

UNIVERSITÄT  
HEIDELBERG  
ZUKUNFT  
SEIT 1386

## ImageBART: Bidirectional Context with Multinomial Diffusion for Autoregressive Image Synthesis

Patrick Esser\*,



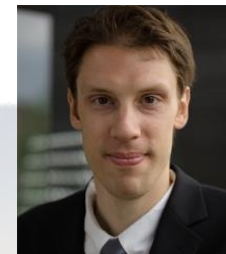
Robin Rombach\*,



Andreas Blattmann\*,

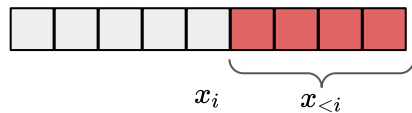


Björn Ommer



\*equal contribution

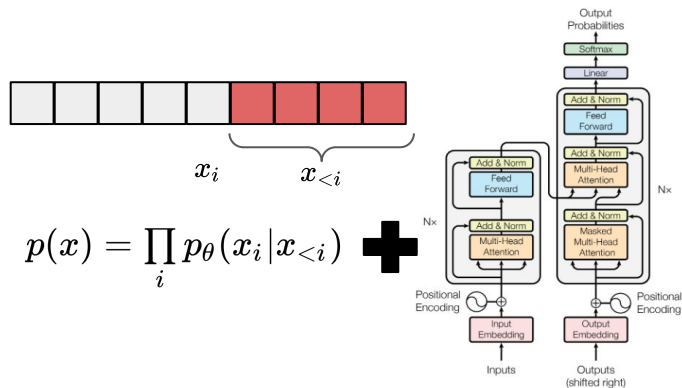
# Autoregressive Generative Modeling



$$p(x) = \prod_i p_\theta(x_i | x_{<i})$$

Sequential Likelihood  
Factorization

# Autoregressive Generative Modeling

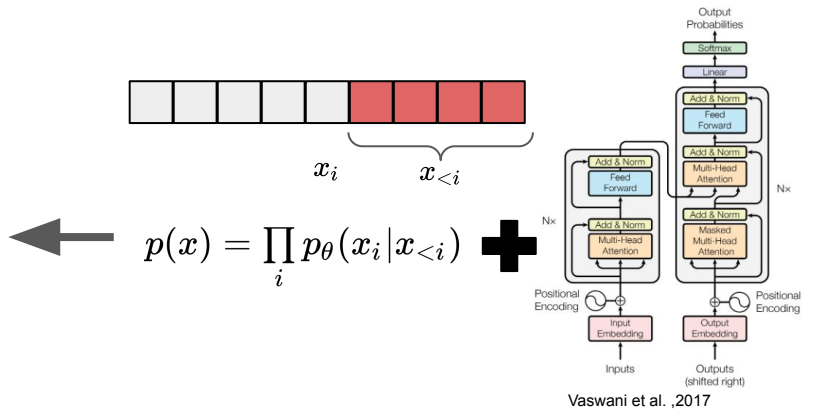
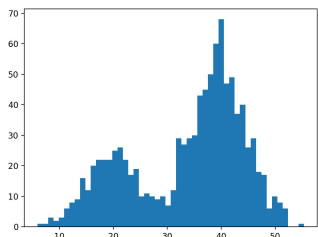


Vaswani et al. ,2017

Sequential Likelihood  
Factorization

Attention Mechanism

# Autoregressive Generative Modeling



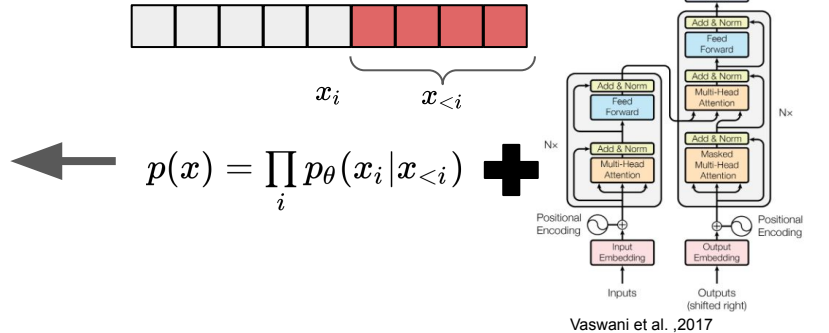
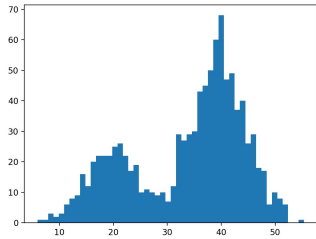
Sequential Likelihood  
Factorization

Attention Mechanism

# Autoregressive Generative Modeling



van den Oord et al. , 2016



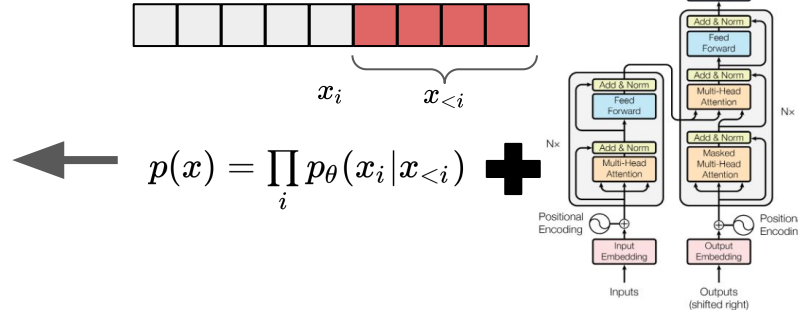
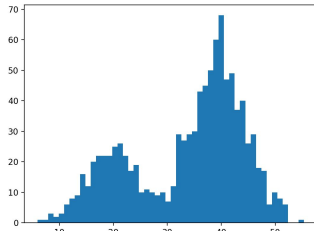
Sequential Likelihood  
Factorization

Attention Mechanism

# Autoregressive Generative Modeling



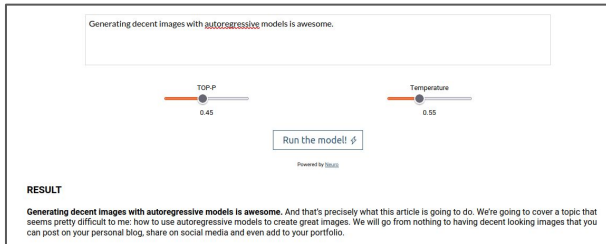
van den Oord et al. , 2016



Vaswani et al. ,2017

Sequential Likelihood  
Factorization

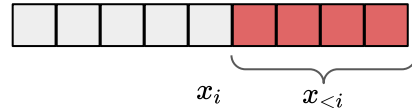
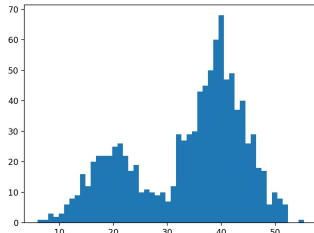
Attention Mechanism



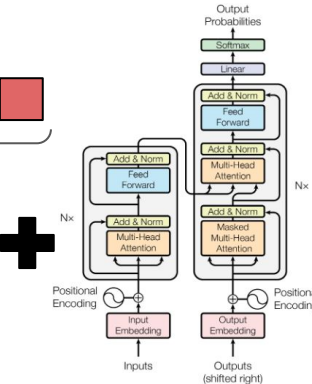
# Autoregressive Generative Modeling



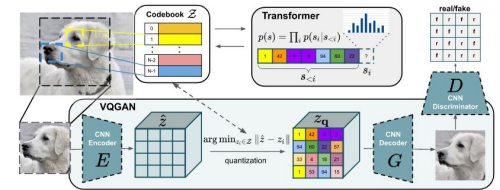
van den Oord et al. , 2016



$$p(x) = \prod_i p_{\theta}(x_i | x_{<i})$$



Vaswani et al. ,2017



Esser et al. , 2020

Sequential Likelihood  
Factorization

Attention Mechanism

Generating decent images with [autoregressive](#) models is awesome.

TOP-P: 0.45      Temperature: 0.95

Run the model!  $\phi$

Powered by

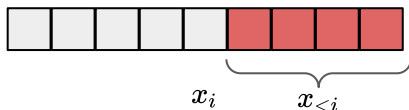
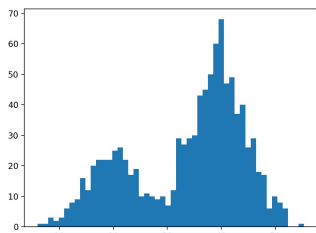
**RESULT**

Generating decent images with autoregressive models is awesome. And that's precisely what this article is going to do. We're going to cover a topic that seems pretty difficult to me: how to use autoregressive models to create great images. We will go from nothing to having decent looking images that you can post on your personal blog, share on social media and even add to your portfolio.

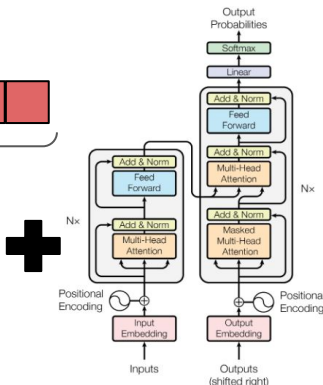
# Autoregressive Generative Modeling



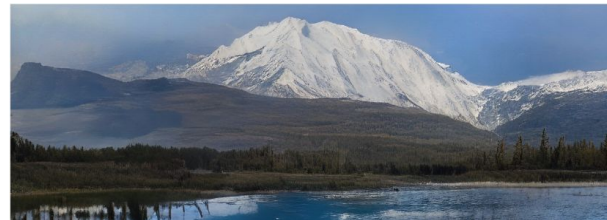
van den Oord et al. , 2016



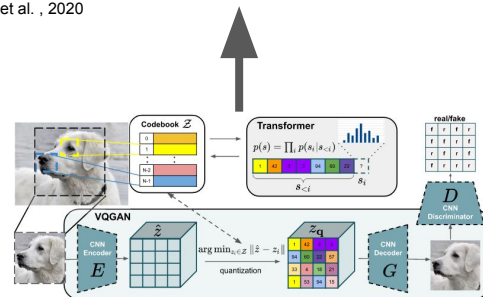
$$p(x) = \prod_i p_\theta(x_i | x_{<i})$$



Vaswani et al. ,2017



Esser et al. , 2020



Esser et al. , 2020

Sequential Likelihood  
Factorization

Attention Mechanism

Generating decent images with autoregressive models is awesome.

TOP-P: 0.45      Temperature: 0.95

Run the model!  $\phi$

Powered by Hugging

**RESULT**

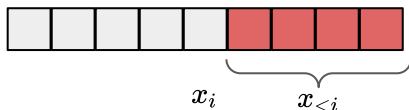
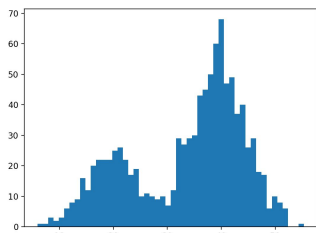
Generating decent images with autoregressive models is awesome. And that's precisely what this article is going to do. We're going to cover a topic that seems pretty difficult to me: how to use autoregressive models to create great images. We will go from nothing to having decent looking images that you can post on your personal blog, share on social media and even add to your portfolio.



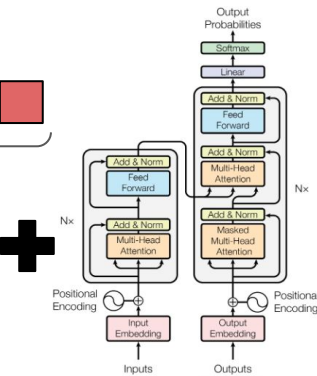
# Autoregressive Generative Modeling



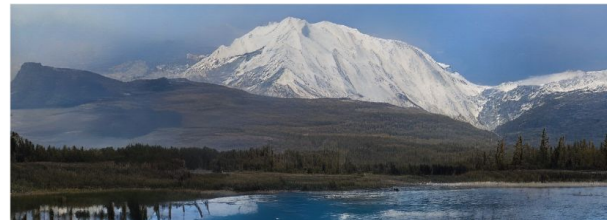
van den Oord et al. , 2016



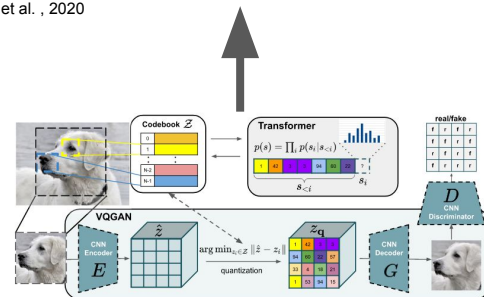
$$p(x) = \prod_i p_\theta(x_i | x_{<i})$$



Vaswani et al. ,2017



Esser et al. , 2020



Esser et al. , 2020

Sequential Likelihood  
Factorization

Attention Mechanism

Generating decent images with autoregressive models is awesome.

TOP-P

0.45

Temperature

0.95

Run the model! ↗

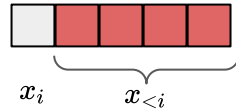
Powered by Hugging

**RESULT**

Generating decent images with autoregressive models is awesome. And that's precisely what this article is going to do. We're going to cover a topic that seems pretty difficult to me: how to use autoregressive models to create great images. We will go from nothing to having decent looking images that you can post on your personal blog, share on social media and even add to your portfolio.

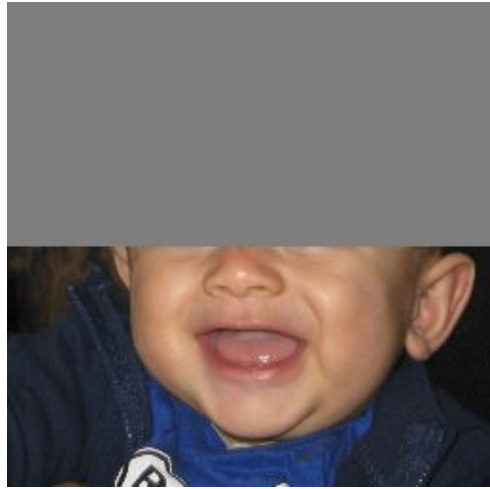


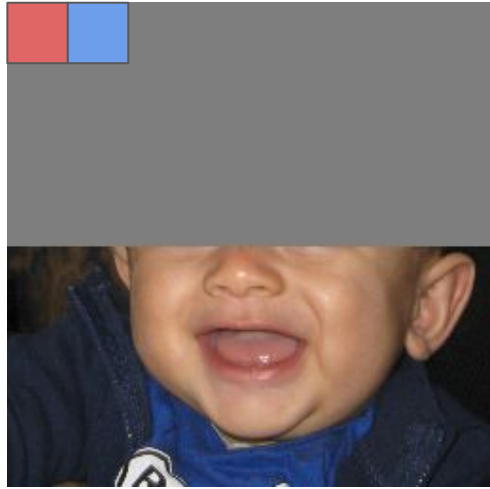
Yu et al. ,2021

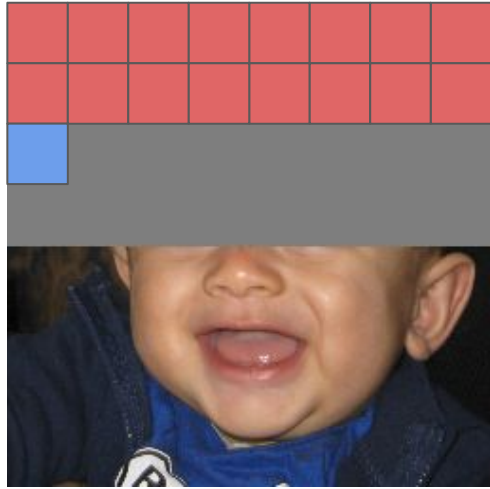


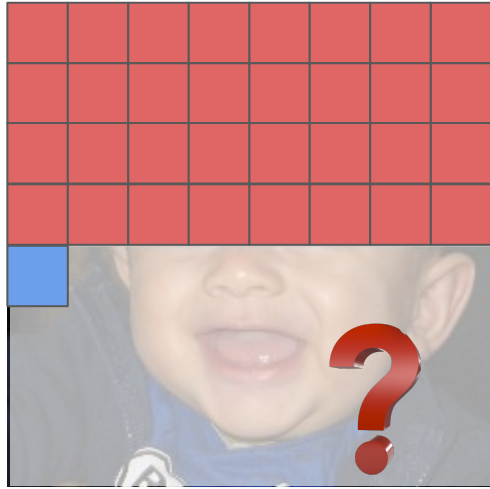
$$p(x) = \prod_i p_\theta(x_i | x_{<i})$$

**Sequential** Likelihood  
Factorization









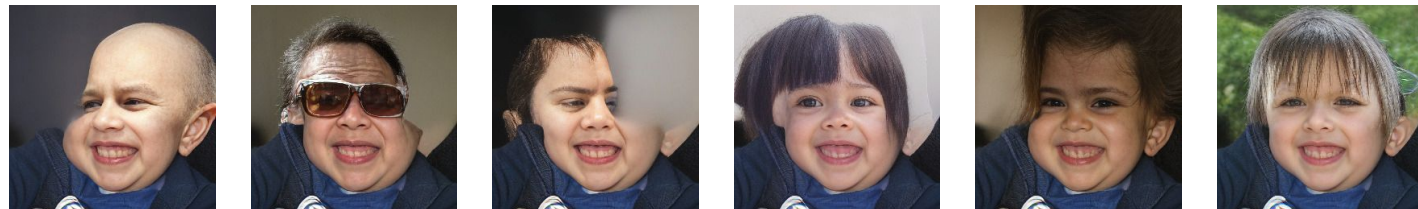


No global image representation





No global image representation



Unrealistic samples

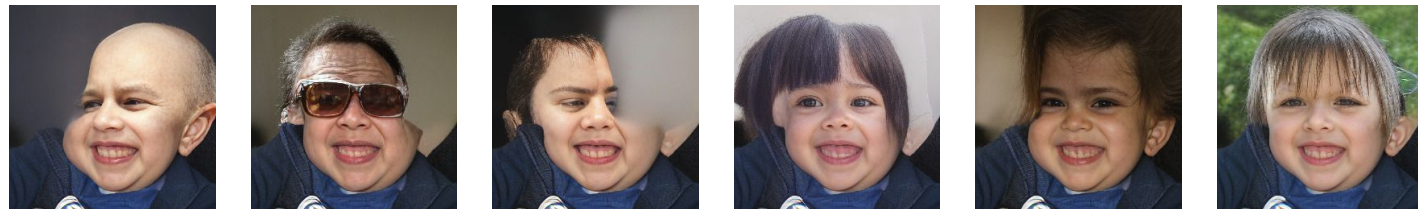




No global image representation



Missing global context



Unrealistic samples

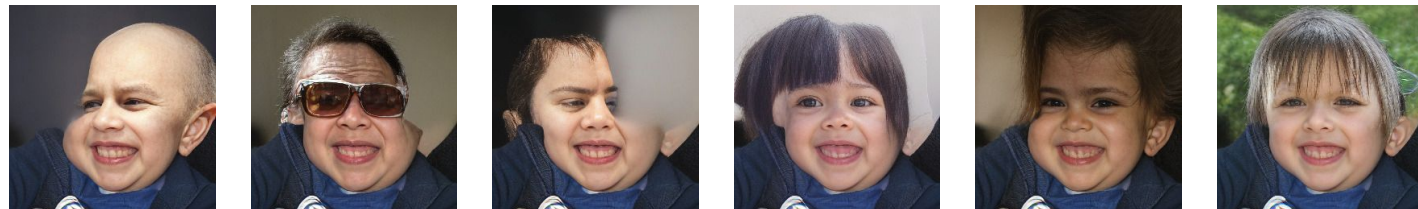


No global image representation



Missing global context

Aggravated conditional image generation



Unrealistic samples



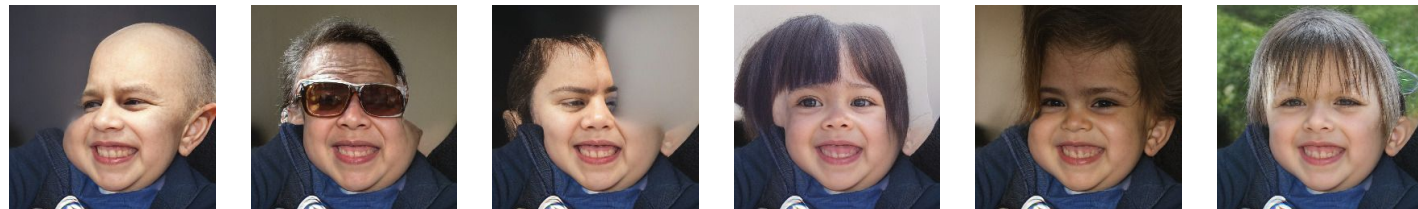
No global image representation



Missing global context

Aggravated conditional image generation

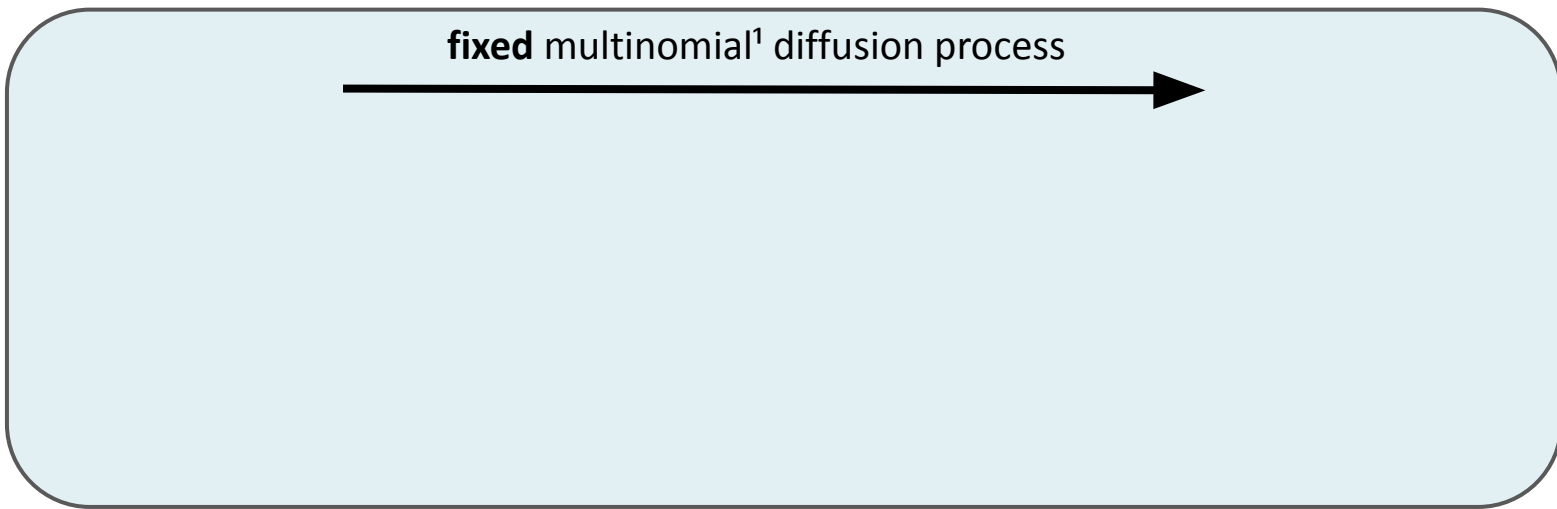
Exposure Bias



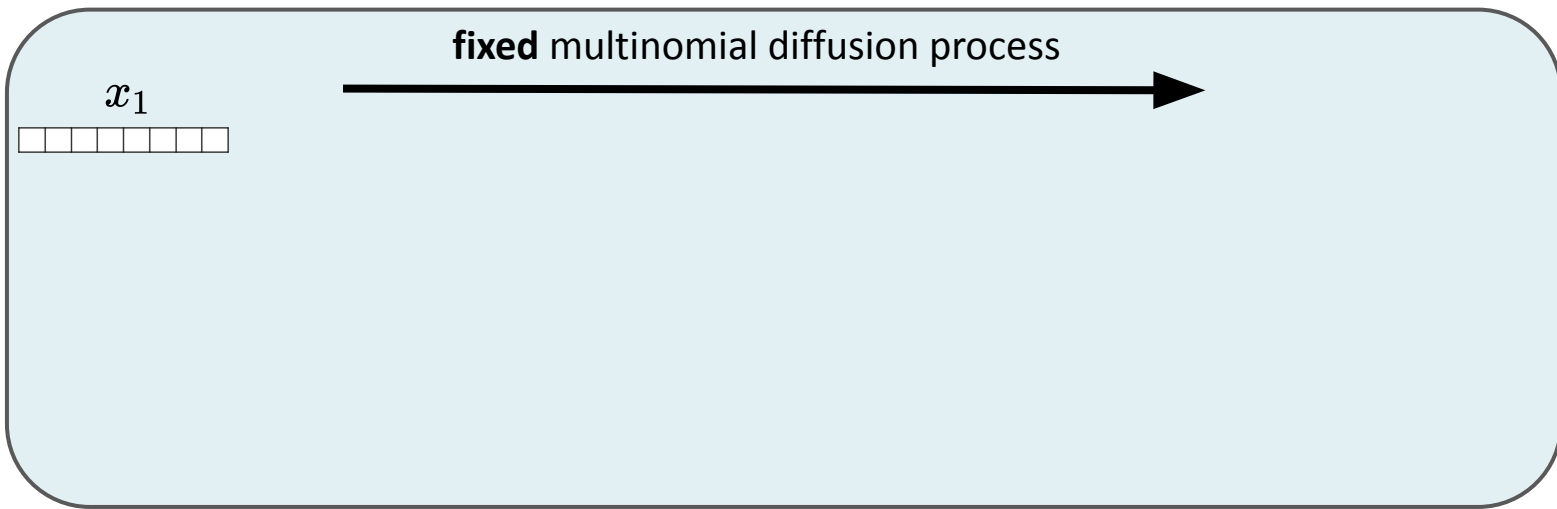
Unrealistic samples

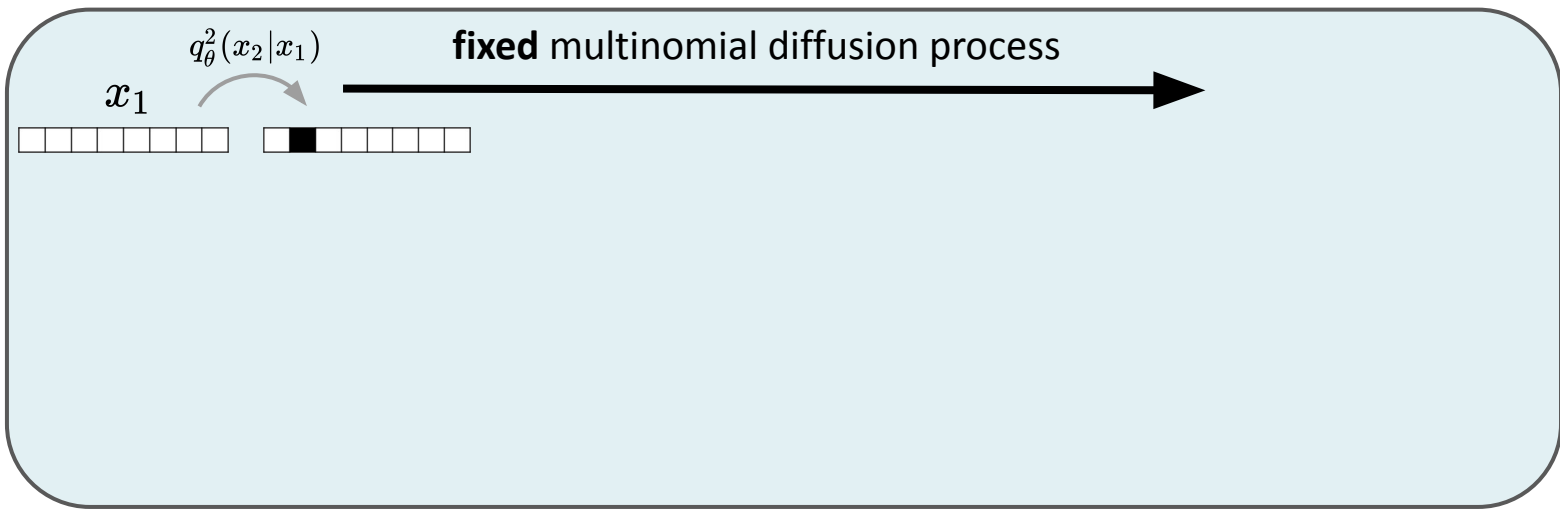
**How can we incorporate global context  
for autoregressive modeling?**

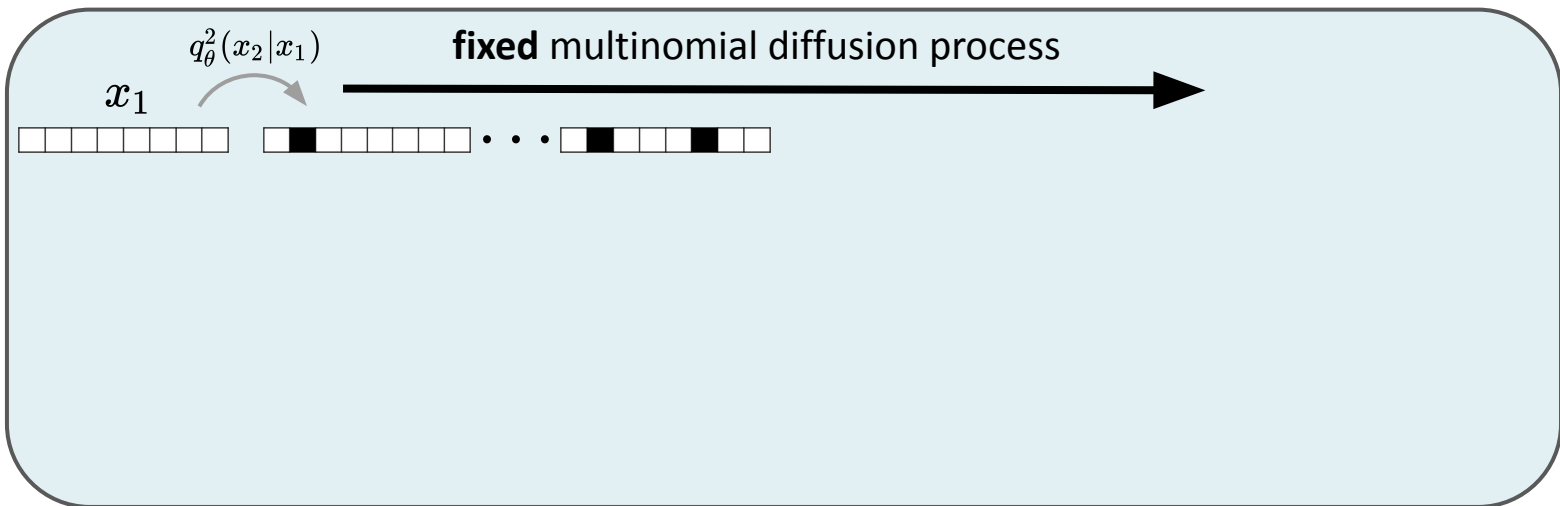
**fixed** multinomial<sup>1</sup> diffusion process



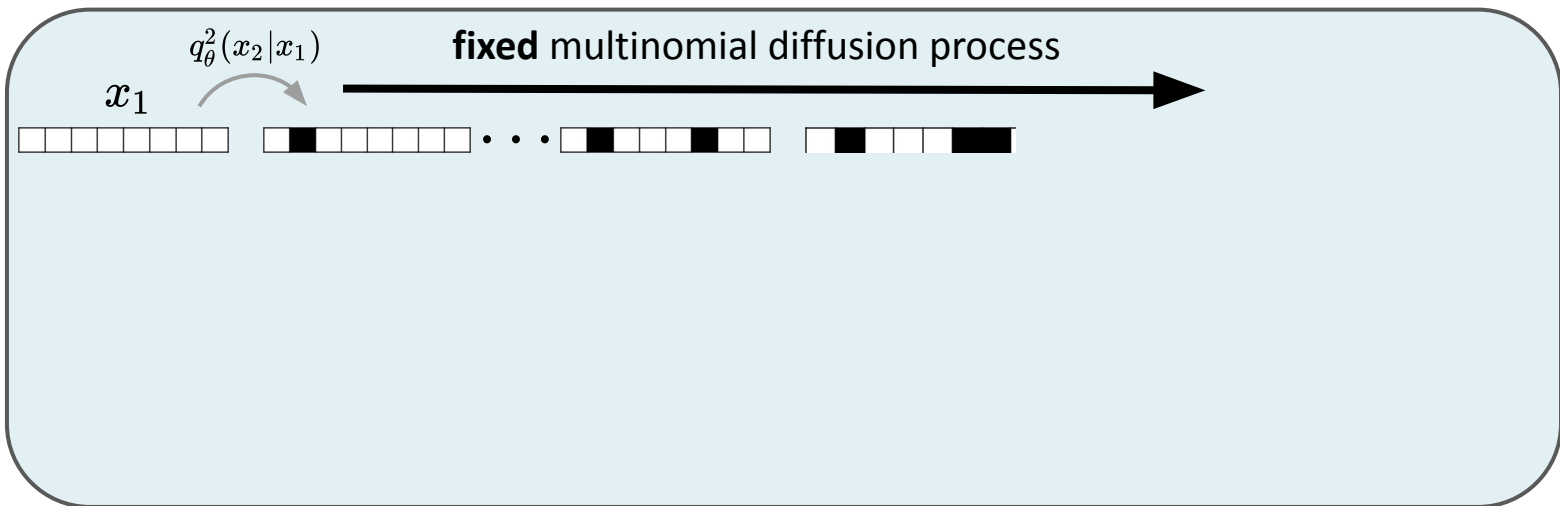
<sup>1</sup> see e.g.: *Argmax Flows and Multinomial Diffusion: Learning Categorical Distributions*, Nielsen et al, 2021

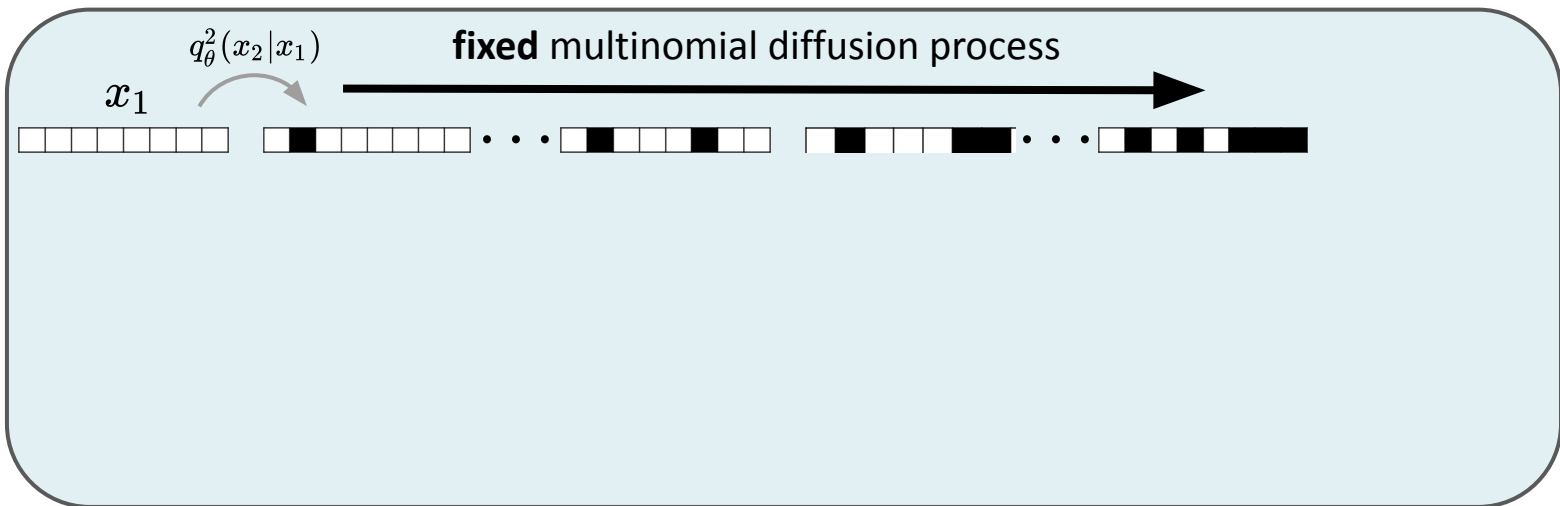


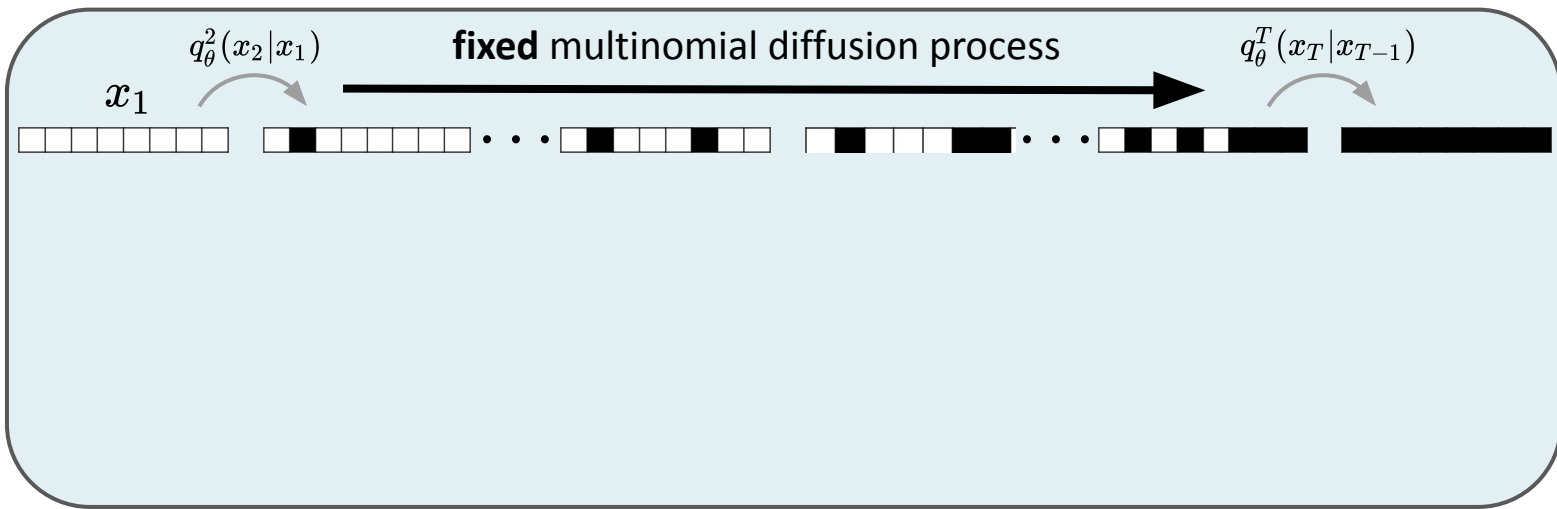


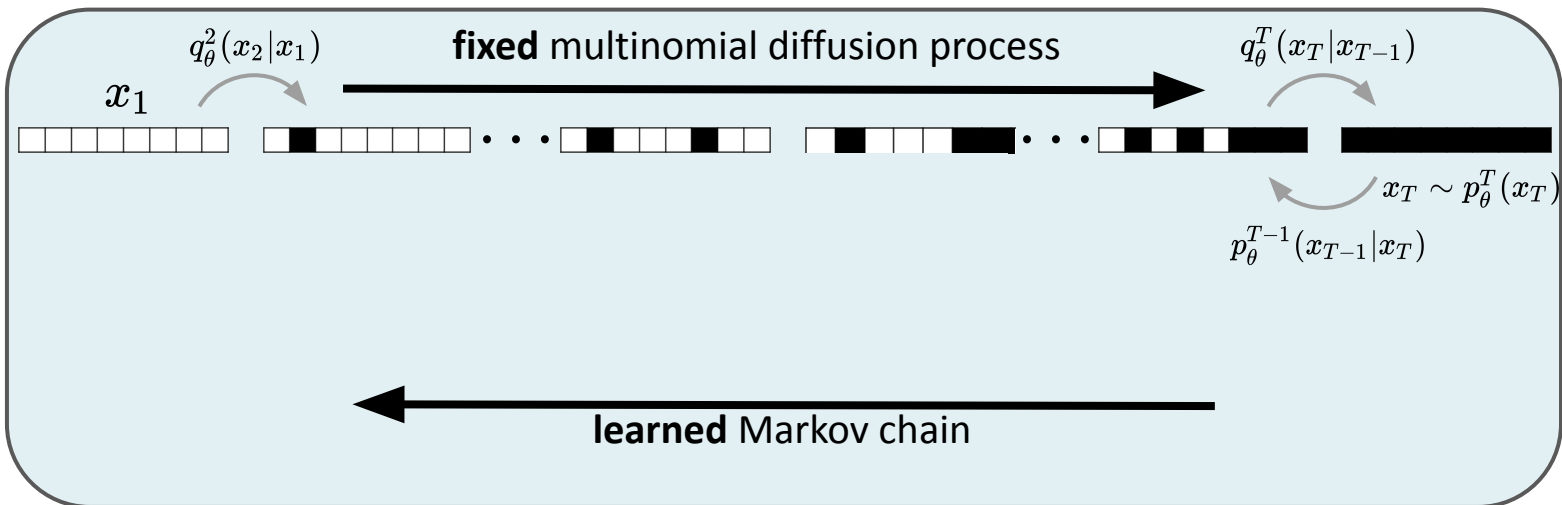


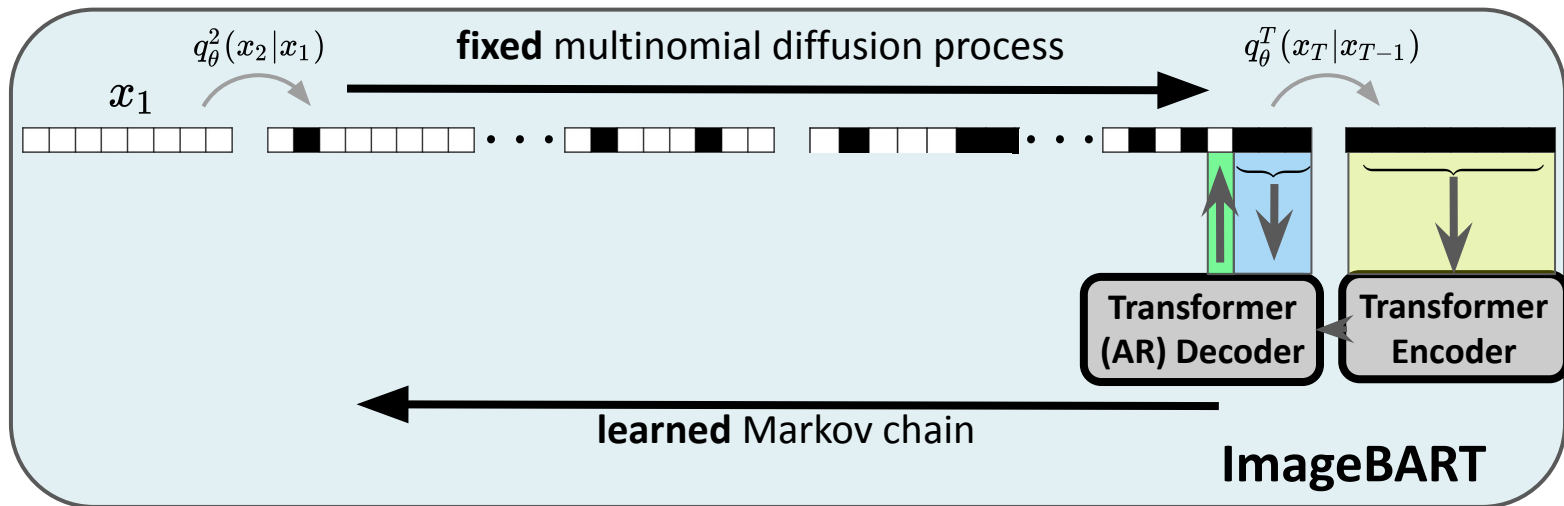


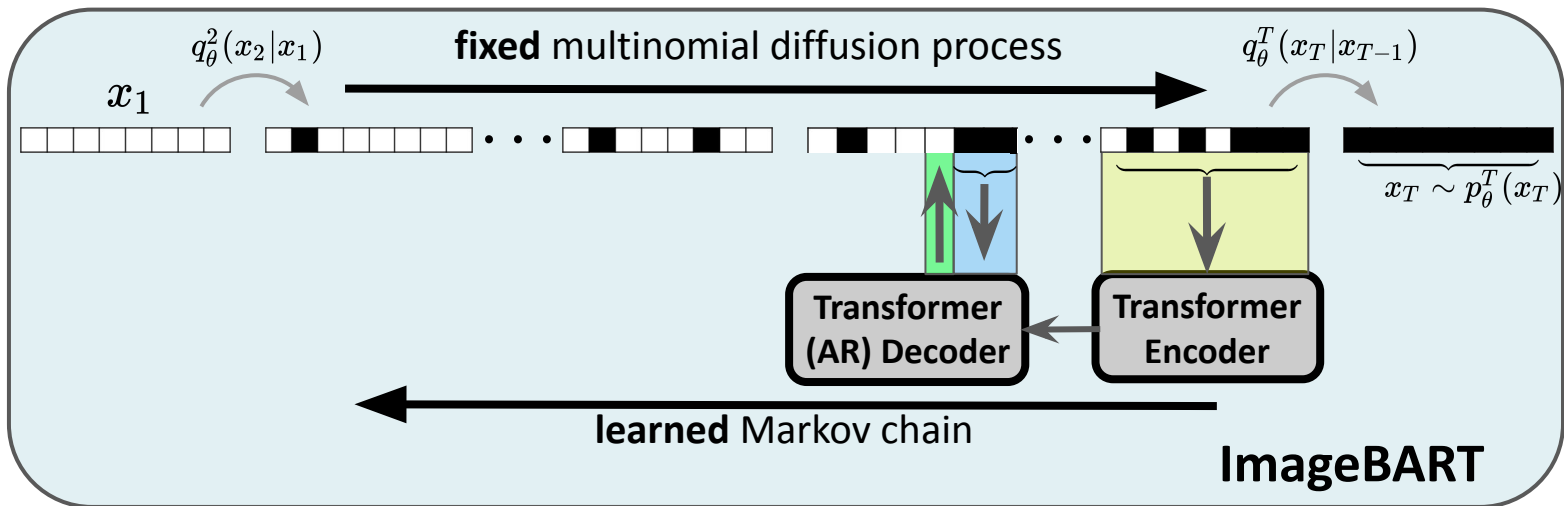


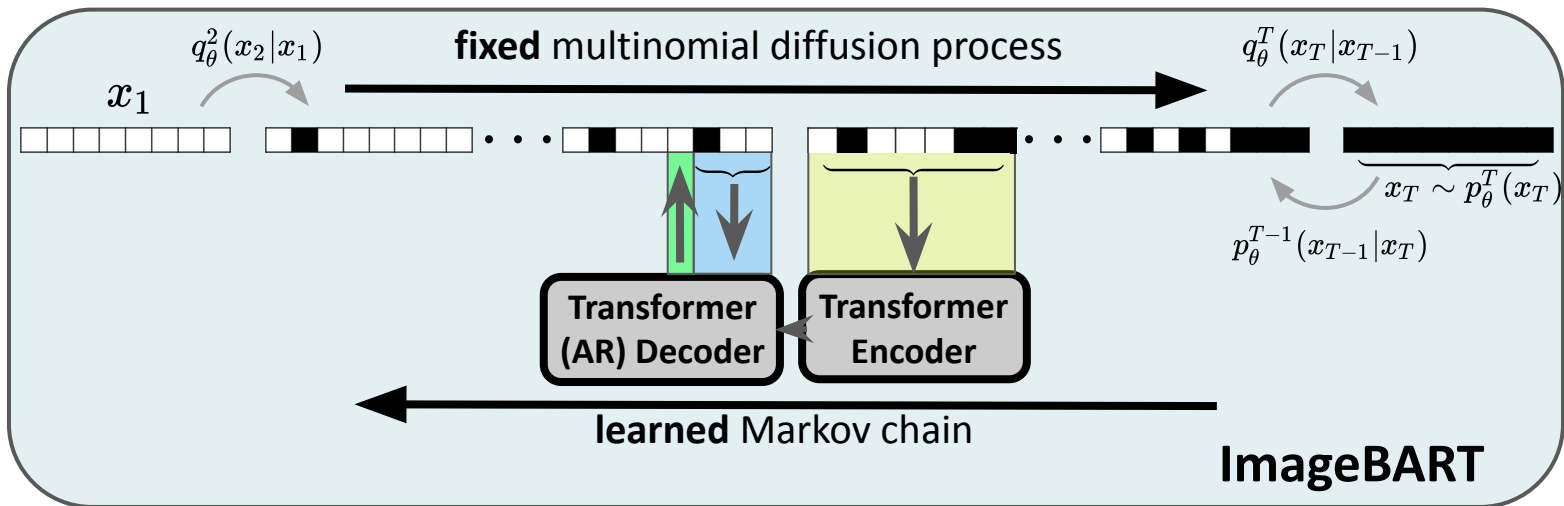


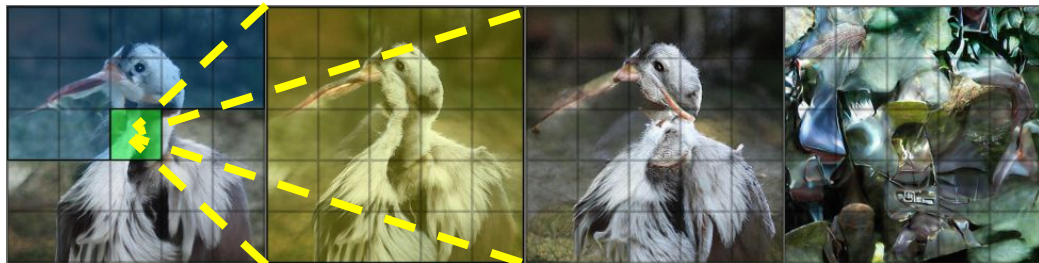
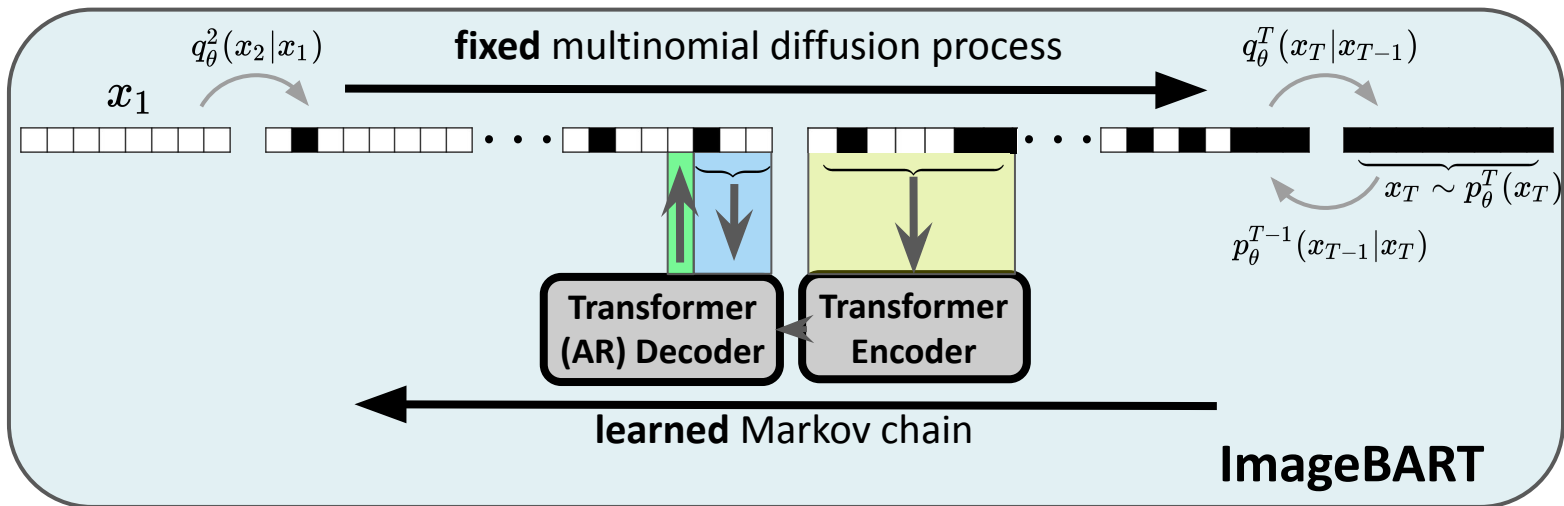




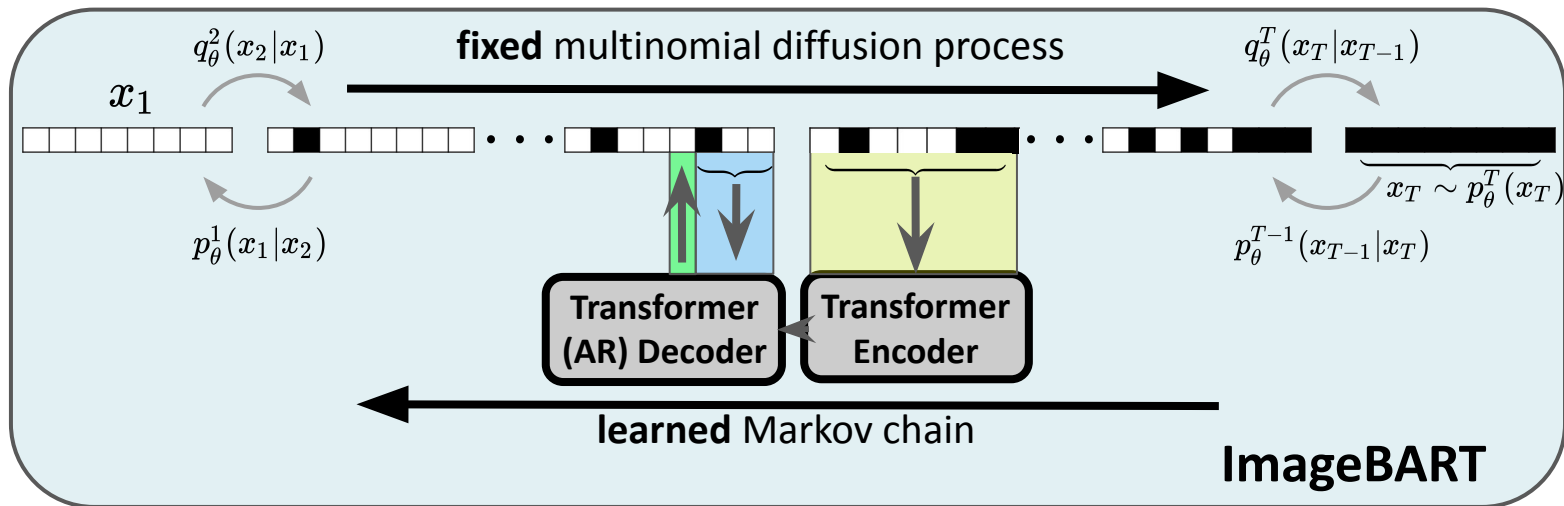


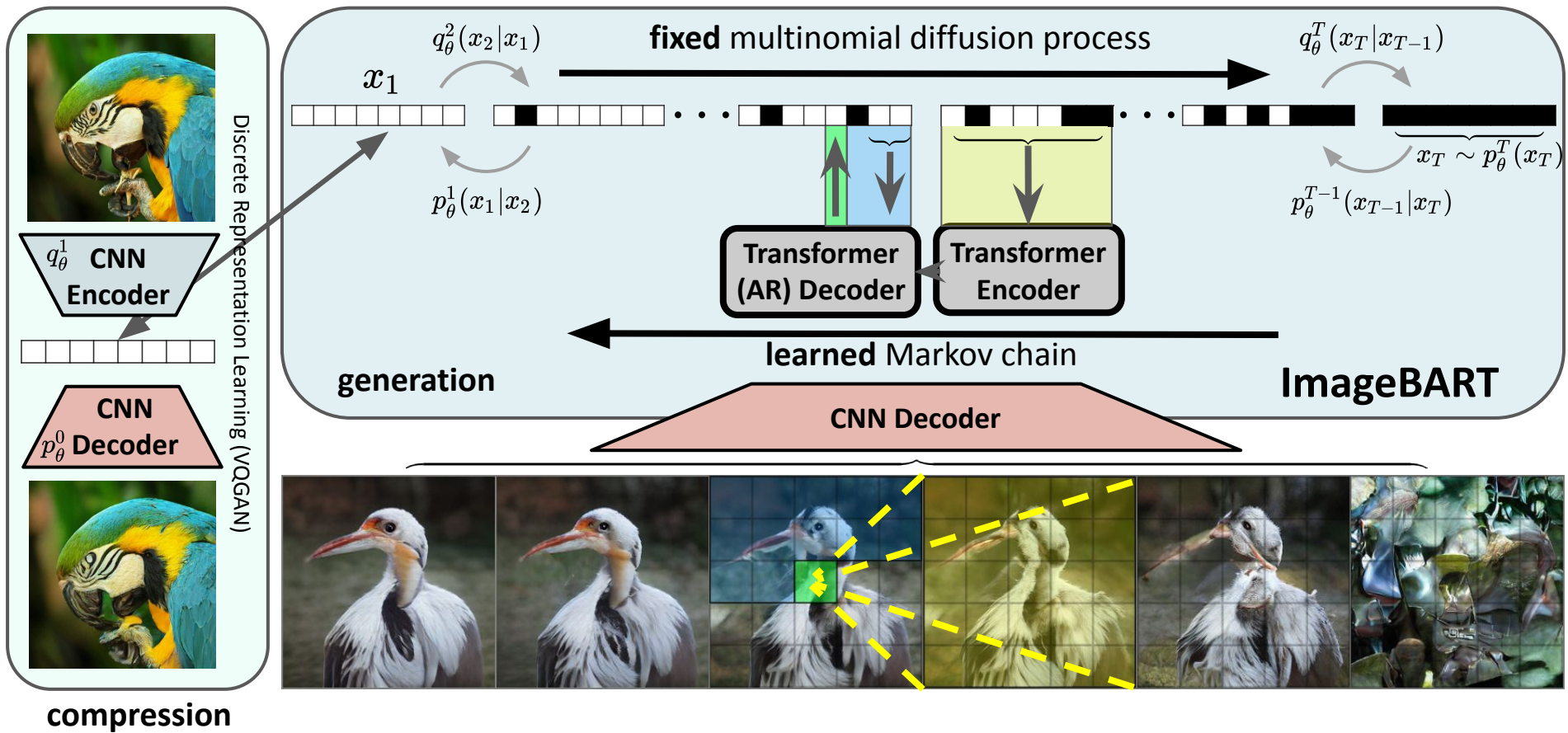








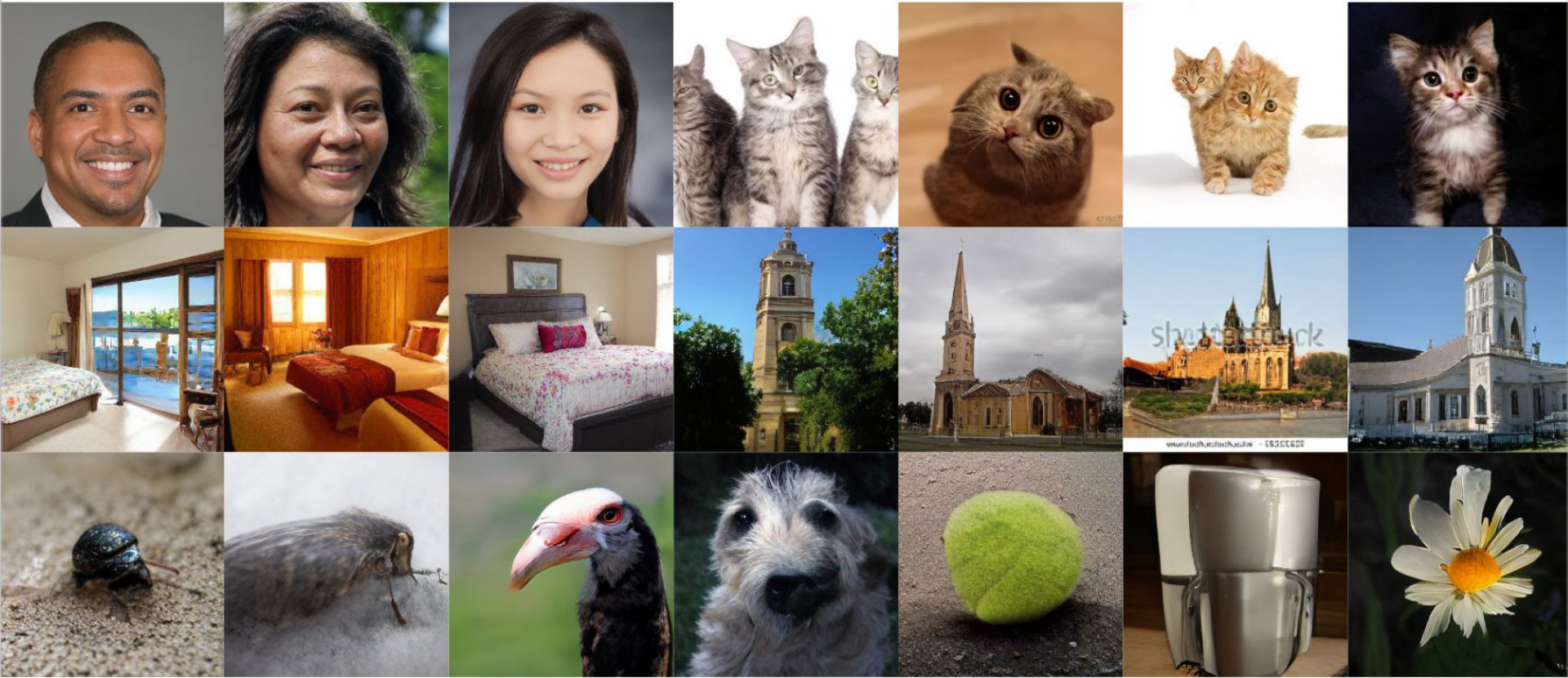




compression

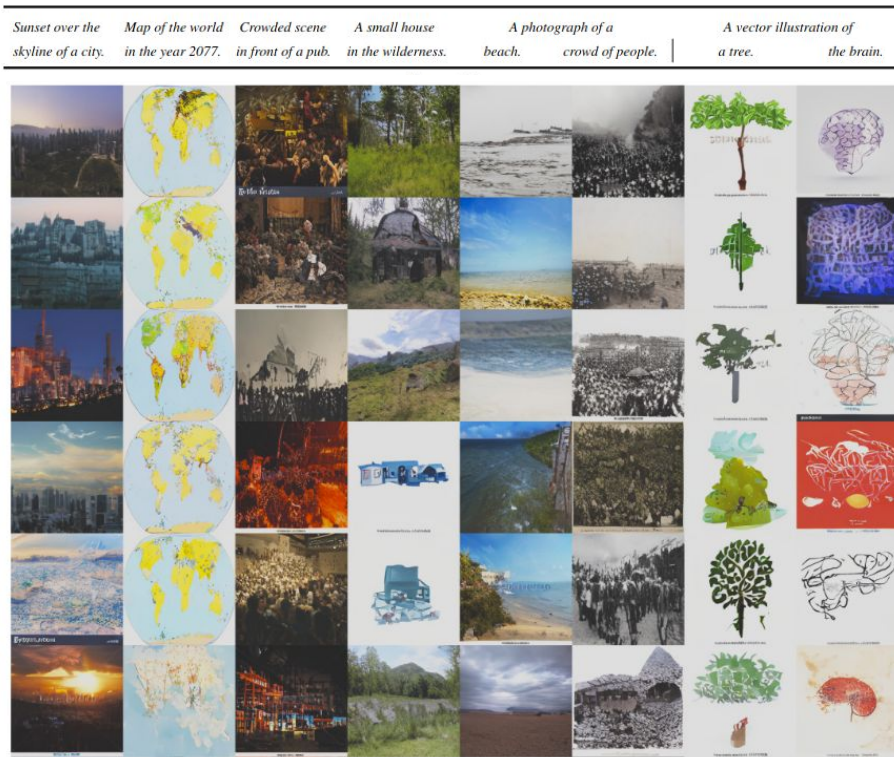
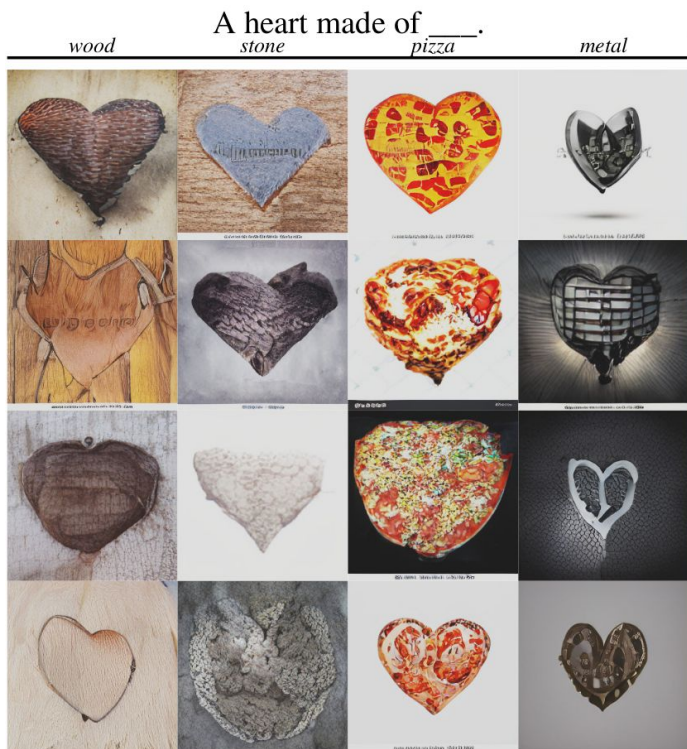
# Basic Image Synthesis

FFHQ, LSUN, ImageNet, ...



# Conditional Image Synthesis

Txt2Img (Conceptual Captions),



# Global Context for Autoregressive Image Synthesis

Masked Input

TT [18]

ImageBART

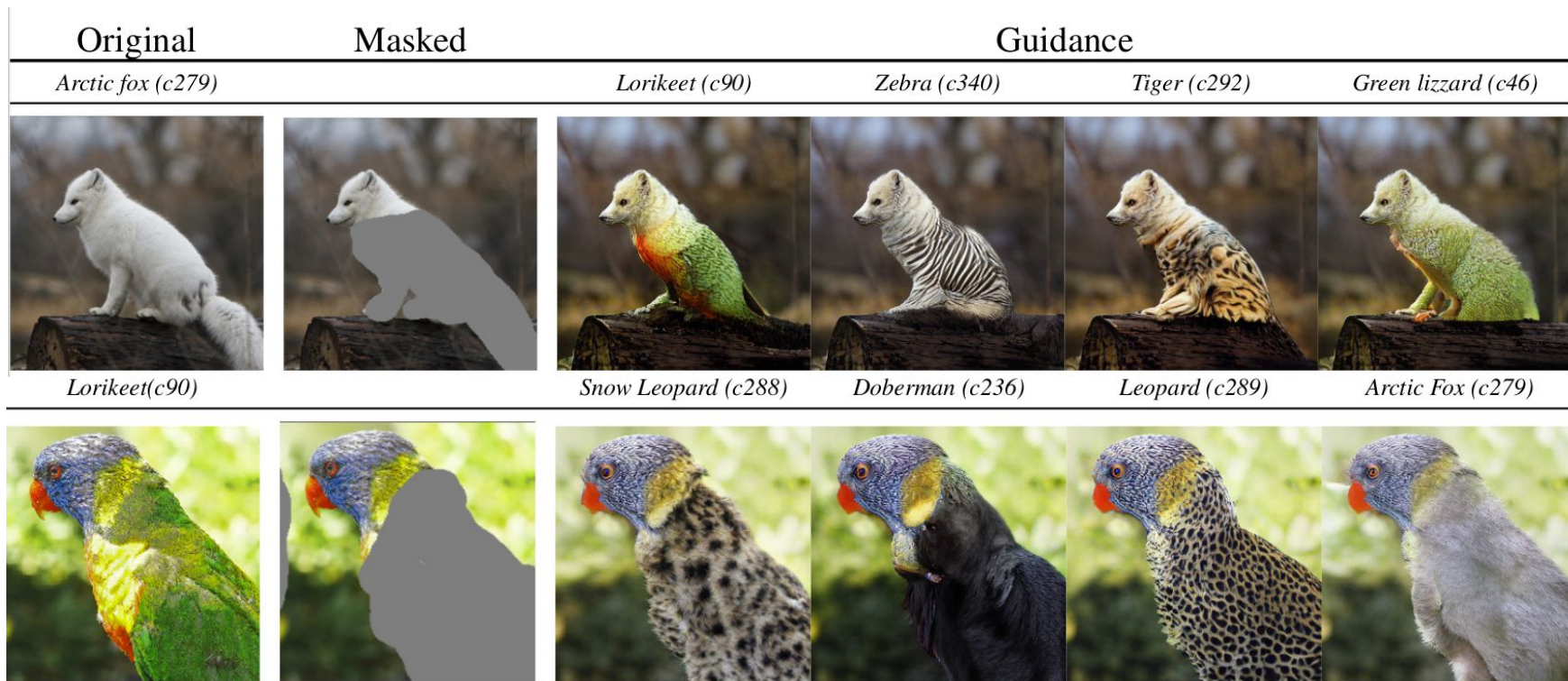


# Arbitrary Image Completion (a.k.a. Inpainting)



# Controllable Inpainting/Modification

e.g. via class labels...



# Controllable Inpainting/Modification

...or text

Original

Masked

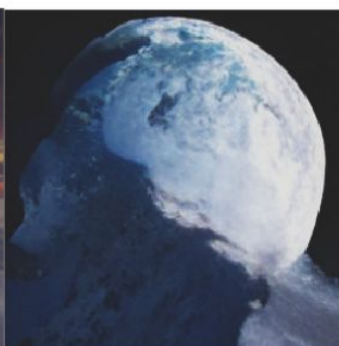
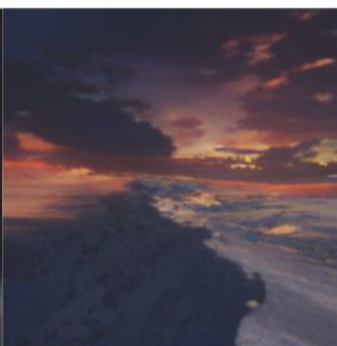
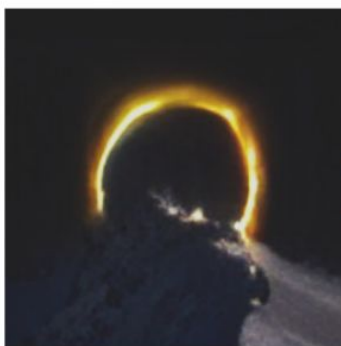
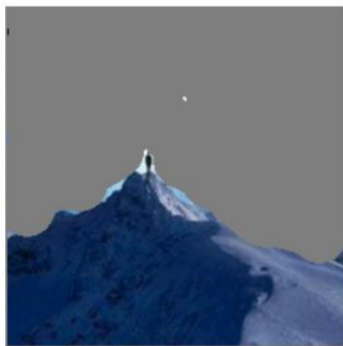
Guidance

*'Man standing on a mountain.'*

*'Solar Eclipse.'*

*'Sunrise.'*

*'Moonlight'*

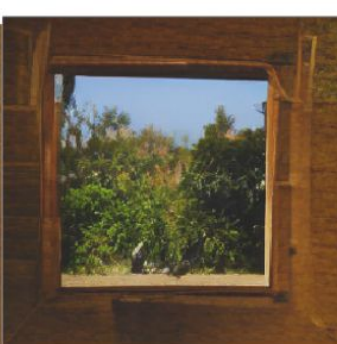


*'The piece of paper.'*

*'A pencil sketch.'*

*'A forest behind the window.'*

*'Oil painting of a cathedral.'*





# Do we really need more steps?

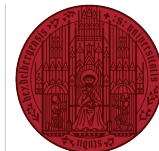
Unconditional Generation			Upper Half Completion		
method	FID ↓	IS ↑	method	FID ↓	IS ↑
TT ( $T = 2$ )	12.44	$3.98 \pm 0.07$	TT ( $T = 2$ )	11.80	$4.48 \pm 0.10$
ImageBART ( $T = 3$ )	12.55	$3.98 \pm 0.07$	ImageBART ( $T = 3$ )	9.25	$4.49 \pm 0.13$
ImageBART ( $T = 5$ )	10.69	$4.27 \pm 0.05$	ImageBART ( $T = 5$ )	6.87	$4.81 \pm 0.13$
ImageBART ( $T = 9$ )	10.81	$4.49 \pm 0.05$	ImageBART ( $T = 9$ )	6.64	$4.86 \pm 0.15$

Code and pretrained models at  
<https://github.com/CompVis/imagebart>

**THE FOXCHAIN**



Thanks!



UNIVERSITÄT  
HEIDELBERG  
ZUKUNFT  
SEIT 1386

# ImageBART: Bidirectional Context with Multinomial Diffusion for Autoregressive Image Synthesis

Patrick Esser\*,



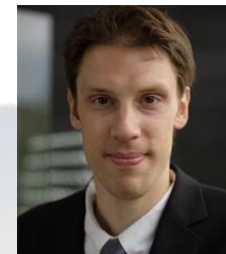
Robin Rombach\*,



Andreas Blattmann\*,



Björn Ommer



\*equal contribution

# Arbitrary Image Completion (a.k.a. Inpainting)

