

## Supplementary Materials

### A. Proof of proposition 1

$$\begin{aligned}
\alpha_{p,q} &= \mathbb{E}_{q(x)} \left[ \min \left( 1, \frac{p(x)}{q(x)} \right) \right] \\
&= \sum_x q(x) \min \left( 1, \frac{p(x)}{q(x)} \right) \\
&= \sum_x \begin{cases} p(x), & \text{if } p(x) < q(x) \\ q(x), & \text{otherwise} \end{cases} \\
&= \sum_x \min(p(x), q(x)) \\
&= \sum_x \frac{p(x) + q(x) - |p(x) - q(x)|}{2} \\
&= \frac{1}{2} \left( \sum_x p(x) + \sum_x q(x) - \sum_x |p(x) - q(x)| \right) \\
&= \frac{1}{2} \left( 1 + 1 - \sum_x |p(x) - q(x)| \right) \\
&= 1 - \frac{1}{2} \sum_x |p(x) - q(x)| \\
&= 1 - TV(p, q)
\end{aligned}$$

### B. Full algorithm of GSD

In this section, we provide detailed pseudocode for the implementation of GSD in Algorithm 1 and 2. Since GSD is built upon SJD, it operates by simply replacing the `VERIFY(·)` part with our `VERIFY_GSD(·)`. For a more detailed implementation, please refer to the attached code.

---

#### Algorithm 1 Grouped Speculative Decoding

---

**Require:** Speculative Length  $L$ , maximum seq length  $N$ , expert model  $p_\theta$ , initial context  $X_{0:n_0}$ , Group size  $G$ , Embedding distance matrix  $M_d$ , thresholds  $d, \delta$

- 1:  $k \leftarrow L, n \leftarrow n_0$
- 2: **while**  $n < N$  **do**
- 3:  $q_{L-k:L}, \hat{X}_{L-k:L} \sim \text{Rand-init}(\cdot)$
- 4: **parallel for**  $j = 0$  to  $L$  ▷ Parallel Verify
- 5:  $p_j \leftarrow p_\theta(\cdot \mid [X_{0:n}, \hat{X}_{0:j}])$
- 6: **end for**
- 7:  $(\hat{X}_{0:k}, k) \leftarrow \text{VERIFY\_GSD}(\hat{X}_{0:L}, p_{0:L}, q_{0:L}, G)$
- 8:  $X_{n:n+k-1} \leftarrow \hat{X}_{0:k}$  ▷ Accept.
- 9:  $q_{0:L-k} \leftarrow p_{k:L}, \hat{X}_{0:L-k} \leftarrow \hat{X}_{k:L}$  ▷ Draft update
- 10:  $n \leftarrow n + k$
- 11: **end while**
- 12: **return**  $X$

---



---

#### Algorithm 2 VERIFY\_GSD( $X, p, q, G$ )

---

**Require:** Draft  $\hat{X}_{0:L}$ , Verifier :  $p_{0:L}(\cdot)$ , Drafter :  $q_{0:L}(\cdot)$ , Group size  $G$ , Embedding distance matrix  $M_d$ , thresholds  $d, \delta$

- 1: **for**  $k = 0$  to  $L$  **do**
- 2:  $p\_sort_{vals}, p\_sort_{idx} \leftarrow \text{sort}(p_k)$
- 3:  $idx \leftarrow \text{find-idx}(p\_sort_{vals}, p_k(\hat{X}_k))$
- 4:
- 5:  $C_{idxs} \leftarrow p\_sort_{idx}[idx - G/2 : idx + G/2]$
- 6:  $C_{vals} \leftarrow p_k[C_{idxs}]$
- 7:
- 8: **for**  $cv, ci$  in  $[C_{vals}, C_{idxs}]$  **do**
- 9: **if**  $|cv - p_k(\hat{X}_k)| > \delta$  **then**  $C_{idxs}.\text{pop}(ci)$
- 10: **if**  $M_d[\hat{X}_k, ci] > d$  **then**  $C_{idxs}.\text{pop}(ci)$
- 11: **end for**
- 12:
- 13:  $p'_C \leftarrow \text{sum}(p_k[C_{idxs}])$
- 14:  $q'_C \leftarrow \text{sum}(q_k[C_{idxs}])$
- 15: **if** not  $r \sim \mathcal{U}[0, 1] \leq \min \left( 1, \frac{p'_C}{q'_C} \right)$  **then**
- 16:  $x \sim [p_k - q_k]_+, \hat{X}_k \leftarrow x$ , **break.**
- 17: **end if**
- 18: **end for**
- 19: **return**  $\hat{X}_{0:k}, k$

---

### C. Additional Results

In this section, we present additional experiments expanding upon the visualizations discussed in the main text.

**Top-1 probabilities** In Fig. 1, we illustrate the visualization of Top-1 probabilities across a wider variety of images. As shown, regardless of the prompts, many images exhibit numerous tokens with low Top-1 probability distributions.

**Visual quality comparison** In Fig. 3, we visually illustrate the differences in generation quality among various methods compared in Table 1. As shown in the figure, our GSD achieves approximately a 4× speed-up while maintaining generation quality comparable to lossless methods such as vanilla AR and SJD. In contrast, the naive lossy method also achieves acceleration but significantly degrades generation quality.

**GSD generation performance** Fig. 2 presents further qualitative results of our method when accelerated by an average factor of 3.6. As demonstrated in the figure, our GSD significantly accelerates AR image decoding while maintaining generation quality across diverse prompts.

## 037 **D. Prompts on Qualitative Experiment**

038 In Figure. 9 on main paper, the prompts for each images are  
039 as follows :

- 040 • *Rusty robot on a skateboard in the hallway of dormitory,*  
041 *photography, 4k, realistic, detailed, bright*
- 042 • *Origami astronaut, walking in the cloud, bright back-*  
043 *ground, realistic, 4k, photography, bright color*
- 044 • *photography, realistic, White cute fluffy dog, skyblue*  
045 *background, very intricate, very detailed, realistic.,*  
046 *bright*
- 047 • *color photo, photography, Face of a young man, very de-*  
048 *tailed, realistic. sharp, film grain, high contrast*
- 049 • *animation art work, cute, cat character, bright color pal-*  
050 *lette*



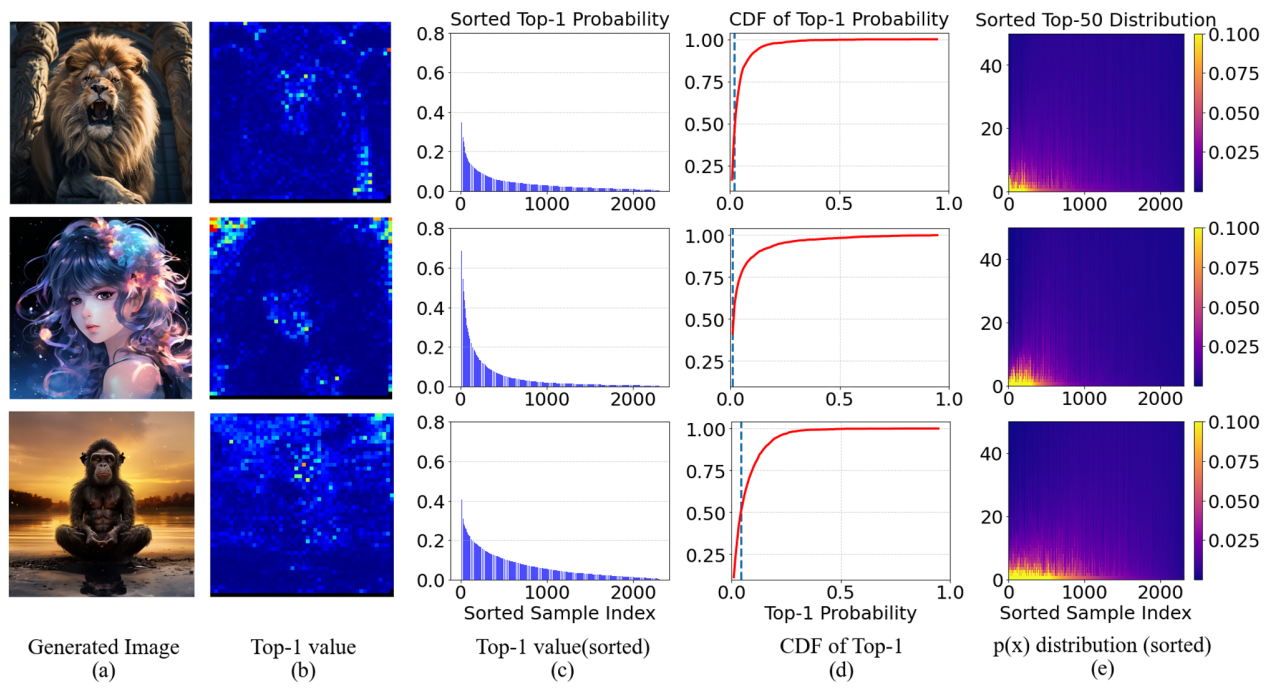


Figure 1. Additional  $p(x)$  visualization.



Figure 2. Qualitative experiment on various prompts. Our GSD shows on average 3.6x NFE acceleration while maintaining image quality





Figure 3. Qualitative comparison between methods in Table 1 of the main paper