This supplementary material consists of five sections. In Section A, we provide the pseudocode for Sieve. In Section B, we present the per-task performance gain of Sieve and Sieve+CLIPScore relative to CLIPScore on *large* scale. In Section C, we present Sieve's capability to augment other data pruning methods: Section C.1, investigating the effect of filtering via text-spotting on Sieve, and Section C.2 investigating the effect of distribution alignment with ImageNet on Sieve. Finally, in Section D, we show the superiority of our sentence transformer for textual semantic clustering compared to CLIP and BLIP text encoders.

## A. Pseudocode

---

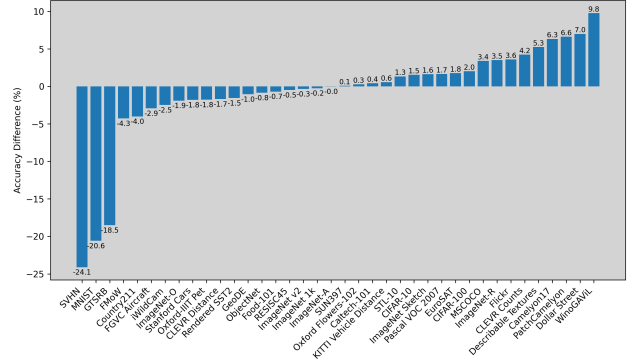**Algorithm 1** Sieve Pseudocode to filter image-text datasets.

---

**Require:** Dataset $\mathcal{D}$, Fraction $k$ fraction to prune, Caption generator $G$, Sentence transformer $S$, Set of medium phrases $\mathcal{M}$, Number of captions to generate $r$
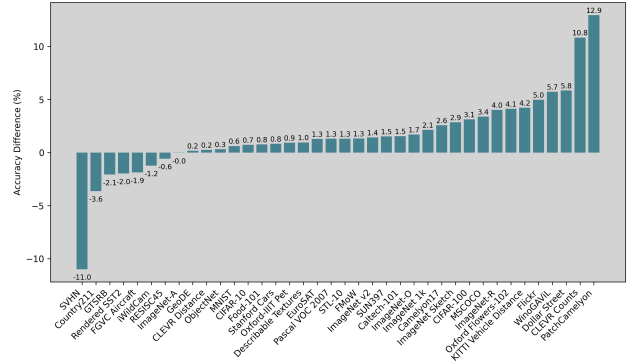
**Ensure:** Pruned dataset $\mathcal{D}_{\text{Sieve}}$

1: **for** each $I, T$ in $\mathcal{D}$ **do**   ▷ Loop over image-text pairs in the dataset
2:     $e \leftarrow S(T)$     ▷ Obtain sentence embedding of text label
3:     $\mathcal{V} \leftarrow \{\}$         ▷ Initialize empty set of scores
4:     $\mathcal{T}^G \leftarrow G(I, r)$ ▷ Generate a set, $\mathcal{T}^G$, of $r$ captions
5:     **for** each $T^G$ in $\mathcal{T}^G$ **do**       ▷ Loop over generated captions
6:         $T^{G'} \leftarrow M(T^G)$       ▷ Mask medium words
7:         $e^G \leftarrow S(T^{G'})$   ▷ Obtain sentence embeddings of generated caption
8:         $v \leftarrow \langle e^G, e \rangle$     ▷ Compute sentence similarity
9:         $\mathcal{V} + = v$             ▷ Append to set of scores
10:     **end for**
11:     $f_{\text{Sieve}}(I, T) \leftarrow \max(\mathcal{V})$   ▷ Obtain maximum score
12: **end for**
13: Rank $\mathcal{D}$ by $f_{\text{Sieve}}(I, T)$ in descending order to obtain $\text{rank}_{\text{Sieve}}(x)$
14: $\mathcal{D}_{\text{Sieve}} \leftarrow$ top $k\%$ of $\mathcal{D}$ based on $\text{rank}_{\text{Sieve}}(x)$
15: **return** $\mathcal{D}_{\text{Sieve}}$

---

## B. Per-task Performance (Large)

Figure 6 shows the change in accuracy introduced by Sieve as well as Sieve+CLIPScore on each task compared to CLIPScore on *large* scale. We observe similar patterns of relative gain for Sieve as observed on the *medium* scale, where Sieve achieves superior performance on retrieval and on medical diagnosis tasks relative to CLIPScore. In addition, when fused with CLIPScore, retrieval performance is boosted while the limited performance of Sieve on tasks that require OCR like SVHN [25] and MNIST [19] is significantly improved. Finally, Sieve with and without fusion



(a) Sieve gain over CLIPScore on *large* scale pool.



(b) Sieve+CLIPScore gain over CLIPScore on *large* scale pool.

Figure 6. The relative performance gain of Sieve and Sieve+CLIPScore relative to CLIPScore on 38 downstream tasks on "large" scale pool.

with CLIPScore achieves a large boost on DollarStreet [35], a dataset that shows pictures of household items from families of diverse ethnic and economic backgrounds. This boost demonstrates that CLIP models pretrained with Sieve filtered image-text datasets is better at interpreting a wide array of real-world scenes and objects.

## C. Combining Sieve with Other Pruning Methods

### C.1. Sieve ∩ Text Spotting

First, we study the effect of Sieve on Text-spotting. Text-spotting involves detecting and recognizing text in images and filtering image-text pairs with high overlap between spotted text (text detected in image) and alt-text (associated label of image) [30]. T-MARS [24] uses text-spotting to reduce the CLIPScore of image-text pairs with high intersection between text in images and their corresponding text. For text spotting, we first utilize a text detector, to detect and compute the percentage of image covered by text. Next, we rank samples in ascending order of the percentage of pixels covered by text. Finally, we keep the top 80% of the sam-

| Filter | ImageNet | ImageNet dist. shift | VTAB | Ret. | Avg. |
|---|---|---|---|---|---|
| CLIPScore | 27.30 | 23.00 | 33.80 | 25.10 | 32.80 |
| Sieve | **29.60** | **24.93** | **35.07** | **28.57** | **34.03** |
| CLIPScore ∩ Text Spotting | 29.75 | 24.10 | **35.65** | 24.95 | **34.05** |
| Sieve ∩ Text Spotting | **30.10** | **25.05** | 34.15 | **28.35** | 33.90 |

Table 6. Intersection of samples ranked by Sieve and samples kept by Text spotting filter on *medium* scale.

| Filter | ImageNet | ImageNet dist. shift | VTAB | Ret. | Avg. |
|---|---|---|---|---|---|
| CLIPScore | 57.8 | 47.4 | 53.8 | 46.6 | 52.9 |
| Sieve | 57.3 | 47.8 | 52.0 | 52.0 | 52.3 |
| Sieve+CLIPScore | 59.7 | 49.1 | **54.8** | 51.1 | **54.6** |
| CLIPScore ∩ ImageNet | 63.1 | 50.8 | 54.6 | 49.8 | 53.7 |
| Sieve ∩ ImageNet | 61.2 | 49.2 | 51.3 | 51.4 | 51.4 |
| Sieve+CLIPScore ∩ ImageNet | **63.8** | **51.4** | 53.1 | **53.3** | 53.6 |

Table 7. Intersection of ImageNet-based filtering on *large* scale. We achieve state-of-the-art performance on ImageNet, ImageNet out-of-distribution and retrieval tasks.

| Top-$k$% | Sieve | Sieve + CLIPScore |
|---|---|---|
| 10 | 16.00% | 27.18% |
| 20 | 29.99% | 44.00% |
| 30 | 39.54% | 56.93% |

Table 8. Intersection-Over-Union between the unique ids of top-$k$ samples selected by CLIPScore and 1) Sieve, 2) Sieve+CLIPScore on *medium* scale

ples and intersect either with top 30% CLIPScore or top 20% top Sieve samples. In Table 6, we observe that CLIPScore filtering gains from Text-spotting more than Sieve. In addition, the performance of CLIPScore ∩ Text-spotting is close to Sieve on its own. It is important to note that we start with top 20% of the samples (24M) for Sieve, we end up with fewer samples for pretraining with Sieve (19.2M) compared to top 30% from CLIPScore (24M). For future work, we are planning to implement text-spotting tightly coupled with Sieve by masking detected text in images prior to the captioning step.

## C.2. Sieve ∩ ImageNet

One of the common approaches to pruning large-scale noisy image-text datasets is by sampling image-text pairs where the image is clustered near pretrained CLIP embeddings of ImageNet training set samples [9]. This distribution alignment can result in improved visual representations at the expense of generalization. We generate a pretraining dataset using Sieve and then investigate the utility of only keep-

ing Sieve samples intersecting with ImageNet samples or Sieve+CLIPScore (fused) samples intersecting with ImageNet samples. Looking at Table 7, compared to Sieve, we observe that the intersection of Sieve and ImageNet leads to significant improvements on ImageNet but at the expense of average performance. We also observe that compared to CLIPScore ∩ ImageNet, Sieve+CLIPScore ∩ ImageNet leads to the best performance on ImageNet, ImageNet out-of-distribution shift and a significant improvement on retrieval tasks (+3.5%).

## D. Textual Semantic Clustering

We investigate the textual semantic clustering abilities of three models; Sentence Transformer, CLIP text encoder and BLIP text encoder. We create 3 groups of sentences, where each group contains 3 sentences describing similar objects. The goal is to study which text encoder is best at distinguishing between sentences from the same group. The best embeddings space for estimating the alignment between generated caption and alt-text, should have a representation space, where semantically similar sentences are close to each other, while semantically-distinct sentences are far away. In Figure 7, we depict the cosine similarity values between the embeddings of all sentences from the 3 groups, where sentences from the same group have the same color. We observe that compared to CLIP and BLIP text embedding space, Sentence transformer has the best intra and inter-cluster values and therefore results in the best alignment score.

| Scale | Filtering | Dataset Size | ImageNet | ImageNet dist. shifts | VTAB | Retrieval | Average over 38 datasets |
|---|---|---|---|---|---|---|---|
| | No Filtering | 12.8M | 2.5 | 3.3 | 14.5 | 11.4 | 13.2 |
| | Basic Filtering | 3.0M | 3.8 | 4.3 | 15.0 | 11.8 | 14.2 |
| Small | LAION Filtering | 1.3M | 3.1 | 4.0 | 13.6 | 9.2 | 13.3 |
| (12.8 Million) | CLIPScore | 3.8M | 5.1 | 5.5 | 19.0 | 11.9 | 17.3 |
| | Sieve | 3.0M | 6.0 | 6.2 | 18.2 | 13.6 | <u>17.4</u> |
| | Sieve+CLIPScore | 2.4M | 6.0 | 6.0 | 19.8 | 13.3 | **18.0** |

Table 9. Zero-shot performance of CLIP models pretrained using filtering strategies on *small* scale pools of the DataComp benchmark.
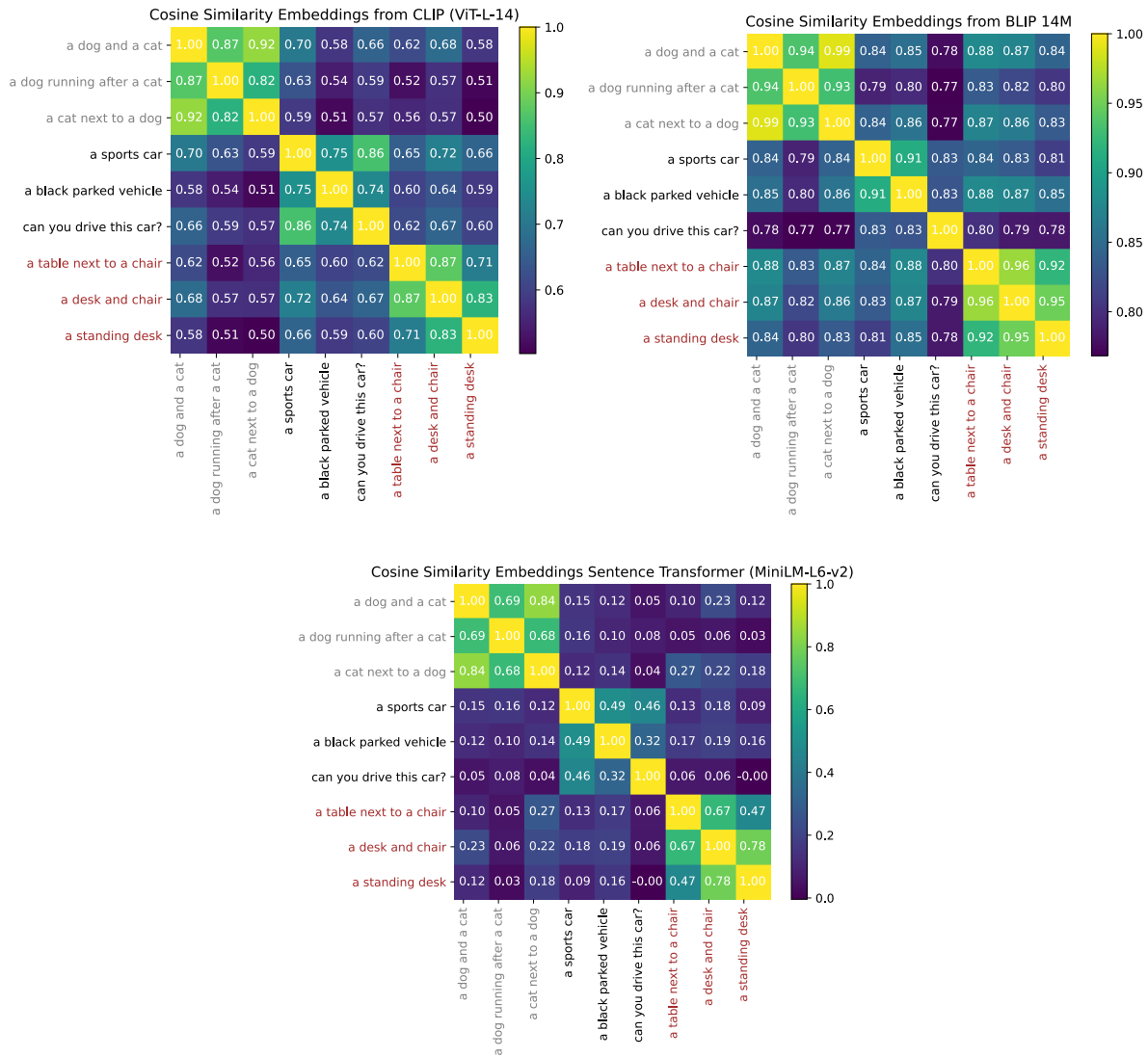


Figure 7. Confusion matrix of cosine similarity illustrating the performance of each text encoder in similarity. We show 9 sentences split into 3 groups of consecutive sentences, where sentences within each group describe similar concepts. We observe that Sentence Transformer has better semantic textual clustering compared to CLIP and BLIP text encoders.