

Prompt to Prompt

Image Editing with Cross-Attention Control

Yuanzhi Zhu

Content

- Motivation & Recap
- Prompt-to-Prompt
- Extensions

Prompt-to-Prompt Image Editing with Cross Attention Control

Motivation

Text-2-image editing is sensitive to text prompts

- text influence the high level semantic only



photo of a cat riding a bike



photo of a cat on a bike

A spatial mask to localize the edit

- hard to draw and
- ignoring the original structure & content

Toward Mask-free Text-2-Image Editing 😊



Origin



Mask



Inpainting Example

Prompt: girl with red hair

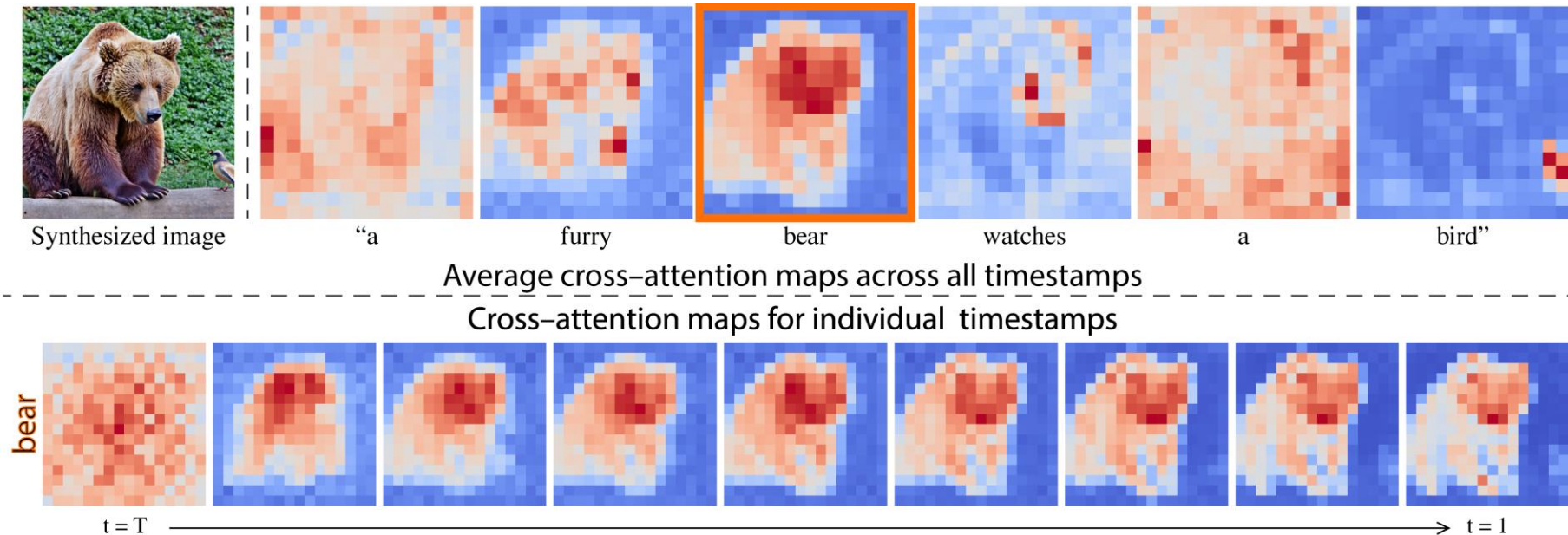
[Prompt-to-Prompt](#)

[Prompt-to-Prompt Image Editing with Cross-Attention Control | OpenReview](#)

[Attend-and-Excite/explain.ipynb at main · AttendAndExcite/Attend-and-Excite \(github.com\)](#)

Prompt-to-Prompt Image Editing with Cross Attention Control

Key Observation: spatial information in the cross-attention maps



Prompt-to-Prompt Image Editing with Cross Attention Control

Model Architecture & Location of Cross Attention

```
print(is_cross, place_in_unet, attn.shape)
```

```
False down torch.Size([16, 4096, 4096])
```

```
True down torch.Size([16, 4096, 77])
```

```
False down torch.Size([16, 4096, 4096])
```

```
True down torch.Size([16, 4096, 77])
```

```
False down torch.Size([16, 1024, 1024])
```

```
True down torch.Size([16, 1024, 77])
```

```
False down torch.Size([16, 1024, 1024])
```

```
True down torch.Size([16, 1024, 77])
```

```
False down torch.Size([16, 256, 256])
```

```
True down torch.Size([16, 256, 77])
```

```
False down torch.Size([16, 256, 256])
```

```
True down torch.Size([16, 256, 77])
```

```
False mid torch.Size([16, 64, 64])
```

```
True mid torch.Size([16, 64, 77])
```

```
False up torch.Size([16, 256, 256])
```

```
True up torch.Size([16, 256, 77])
```

```
False up torch.Size([16, 256, 256])
```

```
True up torch.Size([16, 256, 77])
```

```
False up torch.Size([16, 256, 256])
```

```
True up torch.Size([16, 256, 77])
```

```
False up torch.Size([16, 1024, 1024])
```

```
True up torch.Size([16, 1024, 77])
```

```
False up torch.Size([16, 1024, 1024])
```

```
True up torch.Size([16, 1024, 77])
```

```
False up torch.Size([16, 1024, 1024])
```

```
True up torch.Size([16, 1024, 77])
```

```
False up torch.Size([16, 4096, 4096])
```

```
True up torch.Size([16, 4096, 77])
```

```
False up torch.Size([16, 4096, 4096])
```

```
True up torch.Size([16, 4096, 77])
```

```
False up torch.Size([16, 4096, 4096])
```

```
True up torch.Size([16, 4096, 77])
```

16=2*8

2: cfg

8: heads

Prompt-to-Prompt Image Editing with Cross Attention Control

Attentions in stable diffusion

```
from ldm.modules.attention import SpatialTransformer
```

```
SpatialTransformer → [BasicTransformerBlock(...) for d in range(depth)]
```

```
def _forward(self, x, context=None):  
    x = self.attn1(self.norm1(x), context=context if self.disable_self_attn else None) + x  
    x = self.attn2(self.norm2(x), context=context) + x  
    x = self.ff(self.norm3(x)) + x  
    return x
```

→self
→cross

```
# x.shape = b (h w) c  
# context.shape = b n_token dim
```

Prompt-to-Prompt Image Editing with Cross Attention Control

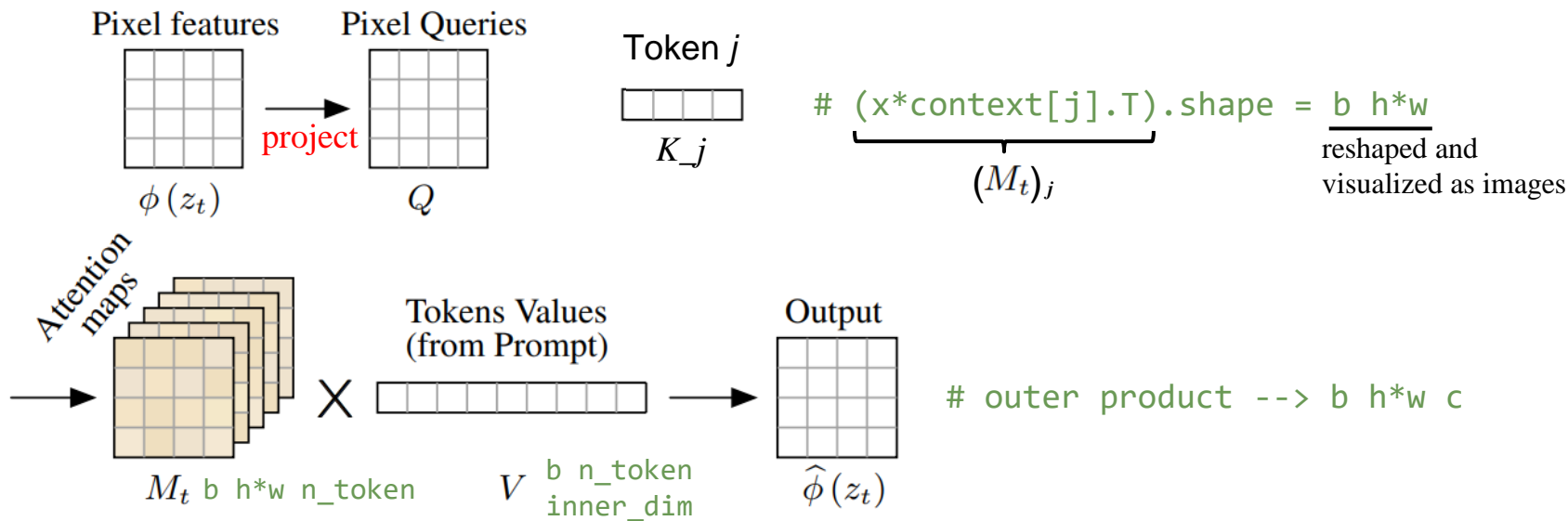
$$M = \text{Softmax} \left(\frac{QK^T}{\sqrt{d}} \right)$$

project

x.shape = b (h w) c --> b h*w inner_dim
 # context.shape = b n_token dim --> b n_token inner_dim

cross-attention

```
# attention, what we cannot get enough of
sim = einsum('b i d, b j d -> b i j', q, k) * self.scale
attn = sim.softmax(dim=-1) # normalizes values along axis -1 (j)
out = einsum('b i j, b j d -> b i d', attn, v)
```

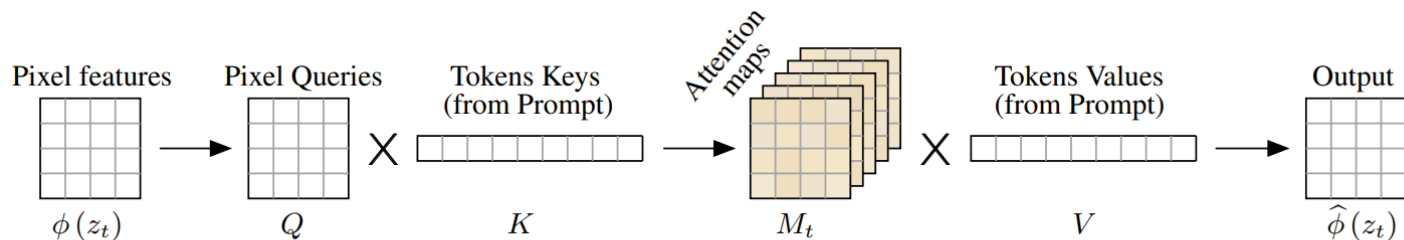


Content

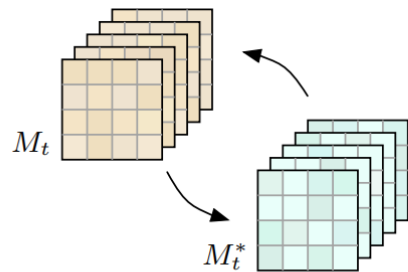
- Motivation & Recap
- Prompt-to-Prompt
- Extensions

Prompt-to-Prompt Image Editing with Cross Attention Control

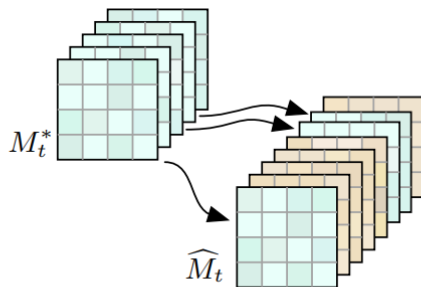
$M = \text{Softmax} \left(\frac{QK^T}{\sqrt{d}} \right)$ we can inject the attention maps M that were obtained from the generation with the original prompt P , into a second generation with the modified prompt P^*



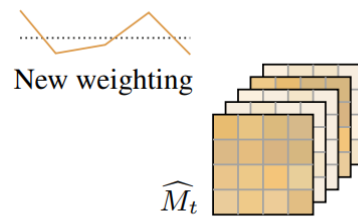
Text to Image Cross Attention
Cross Attention Control



Word Swap



Adding a New Phrase



Attention Re-weighting

Prompt-to-Prompt Image Editing with Cross Attention Control

Algorithm 1: Prompt-to-Prompt image editing

1 **Input:** A source prompt \mathcal{P} , a target prompt \mathcal{P}^* , and a random seed s .
2 **Optional for local editing:** w and w^* , words in \mathcal{P} and \mathcal{P}^* , specifying the editing region.
3 **Output:** A source image x_{src} and an edited image x_{dst} .
4 $z_T \sim N(0, I)$ a unit Gaussian random variable with random seed s ;
5 $z_T^* \leftarrow z_T$;
6 **for** $t = T, T - 1, \dots, 1$ **do**
7 $z_{t-1}, M_t \leftarrow DM(z_t, \mathcal{P}, t, s)$;
8 $M_t^* \leftarrow DM(z_t^*, \mathcal{P}^*, t, s)$;
9 $\widehat{M}_t \leftarrow Edit(M_t, M_t^*, t)$;
10 $z_{t-1}^* \leftarrow DM(z_t^*, \mathcal{P}^*, t, s)\{M \leftarrow \widehat{M}_t\}$;
11 **if local then**
12 $\alpha \leftarrow B(\overline{M}_{t,w}) \cup B(\overline{M}_{t,w}^*)$;
13 $z_{t-1}^* \leftarrow (1 - \alpha) \odot z_{t-1} + \alpha \odot z_{t-1}^*$;
14 **end**
15 **end**
16 **Return** (z_0, z_0^*)

Word Swap

$$Edit(M_t, M_t^*, t) := \begin{cases} M_t^* & \text{if } t < \tau \\ M_t & \text{otherwise} \end{cases}$$

Adding a New Phrase

$$(Edit(M_t, M_t^*, t))_{i,j} := \begin{cases} (M_t^*)_{i,j} & \text{if } A(j) = None \\ (M_t)_{i,A(j)} & \text{otherwise.} \end{cases}$$

i : pixel value; j : text token

Attention Re-weighting

$$(Edit(M_t, M_t^*, t))_{i,j} := \begin{cases} c \cdot (M_t)_{i,j} & \text{if } j = j^* \\ (M_t)_{i,j} & \text{otherwise} \end{cases}$$

Prompt-to-Prompt Image Editing with Cross Attention Control

How to implement?

https://github.com/google/prompt-to-prompt/blob/main/ptp_utils.py#L173

```
### prompt to prompt setup
controller = _setup_attention_controller(attention_control_type, prompts,
                                       cross_replace_steps=cross_replace_steps, self_replace_steps=self_replace_steps,
                                       LocalBlend_pair=LocalBlend_pair,
                                       Reweightwords=Reweightwords, Reweightcales=Reweightcales,
                                       sampling_steps=steps, tokenizer=model.cond_stage_model.tokenizer)
ptp_utils.register_attention_control(model, controller)

def register_recr(net_, count, place_in_unet):
    ## modify the forward function of the cross attention module
    if net_.__class__.__name__ == 'CrossAttention':
        net_.forward = ca_forward(net_, place_in_unet)
        return count + 1
```

Prompt-to-Prompt Image Editing with Cross Attention Control

ca_forward:

```
sim = einsum('2 i d, 2 j d -> 2 i j', q, k) * self.scale
attn = sim.softmax(dim=-1) # normalizes values along axis -1 (j)
attn = controller(attn, is_cross, place_in_unet)
out = einsum('2 i j, 2 j d -> 2 i d', attn, v)
```

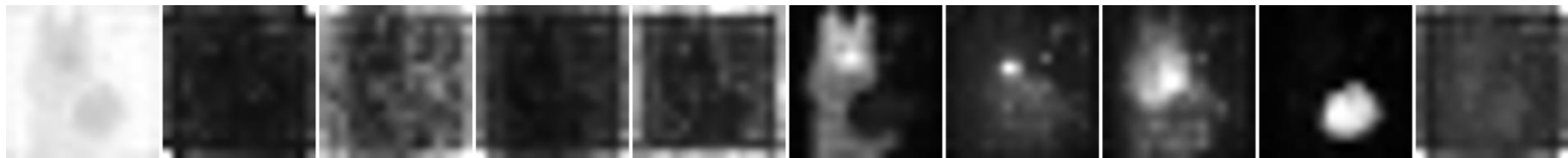
https://github.com/google/prompt-to-prompt/blob/main/ptp_utils.py#L203

```
1 def __call__(self, attn, is_cross: bool, place_in_unet: str):
2     if self.cur_att_layer >= self.num_uncond_att_layers:
3         if LOW_RESOURCE:
4             attn = self.forward(attn, is_cross, place_in_unet)
5         else:
6             h = attn.shape[0]
7             attn[h // 2:] = self.forward(attn[h // 2:], is_cross, place_in_unet)
8     self.cur_att_layer += 1
9     if self.cur_att_layer == self.num_att_layers + self.num_uncond_att_layers:
10        self.cur_att_layer = 0
11        self.cur_step += 1
12        self.between_steps()
13    return attn
```

Visualization of Attention Map

1. Store all the attention maps `self.step_store[key].append(attn)`
2. Aggregate attention map in each sampling step `self.attention_store[key][i] += self.step_store[key][i]`
3. Aggregate attention map at **given resolution** for **all attention heads** and **at all places**
`attention_maps.keys(): ['down_cross', 'mid_cross', 'up_cross', 'down_self', 'mid_self', 'up_self']`

```
1 attention_maps = attention_store.get_average_attention()
2 num_pixels = res ** 2
3 for location in from_where:
4     for item in attention_maps[f"{location}_{'cross' if is_cross else 'self'}"]:
5         if item.shape[1] == num_pixels:
6             cross_maps = item.reshape(len(prompts), -1, res, res, item.shape[-1])[select]
7             out.append(cross_maps)
8 out = torch.cat(out, dim=0)
9 out = out.sum(0) / out.shape[0]
```



<startoftext>

a

painting

of

a

squirrel

eating

a

burger

<endoftext>

Prompt-to-Prompt Image Editing with Cross Attention Control

ca_forward:

```
sim = einsum('2 i d, 2 j d -> 2 i j', q, k) * self.scale
attn = sim.softmax(dim=-1) # normalizes values along axis -1 (j)
attn = controller(attn, is_cross, place_in_unet)
out = einsum('2 i j, 2 j d -> 2 i d', attn, v)
```

https://github.com/google/prompt-to-prompt/blob/main/ptp_utils.py#L203

```
### forward method of controller (AttentionControlEdit)
def forward(self, attn, is_cross: bool, place_in_unet: str):
    super(AttentionControlEdit, self).forward(attn, is_cross, place_in_unet)
    if is_cross or (self.num_self_replace[0] <= self.cur_step < self.num_self_replace[1]):
        h = attn.shape[0] // (self.batch_size)
        attn = attn.reshape(self.batch_size, h, *attn.shape[1:])
        attn_base, attn_repalce = attn[0], attn[1:]
        if is_cross:
            alpha_words = self.cross_replace_alpha[self.cur_step]
            attn_repalce_new = self.replace_cross_attention(attn_base, attn_repalce) * alpha_words \
                + (1 - alpha_words) * attn_repalce
            attn[1:] = attn_repalce_new
        else:
            attn[1:] = self.replace_self_attention(attn_base, attn_repalce)
        attn = attn.reshape(self.batch_size * h, *attn.shape[2:])
    return attn
```

Prompt-to-Prompt Image Editing with Cross Attention Control

How to implement?

https://github.com/google/prompt-to-prompt/blob/main/ptp_utils.py#L74

```
img, pred_x0 = self.p_sample_ddim(...)
### apply local_blend to img
if controller:
    img = controller.step_callback(img)
```

```
1 class LocalBlend:
2     def __call__(self, x_t, attention_store):
3         k = 1
4         maps = attention_store["down_cross"][2:4] + attention_store["up_cross"][:3]
5         maps = [item.reshape(self.alpha_layers.shape[0], -1, 1, 16, 16, MAX_NUM_WORDS) for item in maps]
6         maps = torch.cat(maps, dim=1)
7         maps = (maps * self.alpha_layers).sum(-1).mean(1)
8         mask = nnf.max_pool2d(maps, (k * 2 + 1, k * 2 + 1), (1, 1), padding=(k, k))
9         mask = nnf.interpolate(mask, size=(x_t.shape[2:]))
10        mask = mask / mask.max(2, keepdims=True)[0].max(3, keepdims=True)[0]
11        mask = mask.gt(self.threshold)
12        mask = (mask[:1] + mask[1:]).float()
13        x_t = x_t[:1] + mask * (x_t - x_t[:1]) # x_t[1:] = mask * x_t[1:] + (1-mask) * x_t[:1]
14        return x_t
```


Prompt-to-Prompt Image Editing with Cross Attention Control

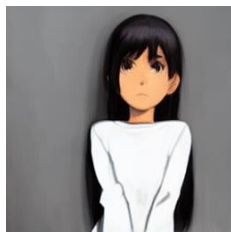
a small girl with blue shirt sitting in front of a mirror, red hair

girl, red hair

a cat(lion) with a hat is lying on a beach chair.

photo of a cat riding a bike(car)

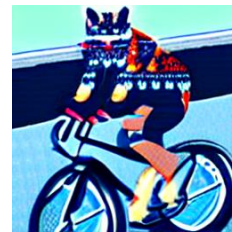
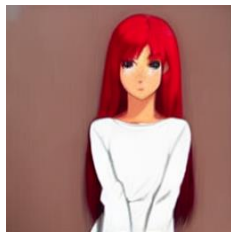
original



w/o ptp



w/ ptp



Prompt-to-Prompt Image Editing with Cross Attention Control

photo of a cat on a bike(car)

cross_replace_steps
self_replace_steps

0.5/0.5



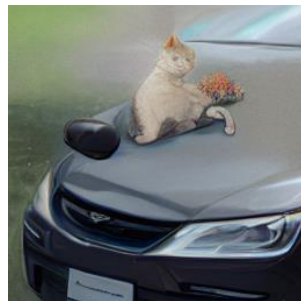
0.8/0.2



0.8/0.3



0.8/0.4



0.8/0.5



0.0/0.0



0.8/0.6



0.5/0.4



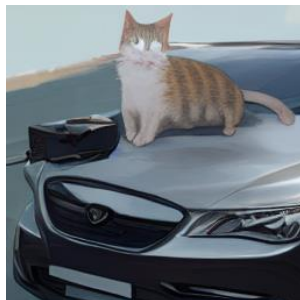
0.6/0.4



0.7/0.4



0.8/0.4



w/o LocalBlend

Prompt-to-Prompt Image Editing with Cross Attention Control

Conclusions

Pros:

- ✓ Training free
- ✓ Mask free
- ✓ Flexible

Cons:

- ✗ Precise control
- ✗ Stable hyper-parameter
- ✗ Natural language instruction!

Content

- Motivation & Recap
- Prompt-to-Prompt
- Extensions

Directed Diffusion: Direct Control of Object Placement through Attention Guidance

Prompt: A painting of a tiger, on the wall in the living room [, in the upper left of the image]



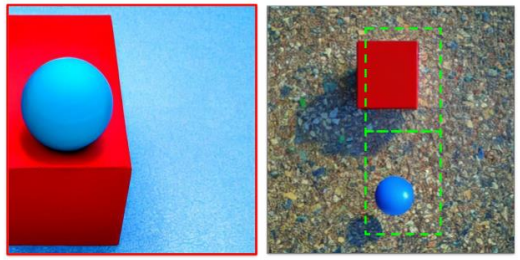
Prompt: A dog sitting next to a mirror



Prompt: A car on a bridge [, in the upper left of the image]



Prompt: A red cube above a blue sphere



SD

DD(1st Q)

DD(2nd Q)

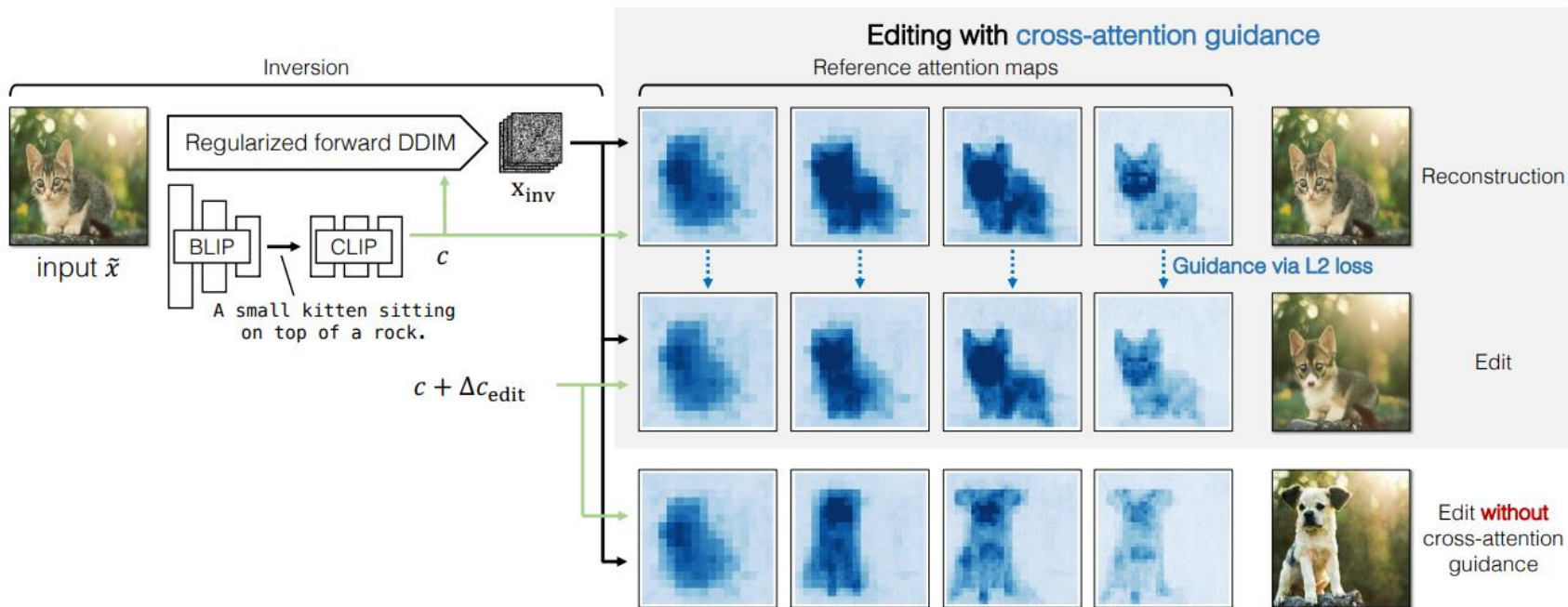
DD(3rd Q)

DD(4th Q)

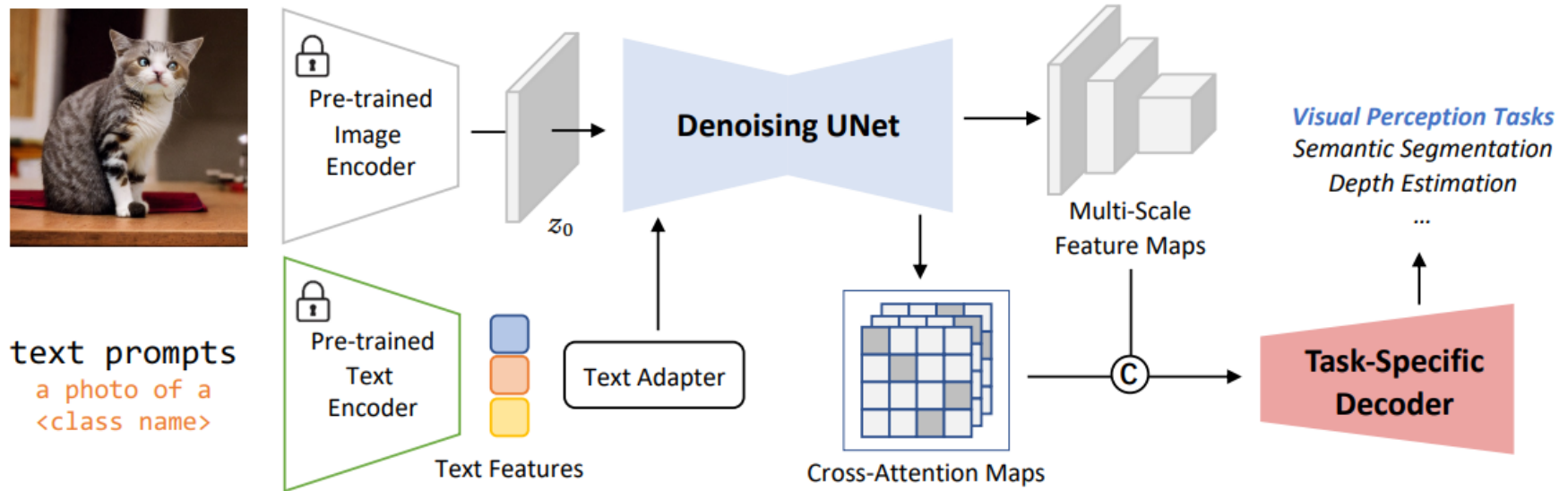
SD

DD

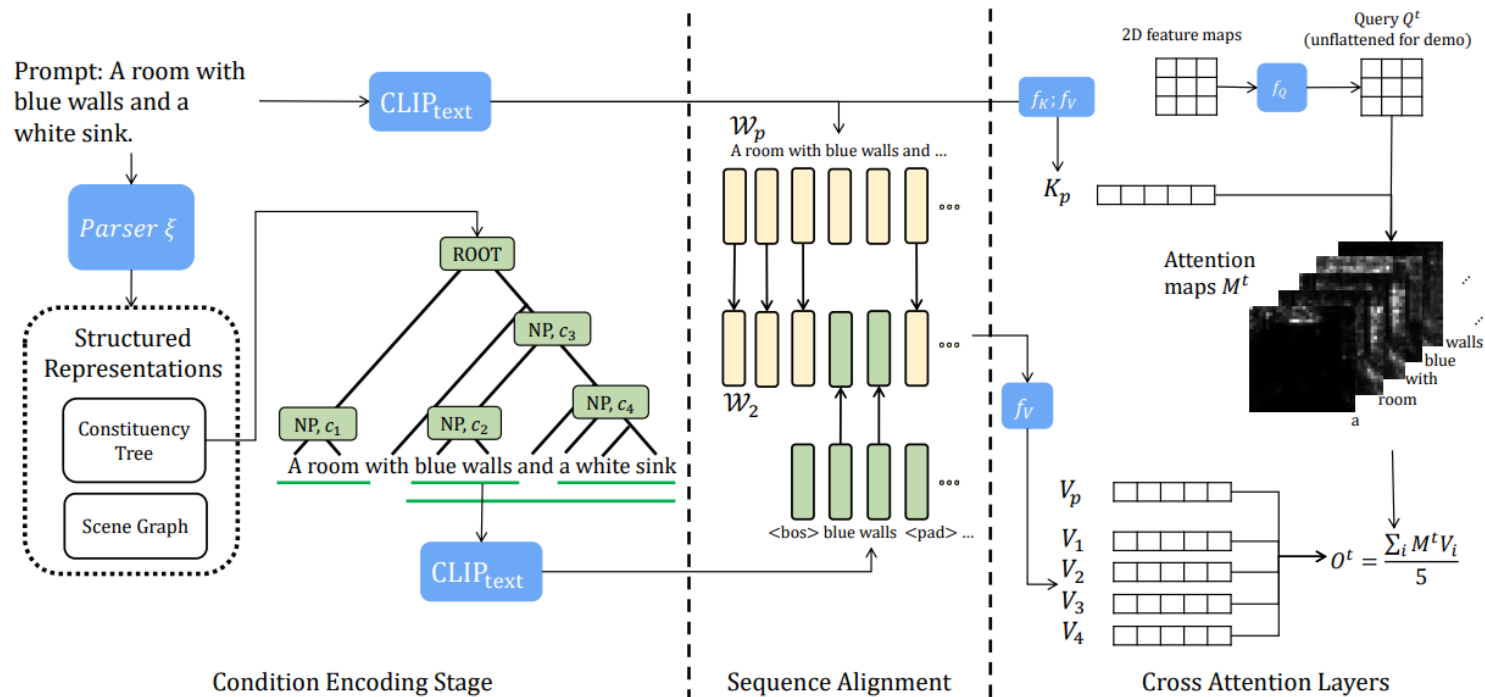
Zero-shot Image-to-Image Translation



Unleashing Text-to-Image Diffusion Models for Visual Perception

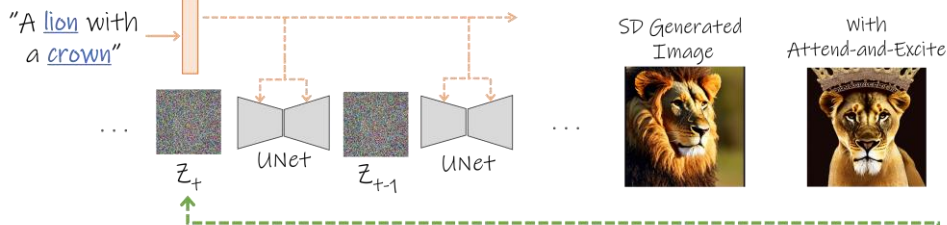


Training-Free Structured Diffusion Guidance for Compositional Text-to-Image Synthesis

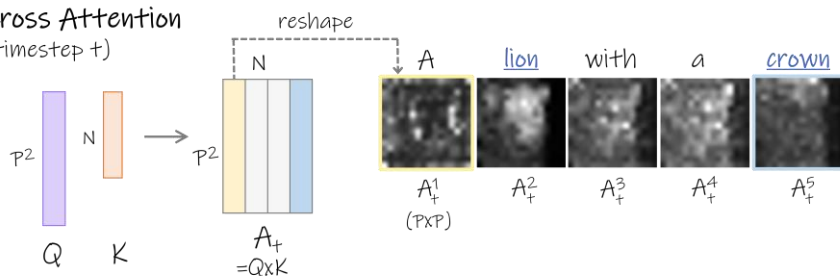


Attend-and-Excite: Attention-Based Semantic Guidance for Text-to-Image Diffusion Models

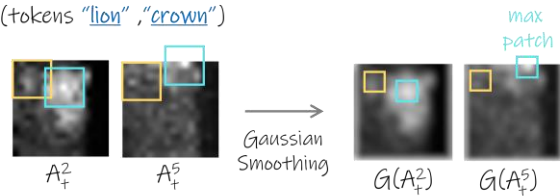
DDPM Process



Cross Attention
(timestep t)



Loss Computation
(tokens "lion", "crown")



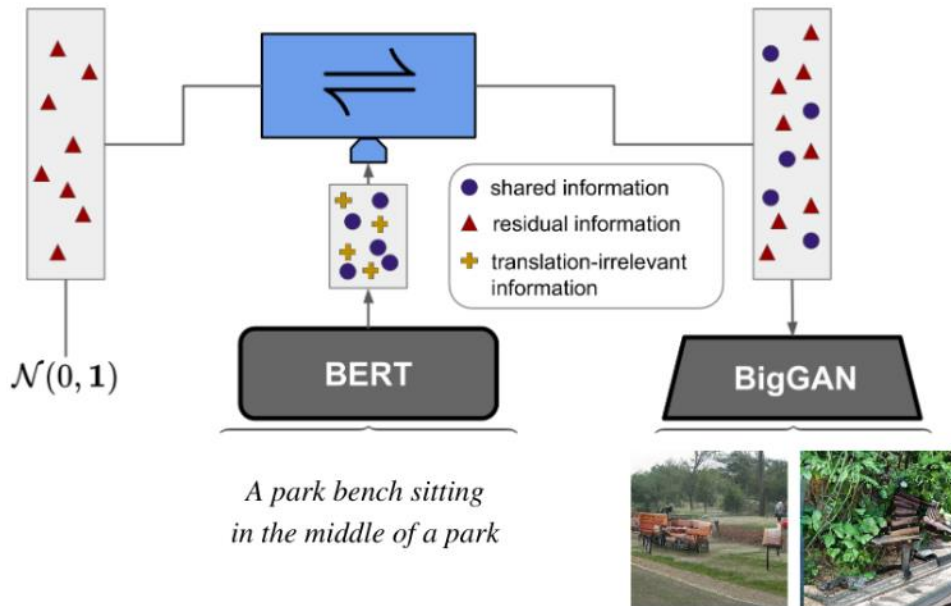
$$L_2 = 1 - \max G(A_t^2)$$

$$L_5 = 1 - \max G(A_t^5)$$

$$\text{Loss: } L = \max(L_2, L_5)$$

$$\text{Update: } z'_t = z_t - \alpha \nabla_{z_t} L$$

Network-to-Network Translation with Conditional Invertible Neural Networks



source domain

target domain

A blue bird sitting on top of a field



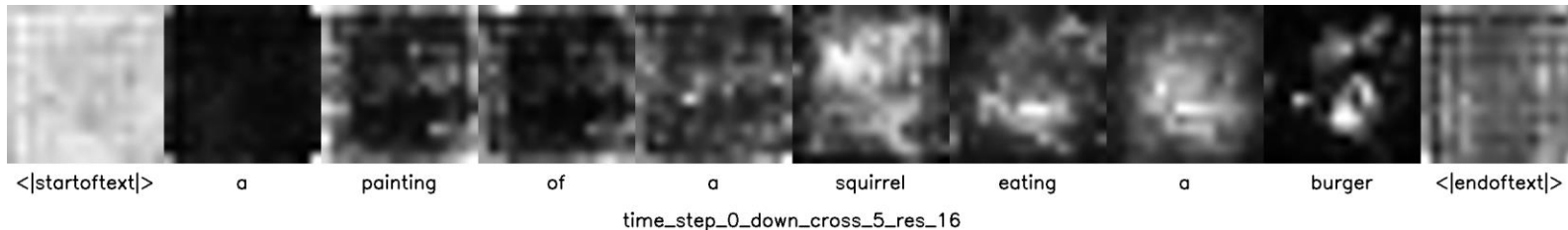
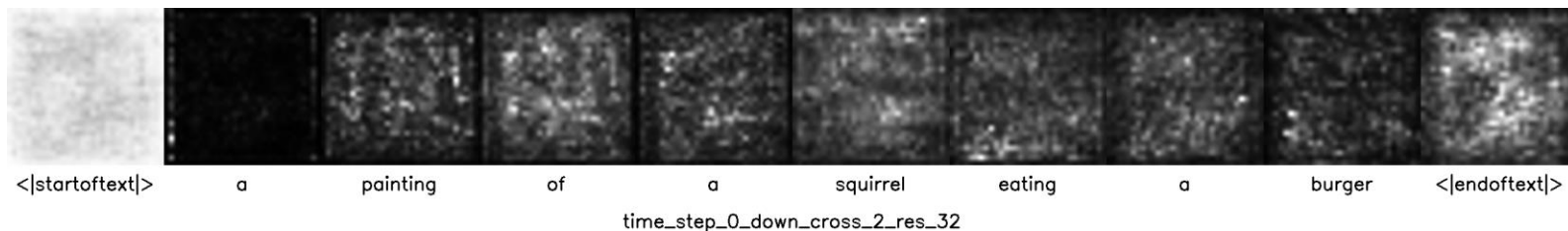
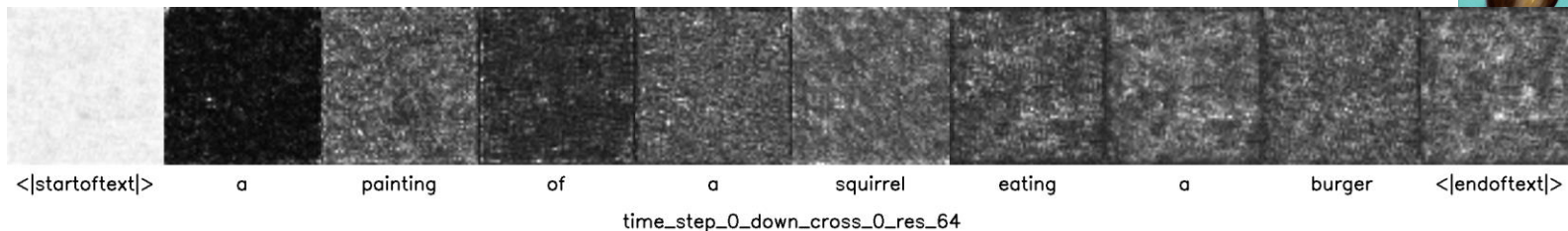
A yellow bird is perched on a branch



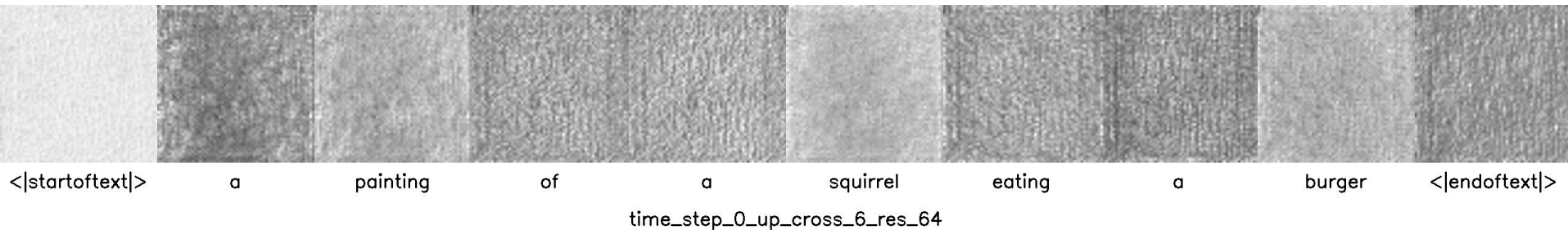
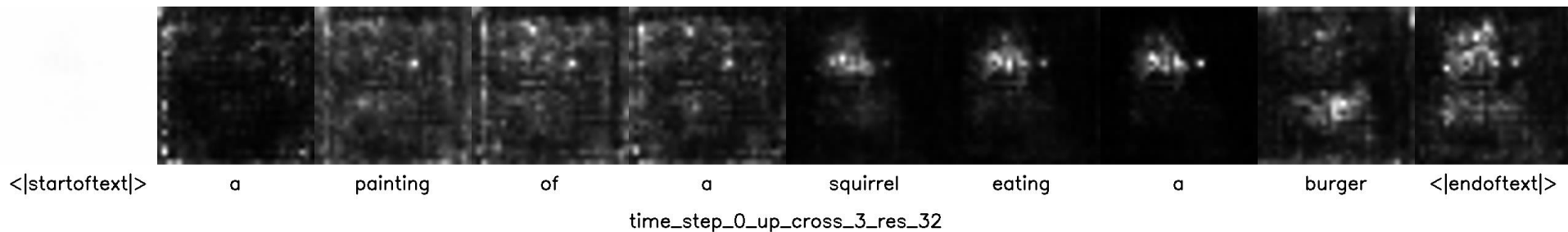
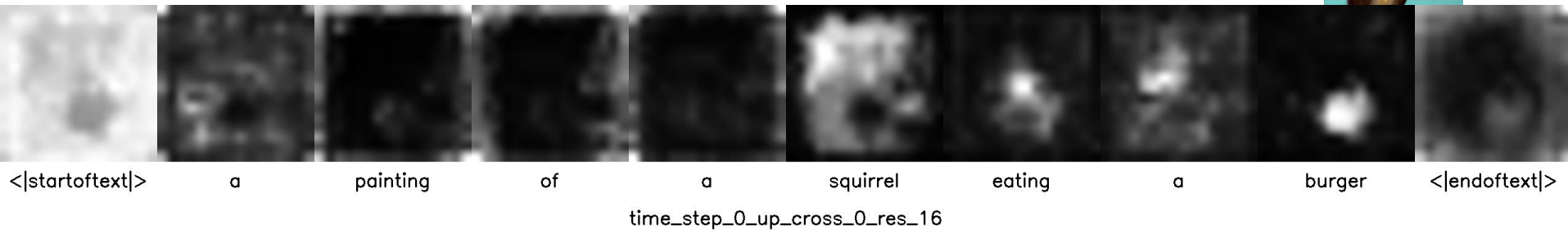
A couple of zebras are standing in a field



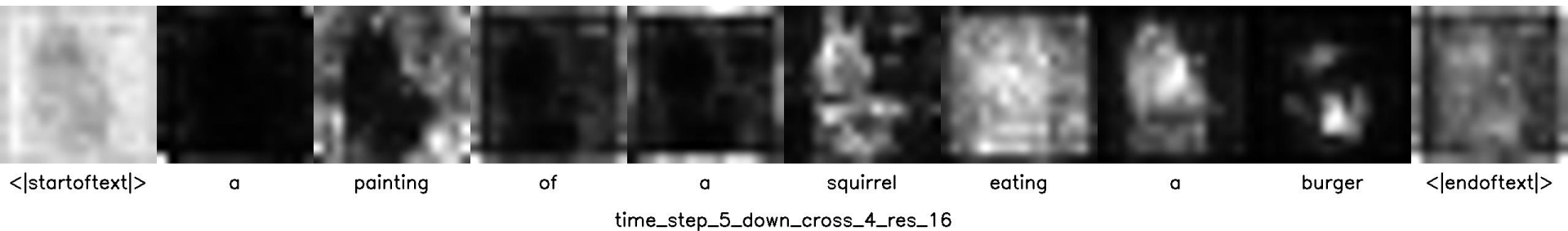
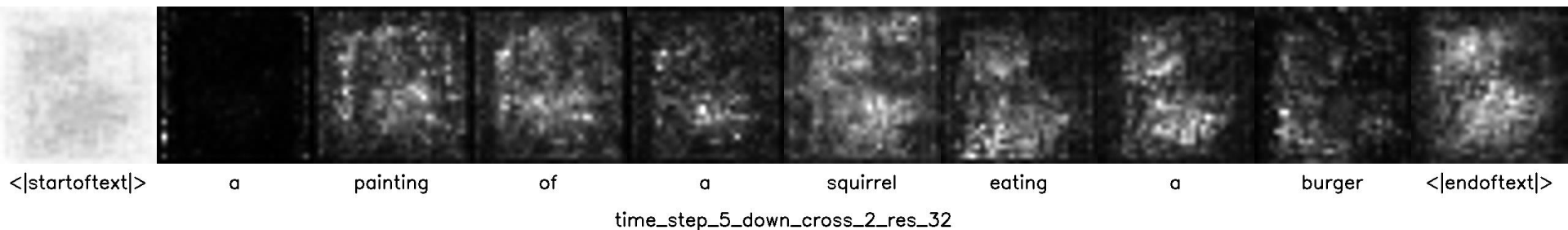
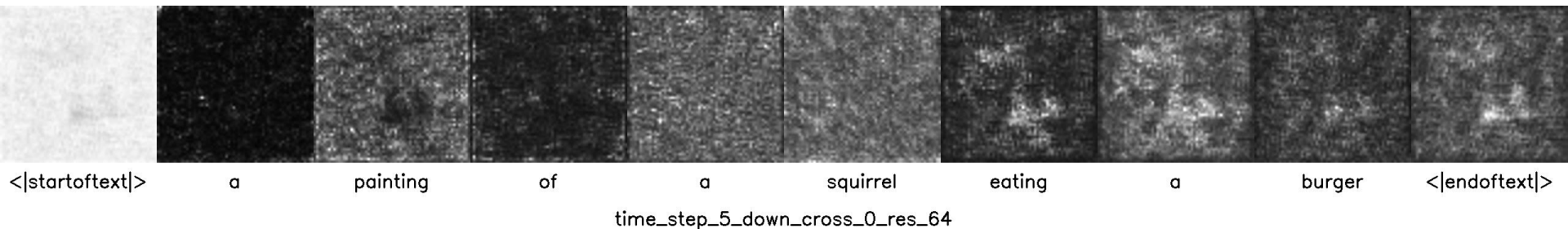
Visualization of Attention Map: timestep 0



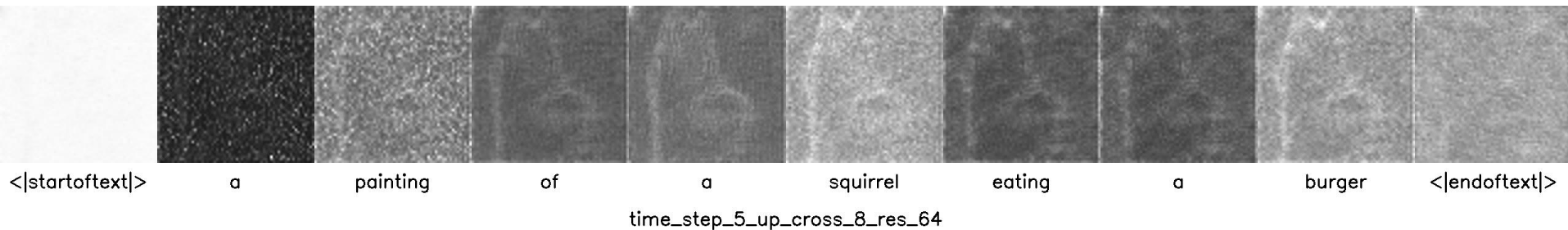
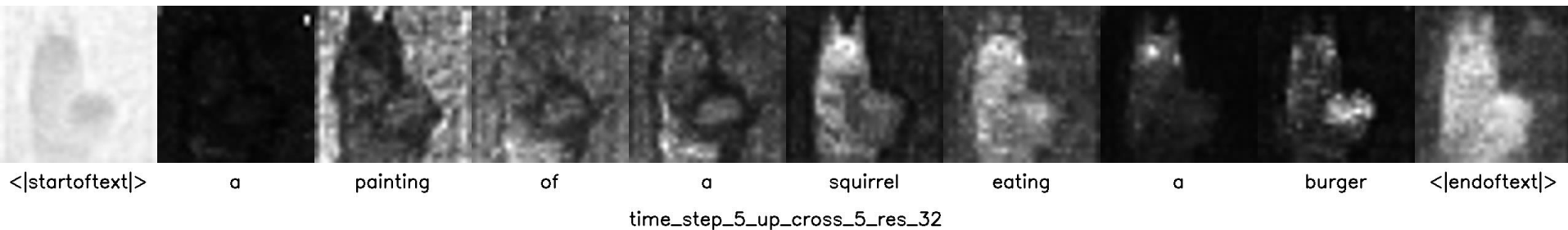
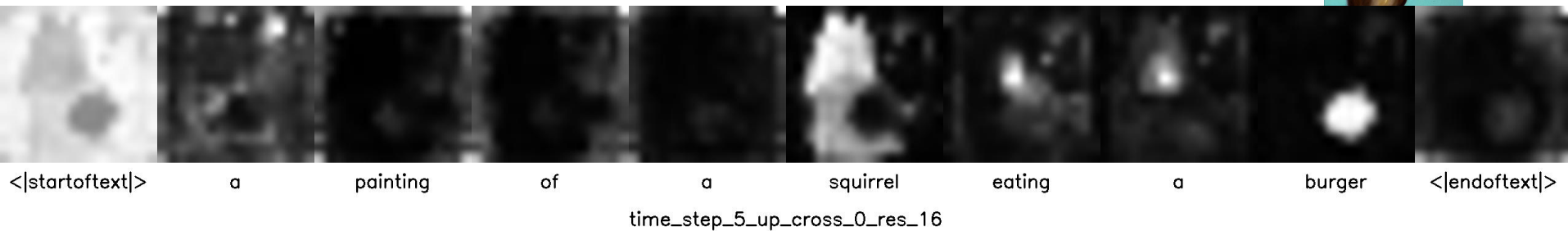
Visualization of Attention Map: timestep 0



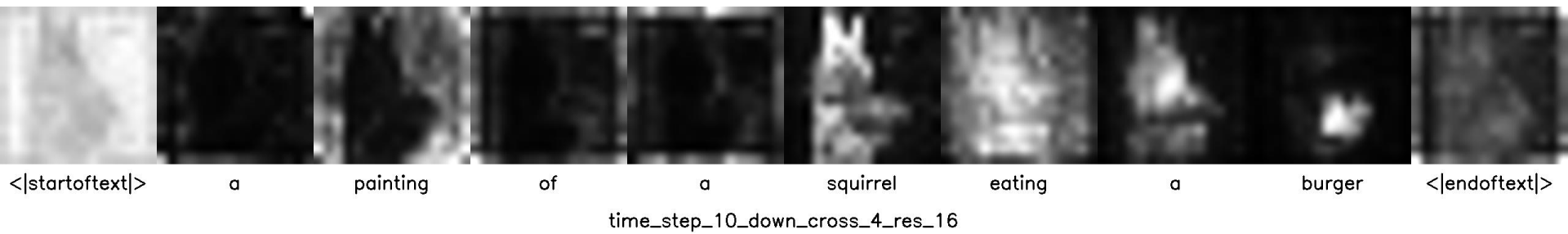
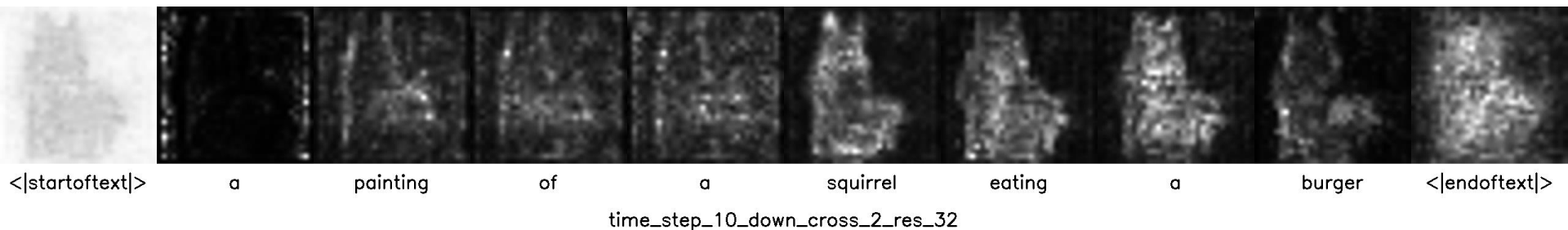
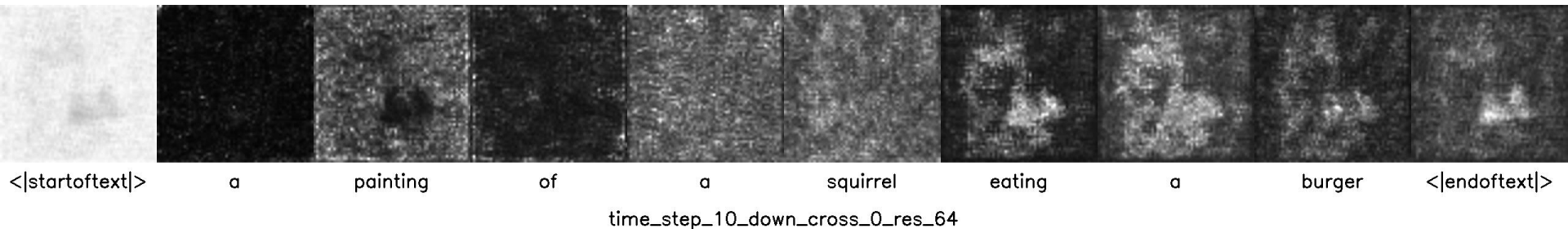
Visualization of Attention Map: timestep 5



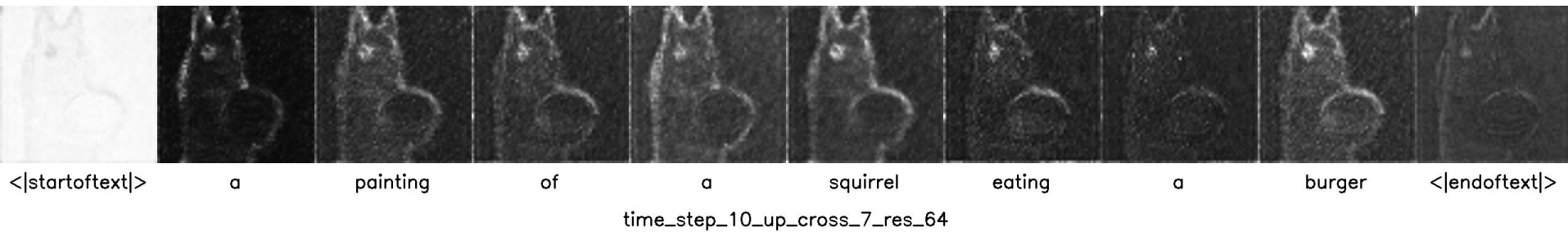
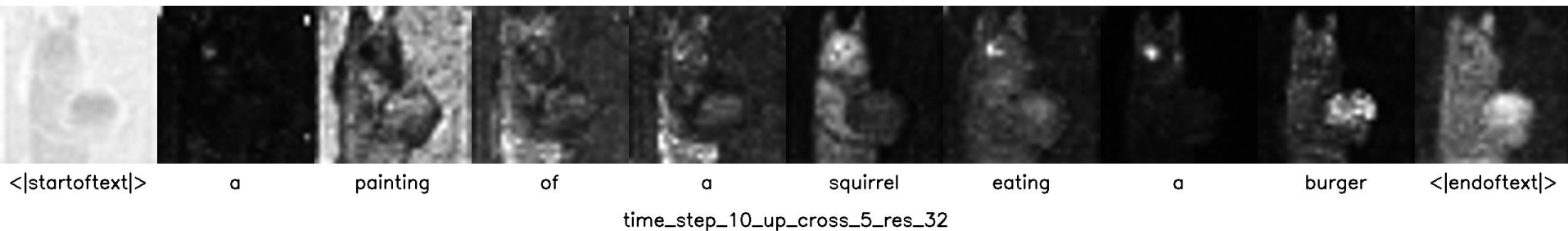
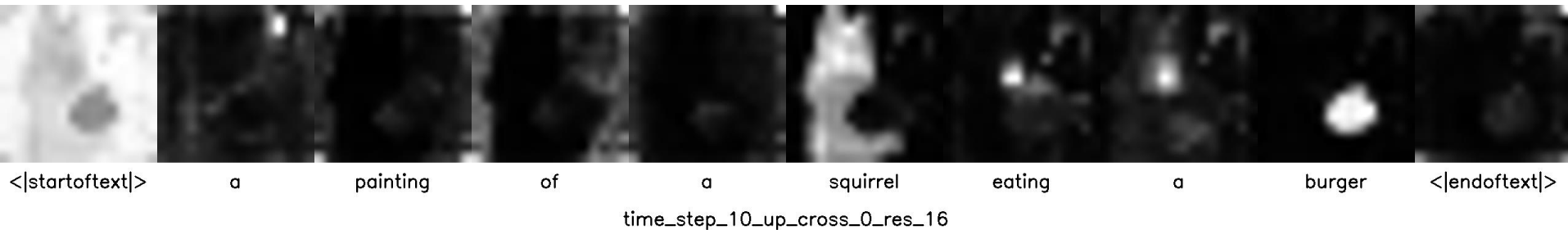
Visualization of Attention Map: timestep 5



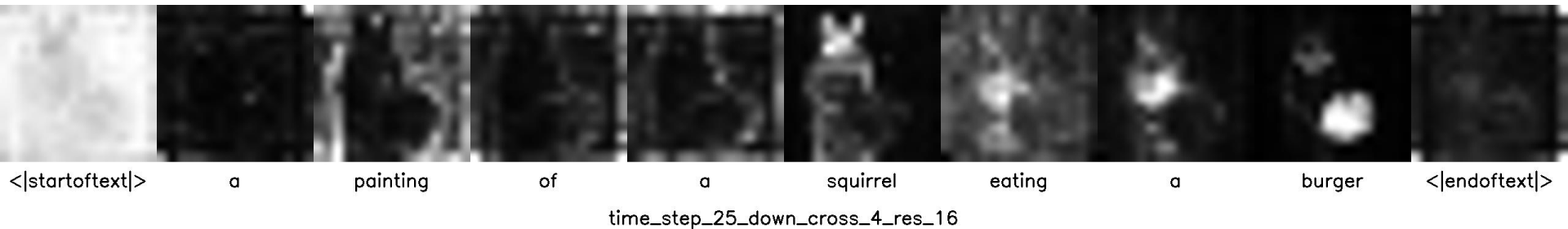
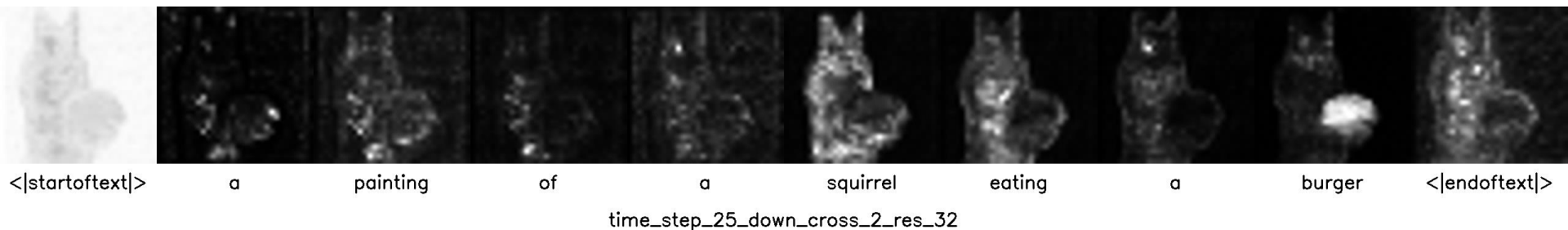
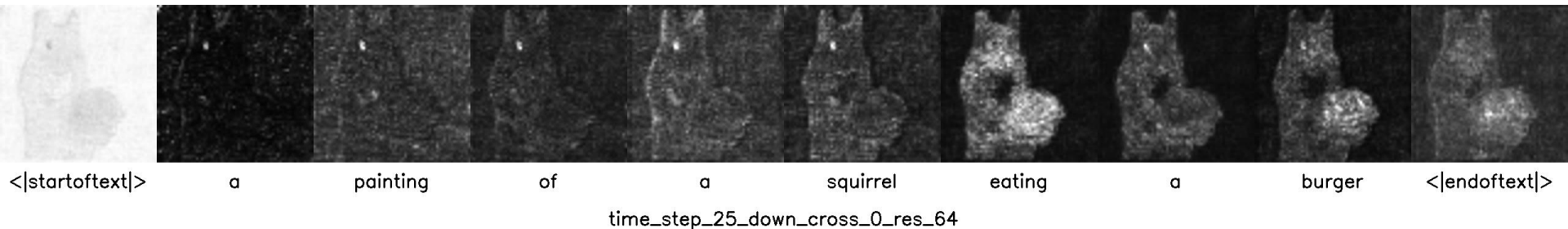
Visualization of Attention Map: timestep 10



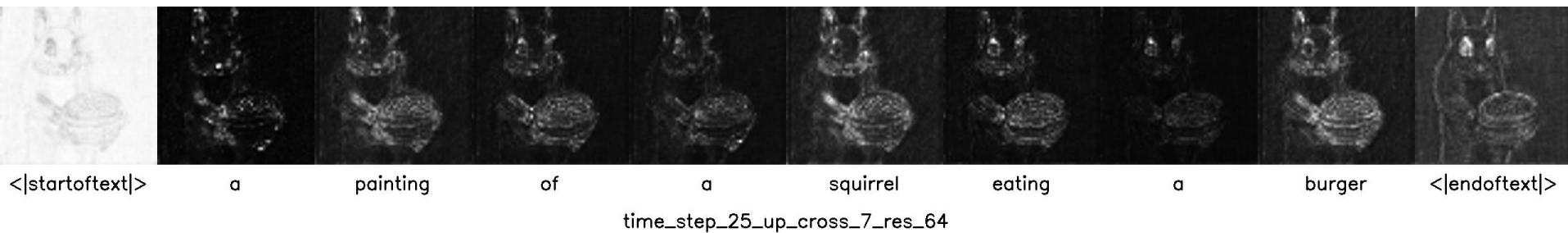
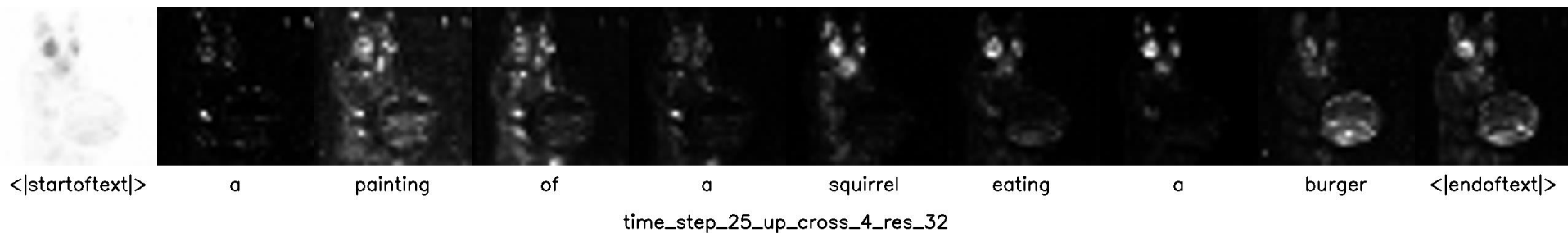
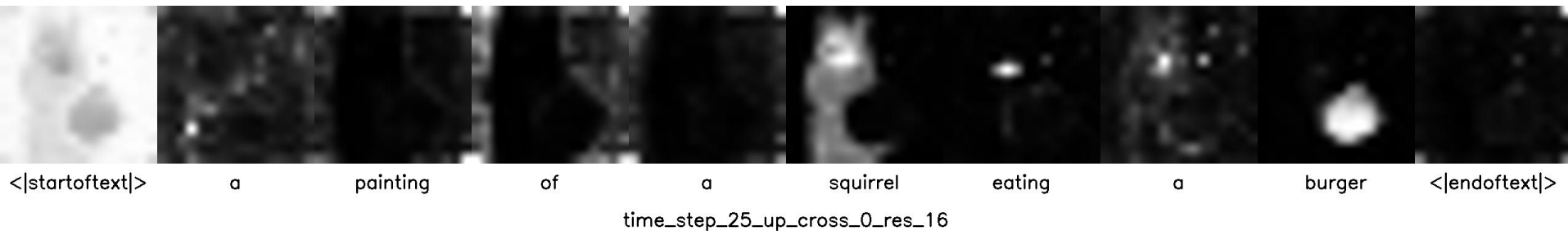
Visualization of Attention Map: timestep 10



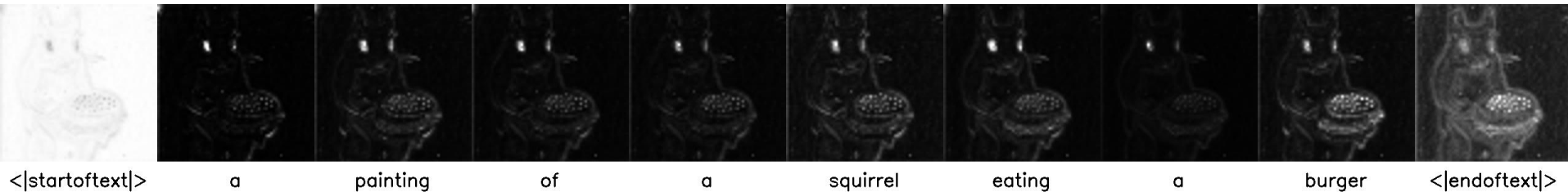
Visualization of Attention Map: timestep 25



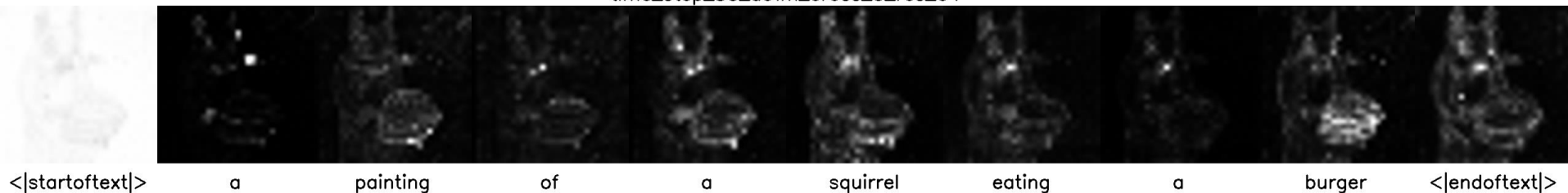
Visualization of Attention Map: timestep 25



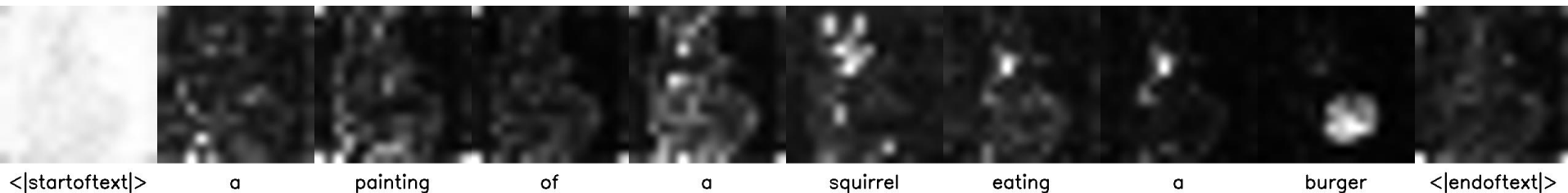
Visualization of Attention Map: timestep 50



time_step_50_down_cross_0_res_64



time_step_50_down_cross_2_res_32



time_step_50_down_cross_4_res_16

Visualization of Attention Map: timestep 50

