# Cornerstones of the Text-to-Pixels Journey

Srikumar Ramalingam

Google Research, NYC
Adjunct Faculty, University of Utah

# Tutorial Speakers



Shobhita Sundaram
MIT

Sadeep Jayasumana
Google Research

Varun Jampani
Stability AI

Dilip Krishnan
Google DeepMind

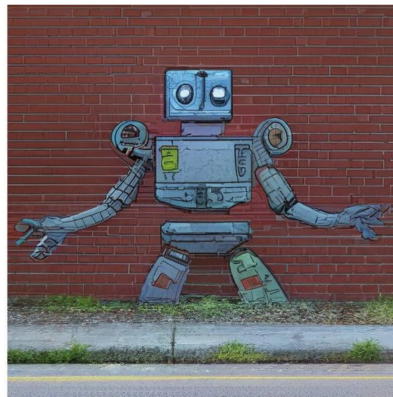Srikumar Ramalingam
Google Research

# Overview

| Time | Speaker | Title |
|------|---------|-------|
| 9:10 -   9.50 | Srikumar Ramalingam | *Cornerstones of the Text-to-Pixels Journey* |
| 9.50 - 10.30 | Shobhita Sundaram | *Image Evaluation Methods* |
| 10.30 - 11.00 | Break | |
| 11.00 - 11.30 | Varun Jampani | *Thinking Slow and Fast: Recent Trends in 3D Generative Models* |
| 11:00 - 12:00 | Dilip Krishnan | *Parallel Decoding and Image Generation* |
| 12:00 - 12:30 | Sadeep Jayasumana | *Structured Prediction Algorithms for Fast Image Generation* |

# Text-to-Image Generation



A robot cooking in the kitchen.

A robot painted as graffiti on a brick wall. a sidewalk is in front of wall, grass is growing out of cracks in the concrete.

A raccoon wearing formal clothes, wearing a tophat. The raccoon is holding a garbage bag.

A hyper-realistic concept art of an alien pyramid landscape, inspired by ArtStation artists.

Jayasumana et al. Markovgen 2023.

# Text to video Generation

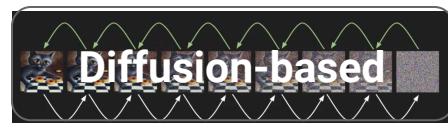# Text-to-3D Generation



https://dreamfusion3d.github.io/

Ben Poole, Ajay Jain, Ben Mildenhall, Jon Barron

# Text-to-Image backbone

Three-quarters front view of a blue 1977 Corvette coming around a curve in a mountain road and looking over a green valley on a cloudy day.

Diffusion-based

Auto-regressive

# Transformers and Diffusion models

**Elon Musk** ✔ ✖
@elonmusk

Subscribe   ...

Who should be President in 2032?

| | |
|---|---|
| **Transformers** | **77.4%** |
| Diffusion | 22.6% |

1,178,197 votes · Final results

# t2i models are centerpieces of many generative models

## Text-to-3D



Use text-to-image and NeRF models as building blocks to generate 3D from text.

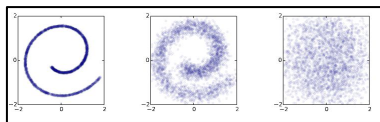https://dreamfusion3d.github.io/

## Text-to-Video



Generate distinct keyframes using text-to-image model, followed by temporal and spatial super-resolution models.

Bar-Tal et al. Lumiere, 2024
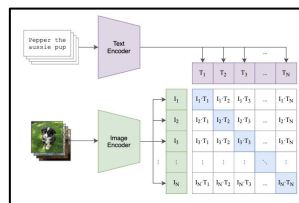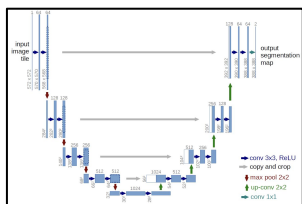
# Pieces of the Text-to-Image Puzzle

**2015**

Diffusion
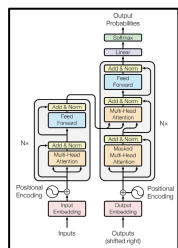


**2020**

CLIP



**2021**

DALL-E





UNet

**2015**



Transformers

**2017**



VQGAN

**2021**

# Pieces of the Text-to-Image Puzzle
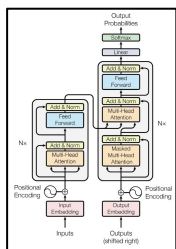
# Image features with similar objects are close



Near duplicates

Sussex spaniel

(red wolf, maned wolf, Canis rufus, Canis niger)   (timber wolf, grey wolf, gray wolf, Canis lupus)

Features corresponding to images containing same semantic objects are close to each other in the embedding space.

# Image-Text Co-Embedding Spaces



Bipartite mapping between image and text embeddings

"artificial intelligence"

"self-driving car"

"ML software"

Image-Text Co-Embedding Space

# Single tower vs. two tower models



Single-tower classification with ResNets or ViTs trained on a chosen set of labels such as in ImageNet.
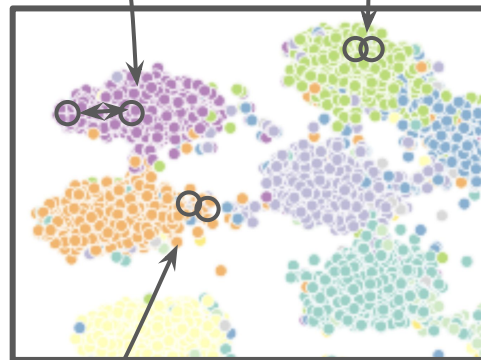
Learning two tower models allows us to use zero-shot classification methods on different classes.

# CLIP/ALIGN



$$L_{i2t} = -\frac{1}{N} \sum_{i}^{N} \log \frac{\exp(x_i^\top y_i / \sigma)}{\sum_{j=1}^{N} \exp(x_i^\top y_j / \sigma)}$$

$$L_{t2i} = -\frac{1}{N} \sum_{i}^{N} \log \frac{\exp(y_i^\top x_i / \sigma)}{\sum_{j=1}^{N} \exp(y_i^\top x_j / \sigma)}$$

$x_i, y_i$  Image and Text normalized embeddings

# Text-Image Coembedding References

[CLIP] [Learning Transferable Visual Models From Natural Language Supervision](#), 2021.

[ALIGN] [Scaling Up Visual and Vision-Language Representation Learning With Noisy Text Supervision](#), 2021.

# Pieces of the Text-to-Image Puzzle



**2015**

Diffusion

**2020**

CLIP

**2021**

DALL-E

UNet

**2015**

Transformers

**2017**

VQGAN

**2021**

# Background: Visual Words

Individual parts of an object reveal a lot of information.

# Background: Visual words



Word #2

Descriptor's
feature space

- Quantize via clustering, let cluster centers be the prototype "words"

- Determine which word to assign to each new image region by finding the closest cluster center.

# Background: Visual words

# Image Tokenization



256 x 256 x 3

Tokenizer

| 9 | 28 | 15 | 78 |
|----|----|----|----|
| 19 | 36 | 27 | 32 |
| 8 | 56 | 68 | 71 |
| 96 | 85 | 49 | 82 |

16 x 16

Detokenizer

256 x 256 x 3

# Key Idea in VQGAN

- Use for CNNs for learning local features and transformers for long range interactions
  - CNNs are used to learn a codebook of context-rich visual parts.
  - Transformers are used to model the long range interactions among the individual visual parts.
- Efficient image generation backbone that allows conditional inputs (similar to ControlNet).
- Default choice in Latent diffusion, MUSE, Parti, Paella, etc.

**Taming transformers for high-resolution image synthesis,**
Patrick Esser*, Robin Rombach*, Björn Ommer

# Overview of VQGAN



Two stage training:
- Learn the encoder, decoder, and codebook.
- Learn the transformer to synthesize images with conditional inputs.

# Codebook

$$x \in \mathbb{R}^{H \times W \times 3}$$

Input Image



| 3861 | 2201 | 743 | 408 |
| 221 | 200 | 4999 | 6021 |
| 421 | 8001 | 7871 | 1213 |
| 7495 | 4259 | 121 | 910 |

Token Image

$$\mathcal{Z} = \{z_k\}_{k=1}^{K} \subset \mathbb{R}^{n_z}$$

Discrete codebook consisting of K vectors



Discrete Token Labels

8192

4021

1

$$z_{\mathbf{q}} \in \mathbb{R}^{h \times w \times n_z}$$

Image represented with codebook entries

# Codebook



$$\hat{z} = E(x) \in \mathbb{R}^{h \times w \times n_z} \qquad \hat{x} = G(z_{\mathbf{q}})$$

$$z_{\mathbf{q}} = \mathbf{q}(\hat{z}) := \left( \arg\min_{z_k \in \mathcal{Z}} \|\hat{z}_{ij} - z_k\| \right) \in \mathbb{R}^{h \times w \times n_z}$$

# Learning the codebook

$$\hat{x} = G(z_{\mathbf{q}}) = G\left(\mathbf{q}(E(x))\right)$$

$$\mathcal{L}_{\text{VQ}}(E, G, \mathcal{Z}) = \|x - \hat{x}\|^2 + \|\text{sg}[E(x)] - z_{\mathbf{q}}\|_2^2$$
$$+ \|\text{sg}[z_{\mathbf{q}}] - E(x)\|_2^2$$

Move the codebook vectors closer to the frozen encoder vectors, and vice versa.

- Reconstruction loss optimizes the encoder and decoder.
- L2 loss to move the encoder outputs towards the codebook entries and another L2 loss to move codebook entries towards the encoder outputs.

[Aaron van den Oord et al. Neural Discrete Representation Learning, 2017]

# Learning a perceptually rich codebook

GAN Loss:

Discriminator wants to maximize this, while the generator wants to minimize this.

$$\mathcal{L}_{\mathrm{GAN}}(\{E, G, \mathcal{Z}\}, D) = \left[\log D(x) + \log(1 - D(\hat{x}))\right]$$

$$\mathcal{Q}^* = \arg\min_{E,G,\mathcal{Z}} \max_{D} \mathbb{E}_{x \sim p(x)} \Big[ \mathcal{L}_{\mathrm{VQ}}(E, G, \mathcal{Z})$$

$$+ \lambda \mathcal{L}_{\mathrm{GAN}}(\{E, G, \mathcal{Z}\}, D) \Big]$$

- Learn the encoder, decoder, and codebook with a perceptual and GAN loss.

# Feature Codebook References

- [VQGAN]: [Taming Transformers for High-Resolution Image Synthesis](#), 2020.
- [VQVAE]: [Neural Discrete Representation Learning](#), 2018.
- [Video Google: A Text Retrieval Approach to Object Matching in Videos](#), 2003.

# Pieces of the Text-to-Image Puzzle



2015
Diffusion

2020
CLIP

2021
DALL-E

UNet
2015

Transformers
2017

VQGAN
2021

# Markov Random Fields (MRFs)



factor nodes

variable nodes



Goal: find most probable interpretation of scene

Minimize an energy function:

$$E(\mathbf{x}) = \mathrm{unary\_cost} + \mathrm{pairwise\_cost}$$

- Solve using using graph cuts or BP

# Model Hierarchy (MRFs -> CNNs -> Transformers)



factor nodes

variable nodes



Extract features from patches hierarchically



Long-range interaction!

MRFs with 4 or 8-neighborhood were solved efficiently using graph cuts and belief propagation.

CNNs are very good at extracting local features!

Transformers allow long range interactions!

Graphcuts
**1999**

AlexNet
**2012**

Transformers
**2017**

# Vision Transformer



**Vision Transformer (ViT)**

**Transformer Encoder**

# Conditioned Synthesis using Transformers

With the encoder, decoder, and codebook, we can treat the image synthesis problem as sequence prediction problem.

**[..,743, 408, 221, 200, ….]**



| 3861 | 2201 | 743 | 408 |
| 221 | 200 | 4999 | 6021 |
| 421 | 8001 | 7871 | 1213 |
| 7495 | 4259 | 121 | 910 |

Token Image

- Based on some ordering, the token prediction can be achieved auto-regressively by feeding the previous tokens.
- To provide conditional inputs, we can learn another codebook if it has spatial extent to generate token indices for conditions.

# Different ordering of tokens for image synthesis



- The ordering is trivial for language tasks, whereas there is no easy way to fix the ordering for images.

# Class conditioned Image Synthesis



256x256 images conditioned on ImageNet

# Conditioned Image Synthesis



Depth -> Image

Low res. -> High res. (Superresolution)

Semantic -> Image

Edge -> Image

# Efficient Text-to-Image Generation using Muse

# MarkovGen: MRFs to speedup Muse

$$E(\mathbf{x}) = \text{unary\_cost} + \text{pairwise\_cost}$$



Label compatibility
$c(x_i, x_j)$

8192

4021

1

1902

Spatial
compatibility
$s(i, j)$

*Detokenization* ↑

| 4021 | 2351 | 743 | 408 |
| 221 | 1902 | 4999 | 600 |
| 420 | 8001 | 6421 | 1213 |
| 7495 | 2001 | 121 | 900 |

**Imperfect token Image**

*MRF Inference* →

↑ *Detokenization*

| 402 | 2351 | 743 | 408 |
| 92 | 1902 | 204 | 600 |
| 420 | 8001 | 6421 | 789 |
| 7495 | 2001 | 121 | 800 |

**Fixed token image**

# MRF: Model Formulation



$$E(\mathbf{x}) = \text{unary\_cost} + \text{pairwise\_cost}$$

Unary Cost
- $\text{cost}(X_i = l) = ?$
- You pay a penalty if your label doesn't agree with the classifier.

Pairwise cost
- $\text{cost}(X_i = l', X_j = l'') = ?$
- You pay a penalty if you assign *"incompatible"* labels to two *"neighboring"* tokens.

$$\text{cost}(X_i = l) = -\text{logit}_i(l)$$

$$\text{cost}(X_i = l', X_j = l'') = -c(l', l'')s(i, j)$$

# Speedup over Muse without quality loss.



Full Muse: All steps

Early Exit Muse: Fewer steps
1.5x faster

MarkovGen: Fewer steps + MRF
1.5x faster

A robot cooking in the kitchen

A robot painted as graffiti on a brick wall. a sidewalk is in front of the wall, and grass is growing out of cracks in the concrete.

| Model | Time (ms) |
|---|---|
| Muse base (single step) | 10.40 |
| Muse super-resolution (single step) | 24.00 |
| MRF inference on base | 0.29 |
| MRF inference on super-resolution | 0.29 |
| T5-XXL inference | 0.30 |
| Detokenizer | 0.15 |
| Muse | 442.05 |
| MarkovGen (ours) | 281.03 |

# MRF and Transformers References

- Masked generative image transformer. In: CVPR (2022)
- Muse:Text-to-image generation via masked generative transformers. ICML (2023)
- Markovgen: Structured prediction for efficient text-to-image generation (2023)
- Hierarchical text-conditional image generation with clip latents. preprint (2022)
- Photorealistic text-to-image diffusion models with deep language understanding. preprint (2022),
- Scaling autoregressive models for content-rich text-to-image generation. In: ICML (2022)

# Pieces of the Text-to-Image Puzzle



**2015**

## Diffusion

**2020**

# CLIP

**2021**

# DALL-E

UNet
**2015**

Transformers
**2017**

VQGAN
**2021**

# Basic idea -> Diffusion Model



Forward diffusion

Reverse diffusion

[Deep unsupervised learning of nonlinear thermodynamics, Sohl-Dickstein et al. 2015]

# Diffusion Models



"a hedgehog using a calculator"

"a corgi wearing a red bowtie and a purple party hat"

"robots meditating in a vipassana retreat"

"a fall landscape with a small cottage next to a lake"

"a surrealist dream-like oil painting by salvador dalí of a cat playing checkers"

"a professional photo of a sunset behind the grand canyon"

"a high-quality oil painting of a psychedelic hamster dragon"

"an illustration of albert einstein wearing a superhero costume"

"a boat in the canals of venice"

"a painting of a fox in the style of starry night"

"a red cube on top of a blue cube"

"a stained glass window of a panda eating bamboo"

"a crayon drawing of a space elevator"

"a futuristic city in synthwave style"

"a pixel art corgi pizza"

"a fog rolling into new york"

[Nichol et al. GLIDE 2021]

# Background: Diffusion models

"*Systematically and slowly destroy structure in a data distribution through an iterative forward diffusion process.*



*We then learn a reverse diffusion process that restores structure in data, yielding a highly flexible and tractable generative model of the data.*"

**[Deep unsupervised learning of nonlinear thermodynamics, Sohl-Dickstein et al. 2015]**

# Background: Diffusion models



- <u>While training</u> we start with clean images from the dataset, add noise and try to predict the added noise.
- <u>While sampling</u>, we start with noise and iteratively denoise the image to generate an image.

# Diffusion model

Mean squared error loss: $||\epsilon - pred||^2$



noise $\epsilon$

image $x_0$

Noised image $x_t$

Diffusion model

$pred$

[Nichol et al. GLIDE 2021]

# Training Diffusion models



$x_0 \sim q(x_0)$  $x_1$  $x_2$  $x_T$

Markov chain of latent variables by progressively adding Gaussian noise.

Sample an
image from
the data
distribution

# Training Diffusion models



$x_0 \sim q(x_0)$     $x_1$        $x_2$                                $x_T$

Sample an image from the data distribution

Markov chain of latent variables by progressively adding Gaussian noise.

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) := \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

# Training diffusion models

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) := \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I})$$



- We are somewhat shrinking the mean and moving it towards the 0.

- If the total noise added is large enough, and if each step adds small enough noise, then $x_T$ can be approximated by $\mathcal{N}(0, \mathcal{I})$.

# Training Diffusion Models



$x_0 \sim q(x_0)$    $x_1$        $x_2$

Sample an image from the data distribution

Markov chain of latent variables by progressively adding Gaussian noise.

$$\alpha_t := 1 - \beta_t \qquad \bar{\alpha}_t := \prod_{s=1}^{t} \alpha_s$$

$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})$$

$$(1 - \alpha_t) < 1, \sqrt{\alpha} < 1$$

# Loss Function



$$L_{simple} = E_{t\sim[1,T],x_0\sim q(x_0),\epsilon\sim\mathcal{N}(0,I)}[||\epsilon - \epsilon_\theta(x_t,t)||^2]$$

# Sampling and Training pseudocode

**Algorithm 1** Training

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

**Algorithm 2** Sampling

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

# Classifier Guidance



With no guidance

With classifier guidance

$x_T$

$y$  Label -> "cats"

[Diffusion Models Beats GANs on Image Synthesis (Dhariwal & Nichol 2021)]

# Classifier Guidance



$y$

E.g., cats

$x_t$

Diffusion model

$f(x_t, y)$

Image Classifier

$p(y|x_t)$

- Classifier trained with noisy images

Final Prediction is given by:

$$f(x_t, y) + s . \nabla_{x_t} p(y|x_t)$$

Scale factor for tradeoff.

$$\nabla_{x_t} p(y|x_t)$$

[Diffusion Models Beats GANs on Image Synthesis (Dhariwal & Nichol 2021)]

# Classifier-Free guidance



$y$

$x_t$

Same Model

$f(x_t, y)$

$f(x_t, \text{-})$

Final Prediction is given by:

$$f(x_t, y) + s.(f(x_t, y) - f(x_t, \text{-}))$$

Scale factor for tradeoff.

[Classifier-Free Diffusion Guidance (Ho & Salimans 2021)]

# CLIP Guidance



Final Prediction is given by:

$$f(x, y) + s. \nabla_{x_t} c_i(x_t)^T c_t(y)$$

Scale factor for tradeoff.

$$\nabla_{x_t} p(y|x_t)$$

$x_t$

$y$

$f(x_t, y)$

CLIP

$c_i(x_t)^T c_t(y)$

$\nabla_{x_t} c_i(x_t)^T c_t(y)$

Diffusion model

- CLIP trained with noisy images

# CLIP verses classifier-free guidance



CLIP Guidance



Classifier-Free Guidance

# Comparison



|  |  |  |  |  |
|---|---|---|---|---|
| "a green train is coming down the tracks" | "a group of skiers are preparing to ski down a mountain." | "a small kitchen with a low ceiling" | "a group of elephants walking in muddy water." | "a living area with a television and a table" |

Real Image

XMC-GAN

DALL-E

GLIDE (CLIP Guid.)

GLIDE (CF Guid.)

# References for Diffusion Models

- [Deep unsupervised learning of nonlinear thermodynamics](#), (Sohl-Dickstein et al. 2015).
- [Denoising Diffusion Probabilistic Models](#) (Ho et al. 2020)
- [Diffusion Models Beats GANs on Image Synthesis](#), (Dhariwal & Nichol 2021)
- Classifier-Free Diffusion Guidance (Ho & Salimans 2021)
- [Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding](#)
- Improved Denoising Diffusion Probabilistic Models (Nichol & Dhariwal 2021)
- GLIDE: Towards Photorealistic Image Generation and Editing with Text-Guided Diffusion Models (Ramesh et al. 2022)
- [Understanding Diffusion Models: A Unified Perspective](#) (Luo et al 2022)

# Discussion

- Larger datasets and GPU/TPU usage led to visually stunning generation results.
    - From 1.2M ImageNet to 5B Laion dataset
    - Hundreds of GPU hours for training
- Going forward, it is extremely important to cut costs of these inference algorithms
    - Hinted the use of parallel decoding and MRF methods for cutting down the costs
        - More detailed algorithms will be presented by Dilip and Sadeep
- Progress in generation hinges on evaluation methods
    - Shobhita will present new evaluation methods