

# BengaliDiff: Diffusion Model for One-Shot Bengali Font Generation

Md Bilayet Hossain<sup>1,2</sup>, Honghui Yuan<sup>1</sup>, Shabnur Anonna Akhy<sup>1,2</sup> and Prof.Keiji Yanai<sup>1</sup>

Department of Informatics, The University of Electro-Communications, Tokyo, Japan<sup>1</sup>  
Department of Computer Science and Engineering, Daffodil International University,  
Dhaka, Bangladesh<sup>2</sup>



国立大学法人

電気通信大学

The University of Electro-Communications



YANAI Lab

# Motivation

- Bengali is spoken by over 200 million people but has very **few high-quality fonts**.
- Existing font generation models (like for Chinese or Latin) don't work well for Bengali.
- Creating Bengali fonts by hand is **slow and complex**.
- We aim to build an automatic system that can **generate Bengali fonts** from just a few examples.

# Challenges

- **Challenge:** Existing models (e.g. **FontDiffuser** [1]) fail to generate on Bengali font.
- **FontDiffuser** struggles with:
  - Distorted matras and ligatures
  - Inconsistent stroke thickness
  - Style mismatch and loss of detail



**Figure 1:** The results generated by FontDiffuser using Bengali font as the input image.

# Base Architecture: FontDiffuser

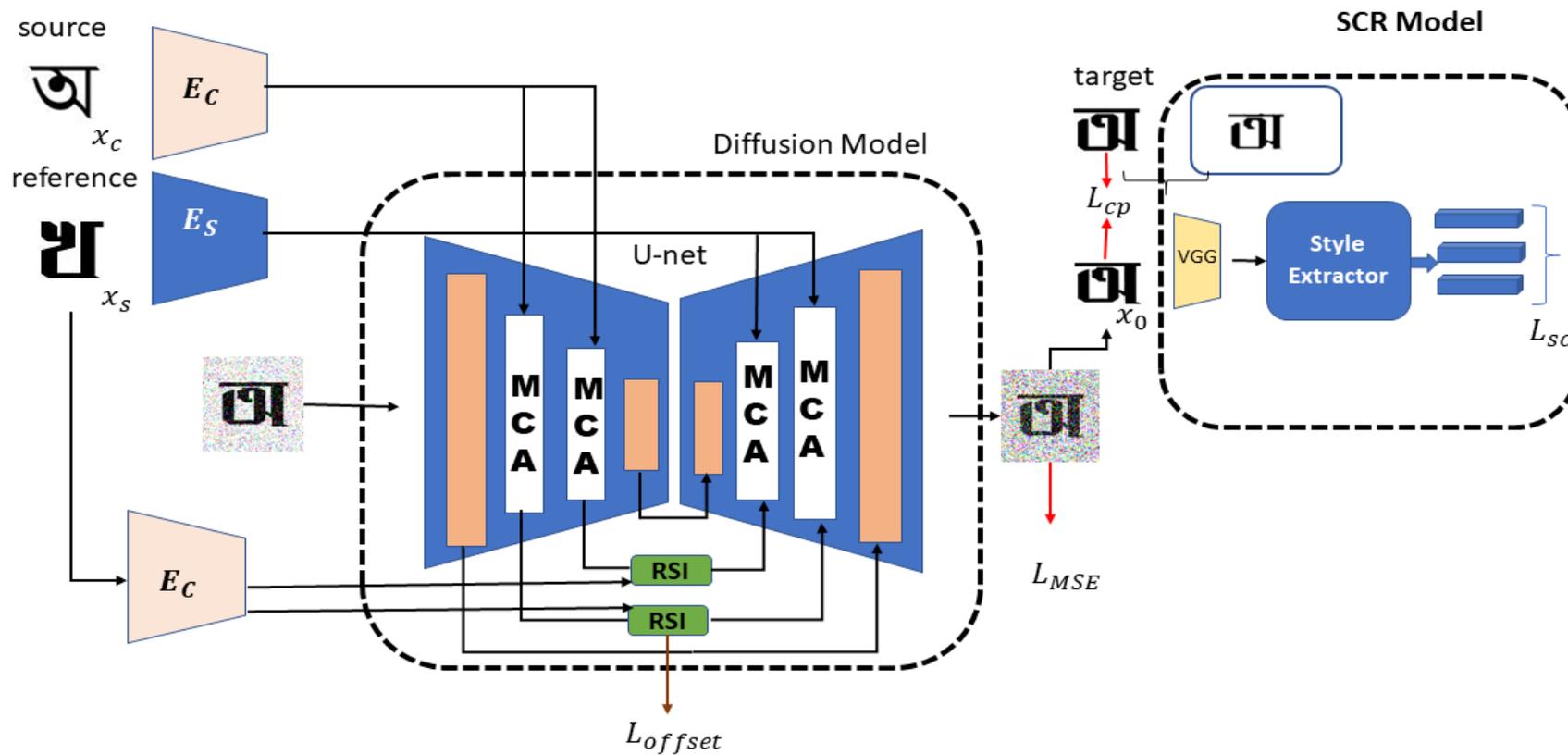


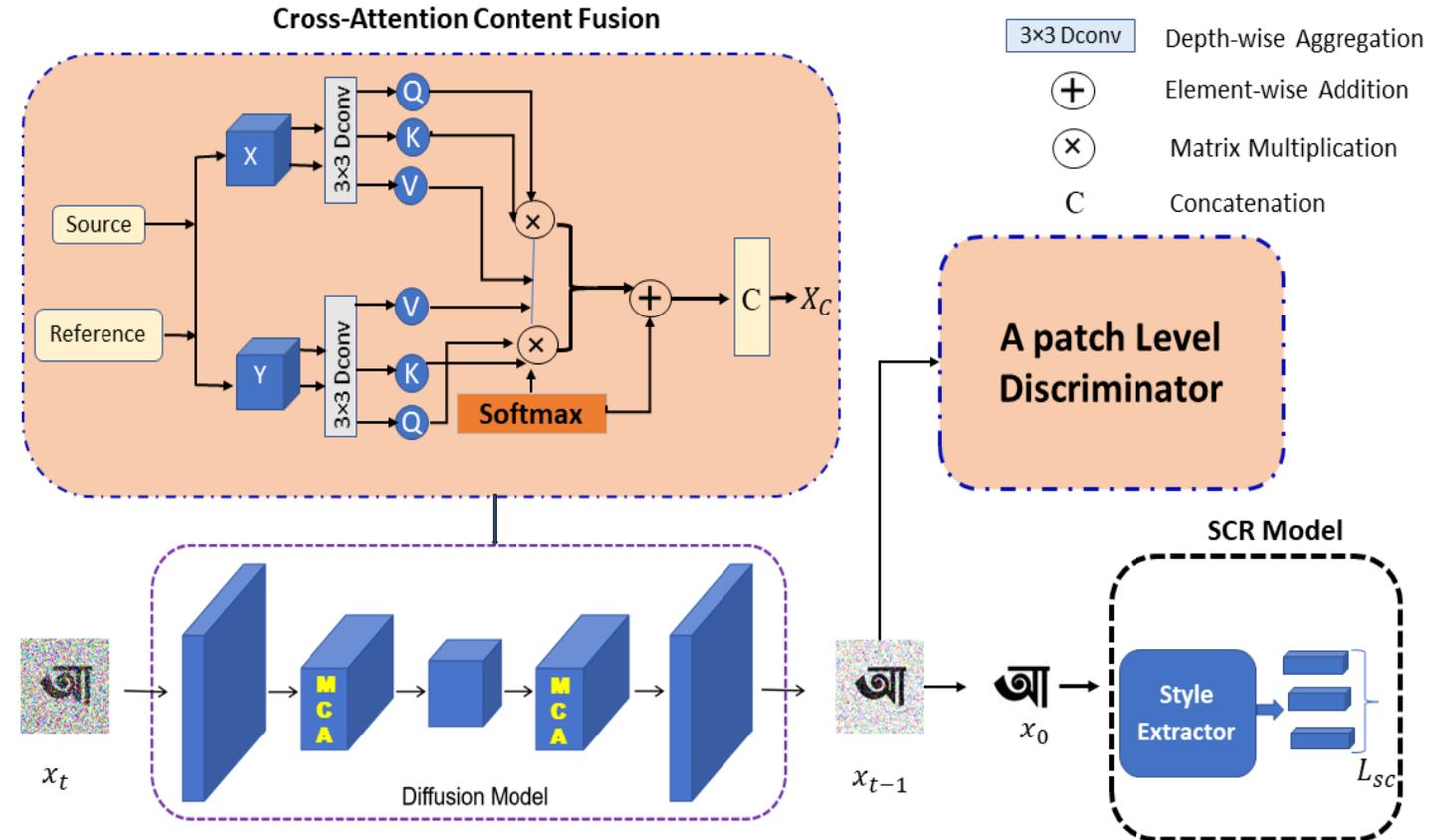
Figure 2: The framework of the base method FontDiffuser [1] .

- **FontDiffuser uses U-Net with:**
  - Content Encoder (structure)
  - Style Encoder (appearance)

- **Two key modules:**
  - MCA: Multi-Scale Content Aggregation
  - RSI: Reference Structure Interaction

# Our Proposed Method “BengaliDiff: Overview”

- **First diffusion-based model** tailored for Bengali font generation.
- **Our Key innovations:**
  - Cross-Attention Content Fusion
  - Patch-level Adversarial Discriminator
- Synthesizes high-quality fonts from few references



**Figure 3:** Overview of our proposed method. The Cross-Attention Content Fusion and A Patch Level Discriminator in the blue dashed line boxes.

# Qualitative Results Comparison with FontDiffuser and Our Method

|                |   |   |   |   |    |   |
|----------------|---|---|---|---|----|---|
| Reference:     | আ |   |   | ঞ |    |   |
| Source:        | অ | খ | ঞ | য | বি | ঞ |
| Font Diffuser: | অ | খ | ঞ | য | বি | ঞ |
| Ours:          | অ | খ | ঞ | য | বি | ঞ |
| Target:        | অ | খ | ঞ | য | বি | ঞ |

## Seen Font Unseen Characters:

- BengaliDiff **outperforms** FontDiffuser
- Achieves **better content preservation**
- **Sharper stroke** details
- More accurate **style transfer**

**Figure 4:** Qualitative Results Comparisons on SFUC between our method and the previous state-of-the-art method of FontDiffuser.

# Qualitative Results Comparison with FontDiffuser and Our Method

|                |    |    |      |    |    |    |
|----------------|----|----|------|----|----|----|
| Reference:     | ব  |    |      | খি |    |    |
| Source:        | বু | পি | স্কু | ডি | যী | জী |
| Font Diffuser: | বু | পি | স্কু | ডি | যী | জী |
| Ours:          | বু | পি | স্কু | ডি | যী | জী |
| Target:        | বু | পি | স্কু | ডি | যী | জী |

## Unseen Font Unseen Characters:

- BengaliDiff maintains **structural accuracy**
- **Style transfer** can still be slightly

**Figure 5:** Qualitative Results Comparisons on UFUC between our method and the previous state-of-the-art method of FontDiffuser.

# Quantitative Results Comparison with FontDiffuser and Our Method

**Table 2: Quantitative results Comparisons with FontDiffuser and BengaliDiff**

| Model               | FID ↓         | LPIPS ↓       | L1 ↓          | RMSE ↓        | SSIM ↑        |
|---------------------|---------------|---------------|---------------|---------------|---------------|
| FontDiffuser (SFUC) | 0.8685        | 0.3568        | <b>0.2377</b> | <b>0.3151</b> | <b>0.7131</b> |
| BengaliDiff (SFUC)  | <b>0.7063</b> | <b>0.3223</b> | 0.2547        | 0.3240        | 0.6751        |
| FontDiffuser (UFUC) | 0.9573        | 0.3738        | <b>0.2419</b> | <b>0.3142</b> | <b>0.6855</b> |
| BengaliDiff (UFUC)  | <b>0.7280</b> | <b>0.3406</b> | 0.2706        | 0.3357        | 0.6491        |

- BengaliDiff achieves **lower FID and LPIPS** scores, indicating more realistic and perceptually similar images.

# Ablation Study

## - Tested 3 configurations:

- Baseline (FontDiffuser)
- + Cross-Attention (CA)
- + CA + Discriminator (Full BengaliDiff)

- Each module improves **visual fidelity and structure**

| Source | Reference | Baseline  | +CA   | +CA+D   | Target  |
|--------|-----------|---|---|---|---|
| শ্বা   | শ্বা      |  |  |  |  |
| শ্বো   | শ্বো      |  |  |  |  |
| শ্বৌ   | শ্বৌ      |  |  |  |  |

**Figure 6:** Qualitative evaluation results of ablation studies. An illustration of several modules. CA and D represent Cross-attention and Discriminator, respectively. Red boxes represent the missing strokes, while green represents the corresponding improvements.

# Conclusion

- **BengaliDiff** generates high-quality Bengali fonts from few samples.
- Preserves **complex character structure** and ensures **style consistency**.
- **Cross-attention and discriminator** improve style and structure.
- Outperforms FontDiffuser in seen font test sets.
- Useful for **OCR, typography, and font design**.

# Limitations & Future Work

## Limitations

- Trained only on digital fonts
- Poor generalization to artistic/handwritten styles
- Limited control over font style

## Future Work

- Add support for handwritten glyphs
- Improve style transfer to unseen
- Enable editable style vectors
- Support dynamic ligatures

**Thank you**

# Question & Answer

**Please Feel Free to Ask Question**

Email: [h2495009@gl.cc.uec.ac.jp](mailto:h2495009@gl.cc.uec.ac.jp)