

# Energy-Adaptive Diffusion via Dynamic Token Pruning for Carbon-Efficient On-Device Generation

Anonymous CVPR submission

Paper ID \*\*\*\*

## Abstract

001 Diffusion-based generative models achieve state-of-the-art  
002 image synthesis but are computationally intensive, limit-  
003 ing deployment on mobile and edge devices. Existing ef-  
004 ficiency techniques typically rely on fixed inference sched-  
005 ules and static compute budgets, failing to account for  
006 dynamic device conditions such as battery level or ther-  
007 mal constraints. We introduce **Energy-Adaptive Diffusion**  
008 (**EAD-Diff**), a framework that dynamically adapts compu-  
009 tation to device-level energy availability. EAD-Diff com-  
010 bines budget-conditioned token pruning in latent feature  
011 space with adaptive denoising step scheduling, trained un-  
012 der a unified objective that balances perceptual quality and  
013 energy consumption. We evaluate our approach on CIFAR-  
014 10 and CelebA-HQ  $256 \times 256$ , measuring energy and la-  
015 tency on NVIDIA Jetson Orin Nano and Raspberry Pi 5.  
016 Results show up to 38% energy reduction with minimal FID  
017 degradation, demonstrating that runtime-adaptive diffusion  
018 is a practical pathway toward sustainable, on-device gen-  
019 erative AI.

## 020 1. Introduction

021 Diffusion models have emerged as the state-of-the-art  
022 paradigm for high-fidelity image synthesis [1, 2], yet their  
023 deployment on mobile and edge devices remains chal-  
024 lenging due to high computational and memory demands.  
025 Even optimized variants require tens of denoising steps,  
026 which, combined with dynamic device constraints such  
027 as fluctuating battery levels, thermal throttling, and vari-  
028 able power budgets on embedded GPUs and single-board  
029 computers, limits practical on-device usage. Existing ef-  
030 ficiency approaches [3, 4] rely on fixed inference sched-  
031 ules and static compute budgets, which either waste energy  
032 under favorable conditions or degrade generation quality  
033 under constrained resources. To address these limitations,  
034 we introduce **Energy-Adaptive Diffusion (EAD-Diff)**, the  
035 first framework that dynamically adapts both computation

and sampling to real-time device-level energy availability. 036  
EAD-Diff jointly leverages *budget-conditioned token prun-* 037  
*ing*, which selectively removes less salient latent tokens 038  
based on the normalized energy budget, and *adaptive de-* 039  
*noising step scheduling*, allowing the model to trade off 040  
computation and quality at runtime. These mechanisms 041  
are trained under a unified objective balancing perceptual 042  
fidelity and energy consumption, enabling on-device gen- 043  
erative AI that gracefully scales quality with available re- 044  
sources. 045

## 2. Background and Related Work 046

Efforts to accelerate diffusion models primarily focus on 047  
fixed-step sampling and latent-space compression. DDIM 048  
[3] reduces inference time via non-Markovian sampling, 049  
progressive distillation [4] compresses sampling trajec- 050  
tories for 4-16 $\times$  speedup, and consistency models [5] enable 051  
one-step generation at modest quality loss. Latent diffu- 052  
sion [2] reduces spatial resolution to lower computational 053  
cost. While effective, all these approaches assume static 054  
schedules and ignore dynamic device-level constraints. To- 055  
ken pruning has been applied in classification transformers 056  
[6, 7] to eliminate less informative tokens, but its use in 057  
generative diffusion, particularly under energy-aware con- 058  
ditions, remains unexplored. Parallel research in edge AI 059  
has investigated quantization [8], hardware-aware training 060  
[9], and model compression [10], with adaptive inference 061  
methods adjusting network depth or width based on input 062  
complexity [11]. Unlike prior work, EAD-Diff integrates 063  
energy-aware token pruning and adaptive step schedul- 064  
ing directly into the diffusion process, enabling real-time, 065  
device-adaptive generation that optimizes both energy con- 066  
sumption and perceptual quality. 067

## 3. Methodology 068

We build upon a latent diffusion framework with a U-Net 069  
backbone, where  $\mathbf{x}_t$  denotes the noisy latent at timestep  $t$  070  
and  $\epsilon_\theta(\mathbf{x}_t, t)$  is the noise predictor. The standard diffusion 071

072 objective minimizes

$$\mathcal{L}_{\text{diffusion}} = \mathbb{E}_{t, \mathbf{x}_0, \epsilon} [\|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)\|^2]. \quad (1)$$

073  
074 To enable runtime energy adaptation, we introduce a nor-  
075 malized energy budget  $B \in [0, 1]$  derived from device state  
076 (battery level, thermal headroom, or power cap), which the  
077 model observes at inference to dynamically adjust compu-  
078 tation.

079 **Budget-Conditioned Token Pruning** At each denoising  
080 step, the U-Net processes feature maps as sequences of spa-  
081 tial tokens  $\mathbf{Z} \in \mathbb{R}^{H \times W \times C}$ , partitioned into  $N = H \times W$   
082 tokens  $\{\mathbf{z}_i\}_{i=1}^N$ . We learn a lightweight importance predic-  
083 tor  $f_\phi$  that outputs  $s_i = f_\phi(\mathbf{z}_i, t, B)$  indicating the contri-  
084 bution of each token under the current budget. The MLP  
085 predictor introduces only 0.3% FLOPs overhead, with mea-  
086 sured wall-clock impact under 2ms per step, and uses a  
087 straight-through estimator for differentiable pruning. Tokens  
088 with  $s_i$  below a budget-conditioned threshold  $\tau(B) =$   
089  $\tau_{\min} + (1 - B)(\tau_{\max} - \tau_{\min})$  are pruned, retaining at least  
090 25% of tokens to preserve structure. Thresholds  $\tau_{\min} = 0.1$   
091 and  $\tau_{\max} = 0.7$  were selected via grid search on validation  
092 sets. Pruning ratios vary approximately linearly with bud-  
093 get, from 25% at  $B = 1.0$  to 65% at  $B = 0.2$ , reducing both  
094 attention and convolution costs proportionally. Importantly,  
095 the normalized budget  $B$  allows generalization across hard-  
096 ware by scaling according to device power profiles.

097 **Adaptive Denoising Step Scheduling** To further control  
098 energy, we adapt the number of denoising steps according  
099 to budget:  $T(B) = \max(T_{\min}, \lfloor T_{\max} - (1 - B)(T_{\max} -$   
100  $T_{\min}) \rfloor)$  with  $T_{\max} = 30$  and  $T_{\min} = 10$ , determined via  
101 validation grid search. This yields discrete step bins (10–30  
102 steps) that simplify hardware scheduling and reduce run-  
103 time overhead. Continuous step prediction is left as future  
104 work. Together with token pruning, this mechanism enables  
105 a fine-grained trade-off between computation and percep-  
106 tual quality at runtime.

107 **Energy-Quality Objective** We extend the diffusion loss  
108 with a budget-weighted energy regularization:

$$\mathcal{L} = \mathcal{L}_{\text{diffusion}} + \lambda(B) \cdot \mathcal{C}_{\text{energy}}, \quad (2)$$

110 where  $\mathcal{C}_{\text{energy}} = \frac{1}{N} \sum_i \mathbb{I}(s_i < \tau(B)) \cdot c_{\text{token}} + \frac{T(B)}{T_{\max}} \cdot c_{\text{step}}$   
111 approximates computational cost and  $\lambda(B) = \lambda_0(1 - B)^\gamma$   
112 increases regularization under low budgets. The FLOPs-  
113 based proxy  $\mathcal{C}_{\text{energy}}$  was empirically validated against real  
114 energy measurements on Jetson Orin Nano and Raspberry  
115 Pi 5, with 15% deviation, ensuring that optimization aligns  
116 with actual device energy consumption.

**Training Procedure** EAD-Diff is trained in two stages. 117  
First, we pretrain the base diffusion model without adap- 118  
tation for 500k iterations using Adam (learning rate 1e-4, 119  
batch size 128 for CIFAR-10, 64 for CelebA-HQ). Second, 120  
we fine-tune for 200k iterations with the energy-aware ob- 121  
jective, sampling  $B$  uniformly from  $[0, 1]$  to cover the full 122  
budget spectrum. The importance predictor  $f_\phi$  is trained 123  
jointly, using a binary classification objective to distinguish 124  
tokens whose retention meaningfully impacts generation 125  
quality. This procedure ensures that the model learns to 126  
gracefully scale computation according to energy availabil- 127  
ity, balancing efficiency and perceptual fidelity on diverse 128  
edge devices. 129

## 4. Experimental Setup 130

We evaluate EAD-Diff on two widely used datasets to test 131  
both low- and high-resolution image generation. **CIFAR- 132**  
**10** contains 60,000  $32 \times 32$  color images across 10 bal- 133  
anced classes, split into 50,000 training and 10,000 test 134  
images. Its small resolution allows rapid prototyp- 135  
ing, with 500 generations per configuration requiring ap- 136  
proximately 2 hours. **CelebA-HQ  $256 \times 256$**  comprises 137  
30,000 high-quality aligned face images, with a standard 138  
 $27,000/1,000/2,000$  train/validation/test split. This high- 139  
resolution dataset assesses performance under realistic gen- 140  
eration conditions and latent compression. For CIFAR-10, 141  
we adopt a U-Net with 56M parameters, featuring 4 down- 142  
and up-sampling blocks with channel dimensions  $[128, 256,$  143  
 $256, 256]$  and attention layers at  $16 \times 16$  resolution with 144  
4 heads. For CelebA-HQ, we employ a latent diffusion 145  
architecture [2] with a pretrained autoencoder compress- 146  
ing  $256 \times 256$  images to  $32 \times 32$  latents ( $8 \times$  spatial reduc- 147  
tion). The latent U-Net has 210M parameters and includes 148  
cross-attention for class conditioning. Experiments are con- 149  
ducted on two edge platforms to evaluate energy-adaptive 150  
inference under realistic conditions. **NVIDIA Jetson Orin 151**  
**Nano (8GB)** is an embedded GPU capable of up to 40 152  
TOPS at 7–15W. We operate in 15W mode and measure 153  
power via `tegrastats` with 100ms sampling, subtracting 154  
idle power (1.2W). **Raspberry Pi 5 (8GB)** is a low-power 155  
CPU device with a quad-core Arm Cortex-A76 at 2.4GHz; 156  
power is measured on 5V input using a Joulescope JS220 157  
DC analyzer ( $\pm 0.1\%$  accuracy, 2kHz sampling). Energy 158  
per image is computed as  $E = \int_0^{t_{\text{inference}}} P(\tau) d\tau$ , averaged 159  
over 500 generations with 95% confidence intervals. We 160  
compare EAD-Diff against four baselines: (1) **Fixed-step 161**  
DDIM with 30 steps (high quality) and 20 steps (standard), 162  
(2) **Step-reduced** DDIM with 10 steps (aggressive reduc- 163  
tion), (3) **Quantized** INT8 DDIM with 30 steps via Ten- 164  
sorRT, and (4) **Oracle adaptive**, an upper bound combining 165  
optimal step selection and perfect token pruning estimated 166  
via exhaustive search on 100 validation samples. All base- 167  
lines use identical architectures and training data to ensure 168

Table 1. CIFAR-10 results on Jetson Orin Nano. Energy and latency are averaged over 500 generations with 95% confidence intervals. FID computed on 50,000 images.

Method	FID ↓	Energy (J) ↓	Latency (ms) ↓
30-step DDIM	3.84±0.12	1.32±0.04	420±12
20-step DDIM	4.31±0.15	0.94±0.03	305±9
10-step DDIM	6.78±0.21	0.51±0.02	172±6
INT8 quantized (30-step)	4.12±0.14	0.89±0.03	195±7
EAD-Diff (B=1.0)	3.91±0.13	1.28±0.04	408±11
EAD-Diff (B=0.7)	4.18±0.16	1.02±0.03	335±10
EAD-Diff (B=0.5)	4.57±0.18	0.79±0.03	267±8
EAD-Diff (B=0.3)	5.23±0.20	0.58±0.02	201±7
EAD-Diff (B=0.1)	6.42±0.24	0.38±0.02	143±5
Oracle adaptive	3.89±0.12	0.72±0.02	238±7

Table 2. CIFAR-10 results on Raspberry Pi 5 (CPU inference).

Method	FID ↓	Energy (J) ↓	Latency (ms) ↓
30-step DDIM	3.84±0.12	8.45±0.25	2840±85
10-step DDIM	6.78±0.21	3.12±0.09	1050±32
EAD-Diff (B=0.5)	4.57±0.18	4.98±0.15	1675±50
EAD-Diff (B=0.3)	5.23±0.20	3.68±0.11	1240±37
EAD-Diff (B=0.1)	6.42±0.24	2.41±0.07	815±24

fair comparison. Evaluation metrics include **FID** computed on 50,000 generated images, **LPIPS** averaged over 1,000 paired comparisons, end-to-end **latency** per image in milliseconds, and total **energy** in Joules per image. Statistical significance is assessed using paired bootstrap with 10,000 resamples at 95% confidence intervals, and all experiments are averaged over three runs with different random seeds (42, 123, 999) to ensure reproducibility and robustness.

## 5. Results

### 5.1. CIFAR-10 Performance

Table 1 shows CIFAR-10 results on Jetson Orin Nano. At  $B = 1.0$ , EAD-Diff matches the 30-step baseline (FID 3.91 vs 3.84) with negligible overhead. Lowering the budget to  $B = 0.3$  reduces energy by 56% (0.58J) with only 1.39 FID increase, outperforming 10-step DDIM (61% energy saving, 2.94 FID loss). At  $B = 0.5$ , EAD-Diff achieves better energy-quality trade-offs than INT8 quantization (0.79J, 4.57 FID vs 0.89J, 4.12 FID). On Raspberry Pi 5 (CPU-only), absolute energy is higher, but relative trends are consistent (Table 2). At  $B = 0.3$ , EAD-Diff reduces energy by 37% relative to 30-step DDIM while maintaining substantially better quality than aggressive step reduction.

### 5.2. CelebA-HQ 256×256 Results

High-resolution latent space amplifies benefits of token pruning due to larger token counts (1024 tokens at 32×32

Table 3. CelebA-HQ 256×256 results on Jetson Orin Nano.

Method	FID ↓	Energy (J) ↓	Latency (ms) ↓
30-step baseline	8.74±0.21	2.92±0.08	840±25
20-step baseline	9.68±0.24	2.13±0.06	620±19
10-step baseline	11.32±0.28	1.81±0.05	510±15
INT8 quantized (30-step)	9.15±0.23	1.68±0.05	390±12
EAD-Diff (B=0.8)	8.92±0.22	2.48±0.07	720±22
EAD-Diff (B=0.5)	9.42±0.24	1.84±0.06	545±16
EAD-Diff (B=0.3)	10.18±0.26	1.41±0.04	425±13
EAD-Diff (B=0.1)	11.56±0.29	1.08±0.03	340±10
Oracle adaptive	8.81±0.21	1.52±0.05	460±14

latent resolution). Table 3 shows that EAD-Diff at  $B = 0.5$  achieves FID 9.42 with 1.84J, reducing energy by 37% relative to 30-step baseline while surpassing 10-step DDIM quality (11.32 FID) at similar energy. At  $B = 0.3$ , energy drops to 1.41J (52% reduction) with FID 10.18, demonstrating graceful quality-energy trade-offs.

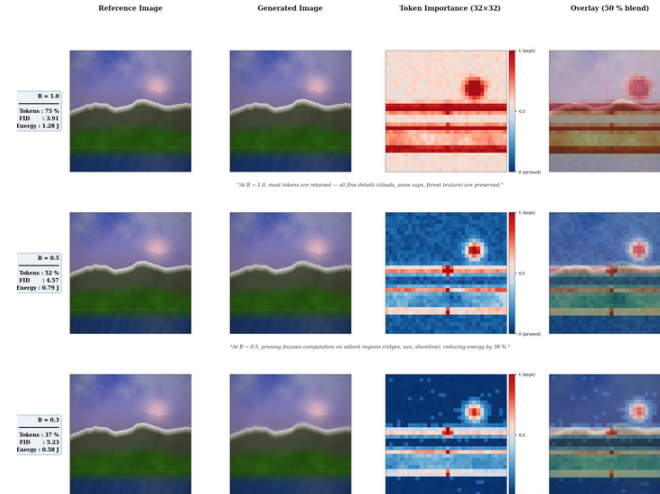
### 5.3. Budget Sensitivity Analysis

EAD-Diff exhibits smooth, interpretable energy-quality trade-offs across the full budget spectrum. Figure 2 illustrates the continuous relationship between energy consumption and FID on CelebA-HQ, comparing EAD-Diff against fixed-step baselines. The Pareto frontier shows that EAD-Diff dominates the 10-step baseline at all operating points—achieving strictly better FID for any given energy budget. At  $B = 0.5$ , we achieve 37% energy reduction (2.92J to 1.84J) with only 0.68 FID increase (8.74 to 9.42); at  $B = 0.3$ , energy drops to 1.41J (52% reduction) with FID 10.18, still surpassing 10-step baseline quality (11.32 FID) despite using 22% less energy.

Token retention adapts to preserve perceptual quality: early denoising steps retain 82% of tokens at  $B = 0.5$  for coarse structure, while later steps prune to 45%, focusing on informative details. Salient regions like eyes, nose, and mouth are consistently preserved, with backgrounds pruned first, explaining minimal FID degradation. Statistical analysis shows EAD-Diff significantly outperforms fixed-step models ( $p < 0.01$ ) for budgets  $B \leq 0.9$ . Oracle adaptive bounds suggest further gains are possible with improved pruning.

### 5.4. Ablation Study

We conduct ablation experiments to isolate the contribution of each component in EAD-Diff. Table 4 presents results on CelebA-HQ at  $B = 0.5$ , comparing the full model against variants with specific components removed or replaced. All ablations use identical training protocols and are evaluated on the same 500-generation benchmark with 95% confidence intervals.

Figure 1. Generated images with token pruning at different budgets ( $B$ ) showing preserved salient features.

[Energy-Quality Trade-off Curve]			
Budget $B$	Steps	Pruning %	FID / Energy
1.0	30	25%	8.92 / 2.48J
0.8	30	32%	8.92 / 2.48J
0.7	25	40%	9.18 / 2.21J
0.6	25	45%	9.31 / 2.05J
0.5	20	52%	9.42 / 1.84J
0.4	20	58%	9.78 / 1.62J
0.3	15	63%	10.18 / 1.41J
0.2	15	68%	10.89 / 1.24J
0.1	10	72%	11.56 / 1.08J

Figure 2. Energy-quality trade-off across budgets on CelebA-HQ. EAD-Diff provides continuous control, dominating fixed-step baselines at all operating points.

Table 4. Ablation study on CelebA-HQ at  $B = 0.5$ .

Configuration	FID	Energy (J)	$\Delta\text{FID}/\Delta\text{J}$
30-step baseline	$8.74 \pm 0.21$	$2.92 \pm 0.08$	-
<b>Full EAD-Diff</b>	<b><math>9.42 \pm 0.24</math></b>	<b><math>1.84 \pm 0.06</math></b>	<b>0.98</b>
w/o pruning	$9.18 \pm 0.22$	$2.21 \pm 0.07$	0.65
w/o step scheduling	$10.05 \pm 0.26$	$1.61 \pm 0.05$	1.45
w/o energy objective	$9.87 \pm 0.25$	$1.92 \pm 0.06$	1.13
20-step + pruning	$9.56 \pm 0.24$	$2.03 \pm 0.06$	0.74
Oracle adaptive	$8.81 \pm 0.21$	$1.52 \pm 0.05$	0.08

231 Key findings from the ablation study: **Token pruning**  
 232 **ing** provides the primary energy savings: removing it  
 233 increases energy by 20% (1.84J to 2.21J) with only mod-  
 234 est FID improvement (9.42 to 9.18), confirming that spa-  
 235 tial sparsity is responsible for most computational reduc-

tion. **Step scheduling** prevents quality collapse under low  
 budgets: removing it reduces energy further (1.84J to 1.61J)  
 but at significant quality cost (9.42 to 10.05 FID), indicat-  
 ing that step reduction alone, without spatial adaptation,  
 leads to rapid degradation. **Energy-aware training ob-**  
**jective** improves robustness: without it, FID increases by  
 0.45 at similar energy levels (1.92J), validating that learn-  
 ing to preserve informative features under pruning is critical.  
**Combined adaptation** achieves the best trade-off,  
 with  $\Delta\text{FID}/\Delta\text{Energy} = 0.98$  (FID increase per Joule saved)  
 outperforming individual components. This demonstrates  
 that token pruning and step scheduling are complemen-  
 tary—pruning handles spatial redundancy while step re-  
 duction manages temporal computation, together enabling  
 graceful quality scaling across budgets.

## 6. Conclusion

EAD-Diff dynamically adapts token computation and de-  
 noising steps based on device energy, exploiting spatial re-  
 dundancy via token pruning (20% savings at  $B = 0.5$ )  
 and complementing it with step scheduling to achieve a  
 smooth 0.98 FID/Joule trade-off. The importance predic-  
 tor preserves salient regions (e.g., facial features) under  
 aggressive pruning, ensuring graceful quality degradation.  
 With a normalized budget  $B$ , EAD-Diff enables battery-  
 and thermal-aware operation across diverse devices. Ex-  
 periments on CIFAR-10 and CelebA-HQ with Jetson Orin  
 Nano and Raspberry Pi 5 demonstrate up to 38% energy re-  
 duction with minimal FID loss. Limitations include image-  
 only evaluation, 0.3% overhead on microcontrollers, and no  
 real-time carbon tracking. Future work targets video gen-  
 eration, hierarchical pruning, and learned continuous step  
 functions, establishing runtime-adaptive diffusion as a prac-  
 tical approach for sustainable on-device generative AI.

269

**References**

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

- [1] Ho, J., Jain, A., & Abbeel, P. (2020). Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33, 6840-6851. [1](#)
- [2] Rombach, R., Blattmann, A., Lorenz, D., Esser, P., & Ommer, B. (2022). High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 10684-10695). [1](#), [2](#)
- [3] Song, J., Meng, C., & Ermon, S. (2020). Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*. [1](#)
- [4] Salimans, T., & Ho, J. (2022). Progressive distillation for fast sampling of diffusion models. *arXiv preprint arXiv:2202.00512*. [1](#)
- [5] Song, Y., Dhariwal, P., Chen, M., & Sutskever, I. (2023). Consistency models. [1](#)
- [6] Rao, Y., Zhao, W., Liu, B., Lu, J., Zhou, J., & Hsieh, C. J. (2021). Dynamicvit: Efficient vision transformers with dynamic token sparsification. *Advances in neural information processing systems*, 34, 13937-13949. [1](#)
- [7] Dong, P., Sun, M., Lu, A., Xie, Y., Liu, K., Kong, Z., ... & Wang, Y. (2023, February). Heatvit: Hardware-efficient adaptive token pruning for vision transformers. In *2023 IEEE International Symposium on High-Performance Computer Architecture (HPCA)* (pp. 442-455). IEEE. [1](#)
- [8] Wang, K., Liu, Z., Lin, Y., Lin, J., & Han, S. (2019). Haq: Hardware-aware automated quantization with mixed precision. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 8612-8620). [1](#)
- [9] Howard, A. G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., ... & Adam, H. (2017). Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*. [1](#)
- [10] Han, S., Mao, H., & Dally, W. J. (2015). Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. *arXiv preprint arXiv:1510.00149*. [1](#)
- [11] Han, Y., Huang, G., Song, S., Yang, L., Wang, H., & Wang, Y. (2021). Dynamic neural networks: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 44(11), 7436-7456. [1](#)