# Text2QR: Harmonizing Aesthetic Customization and Scanning Robustness for Text-Guided QR Code Generation

Guangyang Wu    Xiaohong Liu[†]    Jun Jia    Xuehao Cui    Guangtao Zhai
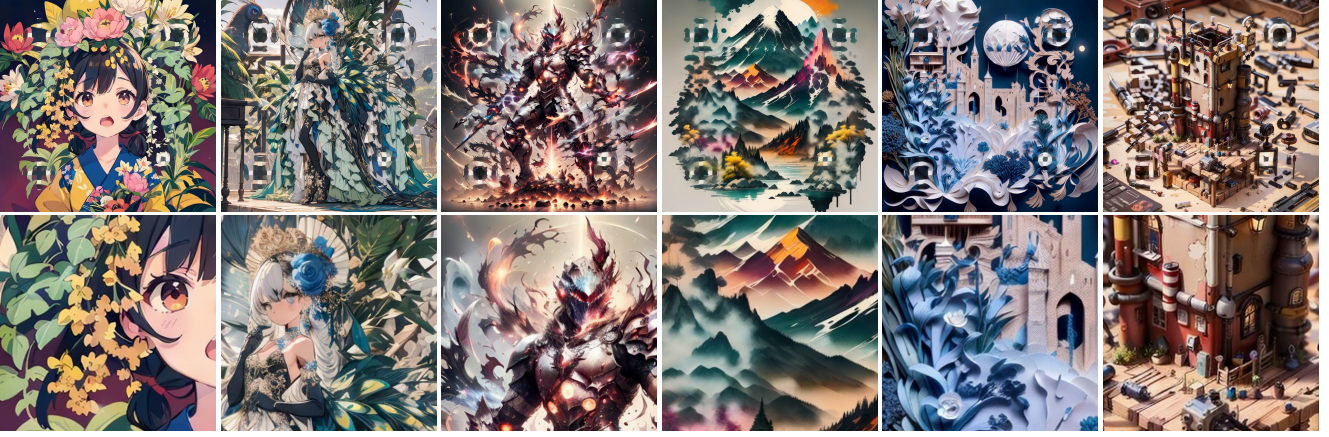Shanghai Jiao Tong University

Figure 1. Aesthetic QR codes (first row) and their zoom-in counterparts (second row) generated by Text2QR. Our QR codes not only exhibit exceptional scanning robustness but also showcase allure and intricate details, accommodating a diverse range of customized styles.

## Abstract

*In the digital era, QR codes serve as a linchpin connecting virtual and physical realms. Their pervasive integration across various applications highlights the demand for aesthetically pleasing codes without compromised scannability. However, prevailing methods grapple with the intrinsic challenge of balancing customization and scannability. Notably, stable-diffusion models have ushered in an epoch of high-quality, customizable content generation. This paper introduces Text2QR, a pioneering approach leveraging these advancements to address a fundamental challenge: concurrently achieving user-defined aesthetics and scanning robustness. To ensure stable generation of aesthetic QR codes, we introduce the QR Aesthetic Blueprint (QAB) module, generating a blueprint image exerting control over the entire generation process. Subsequently, the Scannability Enhancing Latent Refinement (SELR) process refines the output iteratively in the latent space, enhancing scanning robustness. This approach harnesses the potent generation capabilities of stable-diffusion models, navigating the trade-off between image aesthetics and QR code scannability. Our experiments demonstrate the seamless fusion of vi-sual appeal with the practical utility of aesthetic QR codes, markedly outperforming prior methods. Codes are available at* https://github.com/mulns/Text2QR

## 1. Introduction

In an age where digital interaction seamlessly converges with the physical world, Quick Response (QR) codes serve as vital conduits connecting these realms [2, 4, 12, 33, 39]. These ubiquitous two-dimensional codes have found extensive utility, bridging the divide between the physical and digital domains, yet their appearance remains stark, consisting of black and white modules engineered primarily for functional efficiency rather than aesthetic allure.

Amidst a growing consensus among users and stakeholders, a desire has arisen for QR codes that not only fulfill their core functions but also captivate with their visual appeal. The simplicity of QR codes, while undeniably efficient, is increasingly viewed as a missed opportunity to seamlessly integrate them into the modern visual landscape. The demand for aesthetically pleasing QR codes has proliferated, transcending boundaries into marketing, advertising, and artistic domains [2, 4, 33, 39].

Early techniques centered on image-to-image transfor-

---

† Corresponding author.

1

mations, utilizing reshuffling [33], fusion [12, 40], and style transfer methods [33, 39]. Although effective in generating predefined image styles, these approaches struggled to accommodate the diverse stylistic preferences of users, leaving a gap for a unified solution that addresses both customization and consistency. Recent advancements in the intersection of image generation and control have marked a transformative era. stable-diffusion models [28, 43] have emerged as robust engines for producing high-quality, versatile, and dynamically ranged content. Concurrently, an innovative approach for aesthetic QR code generation surfaced [10], leveraging ControlNet's capability to modulate luminance and darkness relationships within QR codes. However, this approach encountered a critical challenge, often exhibiting instability, necessitating the incorporation of auxiliary control models and manual parameter adjustments [10] to ensure both scannability and content quality.

In addressing these challenges, we introduce the innovative Text2QR pipeline, providing a solution for seamlessly generating QR codes that balance user-defined aesthetics and robust scannability. Our framework unfolds through three key steps: **(1)** Users initiate the process by generating their preferred images using the stable-diffusion model, while simultaneously encoding their desired message into a QR Code. **(2)** The synergy begins with the blending of these images in the QR Aesthetic Blueprint (QAB) module. This module generates a blueprint image, incorporating content from the pre-generated image (guidance image) and accurately reflecting the encoded message within the QR code. The blueprint image is then fed into ControlNet, guiding the stable-diffusion models to preserve user-defined aesthetics and maintain desired relationships among light and dark blocks of the QR code. While the generated results may pose decoding challenges, they exhibit a substantially improved distribution of light and dark blocks while remaining consistent with user preferences. **(3)** Subsequent to this stage, we construct an energy equation to quantify content and message consistency in the generated results. Optimizing this energy equation through gradient ascent iterations on latent codes gradually enhances scan robustness while preserving content consistency. Finally, the output QR code excels in both aesthetic appeal and scannability, achieving the delicate balance between user-defined customization and robust utility.

The contributions of this work can be summarized as:
• An integrated pipeline, Text2QR, that harmonizes user-defined aesthetics and robust scannability in QR code generation.
• The introduction of the QR Aesthetic Blueprint (QAB) for creating template images and the Scannability-Enhancing Latent Refinement (SELR) process for optimizing scan robustness while maintaining aesthetics.
• Superior performance compared to existing techniques,

establishing Text2QR as a state-of-the-art solution for QR code generation that excels in both visual quality and scanning robustness.

## 2. Related Works

**Aesthetic 2D Barcode.** In the era of digital interaction, QR codes play a pivotal role in bridging the virtual and physical realms. A variety of aesthetically pleasing 2D barcodes have been proposed as alternatives to the less appealing QR codes. Halftone QR codes, proposed by Chu et al. [5], rearrange the black/white modules of QR codes into an outline that semantically matches an input image. QR Image [12, 40] leverages the redundancy in the coding rules of QR codes to embed color images within QR codes. Recently, Su et al. [32, 33] have combined QR codes with style transfer to create artistic QR codes. These methods adhere to the standard QR code encoding rules and can be scanned and decoded by a common mobile phone scanner. To minimize the visibility of the locating patterns, Chen et al. [3, 4, 21] have designed encoding rules to satisfy the sensitivity of human visual system, making the locating patterns less noticeable. Additionally, TPVM [11] hides QR codes in videos, utilizing the frame rate difference between screens and human eyes. Similarly, invisible information hiding is applied to make information invisible but decodable after camera shooting [8, 9, 13, 14, 34, 36].

**Diffusion Based Generative Models.** Deep learning-based image processing [16–20, 29–31, 35, 37, 38] and generation methods [22, 24, 27, 42] has been fastly developed recently. Diffusion models such as GLIDE [24], DALLE-2 [27], Latent Diffusion [28], and Stable Diffusion [28], are proposed as a novel kind of generative model. These models create images through iterative denoising of initial random Gaussian noise and are able to outperform existing methods in many generative tasks. Of these, Stable Diffusion [28] is particularly innovative, as it transitions the denoising process from the image domain to a variational autoencoder's latent space, leading to significant reductions in data dimensions and training time. Alongside these advancements, various recent studies have presented methods for introducing diverse conditions to control the diffusion process. ControlNet [43] and T2I-Adapter [23] focus on structural control, with the former introducing an adapter mirroring stable diffusion's structure and trained under structural conditions, while the latter fine-tunes a lightweight adapter for detailed control over the produced scenes and content from the diffusion model. Instead, BLIP-Diffusion [15] and SeeCoder [41] aim to achieve controllable results based on image style. BLIP-Diffusion [15] extracts multi-modal topic representations and combines them with text prompts. While SeeCoder [41] discards text prompts and use reference images as control parameters.

## 3. Preliminary

Prior to presenting our method, we elucidate the process by which a QR code scanner decodes binary information from an aesthetic QR code image. Given an colored image featuring a QR code, we initially transform it into a grayscale representation by extracting its luminance channel (Y-channel of the YCbCr color space), denoted as $I \in \mathbb{R}^{H \times W}$, encompassing $L$ gray levels (typically 256). The scanner initially locates the Finder and Alignment markers [25, 40] to identify the QR code region and extract essential information such as the number of modules and module size. Let the QR code encompass $n \times n$ modules, each of size $a \times a$ pixels, where $n \cdot a \leq \min(H, W)$. Using the marker information, we construct a grid comprising $n^2$ modules denoted as $M_k, k \in [1, 2, \ldots, n^2]$. This grid divides image $I$ into $n^2$ patches, represented as $I_{M_k} \in \mathbb{R}^{a \times a}$.

The $k$-th module is decoded into a 1-bit information $\tilde{I}_k$, represented as 0 or 1, where $\tilde{I} \in \mathbb{R}^{n \times n}$ is the resulting binary image. Typically, scanners sample pixels within a central subregion of each module [39, 40]. Let $\theta$ be a square region with a size of $x \times x$ centered on module $M_k$, and $\mathbf{p} \in \{1, 2, \ldots, H\} \times \{1, 2, \ldots, W\}$ denotes pixel coordinates of $I$. The decoded binary value $\tilde{I}_k$ by a scanner is expressed as:

$$v_k = \frac{1}{x^2} \sum_{\mathbf{p} \in \theta} I_{M_k}(\mathbf{p}); \quad \tilde{I}_k = \begin{cases} 0, \text{if } v_k \leq \mathcal{T}_b, \\ 1, \text{if } v_k \geq \mathcal{T}_w, \\ -1, \text{otherwise.} \end{cases} \quad (1)$$

Here, $\mathcal{T}_b$ and $\mathcal{T}_w$ are thresholds for binarization. To account for symmetry, we set $\mathcal{T}_b = L \cdot (1 - \eta)/2$ and $\mathcal{T}_w = L \cdot (1 + \eta)/2$, where $\eta \in (0, 1)$. The hyperparameter $\eta$ governs the strictness of the binarization process.

Conventional QR code markers, traditionally characterized by square patterns, have shown adaptability to diverse styles while maintaining readability for conventional QR code scanners, as indicated by recent studies [3, 10, 21]. Achieving this adaptability involves specific pixel ratios, such as 1:1:3:1:1 for black and white modules, as detailed in [10, 25]. This flexibility includes preserving a cross center region to convey relevant information. Moreover, within the data regions of QR codes, maintaining binary results despite variations in sampled pixel colors is crucial for scanning robustness. As shown in Figure 2, even when colors and shapes subtly blend and vary, the sampled pixels consistently yield binary results aligned with ideal QR codes, ensuring robust decoding by standard QR code readers.

To facilitate analysis, we define the probability $e(I) = p(\tilde{I}_k = \mathcal{M}_k)$ to assess the error level of image $I$ with respect to the code target $\mathcal{M} \in \mathbb{R}^{n \times n}$. Notably, the function $e$ only characterizes the error proportion within data regions, omitting Finder and Alignment regions in a QR code.
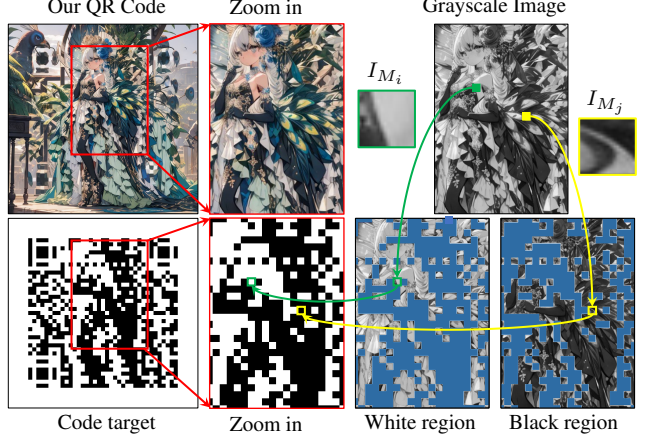


Figure 2. Illustration of preserving scanning-robustness. Each module in our QR code (e.g. $I_{M_i}$ and $I_{M_j}$) is correspondingly mapped to white (green arrow) or black (yellow arrow) blocks, collectively forming a standard QR code target We use blue masks to filter the white and black modules for better visualization.

## 4. Method
### 4.1. Overall

Figure 3 illustrates the overall structure of Text2QR, which is grounded in the Stable Diffusion (SD) model denoted as $\mathcal{G}$. The powerful customization capability enables the SD model to generate a user-preferred image $I^g = \mathcal{G}(c, z_0)$ with customized prompts $c$ and input noise $z_0$. Simultaneously, the input message is encoded into a QR code target $\mathcal{M}$, comprising $n \times n$ binary values (1 for white, 0 for black), representing the ideal color of each module.

Given $I^g$ and $\mathcal{M}$, Text2QR is designed to yield an aesthetically pleasing QR code denoted as $Q$. In pursuit of visual allure, $Q$ faithfully mirrors the semantic content, aesthetic style, and figure layout inherent in $I^g$. Simultaneously, for practical functionality, $Q$ is engineered to seamlessly reveal the encoded message upon scanning, adaptive to any standard QR code reader. The architectural framework of our pipeline unfolds across three distinct stages:

In the first stage, users prepare $I^g$ and $\mathcal{M}$, recording the associated parameters ($c$ and $z_0$). During the second stage, a pivotal step entails the seamless integration of information encapsulated within $I^g$ and $\mathcal{M}$ to formulate a comprehensive blueprint image, denoted as $I^b$, through the innovative QR Aesthetic Blueprint (QAB) module. Subsequently, $I^b$ undergoes processing in a ControlNet $\mathcal{C}$ to exert influence on the SD model. This influence involves adjusting intermediate features through a controlled process defined as:

$$I^s = \mathcal{G}(c, z_0 | \mathcal{C}(I^b, c, z_0)). \quad (2)$$

This integration ensures a synergistic output $I^s$ that harmoniously balances the aesthetic preferences derived from $I^g$ with the structural constraints imposed by $\mathcal{M}$. In the con-
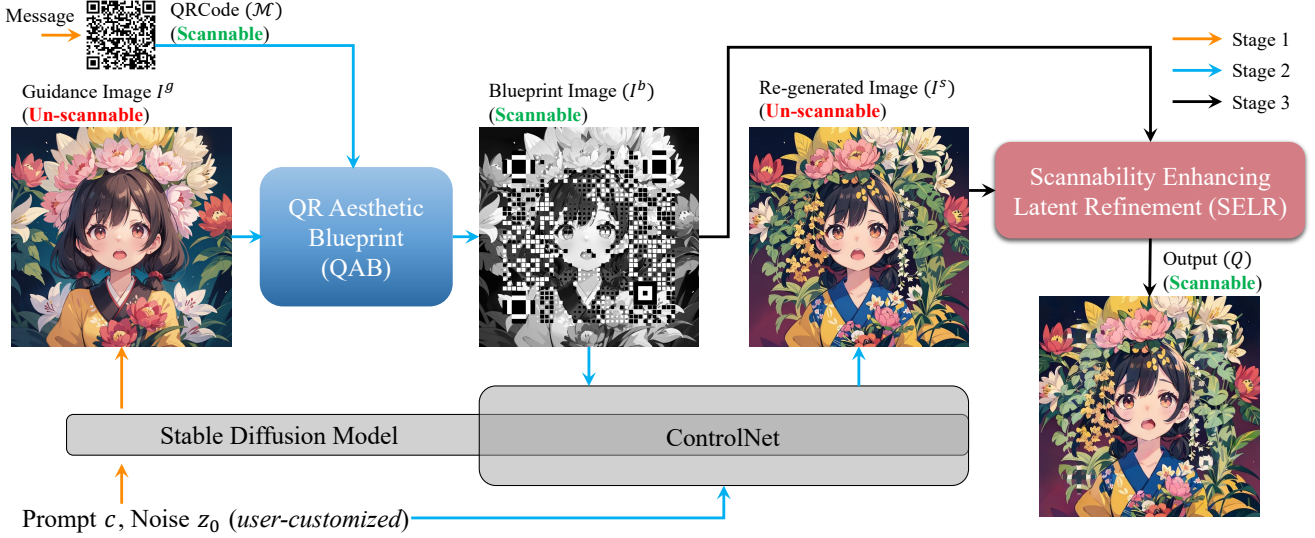
3

Figure 3. Overall Structure of the Text2QR. The pipeline consists of three stages, denoted with orange, blue and black lines. We propose the QAB module for generating a blueprint image used as controlling guidance, and propose the SELR module for refining the controlled output to enhance its scanning robustness.

clusive stage, $I^s$ undergoes iterative fine-tuning through the Scannability Enhancing Latent Refinement (SELR) module. This refines the scanning robustness of $I^s$ while meticulously preserving its aesthetic qualities. The output of this process is an aesthetically impressive QR code $Q$. Subsequent sections delve into detailed expositions on both QAB and SELR modules.

## 4.2. QR Aesthetic Blueprint

The module aims to create a scannable blueprint by integrating QR code information and guidance image details. Initially, we extract the luminance channel, denoted as $I^g_y$, from the guidance image $I^g$. To ensure comparable distributions, we preprocess $I^g_y$ and $\mathcal{M}$ using histogram polarization for luminance adjustment and a module reorganization method for pixel rearrangement, respectively. Finally, the Adaptive-Halftone method is applied to blend them, yielding the blueprint image $I^b$.

**Histogram polarization.** The primary aim of this module is to harmonize the histogram distribution of $I^g_y$ with that of the QR code. This process enhances the contrast of $I^g_y$, yielding a high-contrast grayscale image $I^{hc}$. The histogram polarization operation is represented by a look-up table $\mathcal{H}$, which maps pixel values from one gray level to another. For each pixel $\mathbf{p}$, let $\tau = I^g_y(\mathbf{p})$ and $\tau' = I^{hc}(\mathbf{p})$, we express this transformation as follows:

$$\tau' = \mathcal{H}(\tau). \tag{3}$$

Let $n_\tau$ denote the occurrences of gray level $\tau \in [0, L)$, we introduce the Cumulative Distribution Function (CDF) cor-
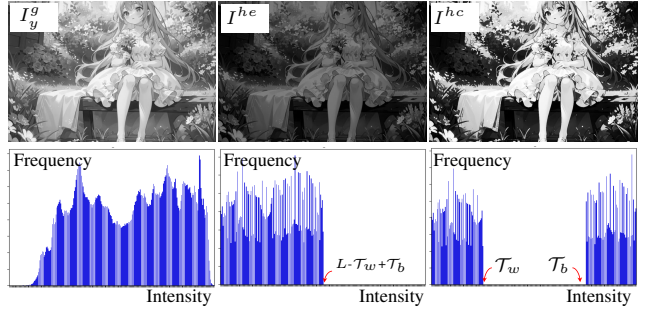


Figure 4. Visualization of the process of Histogram Polarization.

responding to gray level $\tau$ as:

$$\text{cdf}(\tau) = \sum_{i=0}^{\tau} \frac{n_i}{H \times W}. \tag{4}$$

The objective is to generate $I^{hc}$ with a flat histogram in the data range $[0, \mathcal{T}_b) \bigcup [\mathcal{T}_w, L)$, while excluding occurrences in the data range $[\mathcal{T}_b, \mathcal{T}_w)$. To achieve this, we first create a new image $I^{he}$, with a linearized CDF across the value range $[0, L - \mathcal{T}_w + \mathcal{T}_b]$. Let $\tilde{\tau} = I^{he}(\mathbf{p})$, we have:

$$\tilde{\tau} = (L - \mathcal{T}_w + \mathcal{T}_b) \cdot \text{cdf}(\tau). \tag{5}$$

Subsequently, we shift the pixels within the value range $[\mathcal{T}_b, L)$ by adding $\mathcal{T}_w - \mathcal{T}_b$ to obtain $I^{hc}$:

$$\tau' = \mathcal{H}(\tau) = \begin{cases} \tilde{\tau}, & \text{if } \tilde{\tau} < \mathcal{T}_b, \\ \tilde{\tau} + \mathcal{T}_w - \mathcal{T}_b, & \text{if } \tilde{\tau} \geq \mathcal{T}_b. \end{cases} \tag{6}$$

Figure 4 shows visualization of these processes, illustrating

Figure 5. Comparison of blueprint images and their corresponding ControlNet output. Utilizing a pure QR code as the blueprint (first column) yields a low error level $e$ but lacks semantic features. Employing a fixed size of $u = \frac{a}{3}$ (second column) leads to a substantial error level. Our Adaptive-Halftone blending method preserves realistic image content with a minimal error level.

the transformation of the image with a polarized histogram and high-contrast luminance.

**Module reorganization.** To blend the QR code $\mathcal{M}$ with $I^{hc}$, we first binarize the $I^{hc}$ to a binary image $I^{bin}$. This binary image guides the module reorganization method, denoted as $\mathcal{E}_r$, which rearranges the modules of $\mathcal{M}$ while keeping the encoded information. The process can be formulated as:

$$I^{bin}(\mathbf{p}) = \begin{cases} 0, & \text{if } I^{hc}(\mathbf{p}) < \mathcal{T}_b, \\ 1, & \text{if } I^{hc}(\mathbf{p}) > \mathcal{T}_w, \end{cases} \quad (7)$$

$$\mathcal{M}^r = \mathcal{E}_r(\mathcal{M}, I^{bin}). \quad (8)$$

**Adaptive-Halftone blending.** Considering the $k$-th module region $M_k$, we input an image patch after histogram polarization $I^{hc}_{M_k}$ and a target value $\mathcal{M}^r_k = 0$ or 1. Our goal is to obtain the blueprint image $I^b$, where $I^b_{M_k}$ can be *decoded to the correct information while preserving as much image content as possible.*

To achieve this objective, we introduce a novel blending method, Adaptive-Halftone blending. Specifically, for each module $M_k$, let $\theta_k$ be a square region of size $u \times u$ centered on the image patch $I^{hc}_{M_k}$ ($u \le a$). We fill the $I^{hc}_{\theta_k}$ with value $\mathcal{M}^r_k$ to generate $I^b_{M_k}$. The square region size $u$ is optimized by minimizing the code distance within this module. The simulating decoded value of this module corresponding to $u$ is defined as $E_k(I^b_{M_k}|u)$:

$$E_k(I^b_{M_k}|u) = \frac{1}{a^2}\left[\sum_{\mathbf{p}\in\theta_k} L \cdot \mathcal{M}^r_k + \sum_{\mathbf{p}\notin\theta_k} I^b_{M_k}(\mathbf{p})\right] \quad (9)$$
$$= \frac{1}{a^2}\left[u^2 \cdot L \cdot \mathcal{M}^r_k + \sum_{\mathbf{p}\notin\theta_k} I^b_{M_k}(\mathbf{p})\right].$$
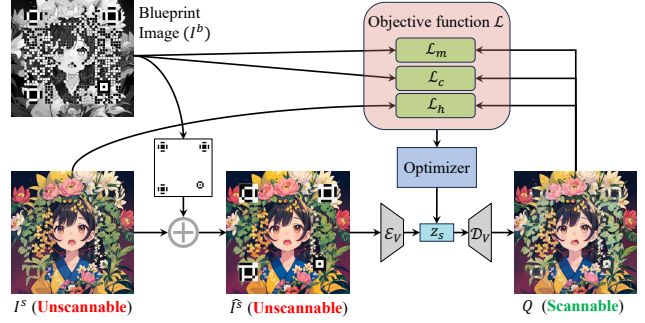


Figure 6. The recurrent pipeline of SELR.

Particularly, $E_k(I^b_{M_k}|u) = L \cdot \mathcal{M}^r_k$ when $u = a$. The objective is to minimize the code distance:

$$u_k = \arg\min_u \|E_k(I^b_{M_k}|u) - L \cdot (\eta \cdot \mathcal{M}^r_k + \frac{1-\eta}{2})\|. \quad (10)$$

According to the definition of thresholds $\mathcal{T}_b$ and $\mathcal{T}_w$, Equation 10 can be further simplified as:

$$s_k = \begin{cases} \arg\min_s \|E_k(I^b_{M_k}|u) - \mathcal{T}_b\|, \text{if } \mathcal{M}^r_k = 0, \\ \arg\min_s \|E_k(I^b_{M_k}|u) - \mathcal{T}_w\|, \text{if } \mathcal{M}^r_k = 1. \end{cases} \quad (11)$$

Having populated each module with a co-centered square block of adaptable size, we proceed to affix markers, including Finders and Alignment markers, onto the finalized blueprint image $I^b$. Figure 5 showcases diverse $I^s$ generated from distinct blueprint images. Our method dynamically adjusts the size of the central block for each module in $I^b$. This adjustment involves shrinking the block size when $I^{hc}$ effectively encapsulates the module's information to preserve more image content. Conversely, it enlarges the block size when a more pronounced control signal is necessary to ensure the module's scannability in $I^s$.

### 4.3. Scannability Enhancing Latent Refinement

While the re-generated image $I^s$ adheres to the structural constraints imposed by $\mathcal{M}^r$, it often lacks scannability due to the presence of numerous error modules. Addressing this issue, the Scannability Enhancing Latent Refinement (SELR) module offers a meticulous refinement process to enhance scanning robustness. The markers, encompassing finder and alignment patterns, are pivotal for determining the location and angle of a QR code, thereby influencing its scannability. Hence, we integrate their appearances onto $I^s$ before refinement, denoted as $\widehat{I^s}$. As shown in Figure 6, we encode the augmented image $\widehat{I^s}$ into a latent code $z_s$ using the encoder of a pre-trained Variational AutoEncoder (VAE) model, denoted as $\mathcal{E}_V$. The total objective function $\mathcal{L}$ is defined as the weighted sum of three terms: marker loss

$\mathcal{L}_m$, code loss $\mathcal{L}_c$, and harmonizing loss $\mathcal{L}_h$:

$$\begin{aligned}
\mathcal{L}(z) = {} & \lambda_1 \mathcal{L}_m(\mathcal{D}_V(z), I^b) \\
& + \lambda_2 \mathcal{L}_c(\mathcal{D}_V(z), I^b) \\
& + \lambda_3 \mathcal{L}_h(\mathcal{D}_V(z), I^s),
\end{aligned} \qquad (12)$$

where $\lambda_1$ to $\lambda_3$ are used to balance the multiple objectives. Here, $\mathcal{D}_V$ represents the decoder in the VAE model. The latent feature $z$ is initialized with $z_s$ and fine-tuned through an optimization process to minimize the total objective function, thereby controlling the scannability and aesthetic quality of the generated QR code $Q = \mathcal{D}_V(z)$. The code loss $\mathcal{L}_c$ is derived from the methodology proposed in [33]. This approach employs the SSLayer to extract module values and computes the module-based code loss through a competitive mechanism.

**Marker loss.** In line with previous discussions, scanners identify QR codes based on specific pixel ratios in marker regions. Consequently, our strategy centers on constraining the cross-center region of the marker, recognizing its crucial role in preserving scannability. To implement this, we introduce a binary mask, $\mathcal{K}_{cc}$, tailored to filter the cross-center region of markers. The aim is to safeguard the essential marker features against potential compromise due to aesthetic customization. Formally, the marker loss function $\mathcal{L}_m$ is defined as follows:

$$\mathcal{L}_m(Q, I^b) = \mathcal{K}_{cc} \cdot \parallel Q_y - I^b \parallel^2, \qquad (13)$$

where $Q_y$ denotes the luminance channel of the QR code $Q$. This formulation ensures marker integrity preservation while allowing for aesthetic modifications in non-marker regions of the QR code.

**Harmonizing loss.** Having addressed the marker and code-related scannability concerns, our approach further ensures the preservation of aesthetic qualities through a harmonizing loss. To maintain the intrinsic aesthetic style of the generated QR code $Q$, we employ a harmonizing loss focused on optimizing visual quality while upholding its original appeal. This loss function computes the $L^2$-Wasserstein distance, denoted as $W_2$, between the feature maps of $Q$ and $I^s$. Specifically, feature map $f_i$ is extracted from $i$-th layer of a pre-trained VGG-19 network ($i \in [1, 6, 11, 18, 25]$). The loss is formulated as:

$$\mathcal{L}_h = \sum_i W_2(f_i(Q), f_i(I^s)), \qquad (14)$$

where $f_i$ denotes the feature map extracted from the $i$-th layer ($i \in [1, 6, 11, 18, 25]$) of the VGG network. The $L^2$-Wasserstein distance, assuming the feature distributions approximate Gaussian distributions described by means and co-variances, can be expressed in a closed form [1]. Let $P_1$ and $P_2$ be Gaussian measures on $\mathbb{R}^n$ with means $\mu_1$

and $\mu_2 \in \mathbb{R}^n$ and non-singular covariance matrices $C_1$ and $C_2 \in \mathbb{R}^{n \times n}$, respectively. The $L^2$-Wasserstein distance $W_2(P_1, P_2)$ is given by:

$$\begin{aligned}
A &= trace(C_1 + C_2 - 2(\sqrt{C_1} C_2 \sqrt{C_1})^{\frac{1}{2}}), \\
W_2(P_1, P_2) &= \sqrt{\parallel \mu_1 - \mu_2 \parallel^2 + A}.
\end{aligned} \qquad (15)$$

The integration of the harmonizing loss ensures that the optimized output not only meets the functional requirements but also preserves aesthetic qualities. In essence, the SELR module leverages marker, code, and harmonizing losses to optimize both the scannability and aesthetic appeal of the generated QR code.

# 5. Experiments

We evaluate the performance of our QR codes in two aspects, aesthetic quality and scanning-robustness.
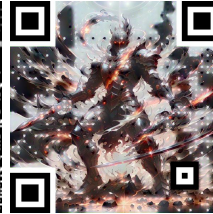
## 5.1. Implementation

We implement our program in PyTorch and conduct experiments on a NVIDIA GeForce 3090 GPU. For scanning-robustness assessment, we display QR codes on a 27-inch, 144Hz IPS-panel monitor. Default settings include $\eta = 0.6$ (following the work [33]) and a controlling strength of 1.4 for the pre-trained ControlNet (i.e., QR-Monster [10]). The VAE model aligns with the SD model, featuring frozen parameters, and VGG-19 that is pre-trained on MS-COCO extracts feature maps in SELR. During refinement, we employ Adam optimizer for 400 iterations with the learning rate set to 0.002 and default weights $\lambda_1$, $\lambda_2$ and $\lambda_3$ set to 1.0. QR codes ($\mathcal{M}$) are generated in version 5 of size $592 \times 592$ (i.e., $37 \times 37$ modules, each of $16 \times 16$ [33]). In this paper, we conduct comparisons on dataset comprising 100 generated images of $1,024 \times 1,024$, span various visual content and artistic styles.
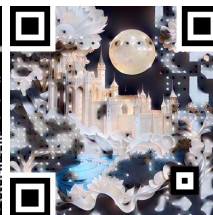
Table 1. Average scanning success rates are assessed across various scanners, considering different sizes and angles. "Scanner" denotes the native scanner of each system. We compare the accuracy (%) of our method and ArtCoder [33].

| Mobile Phone | APPs | Accuracy (%) | | | | | |
|---|---|---|---|---|---|---|---|
| | | $(3cm)^2$ | | $(5cm)^2$ | | $(7cm)^2$ | |
| | | $45°$ | $90°$ | $45°$ | $90°$ | $45°$ | $90°$ |
| iPhone 14Pro | Scanner | 100 | 100 | 100 | 100 | 100 | 100 |
| | TikTok | 100 | 100 | 100 | 100 | 100 | 100 |
| | WeChat | 100 | 100 | 100 | 100 | 96 | 96 |
| Huawei P40 | Scanner | 100 | 100 | 100 | 100 | 100 | 100 |
| | TikTok | 100 | 100 | 100 | 100 | 100 | 100 |
| | WeChat | 100 | 100 | 96 | 100 | 96 | 100 |

Table 2. Visual comparison of different methods. More results can be found in the supplementary material.

| Input | QArt [6] | Halftone QR [5] | ArtCoder [33] | Quick QR [26] | Text2QR (ours) |
|---|---|---|---|---|---|



## 5.2. Scanning Robustness

We assess the performance of our QR codes across various mobile devices and readers in this paper. Initially, we generate a set of 20 aesthetic QR codes with a resolution of 512 × 512. These codes are displayed on the screen in three commonly used sizes: 3cm × 3cm, 5cm × 5cm, and 7cm × 7cm. Positioned at a distance of 20cm, we scan each code using different mobile phones and apps, varying scanning angles. We record the average number of successful scans in 50 attempts, defining success as decoding within 3 seconds. Table 1 presents experimental results indicating that the average success rates consistently exceed 96%. It is worth noting that even in the case where decoding exceeds 3

seconds, our QR codes are still decodeable eventually. This robust performance demonstrates the reliability of our QR codes for real-world applications.

## 5.3. Aesthetic Quality

Although scanning robustness is preserved well, we also concern the visual appeal of the QR code.

**Comparison with existing methods.** We benchmark our methods against various aesthetic QR code approaches, including QArt [6], Halftone QR code [5], ArtCoder [33], and Quick QR [26], as detailed in Table 2. For ArtCoder that utilizes the neural-style transfer technique, we designate $I^g$ as both the content and style target. For QArt, Halftone
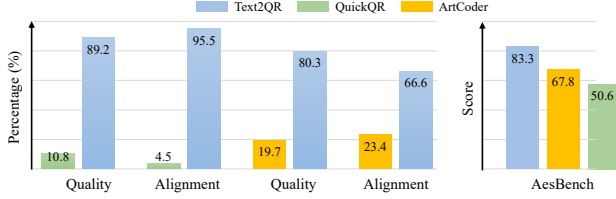
Figure 7. Statistical results of user study (left) and scores of AesBench [7] (right).
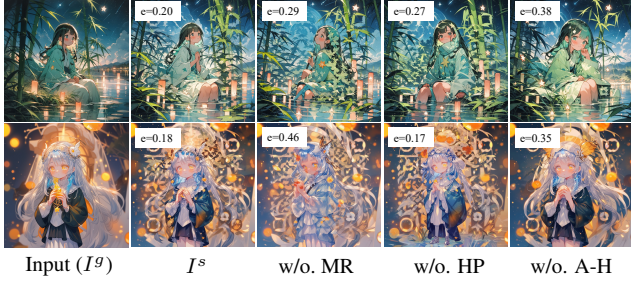


Figure 8. QAB Ablation Study: We assess the impact of Module Reorganization (MR), Histogram Polarization (HP), and Adaptive-Halftone Blending (A-H) on the generated $I^s$. Our result exhibits high consistency with the customized input $I^g$ and achieves a lower error level. (Note: These images are not scannable.)

QR code, and Quick QR, we employ $I^g$ as their reference input images. ArtCoder's results exhibit visible, undesired round spots, indicating repaired modules that can be distractable. Quick QR, on the other hand, yields inconsistent outcomes with the customized input. In contrast, our QR codes seamlessly integrate with the customized input image, featuring modules that are nearly invisible, ensuring superior aesthetic quality characterized by personalization, diversity, and artistic appeal.

In Figure 7, we showcase the outcomes of a user study involving 24 subjects comparing 200 generated QR-code images from various methods, where the ratio values indicate the percentages of participants preferring the corresponding model. Concurrently, we also employ AesBench [7] (ranges from 0 to 100, the higher the better) as an Aesthetic Assessment Metric to systematically score the different methods. Our method achieves superior performance across all evaluated aspects.

## 5.4. Ablation Study

We conduct a comprehensive ablation study, validating the necessity of each module in Text2QR.

**QAB Module.** In Text2QR, the QAB module generates a guidance blueprint image $I^b$ for the SD model to produce a high-quality aesthetic image $I^s$. Our objective is to ensure that $I^s$ not only shares a similar aesthetic style with $I^b$ but also maintains a low error level. This is achieved through
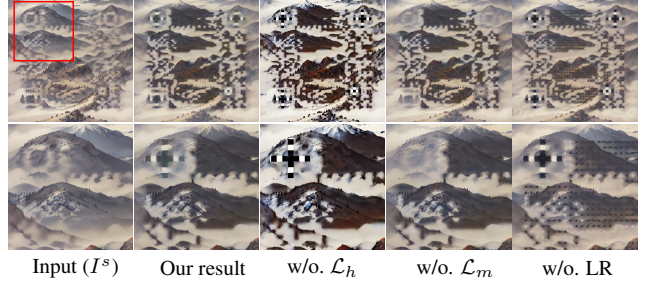


Figure 9. SELR Ablation Study: The first row presents results with various losses and without Latent Refinement (LR), while the second row zooms in on the red box region.

module reshuffling, histogram polarization, and Adaptive-Halftone blending steps. Figure 8 illustrates the impact of these steps on the error level $e$ and aesthetic quality of $I^s$. Our results demonstrate a high consistency between $I^s$ and the customized input $I^g$, accompanied by a notably lower error level.

**SELR Module.** During the SELR process, we initialize $Q$ with $I^s$ and iteratively refine its latent code. The marker loss and code loss actively enhance scanning robustness, while the harmonizing loss meticulously controls aesthetic quality. Figure 9 illustrates the importance of SELR module. Results refined directly on $Q$ (denoted as "w/o. LR") display visible, undesired round spots. Outputs without the marker loss resemble standard outputs but are unscannable. Omitting the harmonizing loss results in outputs with discordant appearances, underscoring its crucial role in achieving harmonious results. In conclusion, SELR refines the QR code result, ensuring a harmonious blend of functionality and aesthetic quality.

## 6. Conclusion

In summary, Text2QR utilizes the Stable-Diffusion (SD) model to effectively address the dual challenge of achieving user-defined aesthetics and scanning robustness in QR code generation. The strategic integration of the QR Aesthetic Blueprint (QAB) module that ensures generation stability and the Scannability Enhancing Latent Refinement (SELR) process that iteratively operates in the latent space, enhancing scanning robustness of the output. This innovative approach adeptly balances image aesthetics and scanning robustness, showcasing visual appeal and practical utility and surpassing previous approaches by a large margin, marking a substantial advancement in QR code generation.

## 7. Acknowledgment

# References

[1] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein Generative Adversarial Networks. In *Proceedings of the 34th International Conference on Machine Learning*, pages 214–223, 2017. 6

[2] Changsheng Chen, Wenjian Huang, Baojian Zhou, Chenchen Liu, and Wai Ho Mow. PiCode: A New Picture-Embedding 2D Barcode. *IEEE Transactions on Image Processing*, 25(8):3444–3458, 2016. 1

[3] Changsheng Chen, Wenjian Huang, Lin Zhang, and Wai Ho Mow. Robust and Unobtrusive Display-to-Camera Communications via Blue Channel Embedding. *IEEE Transactions on Image Processing*, 28(1):156–169, 2018. 2, 3

[4] Changsheng Chen, Baojian Zhou, and Wai Ho Mow. RA Code: A Robust and Aesthetic Code for Resolution-Constrained Applications. *IEEE Transactions on Circuits and Systems for Video Technology*, 28(11):3300–3312, 2018. 1, 2

[5] Hung-Kuo Chu, Chia-Sheng Chang, Ruen-Rone Lee, and Niloy J Mitra. Halftone QR Codes. *ACM Transactions on Graphics (TOG)*, 32(6):1–8, 2013. 2, 7

[6] Russ Cox. Qartcodes. https://research.swtch.com/qart, 2012. 7

[7] Huang et al. Aesbench: An expert benchmark for multimodal large language models on image aesthetics perception. *arXiv*, 2024. 8

[8] Han Fang, Weiming Zhang, Hang Zhou, Hao Cui, and Nenghai Yu. Screen-Shooting Resilient Watermarking. *IEEE Transactions on Information Forensics and Security*, 14(6): 1403–1418, 2018. 2

[9] Han Fang, Dongdong Chen, Feng Wang, Zehua Ma, Honggu Liu, Wenbo Zhou, Weiming Zhang, and Neng-Hai Yu. TERA: Screen-to-Camera Image Code with Transparency, Efficiency, Robustness and Adaptability. *IEEE Transactions on Multimedia*, pages 1–1, 2021. 2

[10] Anthony Fu. Stylistic qr code with stable diffusion. https://antfu.me/posts/ai-qrcode, 2023. 2, 3, 6

[11] Zhongpai Gao, Guangtao Zhai, and Chunjia Hu. The Invisible QR Code. In *Proceedings of the 23rd ACM International Conference on Multimedia*, pages 1047–1050, 2015. 2

[12] Gonzalo J Garateguy, Gonzalo R Arce, Daniel L Lau, and Ofelia P Villarreal. QR Images: Optimized Image Embedding in QR Codes. *IEEE Transactions on Image Processing*, 23(7):2842–2853, 2014. 1, 2

[13] Jun Jia, Zhongpai Gao, Kang Chen, Menghan Hu, Xiongkuo Min, Guangtao Zhai, and Xiaokang Yang. RIHOOP: Robust Invisible Hyperlinks in Offline and Online Photographs. *IEEE Transactions on Cybernetics*, pages 1–13, 2020. 2

[14] Jun Jia, Zhongpai Gao, Dandan Zhu, Xiongkuo Min, Guangtao Zhai, and Xiaokang Yang. Learning invisible markers for hidden codes in offline-to-online photography. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2273–2282, 2022. 2

[15] Dongxu Li, Junnan Li, and Steven CH Hoi. Blip-diffusion: Pre-trained subject representation for controllable text-to-image generation and editing. *arXiv preprint arXiv:2305.14720*, 2023. 2

[16] Wenhao Li, Guangyang Wu, Wenyi Wang, Peiran Ren, and Xiaohong Liu. Fastllve: Real-time low-light video enhancement with intensity-aware look-up table. In *Proceedings of the 31st ACM International Conference on Multimedia*, page 8134–8144, New York, NY, USA, 2023. Association for Computing Machinery. 2

[17] Xiaohong Liu, Yongrui Ma, Zhihao Shi, and Jun Chen. Griddehazenet: Attention-based multi-scale network for image dehazing. In *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*, pages 7313–7322. IEEE, 2019.

[18] Xiaohong Liu, Lingshi Kong, Yang Zhou, Jiying Zhao, and Jun Chen. End-to-end trainable video super-resolution based on a new mechanism for implicit motion estimation and compensation. In *IEEE Winter Conference on Applications of Computer Vision, WACV 2020, Snowmass Village, CO, USA, March 1-5, 2020*, pages 2405–2414. IEEE, 2020.

[19] Xiaohong Liu, Kangdi Shi, Zhe Wang, and Jun Chen. Exploit camera raw data for video super- resolution via hidden markov model inference. *IEEE Trans. Image Process.*, 30: 2127–2140, 2021.

[20] Xiaohong Liu, Zhihao Shi, Zijun Wu, Jun Chen, and Guangtao Zhai. Griddehazenet+: An enhanced multi-scale network with intra-task knowledge transfer for single image dehazing. *IEEE Trans. Intell. Transp. Syst.*, 24(1):870–884, 2023. 2

[21] Zehua Ma, Xi Yang, Han Fang, Weiming Zhang, and Nenghai Yu. Oacode: Overall aesthetic 2d barcode on screen. *IEEE Transactions on Multimedia*, 2023. 2, 3

[22] Maxim Maximov, Ismail Elezi, and Laura Leal-Taixé. Ciagan: Conditional identity anonymization generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5447–5456, 2020. 2

[23] Chong Mou, Xintao Wang, Liangbin Xie, Jian Zhang, Zhongang Qi, Ying Shan, and Xiaohu Qie. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. *arXiv preprint arXiv:2302.08453*, 2023. 2

[24] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021. 2

[25] Sean Owen. Zxing ("zebra crossing") barcode scanning library for java, android. https://github.com/zxing/zxing, 2013. 3

[26] Pixel ML, Inc. Quick qr art. https://quickqr.art, 2023. 7

[27] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1 (2):3, 2022. 2

[28] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 2

[29] Zhihao Shi, Xiaohong Liu, Chengqi Li, Linhui Dai, Jun Chen, Timothy N. Davidson, and Jiying Zhao. Learning for unconstrained space-time video super-resolution. *IEEE Trans. Broadcast.*, 68(2):345–358, 2022. 2

[30] Zhihao Shi, Xiaohong Liu, Kangdi Shi, Linhui Dai, and Jun Chen. Video frame interpolation via generalized deformable convolution. *IEEE Trans. Multim.*, 24:426–439, 2022.

[31] Zhihao Shi, Xiangyu Xu, Xiaohong Liu, Jun Chen, and Ming-Hsuan Yang. Video frame interpolation transformer. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 17461–17470. IEEE, 2022. 2

[32] Hao Su, Jianwei Niu, Xuefeng Liu, Qingfeng Li, Ji Wan, and Mingliang Xu. Q-Art Code: Generating Scanning-robust Art-style QR Codes by Deformable Convolution. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 722–730, 2021. 2

[33] Hao Su, Jianwei Niu, Xuefeng Liu, Qingfeng Li, Ji Wan, Mingliang Xu, and Tao Ren. Artcoder: an end-to-end method for generating scanning-robust stylized qr codes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2277–2286, 2021. 1, 2, 6, 7

[34] Matthew Tancik, Ben Mildenhall, and Ren Ng. Stegastamp: Invisible Hyperlinks in Physical Photographs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2117–2126, 2020. 2

[35] Wenyi Wang, Guangyang Wu, Weitong Cai, Liaoyuan Zeng, and Jianwen Chen. Robust prior-based single image super resolution under multiple gaussian degradations. *IEEE Access*, 8:74195–74204, 2020. 2

[36] Eric Wengrowski and Kristin Dana. Light Field Messaging with Deep Photographic Steganography. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1515–1524, 2019. 2

[37] Guangyang Wu, Lili Zhao, Wenyi Wang, Liaoyuan Zeng, and Jianwen Chen. Pred: A parallel network for handling multiple degradations via single model in single image super-resolution. In *2019 IEEE International Conference on Image Processing (ICIP)*, pages 2881–2885, 2019. 2

[38] Guangyang Wu, Xiaohong Liu, Kunming Luo, Xi Liu, Qingqing Zheng, Shuaicheng Liu, Xinyang Jiang, Guangtao Zhai, and Wenyi Wang. Accflow: Backward accumulation for long-range optical flow. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 12119–12128, 2023. 2

[39] Mingliang Xu, Hao Su, Yafei Li, Xi Li, Jing Liao, Jianwei Niu, Pei Lv, and Bing Zhou. Stylized aesthetic QR code. *IEEE Trans. Multim.*, 21(8):1960–1970, 2019. 1, 2, 3

[40] Mingliang Xu, Qingfeng Li, Jianwei Niu, Hao Su, Xiting Liu, Weiwei Xu, Pei Lv, Bing Zhou, and Yi Yang. ART-UP: A novel method for generating scanning-robust aesthetic QR codes. *ACM Trans. Multim. Comput. Commun. Appl.*, 17(1): 25:1–25:23, 2021. 2, 3

[41] Xingqian Xu, Jiayi Guo, Zhangyang Wang, Gao Huang, Irfan Essa, and Humphrey Shi. Prompt-Free Diffusion: Taking" Text" out of Text-to-Image Diffusion Models. *arXiv preprint arXiv:2305.16223*, 2023. 2

[42] Liming Zhai, Qing Guo, Xiaofei Xie, Lei Ma, Yi Estelle Wang, and Yang Liu. A3gan: Attribute-aware anonymization networks for face de-identification. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 5303–5313, 2022. 2

[43] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847, 2023. 2