

# Rewriting a Deep Generative Model<sup>1</sup>

September 23, 2020

---

<sup>1</sup> "Rewriting a Deep Generative Model" by Bau et al., 2020

# Overview

- ▶ *Rewriting a generator* means changing its weights in such a way that it now generates a different distribution



**Figure:** An example of rewriting: generator now generates churches with trees on top of their domes

- ▶ Authors provide a different perspective on NN layers as key-value memory banks
- ▶ They rewrite G by inserting a new key-value pair into one of its layers

# An old way to do rewriting

- ▶ Imagine that we want to rewrite  $G(z, \theta_0)$
- ▶ For this, we first generate several  $x_i = G(z_i, \theta_0)$  for different  $z_i$
- ▶ Then we update each  $x_i \mapsto x_{*i}$  to describe our new rule
  - ▶ For example, photoshop a hat on a horse' head

(a) Copy



(b) Paste



## An old way to do rewriting (continued)

Now, find the “rewritten”  $G$  parameters  $\theta_1$  by minimizing:

$$\theta_1 = \arg \min_{\theta} \mathcal{L}_{\text{smooth}}(\theta) + \lambda \mathcal{L}_{\text{constraint}}(\theta) \quad (1)$$

$\mathcal{L}_{\text{smooth}}$  helps us not to deviate from old images:

$$\mathcal{L}_{\text{smooth}}(\theta) \triangleq \mathbb{E}_z [\ell(G(z; \theta_0), G(z; \theta))] \quad (2)$$

$\mathcal{L}_{\text{constraint}}$  helps us to enforce the rule:

$$\mathcal{L}_{\text{constraint}}(\theta) \triangleq \sum_i \ell(x_{*i}, G(z_i; \theta)) \quad (3)$$

Unfortunately,  $G$  overfits and does not generalize too well

But authors found that this strategy works much better if we optimize over a single layer instead of all the parameters

## NN layers are memory banks

- ▶ We can see a linear layer  $y = Wx$  as an associative memory (a 50-year old idea):

$$v_i \approx Wk_i \quad (4)$$

- ▶ So, instead of input-output pairs we now have key-value pairs
- ▶ I.e., given an input  $k_i$ , layer  $W$  “recalls” output  $v_i$ .
- ▶ From this perspective, rewriting can be seen as overwriting key  $k_i$  with new value  $v_{*i}$  (or inserting a new one)
- ▶ For a convolutional layer, we just convert  $n_{\text{in}} \times k \times k \times n_{\text{out}}$  tensor into  $(n_{\text{in}} \cdot k \cdot k) \times (n_{\text{out}})$  matrix

## Rewriting by memory update in a linear case

- ▶ We now can rewrite  $G$  by overwriting one of its keys  $k_*$  in layer  $W^\ell$
- ▶ For this, we solve the optimization problem:

$$W_1 = \arg \min_W \|V - WK\|^2 \quad \text{s.t. } v_* = W_1 k_* \quad (5)$$

- ▶ It has an analytical solution, but authors develop a different perspective to be able to generalize to non-linear transforms

## Rewriting by memory update in a *non-linear* case

- ▶ NN layers are usually non-linear transforms:  $y = f(x, W)$
- ▶ This still can be seen as an associative memory  $v_i = f(k_i, W)$
- ▶ In this case, there is no analytical solution to update the memory
- ▶ But for a linear case, we can show that

$$W_1 = W_0 + \Lambda (C^{-1}k_*)^\top \quad (6)$$

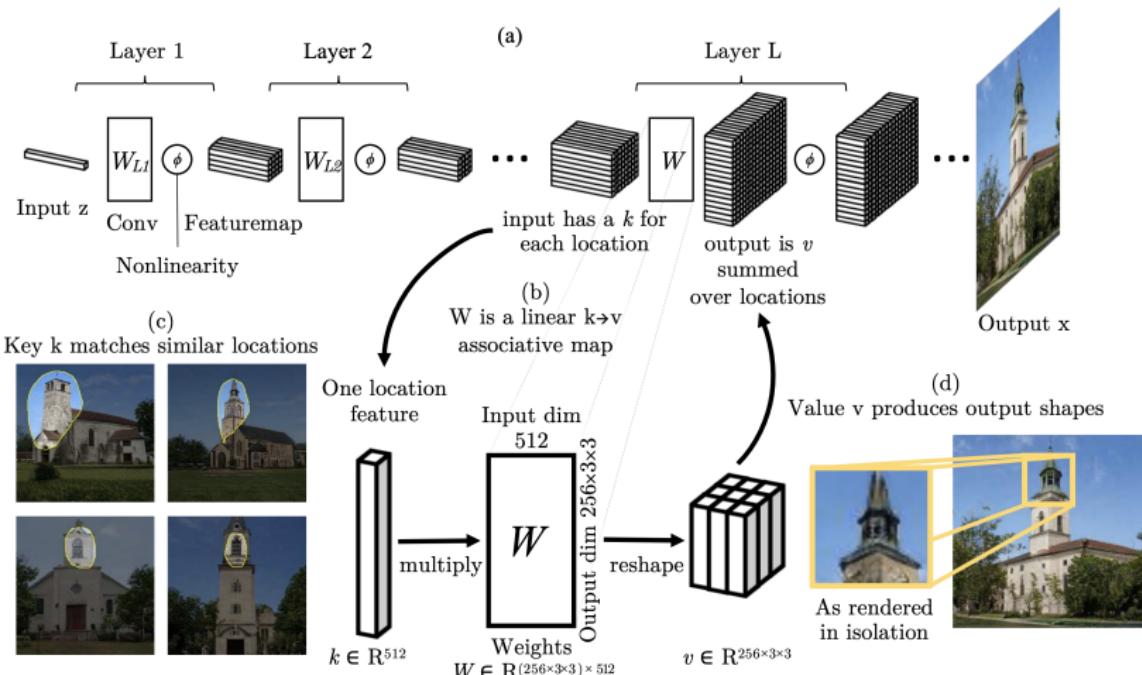
where

$$\Lambda_1 = \arg \min_{\Lambda \in \mathbb{R}^M} \|v_* - [W_0 + \Lambda(C^{-1}k_*)^\top]k_*\| \quad (7)$$

- ▶ Authors generalize this procedure to a non-linear case by finding  $\Lambda_1$  via:

$$\Lambda_1 = \arg \min_{\Lambda \in \mathbb{R}^M} \|v_* - f(k_*; W_0 + \Lambda(C^{-1}k_*)^\top)\| \quad (8)$$

# Rewriting procedure visualization



# Qualitative results



## Final thoughts

- ▶ Rewriting shows a powerful property of modern GANs
- ▶ But generator does not know how to blend a new value nicely with the novel context
  - ▶ That's why for many samples results look clumsily "copy-pasted"
  - ▶ How can we alleviate this?
  - ▶ (Results still look impressive)
- ▶ Associative memory perspective is very interesting since this is how human memory works