

Lecture 3.3

Posters and

Intro to Diffusion Models

Dimitrios Papadopoulos
Associate Professor, DTU Compute

Yesterday

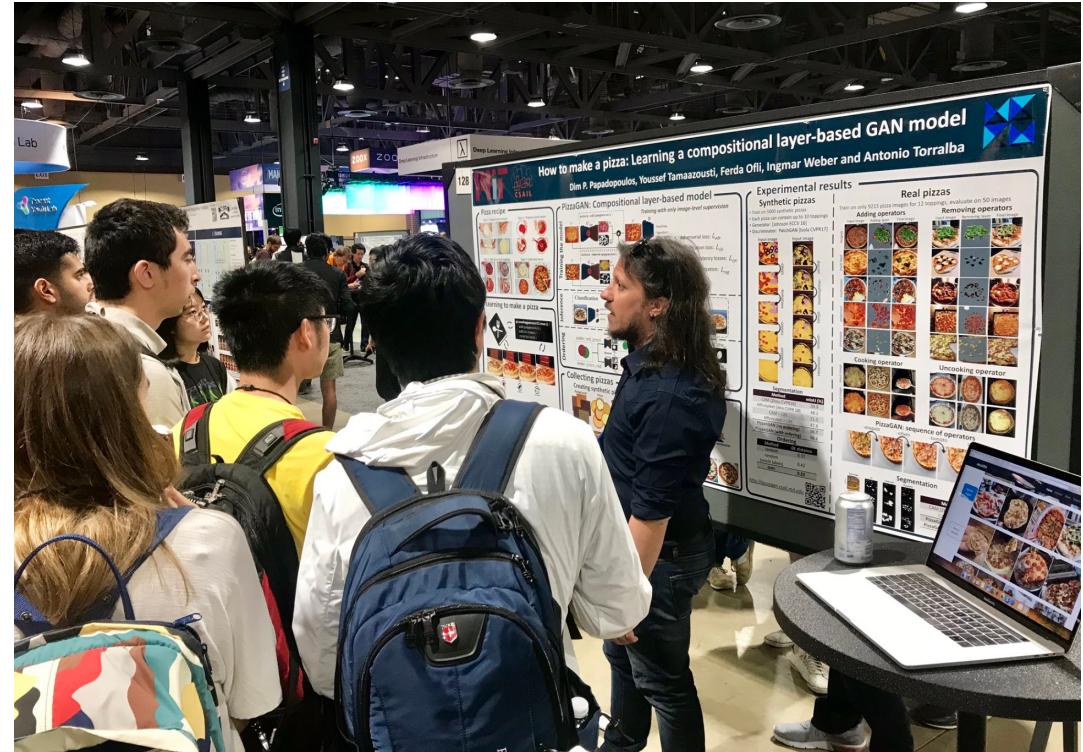
- Progressive GANs, StyleGAN, StyleGAN2, StyleGAN2-ADA
- Interpolation in latent space
- Reconstructing real images
- Learn and apply latent directions
- Generate images from text prompts (DALLE, DALLE2, Imagen)
- CLIP + StyleGAN
- Project 2

Today

- How to make a poster (Tips/ Feedback)
- Quick introduction to diffusion models
- Poster session and feedback to other groups

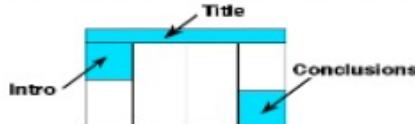


Poster sessions



Section 1 (sizes):

- Posters boards are 48" tall and 96" wide, but we recommend you leave a little border since you may not be able to pin at the vertical edge. Since PowerPoint does not let one define such a large paper size, this template is designed to be printed at 200%, yielding a 46" x94" poster. You can scale it up or down a bit (e.g. 42" is a common paper size at FexEd). Note there is no direct international A0.. A1 equivalent. The poster size is approximately three A0 boards next to each other, i.e., each column in this example is about one A0 board.
- Ideally you want to keep it very readable: this is not your paper, it is a poster. 32pt here (64 final printing) is good for most text:
 - Sub-bullets are 28 here (56 final)
 - Don't use smaller than 24pt in this template (which is 48pt in final printing at 200%)
 - Insert plenty of graphics and any math you need
- When inserting graphics or equations, keep the resolution high (remember this will be printed at 200%). If you can see blocking artifacts at 400% magnification in PowerPoint, consider finding better graphics. This is an example of BAD/LOW RES GRAPHICS



- Leave enough margin for pushpin and remember many big plotters cannot get within .5" of the actual paper edge.
- You are free to use colored backgrounds and such but they generally reduce readability.
- You are free to use what ever fonts you like.
 - San Serif fonts like Arial are more readable from a distance,
 - Serif fonts like times may look more consistent with your mathematics

Section 2 (layout):

- Remember the poster session will be crowded so design the poster to be read in columns so people can read what is in front of them and move left to right to get the whole story.
- The poster should use photos, figures, and tables to tell the story of the study. For clarity, present the information in a sequence that is easy to follow.
- There is often way too much text in a poster - there definitely is in this template! Posters primarily are visual presentations; the text should support the graphics. Look critically at the layout. Some poster 'experts' suggest that if there is about 20-25% text, 40-45% graphics and 30-40% empty space, you are doing well.



Section 3:

- Include more figures than are in the paper so you can talk to them. Include things that are not in the paper and then encourage them to read the paper. Don't try to just put all the paper here.
- If it looks like a cut/paste of the paper, people skip that poster since they can read the papers after the conference. Many people find it better to spend time talking with poster presenters that have more to offer than just redoing the paper content paper in big fonts.
- People will likely have already seen your posted video, so the poster can serve as a talking point for you.
- Remember Poster boards look like this.. This is your canvas. Paint us a picture of your work.



Summary/Conclusion

- Summarize your contributions
- Summarize your results (if applicable)
- You can add in links to additional videos, code, or project website (or QR code)

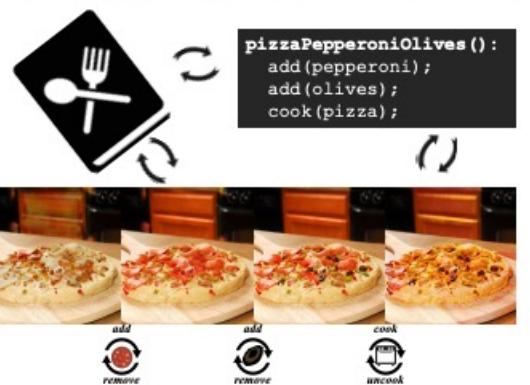
References

How to make a pizza: Learning a compositional layer-based GAN model

Dim P. Papadopoulos, Youssef Tamaazousti, Ferda Ofli, Ingmar Weber and Antonio Torralba

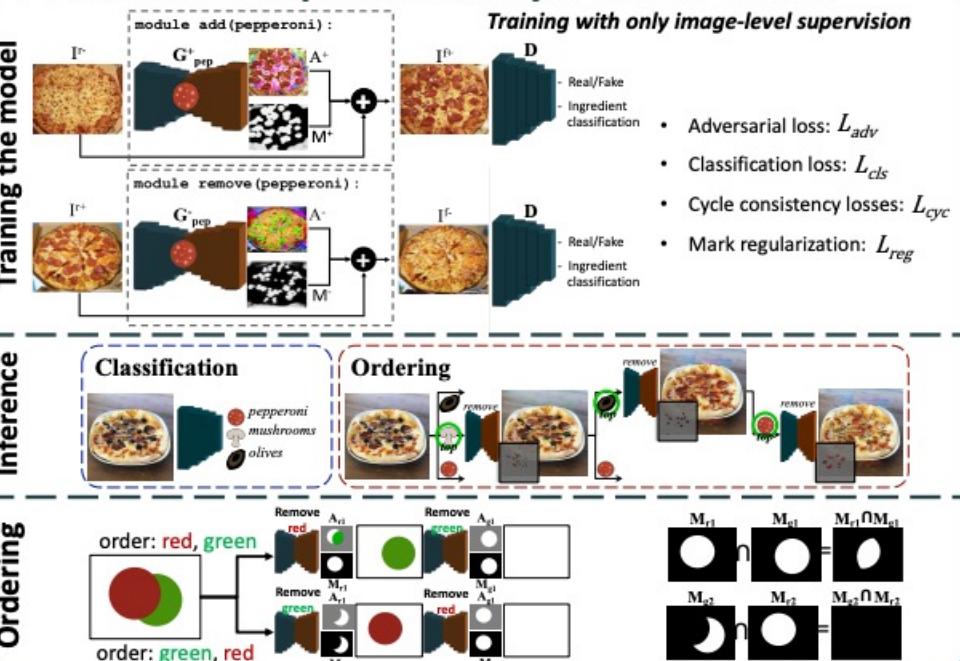


Learning to make a pizza



- Build a generative model to mirror the recipe procedure.
 - Learn composable modules that carry on different cooking operations (e.g. add/remove ingredient).
 - Decompose an image into an ordered sequence of layers.

- PizzaGAN: Compositional layer-based model



- Collecting pizzas



Experimental results

Synthetic pizzas

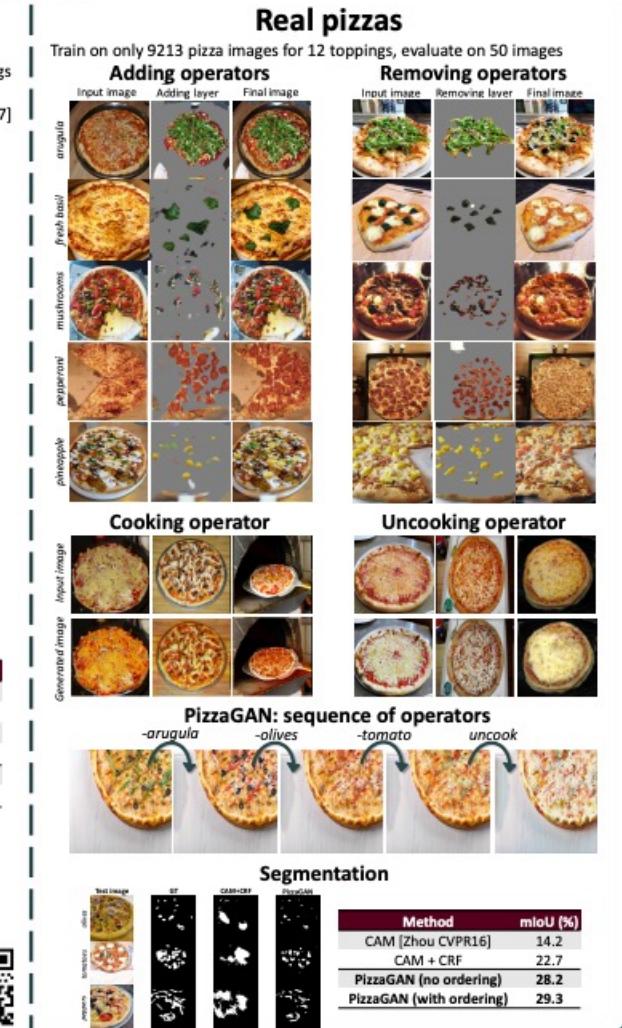
- Train on 5000 synthetic pizzas
 - Each pizza can contain up to 10 toppings
 - Generator: [Johnson ECCV 16]
 - Discriminator: PatchGAN [Isola CVPR17]



Segmentation		
Method	mIoU	
CAM [Zhou CVPR16]	39.3	
AffinityNet [Ahn CVPR 18]	48.3	
CAM + CRF	51.3	
AffinityNet + CRF	47.3	
PizzanGAN (no ordering)	56.3	
PizzanGAN (with ordering)	58.3	
Ordering		

Ordering	
Method	DL distance
random	0.91
random (oracle labels)	0.42
ours	0.33

<http://pizzagan.csail.mit.edu>

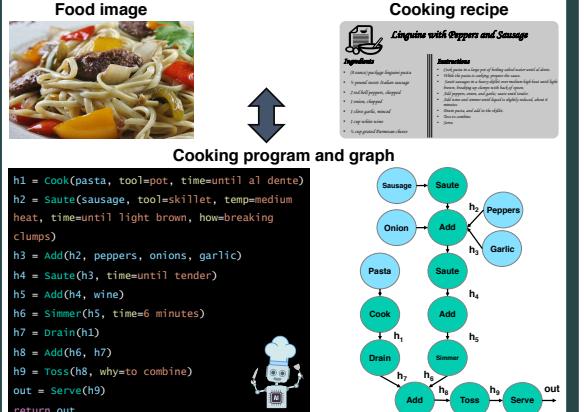


Introduction

Food understanding

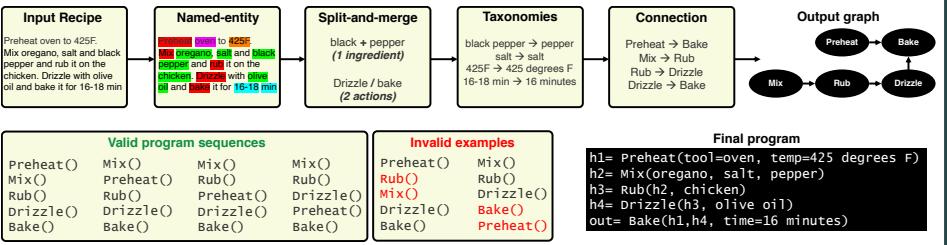


Program representation



- ✓ provide a non-ambiguous and structured presentation
 - ✓ capture cooking semantics and action relationships
 - ✓ can be easily manipulated by users
 - ✓ can be potentially executed by agents

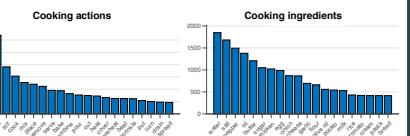
Cooking programs ▶



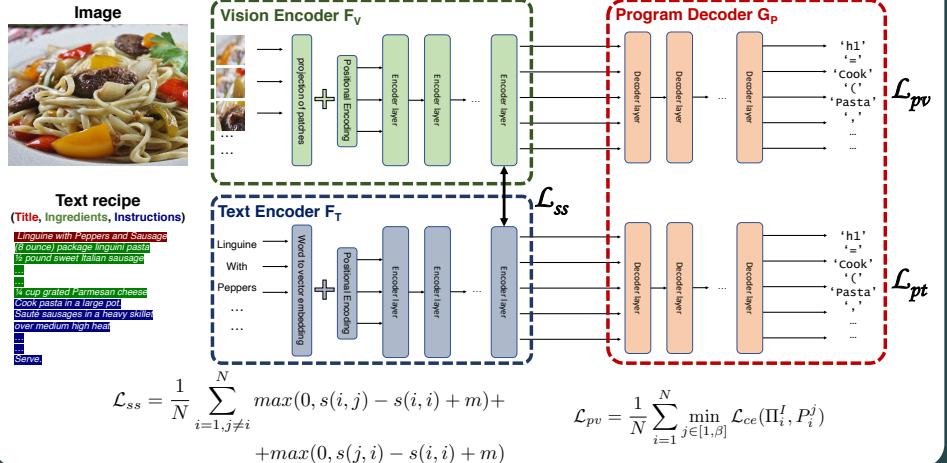
Program collection

- **3,708 programs** from Recipe1M dataset [Salvador CVPR17]
 - **42,473 annotated sentences** with **478,285 tagged words**
 - **Total annotation time:** 4 minutes per recipe
 - **Human agreement:** 97.9% Named-entity and 92.7% Connection

Sentence Length (words)	Number of Sentences
1	~5800
2	~1200
3	~500
4	~200
5	~100
6	~50
7	~20
8	~10
9	~5
10	~3
11	~2
12	~1
13	~1
14	~1
15	~1
16	~1
17	~1
18	~1
19	~1
20	~1
21	~1
22	~1
23	~1
24	~1
25	~1
26	~1
27	~1
28	~1
29	~1
30	~1
31	~1
32	~1
33	~1
34	~1
35	~1
36	~1
37	~1
38	~1
39	~1
40	~1
41	~1
42	~1
43	~1
44	~1
45	~1
46	~1
47	~1
48	~1
49	~1
50	~1
51	~1
52	~1
53	~1
54	~1
55	~1
56	~1
57	~1
58	~1
59	~1
60	~1
61	~1
62	~1
63	~1
64	~1
65	~1
66	~1
67	~1
68	~1
69	~1
70	~1
71	~1
72	~1
73	~1
74	~1
75	~1
76	~1
77	~1
78	~1
79	~1
80	~1
81	~1
82	~1
83	~1
84	~1
85	~1
86	~1
87	~1
88	~1
89	~1
90	~1
91	~1
92	~1
93	~1
94	~1
95	~1
96	~1
97	~1
98	~1
99	~1
100	~1
101	~1
102	~1
103	~1
104	~1
105	~1
106	~1
107	~1
108	~1
109	~1
110	~1
111	~1
112	~1
113	~1
114	~1
115	~1
116	~1
117	~1
118	~1
119	~1
120	~1
121	~1
122	~1
123	~1
124	~1
125	~1
126	~1
127	~1
128	~1
129	~1
130	~1
131	~1
132	~1
133	~1
134	~1
135	~1
136	~1
137	~1
138	~1
139	~1
140	~1
141	~1
142	~1
143	~1
144	~1
145	~1
146	~1
147	~1
148	~1
149	~1
150	~1
151	~1
152	~1
153	~1
154	~1
155	~1
156	~1
157	~1
158	~1
159	~1
160	~1
161	~1
162	~1
163	~1
164	~1
165	~1
166	~1
167	~1
168	~1
169	~1
170	~1
171	~1
172	~1
173	~1
174	~1
175	~1
176	~1
177	~1
178	~1
179	~1
180	~1
181	~1
182	~1
183	~1
184	~1
185	~1
186	~1
187	~1
188	~1
189	~1
190	~1
191	~1
192	~1
193	~1
194	~1
195	~1
196	~1
197	~1
198	~1
199	~1
200	~1
201	~1
202	~1
203	~1
204	~1
205	~1
206	~1
207	~1
208	~1
209	~1
210	~1
211	~1
212	~1
213	~1
214	~1
215	~1
216	~1
217	~1
218	~1
219	~1
220	~1
221	~1
222	~1
223	~1
224	~1
225	~1
226	~1
227	~1
228	~1
229	~1
230	~1
231	~1
232	~1
233	~1
234	~1
235	~1
236	~1
237	~1
238	~1
239	~1
240	~1
241	~1
242	~1
243	~1
244	~1
245	~1
246	~1
247	~1
248	~1
249	~1
250	~1
251	~1
252	~1
253	~1
254	~1
255	~1
256	~1
257	~1
258	~1
259	~1
260	~1
261	~1
262	~1
263	~1
264	~1
265	~1
266	~1
267	~1
268	~1
269	~1
270	~1
271	~1
272	~1
273	~1
274	~1
275	~1
276	~1
277	~1
278	~1
279	~1
280	~1
281	~1
282	~1
283	~1
284	~1
285	~1
286	~1
287	~1
288	~1
289	~1
290	~1
291	~1
292	~1
293	~1
294	~1
295	~1
296	~1
297	~1
298	~1
299	~1
300	~1
301	~1
302	~1
303	~1
304	~1
305	~1
306	~1
307	~1
308	~1
309	~1
310	~1
311	~1
312	~1
313	~1
314	~1
315	~1
316	~1
317	~1
318	~1
319	~1
320	~1
321	~1
322	~1
323	~1
324	~1
325	~1
326	~1
327	~1
328	~1
329	~1
330	~1
331	~1
332	~1
333	~1
334	~1
335	~1
336	~1
337	~1
338	~1
339	~1
340	~1
341	~1
342	~1
343	~1
344	~1
345	~1
346	~1
347	~1
348	~1
349	~1
350	~1
351	~1
352	~1
353	~1
354	~1
355	~1
356	~1
357	~1
358	~1
359	~1
360	~1
361	~1
362	~1
363	~1
364	~1
365	~1
366	~1
367	~1
368	~1
369	~1
370	~1
371	~1
372	~1
373	~1
374	~1
375	~1
376	~1
377	~1
378	~1
379	~1
380	~1
381	~1
382	~1
383	~1
384	~1
385	~1
386	~1
387	~1
388	~1
389	~1
390	~1
391	~1
392	~1
393	~1
394	~1
395	~1
396	~1
397	~1
398	~1
399	~1
400	~1
401	~1
402	~1
403	~1
404	~1
405	~1
406	~1
407	~1
408	~1
409	~1
410	~1
411	~1
412	~1
413	~1
414	~1
415	~1
416	~1
417	~1
418	~1
419	~1
420	~1
421	~1
422	~1
423	~1
424	~1
425	~1
426	~1
427	~1
428	~1
429	~1
430	~1
431	~1
432	~1
433	~1
434	~1
435	~1
436	~1
437	~1
438	~1
439	~1
440	~1
441	~1
442	~1
443	~1
444	~1
445	~1
446	~1
447	~1
448	~1
449	~1
450	~1
451	~1
452	~1
453	~1
454	~1
455	~1
456	~1
457	~1
458	~1
459	~1
460	~1
461	~1
462	~1
463	~1
464	~1
465	~1
466	~1
467	~1
468	~1
469	~1
470	~1
471	~1
472	~1
473	~1
474	~1
475	~1
476	~1
477	~1
478	~1
479	~1
480	~1
481	~1
482	~1
483	~1
484	~1
485	~1
486	~1
487	~1
488	~1
489	~1
490	~1
491	~1
492	~1
493	~1
494	~1
495	~1
496	~1
497	~1
498	~1
499	~1
500	~1



From food images and recipes to programs



Experiments

- Recipe1M** [Salvador CVPR17], 340k recipes with images
Vision encoder: ViT-B/16 [Dosovitskiy ICLR21]
Text encoder and Program decoder: Transformer [Vaswani NeurIPS17]

Image-to-recipe retrieval



50 epochs						
Adam optimizer						
0.1 step decay every 20 epochs						
image-to-image				recipeto-image		
R _{E1}	R _{E2}	R _{E10}	R@10	medR _{E1}	R _{E1}	R@10
5	24.0	54.4	67.8	4.5	23.8	54.2
39.1	68.1	80.3	3.0	38.9	68.4	80.5
36.6	65.2	76.8	3.0	36.5	65.8	76.9
43.6	75.6	85.0	2.0	44.9	76.2	85.4
53.5	81.8	92.0	1.0	53.1	82.0	89.6
58.6	85.7	91.7	1.0	58.2	85.5	92.0
66.9	90.9	95.1	1.0	66.8	89.8	94.6
image-to-image				recipeto-image		
R _{E1}	R _{E2}	R _{E10}	R@10	medR _{E1}	R _{E1}	R@10
5	5.0	51.0	65.0	5.1	52.0	65.0
25.6	53.7	66.9	4.6	25.7	53.9	67.1
39.8	69.0	77.4	1.0	40.2	68.1	78.7
39.1	71.0	81.7	2.0	40.6	72.6	83.3
48.2	75.8	83.6	1.9	48.4	76.1	83.7
49.7	79.5	86.3	1.9	50.1	80.2	86.7
60.2	87.5	93.0	1.0	52.8	80.2	87.6
54.0	81.8	87.8	1.0	54.9	81.9	89.0
60.2	84.0	89.7	1.0	60.3	87.6	93.2
60.0	87.6	97.4	1.0	60.3	87.6	93.2
63.2	88.3	93.1	—	—	—	—
66.9	90.9	95.1	1.0	66.8	89.8	94.6

Program prediction



Image generation



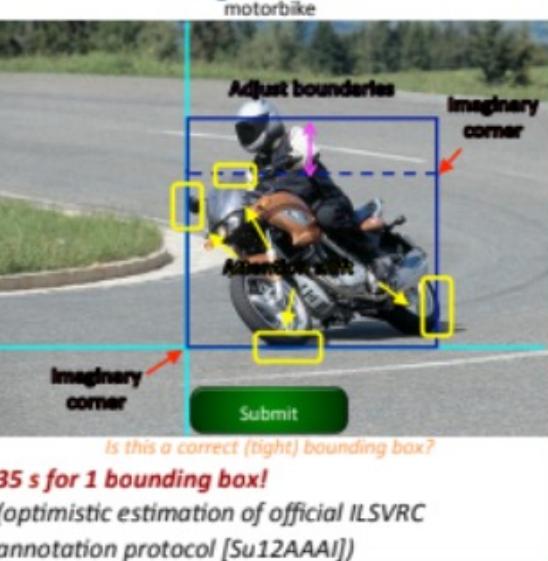


Extreme clicking for efficient object annotation

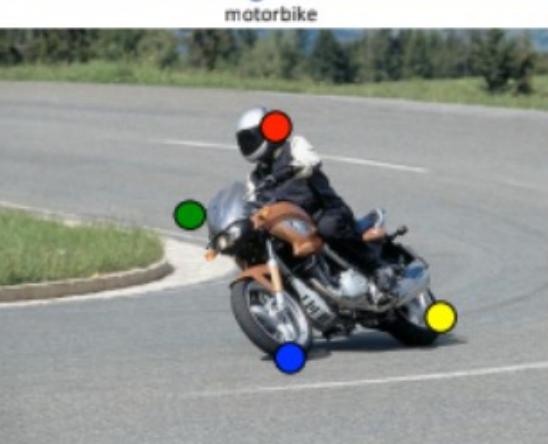
Dim P. Papadopoulos, Jasper R. R. Uijlings, Frank Keller and Vittorio Ferrari



Draw bounding boxes

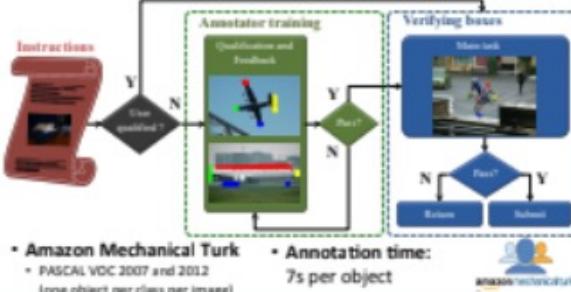


Extreme clicking



- Click on four extreme points: **left-most, right-most, top, bottom**
- Attention is on single point at a time
- Physical points on object, not imaginary
- No rectangle is involved, neither real nor imaginary
- Single task, no task-switching

Crowd-source extreme clicks



Create object segmentation

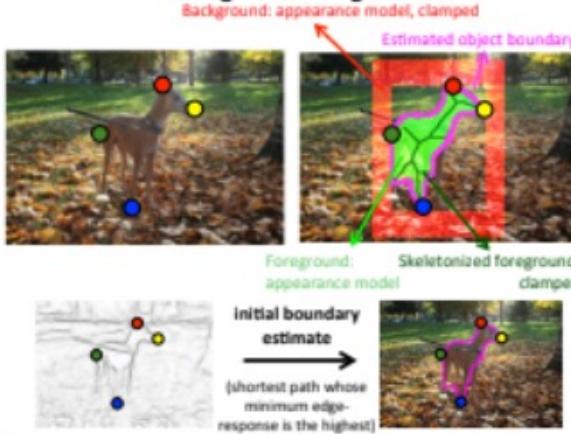
Boxes to segments using GrabCut

$$E(L) = \sum_p U(l_p) + \sum_{p,q} V(l_p, l_q) \quad [\text{Rother SIGGRAPH 04}]$$

Background: appearance model, clamped



Extreme clicks to segments using GrabCut



Extreme clicking results

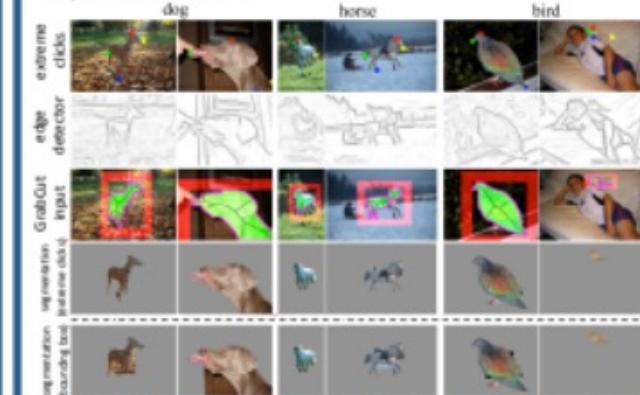
- PASCAL VOC 2007 and 2012 (1 instance per class per image)
- Quality drawing by using expert agreement: boxes vs boxed segment
- Detector: Fast-RCNN [Girshick ICCV 15]

Dataset	Approach	Annotation quality		Detection performance (mAP)	Annotation time
		mIoU	AlexNet		
VOC	Extreme clicks	88	56	66	7.0
2007	PASCAL GT boxes	88	56	66	34.5
VOC	Extreme clicks	87	52	62	7.2
2012	PASCAL GT boxes	87	52	62	34.5

5 times faster without any compromise on quality

Results on object segmentation

Qualitative results



Quantitative results

Approach	mIoU	
	VOC 2007	VOC 2012
original GrabCut	37.3	35.8
optimized GrabCut	74.4	71.0
GrabCut and extreme clicks	78.1	72.7

extreme clicks into GrabCut improve object segmentation by 2-4%

[Hariharan ICCV 2011]

DeepLab on PASCAL VOC 2012

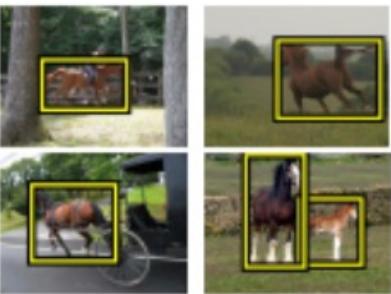
	Full supervision	Segments from GT Boxes	Segments from extreme clicks
mIoU	59.9	55.8	58.4

extreme clicks strong for weakly supervised semantic segmentation

Introduction

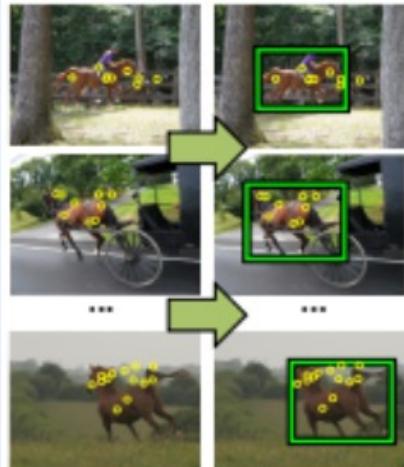
Training object detectors

Current way: draw bounding-boxes



- time consuming (26s-42s per box) [Su AAAI 2012]
- need detailed annotation guidelines

New way: bounding-boxes from eye-tracking data



- + annotation time (1s per image)
- + reduce annotation time by 6.8x
- + simple annotation guidelines
- + correct localizations in half images

Eye tracking dataset

Data

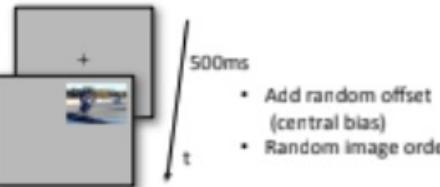
- Large scale (6270 images)
- Pascal VOC 12 : train+val images of 10 classes
- 5 distinct viewers (28 in total) for each image



- 178,000 fixations (5.7 per viewer per image)
- Mean response time = 889 ms/image
- Fixations on target objects = 75.2%

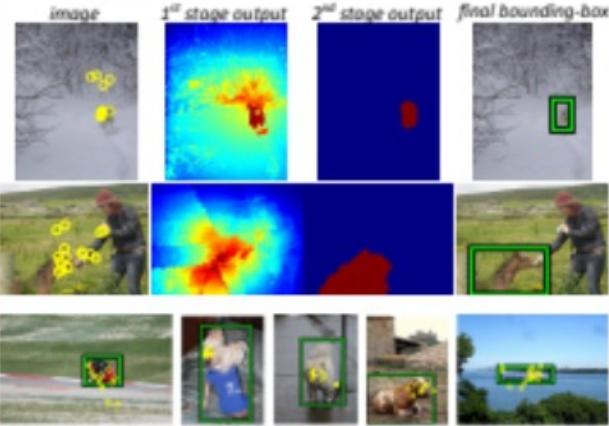
Experiment

- Visual search paradigm**
 - more fixations on target objects
 - faster than free-viewing
- Pairs of classes**
 - two-alternative forced choice object discrimination
 - pair classes with similar background (e.g. cat, dog)



- Add random offset (central bias)
- Random image order

Results



From fixations to bounding-boxes

Bounding-box estimation as figure-ground superpixel labeling

\mathcal{R}_{bb+fix}

- small subset (7%)
- learn to predict bounding-boxes from fixations



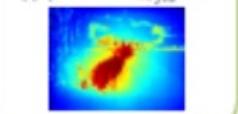
\mathcal{R}_{fix}

- derive new bounding-boxes from fixations



Initial object segmentation

- Train a superpixel classifier in \mathcal{R}_{bb+fix} (linear SVM + Platt scaling)
- Apply model in \mathcal{R}_{fix} set



Segmentation refinement

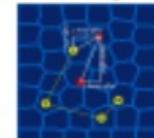
- Grabcut-like energy minimization



Features

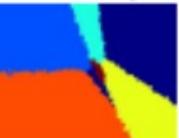
Fixation position

Visual search increases fixations on [or near] the target object



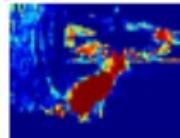
Fixation timing

Timing matters: the longer or the later, the more significant



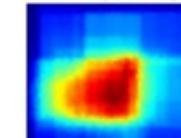
Fixation appearance

Learn color distribution of **fg** and **bg** superpixels from fixations



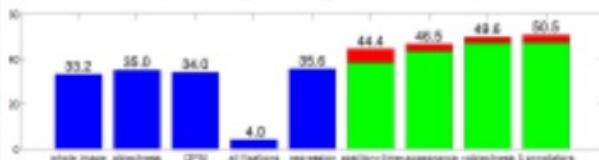
Objectness

Probability that a window contains object of **any** class



Quantitative results

- 10 Pascal VOC 12 classes, 6270 images in train+val
- Evaluate predicted bounding-boxes in \mathcal{R}_{fix}
- Performance: percentage of images with correct predictions



- + all feature types contribute
- + full model outperforms all baselines
- + segmentation refinement always helps (+ 3-5%)

Train DPM detector from fixations [Felzenszwalb PAMI10]

- Pascal VOC 12: train on train+val, test on test set (10991 images)

All ground-truth bounding boxes



mAP = 25.5%

All predicted bounding boxes



mAP = 12.5%

6.8x fewer ground-truth



mAP = 13.7%

*Given same annotation time,
get comparable mAP performance*

A quick introduction to diffusion models



DALLE

TEXT PROMPT

an armchair in the shape of an avocado. an armchair imitating an avocado.

AI-GENERATED
IMAGES



<https://openai.com/blog/dall-e/>

DALLE

TEXT PROMPT

a soap dispenser in the style of a pikachu. a soap dispenser imitating a pikachu.

AI-GENERATED IMAGES



<https://openai.com/blog/dall-e/>

DALLE 2

TEXT DESCRIPTION

An astronaut **Teddy bears** A bowl of

soup

mixing sparkling chemicals as mad
scientists shopping for groceries working
on new AI research

as kids' crayon art **on the moon in the**
1980s underwater with **1990s technology**



DALL-E 2



Imagen by Google AI



A photo of a Corgi dog riding a bike in Times Square. It is wearing sunglasses and a beach hat.

<https://Imagen.research.google/>

Make-A-Video by Meta



<https://makeavideo.studio/>

Dreambooth



Input images



in the Acropolis



swimming



sleeping



getting a haircut

Dreamfusion



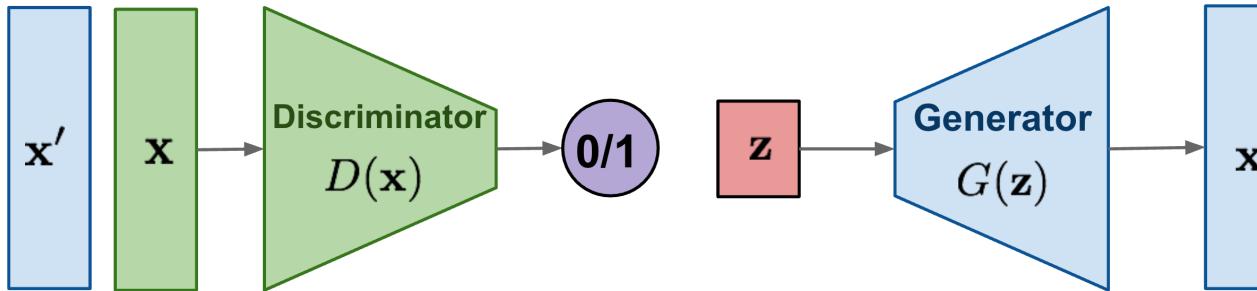
More diffusion models

- Resources
 - Introductory Posts
 - Introductory Papers
 - Introductory Videos
 - Introductory Lectures
 - Tutorial and Jupyter Notebook
- Papers
 - Survey
 - Vision
 - Generation
 - Classification
 - Segmentation
 - Image Translation
 - Inverse Problems
 - Medical Imaging
 - Multi-modal Learning
 - 3D Vision
 - Adversarial Attack
 - Miscellany
- Audio
 - Generation
 - Conversion
 - Enhancement
 - Separation
 - Text-to-Speech
 - Miscellany
- Natural Language
- Tabular and Time Series
 - Generation
 - Forecasting
 - Imputation
 - Miscellany
- Graph
 - Generation
 - Molecular and Material Generation
- Reinforcement Learning
- Theory
- Applications

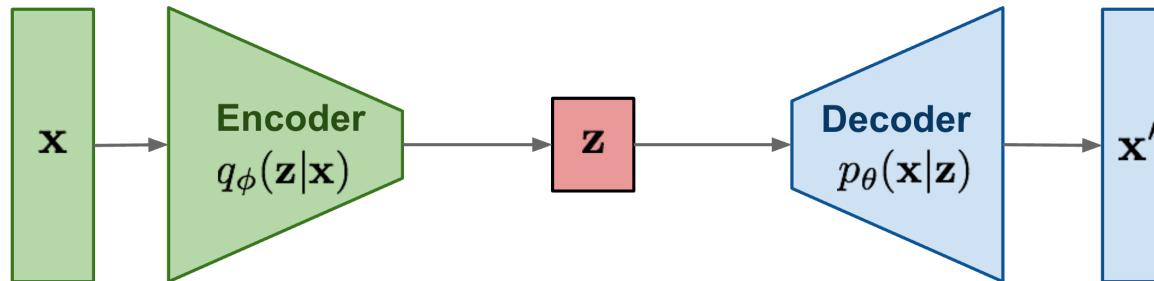
<https://github.com/heejkoo/Awesome-Diffusion-Models>

Generative models

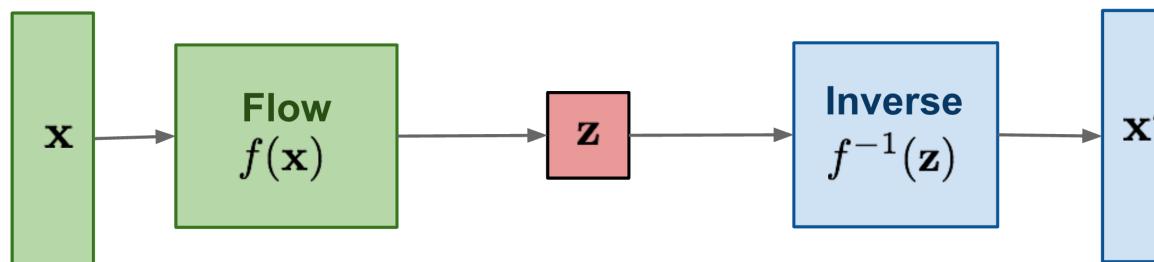
GAN: Adversarial training



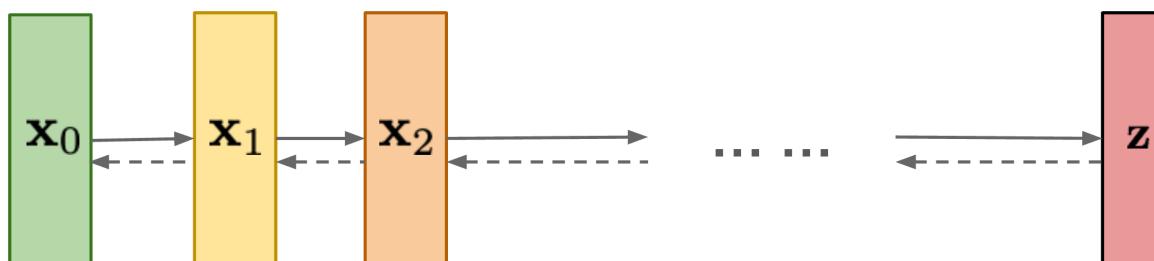
VAE: maximize variational lower bound



Flow-based models:
Invertible transform of distributions



Diffusion models:
Gradually add Gaussian noise and then reverse



Deep Unsupervised Learning using Nonequilibrium Thermodynamics

Jascha Sohl-Dickstein

Stanford University

JASCHA@STANFORD.EDU

Eric A. Weiss

University of California, Berkeley

EAWEISS@BERKELEY.EDU

Niru Maheswaranathan

Stanford University

NIRUM@STANFORD.EDU

Surya Ganguli

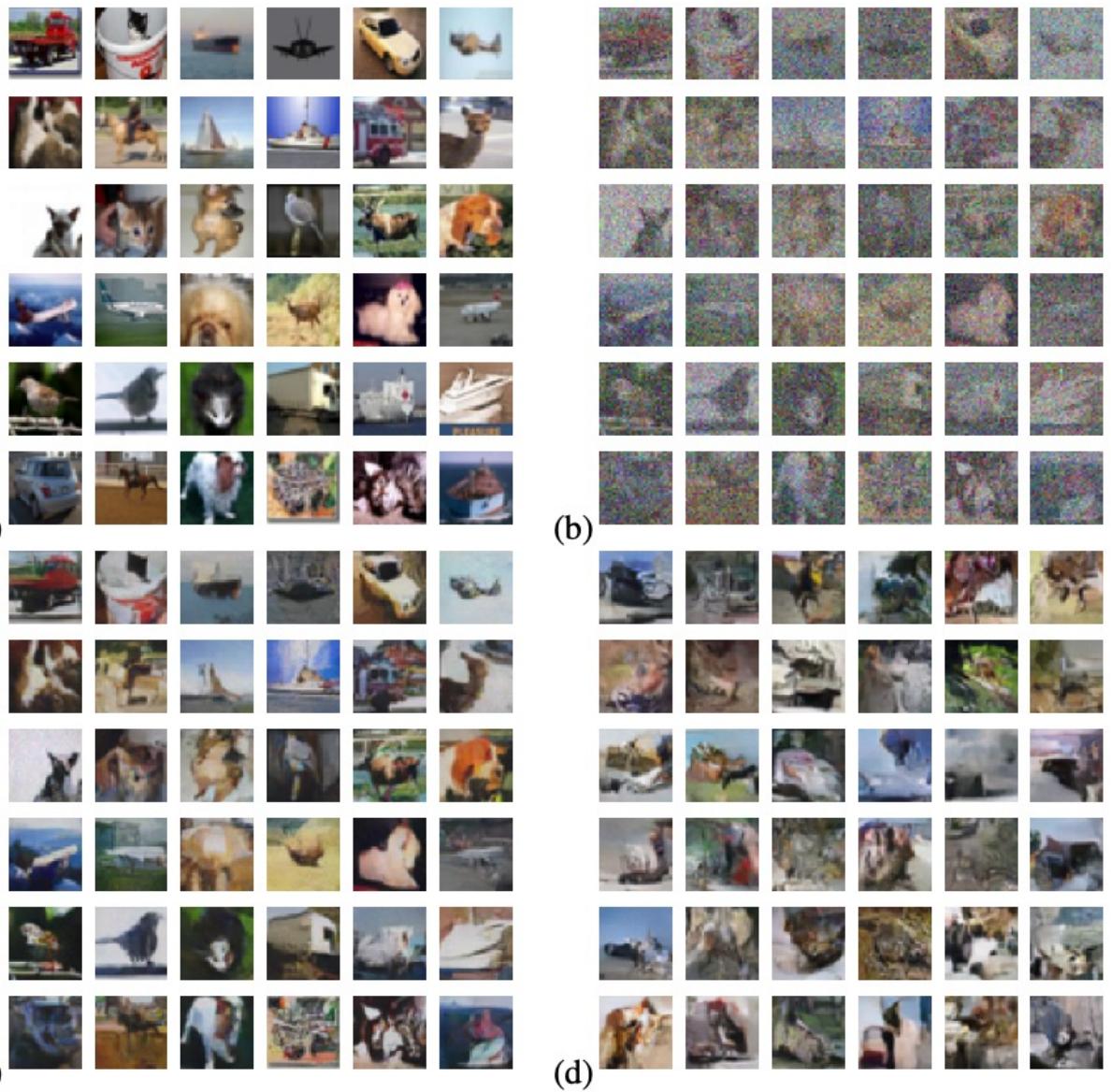
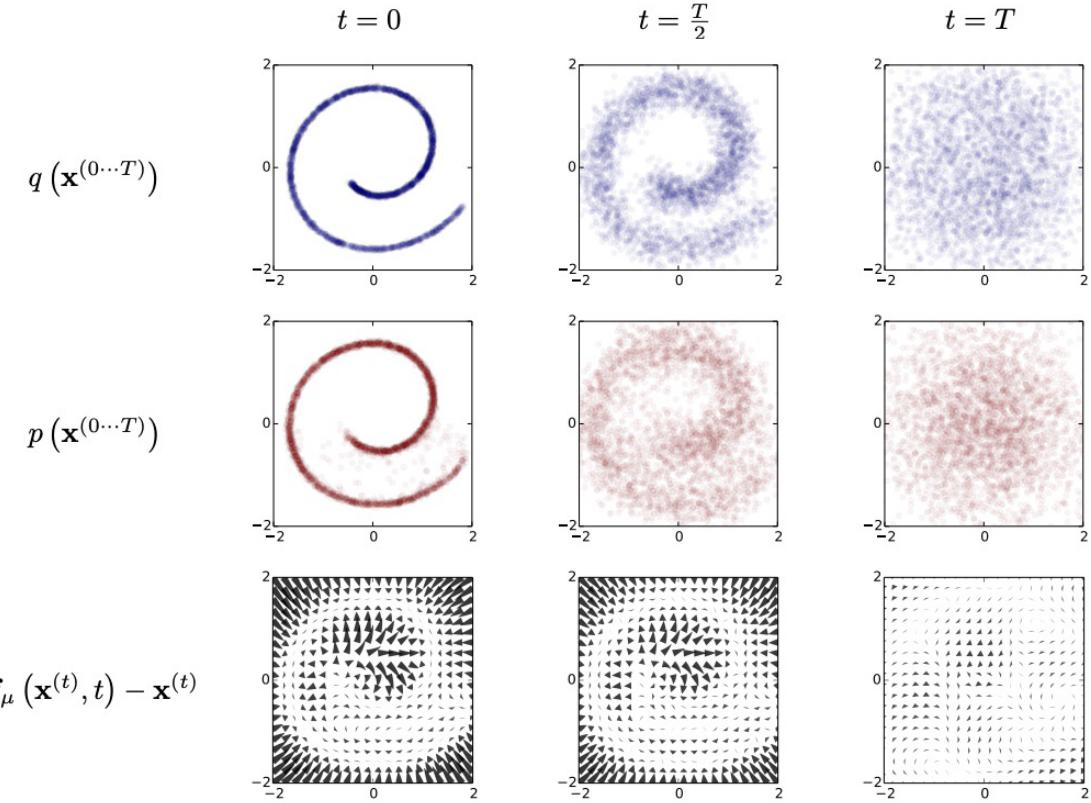
Stanford University

SGANGULI@STANFORD.EDU

Abstract

A central problem in machine learning involves

these models are unable to aptly describe structure in rich datasets. On the other hand, models that are *flexible* can be molded to fit structure in arbitrary data. For example, we



Denoising Diffusion Probabilistic Models

Denoising Diffusion Probabilistic Models

Jonathan Ho

UC Berkeley

jonathanho@berkeley.edu

Ajay Jain

UC Berkeley

ajayj@berkeley.edu

Pieter Abbeel

UC Berkeley

pabbeel@cs.berkeley.edu

Abstract

We present high quality image synthesis results using diffusion probabilistic models, a class of latent variable models inspired by considerations from nonequilibrium thermodynamics. Our best results are obtained by training on a weighted variational bound designed according to a novel connection between diffusion probabilistic models and denoising score matching with Langevin dynamics, and our models naturally admit a progressive lossy decompression scheme that can be interpreted as a generalization of autoregressive decoding. On the unconditional CIFAR10 dataset, we obtain an Inception score of 9.46 and a state-of-the-art FID score of 3.17. On 256x256 LSUN, we obtain sample quality similar to ProgressiveGAN. Our implementation is available at <https://github.com/hojonathanho/diffusion>.

1 Introduction

Deep generative models of all kinds have recently exhibited high quality samples in a wide variety of data modalities. Generative adversarial networks (GANs), autoregressive models, flows, and variational autoencoders (VAEs) have synthesized striking image and audio samples [14, 27, 3, 58, 38, 25, 10, 32, 44, 57, 26, 33, 45], and there have been remarkable advances in energy-based modeling and score matching that have produced images comparable to those of GANs [11, 55].

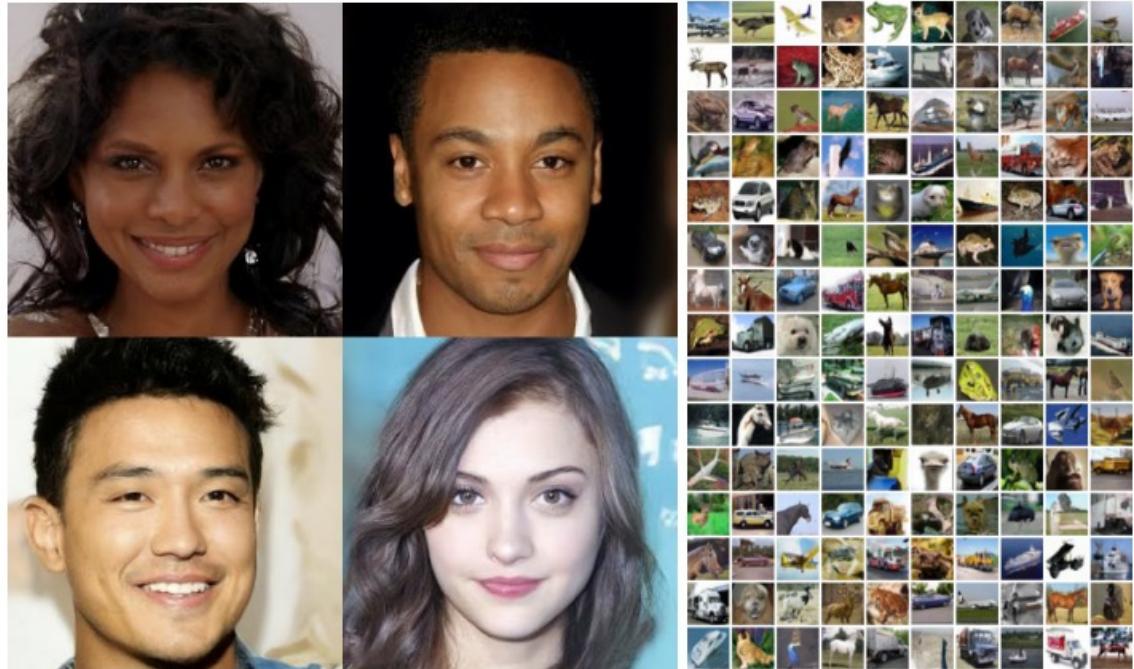
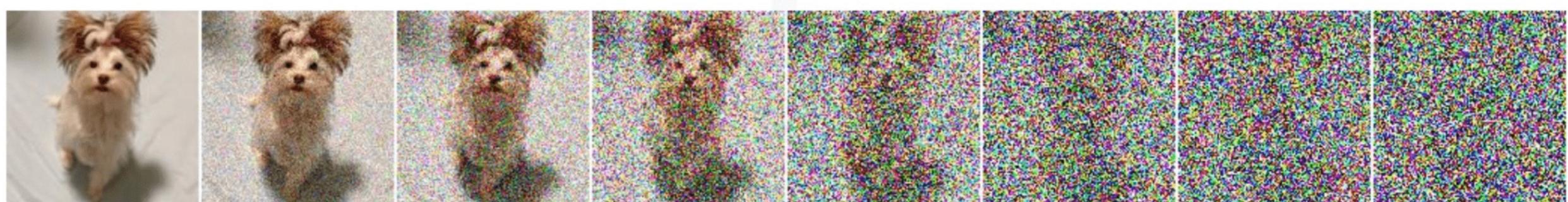


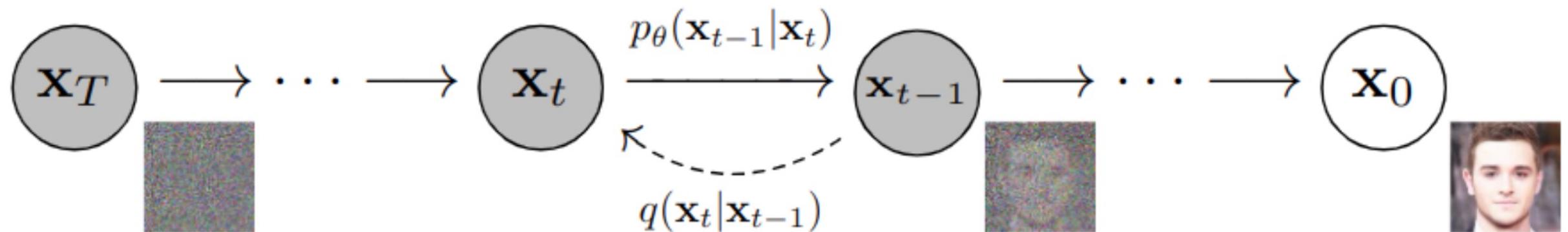
Figure 1: Generated samples on CelebA-HQ 256 × 256 (left) and unconditional CIFAR10 (right)

Denoising Diffusion Probabilistic Models

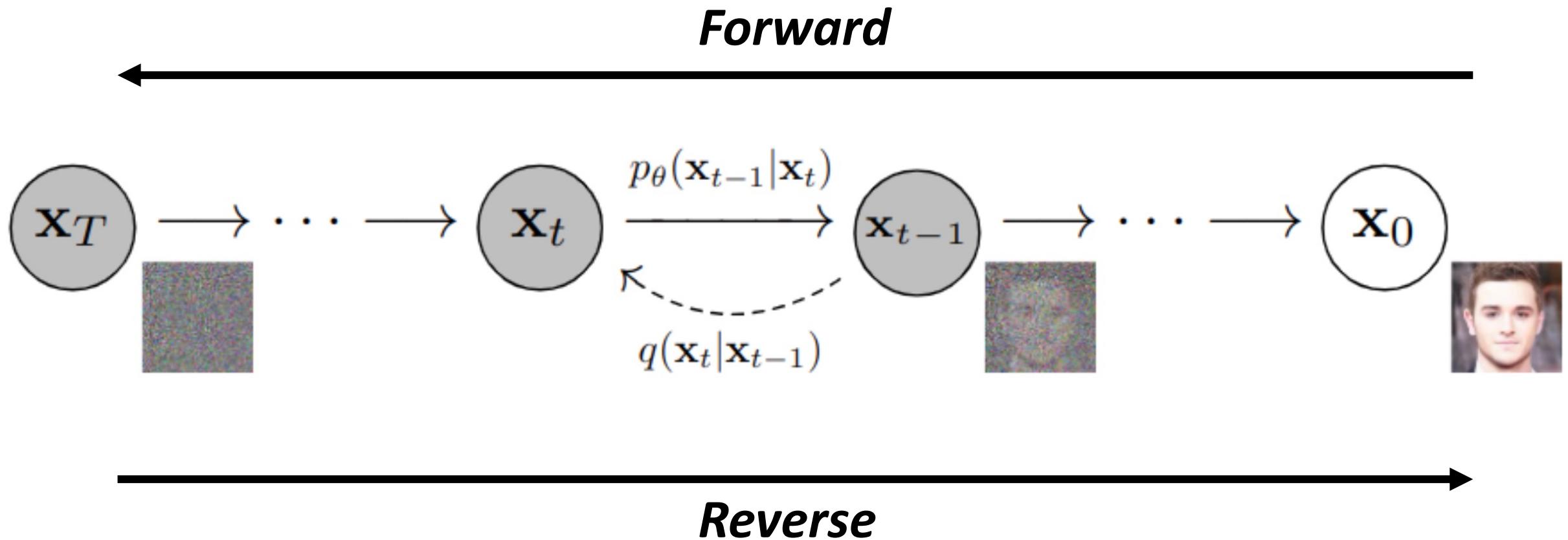


<https://huggingface.co/blog/annotated-diffusion>

Algorithm



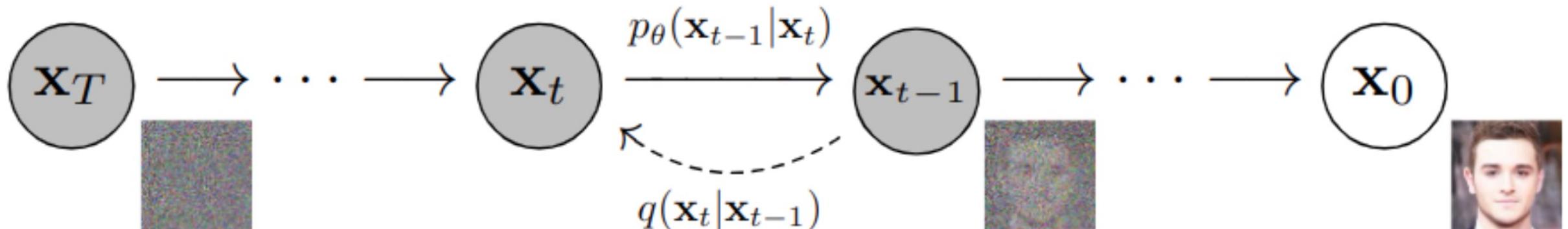
Forward and reverse process



Forward and reverse process

Forward

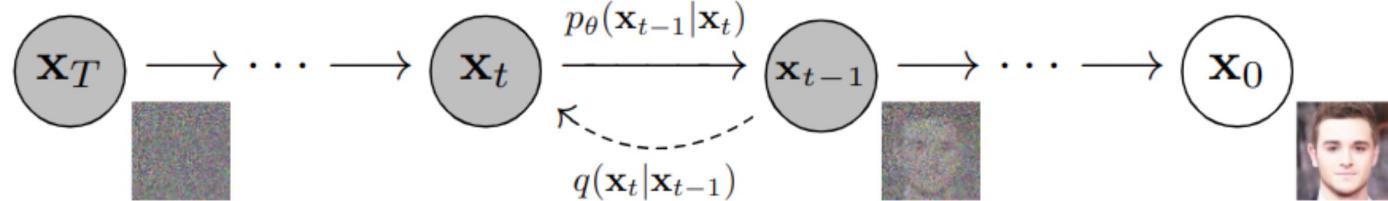
Fixed process that gradually adds Gaussian noise to an image, until you end up with pure noise.



Reverse

Learned reverse denoising diffusion process where a neural network is trained to gradually denoise an image starting from pure noise, until you end up with an actual image

Forward process



Forward diffusion process

Given a schedule $\beta_1 < \beta_2 < \dots < \beta_T$,

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t I) \quad (1)$$

$$q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1})$$

We define $\alpha_t = 1 - \beta_t$ and $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$, then we have

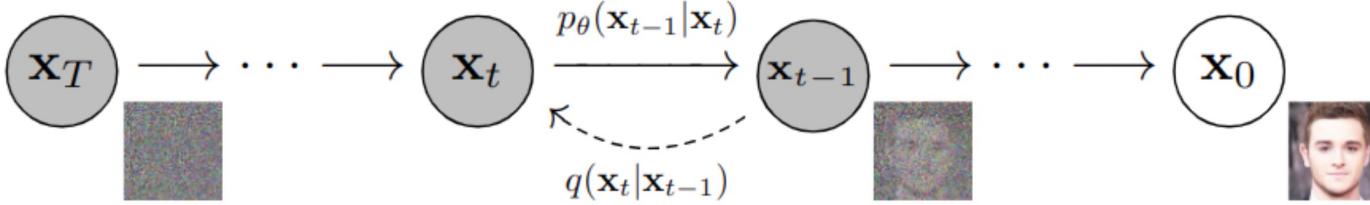
$$\begin{aligned} x_t &= \sqrt{\alpha_t}x_{t-1} + \sqrt{1 - \alpha_t}\epsilon_{t-1}, \text{ with } \epsilon_{t-1} \sim \mathcal{N}(0, I) \\ &= \sqrt{\alpha_t\alpha_{t-1}}x_{t-2} + \sqrt{\alpha_t(1 - \alpha_{t-1})}\epsilon_{t-2} + \sqrt{1 - \alpha_t}\epsilon_{t-1} \\ &= \sqrt{\alpha_t\alpha_{t-1}}x_{t-2} + \sqrt{1 - \alpha_t\alpha_{t-1}}\tilde{\epsilon}_t \end{aligned} \quad (2)$$

- Diffusion schedule (beta values)
- DDPM T=1000

Hence, we have

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon \quad (3)$$

Approximate the reverse diffusion



Training: to approximate **the reversed diffusion** $q(x_{t-1}|x_t)$ by a neural network given by $p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \beta_t I)$ and $p(x_T) \sim \mathcal{N}(0, I)$, we maximize the usual Variational bound:

$$\mathbb{E}_{q(x_0)} \ln p_\theta(x_0) \geq L_T + \sum_{t=2}^T L_{t-1} + L_0 \text{ with, } L_{t-1} = \mathbb{E}_q \left[\frac{1}{2\sigma_t^2} \|\mu_\theta(x_t, t) - \mu(x_t, \epsilon, t)\|^2 \right] \quad (17)$$

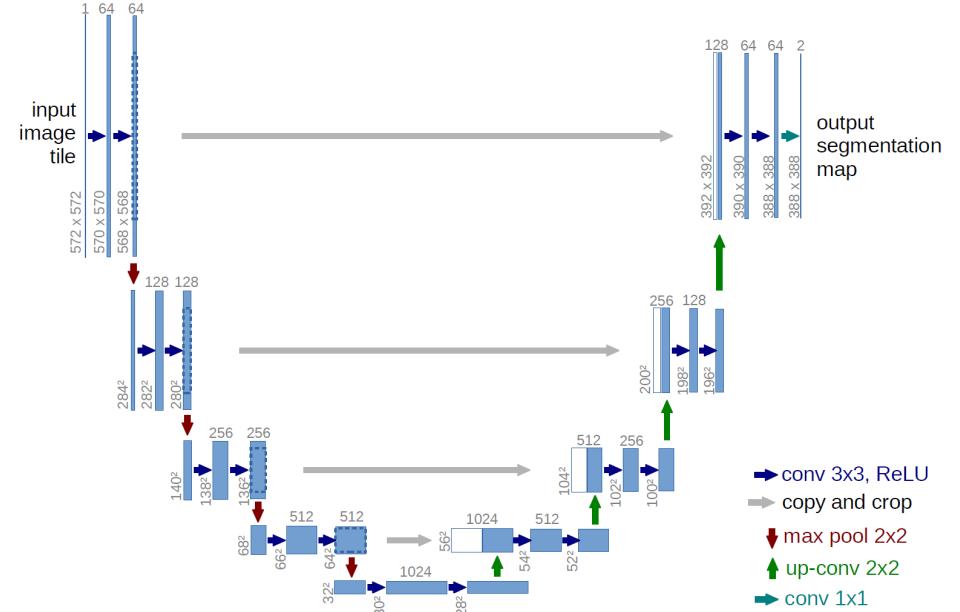
With the change of variable:

$$\mu_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t, t) \right), \quad (18)$$

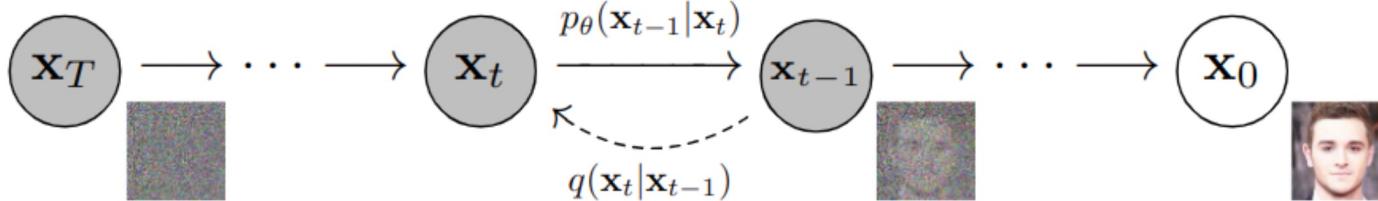
ignoring the prefactor and sampling τ instead of summing over all t , the loss is finally:

$$\ell(\theta) = \mathbb{E}_\tau \mathbb{E}_\epsilon [\|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_\tau} x_0 + \sqrt{1 - \bar{\alpha}_\tau} \epsilon, \tau)\|^2] \quad (19)$$

- Train denoising U-Nets
- MSE Loss
- Predict noise at each step



Training and Sampling



Algorithm 1 Training

```

1: repeat
2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 
3:    $t \sim \text{Uniform}(\{1, \dots, T\})$ 
4:    $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:   Take gradient descent step on
         $\nabla_{\theta} \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t)\|^2$ 
6: until converged

```

- we take a random sample \mathbf{x}_0 from the real unknown and possibly complex data distribution $q(\mathbf{x}_0)$
- we sample a noise level t uniformly between 1 and T (i.e., a random time step)
- we sample some noise from a Gaussian distribution and corrupt the input by this noise at level t (using the nice property defined above)
- the neural network is trained to predict this noise based on the corrupted image \mathbf{x}_t (i.e. noise applied on \mathbf{x}_0 based on known schedule β_t)

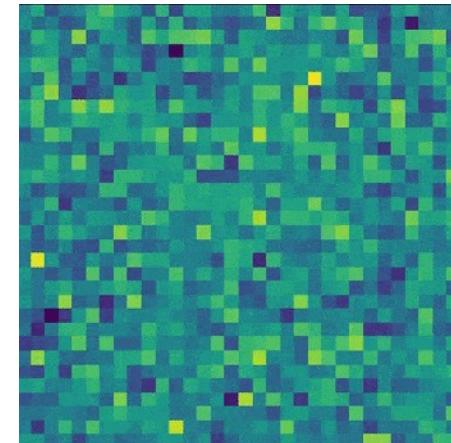
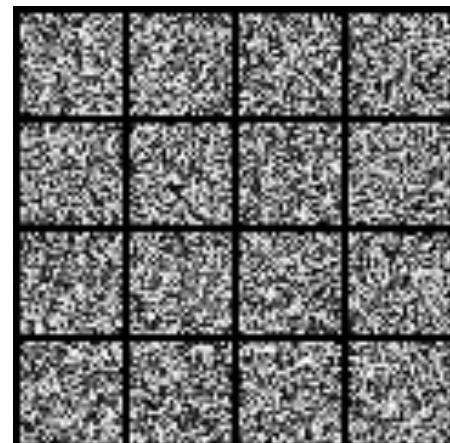
In reality, all of this is done on batches of data, as one uses stochastic gradient descent to optimize neural networks.

Algorithm 2 Sampling

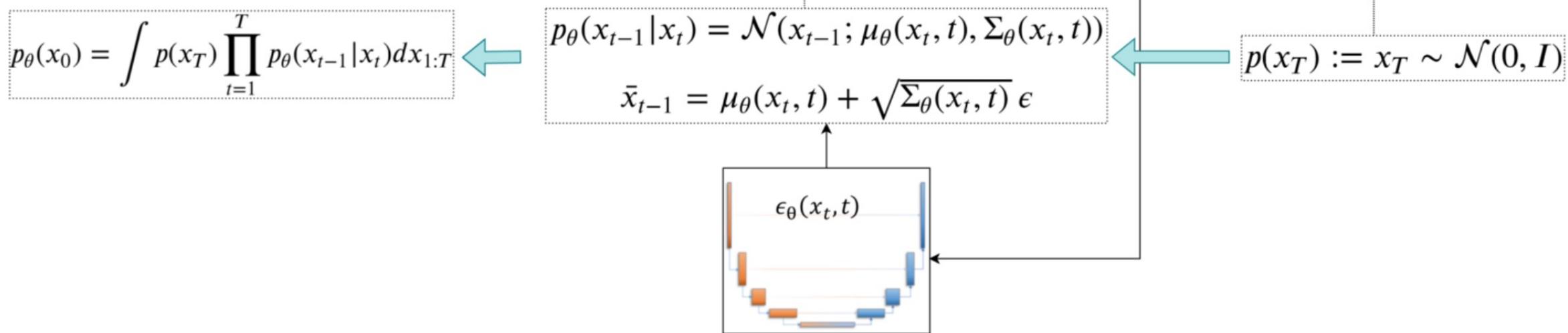
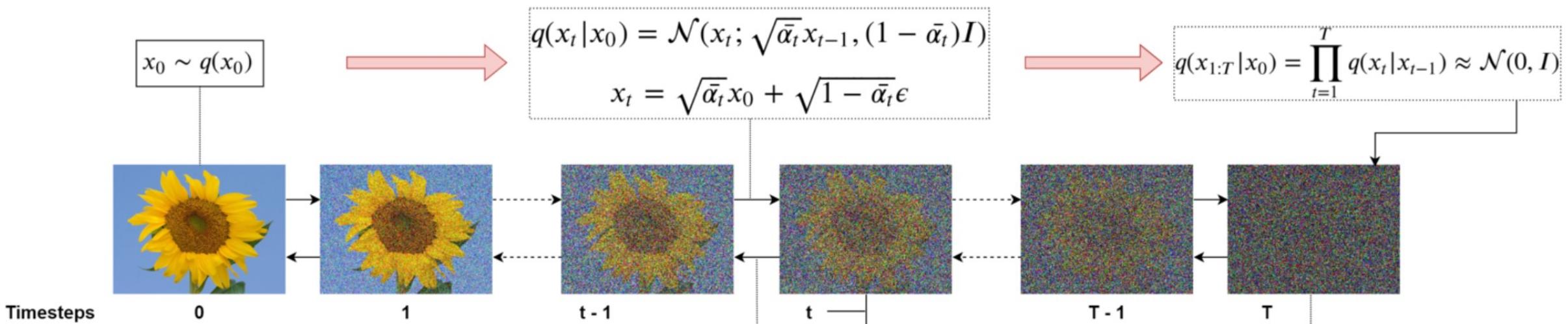
```

1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 

```

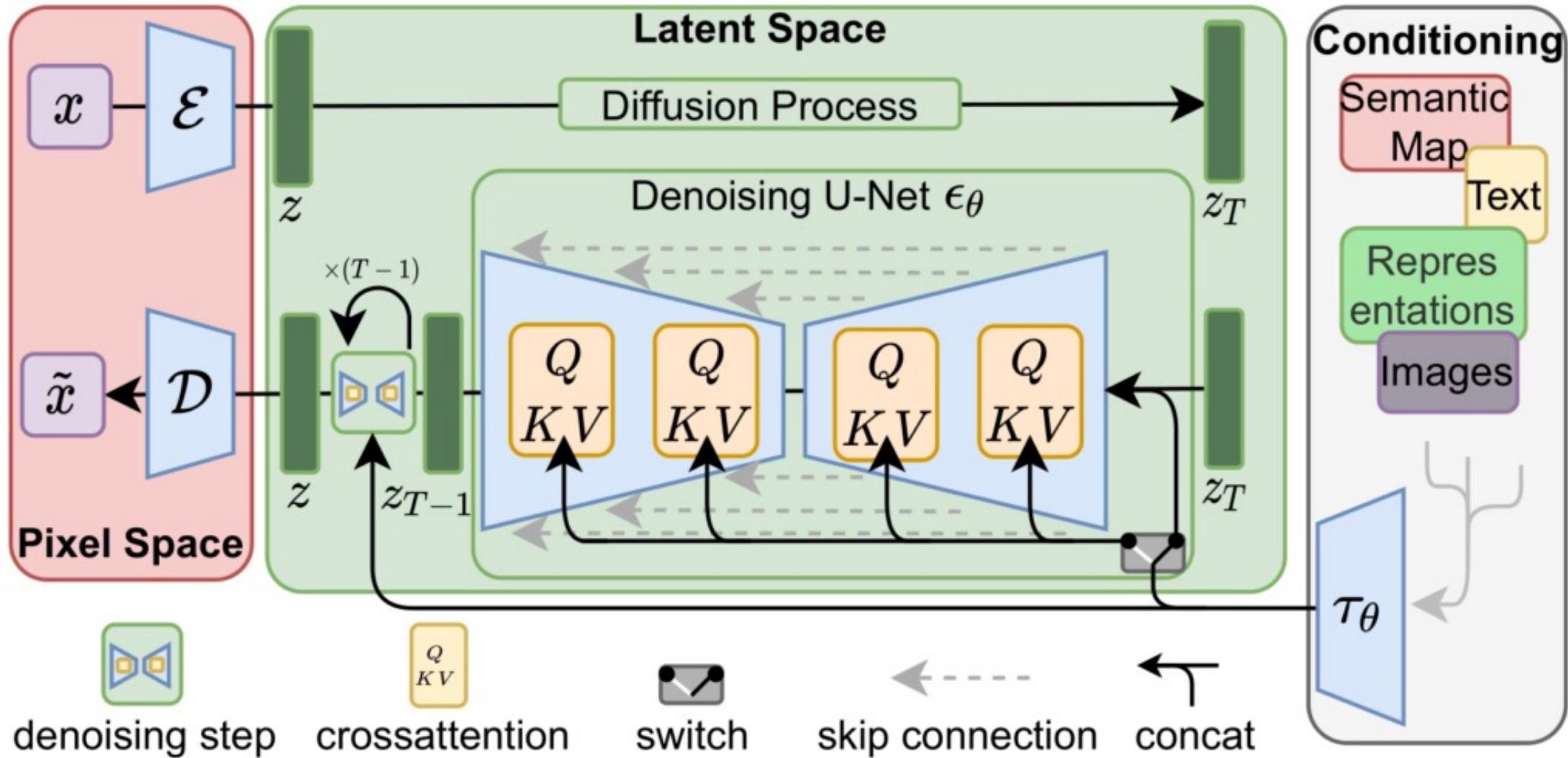


Forward Diffusion Process



Reverse Diffusion Process

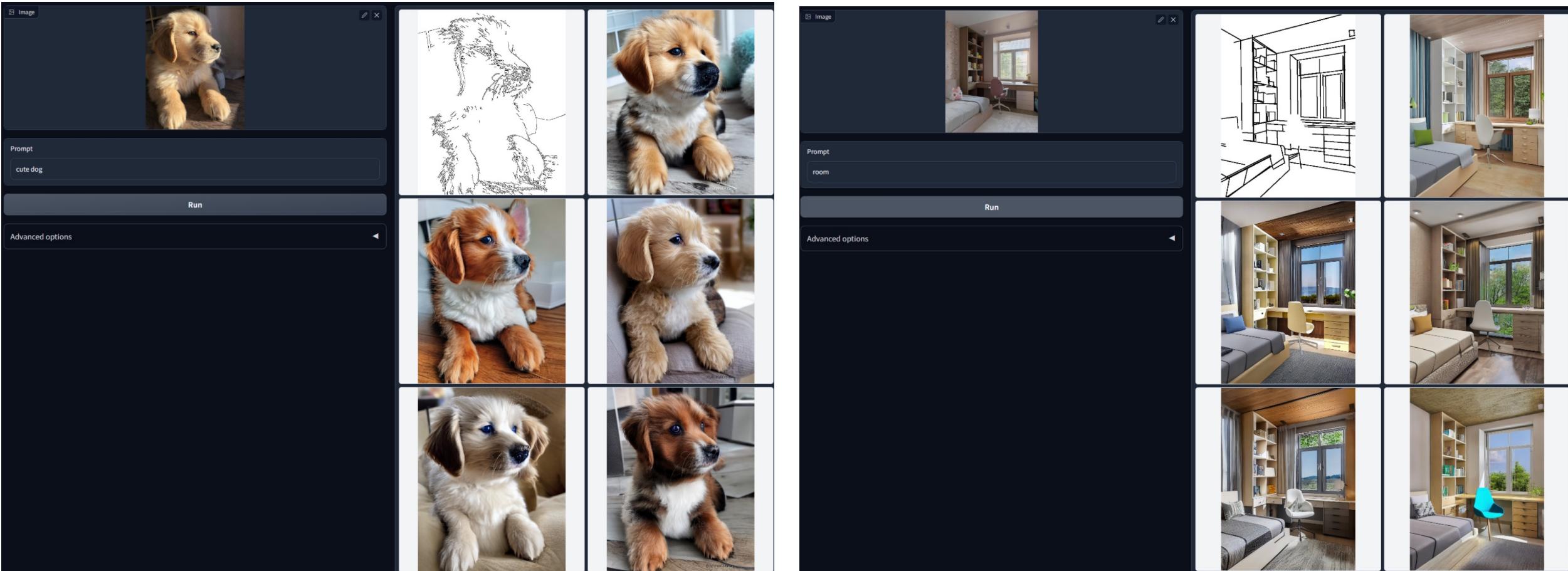
Latent diffusion models (pixel vs feature space)



Latent diffusion models (pixel vs feature space)



ControlNet



ControlNet

