tokens can be generated—a key usability metric for LLMs.

Despite this expensive cache-fetch at each step, tokens generally only attend to a small part of the sequence at a time (Vig, 2019; Yun et al., 2020). Building upon this understanding, we show that it is possible to predict tokens with high attention scores without additional training, by selecting *components* of the query and key tensors. We also outline a better approximation scheme to achieve sparse attention, where attention mass is reallocated to a mean-value vector. Following these observations, we present **SparQ (Sparse Query) Attention**, a practical technique for improving the memory bandwidth efficiency of transformer inference by 1) predicting high-scoring tokens using query components, 2) fetching these from the KV cache and 3) reallocating the attention score of low-scoring tokens to a mean-value vector.

2 COMPUTATIONAL EFFICIENCY

In this section we provide a straightforward framework to understand the computational efficiency of sequence generation using transformer models (similar to the modelling introduced by Kaplan et al. (2020)) and use it to motivate transfer-efficient attention mechanisms.

Arithmetic intensity A compute unit capable of r_A scalar arithmetic operations per second is connected to a memory via an interface that can transfer r_M scalar elements per second, processing a workload requiring \mathcal{A} arithmetic operations and \mathcal{M} transfers. Assuming concurrent compute and data transfer, when the *arithmetic intensity* \mathcal{A}/\mathcal{M} of the workload is less than the ratio r_A/r_M , execution time is limited by r_M .

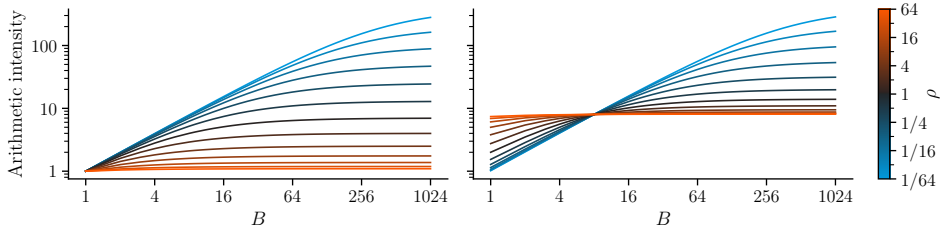


Figure 1: Relationship between $\rho = S/(gd_m)$, batch size B and arithmetic intensity during sequence generation. *Left*: Multi-head attention. *Right*: GQA ($g=8$). ML hardware provides $r_A/r_M > 200$, making memory bandwidth the limiting factor in many practical scenarios (see Appendix D).

Sequence generation Consider a full transformer layer, with N parameters, batch size B , and C elements in the attention KV cache per batch element. We assume Grouped Query Attention (GQA) (Ainslie et al., 2023) with g grouped-query heads ($g = 1$ for standard multi-head attention). This implies the arithmetic intensity:

$$\frac{\mathcal{A}}{\mathcal{M}} = \frac{BN + BCg}{N + BC} = \frac{N + Cg}{N/B + C} \tag{1}$$

We can increase arithmetic intensity by making B large, causing \mathcal{A}/\mathcal{M} to approach $N/C + g$. Hence the limiting factor for large-batch transformer inference is the ratio of the KV cache size per-item to the size of the model. An alternative formulation for a standard transformer with model dimension d_m and sequence-length S , has $N = 12(d_m)^2$ and $C = 2 S d_m/g$, giving:

$$\frac{\mathcal{A}}{\mathcal{M}} = \frac{6 + \rho g}{6/B + \rho} \tag{2}$$

where $\rho = S/(gd_m)$. The value ρ underlies the KV cache-model size relationship outlined above, determining the point at which the model becomes memory bandwidth bound (Figure 1). Since long sequences are desirable, data transfer is the performance-limiting factor, motivating the search for transfer-efficient approximations to full attention.

3 APPROXIMATING ATTENTION

A single attention *query head* produces the output $\mathbf{y} = \mathbf{s} \cdot \mathbf{V}$, given *attention scores* $\mathbf{s} \in (0, 1)^S$:

$$\mathbf{s} = \text{softmax}\left(\frac{\mathbf{q} \cdot \mathbf{K}^\top}{\sqrt{d_h}}\right) \tag{3}$$

where $\mathbf{q} \in \mathbb{R}^{d_h}$ is the query, and $\mathbf{K}, \mathbf{V} \in \mathbb{R}^{S \times d_h}$ are the key and value caches respectively. In multi-head attention, the number of scalar elements transferred when calculating \mathbf{y} is given by

$$\mathcal{M}_{\text{base}} = 2 S d_h + 2 d_h \tag{4}$$

where the first term corresponds to reading the \mathbf{K} and \mathbf{V} caches and the second term corresponds to writing the current \mathbf{k} and \mathbf{v} to memory.

Attention scores sparsity (Sheng et al., 2023) Due to the normalising effect of the softmax function, the resulting \mathbf{s} vector is *sparse* (see Figures 2a and 2b), i.e. we can find a boolean mask $\mathbf{m}_s \in \{0, 1\}^S$ corresponding to the top- k elements in \mathbf{s} ($k \ll S$) such that:

$$\mathbf{y}_1 = (\mathbf{s} \circ \mathbf{m}_s) \cdot \mathbf{V} \approx \mathbf{s} \cdot \mathbf{V} \tag{5}$$

As a result, only the values \mathbf{v}_i corresponding to the non-zero elements of \mathbf{m}_s need to be fetched from memory. However, we still must fetch the full \mathbf{K} from memory to calculate the attention scores \mathbf{s} .

Mean value reallocation We note that \mathbf{v}_i vectors from the same sequence exhibit a high degree of auto-correlation (see Table E1). We therefore introduce an additional correction term to improve our approximation, using a running-mean value vector $\bar{\mathbf{v}} = \frac{1}{S} \sum_{i=1}^S \mathbf{v}_i$:

$$\mathbf{y}_2 = (\mathbf{s} \circ \mathbf{m}_s) \cdot \mathbf{V} + (1 - \mathbf{s} \cdot \mathbf{m}_s) \bar{\mathbf{v}} \tag{6}$$

This introduces a minimal additional overhead compared to Equation (5) as the mean vector $\bar{\mathbf{v}}$ needs to be updated and written back to memory at each step.

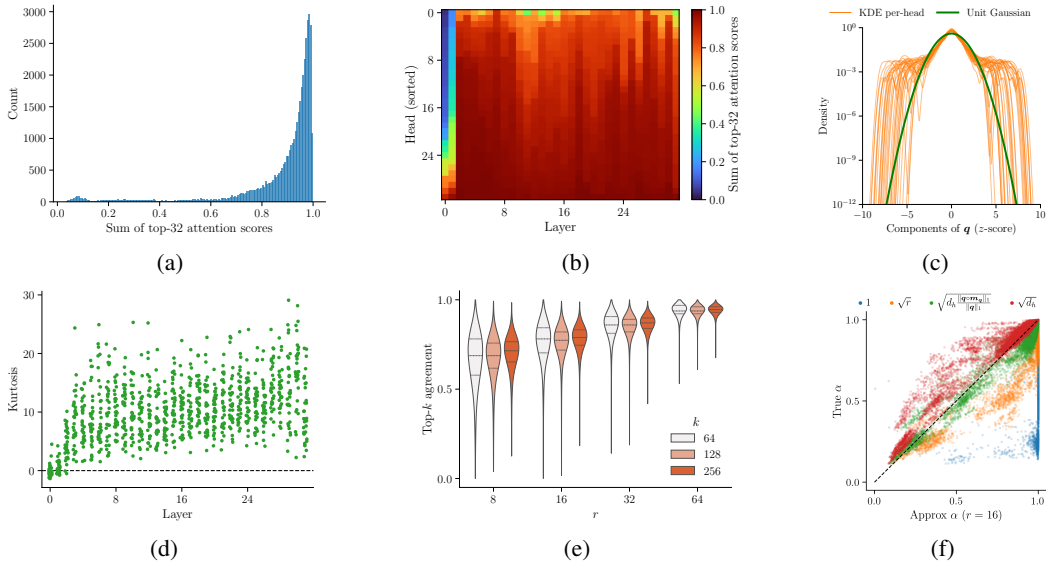


Figure 2: Statistics of Llama 2 7B over 40 SQuAD queries, for all 32 layers \times 32 heads unless noted. (a) Sum softmax output allocated to the 32 highest-scoring positions, demonstrating natural attention sparsity; (b) for each head. (c) Kernel density estimate (Rosenblatt, 1956) of components of q in layer 16, showing heavy tails. (d) Fisher Kurtosis of q components, showing that the query vector is leptokurtic for most heads. (e) Top- k agreement, the proportion of the top- k positions that are correctly predicted by an approximated softmax for various r and k . (f) Agreement between the coverage α based on estimated scores versus the true mass of the top 128 scores, for different softmax temperatures (a point for each example \times head), showing the importance of temperature.

Query sparsity We now consider approximations to the mask m_s via approximate attention scores \hat{s} without using the full matrix K . Here, we observe that the query vector q is *heavy-tailed* (see Figures 2c and 2d), indicating that basis-aligned sparsity (masking) should be effective. We introduce a per-query boolean mask $m_q \in \{0, 1\}^{d_h}$ corresponding to the top- r components of q . The scores are then approximated as:

$$\hat{s} = \text{softmax}\left(\frac{(q \circ m_q) \cdot K^\top}{\tau}\right) \tag{7}$$

where τ is the softmax temperature. Only the components of K corresponding to non-zero elements of m_q need to be fetched from memory. The top- k mask $m_{\hat{s}} \in \{0, 1\}^S$ can then be calculated using \hat{s} (showing good agreement with m_s , see Figure 2e), and the approximate attention output is obtained as:

$$y_3 = \text{softmax}\left(\frac{q \cdot K^\top}{\sqrt{d_h}} + \log(m_{\hat{s}} + \epsilon)\right) \cdot V \tag{8}$$

with $\epsilon \rightarrow 0$. Only the k_i, v_i pairs corresponding to non-zero $[m_{\hat{s}}]_i$ need to be fetched from memory.

Mean value reallocation with query sparsity Finally, we reintroduce the mean value reallocation improvement of Equation (6). As we do not have access to the full scores s , we approximate the weighted sum using approximate scores \hat{s} . Since the query-key dot product is performed over only r dimensions, the softmax temperature τ from Equation (7) must be carefully chosen. If r components were chosen randomly, the appropriate temperature would be \sqrt{r} . On the other hand, if the top- r components were the only non-zero elements of the query vector, the appropriate temperature would remain $\sqrt{d_h}$. As a balance between the two extremes, we have found $\tau = \sqrt{d_h} \frac{\|q \circ m_q\|_1}{\|q\|_1}$ to yield a good approximation (see Figure 2f).

The final attention output is then calculated as a weighted sum $y = \alpha y_3 + (1 - \alpha)\bar{v}$, where $\alpha = m_{\hat{s}} \cdot \hat{s}$ is the relative weight of the top- k terms.

Table 1: Results for the largest models tested. SQuAD and TriviaQA measure performance in accuracy. CNN/DailyMail uses ROUGE-L score. WikiText measures perplexity in bits per character (BPC) and Repetition counts the number of characters before generation diverges. The best score for a model, task and compression setting is in bold. Median standard errors across all models and sparsity settings are: SQuAD 0.7, TriviaQA 0.7, CNN/DailyMail 0.4, WikiText 0.006, Repetition 2.

| Dataset Name | | SQuAD \uparrow | | | TriviaQA \uparrow | | | CNN/DailyMail \uparrow | | | WikiText \downarrow | | | Repetition \uparrow | | |
|----------------|------------------|------------------|-------------|-------------|---------------------|-------------|-------------|--------------------------|-------------|-------------|-----------------------|-------------|-------------|-----------------------|------------|------------|
| Compression | | 1 | 1/2 | 1/8 | 1 | 1/2 | 1/8 | 1 | 1/2 | 1/8 | 1 | 1/2 | 1/8 | 1 | 1/2 | 1/8 |
| Llama 2 13B | LM- ∞ | | 50.0 | 32.4 | | 73.4 | 69.0 | | 16.8 | 15.1 | | 0.64 | 0.69 | | 76 | 29 |
| | H ₂ O | 80.8 | 73.2 | 64.1 | 78.7 | 78.5 | 78.4 | 22.1 | 22.2 | 20.8 | 0.61 | 0.61 | 0.63 | 229 | 61 | 26 |
| | SparQ | | 80.7 | 78.0 | | 78.8 | 78.2 | | 22.5 | 22.2 | | 0.61 | 0.64 | | 227 | 190 |
| Mistral 7B | LM- ∞ | | 51.0 | 31.6 | | 75.8 | 72.8 | | 18.0 | 16.8 | | 0.65 | 0.70 | | 81 | 20 |
| | H ₂ O | 81.0 | 71.2 | 59.2 | 80.9 | 80.8 | 80.6 | 23.7 | 23.5 | 23.4 | 0.62 | 0.63 | 0.65 | 231 | 38 | 14 |
| | SparQ | | 80.9 | 77.5 | | 80.8 | 79.0 | | 23.5 | 23.0 | | 0.63 | 0.65 | | 209 | 201 |
| Pythia 6.9B | LM- ∞ | | 38.5 | 18.9 | | 41.6 | 32.0 | | 14.9 | 14.1 | | 0.71 | 0.77 | | 64 | 18 |
| | H ₂ O | 57.8 | 52.9 | 46.6 | 52.6 | 52.6 | 52.3 | 20.2 | 20.3 | 18.9 | 0.68 | 0.69 | 0.71 | 150 | 47 | 19 |
| | SparQ | | 58.0 | 57.1 | | 52.4 | 51.7 | | 20.6 | 20.6 | | 0.68 | 0.70 | | 151 | 144 |

SparQ Attention Following this analysis, our method (see Appendix A), consists of three steps:

Step 1: Find the indices of r largest components of $|\mathbf{q}|$ and only fetch \mathbf{K} along the corresponding dimensions. Calculate *approximate* attention scores $\hat{\mathbf{s}}$ using the sliced query and keys.

Step 2: Find the top- k positions in the approximate attention scores and fetch the corresponding *full* key and value vectors. Calculate the output of the attention layer using the top- k keys and values.

Step 3: Estimate the total score α assigned to the top- k positions using $\hat{\mathbf{s}}$. Use this total score to interpolate between the attention output from the top- k positions, and a *mean value* vector, $\bar{\mathbf{v}}$.

The memory transfer of the SparQ Attention algorithm for a single attention head forward-pass:

$$\mathcal{M}_{\text{SparQ}} = S r + 2 k d_h + 4 d_h \quad (9)$$

where the first term corresponds to reading r rows of \mathbf{K} , the second term corresponds to reading the top- k columns of \mathbf{K} and \mathbf{V} and the third term corresponds to transfers associated with writing the current \mathbf{k} and \mathbf{v} , in addition to reading and writing $\bar{\mathbf{v}}$. The compression ratio is defined as $\mathcal{M}_{\text{SparQ}}/\mathcal{M}_{\text{base}}$.

4 EXPERIMENTS AND RESULTS

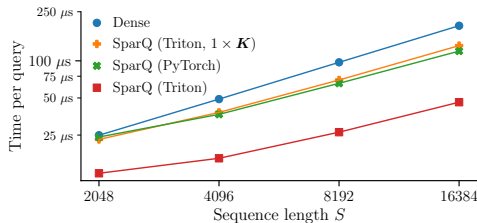
We compare SparQ Attention against two alternative sparse attention methods: *H₂O* (Zhang et al., 2023), an eviction scheme which iteratively removes tokens from the KV cache that are predicted to not influence future token generation, *FlexGen*, in which the exact top- k highest attention scores are calculated before transferring the corresponding columns of \mathbf{V} , and *LM-Infinite* (Han et al., 2023), which only transfers a fixed number of initial tokens and recent tokens from the sequence.

We evaluate these different methods over a range of compression ratios on five different NLP tasks, with datasets adapted to generate sequence lengths between 1k-2k tokens (see Appendix C). Our experiments span six models: Llama 2 {7B, 13B}, Mistral 7B, and Pythia {1.4B, 2.8B, 6.9B} (Touvron et al., 2023; Jiang et al., 2023; Biderman et al., 2023). Results from Llama 2 13B and Mistral 7B are presented in Table 1, with further results in Figures B1 and B2, showing that:

- SparQ Attention performance is robust across all tasks and model sizes tested. Compression ratios of $2\times$ to $8\times$ are readily achievable with little to no loss in task performance.
- The simple recipe of setting $k = 128$ and tuning r to set the compression/performance trade-off seems generally robust over all models and tasks considered (see Figure F3).
- Certain tasks are more challenging for H₂O (Repetition, SQuAD), while others are more forgiving (TriviaQA, WikiText-103).
- LM-Infinite degrades performance across all tasks, demonstrating that the tasks do not permit the trivial solution of discarding the long input sequence.

| Kernel | A100 (40GB) | A10G |
|---------------------------------|-----------------------------|-----------------------------|
| Dense | 49 μ s (1 \times) | 128 μ s (1 \times) |
| SparQ (Triton, 1 \times K) | 38 μ s (1.28 \times) | 79 μ s (1.63 \times) |
| SparQ (PyTorch) | 37 μ s (1.33 \times) | 78 μ s (1.63 \times) |
| SparQ (Triton) | 16 μ s (3.02 \times) | 31 μ s (4.17 \times) |

(a)



(b)

Figure 3: GPU performance, with batch size 64, $r = 32$, $k = 128$. (a) Microbenchmark results with sequence length $S = 4096$. (b) Scaling sequence length on A100 (40GB) GPU.

5 BENCHMARKING

The results above use a theoretical cost model of total memory transfers, allowing us to evaluate SparQ Attention independently of a specific hardware setup. To validate this approach, we performed a set of microbenchmarks of an attention operation in isolation.

SparQ Attention benefits from two optimisations. The first is to store K twice, in both d_h -contiguous and S -contiguous layouts, since this allows for an efficient gather (indexing) on either axis, at the cost of 50% extra memory usage. The second optimisation is to use a fused gather-then-matmul operation to avoid writing the result of the gather to memory.

We tested multiple implementations of baseline and SparQ Attention on IPU using the Poplar C++ interface and GPU using PyTorch (Paszke et al., 2019). In all cases, we used the Llama 2 7B shape parameters: 32 heads, $d_h = 128$. The implementations tested were: **Dense** baseline, choosing the faster of a plain PyTorch implementation and the built-in `scaled_dot_product_attention`, **SparQ (Triton)**, storing K twice and using fused gather-then-matmul kernels written using Triton (Tillet et al., 2019), **SparQ (PyTorch)**, with no Triton and **SparQ (Triton, 1 \times K)**, storing K in d_h -contiguous layout only, for no additional memory cost.

Our achieved GPU speed-ups are presented in Figure 3a, and the performance trend with sequence length is shown in Figure 3b. Standard error for all results given is $< 1\%$ of the mean. See Appendix G for further details. These microbenchmark results show that the theoretical benefits of SparQ Attention can yield substantial wall-clock time speedups on current hardware. Further work is needed to show improvements for small batches, and to investigate alternatives to storing K twice.

6 CONCLUSION

In this work, we explore the scalability of attention mechanisms in modern language models, and find that in many realistic settings, transferring the KV cache from memory creates an LLM inference bottleneck. By analysing the statistics of tensors within attention, we have identified opportunities to approximate attention by sparsely accessing the KV cache, reducing total data transfer.

These opportunities have been realised in SparQ Attention, a novel technique for unlocking faster inference for pre-trained LLMs, without any fine-tuning or modifications to the weights of the model. Our proposed technique modifies the attention mechanism by predicting keys that yield large attention scores, permitting only the relevant tokens from the KV cache to be transferred on every generation step. This allows for pre-trained models to be executed more efficiently.

Related work Sparse attention as an architectural change has been explored by Child et al. (2019); Ren et al. (2021); Beltagy et al. (2020); Zaheer et al. (2020); Kitaev et al. (2020); Tay et al. (2020); Xiao et al. (2023). Compared with existing post-training techniques (Sheng et al., 2023; Mao et al., 2023; Chen et al., 2021), SparQ Attention relies on component-wise sparsity of q , K and introduces V reallocation.

The full paper is available at <https://arxiv.org/abs/2312.04985>.

ACKNOWLEDGEMENTS

We would like to thank Daniel Justus, Paul Balana and Andrew Fitzgibbon for their helpful input and feedback on this work.

REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Joshua Ainslie, James Lee-Thorp, Michiel de Jong, Yury Zemlyanskiy, Federico Lebr3n, and Sumit Sanghai. GQA: Training generalized multi-query transformer models from multi-head checkpoints. *arXiv preprint arXiv:2305.13245*, 2023.
- Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*, 2020.
- Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle OBrien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pp. 2397–2430. PMLR, 2023.
- Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. Improving language models by retrieving from trillions of tokens. In *International conference on machine learning*, pp. 2206–2240. PMLR, 2022.
- Beidi Chen, Tri Dao, Eric Winsor, Zhao Song, Atri Rudra, and Christopher R3. Scatterbrain: Unifying sparse and low-rank attention. *Advances in Neural Information Processing Systems*, 34: 17413–17426, 2021.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023.
- Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse transformers. *arXiv preprint arXiv:1904.10509*, 2019.
- Graphcore. Bow-2000 datasheet. (Online: accessed 25 January 2024), March 2023. URL <https://docs.graphcore.ai/projects/bow-2000-datasheet>.
- Chi Han, Qifan Wang, Wenhan Xiong, Yu Chen, Heng Ji, and Sinong Wang. LM-infinite: Simple on-the-fly length generalization for large language models. *arXiv preprint arXiv:2308.16137*, 2023.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension. *arXiv preprint arXiv:1705.03551*, 2017.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- Andrej Karpathy. The unreasonable effectiveness of recurrent neural networks. (Online: accessed 27 January 2024), 2015. URL <https://github.com/karpathy/char-rnn>.
- Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. *arXiv preprint arXiv:2001.04451*, 2020.

- Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pp. 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics.
- Yuzhen Mao, Martin Ester, and Ke Li. Iceformer: Accelerated inference with long-sequence transformers on CPUs. In *Third Workshop on Efficient Natural Language and Speech Processing (ENLSP-III): Towards the Future of Large Language Models and their Emerging Descendants*, 2023.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. *arXiv preprint arXiv:1609.07843*, 2016.
- NVIDIA. NVIDIA A10 datasheet. (Online: accessed 22 January 2024), March 2022. URL <https://www.nvidia.com/content/dam/en-zz/Solutions/Data-Center/a10/pdf/datasheet-new/nvidia-a10-datasheet.pdf>.
- NVIDIA. NVIDIA H100 datasheet. (Online: accessed 22 January 2024), July 2023. URL <https://www.nvidia.com/en-gb/data-center/h100/>.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- Reiner Pope, Sholto Douglas, Aakanksha Chowdhery, Jacob Devlin, James Bradbury, Jonathan Heek, Kefan Xiao, Shivani Agrawal, and Jeff Dean. Efficiently scaling transformer inference. *Proceedings of Machine Learning and Systems*, 5, 2023.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*, 2016.
- Hongyu Ren, Hanjun Dai, Zihang Dai, Mengjiao Yang, Jure Leskovec, Dale Schuurmans, and Bo Dai. Combiner: Full attention transformer with sparse computation cost. *Advances in Neural Information Processing Systems*, 34:22470–22482, 2021.
- Alexis Roche, Grégoire Malandain, Xavier Pennec, and Nicholas Ayache. The correlation ratio as a new similarity measure for multimodal image registration. In *Medical Image Computing and Computer-Assisted Intervention MICCAI98: First International Conference Cambridge, MA, USA, October 11–13, 1998 Proceedings 1*, pp. 1115–1124. Springer, 1998.
- Murray Rosenblatt. Remarks on Some Nonparametric Estimates of a Density Function. *The Annals of Mathematical Statistics*, 27(3):832 – 837, 1956.
- Abigail See, Peter J Liu, and Christopher D Manning. Get to the point: Summarization with pointer-generator networks. *arXiv preprint arXiv:1704.04368*, 2017.
- Ying Sheng, Lianmin Zheng, Binhang Yuan, Zhuohan Li, Max Ryabinin, Beidi Chen, Percy Liang, Christopher Ré, Ion Stoica, and Ce Zhang. FlexGen: high-throughput generative inference of large language models with a single GPU. In *International Conference on Machine Learning*, pp. 31094–31116. PMLR, 2023.
- Yi Tay, Dara Bahri, Liu Yang, Donald Metzler, and Da-Cheng Juan. Sparse sinkhorn attention. In *International Conference on Machine Learning*, pp. 9438–9447. PMLR, 2020.
- Philippe Tillet, Hsiang-Tsung Kung, and David Cox. Triton: an intermediate language and compiler for tiled neural network computations. In *Proceedings of the 3rd ACM SIGPLAN International Workshop on Machine Learning and Programming Languages*, pp. 10–19, 2019.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shrutvi Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Jesse Vig. A multiscale visualization of attention in the transformer model. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pp. 37–42, 01 2019.

Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming language models with attention sinks. *arXiv preprint arXiv:2309.17453*, 2023.

Chulhee Yun, Yin-Wen Chang, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank Reddi, and Sanjiv Kumar. $O(n)$ connections are expressive enough: Universal approximability of sparse transformers. *Advances in Neural Information Processing Systems*, 33:13783–13794, 2020.

Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. Big bird: Transformers for longer sequences. *Advances in Neural Information Processing Systems*, 33:17283–17297, 2020.

Zhenyu Zhang, Ying Sheng, Tianyi Zhou, Tianlong Chen, Lianmin Zheng, Ruisi Cai, Zhao Song, Yuandong Tian, Christopher Ré, Clark Barrett, et al. H₂O: Heavy-hitter oracle for efficient generative inference of large language models. *arXiv preprint arXiv:2306.14048*, 2023.

A SPARQ ATTENTION

A.1 SCHEMATIC

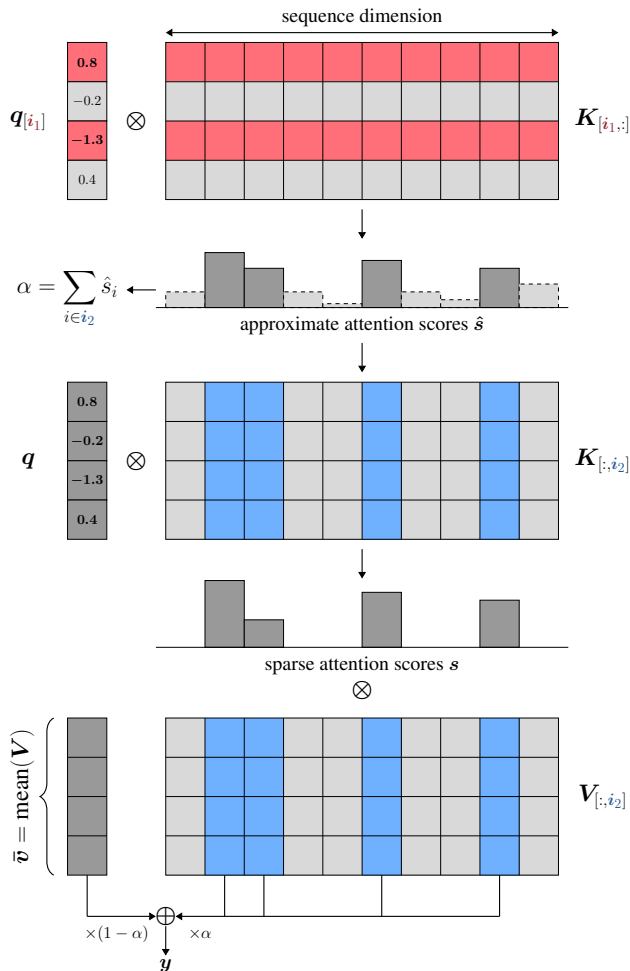


Figure A1: SparQ Attention for a single attention head. The algorithm consists of three steps. First, we find the r largest components of the incoming query vector and gather the corresponding components along the hidden dimension of the key cache K . This allows us to approximate the full attention scores (\hat{s}). In the second step, we identify the top- k largest scores in the approximation and proceed to gather the corresponding full key and value vectors from the cache. As a final step, to compensate for the missing value vectors, we additionally maintain and fetch the running mean value vector \bar{v} and reassign it the leftover mass based on approximate score weightings. The attention output is then calculated as usual using the top- k fetched key and value pairs, together with \bar{v} .

A.2 ALGORITHM

Algorithm 1 SparQ Attention

Input: $\mathbf{q} \in \mathbb{R}^{d_h}$, $\mathbf{K} \in \mathbb{R}^{S \times d_h}$, $\mathbf{V} \in \mathbb{R}^{S \times d_h}$, $\bar{\mathbf{v}} \in \mathbb{R}^{d_h}$, $r \in \mathbb{N}$, $k \in \mathbb{N}$, $l \in \mathbb{N}$

$$\mathbf{i}_1 \leftarrow \text{argtopk}(|\mathbf{q}|, r) \quad \# \text{Indices of top } r \text{ elements of } |\mathbf{q}|$$

$$\tau \leftarrow \sqrt{d_h \cdot \frac{\|\mathbf{q}_{[\mathbf{i}_1]}\|_1}{\|\mathbf{q}\|_1}} \quad \# \text{Softmax temperature, weighted by L1 coverage}$$

$$\hat{\mathbf{s}} \leftarrow \text{softmax}(\mathbf{q}_{[\mathbf{i}_1]} \cdot \mathbf{K}_{[\mathbf{i}_1, :]}^\top / \tau) \quad \# \text{Approximate attention scores (all positions)}$$

$$\mathbf{m} \leftarrow [1 \ i > S - l \ 0]_{i=1}^S \quad \# \text{Local mask of last } l \text{ positions}$$

$$\mathbf{i}_2 \leftarrow \text{argtopk}(\hat{\mathbf{s}} + \mathbf{m}, k) \quad \# \text{Indices of top } k \text{ approximate scores or local}$$

$$\alpha \leftarrow \text{sum}(\hat{\mathbf{s}}_{[\mathbf{i}_2]}) \quad \# \text{Total approximate score of top } k$$

$$\mathbf{s} \leftarrow \text{softmax}(\mathbf{q} \cdot \mathbf{K}_{[:, \mathbf{i}_2]}^\top / \sqrt{d_h}) \quad \# \text{Final attention scores (top } k \text{ positions)}$$

$$\mathbf{y} \leftarrow \alpha \mathbf{s} \cdot \mathbf{V}_{[:, \mathbf{i}_2]} + (1 - \alpha) \bar{\mathbf{v}} \quad \# \text{Mixed scores and values, interpolating with } \bar{\mathbf{v}}$$

return \mathbf{y}

A.3 CODE

```

from torch import softmax, sqrt, tensor, topk

def gather(t, dim, i):
    dim += (dim < 0) * t.ndim
    return t.gather(dim, i.expand(*t.shape[:dim], i.shape[dim], *t.shape[dim + 1 :]))

def attn(Q, K, V, M):
    s = (Q @ K.transpose(-1, -2)) / sqrt(tensor(Q.shape[-1])) + M
    y = softmax(s, dim=-1) @ V
    return y

def sparq_attn(Q, K, V, V_mean, M, r, k):
    # Q -- (batch_size, n_kv_heads, n_heads // n_kv_heads, 1, head_size)
    # K, V -- (batch_size, n_kv_heads, 1, seq_len, head_size)

    # 1. Approximate attention scores using r largest components of Q
    i1 = topk(abs(Q).sum(dim=2, keepdim=True), r, -1).indices
    Q_hat, K_hat = gather(Q, -1, i1), gather(K, -1, i1)
    scale = sqrt(
        Q.shape[-1]
        * abs(Q_hat).sum(dim=-1, keepdim=True)
        / abs(Q).sum(dim=-1, keepdim=True)
    )
    s_hat = softmax(Q_hat @ K_hat.transpose(-1, -2) / scale + M, dim=-1)

    # 2. Gather top k positions based on approximate attention scores & run attention
    i2 = topk(s_hat.sum(dim=2, keepdim=True), k, -1).indices
    iKV = i2[... , 0, :, None]
    K, V, M = gather(K, -2, iKV), gather(V, -2, iKV), gather(M, -1, i2)
    y_ = attn(Q, K, V, M)

    # 3. Estimate the total score of the top k, and interpolate with V_mean
    alpha = gather(s_hat, -1, i2).sum(-1, keepdim=True)
    y = alpha * y_ + (1 - alpha) * V_mean
    return y

```

B DETAILED RESULTS

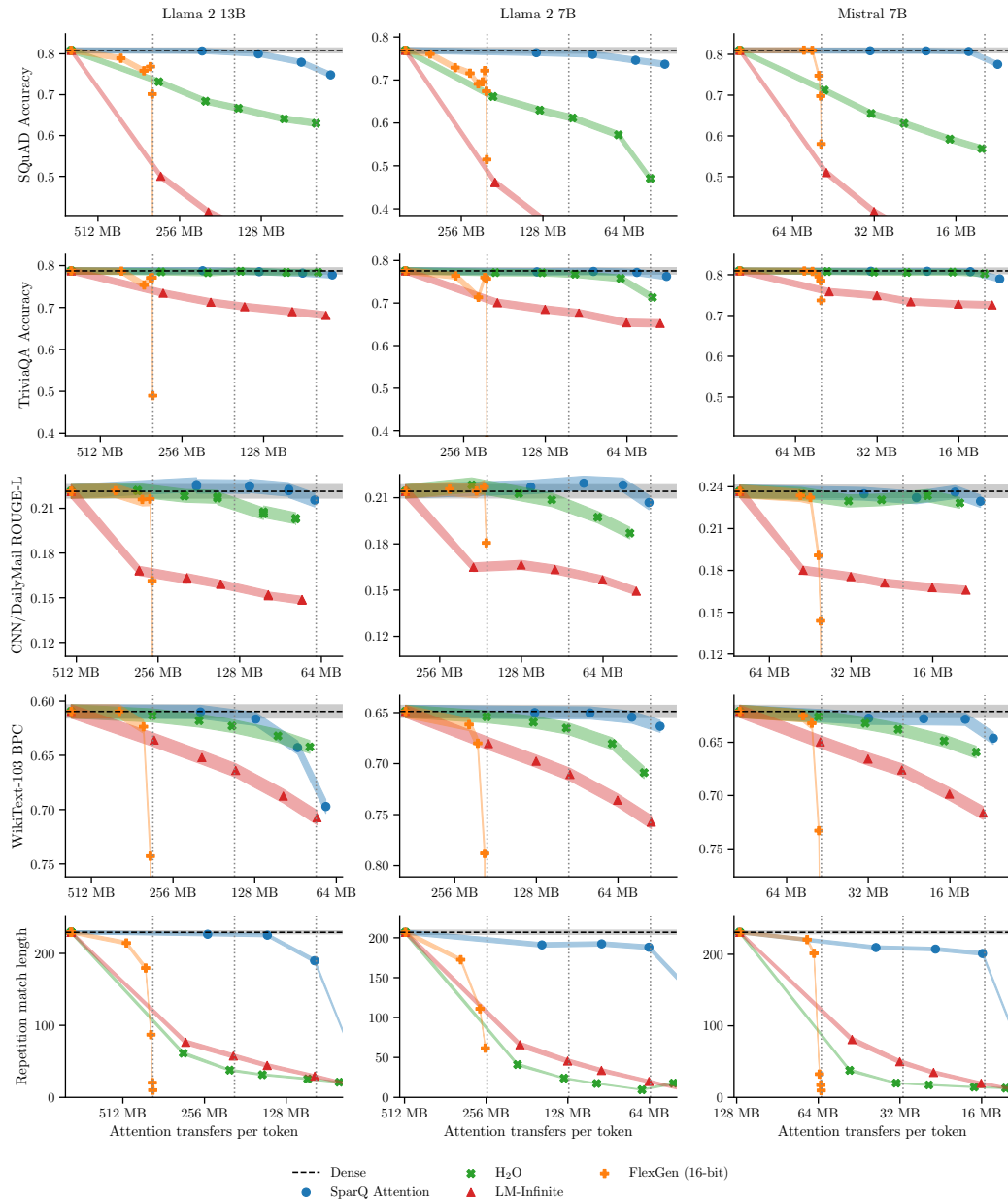


Figure B1: Compression versus performance trade-off curves over all tasks and multiple models. The y-axis minimum is set to $(0.5, 0.5, 0.5, 1.25, 0.0) \times$ the dense baseline for the tasks, reading top-to-bottom, in order to give a consistent view of the performance loss across models. Vertical dotted lines show $(1/2)$, $(1/4)$ and $(1/8)$ compression versus dense. Shaded lines show ± 1 standard error of the mean (uncertainty due to a finite test set).

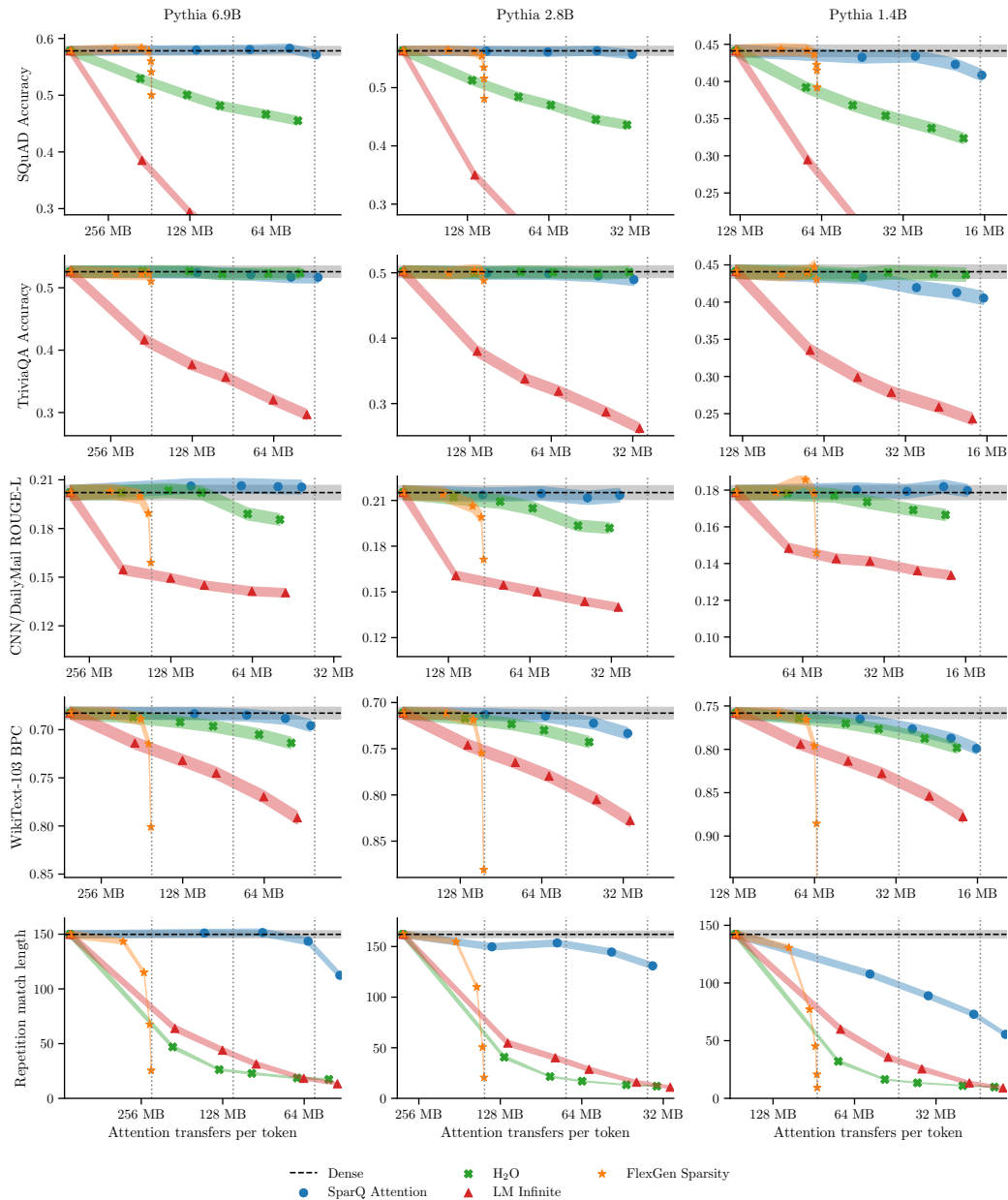


Figure B2: Compression versus performance trade-off curves for Pythia models (see Figure B1).

C TASKS AND DATASETS

In order to evaluate our method on a spectrum of relevant NLP tasks that present a particular challenge to sparse attention techniques, our evaluation setup consists of various tasks requiring information retrieval and reasoning over long input sequences. This includes question answering, summarisation, perplexity/bits-per-character (BPC), and text repetition. For this, we adapted standard downstream tasks and datasets to generate examples of sequence lengths between 1k and 2k tokens. As we wanted to define the tasks independently of the selected models, our examples were chosen to have sequence lengths between 4000 and 8000 characters, roughly giving the desired lengths in tokens.

For question answering, we use the SQuAD (Rajpurkar et al., 2016) and TriviaQA (Joshi et al., 2017) datasets in the *open-book* setting. In order to construct the SQuAD examples, we augment the provided context (i.e. the standard SQuAD input sequence required to answer the question) with seven additional “*confusion contexts*” from unrelated questions. This ensures that the examples have a large sequence length, while making the task harder as the model needs to distinguish the relevant information from the context from the unrelated paragraphs. We use SQuAD v1.1, as it does not include unanswerable questions included in SQuAD v2.0, since we aim to measure the model’s ability to extract useful information from the KV cache. For both question answering tasks we use exact string match accuracy as the evaluation metric.

Summarisation is evaluated on the CNN/DailyMail dataset (See et al., 2017) using the ROUGE-L F-score (Lin, 2004) as the metric. We use the WikiText-103 dataset (Merity et al., 2016) with bits per character (BPC) for evaluating language modelling performance. We quote performance for sub-word language modelling in BPC, to account for any differences in vocabulary across models.

Finally, we construct an artificial “Text Repetition” task to evaluate the capability of the model to repeat sentences from its context verbatim. Such a task can commonly appear in a dialogue setting where the LLM agent is required to retrieve a piece of text from a possibly long context provided, and can be challenging for sparse attention techniques. We construct examples using the Tiny-Shakespeare dataset (Karpathy, 2015) by chunking the text into contexts of the appropriate size, appending them with the prompts containing a subset of the context, and evaluating the output exact character length match with the continuation from the context.

D ARITHMETIC INTENSITY

Following the framework of Section 2, we provide concrete examples of arithmetic intensity for various models and the implications for execution modern machine learning hardware. We observe from Equation (2) that the arithmetic intensity as batch size increases approaches $g + 6/\rho$. For example:

| Model | g | d_m | S | $\rho = S/(gd_m)$ | Max \mathcal{A}/\mathcal{M} |
|-------------|-----|-------|-------|-------------------|-------------------------------|
| Llama 2 7B | 1 | 4096 | 4096 | 1 | 7 |
| Llama 2 70B | 8 | 8192 | 4096 | 1/16 | 104 |
| Llama 2 70B | 8 | 8192 | 16384 | 1/4 | 32 |

Hardware Properties of selected machine learning hardware.¹ Note that $r_{\mathcal{A}}$ is the number of multiply-adds per second and $r_{\mathcal{M}}$ the number of data elements transferred per second.

| Name | Memory technology | $r_{\mathcal{A}}/10^{12}$ | $r_{\mathcal{M}}/10^{12}$ | $r_{\mathcal{A}}/r_{\mathcal{M}}$ |
|--------------------|-------------------|---------------------------|---------------------------|-----------------------------------|
| Bow IPU (FP16) | SRAM | 175 | 5.5 | 32 |
| A10 GPU (INT8) | GDDR | 125 | 0.6 | 210 |
| H100 SXM GPU (FP8) | HBM | 990 | 3.35 | 295 |

Comparing $r_{\mathcal{A}}/r_{\mathcal{M}}$ for this hardware to the arithmetic intensity achievable for standard transformer models, it’s clear that sequence generation will hit a data transfer bottleneck.

¹For IPU (Graphcore, 2023), we use the exchange memory bandwidth of 11 TB/s. A10 (NVIDIA, 2022). H100 (NVIDIA, 2023).

Table E1: *Excess* correlation ratio η (Roche et al., 1998) along axes of \mathbf{V} (excess: subtract $d^{-0.5}$, so uniform random data = 0.0). This demonstrates substantial auto-correlation along the sequence axis. Calculated for Llama 7B over 40 SQuAD examples.

| | B | S | Layer | Head | d_h |
|-------------------|-------|--------------|-------|------|-------|
| $\eta - d^{-0.5}$ | 0.143 | 0.256 | 0.0 | 0.0 | 0.0 |

E ATTENTION SPARSITY ANALYSIS

In order to understand how to approximate attention in pre-trained transformers, we analysed the queries, values and intermediate *scores* vector (softmax output). We took 40 examples from our SQuAD 1-shot task, and generated the first completion token using the dense Llama 2 7B model, capturing the \mathbf{q} vector and \mathbf{K}, \mathbf{V} matrices from every layer and attention head, showing derived statistics in Figures 2, E1 and E2 and Table E1.

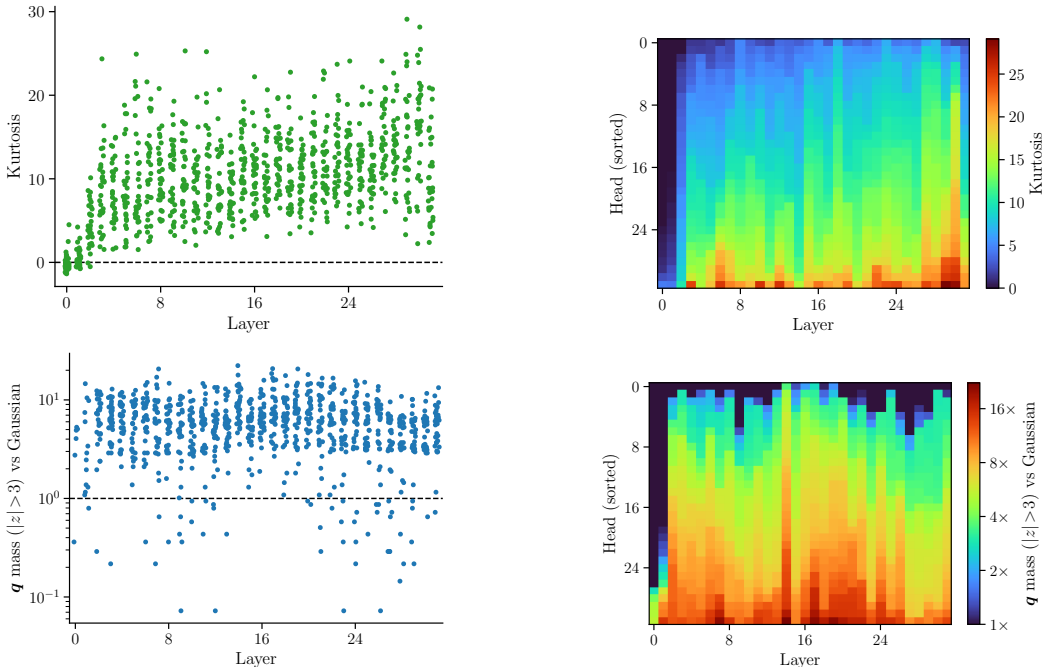


Figure E1: Statistics of components of \mathbf{q} for each head, as a function of layer. (Top) Kurtosis (Fisher), indicating that most heads have heavy-tailed \mathbf{q} . (Bottom) z -value mass, normalised by that of a Gaussian (0.3%), showing that most heads are outlier-heavy. All Llama 2 7B, measured over 40 SQuAD examples.

In Figures 2c and 2d we show that elements of the query vectors are not normally distributed, but have high sample kurtosis values. If compared to a normal distribution, the combined mass of the elements with absolute z -score exceeding 3.0 is up to $20\times$ higher. This leads us to theorise that query vectors in a pre-trained model inherently encode information **sparsely** using the tails. Therefore, the magnitude based sparsity we induce in the first stage of the algorithm does not significantly harm the approximation of the attention mappings.

We validate this claim by comparing the correspondence between the exact and approximated attention scores. SparQ Attention uses the approximate attention scores to only choose the tokens that are important for the next generation step. The actual values of the approximate scores are not relevant, as these scores are not multiplied with value vectors and thus the property of interest to us is whether the top- k indices in the approximate scores match those of the exact counterpart. This

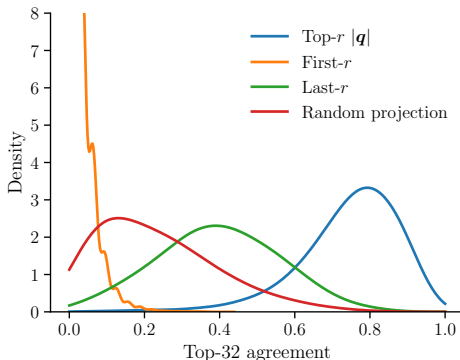


Figure E2: Top- k agreement between approximate and true scores (Llama 2 7B, measured over 40 SQuAD examples). Top- k agreement is the proportion of the top- k positions that are correctly predicted by an approximated softmax, using a projection of \mathbf{q} , either component-wise or a random low-rank projection.

can be measured on a scale from 0 to 1, where 1 means top- k indices are identical between the approximation and the exact scores and 0 means these sets do not overlap. We call this measure *top-k correspondence*. Figure E2 provides an overview how the choice of rank and k affects the top- k correspondence aggregated over all attention heads of the model. We see that the query vector sparsity of 50% and 75% maintain high top- k correspondence to the exact attention scores, which is consistently maintained over various values of k .

It is useful to drop positions in \mathbf{V} given attention scores, but this can save at most half of the data transfers, since the whole of \mathbf{K} is needed to calculate these scores. We propose approximating these scores using a subset of the components of \mathbf{K} . To test such an approximation, we measure the proportion of *overlap* between the top 32 positions in the approximated and true scores. If overlap is high, we can use the approximation to avoid transferring the whole \mathbf{K} matrix, instead only transferring some components of \mathbf{K} for all positions, then all components of \mathbf{K} for some positions.

Our hypothesis is that the r largest-magnitude components of \mathbf{q} are most useful to predicting the score, $\mathbf{q}\mathbf{K}^\top$. The coverage of this technique against an arbitrary-component baseline is shown in Figure E2. These results show that it is possible to achieve reasonably high overlap even using $r = d_h/8$, but that some later layers are harder to predict. Using the top- r components outperforms the first r baseline considerably.

F ABLATIONS AND ADDITIONAL EXPERIMENTS

Key cache compression The first step in SparQ Attention involves reading r components of the key cache to approximately determine which keys yield the highest attention scores. To examine the practical trade-off of the approximation we look at how SparQ Attention performs when compared to a theoretical upper-bounding “oracle” which provides the exact top- k keys without any data transfer. The results in Figure F1 show that SparQ Attention retains comparable performance to the oracle for a wide range of compression ratios, and attains considerably higher performance than a baseline compression scheme, in which a random low rank projection of \mathbf{K} is transferred from memory.

Approximate softmax temperature To empirically support our statistical analysis of α agreement shown in Figure 2 (Bottom Right), we evaluate a number of different viable temperature settings, including the square root of the head dimension ($\tau = \sqrt{d_h}$), the square root of the rank ($\tau = \sqrt{r}$), and our own proposed temperature, defined as $\tau = \sqrt{d_h \|\mathbf{q} \circ \mathbf{m}_q\|_1 / \|\mathbf{q}\|_1}$. We also consider the scenario where we do not reallocate mass to mean value ($\alpha = 0$), which corresponds to the limit of the temperature tending towards 0. We find that our proposed temperature performs best, as shown in Figure F2.

Hyperparameter selection The reduction of data transfer attained by SparQ Attention is controlled by its two hyperparameters, k and r . Reducing either of these variables will improve the bandwidth efficiency, but can negatively impact task performance. Figure F3 shows the relationship between k and r on both of these factors. Based on these results, we propose a simple recipe of setting $k = 128$ and tuning r to maintain a good trade-off between data transfer and task performance for a range of models and tasks.

Sequence length scaling The sequence lengths of different examples in our main tasks vary between 1k and 2k tokens, whereas many LLMs support sequence lengths far greater than this. We developed a variation of the SQuAD task that increases the task difficulty, as well as the sequence length by increasing the number of confusion contexts present in the prompt in Figure F4, which is akin to increasing the number of retrieved documents with a retrieval augmented generation system (Borgeaud et al., 2022). We test SparQ Attention and H₂O in this setting using Vicuna (Chiang et al., 2023), a descendent of Llama 2 that has been adapted and for longer sequences. Both SparQ Attention and H₂O are configured to maintain a fixed compression ratio versus the dense baseline (keeping $r = 32$ and modifying k to maintain 1/4 compression), showing that SparQ Attention is scalable to large sequences.

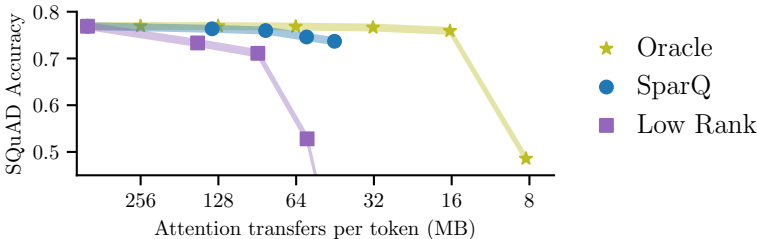


Figure F1: SQuAD 1-shot accuracy with Llama 2 7B of SparQ Attention and a random low rank compression scheme against an oracle top- k selector.

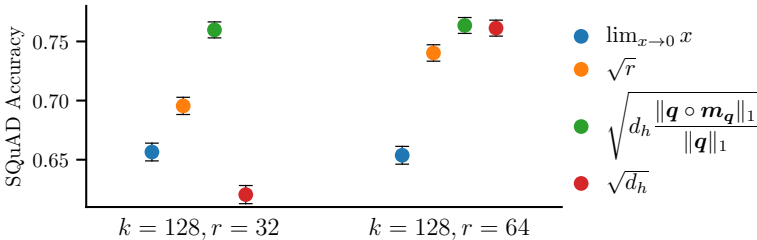


Figure F2: Comparison of different softmax temperatures for approximate attention scores for two different hyperparameter configurations (Llama 2 7B SQuAD 1-shot performance).

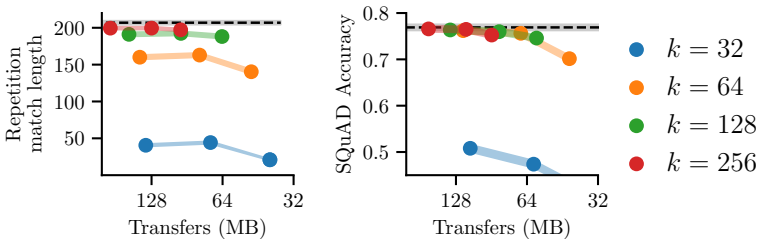


Figure F3: Results for Repetition and SQuAD tasks with $r \in \{16, 32, 64\}$.

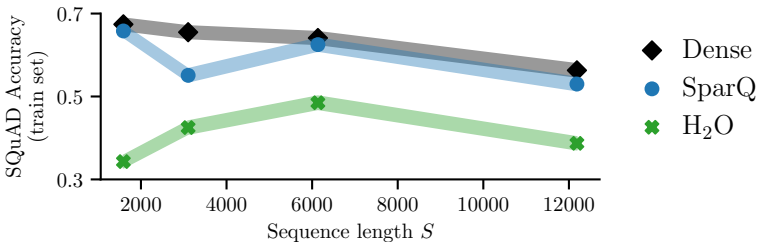


Figure F4: SQuAD performance vs input sequence length. The compression ratio is fixed at 1/4. Uses Vicuna 1.5 7B with 16k maximum sequence length against our SQuAD (train) task with 7 (default) to 63 confusion contexts to increase the sequence length.

G BENCHMARKING DETAIL

Benchmarking code is made available from:

<https://github.com/graphcore-research/llm-inference-research/tree/2024-01-paper>.

IPU measurements We tested custom fully-fused Poplar implementations of both dense attention and SparQ Attention, compiled using Poplar SDK 3.3.0+1403. On initialisation, we fill large \mathbf{K} and \mathbf{V} tensors with values $\sim N(0, 1)$ in streaming memory. On each benchmarking (outer) iteration, we first randomise the contents of a \mathbf{q} in local memory, then perform multiple inner repeats of the attention op being profiled. We use 4 inner repeats for dense attention, otherwise $1024/\text{batch_size}$, chosen because dense attention is much slower, and we swept a wide range of settings. We ran an outer loop of 2 warm-up iterations followed by 10 timed iterations, reporting the mean and standard error. The sweep covered $S \in [1024, 2048, \dots, 65536]$, $\text{batch_size} \in [1, 4, 16, 64]$, SparQ Attention $r \in [16, 32, 64]$ and $k \in [64, 128, 256, 512]$.

GPU measurements All experiments use PyTorch 2.1.2+cu121 on Ubuntu AWS instances. To set up the experiment, we initialise the large \mathbf{K} and \mathbf{V} tensors with values $\sim N(0, 1)$. On each step, we draw $\mathbf{q} \sim N(0, 1)$, run `torch.cuda.synchronize` before starting a host-side wall-clock timer, run the op, and synchronize again before stopping the timer. We run 20 warm-up iterations followed by 200 timed iterations, reporting mean and standard error. For dense baseline implementations, we tested a vanilla PyTorch implementation, with/without `torch.compile` and `torch.nn.functional.scaled_dot_product_attention`, selecting each backend (math, flash, mem_efficient) manually. For SparQ Attention implementations, we tested vanilla PyTorch (lightly hand-optimised from Appendix A.3), with/without `torch.compile`. We also toggled fused gather-matmul kernels written in Triton, and whether \mathbf{K} was stored twice in S -contiguous (for **Step 1**) and d_h -contiguous (for **Step 2**) layouts, or only once in d_h -contiguous layout. We tested $S \in [1024, 2048, 4096, 8192, 16384]$, $\text{batch_size} \in [1, 4, 16, 64]$, SparQ Attention $r \in [16, 32, 64]$ and $k \in [64, 128, 256, 512]$.

Additional results In addition to the headline results shared in Section 5 and Figure 3b, we give an aggregate picture of the trends in Figure G1. Since the number and dimension of heads is fixed, the x-axis is proportional to the size of the input tensors. On IPU (M2000), strong speedups are available across a range of input sizes, principally depending on r , but also on k (not shown). On GPU, sufficient input size is required to observe a speedup over the dense baseline, with the more bandwidth-limited A10G reaching speedups sooner. While part of this effect can be linked to the fundamental additional complexity of SparQ Attention, we anticipate that small input sizes could be accelerated considerably with additional kernel fusion. With an appropriate limit to sequence length, SparQ Attention could even be fused into a single CUDA kernel.

Storing \mathbf{K} twice One limitation of a theoretical model of data transfer is that it does not account for the granularity of memory access. Since the \mathbf{K} matrix is indexed on different axes in **Step 1** and **Step 2** of SparQ Attention, a naive implementation would fetch non-contiguous elements in one of

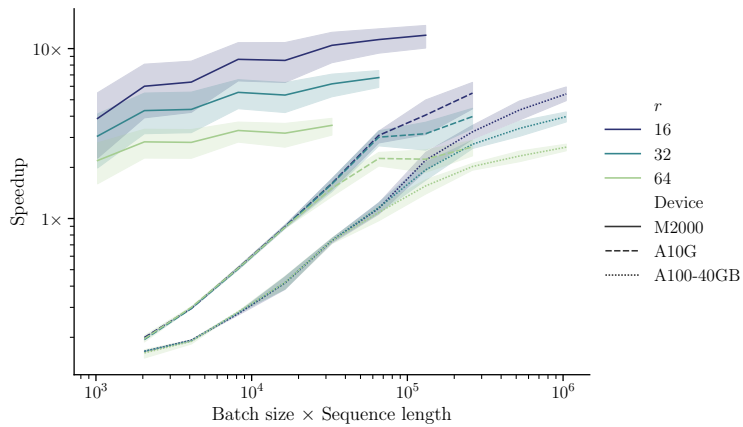


Figure G1: SparQ speedup over the dense baseline, across a range of batch size (1-64), sequence length (1024-65536) and k (64-512), for different devices. We note that for both GPUs, the number of KV elements is a limiting factor for the achieved speedup, and that this could be improved by writing a fully fused SparQ Attention kernel.

the two steps. To mitigate this, we propose storing \mathbf{K} twice, once in S -major format and once in d_h -major format. This increases KV cache memory usage by 50%, but uses only a small amount of extra bandwidth to write \mathbf{k} twice. This extra write is non-contiguous, but small, so should not form a bottleneck.

H METHODOLOGY

We provide a comprehensive set of hyperparameters for reference in Table H1. Typical compression ratios for settings of (r, k) are given in Table H2.

We use our own implementation of H₂O (Zhang et al., 2023), which differs from the authors’ implementation in that it uses a fixed cache size k , rather than a ratio of the current sequence length. To validate that these implementations are sufficiently similar, we ran their implementation through our harness on a small model and sample size. On SQuAD 1-shot, with Pythia-1.4B, using $k = 256$, $l = 64$, our implementation was correct for 60 of 200 examples, theirs for 57 (the dense baseline achieved 74). Perhaps more importantly, we found that of the 79 times that either output differed from dense, 41 occurrences showed a 20-character prefix match between our implementation and theirs. The fact that the two implementations often generate the same errors (despite minor implementation differences) reassures us that our results should be a fair representation of H₂O.

| | | |
|-----------------|----------------------|---|
| Dense model | Family | Llama 2 (13B, 7B), Mistral (7B), Pythia (6.9B, 2.8B, 1.4B) |
| | d_h | {80, 128} |
| | Max S | {2048, 4096} |
| Tasks | Question Answering | SQuAD 1-shot (4000 samples) TriviaQA 0-shot (2992 samples) |
| | Summarisation | CNN/DailyMail 0-shot (500 samples) |
| | Language Modelling | WikiText-103 LM (500 samples) |
| | Artificial | Repetition (1000 samples) |
| Baselines | Eviction | keep $(k - l)$ tokens with highest $\text{score}(n) = \sum_i s_{in}$ and the most recent $l = k/4$ $k \in \{192, 256, 384, 512, 768\}$ |
| | LM-Infinite | take the first 16 tokens, and most recent $k - 16$ $k \in \{192, 256, 384, 512, 768\}$ |
| | FlexGen | fetch the top- k values based on scores $k \in \{2, 8, 32, 128, 256\}$ |
| SparQ Attention | Rank r | {8, 16, 32, 64} |
| | Number of values k | 128 |
| | Local window l | $k/4$ |

Table H1: Experiment hyperparameters

| Method | k | r | Compression |
|------------------|-----|-----|-------------|
| SparQ Attention | 128 | 16 | 0.13 - 0.17 |
| | | 32 | 0.19 - 0.23 |
| | | 64 | 0.31 - 0.36 |
| H ₂ O | 192 | - | 0.10 - 0.17 |
| | 256 | | 0.13 - 0.22 |
| | 384 | | 0.20 - 0.33 |
| | 512 | | 0.26 - 0.43 |
| | 768 | | 0.39 - 0.65 |

Table H2: Range of compression ratios for different settings of (r, k) , for Llama 2 7B and Pythia 6.9B. The compression ratio achieved varies across models and tasks, based on the sequence length and head size.

H.1 EXAMPLES

We illustrate the task setup with a single example per task, showing the prompt formatting and a cherry-picked example. In each case, we show outputs from a dense Llama 2 13B model, SparQ Attention ($r = 8, k = 128$), H₂O and LM-Infinite ($k = 192$). Where “...” appears, we have truncated the line of text for brevity.

H.1.1 QUESTION ANSWERING (SQUAD 1-SHOT)

```
Title: University of Chicago. Background: Current ...
Title: Harvard University. Background: Harvard has...
Title: Oxygen. Background: In one experiment, Lavo...
Title: Oxygen. Background: Oxygen storage methods ...
Title: Fresno, California. Background: This vibran...
Title: Fresno, California. Background: Before Worl...
Title: Steam engine. Background: The working fluid...
Title: Sky (United Kingdom). Background: While BSk...
From what you've just read about Fresno, California, please answer the
following questions.
Question: Where is Audra McDonald from?
Answer: Fresno
Question: In what year did Roger Rocka's Dinner Theater & Good Company
Players open?
Answer:

### OUTPUT
DENSE: 1978
SPARQ: 1978
H2O: 1979
LM-INFINITE: 1975 (Roger Rock
```

H.1.2 QUESTION ANSWERING (TRIVIAQA 0-SHOT)

```
Apritifs and digestifs ( and) are drinks, typical...
Apritifs
An apéritif is an alcoholic beverage usually serve...
"Apritif" may also refer to a snack that precedes...
"Apritif" is a French word derived from the Latin...
...
...
* Distilled liquors (ouzo, tequila, whisky or akva...
* Liquor cocktails (Black Russian, Rusty Nail, etc...
In certain areas, it is not uncommon for a digesti...
Bitter digestifs typically contain carminative her...
In many countries, people drink alcoholic beverage...
Question: Which aperitif is named for the Paris chemist who created it in
1846?
Answer:

### OUTPUT
DENSE: Dubonnet
SPARQ: Dubonnet
H2O: Dubonnet
LM-INFINITE: Byrrh
```

Note that for Pythia, the prompt “Single-word answer:” was used in place of “Answer:”, as this helped prevent the model from restating the question in the answer (often qualitatively correct, but not a regex match).

H.1.3 SUMMARISATION (CNN/DAILYMAIL)

```
Article: Prince William arrived in China tonight for one of the most high
-profile...
Summary:
```

```
### OUTPUT
DENSE: Prince William arrived in China tonight for one of the most
      high-profile ...
SPARQ: Prince William arrived in China tonight for one of the most
      high-profile ...
H2O: Prince William arrived in China tonight for a high-profile
     visit that will ...
LM-INFINITE: Prince William and Kate Middleton are in Japan for a three-
            day tour. The ro...
```

H.1.4 REPETITION (SHAKESPEARE)

```
you mistake me much;
I do lament the sickness of the king.
...
...
Peace, children, peace! the king doth love you well:
Incapable and shallow innocents,
You cannot guess who caused your father's death.

Boy:
Grandam, we can; for my good uncle Gloucester
Told me, the king, provoked by the queen,
Devised impeachments to imprison him :
And when my uncle told me so, he wept,
And hugg'd me in his arm, and kindly kiss'd my cheek;
...
...
the king doth love you well:
Incapable and shallow innocents,
You cannot guess who caused your father's death.

Boy:
Grandam, we
```

```
### OUTPUT
DENSE: can; for my good uncle Gloucester
SPARQ: can; for my good uncle Gloucester
H2O: can;
LM-INFINITE: 'll not stand to prate, but to the purpose.
```

H.1.5 LANGUAGE MODELLING (WIKITEXT-103)

```
= Mellor hill fort =

Mellor hill fort is a prehistoric site in North West England , that
  dates from ...

= = Location = =

Mellor lies on the western edge of the Peak District in the Metropolitan
  Boroug...

= = Background = =

Until the 19th century little was known about hill forts ; none had been
  excava...
The study of hill forts was popular in the 19th century , with a revival
  in the...

= = History = =
```

There is evidence of human activity on the site pre @-@ dating the Iron Age , a...
A flint dagger was discovered on the site . This type of artefact is rare in Gr...
The hill fort was built in and used throughout the Iron Age , as demonstrated b...

Fragments of glass , possibly Roman in origin , and shards of pottery which date to the 1st and 2nd centuries AD , indicate the site was used in the Romano @-@ British period . However no Roman structures have been discovered , and the nature of Roman activity at the site is a source of speculation . The position of the hilltop indicate that it was easily defended ; however , local finds indicate it was a high @-@ status settlement rather than a military outpost unless a similar feature was located nearby . One reason that Roman structures have not been identified is that the Romano

DENSE: 0.669
SPARQ: 0.673
H2O: 0.685
LM-INFINITE: 0.692