

Generative Memory-Guided Semantic Reasoning Model for Image Inpainting

Xin Feng, Wenjie Pei, Fengjun Li, Fanglin Chen, *Member, IEEE*, David Zhang, *Life Fellow, IEEE* and Guangming Lu, *Member, IEEE*

Abstract—The critical challenge of single image inpainting stems from accurate semantic inference via limited information while maintaining image quality. Typical methods for semantic image inpainting train an encoder-decoder network by learning a one-to-one mapping from the corrupted image to the inpainted version. While such methods perform well on images with small corrupted regions, it is challenging for these methods to deal with images with large corrupted area due to two potential limitations. 1) Such one-to-one mapping paradigm tends to overfit each single training pair of images; 2) The inter-image prior knowledge about the general distribution patterns of visual semantics, which can be transferred across images sharing similar semantics, is not explicitly exploited. In this paper, we propose the Generative Memory-guided Semantic Reasoning Model (*GM-SRM*), which infers the content of corrupted regions based on not only the known regions of the corrupted image, but also the learned inter-image reasoning priors characterizing the generalizable semantic distribution patterns between similar images. In particular, the proposed *GM-SRM* first pre-learns a generative memory from the whole training data to explicitly learn the distribution of different semantic patterns. Then the learned memory are leveraged to retrieve the matching semantics for the current corrupted image to perform semantic reasoning during image inpainting. While the encoder-decoder network is used for guaranteeing the pixel-level content consistency, our generative priors are favorable for performing high-level semantic reasoning, which is particularly effective for inferring semantic content for large corrupted area. Extensive experiments on Paris Street View, CelebA-HQ, and Places2 benchmarks demonstrate that our *GM-SRM* outperforms the state-of-the-art methods for image inpainting in terms of both visual quality and quantitative metrics.

Index Terms—Image inpainting, generative memory, image synthesis, semantic reasoning.

I. INTRODUCTION

Image inpainting aims to infer the content of missing regions given a corrupted image. It serves as the essential technical step for many image processing tasks, such as old-photo restoration [1] and image edit [2]. Image inpainting is challenging in that the predicted content for the missing regions is required to be consistent with the known regions at both the pixel level and the semantic level. Despite the significant progress for image inpainting made in recent years,

image inpainting for images with large corrupted area remains an extremely challenging task.

Compared to the traditional techniques [3], [4] for image inpainting, deep learning methods based on convolutional neural networks have boosted the performance of image inpainting substantially due to its excellent capability of feature learning. Most existing deep learning methods [5], [6] for image inpainting follow the encoder-decoder framework, in which an encoder is designed to extract useful features from the known regions of the input image, and then the decoder infers the content of the corrupted regions based on the encoded features. While such straightforward modeling way performs well for image inpainting with small corrupted regions, it can hardly deal with images with large corrupted area. This is largely because such methods focus on modeling the image-to-image mapping between the input corrupted image and the groundtruth intact image during training, whereas the reasoning from the known regions to the corrupted regions is not explicitly learned.

To fully take advantage of the information of known regions for inferring the content of the corrupted regions, many methods for image inpainting [7]–[22] aim to learn effective prior knowledge from the known regions of the input image, and then infer the content of the corrupted regions based on such prior knowledge. These methods can be classified into four categories by the way of leveraging priors from the known regions. The first type of methods [7]–[9] employs attention mechanism to learn a confidence mask for the corrupted regions based on prior information of known regions, and then infers the pixel values of corrupted regions in a progressive propagation manner from high-confidence area to low-confidence area. The second way of leveraging priors of known regions [10]–[13] is to predict the structure information of the corrupted regions from known regions first, namely the edge (high-frequency) information, then infer the detailed texture information under the guidance of the structure information. The prior knowledge of known regions can also be learned using VAE [23] by estimating the distribution of pixel values of the whole image based on known regions [21], [22]. Then the values of the corrupted regions can be inferred based on the obtained distribution. The last type of methods [14]–[20] performs constraints on the semantic consistency between the known regions and the predicted corrupted regions.

All aforementioned methods focus on modeling the image-to-image mapping and implicitly learn semantic priors from the known regions of current input image to infer the content of the corrupted regions. Such modeling paradigm suffers

Xin Feng, Wenjie Pei, Fengjun Li, Fanglin Chen and Guangming Lu are with the Department of Computer Science, Harbin Institute of Technology at Shenzhen, Shenzhen 518057, China (e-mail: fengx_hit@outlook.com; wenjiecoder@outlook.com; 20s151173@stu.hit.edu.cn; chenfanglin@hit.edu.cn; luguangm@hit.edu.cn).

David Zhang is with the School of Science and Engineering, The Chinese University of Hong Kong at Shenzhen, Shenzhen 518172, China (e-mail: davidzhang@cuhk.edu.cn)

Manuscript received Xxxx xx, xxxx; revised Xxxx xx, xxx.

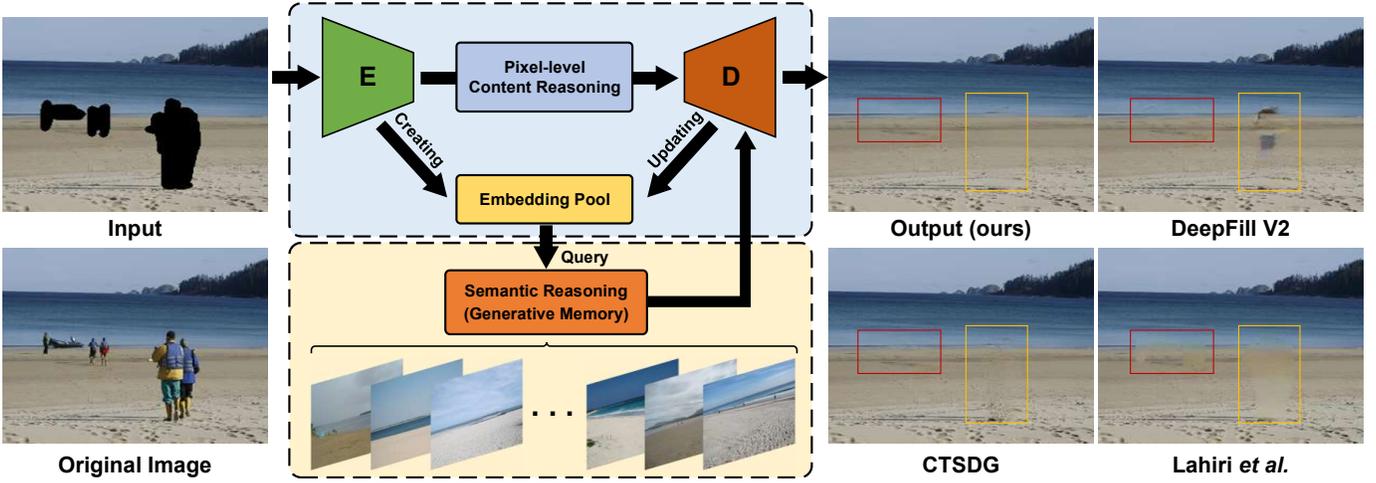


Fig. 1. Given a corrupted image, our *GM-SRM* infers the content for the corrupted area by leveraging two types of content reasoning by the decoder (**D**): 1) pixel-level content reasoning by the encoder-decoder framework, and 2) high-level semantic reasoning by the pre-trained generative memory. Our *GM-SRM* obtains the semantic priors from the whole training data by learning a generative memory, and the learned priors are queried by the embeddings of known information for high-level semantic reasoning. Consequently, our model is able to inpaint more reasonable content than other state-of-the-art methods, such as DeepFill V2 [8] which implicitly learns the inter-image semantic priors, CTSDG [24] which uses a two-stream network to facilitate image inpainting by the structure-constrained texture synthesis and texture-guided structure reconstruction, or Lahiri *et al.* [25] which explicitly learns inter-image priors by simply pre-training the decoder as an image synthesizing network.

from two potential limitations: 1) Such one-to-one learning paradigm tends to overfit each single training pair of images; 2) The inter-image prior knowledge about the general distribution patterns of visual semantics, which can be transferred across images sharing similar semantics, is not explicitly exploited. For instance, similar objects or scenes from different but similar images always share similar textures. Besides, the spatial/structural associations among objects/sub-regions can also be distilled as priors to generalize across similar semantics. Lahiri *et al.* [25] made the first heuristic attempt to explicitly learn such semantic prior by pre-training a plain Generative Adversarial Net (GAN) [26] as the decoder to reconstruct the corrupted image. The pre-trained GAN captures the general semantic distribution, which is used to predict the content of the corrupted regions during decoding. However, such a straightforward way suffers from three main weaknesses for handling large corruptions. 1) The pre-trained decoder is hard to synthesize exquisite details only with semantic reasoning; 2) The manner of semantic query solely relies on one highly coupled embedding to control all types of semantics for different feature scales; 3) Vanilla GAN has limited generative performance to provide complex semantics for largely corrupted regions. As a result, this method only shows effectiveness on images with small corrupted regions, but cannot handle cases with large corrupted regions.

In this paper, we propose the **Generative Memory-guided Semantic Reasoning Model (GM-SRM)**, which infers the content of corrupted regions based on not only the known regions of the corrupted image, but also the learned inter-image reasoning priors characterizing the generalizable semantic distribution patterns between similar images. Similar to most existing methods, the pixel-level reasoning from the known regions of the input image learns the image-to-image mapping following the typical encoder-decoder framework.

To obtain more comprehensive semantic priors, our *GM-SRM* mines the general distribution of semantic patterns that can be transferred across images sharing similar semantic distributions. Specifically, our proposed *GM-SRM* first pre-learns a generative memory from the whole corpus of training data to capture the semantic distributions in a global view. Then the learned memory is leveraged to guide the process of image inpainting: it is favorable for performing high-level semantic reasoning while the typical encoder-decoder framework focuses on guaranteeing mainly the low-level (like pixel-level) semantic consistency. As shown in Figure 1, our model is able to synthesize more reasonable content for the corrupted regions than competing methods. The main contributions of our *GM-SRM* are summarized as follows:

- A novel image inpainting framework is proposed to learn both pixel-level content reasoning and semantic reasoning by the generative prior knowledge. It employs the inter-image priors from other images sharing similar semantic distributions to facilitate inpainting reasonable content in the corrupted area.
- We design a GAN-based generative memory to learn the generative prior knowledge about the general distribution of semantic patterns, which can be seamlessly integrated into the typical encoder-decoder framework.
- We present the Conditional Stochastic Variation mechanism to learn the texture distribution of the known regions, thereby synthesizing rich yet semantically reasonable textures during image inpainting.
- Extensive experiments on three benchmarks of image inpainting, including Paris Street View, CelebA-HQ, and Places2, demonstrate both the quantitative and qualitative superiority of our proposed *GM-SRM* over other state-of-the-art methods for image inpainting, especially in scenarios with large corrupted area.

II. RELATED WORK

There is a substantial amount of work on image inpainting. In this section, we review the most typical methods. Mainstream image inpainting methods can be roughly divided into two categories: non-learning based methods and learning-based methods.

A. Non-learning based Image Inpainting

Earlier research on image inpainting employs the content of existing regions to directly fill in the missing regions, which can be divided into two categories, diffusion-based methods, and patch-based methods. The diffusion-based methods [3], [27] extend contextual pixels from the boundary to the hole along the isophote direction. By imposing various boundary conditions, diffusion-based methods have shown promising performance in filling reasonable content in images with small corruptions. However, diffusion-based inpainting methods often cause boundary-related blurriness and synthesize implausible semantics.

Aiming to solve such problems, the patch-based methods [4], [28]–[31] copy and paste the most similar patch from known regions or external images to replace the corrupted regions. Thus, many patch-based methods contribute to designing the optimal algorithms for patch selection. For instance, Criminisi *et al.* [30] propose an exemplar-based method for texture synthesis and calculate the confidence value for determining the inpainting order of the corrupted region. Barnes *et al.* [4] propose the PatchMatch algorithm for quickly selecting approximate nearest neighbor matches between image patches. Later, Jin *et al.* [31] propose a patch-sparsity-based image inpainting algorithm through facet deduced directional derivative. This type of inpainting methods can produce high-quality results on images with repeated textures. Nevertheless, it generates significant artifacts in the output due to the lack of global semantic understanding and generalization. In summary, non-learning based inpainting methods are difficult to satisfy the demand of practical application nowadays.

B. Learning-based Image Inpainting

With the great success of deep convolutional neural networks (CNNs) in various computer vision tasks, recent research leverages the encoder-decoder CNN framework to facilitate image inpainting. Main CNN-based methods for image inpainting employ known regions as priors to infer missing pixels from the outside to the inside. Pathak *et al.* [6] firstly propose a context-encoder framework to restore corrupted images in latent space. To further cope with irregular corruption, Liu *et al.* introduce the partial convolution (PConv) [7] to avoid blurry artifacts caused by traditional convolution. Inspired by attention mechanism [32], Yu *et al.* propose gated convolution (DeepFill V2) [8] to gradually learn the soft mask rather than the hard mask in PConv, enhancing the inference reliability. Yu *et al.* believe features of known regions and inferred regions should be normalized respectively, thus proposing the region normalization (RN) [33] to improve inpainting performance. Li *et al.* [9] propose a recurrent

framework to iteratively predict the features of corrupted regions and merge them to infer the corrupted content.

Structural prior-guided inpainting methods. To synthesize more plausible results, some research employs structure information as prior knowledge to reconstruct unknown regions. For instance, Liu *et al.* [18] improve the attention mechanism by computing coherent semantic relevance to infer missing content. Nazeri *et al.* [10] predict edge maps to guarantee structure correctness and reuse edge maps to facilitate inpainting plausible content. Yang *et al.* [12] further integrate structural information with semantic information to enhance stability of image inpainting. Xiong *et al.* [34] propose a foreground-aware framework that restores foreground objects of input image. Liu *et al.* [35] propose a mutual encoder-decoder framework to equalize the restoration of structure features and texture features. Xu *et al.* [36] design a two-step framework, which first generates edges inside the missing areas, and then generates inpainted image based on the edges. Wang *et al.* [37] introduce a structure-guided method for video inpainting, which integrates scene structure, texture, and motion to complete the missing regions. Recently, Guo *et al.* [24] propose a two-stream network which casts image inpainting into two collaborative subtasks, structure-constrained texture synthesis and texture-guided structure reconstruction. Though structural prior-based methods have improved the performance to restore large corruption, the error of predicting structural prior often leads to synthesizing unreasonable semantics and reducing the quality of filled content.

Generative prior-guided inpainting methods. Learning the distribution of semantic patterns by generative models leads to promising performance on image inpainting. Promoted by excellent generative representation ability of generative adversarial nets [26] and variation auto-encoders [23], some research attempts to employ them to learn generative prior for improving the quality of synthesized images, and generate diverse content to fill in corrupted regions. Li *et al.* [38] propose a generative adversarial model to learn distribution representation of natural images, and thus can synthesize reasonable content from random noise. Zheng *et al.* [21] introduce the VAE to enhance the diversity of generated content by posterior probability maximization. Later, Zhao *et al.* [22] propose a VAE-based framework to improve the diversity of filled content by cooperation with randomly selecting instance images.

To explicitly learn the generative prior, Lahiri *et al.* [25] pre-train a vanilla gan model as generative prior of the whole dataset, and then search effective information from the distribution space. Kelvin *et al.* [39] propose the GLEAN framework to learn generative priors for synthesizing realistic textures in the image super-resolution task. However, current generative prior-guided methods fail to restore images with complex content, or large corruption due to the weakness of low-level content inference. Hence, our *GM-SRM* combines the pixel-level content reasoning in the typical encoder-decoder framework and high-level semantic reasoning by learning the generative prior from the proposed generative memory respectively.

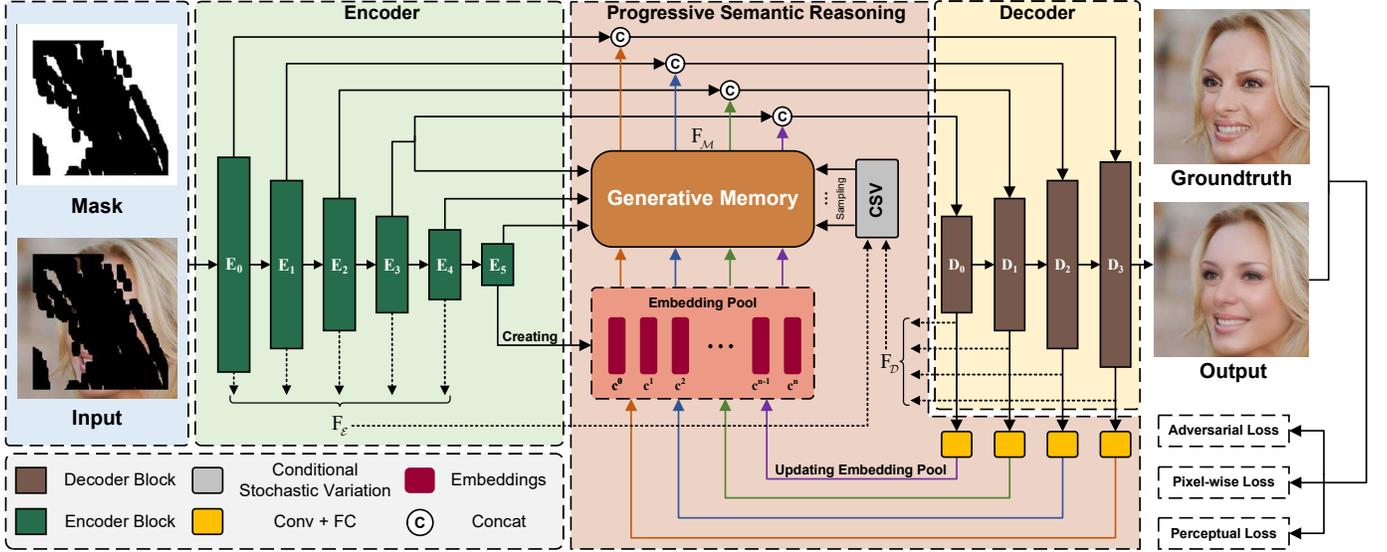


Fig. 2. The architecture of the proposed Generative Memory-guided Semantic Reasoning Model ($GM-SRM$). It consists of three modules: the encoder \mathcal{E} , the generative memory \mathcal{M} , and the decoder \mathcal{D} . The proposed $GM-SRM$ infers the corrupted regions by two types of content reasoning: 1) the generative priors on the distribution of various semantic patterns, which are learned by the proposed generative memory \mathcal{M} from the whole training corpus, and 2) pixel-level content reasoning by the encoder-decoder framework assisted by the semantics from the generative memory. For each input corrupted image, the semantic priors are retrieved from the pre-trained generative memory \mathcal{M} by searching for semantically matched features for the corrupted area using the encoded embedding c , representing known information. The semantic priors by generative memory favor the high-level semantic reasoning for large corruption.

III. GENERATIVE MEMORY-GUIDED SEMANTIC REASONING MODEL

Given a corrupted image, the goal of image inpainting is to infer the content of the missing (corrupted) regions, which is required to be consistent with the known regions at both the pixel level and the semantic level. To this end, we propose the **Generative Memory-guided Semantic Reasoning Model** ($GM-SRM$), which learns a generative memory from the corpus of training data to mine the general distribution patterns of visual semantics within images. Such learned generative memory is then leveraged to guide the process of image inpainting by performing semantic reasoning for the missing regions.

In this section, we will first introduce the overall framework of image inpainting by our model, and then we will elaborate on how to construct the generative memory and perform semantic reasoning to infer the content of the missing regions. Next, we present the conditional stochastic variation, a technique proposed to synthesize rich yet semantically reasonable details during image inpainting. Finally, we will show how to perform supervised learning to train our proposed $GM-SRM$.

A. Overall Framework for Image Inpainting

As illustrated in Figure 2, our $GM-SRM$ mainly consists of three modules: the encoder \mathcal{E} , the generative memory \mathcal{M} , and the decoder \mathcal{D} . The whole model follows the classical encoder-decoder framework to perform image inpainting. The encoder \mathcal{E} is responsible for extracting useful visual features for the known regions of the input corrupted image and meanwhile generating a latent embedding which serves as a query embedding to search for semantically matched visual features from the generative memory \mathcal{M} . The generative memory \mathcal{M} is constructed to learn the general distribution patterns of visual

semantics from training data and help infer the content of the missing regions that is semantically consistent with the known regions in the input image. Thus the generative memory \mathcal{M} summarizes the prior knowledge about the distribution of high-level visual semantics within images from a global view (covering the whole corpus of training data). Such a generative prior knowledge enables our model to explicitly infer more semantically reasonable content for the missing regions than the typical way, which trains the deep model to implicitly learn the consistency between similar semantics. The retrieved semantically matched features from the memory \mathcal{M} , together with the encoded features from the encoder \mathcal{E} , are fed into the decoder \mathcal{D} to synthesize the intact image.

Encoder. Given a corrupted image I with a mask M (with the same size as I) indicating the missing regions, the encoder \mathcal{E} extracts visual features for the known regions of I by:

$$\mathbf{F}_{\mathcal{E}} = \mathcal{E}(I, M), \quad (1)$$

where $\mathbf{F}_{\mathcal{E}}$ denotes the obtained encoded features. As shown in Figure 2, the encoder \mathcal{E} is designed by repetitively stacking a basic residual block to iteratively refine visual features and meanwhile downsample the resolution of feature maps. Such residual block consists of stride-2 convolution layers along with ReLU layers. Besides, we utilize instance normalization (IN) [40] and channel attention (CA) mechanism [41] to optimize the feature learning process. Specifically, the i -th basic residual block in the encoder \mathcal{E} can be mathematically formulated as:

$$\begin{aligned} \mathbf{F}^e &= \text{ReLU}(\text{IN}(\text{Conv}(\mathbf{F}_{\mathcal{E}}^{i-1}))), \\ \mathbf{F}_{\mathcal{E}}^i &= \mathbf{F}^e + \text{CA}(\text{ReLU}(\text{IN}(\text{Conv}(\mathbf{F}^e))), \end{aligned} \quad (2)$$

where $\mathbf{F}_{\mathcal{E}}^i$ is the output feature maps by the i -th basic residual

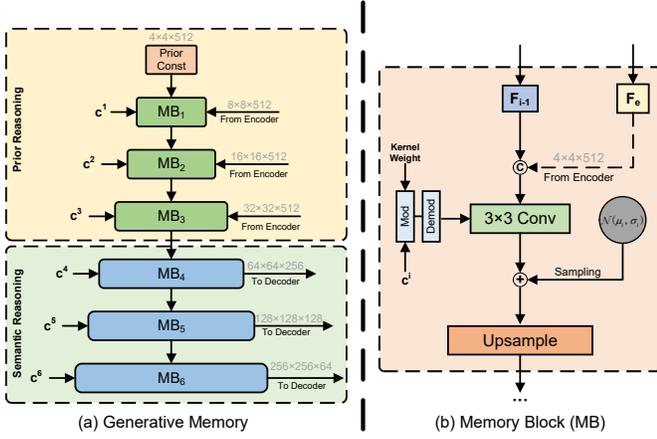


Fig. 3. The structure of the proposed generative memory in *GM-SRM*. c^i denotes the latent embedding, F_e and F_{i-1} represent the features from the encoder and last memory block respectively, and \odot is the feature concatenation operation.

block of \mathcal{E} , and \mathbf{F}^e is the intermediate features in the residual connection.

Besides learning visual features for the known regions of the input corrupted image, the encoder \mathcal{E} also generates the latent embeddings to establish an updating embedding pool based on the encoded features, and the query embeddings are the representation of progressively renovated known information. The embedding pool performs feature query by embeddings on the generative memory \mathcal{M} to search for semantically matched features for the missing regions:

$$\mathbf{c} = f_c(\mathbf{F}_{\mathcal{E}}), \quad (3)$$

where the latent embedding \mathbf{c} is a vector obtained by a non-linear mapping function f_c performed by a convolutional layer and a fully connected layer from the encoded features $\mathbf{F}_{\mathcal{E}}$.

Generative Memory. The generative memory \mathcal{M} is designed as a distribution learner by a StyleGAN-based generative network to infer semantics in corrupted regions from the latent vectorial embedding distilled by the features of known regions. The main goal of the generative memory \mathcal{M} is to synthesize the matched feature maps that are semantically consistent with the known regions, under the supervised learning for the whole *GM-SRM* model. It takes the updating latent embedding \mathbf{c} as well as the encoded features $\mathbf{F}_{\mathcal{E}}$ as input and generates the corresponding feature maps. As a result, the generative memory \mathcal{M} learns a mapping from the latent coding space to the reasoned semantic features. Formally, the reasoned semantic features $\mathbf{F}_{\mathcal{M}}$ corresponding to the latent embedding \mathbf{c} by the generative memory \mathcal{M} is queried by:

$$\mathbf{F}_{\mathcal{M}} = \mathcal{M}(\mathbf{c}, \mathbf{F}_{\mathcal{E}}). \quad (4)$$

We will elaborate on the model structure of the generative memory \mathcal{M} and describe how to perform progressive semantic reasoning based upon \mathcal{M} in the following Section III-B.

Decoder. The queried features $\mathbf{F}_{\mathcal{M}}$ can be considered as the inferred semantic features as generative priors by the memory \mathcal{M} for the missing regions of the input corrupted image. Both $\mathbf{F}_{\mathcal{M}}$ and the encoded features $\mathbf{F}_{\mathcal{E}}$ for the known regions as

Algorithm 1 Inference Process of *GM-SRM*

Input: $I_{in}, M_{in}, \mathcal{M}, E_i, D_i, \mathbf{c}$.

I_{in} : Input corrupted image;

M_{in} : Mask of corrupted region;

\mathcal{M} : Pre-trained generative memory;

E_i : i -th layer of encoder E ;

D_i : i -th layer of decoder D ;

\mathbf{c} : latent embeddings for memory reasoning.

Output: \hat{I}_{out} , restored image from input I_{in} .

- 1: # Initializing input;
- 2: $\mathbf{F}_{\mathcal{E}}^1 = \text{Concat}[I_{in}, M_{in}]$
- 3: # Encoding, n_e is the layer number of encoder;
- 4: **for** $i = 1 \rightarrow n_e - 1$ **do**
- 5: $\mathbf{F}_{\mathcal{E}}^{i+1} = E_i(\mathbf{F}_{\mathcal{E}}^i)$
- 6: # Estimating the mean μ and the standard variation σ^2 ;
- 7: $\mu_i, \sigma_i^2 = \text{Mean}(\text{Split}(\text{Conv}(\mathbf{F}_{\mathcal{E}}^i)))$
- 8: **end for**
- 9: # Initializing query embedding pool;
- 10: $\mathbf{c}_1 = \text{Linear}(\text{Conv}(\mathbf{F}_{\mathcal{E}}^{n_e}))$
- 11: $\mathbf{F}_{\mathcal{D}}^1 = \mathbf{F}_{\mathcal{D}}^{n_e}$
- 12: # Decoding, n_d is the layer number of decoder;
- 13: **for** $j = 1 \rightarrow n_d - 1$ **do**
- 14: $\mathbf{F}_{\mathcal{D}}^{j+1} = D_j(\mathbf{F}_{\mathcal{D}}^j)$
- 15: # Updating the latent embedding pool;
- 16: $\mathbf{c}_{j+1} = \text{Linear}(\text{Conv}(\mathbf{F}_{\mathcal{D}}^{n_d-j})) + \mathbf{c}_j$
- 17: Sampling noise $_{j+1}$ from distribution $\mathcal{N}(\mu_{n_d-j}, \sigma_{n_d-j}^2)$
- 18: # Semantic reasoning in generative memory;
- 19: $\mathbf{F}_{\mathcal{M}} = \mathcal{M}(\mathbf{F}_{\mathcal{E}}^{n_d-j}, \mathbf{c}_{j+1}, \text{noise}_{j+1})$
- 20: # Inferring by known content and memory query;
- 21: $\mathbf{F}_{\mathcal{D}}^{j+1} = \text{Concat}[\mathbf{F}_{\mathcal{D}}^{j+1}, \mathbf{F}_{\mathcal{M}}]$
- 22: **end for**
- 23: # Synthesizing final restored output image.
- 24: $\hat{I}_{out} = \text{OutConv}(\mathbf{F}_{\mathcal{D}}^{n_d})$

the pixel-level inference cues are fed into the decoder \mathcal{D} to synthesize the output intact image \hat{I} :

$$\hat{I} = \mathcal{D}(\mathbf{F}_{\mathcal{M}}, \mathbf{F}_{\mathcal{E}}). \quad (5)$$

The decoder \mathcal{D} is built in the similar way as the encoder: repetitively stacking a basic residual block to progressively synthesize the final intact image. One major difference from the encoder is that the decoder utilizes the bi-linear interpolation f_{interp} in each basic block to upsample feature maps gradually. Specifically, the i -th basic block of the decoder \mathcal{D} is formulated as:

$$\begin{aligned} \mathbf{F}^d &= \text{ReLU}(\text{IN}(\text{Conv}(f_{\text{interp}}(\text{Concat}[\mathbf{F}_{\mathcal{M}}^i, \mathbf{F}_{\mathcal{E}}^i, \mathbf{F}_{\mathcal{D}}^{i-1}])))), \\ \mathbf{F}_{\mathcal{D}}^i &= \mathbf{F}^d + \text{CA}(\text{ReLU}(\text{IN}(\text{Conv}(\mathbf{F}^d)))), \end{aligned} \quad (6)$$

where $\mathbf{F}_{\mathcal{D}}^i$ is the output feature maps of the i -th basic block in the decoder \mathcal{D} , and \mathbf{F}^d is the intermediate features in the residual connection. The whole inference process of our *GM-SRM* is summarized in Algorithm 1.

B. Semantic Reasoning by Generative Memory

The generative memory \mathcal{M} is proposed to employ adversarial learning to explicitly learn the distribution of various semantic patterns, and facilitate semantic reasoning from limited information. Our *GM-SRM* then leverages the learned prior

knowledge about the semantic distributions by the generative memory to perform high-level semantic reasoning for the missing regions given the information of known regions.

Construction of Generative Memory. As shown in Figure 3, we opt for the generative model of revised StyleGAN [42] due to its excellent performance of synthesizing target image controlled by decoupled latent embeddings via the improved AdaIN operation with the proposed demodulation [42]. Inspired by previous effective attempts to learn generative priors in other tasks [39], we pre-train the generative memory on the whole training data in the similar training way as plain StyleGAN, except that during training GM-SRM we replace the mapping network in the StyleGAN with the mapping function in Equation 3 to learn the underlying correspondence between the latent embedding (representing the semantics of the existing regions) and the reasoned semantic features for the missing regions. After sufficient training on the whole corpus of training data, the generative memory \mathcal{M} is expected to learn a good mapping between the latent space and visual semantic features, thereby summarizing the prior knowledge about the distribution of various semantic patterns from a global view. Once the memory \mathcal{M} is pre-trained, its parameters \mathcal{M} are frozen, and \mathcal{M} is used to perform semantic reasoning as generative prior knowledge during the training of the whole GM-SRM, i.e., to infer the semantic features for the missing regions of the input image based on the latent embedding of the known regions.

Progressive Semantic Reasoning. We leverage the pre-trained generative memory \mathcal{M} to perform progressive semantic reasoning and thereby assist the decoder \mathcal{D} to conduct image inpainting. As shown in Equation 6, the decoder \mathcal{D} employs multiple basic blocks to synthesize the intact image in a coarse-to-fine manner by progressively expanding the resolution of the generated feature maps. Accordingly, the generative memory infers the semantics for the missing regions for each resolution of feature maps to keep pace with the synthesizing process of the decoder \mathcal{D} . As a result, the memory \mathcal{M} and the decoder \mathcal{D} are able to perform image inpainting collaboratively in an iterative way: the inferred semantic features by \mathcal{M} are fed into \mathcal{D} to provide high-level semantic cues for decoding the same resolution of feature maps, while the decoded feature maps are in turn provided for \mathcal{M} to update the embedding pool for the next scale of semantic inferring in larger resolution:

$$\mathbf{c}^i = f_c(\mathbf{F}_{\mathcal{D}}^{i-1}) + \mathbf{c}^{i-1}, \quad (7)$$

where f_c corresponds to the same mapping function as in Equation 3. \mathbf{c}^i denotes the latent embedding to retrieve the semantic features $\mathbf{F}_{\mathcal{M}}^i$ from \mathcal{M} , which is prepared for decoding $\mathbf{F}_{\mathcal{D}}^i$ by the i -th block of \mathcal{D} . As a result, the newly inferred features for the missing regions can be incorporated into the updated latent embedding to predict the semantic features for the next basic block of \mathcal{D} . The image inpainting is performed in such a progressively inferring mechanism to predict the content of missing regions by the generative memory \mathcal{M} and the decoder \mathcal{D} collaboratively.

Compared to the typical methods for image inpainting that implicitly learns the distribution of semantics by modeling

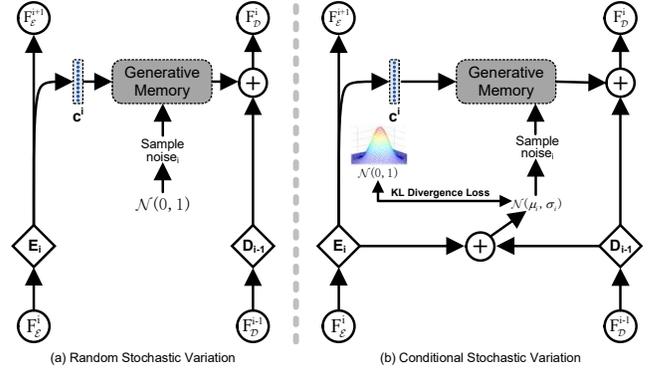


Fig. 4. Unlike typical StyleGAN that simply samples noise randomly from a standard normal distribution (a), our proposed Conditional Stochastic Variation mechanism samples noise conditioned on the encoded features $\mathbf{F}_{\mathcal{E}}^i$ and the decoded features $\mathbf{F}_{\mathcal{D}}^{i-1}$.

the image-to-image mapping, the key benefit of our GM-SRM is that the generative memory explicitly learns the semantic distributions as prior knowledge in a global view from the whole corpus of training data, which can be generalized across different images sharing similar semantic distributions. Such learned prior knowledge about the general semantic distributions enables our model to perform high-level semantic reasoning to infer more semantically reasonable content for the missing regions. As a result, our method is particularly effective in scenarios with large corrupted area.

C. Conditional Stochastic Variation

To simulate the stochastic appearance variations in synthesized images that do not violate the correct semantics and enrich the texture details, random noise is introduced into synthesizing process in StyleGAN in multiple generative layers. Inspired by such Stochastic Variation scheme in StyleGAN, we also incorporate noise as input during image synthesis in the generative memory \mathcal{M} . Unlike the StyleGAN that synthesizes images from random latent embedding without conditioning on any input information, the generated features by \mathcal{M} of our GM-SRM model are required to be semantically consistent with the known regions. Thus we propose the Conditional Stochastic Variation mechanism, which introduces noise conditioned on the semantics of the known regions into the synthesis process of the generative memory \mathcal{M} to enrich the appearance details of synthesized images while keeping the semantics correct.

Noise is introduced into each synthesizing layer of \mathcal{M} to ensure the texture diversity of generated features in each resolution. As shown in Figure 4, the noise in each layer is generated conditioned on the encoded features $\mathbf{F}_{\mathcal{E}}$ by \mathcal{E} in the current layer and the decoded features $\mathbf{F}_{\mathcal{D}}$ by \mathcal{D} in the previous layer. Specifically, a normal distribution is predicted from $\mathbf{F}_{\mathcal{E}}$ and $\mathbf{F}_{\mathcal{D}}$ by a convolution layer and a fully connected layer (FC), then the conditional noise is sampled from such predicted normal distribution. For instance, the noise for the i -th layer of \mathcal{M} is generated by the proposed Conditional Stochastic Variation Mechanism by:

$$\begin{aligned} \mu_i, \sigma_i &= \text{FC}(\text{Conv}(\mathbf{F}_{\mathcal{E}}^i, \mathbf{F}_{\mathcal{D}}^{i-1})), \\ \text{noise}_i &= \text{Sample}(\mathcal{N}(\mu_i, \sigma_i^2)), \end{aligned} \quad (8)$$

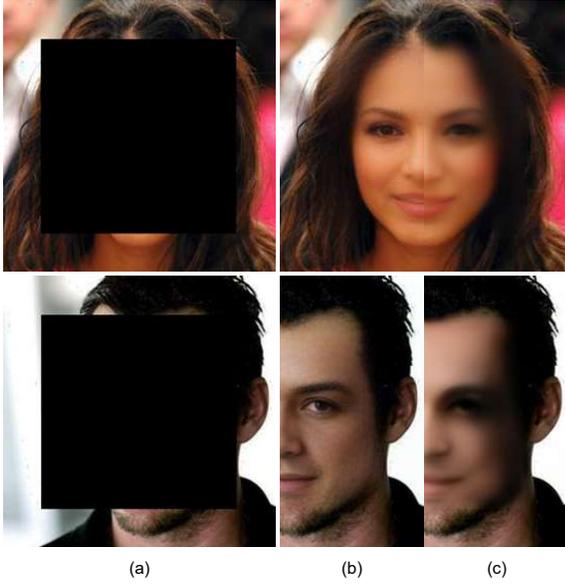


Fig. 5. Comparison of inpainting results between our *GM-SRM* using random noise and our Conditional Stochastic Variation mechanism. (a) Corrupted input images. (b) Inpainted images using the conditional noise by the proposed Conditional Stochastic Variation mechanism. (c) Inpainted images using plain stochastic variation.

where μ_i and σ_i are the mean and standard variation for the predicted normal distribution for the i -th synthesis layer conditioned on the encoded features $F_{\mathcal{E}}^i$ and the decoded features $F_{\mathcal{D}}^i$. To encourage the generated noise distribution to be close to the standard normal distribution to ease the noise-sampling process, we leverage KL divergence to guide the parameter learning in Equation 8, which is similar to VAE [23]:

$$\mathcal{L}_{\text{KL}} = \sum w_i \cdot \text{KL}(\mathcal{N}(\mu_i, \sigma_i^2) \mid \mathcal{N}(0, 1)), \quad (9)$$

where w for different layers are the weights to balance between different synthesis layers, and KL denotes the KL divergence between two distributions. Figure 5 shows two examples that illustrate the comparison between the typical Stochastic Variation mechanism performed by random noise and our proposed Conditional Stochastic Variation Mechanism. These examples reveal that the conditional noise generated by our Conditional Stochastic Variation leads to more consistent semantic content in the missing regions, while keeping rich texture details. In contrast, the random noise tends to result in blurred texture due to completely random nature of introduced noise.

D. Supervised Parameter Learning

Our proposed *GM-SRM* is trained in two steps: the generative memory \mathcal{M} is first pre-trained to learn the prior knowledge about the general distribution patterns of visual semantics. Then the whole *GM-SRM* is trained for image inpainting while keeping the parameters of \mathcal{M} frozen. Since the generative memory \mathcal{M} is trained in the similar way as the training process of StyleGAN, we explicate how to perform supervised learning to train the whole *GM-SRM*.

We employ four types of loss functions to train our *GM-SRM*. Apart from the loss function in Equation 9, the other

three loss functions are presented below:

- **Pixel-wise L1 Reconstruction Loss**, which focuses on pixel-level measurement of known regions and corrupted regions between groundtruth and synthesized images respectively:

$$\begin{aligned} \mathcal{L}_{\text{L1}} &= \|\hat{I} - I_{GT}\|_1, \\ \mathcal{L}_{\text{rec}} &= \mathcal{L}_{\text{known-L1}} + \gamma \mathcal{L}_{\text{corrupted-L1}}, \end{aligned} \quad (10)$$

where \hat{I} is the synthesized image by our *GM-SRM*, and I_{GT} are the corresponding ground-truth image. γ is a hyper-parameter to balance between losses, which is tuned to be 10 on a held-out validation set.

- **Perceptual Loss** [43], which aims to minimize the semantic difference between restored image and the ground-truth image in deep feature space:

$$\mathcal{L}_{\text{perc}}(\hat{I}, I_{GT}) = \sum_{l=1}^L \frac{1}{C_l H_l W_l} \|f_{\text{vgg}}^l(\hat{I}) - f_{\text{vgg}}^l(I_{GT})\|_1, \quad (11)$$

where $f_{\text{vgg}}^l(\hat{I})$ and $f_{\text{vgg}}^l(I_{GT})$ are the extracted feature maps, normalized by feature dimensions $C_l \times H_l \times W_l$, for the generated image \hat{I} and the ground-truth image I_{GT} respectively from the l -th convolution layer of the pre-trained VGG-19 network [44].

- **Conditional Adversarial Loss**, which encourages the synthesized image \hat{I} to be as realistic as the ground-truth image I_{GT} . We employ spectral normalization [45] in discriminator to stabilize the training process:

$$\mathcal{L}_{\text{adv}} = -\mathbb{E}_{B \sim \mathbb{P}_{\text{GM-SRM}}} [D^{sn}(G(I))], \quad (12)$$

where D^{sn} is the spectral normalized discriminator. It is trained by:

$$\begin{aligned} \mathcal{L}_{D^{sn}} &= \mathbb{E}_{B_{GT} \sim \mathbb{P}_{\text{data}}} [1 - D^{sn}(G(I))] \\ &\quad + \mathbb{E}_{B \sim \mathbb{P}_{\text{GM-SRM}}} [D^{sn}(G(I))]. \end{aligned} \quad (13)$$

Overall, the whole *GM-SRM* is trained by:

$$\mathcal{L} = \lambda_1 \mathcal{L}_{\text{rec}} + \lambda_2 \mathcal{L}_{\text{perc}} + \lambda_3 \mathcal{L}_{\text{adv}} + \lambda_4 \mathcal{L}_{\text{KL}}, \quad (14)$$

where $\lambda_1, \lambda_2, \lambda_3$ and λ_4 are hyper-parameters to balance between different losses. In our experiments, we empirically set $\lambda_1=1, \lambda_2=0.1, \lambda_3=0.01$ and $\lambda_4=0.01$.

IV. EXPERIMENTS

In this section, we conduct experiments to quantitatively and qualitatively evaluate the proposed *GM-SRM*. In Section IV-A, we introduce benchmark datasets used in the experiments, evaluation metrics, and implementation details. In Section IV-B, we conduct experiments to compare our model with state-of-the-art methods for image inpainting. Finally, we perform ablation study to investigate the effectiveness of each component in our *GM-SRM* in Section IV-C.

A. Experimental Setting

Evaluation Metrics. We employ five generally used criteria as quantitative measurements to quantify the quality of the restored images in our experiments: 1) the peak signal-to-noise

TABLE I
PERFORMANCE OF DIFFERENT MODELS FOR IRREGULAR IMAGE INPAINTING ON THE PARIS STREET VIEW [46], CELEBA-HQ256 [47], AND PLACES [48] DATASETS. METRICS WITH \uparrow , HIGHER VALUE DENOTES BETTER PERFORMANCE, WHEREAS \downarrow DENOTES LOWER IS BETTER. THE BEST RESULTS OF EACH METRIC ARE HIGHLIGHTED IN BOLD.

Method	Paris Street View				CelebA-HQ				Places2				
	20-30%	30-40%	40-50%	50-60%	20-30%	30-40%	40-50%	50-60%	20-30%	30-40%	40-50%	50-60%	
PSNR \uparrow	PConv [7]	26.85	24.66	22.67	20.52	28.29	26.11	23.90	21.47	25.05	22.99	21.11	19.04
	DeepFill V2 [8]	26.29	23.99	21.86	19.79	27.81	25.52	23.19	20.64	25.15	23.00	20.94	18.68
	EdgeConnect [10]	27.62	25.39	23.19	20.86	28.53	26.14	23.53	20.38	25.75	23.64	21.66	19.46
	LISK [12]	25.04	23.10	21.15	19.00	30.22	27.63	24.95	21.51	27.04	24.68	22.35	19.77
	LBAM [19]	27.54	25.25	23.09	20.84	28.86	26.55	24.22	21.68	25.69	23.46	21.41	19.21
	E2I [36]	26.78	23.80	21.94	19.14	26.28	25.03	23.19	20.20	24.91	23.14	21.09	18.56
	CTSDG [24]	27.33	24.67	22.00	19.98	28.49	25.99	23.77	21.22	26.35	24.11	21.58	19.90
	MEDFE [35]	27.23	24.78	22.30	19.84	25.50	24.10	22.31	20.04	24.84	23.38	21.38	18.27
	Lahiri <i>et al.</i> [25]	27.26	24.71	22.76	20.58	27.96	25.89	24.73	21.26	25.54	23.51	21.60	19.41
GM-SRM (ours)	29.04	26.79	24.51	22.20	30.23	27.91	25.51	22.70	27.09	24.93	22.91	20.60	
SSIM \uparrow	PConv [7]	0.861	0.795	0.711	0.622	0.896	0.847	0.783	0.711	0.844	0.777	0.693	0.611
	DeepFill V2 [8]	0.865	0.801	0.714	0.625	0.894	0.842	0.772	0.693	0.859	0.796	0.714	0.630
	EdgeConnect [10]	0.879	0.821	0.740	0.651	0.905	0.856	0.784	0.689	0.862	0.799	0.717	0.633
	LISK [12]	0.881	0.823	0.744	0.642	0.928	0.886	0.822	0.724	0.882	0.823	0.744	0.655
	LBAM [19]	0.879	0.816	0.732	0.637	0.907	0.860	0.796	0.718	0.861	0.794	0.708	0.619
	E2I [36]	0.865	0.810	0.718	0.620	0.895	0.846	0.779	0.698	0.861	0.798	0.719	0.634
	CTSDG [24]	0.874	0.819	0.728	0.629	0.905	0.889	0.780	0.697	0.869	0.814	0.728	0.649
	MEDFE [35]	0.876	0.811	0.718	0.616	0.884	0.810	0.759	0.666	0.833	0.809	0.738	0.636
	Lahiri <i>et al.</i> [25]	0.871	0.813	0.739	0.651	0.903	0.859	0.802	0.716	0.869	0.811	0.738	0.666
GM-SRM (ours)	0.902	0.852	0.780	0.696	0.925	0.886	0.832	0.760	0.883	0.827	0.756	0.677	
NCC \uparrow	PConv [7]	0.957	0.934	0.895	0.835	0.986	0.978	0.964	0.938	0.957	0.934	0.898	0.838
	DeepFill V2 [8]	0.951	0.924	0.878	0.811	0.985	0.975	0.958	0.926	0.956	0.933	0.894	0.826
	EdgeConnect [10]	0.963	0.943	0.906	0.848	0.987	0.978	0.960	0.920	0.962	0.942	0.908	0.850
	LISK [12]	0.971	0.954	0.919	0.859	0.991	0.985	0.972	0.940	0.972	0.954	0.921	0.862
	LBAM [19]	0.962	0.941	0.904	0.845	0.987	0.980	0.966	0.940	0.962	0.941	0.906	0.848
	E2I [36]	0.956	0.929	0.881	0.820	0.986	0.979	0.961	0.928	0.959	0.936	0.896	0.830
	CTSDG [24]	0.959	0.940	0.900	0.822	0.982	0.979	0.959	0.931	0.960	0.947	0.913	0.835
	MEDFE [35]	0.961	0.937	0.889	0.814	0.976	0.967	0.950	0.917	0.956	0.939	0.906	0.813
	Lahiri <i>et al.</i> [25]	0.964	0.947	0.916	0.854	0.988	0.981	0.969	0.937	0.968	0.952	0.925	0.876
GM-SRM (ours)	0.972	0.958	0.929	0.879	0.991	0.985	0.974	0.952	0.972	0.956	0.929	0.880	
LMSE \downarrow	PConv [7]	0.016	0.025	0.037	0.051	0.007	0.011	0.017	0.026	0.018	0.027	0.038	0.053
	DeepFill V2 [8]	0.019	0.030	0.046	0.064	0.008	0.013	0.020	0.031	0.019	0.029	0.043	0.062
	EdgeConnect [10]	0.013	0.020	0.044	0.064	0.006	0.010	0.017	0.031	0.015	0.023	0.034	0.048
	LISK [12]	0.010	0.016	0.027	0.040	0.004	0.007	0.013	0.024	0.011	0.019	0.030	0.046
	LBAM [19]	0.014	0.022	0.033	0.047	0.006	0.010	0.016	0.025	0.016	0.025	0.037	0.053
	E2I [36]	0.015	0.027	0.031	0.052	0.008	0.012	0.018	0.029	0.016	0.027	0.039	0.051
	CTSDG [24]	0.012	0.022	0.035	0.048	0.009	0.010	0.019	0.028	0.015	0.024	0.035	0.052
	MEDFE [35]	0.014	0.023	0.038	0.055	0.011	0.015	0.021	0.032	0.018	0.025	0.037	0.058
	Lahiri <i>et al.</i> [25]	0.011	0.020	0.029	0.049	0.006	0.009	0.013	0.024	0.014	0.019	0.040	0.046
GM-SRM (ours)	0.009	0.015	0.023	0.034	0.004	0.006	0.011	0.018	0.011	0.017	0.025	0.037	
LPIPS \downarrow	PConv [7]	0.123	0.162	0.208	0.282	0.047	0.078	0.100	0.157	0.098	0.140	0.202	0.279
	DeepFill V2 [8]	0.118	0.154	0.191	0.269	0.054	0.080	0.117	0.164	0.079	0.119	0.178	0.256
	EdgeConnect [10]	0.081	0.114	0.174	0.210	0.045	0.067	0.106	0.165	0.075	0.107	0.164	0.279
	LISK [12]	0.072	0.098	0.168	0.209	0.039	0.051	0.082	0.134	0.076	0.110	0.161	0.247
	LBAM [19]	0.069	0.105	0.143	0.257	0.038	0.057	0.087	0.128	0.074	0.113	0.172	0.250
	E2I [36]	0.094	0.137	0.188	0.260	0.049	0.072	0.107	0.160	0.077	0.111	0.178	0.261
	CTSDG [24]	0.074	0.120	0.171	0.245	0.040	0.068	0.094	0.151	0.080	0.122	0.188	0.253
	MEDFE [35]	0.101	0.127	0.176	0.238	0.052	0.084	0.115	0.166	0.105	0.130	0.190	0.290
	Lahiri <i>et al.</i> [25]	0.095	0.114	0.165	0.286	0.046	0.069	0.101	0.149	0.080	0.124	0.180	0.259
GM-SRM (ours)	0.063	0.098	0.138	0.207	0.031	0.048	0.079	0.128	0.071	0.109	0.161	0.234	

ratio (PSNR), 2) the structural similarity index (SSIM) [49], 3) the normalized cross correlation (NCC) [50], 4) the local mean square error (LMSE) [51], and the learned perceptual image patch similarity (LPIPS) [52]. Additionally, we perform qualitative evaluation by visually comparing the restored results of randomly selected test samples by various models for different degrees of corruption in experiments. As a complement to the standard evaluation metrics, we further conduct user study to compare our results to the state-of-the-art results by human evaluation.

Datasets. We perform experiments on three benchmark datasets for image inpainting:

- **Paris Street View** [46], which is collected from street views of Paris, and we leverage its original splits, 14,900 images for training and 100 images for testing.
- **CelebA-HQ** [47], which contains 30,000 images of human face. We randomly select 3000 images as validation and testing dataset, and leverage remained 27,000 images as training dataset.
- **Places2** [48], which is composed of over 2,000,000 images

from 365 scenes. We select two full categories obtaining 80,000 images and randomly picks 1,500 images from each category as test set respectively. The remaining 74,000 images are employed as training set.

By following previous experience [21], we randomly generate irregular and regular masks for training. As for test, we leverage PConv’s [7] irregular mask and center mask in different degradation ratios.

Implementation Details. We implement our *GM-SRM* in distribution mode with 4 RTX 3090 GPUs under Pytorch framework. Adam [53] is employed for gradient descent optimization with batchsize set to be 8. The initial learning rate is set to be 2×10^{-4} and the training process takes maximally 100 epochs. In our experiments, we resize all images to make the shorter side be 320, and then crop into 256×256 . Random flipping, random cropping, and resizing are used for data augmentation.

B. Comparison with State-of-the-art Methods

We first conduct experiments to compare *GM-SRM* with state-of-the-art methods for image inpainting on three datasets, including Paris StreetView [46], CelebA-HQ [47], and Places2 [48]. In particular, we divide masks by ranging percentage of corrupted size.

Baselines. Concretely, we compare our *GM-SRM* with 1) **PConv** [7], employing partial convolution to cope with irregular corruption; 2) **DeepFill V2** [8], which proposes gated convolution to generalize PConv; 3) **EdgeConnect** [10], which first predicts edge map and then leverages predicted edge maps to facilitate restoration; 4) **LISK** [12], which incorporates structural knowledge to reconstruct corrupted image and structure maps simultaneously; 5) **LBAM** [19], introducing a learnable reverse attention mechanism to fill missing regions; 6) **E2I** [36], which designs a two-step framework to first generate edges inside the missing areas, and then generate inpainted image based on the edges; 7) **MEDFE** [35], which proposes a mutual encoder-decoder framework to reconstruct structure and textures separately, and then fuses them by feature equalization; 8) **CTSDG** [24], which proposes a two-stream network which casts image inpainting into structure-constrained texture synthesis and texture-guided structure reconstruction; 9) **Lahiri et al.** [25], which explicitly learns the generative priors by pre-training a vanilla GAN as the decoder. It should be noted that all these methods except Lahiri *et al.* implicitly learn the distribution of different semantics to perform image inpainting.

Quantitative Evaluation. Table I and Table II present the experimental results of different methods for image inpainting on three benchmarks in terms of PSNR, SSIM, NCC, LMSE, and LPIPS. To quantify model performance, we leverage two types of corrupted masks for testing, regularly center masks and irregular masks. In addition, we divide them into various corrupted ratios, *i.e.*, (20%, 30%), (30%, 40%), (40%, 50%), and (50%, 60%) for irregular mask, and 25% and 50% for center mask. For a fair comparison, we obtain restoration results from officially released source codes and pre-trained

TABLE II
PERFORMANCE OF DIFFERENT MODELS FOR REGULAR IMAGE INPAINTING ON THE PARIS STREET VIEW [46], CELEBA-HQ256 [47], AND PLACES2 [48] DATASETS. METRICS WITH \uparrow , HIGHER VALUE DENOTES BETTER PERFORMANCE, WHEREAS \downarrow DENOTES LOWER IS BETTER. THE BEST RESULTS OF EACH METRIC ARE HIGHLIGHTED IN BOLD.

Method	Paris Street View		CelebA-HQ		Places2		
	25%	50%	25%	50%	25%	50%	
PSNR \uparrow	PConv [7]	23.97	20.02	20.54	14.56	21.69	18.09
	DeepFill V2 [8]	23.32	18.98	25.45	19.86	21.27	17.54
	EdgeConnect [10]	24.91	20.32	23.68	18.83	22.39	18.48
	LISK [12]	22.28	18.32	24.28	18.02	22.52	18.54
	LBAM [19]	24.18	20.14	25.91	20.73	21.79	18.38
	E2I [36]	23.40	19.27	25.59	19.99	21.32	18.17
	CTSDG [24]	24.22	19.67	25.51	20.08	22.27	18.40
	MEDFE [35]	23.43	19.09	25.54	20.26	21.18	17.57
	Lahiri <i>et al.</i> [25]	24.79	19.92	25.51	20.50	22.35	18.36
	GM-SRM(ours)	25.86	21.41	26.97	21.75	23.33	19.51
SSIM \uparrow	PConv [7]	0.828	0.638	0.797	0.596	0.816	0.633
	DeepFill V2 [8]	0.832	0.633	0.876	0.694	0.823	0.638
	EdgeConnect [10]	0.846	0.659	0.860	0.672	0.827	0.647
	LISK [12]	0.828	0.639	0.873	0.675	0.842	0.662
	LBAM [19]	0.832	0.637	0.887	0.719	0.820	0.634
	E2I [36]	0.834	0.636	0.879	0.699	0.825	0.640
	CTSDG [24]	0.831	0.638	0.874	0.691	0.828	0.641
	MEDFE [35]	0.827	0.615	0.880	0.703	0.813	0.622
	Lahiri <i>et al.</i> [25]	0.835	0.649	0.882	0.701	0.830	0.669
	GM-SRM(ours)	0.862	0.695	0.901	0.757	0.845	0.683
NCC \uparrow	PConv [7]	0.924	0.817	0.894	0.691	0.909	0.800
	DeepFill V2 [8]	0.913	0.781	0.974	0.911	0.899	0.781
	EdgeConnect [10]	0.937	0.829	0.962	0.887	0.920	0.815
	LISK [12]	0.939	0.833	0.968	0.871	0.930	0.833
	LBAM [19]	0.928	0.820	0.976	0.926	0.917	0.820
	E2I [36]	0.918	0.794	0.977	0.918	0.903	0.792
	CTSDG [24]	0.925	0.811	0.969	0.911	0.916	0.818
	MEDFE [35]	0.919	0.782	0.974	0.918	0.902	0.781
	Lahiri <i>et al.</i> [25]	0.929	0.830	0.969	0.926	0.919	0.830
	GM-SRM(ours)	0.946	0.860	0.981	0.940	0.932	0.840
LMSE \downarrow	PConv [7]	0.024	0.051	0.023	0.044	0.027	0.054
	DeepFill V2 [8]	0.029	0.069	0.012	0.034	0.034	0.068
	EdgeConnect [10]	0.019	0.046	0.017	0.038	0.024	0.051
	LISK [12]	0.017	0.044	0.015	0.032	0.020	0.047
	LBAM [19]	0.022	0.049	0.011	0.028	0.028	0.055
	E2I [36]	0.025	0.057	0.012	0.033	0.031	0.056
	CTSDG [24]	0.024	0.050	0.013	0.032	0.025	0.052
	MEDFE [35]	0.022	0.060	0.012	0.033	0.030	0.063
	Lahiri <i>et al.</i> [25]	0.026	0.047	0.020	0.034	0.030	0.050
	GM-SRM(ours)	0.015	0.035	0.008	0.022	0.017	0.038
LPIPS \downarrow	PConv [7]	0.110	0.242	0.079	0.341	0.146	0.301
	DeepFill V2 [8]	0.106	0.239	0.056	0.166	0.133	0.277
	EdgeConnect [10]	0.099	0.219	0.072	0.185	0.122	0.263
	LISK [12]	0.096	0.226	0.067	0.298	0.149	0.287
	LBAM [19]	0.100	0.222	0.047	0.135	0.135	0.272
	E2I [36]	0.105	0.233	0.056	0.159	0.131	0.270
	CTSDG [24]	0.103	0.225	0.059	0.160	0.142	0.281
	MEDFE [35]	0.114	0.250	0.052	0.143	0.160	0.312
	Lahiri <i>et al.</i> [25]	0.117	0.277	0.053	0.154	0.151	0.325
	GM-SRM(ours)	0.093	0.218	0.045	0.132	0.120	0.263

models. According to their source code and paper, we re-implement models that do not release official pre-trained models. We leverage the same mask for each test image to evaluate the results of different methods.

As illustrated in Table I, we compare quantitative results of our *GM-SRM* with state-of-the-art methods on three benchmarks for irregular inpainting. On Paris Street View dataset [46], our *GM-SRM* outperforms other competing methods by a large margin on five metrics. As illustrated in Figure 7, our *GM-SRM* achieves significant improvements (higher than 1.3 dB PSNR) for each corruption ratio. For



Fig. 6. Visualization of inpainted images by five state-of-the-art models for image inpainting and our *GM-SRM* on three randomly selected samples from Paris Street View [46] test set. Our model is able to restore higher-quality image than other methods. Best viewed in zoom-in mode.

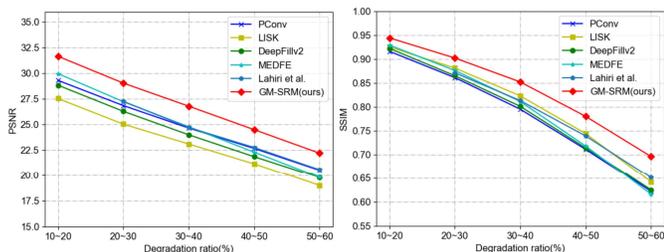


Fig. 7. Quantitative comparison on Paris Street View [46] for varying degradation ratios in terms of PSNR and SSIM. Competing methods includes PConv [7], DeepFillv2 [8], LISK [12], MEDFE [35], Lahiri *et al.* [25], and our *GM-SRM*. Though higher degradation ratio leads to faster decreasing of quantitative metrics, the *GM-SRM* performs more stably than other state-of-the-art methods.

CelebA-HQ dataset, although face structure is relatively fixed, our *GM-SRM* still outperforms other state-of-the-art methods more than 1 dB of PSNR gain for the large hole-to-image area ratio (50%,60%). The performance of our *GM-SRM* on Places2 [48] dataset is also better than other methods, especially in the cases with large corruption ratios. In sum, our *GM-SRM* achieves the best performance in terms of five quantitative metrics for irregular corruption, and outperforms other state-of-the-art methods in various corruption ratios. Comparing with the methods that implicitly learn semantic distribution, namely all methods in comparison except Lahiri *et al.*, our method also explicitly learns the semantic priors by the proposed generative memory, which enables our model to perform high-level semantic reasoning between images with similar semantic distributions. It is worth noting that our *GM-SRM* significantly outperforms Lahiri *et al.* on all datasets, which also explicitly learns the generative priors by simply pre-training the decoder. It reveals the effectiveness of the proposed generative memory in *GM-SRM* for learning semantic priors.

Table II lists experimental results for challenging center masked images. Unlike former experiments that only compare 25% corrupted ratios, we provide results for both 25% and 50% corrupted ratios from continuous center masked cor-

TABLE III
USER STUDY ON 100 RESTORED RESULTS OF PLACES2 [48], AND CELEBA-HQ256 [47]. 50 HUMAN SUBJECTS ARE PERFORMED FOR COMPARISON BETWEEN OUR *GM-SRM* AND STATE-OF-THE-ART METHODS.

Method	Share of the vote		Winning samples	
	CelebA-HQ	Places2	CelebA-HQ	Places2
DeepFill V2 [8]	3.96%	8.92%	1	3
LISK [12]	15.84%	16.60%	4	6
Lahiri <i>et al.</i> [25]	0.28%	4.68%	0	2
MEDFE [35]	19.00%	11.16%	6	7
<i>GM-SRM</i>	60.92%	58.64%	39	32

ruption to evaluate the performance of methods on image inpainting with the small and large corrupted area, respectively. As shown in Table II, *GM-SRM* outperforms other state-of-the-art methods in terms of all five metrics on three benchmark datasets. It demonstrates *GM-SRM* can synthesize more reasonable content than other methods, especially in cases with large corrupted area.

Qualitative Evaluation. We compare the restored results by our *GM-SRM* and other state-of-the-art image inpainting methods on the Paris Street View [46], CelebA-HQ [47], and Places2 [48] in Figure 6, 8 and 9. PConv [7] is specifically proposed to handle irregular corruption, and thus it restores plausible results in Figure 6(b) and Figure 8(b). However, it even cannot fill in reasonable structures when restoring large continuous corruption. DeepFillv2 [8] normalizes feature maps by the soft mask mechanism, and this relieves the limitation of PConv. Although convincing structures are synthesized in corrupted areas, undesired artifacts can still be observed in Figure 8(c). LISK [12] integrates structural information to infer corrupted regions, which may leads to confusing content for the error prediction of image structure(see Figure 8)(d). MEDFE [35] restores structure and texture at the same time. Thus, it predicts stable structure than other competing methods from Figure 8. And yet MEDFE produces unwanted blurriness and color discrepancy(see Figure 9(e)). Conversely, our *GM-SRM* synthesizes higher quality content and more reasonable

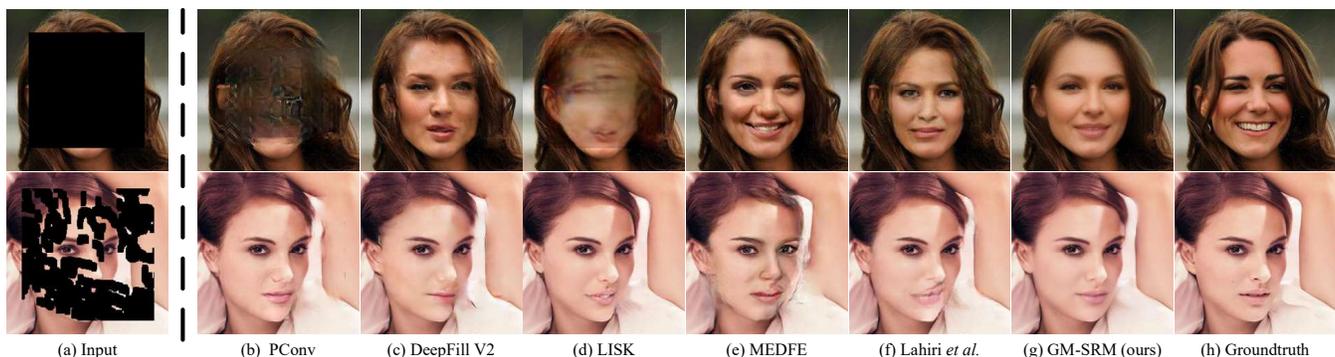


Fig. 8. Visualization of inpainted images by five state-of-the-art models for image inpainting and our *GM-SRM* on two randomly selected samples from CelebA-HQ [47] test set. Our model is able to restore higher-quality image than other methods. Best viewed in zoom-in mode.

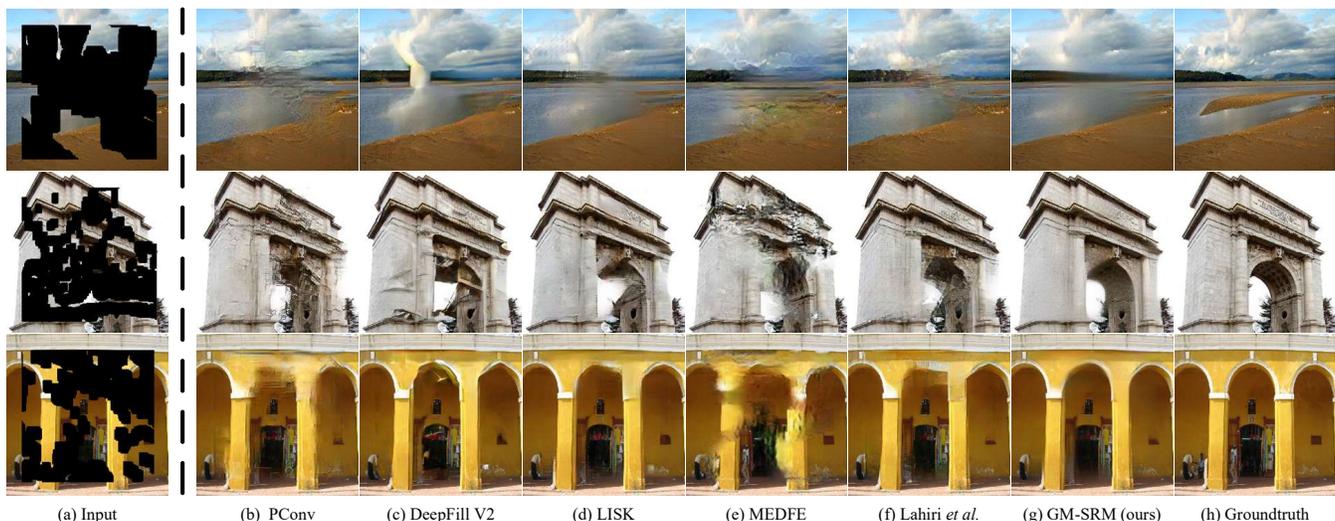


Fig. 9. Visualization of inpainted images by five state-of-the-art models for image inpainting and our *GM-SRM* on three randomly selected samples from Places2 [48] test set. Our model is able to restore higher-quality image than other methods. Best viewed in zoom-in mode.

semantic structures. From Figure 6, 8 and 9, *GM-SRM* avoids undesired blurriness, artifacts and color discrepancy to making restored images more realistic.

Compared with Lahiri *et al.* which also explicitly learns the semantic priors for image inpainting, our *GM-SRM* is able to restore content for the corrupted area that is more realistic in pixel level and more reasonable in semantic level. These qualitative comparisons illustrate the advantages of our model over Lahiri *et al.* in learning the semantic priors.

User Study. Quantitative metrics have their bias for the quality evaluation of restored images. To standardize evaluation process, we further perform user study with another four state-of-the-art methods for image inpainting, including DeepFillv2 [8], LISK [12], Lahiri *et al.* [25] and MEDFE [35]. We randomly select 50 test images of CelebA-HQ [47] and Places2 [48] respectively, and present restoration results of four methods to 50 human subjects for manual ranking of image quality. Table III lists the voting results of this user study. For the CelebA-HQ dataset, our *GM-SRM* reaches 60.92% votes among total $50 \times 50 = 2500$ rankings, which is much higher than other competing methods. In addition, we count winning samples of each method, and our *GM-SRM* wins on 39 test samples and others altogether 11 samples. As for the Places2 dataset, our *GM-SRM* reaches 58.64%

votes among 2500 rankings, which is much higher than other competing methods as well. And our *GM-SRM* wins on 32 test samples and others altogether 18 samples.

C. Investigation on *GM-SRM* by Ablation Study

We conduct ablation study to investigate the effect of different technical components in *GM-SRM*. To such end, we perform experiments on four variants of our *GM-SRM*.

- **Base model**, which only employs an encoder-decoder model as the base network to restore corrupted images. Thus, no generative memory or stochastic variation strategies are leveraged.
- **GM-BM**, which further leverages decoder to integrate features from both generative memory and encoder, and this is equivalent to plug **Generative Memory** in **Base Model** as semantic query unit for image generation.
- **GM-CSV**, which leverages **Generative Memory** with **Conditional Stochastic Variation** strategy to obtain more realist details in restored area.
- **GM-SRM**, which further leverages progressive reasoning strategy to stabilize generated structure and improve image quality. The resulting model is our intact *GM-SRM*.

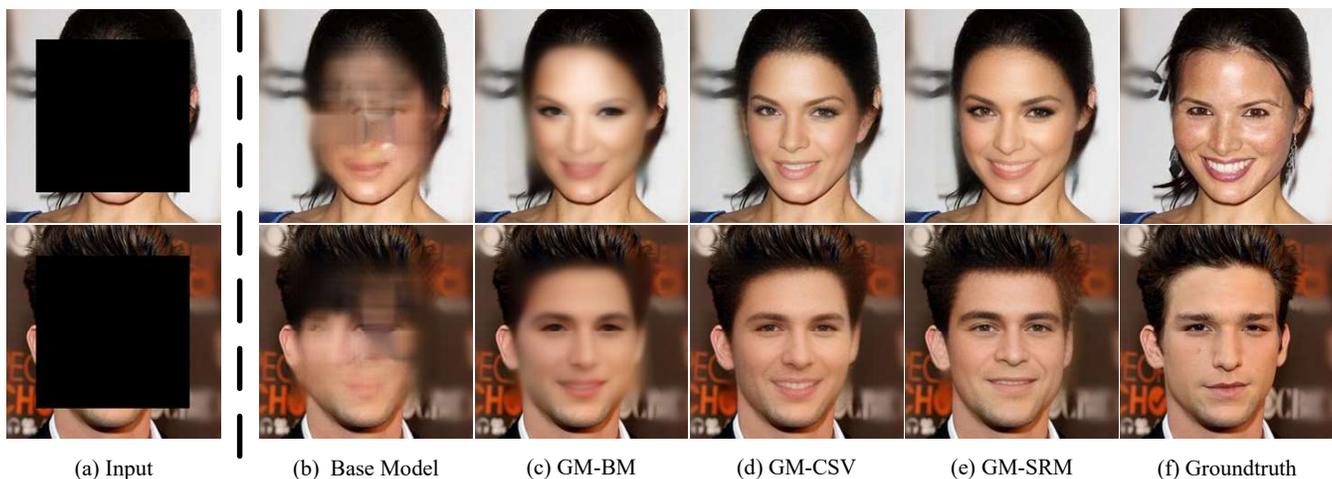


Fig. 10. Visualization of inpainted images from 50% center corruption by four variants of the *GM-SRM* for image inpainting on two samples which are randomly selected from CelebA-HQ test set. Best viewed in zoom-in mode.

TABLE IV
ABLATION EXPERIMENTS OF 50% CORRUPTION RATIO ON
CELEBA-HQ [47] IN TERMS OF PSNR, SSIM, NCC, LMSE, AND LPIPS
TO INVESTIGATE EFFECT OF EACH COMPONENT IN OUR *GM-SRM*.

Method	PSNR	SSIM	NCC	LMSE	LPIPS
Base model	18.39	0.694	0.874	0.032	0.235
GM-BM	20.43	0.729	0.921	0.027	0.143
GM-CSV	21.06	0.734	0.930	0.023	0.139
GM-SRM	21.75	0.757	0.940	0.022	0.132

Table IV illustrates the results of four variants of *GM-SRM* on CelebA-HQ dataset, in terms of PSNR, SSIM, NCC, MSE and LPIPS.

Qualitative ablation evaluation. As illustrated in Figure 10, restoration results of four variants of *GM-SRM* are visualized. The results demonstrate that the visual quality of restored content is increasingly better as the augment of *GM-SRM* with different functional components.

Effect of generative memory. As illustrated in Table IV, the large gap between **Base model** and **GM-BM** demonstrate that our proposed generative memory improves image inpainting performance significantly compared to typical encoder-decoder architecture. Combined pixel-level content reasoning with generative semantic reasoning, generative memory facilitates inpainting reasonable results while coping with large corrupted areas.

Effect of conditional stochastic variation. The conditional stochastic variation improves the texture quality of synthesized images due to learning conditional distribution from known information as a constraint. Intuitive results are shown in Figure 5 and Figure 10.

Effect of progressive reasoning strategy. Comparison between **GM-CSV** and **GM-SRM** demonstrates that progressive reasoning strategy indeed improves the inpainting quality. In Figure 10, the *GM-SRM* makes the inferred structure and texture more realistic. For instance, in Figure 10(e), the *GM-CSV* generates a weird ear, which seems too long, but the *GM-SRM* corrects it. Considering information increment from feature reasoning in different scales, the progressive strategy

makes the inpainting results more reasonable according to the known content.

V. CONCLUSION

In this work, we have presented the Generative Memory-guided Semantic Reasoning Model (*GM-SRM*) for image inpainting, which infers the corrupted content based on not only the known regions of the input image, but also the explicitly learned inter-image reasoning priors characterizing the generalizable semantic distribution patterns between similar images. Our *GM-SRM* employ the encoder-decoder framework to guarantee the pixel-level content consistency, and our pre-trained generative memory is favorable for performing high-level semantic reasoning. Compared to previous implicitly learned semantic priors, our generative priors favor high-level semantic reasoning. As a result, our model is able to synthesize semantically reasonable content for the corrupted regions, especially in the scenarios with large corrupted area.

REFERENCES

- [1] Z. Wan, B. Zhang, D. Chen, P. Zhang, D. Chen, J. Liao, and F. Wen, "Bringing old photos back to life," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 2747–2757.
- [2] D. Bau, H. Strobel, W. Peebles, J. Wulff, B. Zhou, J.-Y. Zhu, and A. Torralba, "Semantic photo manipulation with a generative image prior," *ACM Transactions on Graphics*, vol. 38, no. 4, pp. 1–11, 2019.
- [3] G. Ballester, M. Bertalmio, V. Caselles, G. Sapiro, and J. Verdera, "Filling-in by joint interpolation of vector fields and gray levels," *IEEE Transactions on Image Processing*, vol. 10, no. 8, pp. 1200–1211, 2001.
- [4] C. Barnes, E. Shechtman, A. Finkelstein, and D. B. Goldman, "Patchmatch: A randomized correspondence algorithm for structural image editing," *ACM Trans. Graph.*, vol. 28, no. 3, p. 24, 2009.
- [5] S. J. J. Ren, L. Xu, Q. Yan, and W. Sun, "Shepard convolutional neural networks," *Annual Conference on Neural Information Processing Systems*, 2015.
- [6] D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, and A. A. Efros, "Context encoders: Feature learning by inpainting," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2536–2544.
- [7] G. Liu, F. A. Reda, K. J. Shih, T.-C. Wang, A. Tao, and B. Catanzaro, "Image inpainting for irregular holes using partial convolutions," in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 85–100.

- [8] J. Yu, Z. Lin, J. Yang, X. Shen, X. Lu, and T. S. Huang, "Free-form image inpainting with gated convolution," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 4471–4480.
- [9] J. Li, N. Wang, L. Zhang, B. Du, and D. Tao, "Recurrent feature reasoning for image inpainting," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, June 2020.
- [10] K. Nazeri, E. Ng, T. Joseph, F. Qureshi, and M. Ebrahimi, "Edgeconnect: Structure guided image inpainting using edge prediction," in *The IEEE International Conference on Computer Vision Workshops*, Oct 2019.
- [11] Y. Ren, X. Yu, R. Zhang, T. H. Li, S. Liu, and G. Li, "Structureflow: Image inpainting via structure-aware appearance flow," in *IEEE International Conference on Computer Vision*, 2019.
- [12] Y. S. Jie Yang, Zhiqian Qi, "Learning to incorporate structure knowledge for image inpainting," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 7, 2020, pp. 12 605–12 612.
- [13] J. Li, F. He, L. Zhang, B. Du, and D. Tao, "Progressive reconstruction of visual structure for image inpainting," in *The IEEE International Conference on Computer Vision*, October 2019.
- [14] Y. Song, C. Yang, Z. Lin, X. Liu, Q. Huang, H. Li, and C.-C. J. Kuo, "Contextual-based image inpainting: Infer, match, and translate," in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 3–19.
- [15] Z. Yan, X. Li, M. Li, W. Zuo, and S. Shan, "Shift-net: Image inpainting via deep feature rearrangement," in *The European Conference on Computer Vision*, September 2018.
- [16] J. Yu, Z. Lin, J. Yang, X. Shen, X. Lu, and T. S. Huang, "Generative image inpainting with contextual attention," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5505–5514.
- [17] N. Wang, J. Li, L. Zhang, and B. Du, "Musical: Multi-scale image contextual attention learning for inpainting," in *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence*, 2019, pp. 3748–3754.
- [18] H. Liu, B. Jiang, Y. Xiao, and C. Yang, "Coherent semantic attention for image inpainting," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 4170–4179.
- [19] C. Xie, S. Liu, C. Li, M.-M. Cheng, W. Zuo, X. Liu, S. Wen, and E. Ding, "Image inpainting with learnable bidirectional attention maps," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 8858–8867.
- [20] N. Wang, S. Ma, J. Li, Y. Zhang, and L. Zhang, "Multistage attention network for image inpainting," *Pattern Recognition*, vol. 106, p. 107448, 2020.
- [21] C. Zheng, T.-J. Cham, and J. Cai, "Pluralistic image completion," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1438–1447.
- [22] L. Zhao, Q. Mo, S. Lin, Z. Wang, Z. Zuo, H. Chen, W. Xing, and D. Lu, "Uctgan: Diverse image inpainting based on unsupervised cross-space translation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 5741–5750.
- [23] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.
- [24] X. Guo, H. Yang, and D. Huang, "Image inpainting via conditional texture and structure dual generation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, October 2021, pp. 14 134–14 143.
- [25] A. Lahiri, A. K. Jain, S. Agrawal, P. Mitra, and P. K. Biswas, "Prior guided gan based semantic inpainting," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 13 696–13 705.
- [26] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in Neural Information Processing Systems*, vol. 27. Curran Associates, Inc., 2014.
- [27] M. Bertalmio, G. Sapiro, V. Caselles, and C. Ballester, "Image inpainting," in *Proceedings of the 27th Annual Conference on Computer Graphics and Interactive Techniques*, 2000, pp. 417–424.
- [28] A. A. Efros and W. T. Freeman, "Image quilting for texture synthesis and transfer," in *Proceedings of the 28th Annual Conference on Computer Graphics and Interactive Techniques*, 2001, pp. 341–346.
- [29] V. Kwatra, I. Essa, A. Bobick, and N. Kwatra, "Texture optimization for example-based synthesis," *ACM Transactions on Graphics*, vol. 24, no. 3, pp. 795–802, 2005.
- [30] A. Criminisi, P. Pérez, and K. Toyama, "Region filling and object removal by exemplar-based image inpainting," *IEEE Transactions on image processing*, vol. 13, no. 9, pp. 1200–1212, 2004.
- [31] D. Jin and X. Bai, "Patch-sparsity-based image inpainting through a facet deduced directional derivative," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 29, no. 5, pp. 1310–1324, 2018.
- [32] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems*, vol. 30. Curran Associates, Inc., 2017.
- [33] T. Yu, Z. Guo, X. Jin, S. Wu, Z. Chen, W. Li, Z. Zhang, and S. Liu, "Region normalization for image inpainting," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 07, 2020, pp. 12 733–12 740.
- [34] W. Xiong, J. Yu, Z. Lin, J. Yang, X. Lu, C. Barnes, and J. Luo, "Foreground-aware image inpainting," in *The IEEE Conference on Computer Vision and Pattern Recognition*, June 2019.
- [35] H. Liu, B. Jiang, Y. Song, W. Huang, and C. Yang, "Rethinking image inpainting via a mutual encoder-decoder with feature equalizations," in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*. Springer, 2020, pp. 725–741.
- [36] S. Xu, D. Liu, and Z. Xiong, "E2i: Generative inpainting from edge to image," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 4, pp. 1308–1322, 2020.
- [37] C. Wang, X. Chen, S. Min, J. Wang, and Z.-J. Zha, "Structure-guided deep video inpainting," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 8, pp. 2953–2965, 2020.
- [38] Y. Li, S. Liu, J. Yang, and M.-H. Yang, "Generative face completion," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 3911–3919.
- [39] K. C. Chan, X. Wang, X. Xu, J. Gu, and C. C. Loy, "Glean: Generative latent bank for large-factor image super-resolution," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 14 245–14 254.
- [40] D. Ulyanov, A. Vedaldi, and V. Lempitsky, "Instance normalization: The missing ingredient for fast stylization," *arXiv preprint arXiv:1607.08022*, 2016.
- [41] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7132–7141.
- [42] T. Karras, S. Laine, M. Aittala, J. Hellsten, J. Lehtinen, and T. Aila, "Analyzing and improving the image quality of stylegan," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 8110–8119.
- [43] J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual losses for real-time style transfer and super-resolution," in *European Conference on Computer Vision*. Springer, 2016, pp. 694–711.
- [44] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [45] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida, "Spectral normalization for generative adversarial networks," *arXiv preprint arXiv:1802.05957*, 2018.
- [46] C. Doersch, S. Singh, A. Gupta, J. Sivic, and A. Efros, "What makes paris look like paris?" *ACM Transactions on Graphics*, vol. 31, no. 4, 2012.
- [47] T. Karras, T. Aila, S. Laine, and J. Lehtinen, "Progressive growing of gans for improved quality, stability, and variation," in *International Conference on Learning Representations*, 2018.
- [48] B. Zhou, A. Lapedriza, A. Khosla, A. Oliva, and A. Torralba, "Places: A 10 million image database for scene recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 6, pp. 1452–1464, 2017.
- [49] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE Transactions on Image Processing*, vol. 13, no. 4, pp. 600–612, 2004.
- [50] K. Wei, J. Yang, Y. Fu, D. Wipf, and H. Huang, "Single image reflection removal exploiting misaligned training data and network enhancements," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 8178–8187.
- [51] R. Grosse, M. K. Johnson, E. H. Adelson, and W. T. Freeman, "Ground truth dataset and baseline evaluations for intrinsic image algorithms," in *2009 IEEE 12th International Conference on Computer Vision*. IEEE, 2009, pp. 2335–2342.
- [52] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in *CVPR*, 2018.
- [53] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.