

APPROXIMATE CACHING FOR EFFICIENTLY SERVING TEXT-TO-IMAGE DIFFUSION MODELS

Shubham Agarwal¹, Subrata Mitra^{*1}, Sarthak Chakraborty²,
Srikrishna Karanam¹, Koyel Mukherjee¹, Shiv K. Saini¹

Adobe Research¹, UIUC²



Shubham Agarwal

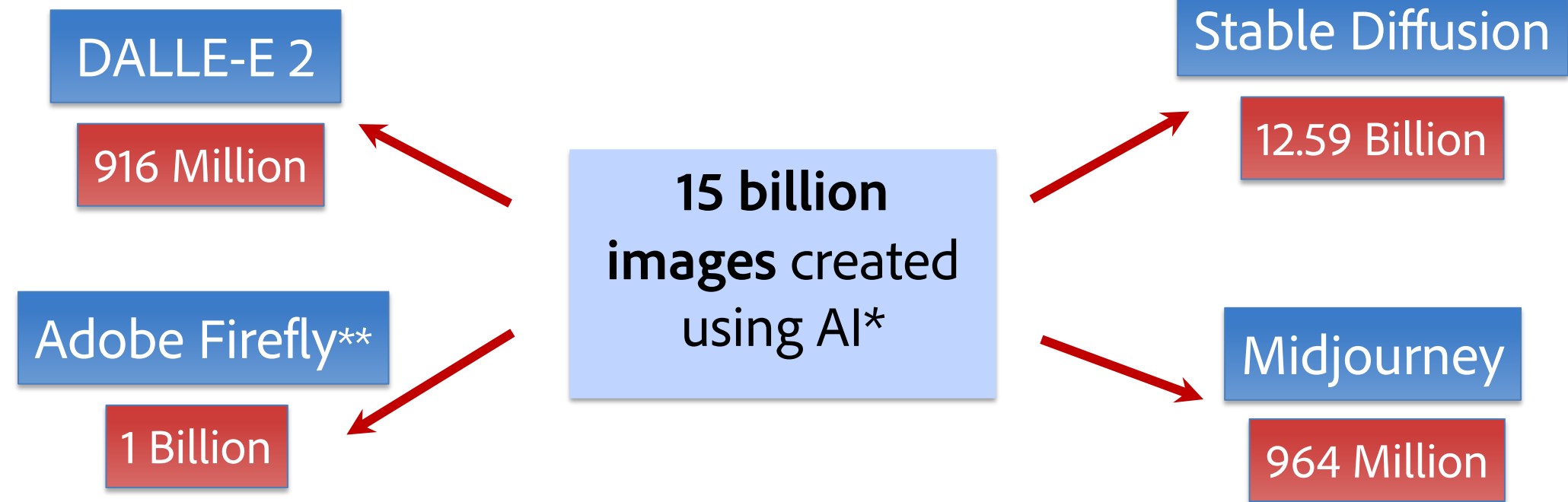
Research Associate, Adobe Research

Networked Systems Design and Implementation
(NSDI 2024)

Popularity of Text-to-Image



AI Has Already Created As Many Images As Photographers Have Taken in 150 Years. Statistics for 2023*



Text Art



Cartoons



Abstract Arts



Flyers and templates

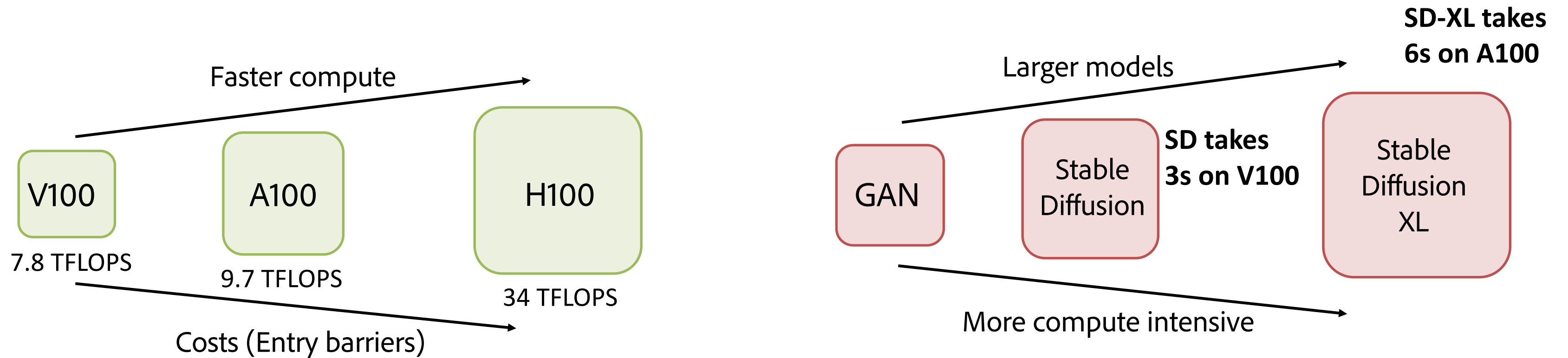
* As of Aug 2023, [<https://journal.everypixel.com/ai-image-statistics>]

** To date, Firefly generated over 6.5 billion images and counting!

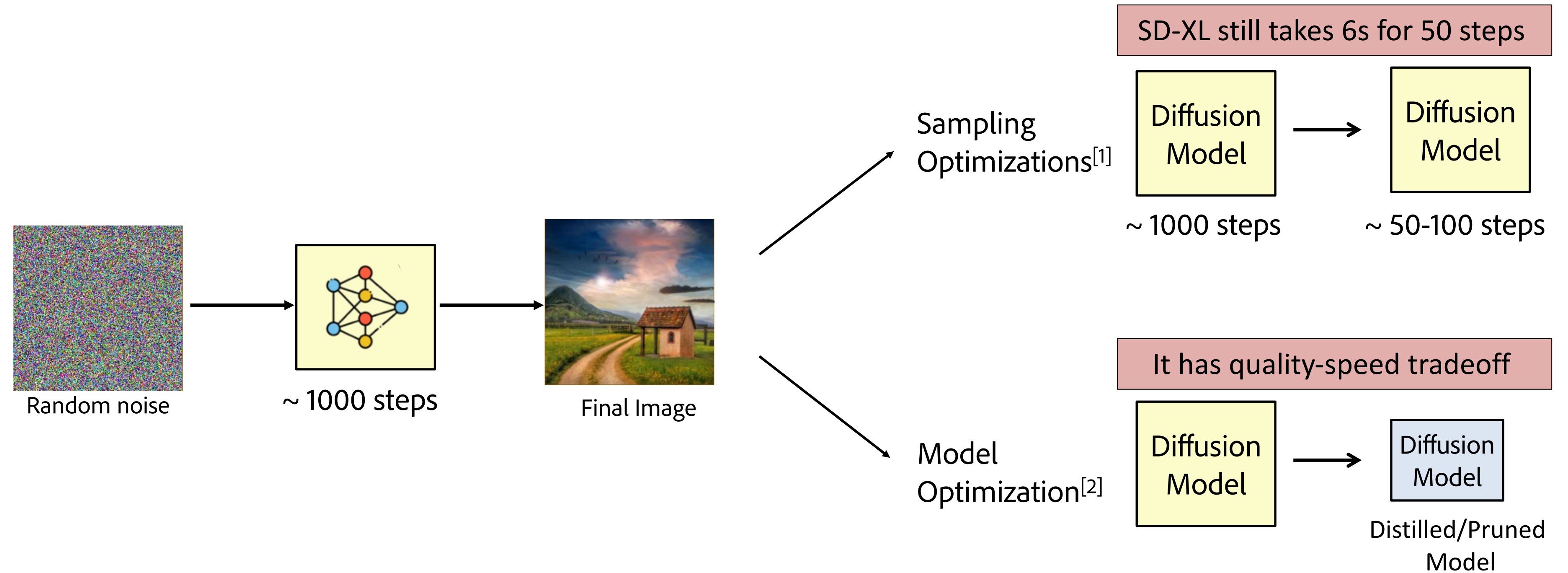
Soaring Costs with Popularity

Evolution of Compute and Model Complexity

Increasing Model Complexity can outpace the Hardware Evolution



Efficiency of Diffusion models

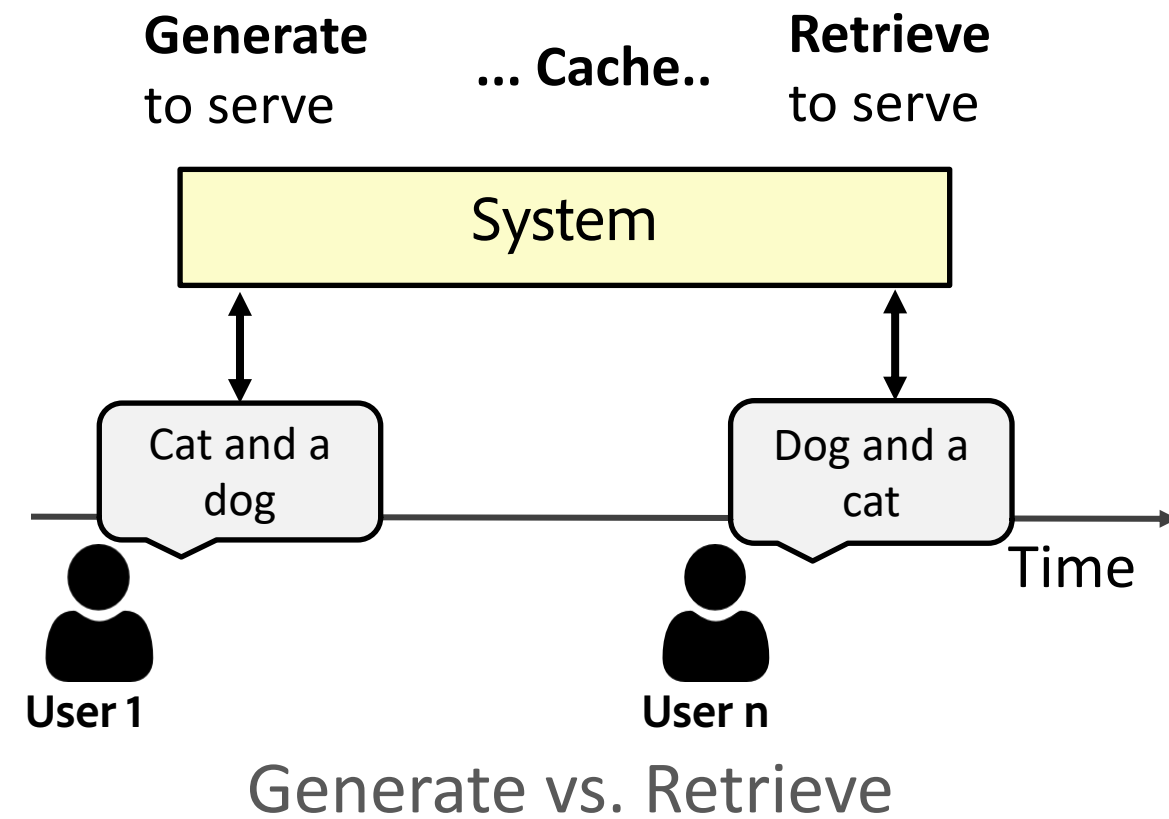


[1] Hongkai Zheng, Weili Nie, Arash Vahdat, Kamyar Azizzadenesheli, and Anima Anandkumar. Fast sampling of diffusion models via operator learning. In ICML 2023.

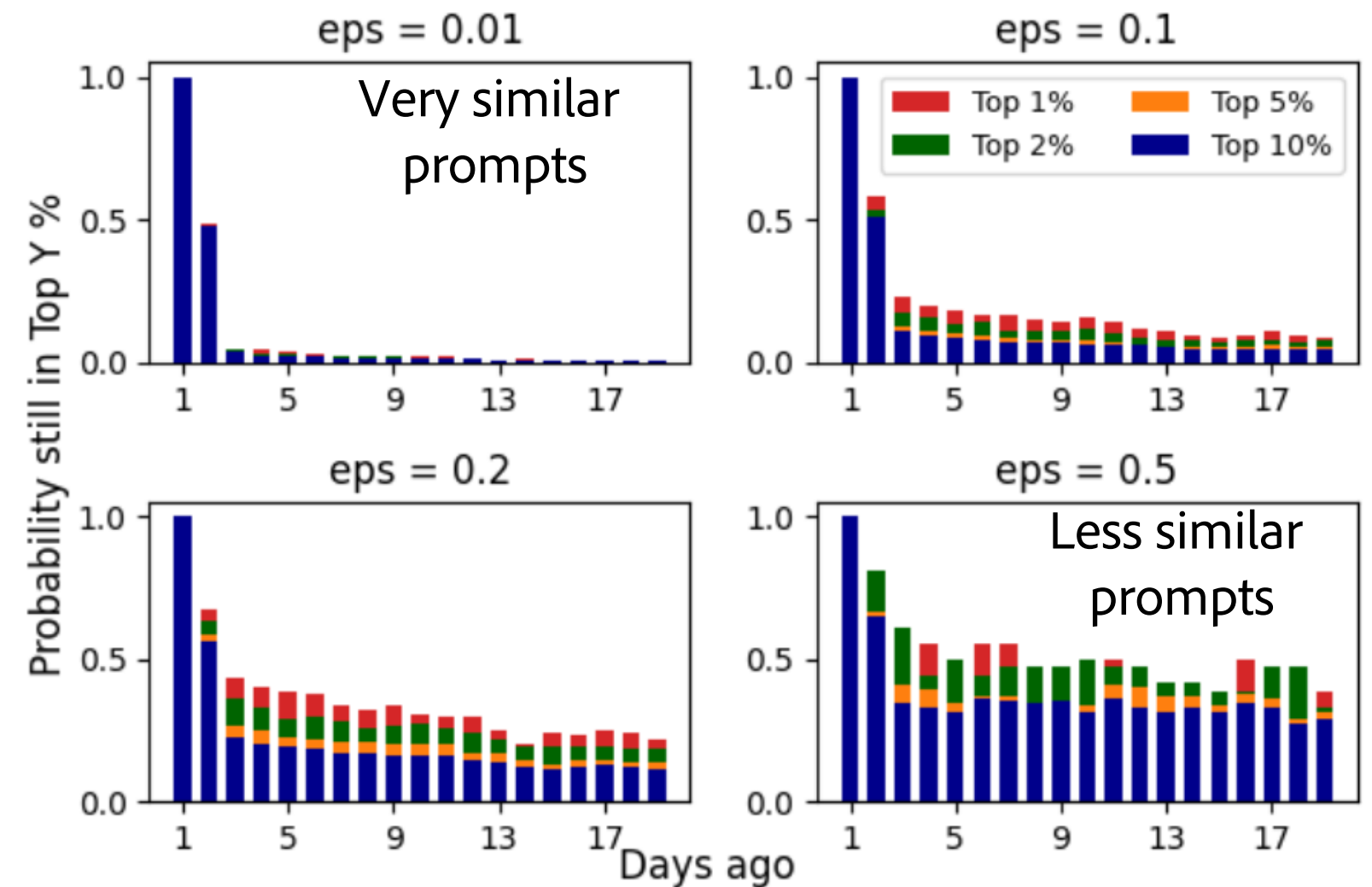
[2] Chenlin Meng, Robin Rombach, Ruiqi Gao, Diederik Kingma, Stefano Ermon, Jonathan Ho, and Tim Salimans. On distillation of guided diffusion models. In CVPR, 2023.

Proposed Caching Technique

Previous works have overlooked the use of text-to-image systems across multiple generations

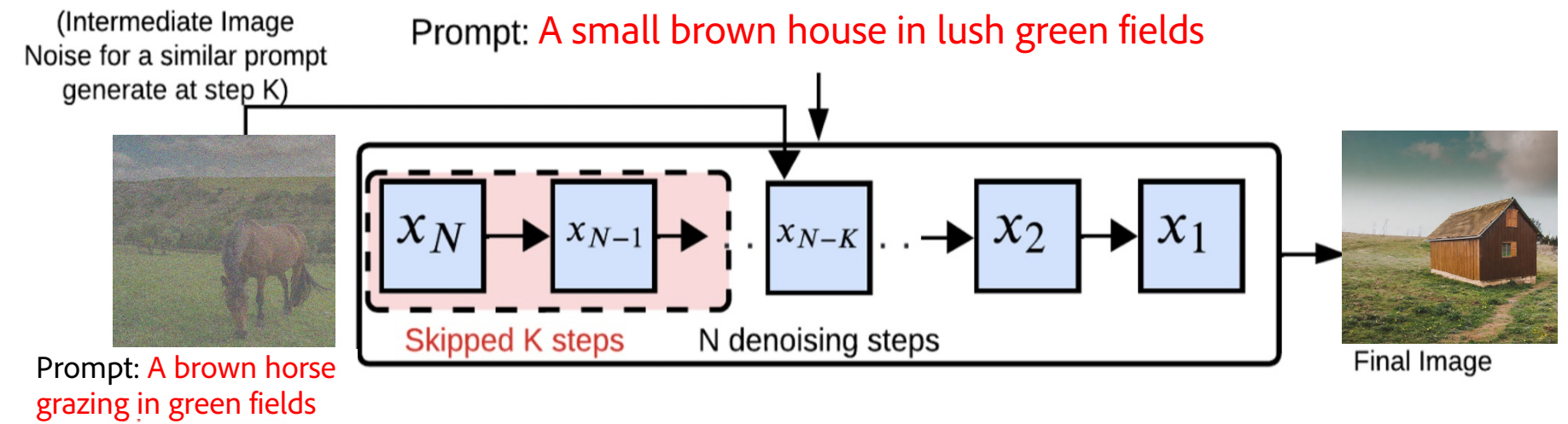
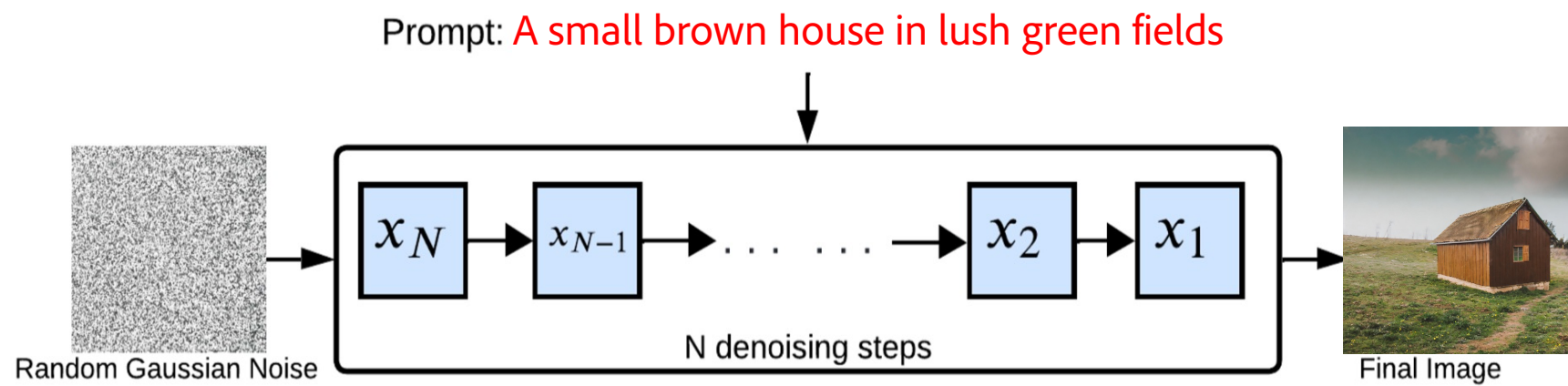


In this work, we reduce the generation time by using a simple-yet-novel approach of **caching**

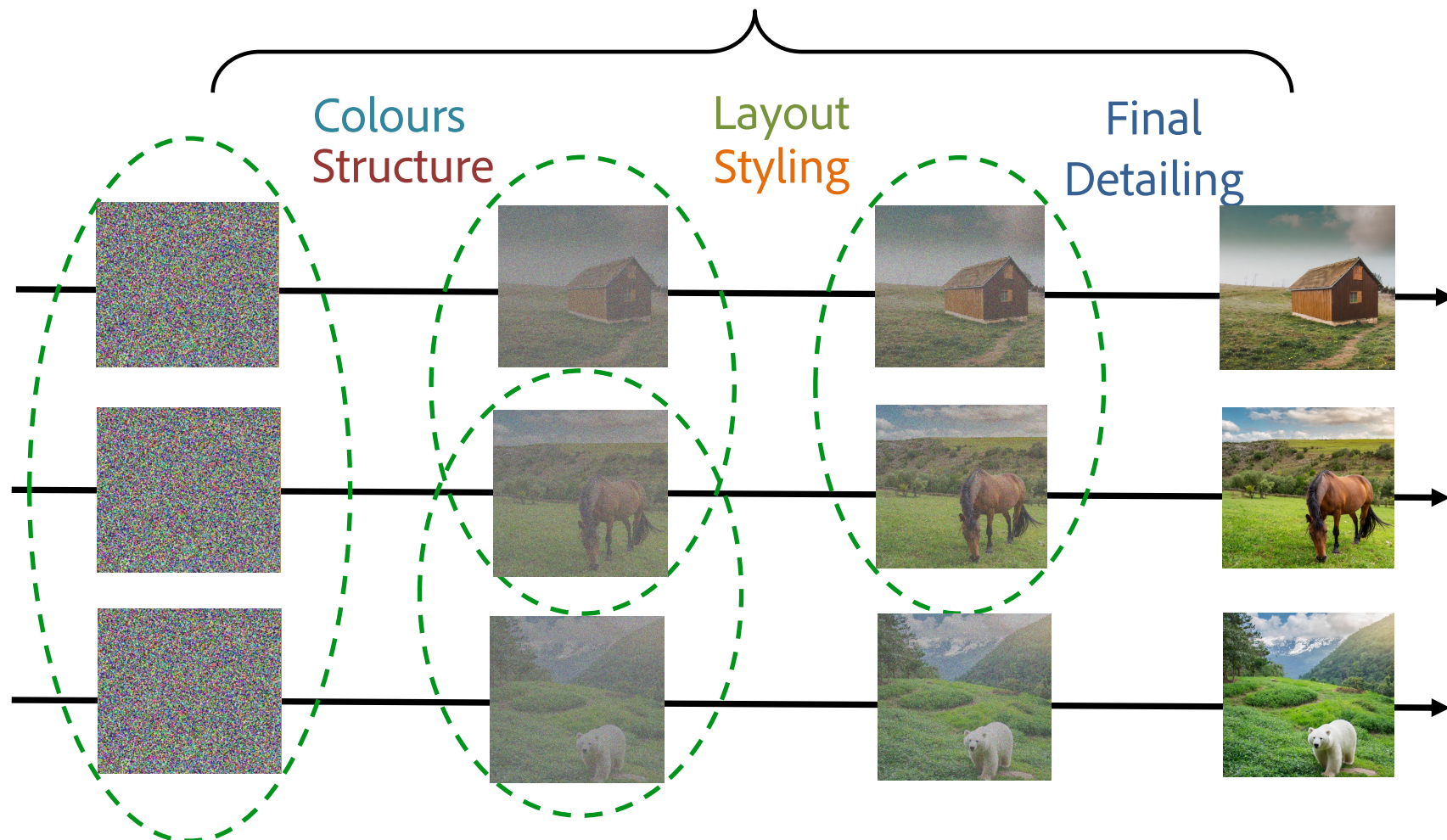


Top Prompt Clusters' Popularity Over Days

Proposed Approximate Caching

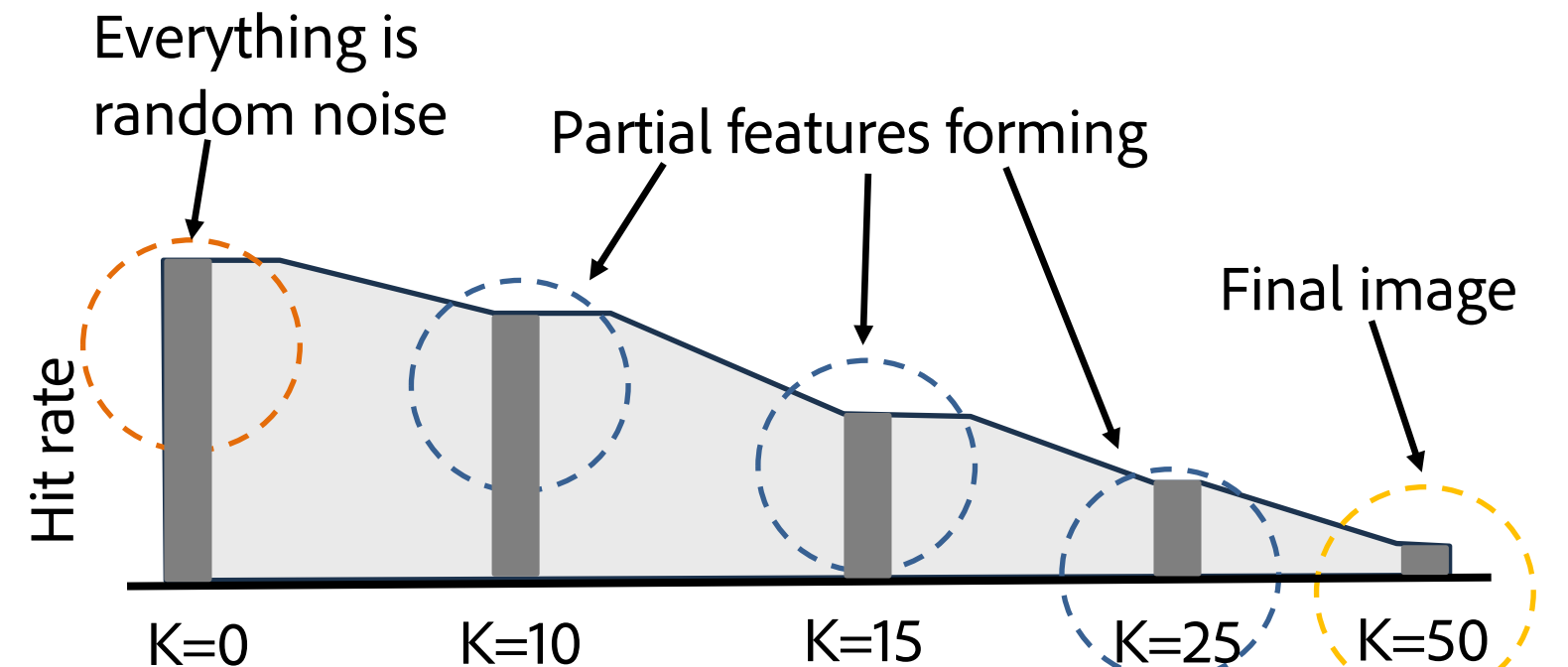


Different aspects of image forms across different steps



Opportunity: Cache and generate from intermediate step


$$\text{Maximize compute savings } f_C = h(K) \cdot \frac{K}{N}$$











Hit rate is different for different steps of generation

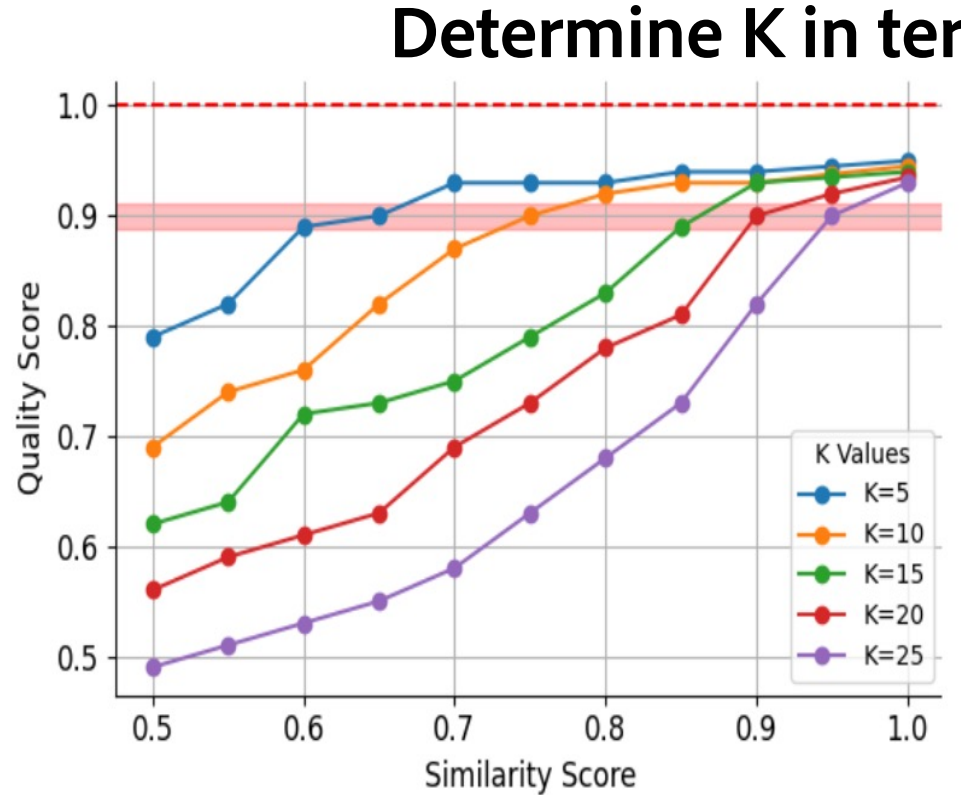
Cache Selection

Cache Selector Heuristic Hypothesis – More steps can be skipped for a prompt if its cache is more similar

Retrieved prompt:  A brown horse grazing in a green field

<div style="color: green; font-size: 2em;">✓</div> <div style="color: red; font-size: 2em;">✗</div>	Worked	<p>"A white horse grazing in a green field"</p>  <p>K = 10</p>	<p>"A brown horse running in a green field"</p>  <p>K = 10</p>	<p>"A brown bear grazing in a green field"</p>  <p>K = 20</p>	<p>"A brown horse grazing near a river"</p>  <p>K = 20</p>
	Did not work	 <p>K = 15</p>	 <p>K = 15</p>	 <p>K = 25</p>	 <p>K = 25</p>

K is Determined by Prompt-Cache Similarity:
 Noise from a Brown Horse Transforms into Different Prompts, Limited by K



- Generate Images for queries at different K values
- Profile Image Quality at different K values for different Prompt-Cache similarity levels

Example output →

```

1 def cache_selector(s):
2   if s > 0.95: k = 25
3   elif s > 0.9: k = 20
4   elif s > 0.85: k =
5     15
6   elif s > 0.75: k =
7     10
8   elif s > 0.65: k = 5
9   else: k = 0
10  return k
  
```

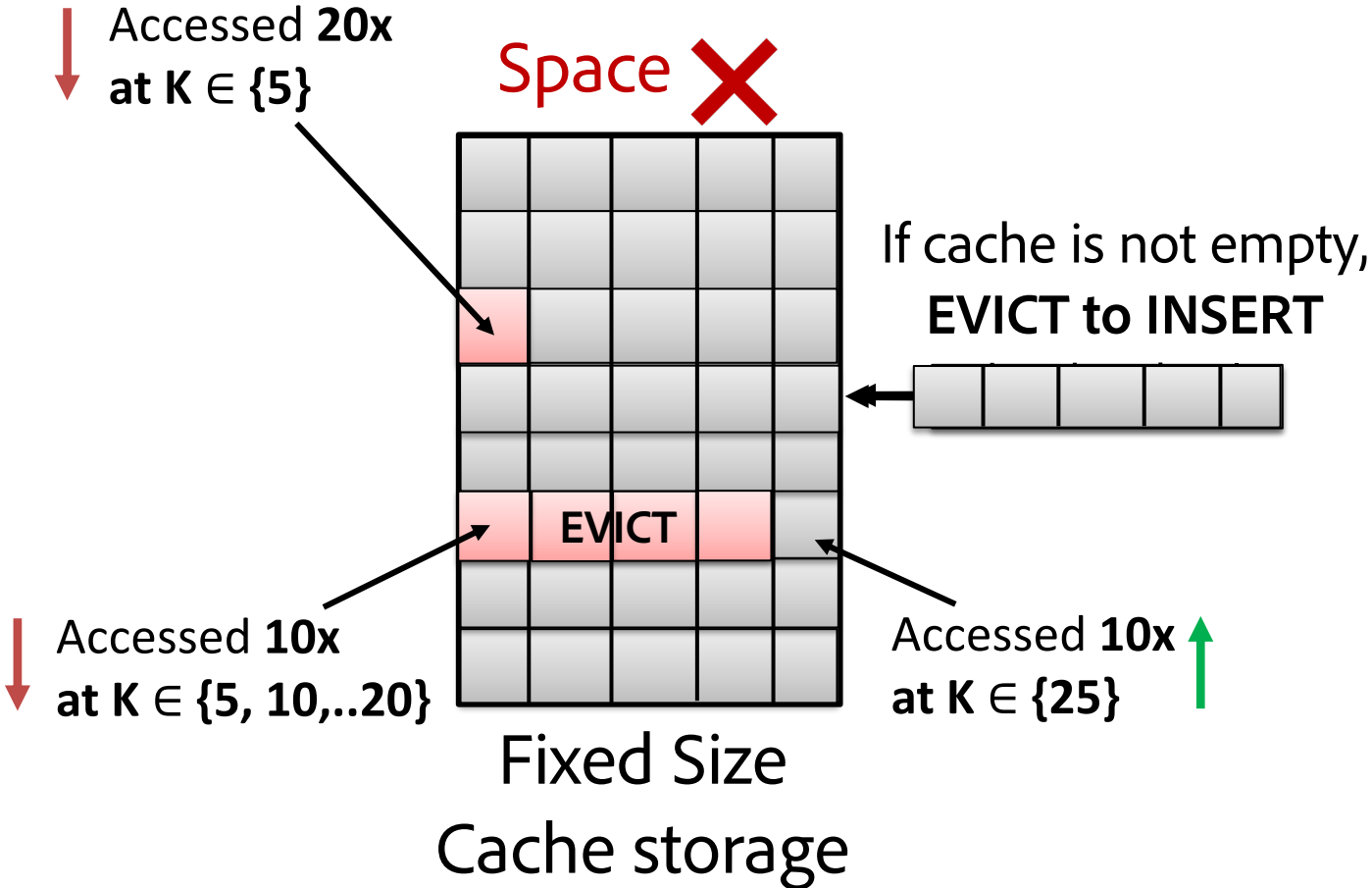
- Map from similarity score to the max K that can be skipped

Cache Management

LCBFU Policy

How to maintain the cache for storing noises ?

Least Computationally Beneficial Frequently Used



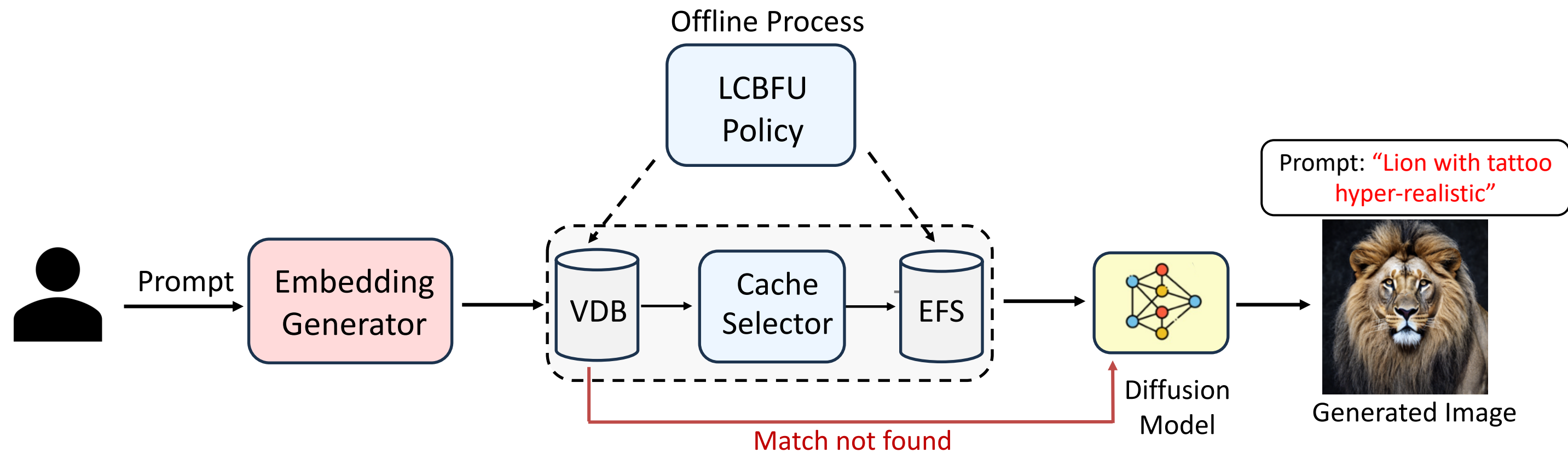
Cache Size(GB)	#noises in cache	FIFO	LRU	LFU	LCBFU
1GB	1500	0.11	0.12	0.12	0.12
10GB	15000	0.13	0.14	0.14	0.15
100GB	150000	0.14	0.16	0.16	0.18
1000GB	1500000	0.17	0.20	0.19	0.23

Compute Savings facilitated by LCBFU compared to other eviction techniques

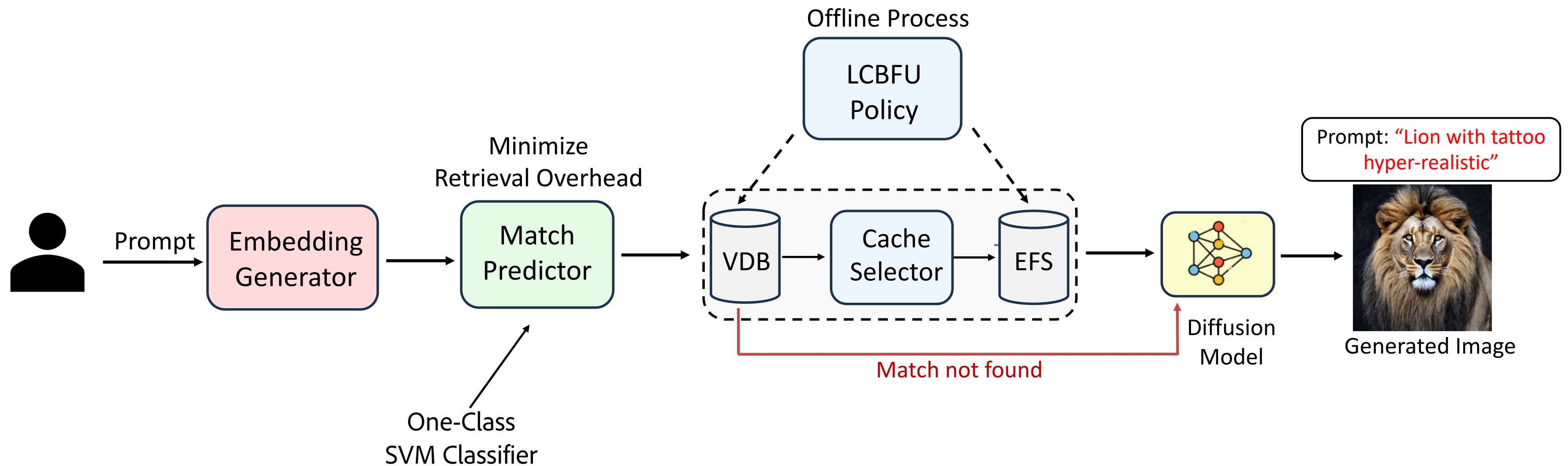
Eviction Policy?

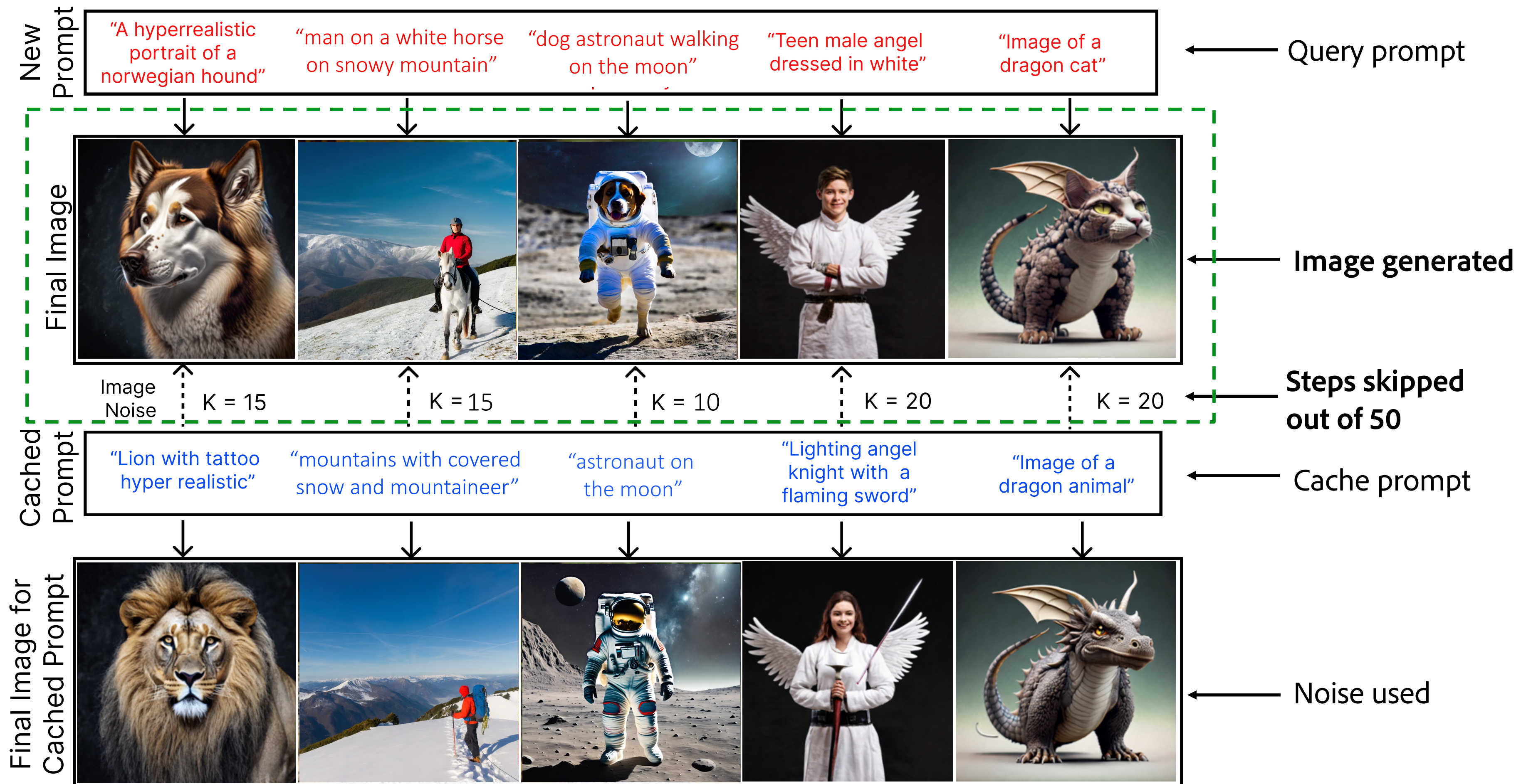
Least Computationally Beneficial Cache (LCBFU)
 Score = K (potential savings) * f (frequency of use)

NIRVANA: Proposed pipeline



NIRVANA: Proposed pipeline



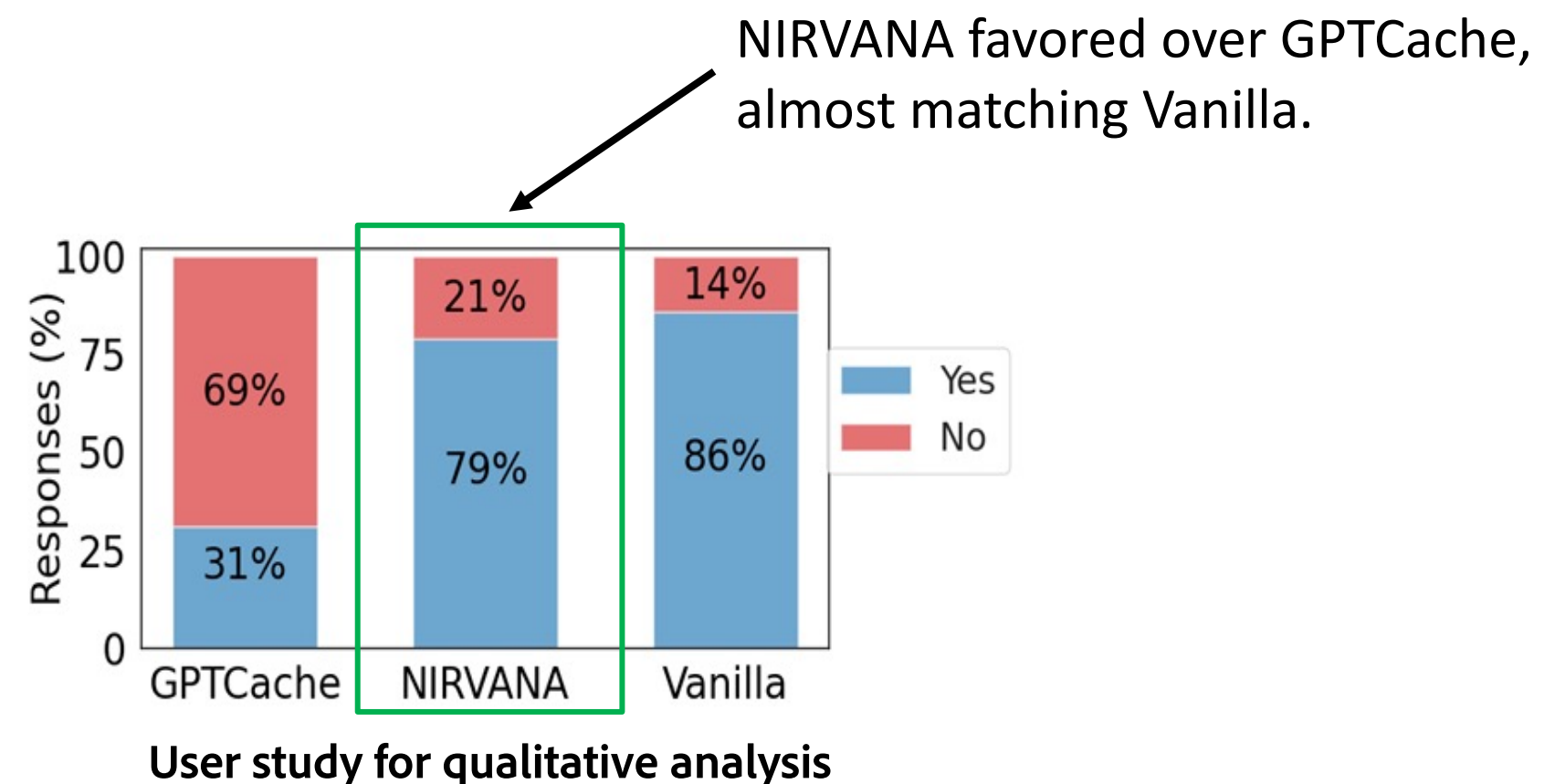


NIRVANA: Results

- We evaluate Nirvana quality using FID, CLIP and PickScore.
- We use real prompt traces from production.
- We conduct a user study with 60 participants.

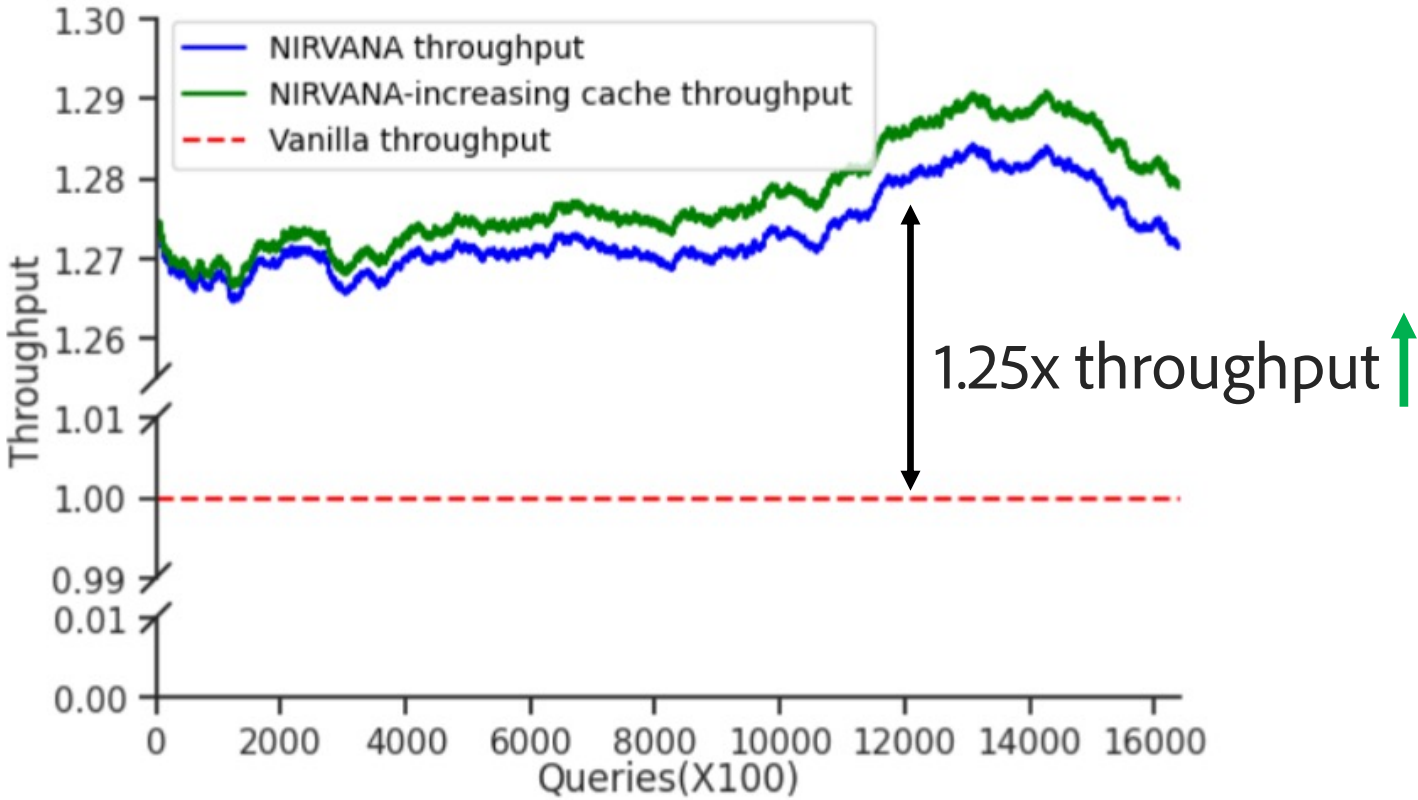
Dataset	Models	Quality		
		FID ↓	CLIP Score ↑	PickScore ↑
DiffusionDB*	GPT-CACHE	7.98	25.84	19.04
	PINECONE	10.92	24.83	18.92
	CRS	8.43	24.05	18.84
	SMALLMODEL	11.14	25.64	18.65
	NIRVANA – w/oMP	4.94	28.65	20.35
	NIRVANA	4.68	28.81	20.41
	VANILLA	6.12-6.92	30.28	20.86
SYSTEM-X	GPT-CACHE	8.15	26.32	19.11
	PINECONE	10.12	24.43	18.83
	CRS	8.38	23.81	18.78
	SMALLMODEL	11.35	25.91	18.92
	NIRVANA – w/oMP	4.48	28.94	20.31
	NIRVANA	4.15	29.12	20.38
	VANILLA	5.42-6.12	30.4	20.71

Comparison of NIRVANA against retrieval-based baselines.
We also compare against a small distilled model.

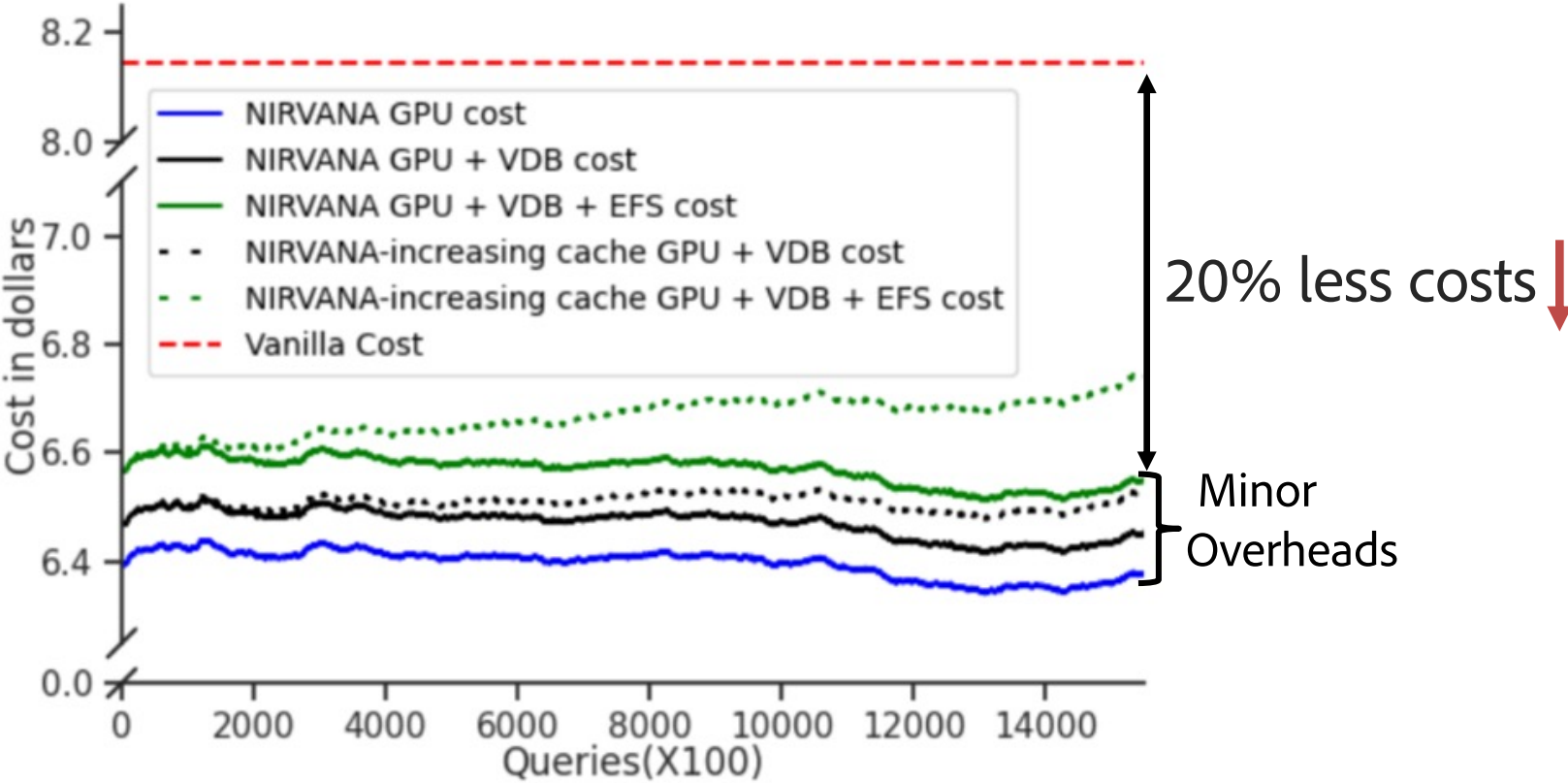


NIRVANA: Results

- We assess the end-to-end speedup and cost reductions realized by NIRVANA.

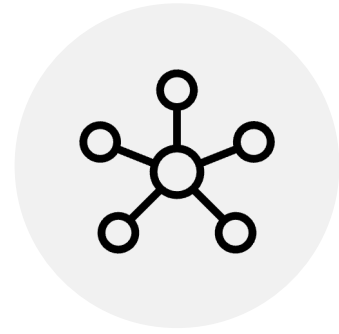


Throughput attained by NIRVANA w.r.t. standard deployment

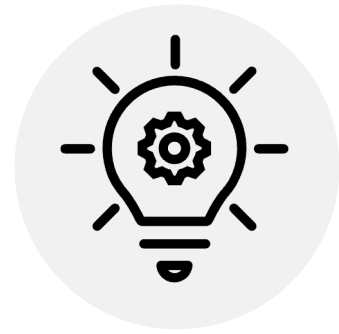


Cost analysis of deploying NIRVANA w.r.t. standard deployment

Analysis and Conclusion



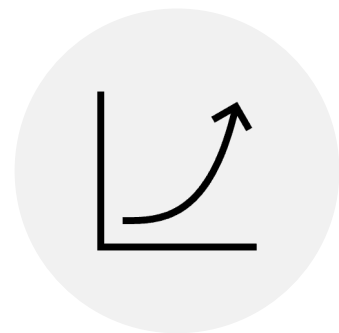
We introduce **NIRVANA: An Approximate Caching Technique** for Diffusion models.



The idea is to **cache and reuse intermediate states** from previous text prompts.



We introduces a **new eviction technique *LCBFU*** for cache management.

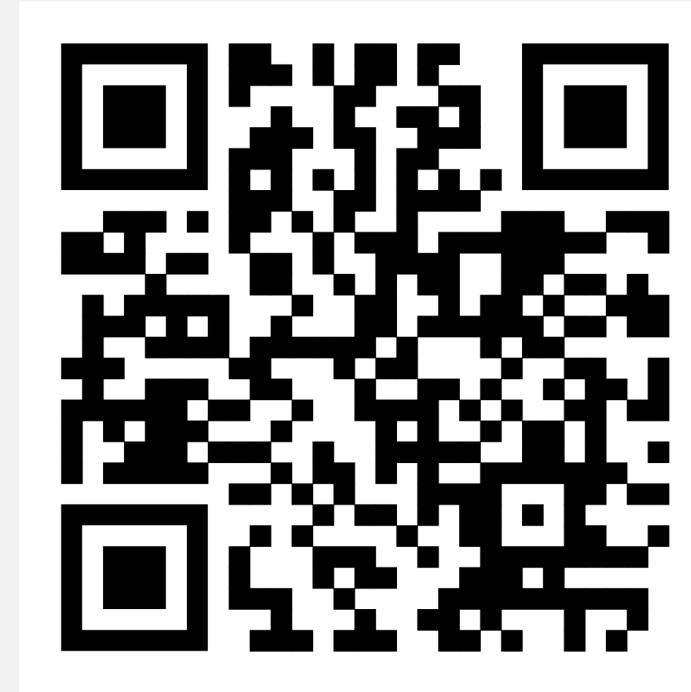


We show the effectiveness of NIRVANA using **two real prompt logs***.

Thank you



NSDI Paper link



Adobe Research

Contact: shagarw@adobe.com