

Fig. 2: Performance of different sharing models

Fig. 2 shows the time to create containers and their size with these policies over the multi-year trace. The figure shows that reuse can improve performance over the baseline by $\sim 5x$. If we impose a constraint on cache capacity (in this case 20 containers) we see improvement of $\sim 2x$.

V. CACHING METRICS

We explore five metrics for each package. **Versions:** Average time between versions. **Popularity:** Average number of stars and forks. **Size:** Total size on disk. **Time:** Total time to install. **Dynamic Count:** An online approach that counts the number of package invocations within a sliding window.

We evaluate caching strategies that order containers based on each metric. Table I shows the average performance of each strategy using our simulator and Binder trace. We compare against a baseline least recently used (LRU) implementation. Our results suggest that individually, LRU is best.

Metric	Size (TB)	Time (Hrs)	Hit rate (%)
LRU	2.71	115.39	84.61
Dynamic	5.11	128.05	83.39
Size	4.90	123.25	83.47
Popularity	5.68	130.11	83.11
Time	4.75	119.37	83.84
Versions	5.89	134.26	82.75

TABLE I: Average cache strategy performance with cache size limits of [2000, 4000, 6000, 8000, 10000] megabytes.

VI. WEIGHTED COMBINATION: CACHERANK

Without an obvious “best” metric, we combined metrics into a single ranking metric (CacheRank), implemented as follows.

```

for met in metrics:
    rk = Rank(met)
    score += met_coeff * (rk - mean(rk)) / std(rk)

```

This approach aims to normalize the different metrics’ scores, efficiently account for overlap, and utilizes input weightings for the factors. We use random values from a Gaussian distribution as parameters for our simulations. Fig. 3 shows the best performance at every cache size limit and the improvement over LRU in terms of size, time, and hit rate.

VII. MRU PROTECTION

CacheRank weights the ranking of each container uniformly throughout the cache. However, this is likely to be flawed as the containers at the front of the cache (most recently

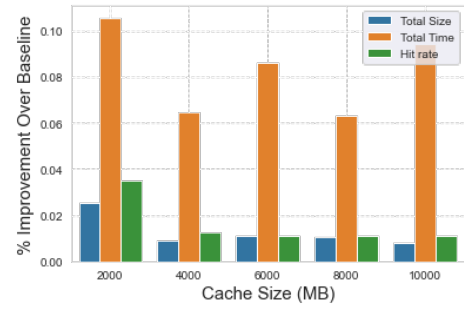


Fig. 3: Improvement over LRU for different cache sizes.

used, MRU) are more impactful proportionally to the least recently used containers. To integrate this observation we explored protecting the MRU containers in the cache. In our simulations, we set the MRU protection parameter *cache_safe* to be the 20th percentile for the cache limit being evaluated. We can see in Fig. 4 that this strategy provides the best performance as the most optimal value of *cache_safe* lies somewhere between 20% and 60% of the capacity.

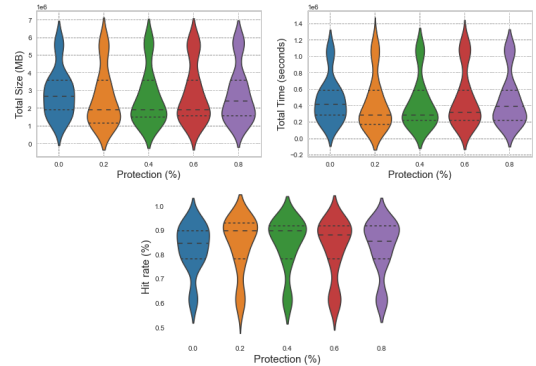


Fig. 4: Distribution and quartiles for protection fractions

With the unprotected sector of the cache, we tested different heuristics to prioritize the removal of containers, using: the metric score as a sole basis for decision, the smallest total size based on a minimum Knapsack algorithm, and a linear combination of both score and size. We found the linear combination outperformed the other methods.

VIII. SUMMARY

We explored container sharing and caching strategies to improve the performance of multi-user computing services. We collected features such as installed package size, installation time, popularity metrics, and version history to create metrics that could be combined with an LRU cache. Our strategies, CacheRank and MRU protection, use these metrics to improve cache performance under space constraints.

REFERENCES

- [1] Binder, <https://mybinder.org>, 2021
- [2] T. Shaffer, K. Chard, D. Thain, “An Empirical Study of Package Dependencies and Lifetimes in Binder Python Containers”. eScience, 2021.
- [3] T. Shaffer et al., “Solving the Container Explosion Problem for Distributed High Throughput Computing.” IPDPS, 2020.