

Object Detection/Segmentation

Describe Faster R-CNN. At a high level, how did it improve on Fast R-CNN and R-CNN?

These three methods address object detection. Faster R-CNN works with 4 stages:

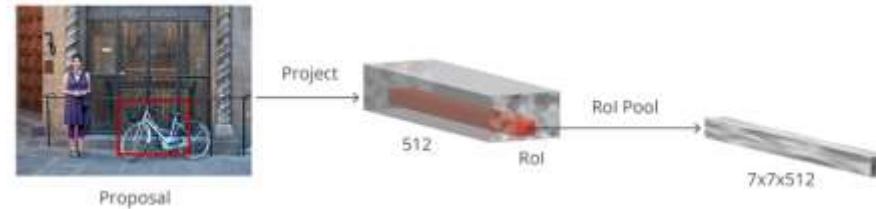
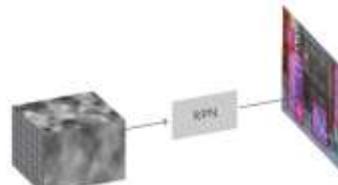
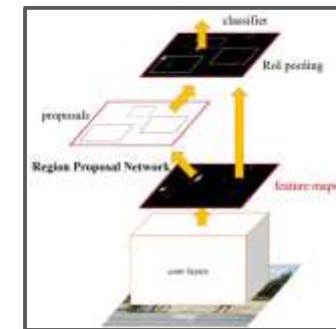
- 1) **Region Proposal Network (RPN)** takes a CNN feature map as input. Then for a large, discrete set of predefined anchors and bounding boxes on the feature map, it outputs
 - a) A "objectness" score (representing foreground/background), formulated as binary classification
 - b) Slight transformations on that bounding box, formulated as regression.

The regular grid of anchor points naturally arises due to the downscaling of the FEN, it doesn't have to be implemented explicitly.

This network is implemented in a fully convolutional way using 1x1 convolutions on the feature map. Only fg objects have a regression loss.
- 2) **Non-maximum suppression (NMS)** tries to remove duplicates, by iterating over the list of proposals sorted by score, and discarding those which have an IoU larger than some threshold with a proposal that has a higher score.
- 3) **A classification head**, which:
 - a) Applies **ROI pooling** to randomly selected fg and bg bboxed feature proposals from stage 1, by projecting from image to feature map, and standardizing the size by splitting the ROI into a grid and max pooling (better than avg pooling in practice).
 - b) Performs a final classification, to label it as one of the target classes, or background. This uses FC layers.
 - c) Performs a final bbox regression loss based on class. This uses FC layers.
- 3) Loss is computed as a classification loss (both fg/bg bboxes) and bbox loss based on L1 dist (on only fg)

Some notes:

- RPN needs to have high recall, since if we don't get proposals here there's no chance they will be classified
- For some applications you can just use RPN to save time (eg if you are only recognizing a single class)
- Comparing the 3 methods:
 - R-CNN does not use a CNN to predict region proposals; it uses a selective search algorithm. Then, it uses a SVM to classify.
 - Fast R-CNN, still uses selective search, but uses ROI pooling and classifies with a CNN.
 - Faster R-CNN adds the RPN, and is thus end-to-end.
 - The time it takes to process an image are approx 50 secs for R-CNN, 2 secs for Fast R-CNN, and 0.2 secs for Faster R-CNN.



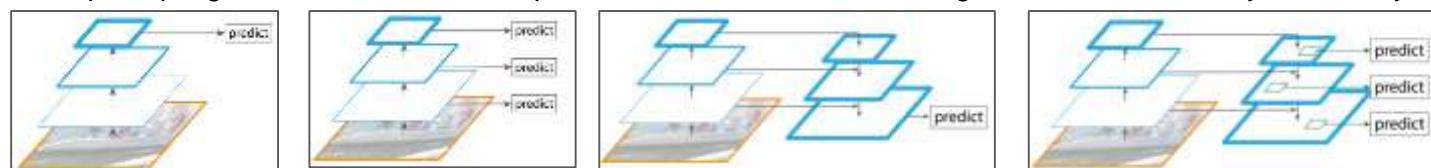
Explain the 1 stage detectors: YOLO, SSD, FPN, and RetinaNet. How do they compare to 2 stage detectors?

YOLO directly regresses bboxes with a CNN + FC layers:

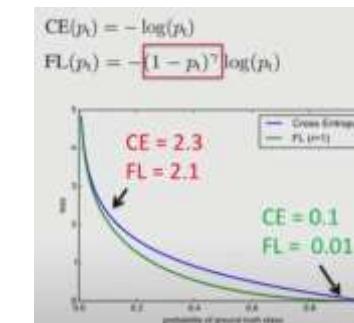
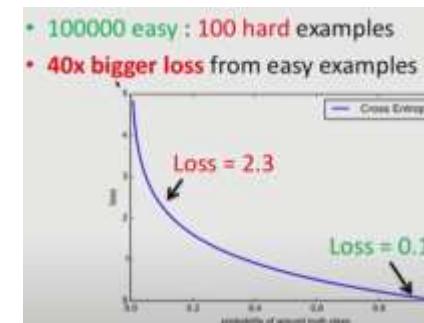
- Input is split to a grid, and the network outputs a fixed number of bboxes per square. For each bbox, there are parameters for (x,y) relative to center of obj, width/height of bbox, confidence, and probabilities for each class.

SSD is similar to YOLO, but predicts independent object detections from multiple feature maps of different resolutions. This helps with detecting small objects.

Feature Pyramid Networks (FPN) implements a U-shaped architecture which enables even better detection of small objects. As in U-Net, the intuition is to combine low-resolution, semantically strong features with high-resolution, semantically weak features via a top-down pathway and lateral connections to get rich semantics at all levels. It uses nearest neighbor upsampling (to “hallucinate” high level semantics, yet at high resolutions) and conv layers for upsampling, while U-Net uses transpose convolutions. From left to right, we show YOLO style, SSD style, UNet style, and FPN style.



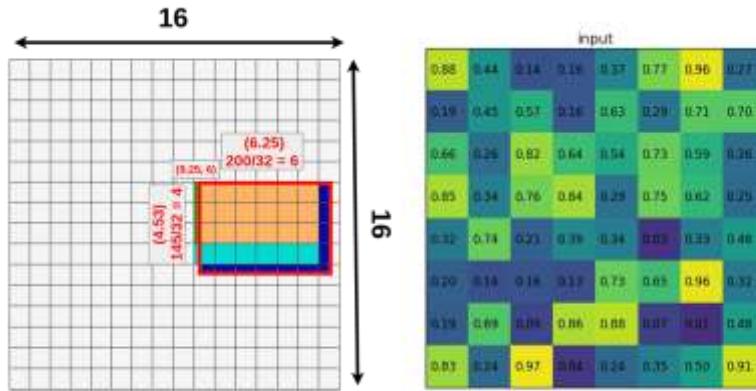
RetinaNet is essentially a combination of a ResNet backbone, FPN neck, and a Focal loss which helps with the fg/bg imbalance issue. Even with Cross Entropy loss's asymptotic behavior, the many easy bg examples can dominate the few hard examples. So, the focal loss adds a simple term to exacerbate the asymptotic behavior.



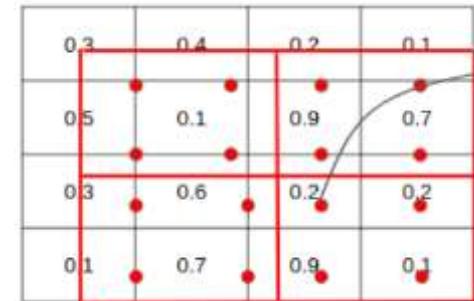
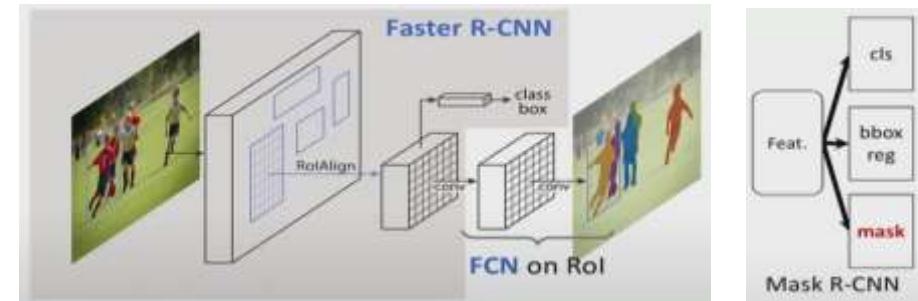
Explain how Mask-RCNN works.

Essentially a Faster R-CNN (ResNet+FPN backbone), with a Fully Convolutional Network (FCN) on each ROI for segmentation.

Uses **RoI Align** instead of RoI (max) Pooling, which does not use any quantization and does not break pixel-to-pixel alignment, by using bilinear interpolation before the max pooling operation.



**RoI (Max)
Pooling**



0.1	0.5	0.48	0.45
0.6	0.68	0.62	0.15
0.32	0.2	0.284	0.15
0.84	0.34	0.34	0.1

2 x 2

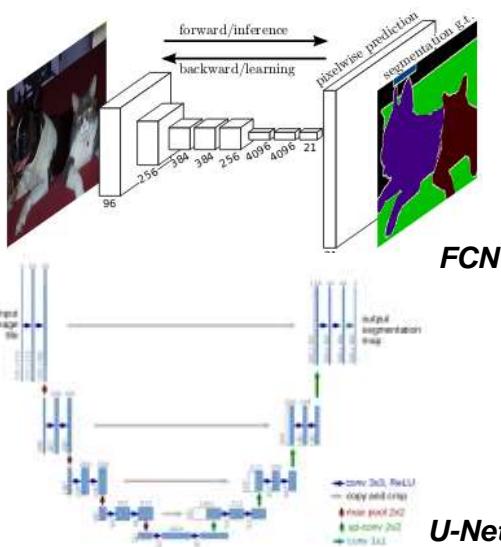
RoI Align

Explain FCN, U-Net, and DeepLab

All three are for semantic segmentation. Here, the key challenge is upsampling, while preserving detail.

Fully Convolutional Network (FCN) stacks a bunch of conv layers in an encoder-decoder fashion.

- No linear layers at the end for classification; only 1x1 convolutions. This allows different size img inputs.
- Upsamples only once in the decoder (using a deconvolution)
- Skip connections are added. The intuition is that this can recover fine-grained spatial information lost in the pooling/downsampling layers for making the final predictions

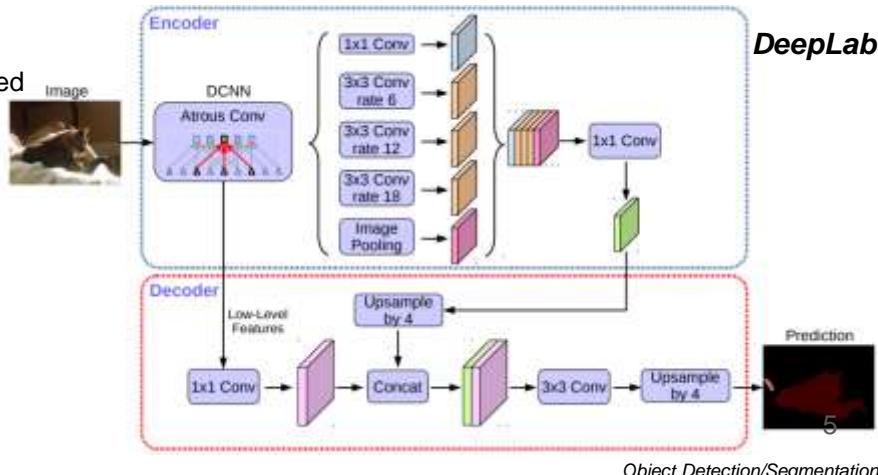
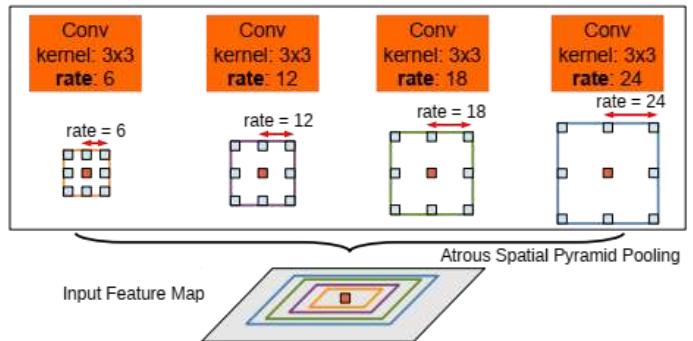


The **U-Net** is similar, but has some modifications:

- Multiple upsampling layers, so that the encoder-decoder is symmetric
- Skip connections are concatenated

The **DeepLabv3+** (from Google) is the best performing of the 3:

- Uses multiple Atrous convolutions, in a depthwise separable manner, at different dilation rates (creating a pyramid) in Encoder
- Upsamples 4x twice using bilinear interpolation, in Decoder
- Original DeepLab used conditional random fields (CRFs), but this was dropped



What losses are used for semantic segmentation?

- Cross entropy
 - Dice loss
 - I.e., IoU loss
 - In practice a “soft dice loss” is used, where A is the probability map mask and the g.t. Mask
 - In general, this is conceptually more “direct”, but cross entropy can be better in implementation because it has nicer gradient properties which do not explode as easily

$$\text{Dice loss} = 1 - \frac{2|A \cap B|}{|A| + |B|}$$

Explain how the Detection Transformer (DETR) works.

- DETR (from facebook) leverages a **transformer encoder-decoder** architecture for object detection
 - Moves away from traditional object detection like Faster-RCNN, which rely on complex pipelines with region proposals, anchor boxes, and NMS
- Transformer decoder attends to encoded features and object queries (fixed learned embeddings), to predict all bounding boxes and classes in parallel (not autoregressive)
 - Final decoder outputs passed to a FFN to predict a detection (class & bbox) or a “no object class”
- **Bipartite Matching Loss (Hungarian Loss)** used to enforce one-to-one mapping between predictions & ground truth
 - Since this method effectively outputs an unordered set, we need to assign GT to predictions before backprop
 - A cost function is computed for all pairings
 - Optimal matching is determined via the Hungarian algorithm (not differentiable, but doesn’t need to be; only a preprocessing step to assign GT)
 - Final loss includes classification loss and Generalized IoU + L1 loss for bounding boxes.

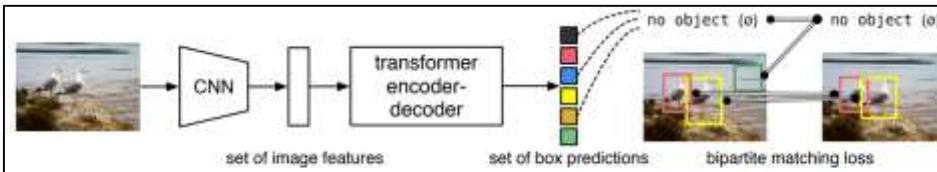


Fig. 1: DETR directly predicts (in parallel) the final set of detections by combining a common CNN with a transformer architecture. During training, bipartite matching uniquely assigns predictions with ground truth boxes. Prediction with no match should yield a “no object” (\emptyset) class prediction.

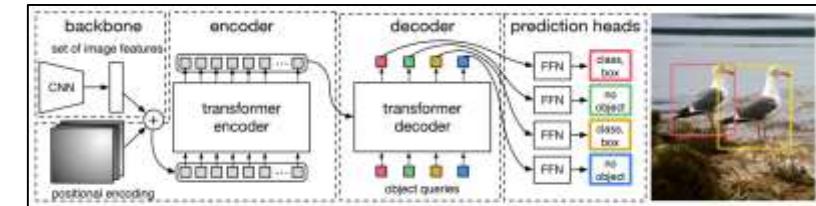


Fig. 2: DETR uses a conventional CNN backbone to learn a 2D representation of an input image. The model flattens it and supplements it with a positional encoding before passing it into a transformer encoder. A transformer decoder then takes as input a small fixed number of learned positional embeddings, which we call *object queries*, and additionally attends to the encoder output. We pass each output embedding of the decoder to a shared feed forward network (FFN) that predicts either a detection (class and bounding box) or a “no object” class.

Explain how Mask2Former works.

- General / “universal” formulation that address any segmentation task (panoptic, instance, & semantic)
 - Predicts a set of binary masks and assigns a single class to each mask. Each mask represents the segmentation of one instance
- Main MaskFormer components for K classes, up to N detections
 - Backbone:** extracts low-resolution image features
 - Pixel Decoder:** upsamples encoded features from backbone to high-resolution per-pixel embeddings
 - Transformer Decoder:** object queries attend to image features. Produces N mask embeddings and class predictions.
 - Final N mask predictions computed by dot product of image features and mask embeddings
 - Matching between set of predictions and ground truth segments done using bipartite Hungarian matching
- Mask2Former improvements:
 - Masked Attention:** Variant of cross attention that only attends within the predicted foreground mask region for each query
 - Leads to more effective learning/convergence; hypothesizes that local features can update query features, and context info can be gathered in self-attention later
 - Implemented by adjusting attention probabilities so that only positive mask regions have nonzero probability
 - Now, different transformer decoder layers can cross-attend to different sizes/resolutions of features
 - Switch order of self and cross/masked attention; now, cross attention goes first
 - Query features to the first self-attention layer are image-independent and do not have signals from the image, thus applying self-attention is unlikely to enrich information
 - Improved training efficiency**
 - Uniformly randomly samples points to calculate bipartite matching step instead using whole mask.
 - In the final loss, sample points in mask for supervision using **importance sampling** (choosing regions with high uncertainty)
 - Remove dropout in residual connections and attention maps
 - Make query features learnable (previously, only query positional embeddings were learnable)

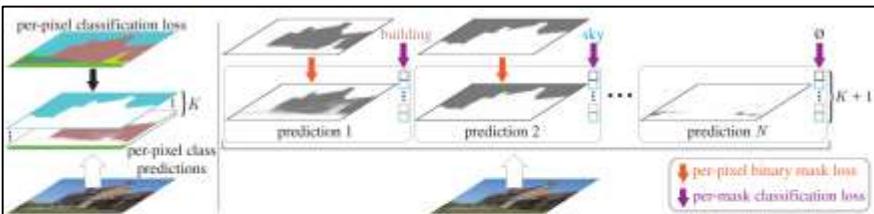


Figure 1: **Per-pixel classification vs. mask classification.** (left) Semantic segmentation with per-pixel classification applies the same classification loss to each location. (right) Mask classification predicts a set of binary masks and assigns a single class to each mask. Each prediction is supervised with a per-pixel binary mask loss and a classification loss. Matching between the set of predictions and ground truth segments can be done either via *bipartite matching* similarly to DETR [4] or by *fixed matching* via direct indexing if the number of predictions and classes match, i.e., if $N = K$.

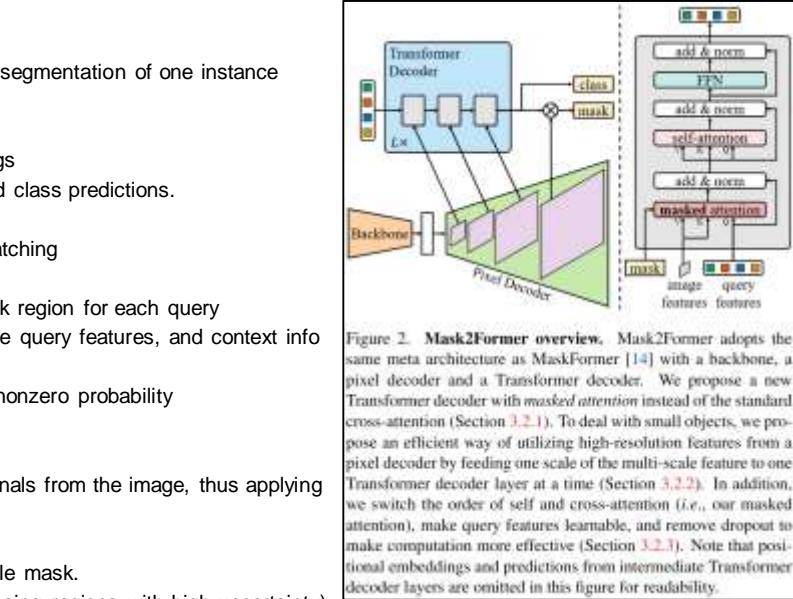
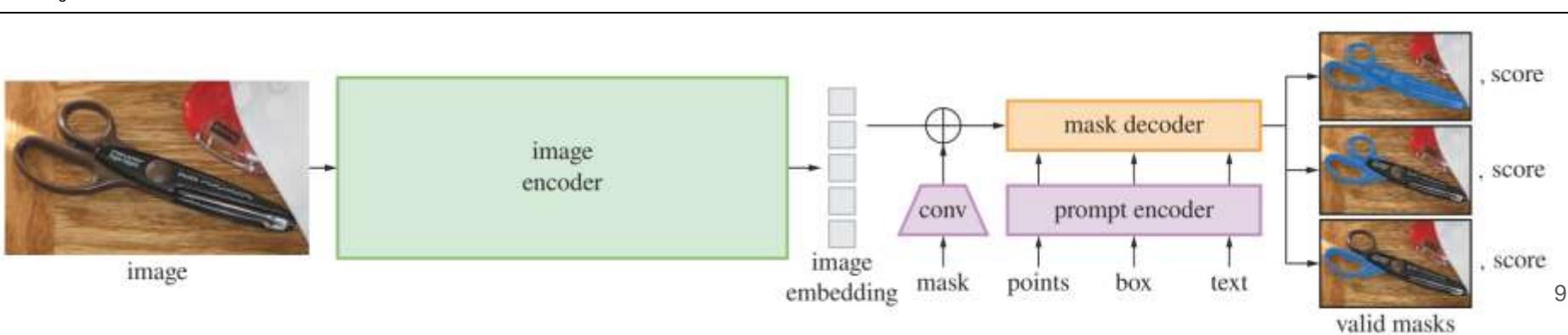
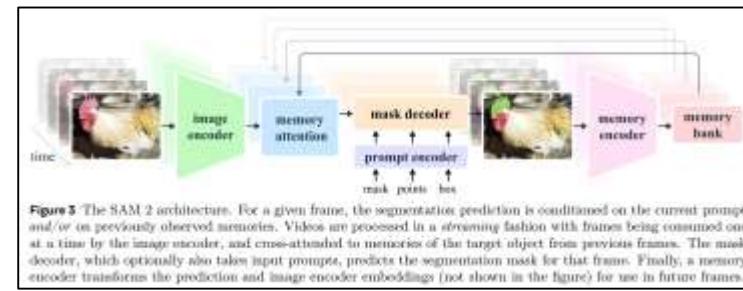
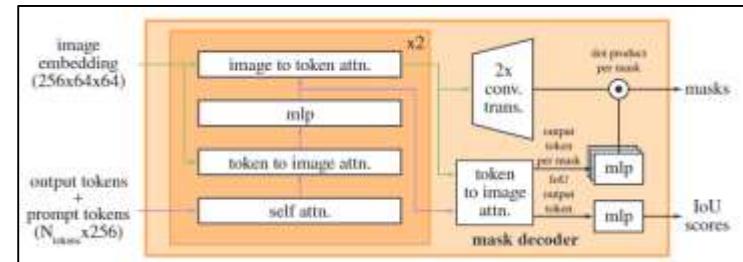


Figure 2: **Mask2Former overview.** We use a backbone to extract image features \mathcal{F} . A pixel decoder gradually upsamples image features to extract per-pixel embeddings $\mathcal{E}_{\text{pixel}}$. A transformer decoder attends to image features and produces N per-segment embeddings \mathcal{Q} . The embeddings independently generate N class predictions with N corresponding mask embeddings $\mathcal{E}_{\text{mask}}$. Then, the model predicts N possibly overlapping binary mask predictions via a dot product between pixel embeddings $\mathcal{E}_{\text{pixel}}$ and mask embeddings $\mathcal{E}_{\text{mask}}$ followed by a sigmoid activation. For semantic segmentation task we can get the final prediction by combining N binary masks with their class predictions using a simple matrix multiplication (see Section 3.4). Note, the dimensions for multiplication \otimes are shown in gray.

How does Segment Anything (SAM) work?

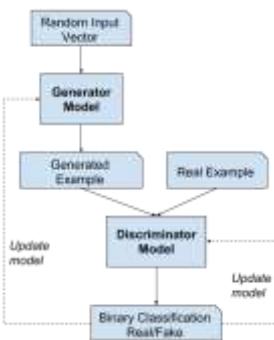
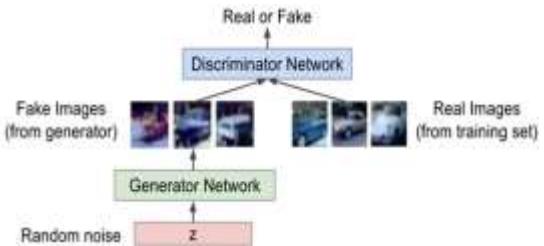
- Given an image and prompt (eg positive/negative points, box, mask, or text), output a valid segmentation mask
 - “Valid” means that even when the prompt is ambiguous, it should still be reasonable.
- Marketed as a **foundation model** for image segmentation
 - Important qualities: **prompt engineerable** and generalizable to solve a range of downstream segmentation problems; easy to fit as a piece in a larger ML system, eg with VR using gaze detection.
- Architecture consists of different types of prompt encoders (outputting 256 dimensional vectors), an image encoder (Masked Autoencoder ViT outputting 256x64x64 feature map), and a Mask decoder similar to MaskFormer
 - Multiple (3) output tokens, to naturally allow for ambiguity. We only backprop on the one with the lowest loss.
 - Mask decoder block has cross-attention in both directions
 - Also predicts IoU scores, supervised by true IoU score
 - Text embedding can be used with clip embeddings
- Training procedure simulates interactive 11-round setups
 - First, either randomly select an initial point or bbox from the GT mask
 - Iteratively make predictions and adding new query point prompts based on the error region between prediction and GT; each new point is foreground or background if the error point is a false negative or false positive respectively. Additionally re-feed mask into network.
 - Only backprop after the 11-rounds.
- Comes with a 1.1B mask dataset
 - Collected with extensive data engine starting with assisted manual annotations
 - Labelers had simple instructions: to just label objects they could name/describe; did not collect names/descriptions
 - Continually bootstrapped SAM model to improve automation; final dataset is composed of all pseudolabels by prompting model with regular grid and filtering for confident masks.
 - All images come from photographers, but is able to zero-shot generalize decently to other domains / image styles
- SAM2 adds streaming memory-based transformer for video segmentation & occlusion prediction head; trained on new video segmentation dataset



Generative Modeling

How do (W)GANs work, and how are they optimized? Can they be conditioned on?

GANs implicitly model probability density by focusing on the ability to sample from an implicit distribution, in an unsupervised way. This is done through a two-player game between a generator and discriminator net. Convergence is reached at “nash equilibrium”, where which each player cannot reduce their cost without changing the other players’ parameters.



1. Gradient ascent on discriminator

$$\max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

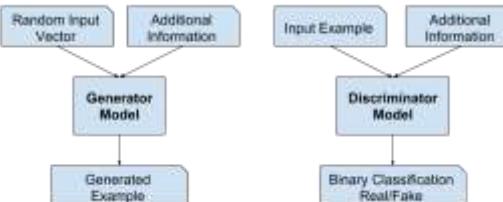
2. Instead: Gradient ascent on generator, different objective

$$\max_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(D_{\theta_d}(G_{\theta_g}(z)))$$

The **Wasserstein GAN (WGAN)** is an alternative formulation which may have better theoretical properties.

- Discriminator becomes a “critic”, which optimizes towards an approximation of the Wasserstein distance between the real and generated distribution.
- Unlike standard GANs, we want to train critic to convergence. In principle, WGAN loss should be informative (unlike standard gans losses, which have no immediate meaning for how training is progressing).
- The critic/generator are supposed to be K-Lipschitz, so gradient clipping or penalty is needed.
- The critic is optimized to maximize the distance between average real and generated outputs; this approximates the wasserstein loss (lines 2-8)
 - Critic outputs “realness score”; lower critic score= fake, higher critic score= real
- Then, the generator is trained to maximize the critic’s score

(W)GANs can be conditioned (cGAN) on by providing both the generator and discriminator with some additional input.



Algorithm 1 WGAN, our proposed algorithm. All experiments in the paper used the default values $\alpha = 0.00005$, $c = 0.01$, $m = 64$, $n_{critic} = 5$.

Require: : α , the learning rate. c , the clipping parameter. m , the batch size. n_{critic} , the number of iterations of the critic per generator iteration.

Requires: : w_0 , initial critic parameters. θ_0 , initial generator’s parameters.

```

1: while  $\theta$  has not converged do
2:   for  $t = 0, \dots, n_{critic}$  do
3:     Sample  $\{x^{(i)}\}_{i=1}^m \sim P_r$  a batch from the real data.
4:     Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples.
5:      $g_w \leftarrow \nabla_w \left[ \frac{1}{m} \sum_{i=1}^m f_w(x^{(i)}) - \frac{1}{m} \sum_{i=1}^m f_w(g_w(z^{(i)})) \right]$ 
6:      $w \leftarrow w + \alpha \cdot RMSProp(w, g_w)$ 
7:      $w \leftarrow clip(w, -c, c)$ 
8:   end for
9:   Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples.
10:   $g_\theta \leftarrow -\nabla_\theta \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)}))$ 
11:   $\theta \leftarrow \theta - \alpha \cdot RMSProp(\theta, g_\theta)$ 
12: end while
  
```

Explain how Pix2Pix, CycleGAN, and conditional GANs work, and what they're useful for.

Both methods perform the task of image-to-image translation. Pix2Pix requires paired data, while CycleGAN does not.

Pix2Pix is a conditional GAN (cGAN); the discriminator is fed both domains to discriminate.

CycleGAN for domains X and Y, learns a mapping $G: X \rightarrow Y$ and inverse mapping $F: Y \rightarrow X$. The outputs of both have their own discriminator. It's enforced that cycles $F(G(x))=x$ and $g(F(y))=y$ hold.

Conditional GAN: both the generator and discriminator are conditioned (fed as input) with some additional information, like class labels or data from other modalities. Since discriminator now has to also ensure that input images belong to the correct class/criteria, these gradients guide the generator to create more specific images of the appropriate class/criteria

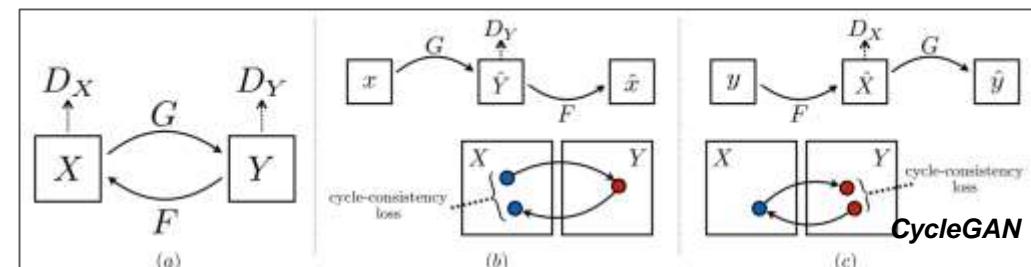
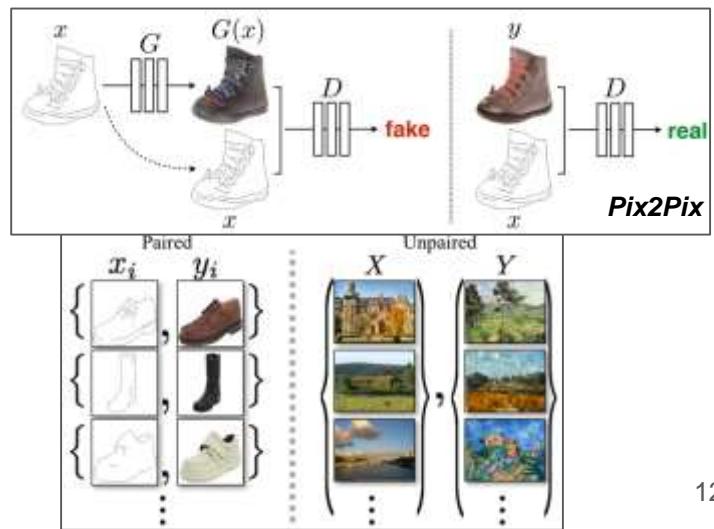
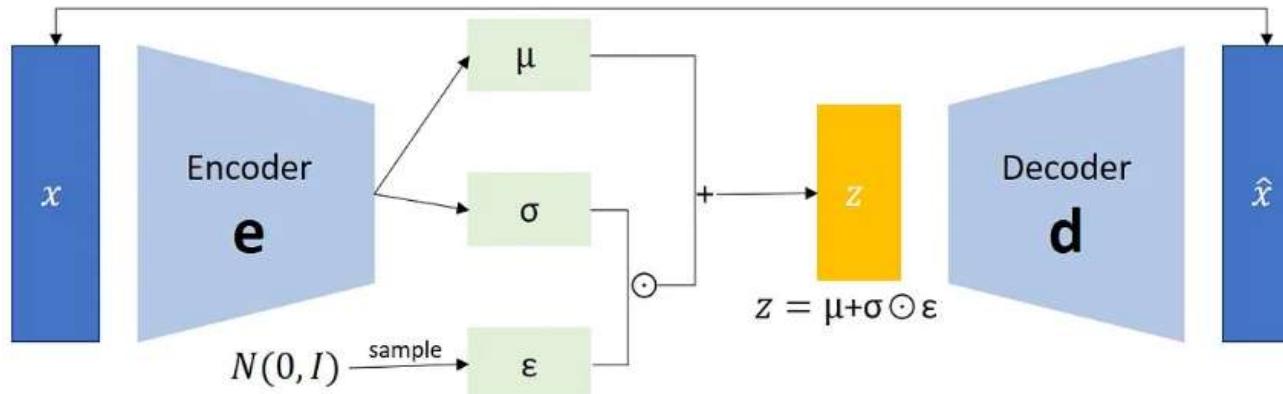


Figure 3: (a) Our model contains two mapping functions $G : X \rightarrow Y$ and $F : Y \rightarrow X$, and associated discriminators D_Y and D_X . D_Y encourages G to translate X into outputs indistinguishable from domain Y , and vice versa for D_X and F . To further regularize the mappings, we introduce two *cycle consistency losses* that capture the intuition that if we translate from one domain to the other and back again we should arrive at where we started: (b) forward cycle-consistency loss: $x \rightarrow G(x) \rightarrow F(G(x)) \approx x$, and (c) backward cycle-consistency loss: $y \rightarrow F(y) \rightarrow G(F(y)) \approx y$



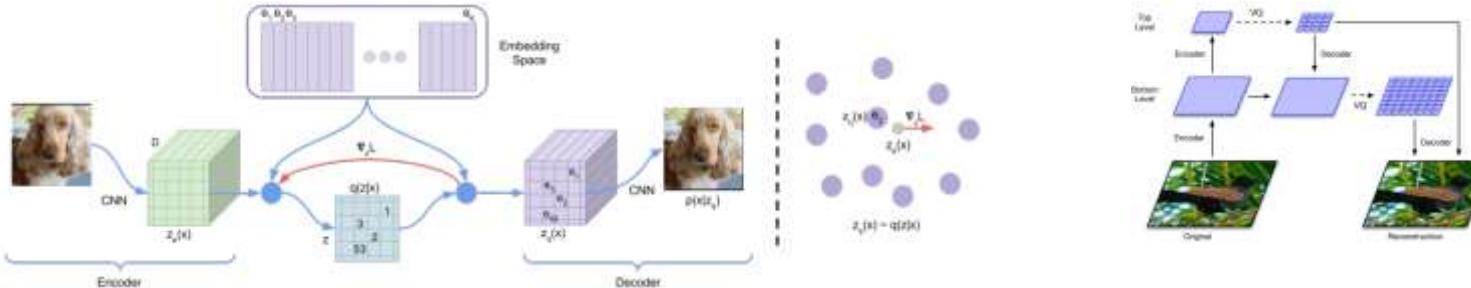
How do VAEs work?

- Encoder $q_\phi(z|x)$ ie posterior, decoder $p_\theta(x|z)$, prior distribution $p(z) = N(0, I)$
- Encoder outputs parameters μ, σ of a gaussian distribution
- Loss involves reconstruction term and KL divergence to regularize to standard normal:
$$\mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|z)] - D_{KL}(q_\phi(z|x)\|p(z))$$
- Conducive to compact, & semantically organized latent space for stochastic sampling
- Reparameterization trick where $\epsilon \sim N(0, I)$ allows gradients to flow through μ, σ
- Can generate new samples from test time by sampling from $N(0, I)$ & passing through decoder
- Weaknesses: produces blurry images due to averaging nature of the reconstruction loss & normal regularization
 - In limit, can suffer to posterior collapse; encoder collapses to the prior distribution, causing latent variables to be uninformative and leading to poor diversity in generated samples



How does VQ-VAE work?

- During training an encoder maps the image to a continuous feature volume $z_e(x)$; each channel is then quantized to the nearest entry in a learned, finite-sized codebook to obtain $z_q(x)$; this is then decoded to reconstruct the original image.
 - The argmin “snapping” operator here is non-differentiable, but in practice everything seems to work fine if you just approximate & pass the decoder gradient directly through to the encoder
- Loss function contains 1) reconstruction loss, 2) “codebook alignment loss” to incentive learned codebook embeddings to be close to the encoder output. and 3) “codebook commitment loss” to encourage encoder outputs to be close to codebook embeddings
 - $\log(p(x|q(x))) + \|\text{sg}[z_e(x)] - e\|_2^2 + \beta\|z_e(x) - \text{sg}[e]\|_2^2$
 - Here, sg stands for “stop gradient” and prevents the gradient from flowing through that part
- The use of discrete latents facilitates the use of autoregressive token-by-token latent priors for speech/audio/images
- At test time, we can learn a new prior $p(z)$ by 1) taking a dataset of images, 2) obtaining their corresponding sequences of codebook tokens with the encoder & codebook, 3) training an autoregressive generative model (eg pixelCNN or transformer) model on those sequences, 4) generating new novel sequences, and 4) feeding the sequences into the decoder for image generation
 - In this paradigm, intuitively pixelCNN would be the “creative director” generating a sequence of semantically meaningful latents and the decoder would carry out the work of realizing those ideas in pixel space
- Some later works improve this by 1) using vision transformers (ViT) for the encoder/decoder, 2) using multiple scales w/ corresponding different scale codebooks, 3) adding styleGAN discriminator losses, 4) adding perceptual loss



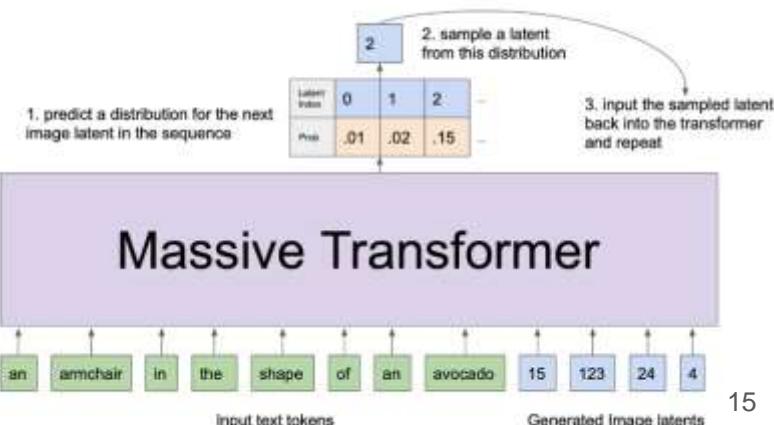
How does DALL-E 1 implement text to image generation?

- The first step involves using a Discrete VAE (dVAE), which is similar to VQ-VAEs, but instead of the encoder trying to match codebook values, it outputs a distribution over codebook vectors for each latent, which is more expressive. (Note, this is not a fundamental change, and in principle DALL-E could have used vanilla VQ-VAEs.)
- To backpropagate through sampling from this categorical distribution, it uses a gumbel softmax relaxation:
 - First, we use a “gumbel” reparameterization trick to allow us to disentangle nn parameters from the source of randomness when categorically sampling.
 - It can be shown that if g_i are iid samples from the Gumbel distribution, then the following allows for random categorical sampling, where $q(e_i|x)$ is the probability prediction of the e_i^{th} codebook vector
$$z = \text{codebook}[\text{argmax}_i[g_i + \log(q(e_i|x))]]$$
 - Second, we use a softmax relaxation with a temperature τ that approaches 0 throughout training, instead of argmax. This allows us to “softly” choose the max via weighted sum in a backpropagatable way
 - When τ is large, this resembles a uniform distribution
 - When $\tau = 1$, this is a softmax over the categorical distribution
 - As $\tau \rightarrow 0$, the differences become exaggerated and it approaches a max
$$y_i = \frac{e^{\frac{g_i + \log(q(e_i|x))}{\tau}}}{\sum_{j=1}^k e^{\frac{g_j + \log(q(e_j|x))}{\tau}}}$$
- After the dVAE is trained, we can train a transformer to predict sequences of text, followed by image latents (which we can subsequently look up in codebook & feed to the dVAE decoder)
 - DALL-E was trained on 250 million text-image pairs scraped from the internet

TEXT PROMPT

an armchair in the shape of an avocado [...]

AI-GENERATED IMAGES



How does a Masked Generative Image Transformer (MaskGIT) work?

- MaskGIT uses VQVAE-style tokenization, but modifies the token generation from sequential AR to an iterative parallel approach by using a bidirectional transformer decoder similar to BERT.
- During training, we learn to predict randomly masked tokens by attending to tokens in all directions
- During inference, we iteratively predict all tokens simultaneously, but only keep the most confident ones; remaining tokens are masked out and re-predicted in the next iteration.
 - In principle the model is able to infer all tokens and generate the entire image in a single pass, but this doesn't work well due to the training task's formulation
- Shown to also natively work well for zero-shot tasks such as class-conditioned image manipulation, inpainting, and outpainting
- Much faster, intuitively like solving a jigsaw puzzle: predicting masked parts of an image in parallel, guided by confidence and contextual understanding

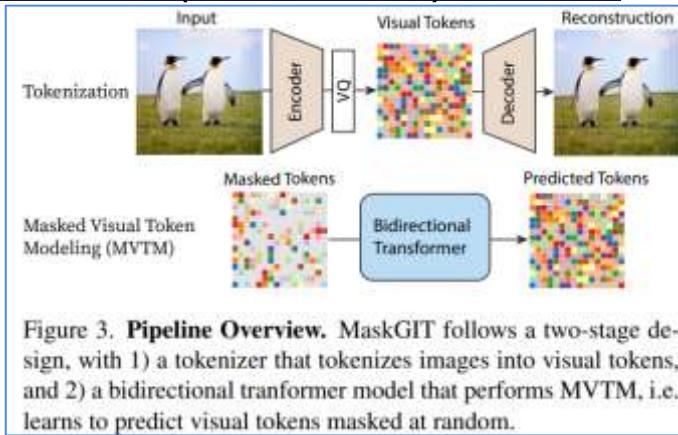
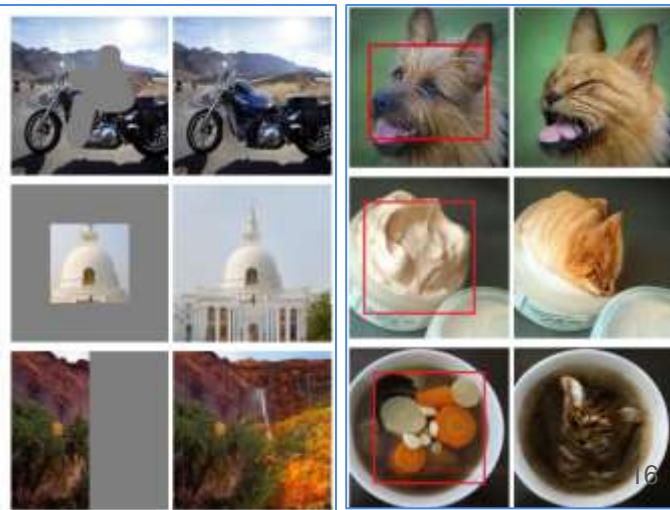
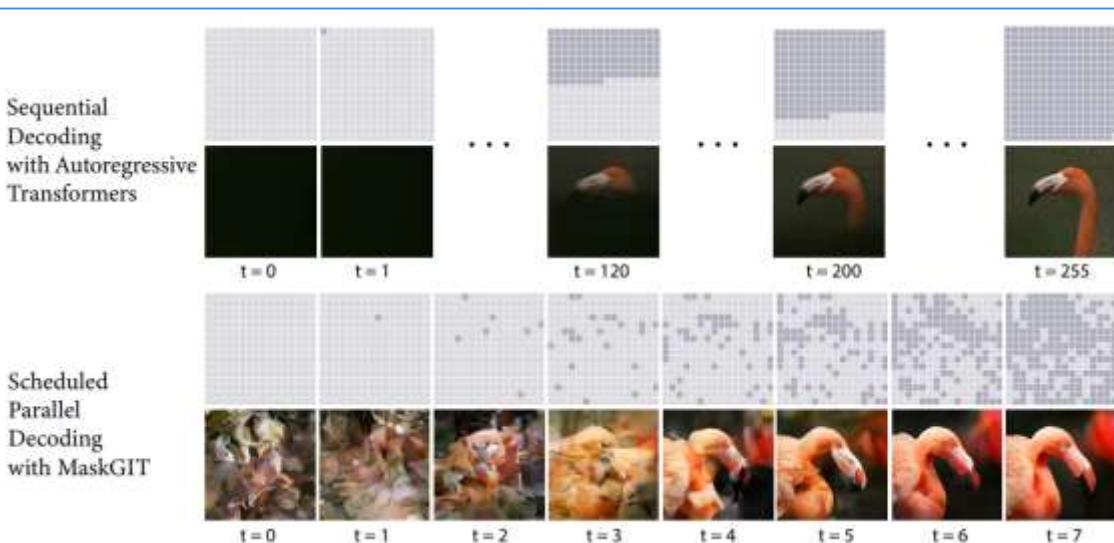


Figure 3. **Pipeline Overview.** MaskGIT follows a two-stage design, with 1) a tokenizer that tokenizes images into visual tokens, and 2) a bidirectional tranformer model that performs MVTM, i.e. learns to predict visual tokens masked at random.



Explain how Visual AutoRegressive (VAR) models work.

- Motivation: VQVAE & similar works perform flattening & process tokens unidirectionally (eg raster-scan or spiral), even though on images there are bidirectional spatial correlations
- VAR reformulates the task by shifting from “next token” to “next scale” prediction; autoregressive unit is an entire token map, instead of just a single token. Leads to SoTA performance and speed
- Multi-scale VQVAE Encoding (algorithm 1 below):
 - Start with an initial feature map f . For K iterations, we successively go from low resolution to high resolution, storing quantized feature map interpolations (ie token maps) & residually updating f by subtracting the codebook vectors interpolated to the original size (ϕ_k conv layers follow the interpolation)
- Multi-scale VQVAE Decoding (algorithm 2 below):
 - Given the encoded list of multi-scale token maps, sum them all by converting them to codebook vectors and interpolating them to the original size. Then, pass that feature volume to the decoder to reconstruct
- Compound loss used, with reconstruction loss + codebook quantization error loss + perceptual loss + styleGAN discriminator loss
- Standard decoder-only transformer used for generation, where in the k^{th} autoregressive step we treat the sequence of tokens (r_1, \dots, r_{k-1}) as the conditional prefix and generate $h_k \times w_k$ tokens in

Algorithm 1: Multi-scale VQVAE Encoding

```
1 Inputs: raw image  $im$ ;  
2 Hyperparameters: steps  $K$ , resolutions  
     $(h_k, w_k)_{k=1}^K$ ;  
3  $f = \mathcal{E}(im)$ ,  $R = []$ ;  
4 for  $k = 1, \dots, K$  do  
5    $r_k = Q(\text{interpolate}(f, h_k, w_k))$ ;  
6    $R = \text{queue\_push}(R, r_k)$ ;  
7    $z_k = \text{lookup}(Z, r_k)$ ;  
8    $z_k = \text{interpolate}(z_k, h_K, w_K)$ ;  
9    $f = f - \phi_k(z_k)$ ;  
10 Return: multi-scale tokens  $R$ ;
```

Algorithm 2: Multi-scale VQVAE Reconstruction

```
1 Inputs: multi-scale token maps  $R$ ;  
2 Hyperparameters: steps  $K$ , resolutions  
     $(h_k, w_k)_{k=1}^K$ ;  
3  $\hat{f} = 0$ ;  
4 for  $k = 1, \dots, K$  do  
5    $r_k = \text{queue\_pop}(R)$ ;  
6    $z_k = \text{lookup}(Z, r_k)$ ;  
7    $z_k = \text{interpolate}(z_k, h_K, w_K)$ ;  
8    $\hat{f} = \hat{f} + \phi_k(z_k)$ ;  
9    $\hat{im} = \mathcal{D}(\hat{f})$ ;  
10 Return: reconstructed image  $\hat{im}$ ;
```



How does diffusion work?

Forward Diffusion:

$$0 < \beta_1 < \beta_2 < \dots < \beta_T < 1 \quad (\text{Known variance schedule})$$

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

$$\mathbf{x}_t = \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \sqrt{\beta_t} \epsilon \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (\text{Variance towards 1 while } \mathbf{x} \text{ towards 0; approaches unit normal})$$

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

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

Reverse Diffusion (Denoising) via Bayes Rule:

$$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\mu}(\mathbf{x}_t, \mathbf{x}_0), \tilde{\sigma}_t \mathbf{I})$$

$$\tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} * \beta_t \quad (\text{Derived result; can show that the reverse itself is close to Gaussian as well})$$

$$\tilde{\mu}_t = \frac{1}{\sqrt{\bar{\alpha}_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_t \right)$$

Objective Function Formulation:

$$\|\epsilon - \epsilon_\theta(\mathbf{x}_t, t)\|^2 = \|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{(1 - \bar{\alpha}_t)} \epsilon, t)\|^2 \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

- Take batch of images
- Forward diffuse to randomly sampled noisy timesteps
- NN predicts the added gaussian noise relative to their respective previous timestep
- NN receives a mixture of image & noise; it's its job to understand the distribution of images and be able to separate out the noise

Test-Time Sampling:

Algorithm 2 Sampling

```

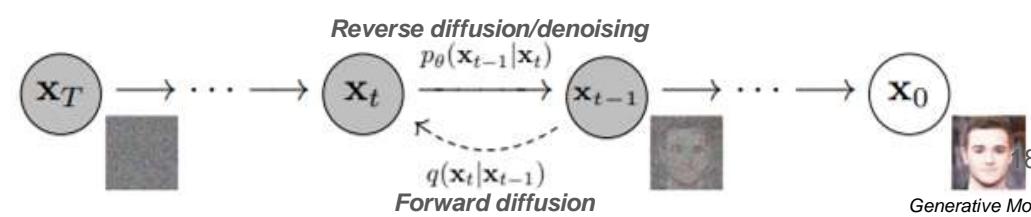
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1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

```

- Note: mathematically, in principle we can directly jump from x_T to x_0 based on the objective function. But this isn't done in practice b/c the network is trained to work on many timesteps, as it empirically outperforms.

Other Details:

- Original arch is U-Net w/ transformer layers & groupnorm
- t encoded as sinusoidal positional embedding & added to feature maps multiple times throughout, where position embeddings dimensions = # of channels
- One of the keys to making this work is the concept of iterative refinement: rather than generating a sample in a single pass through a neural network, the model is applied repeatedly to refine a canvas, and hence the unrolled sampling procedure corresponds to a much "deeper" computation graph (similar to LLMs). The paradigm is train-efficient since we don't need to backprop through everything; generative process not trained end-to-end



Explain how diffusion models can be conditioned, with and without a classifier.

Classifier Guidance:

- If you train a classifier $p_\phi(y|x_t)$ on noisy images, we can use its gradients to guide the diffusion sampling process towards a class-semantic correct image
- Classifier tells us which direction in image space increases probability of the class, and we follow that direction during generation
- No need to retrain the diffusion model, but training classifier to be accurate on noisy images can be challenging

Classifier-Free Guidance

- Involves training a diffusion model to be exposed to both conditioned & unconditioned inputs (via null labels)
- Then, at inference use a linear combination of conditional and unconditional score estimates ($w > 1$)
- Works better than classifier guidance b/c we've constructed the "classifier" from a generative model, which provides much more robust gradients

$$\tilde{\epsilon}_\theta(\mathbf{z}_\lambda, \mathbf{c}) = (1 + w)\epsilon_\theta(\mathbf{z}_\lambda, \mathbf{c}) - w\epsilon_\theta(\mathbf{z}_\lambda)$$

Algorithm 1 Classifier guided diffusion sampling, given a diffusion model $(\mu_\theta(x_t), \Sigma_\theta(x_t))$, classifier $p_\phi(y|x_t)$, and gradient scale s .

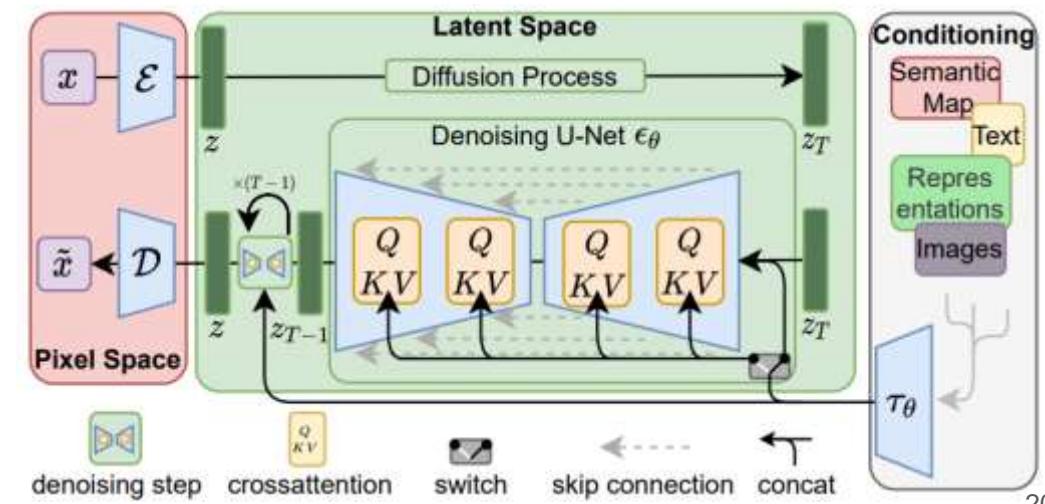
```
Input: class label  $y$ , gradient scale  $s$ 
 $x_T \leftarrow$  sample from  $\mathcal{N}(0, \mathbf{I})$ 
for all  $t$  from  $T$  to 1 do
     $\mu, \Sigma \leftarrow \mu_\theta(x_t), \Sigma_\theta(x_t)$ 
     $x_{t-1} \leftarrow$  sample from  $\mathcal{N}(\mu + s\Sigma \nabla_{x_t} \log p_\phi(y|x_t), \Sigma)$ 
end for
return  $x_0$ 
```



Unconditional (left) vs conditional (right)

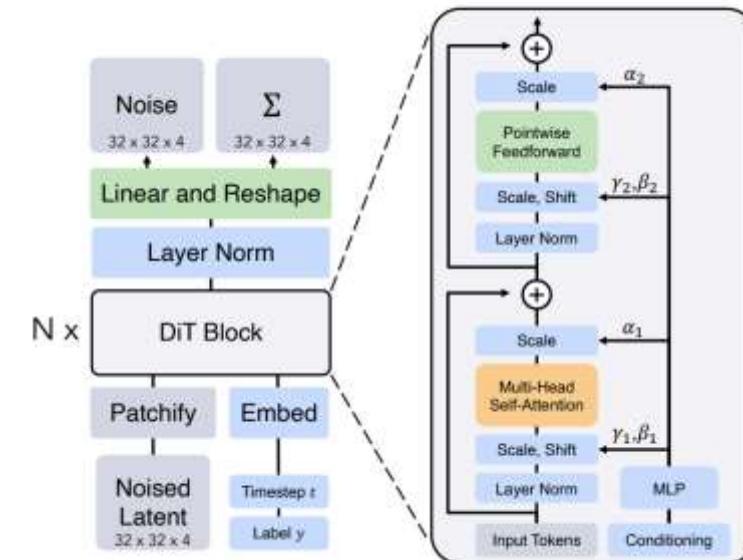
How can diffusion be done in the latent space?

- Two stage recipe:
 - First, train some type of autoencoder (eg vanilla, VAE, VQ-VAE, etc) on input signals
 - Second, train generative diffusion model on latent representations. In this stage, encoder parameters are frozen (and we don't need the decoder until inference)
- Intuition: by extracting a more compact representation that focuses on the perceptually relevant fraction of signal content, and modelling that instead of the original representation, we are able to make relatively modestly sized generative models punch above their weight.



How do Diffusion Transformers Work?

- Transformer-based diffusion model operating on latent representations
 - Latents are from a frozen pretrained VAE encoder, which is generally convolution / U-Net based
- Uses AdaLN-Zero for conditioning, which was found to perform better than cross-attention and in-context tokens
 - Embeds timestamp and class label with separate MLPs into e_t, e_y , then adds to get e
 - Layernorm residual scaled with separate MLPs: $h_{out} = h + \alpha(e) * Block(\gamma(e) * LayerNorm(h) + \beta(e))$



Latent Diffusion Transformer

DiT Block with adaLN-Zero

How can diffusion sampling be made deterministic?

- Forward process (noising) is the same as before
- Reverse process involves first obtaining \widehat{x}_0 , a “one-shot” estimate of x_0
- Then, the update is a weighted combination of the clean image \widehat{x}_0 and the noise direction
- Allows for a deterministic mapping from noise -> image, and image back to noise

$$\widehat{x}_0 = \frac{1}{\sqrt{\alpha_t}} (x_t - \sqrt{1 - \alpha_t} \epsilon_\theta(x_t, t))$$

$$x_{t-1} = \sqrt{\alpha_{t-1}} \cdot \widehat{x}_0 + \sqrt{1 - \alpha_{t-1}} \cdot \epsilon_\theta(x_t, t)$$

How can diffusion be sped up using distillation?

- Train a standard teacher DDIM model
- Perform inference on a dataset of images, and cache some intermediate partially-denoised results
- Train student to directly map those intermediate results

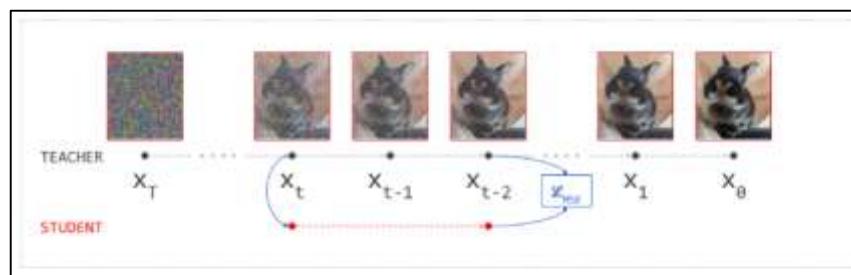


Diagram showing progressive distillation. The student learns to match the result of two sampling steps in one forward pass.

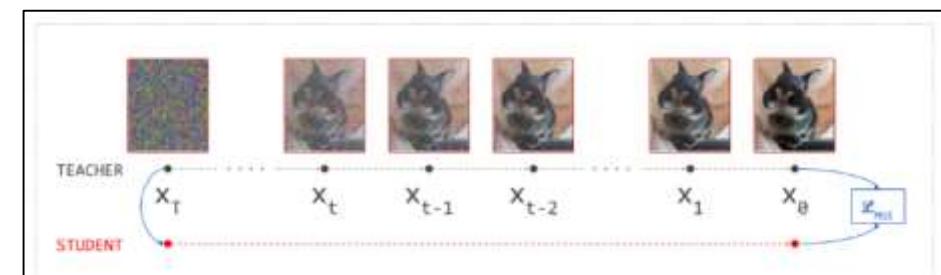


Diagram showing distillation of the diffusion sampling procedure into a single forward pass.

What is ControlNet?

- Method to enable precise spatial control/conditioning in large text-to-image diffusion models, without modifying the original pretrained model
 - Inputs include edge maps, human poses, depth maps, segmentation maps, sketches, etc
- Freezes pretrained model weights and connects it to zero-initialized 1x1 convolutions (so that it avoids injecting noise too early, and promote gradually learning how to modulate network)
- Requires paired data with conditioning, and optionally paired text captions



Figure 1: Controlling Stable Diffusion with learned conditions. ControlNet allows users to add conditions like Canopy edges (top), human pose (bottom), etc., to control the image generation of large pretrained diffusion models. The default results use the prompt "a high-quality, detailed, and professional image". Users can optionally give prompts like the "chef in kitchen".

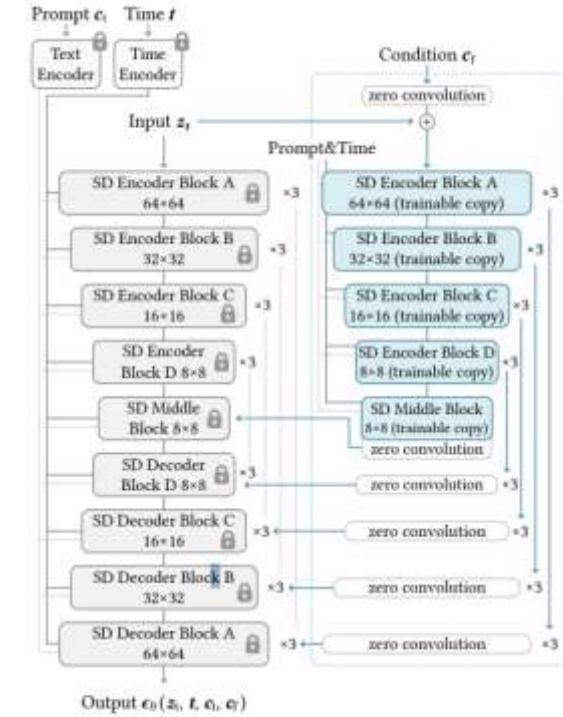


Figure 3: Stable Diffusion's U-net architecture connected with a ControlNet on the encoder blocks and middle block. The locked, gray blocks show the structure of Stable Diffusion V1.5 (or V2.1, as they use the same U-net architecture). The trainable blue blocks and the white zero convolution layers are added to build a ControlNet.

How can diffusion be used to do inpainting?

- Can be done without retraining the diffusion model; during inference timesteps, noise known (unmasked) regions to the appropriate corresponding noise level, while unknown (masked) regions are initialized with random noise and proceed as usual
 - Has additional benefit to not restrict training to a specific distribution of mask, which can limit generalization
 - The RePaint paper additionally applies strategy to include periodic jumps forward/back in diffusion steps to refine results
- Alternatively, you can train the model to directly do inpainting by training on masked images
 - Network trained to see partially noised images, in the masked regions

$$x_t = M \odot x_{\text{noisy}} + (1 - M) \odot x_0$$

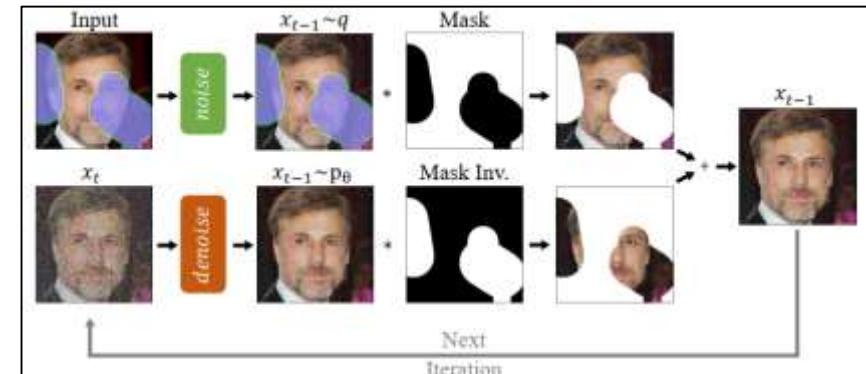
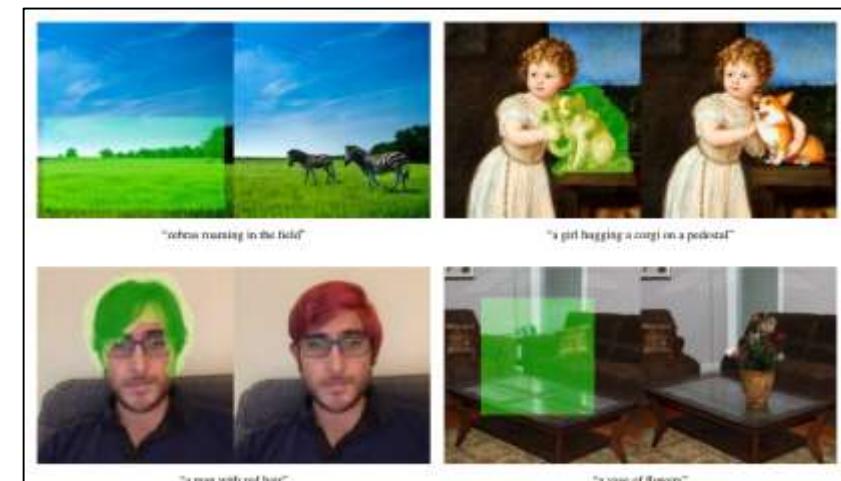


Figure 2. **Overview of our approach.** RePaint modifies the standard denoising process in order to condition on the given image content. In each step, we sample the known region (*top*) from the input and the inpainted part from the DDPM output (*bottom*).

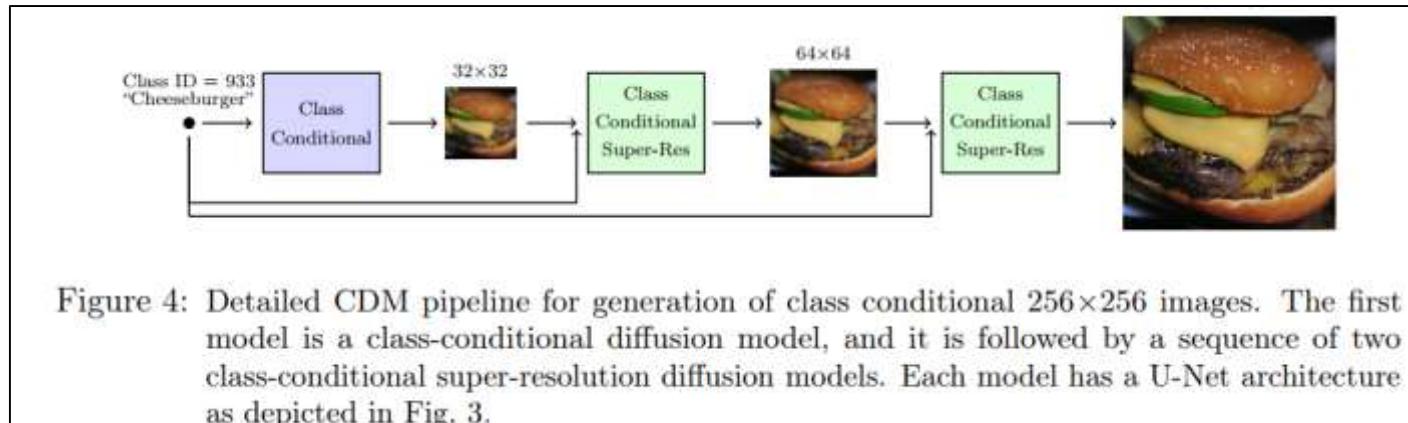
What is GLIDE and how can it perform text-guided inpainting?

- Text-guided diffusion (using classifier-free guidance); outperforms CLIP-guided diffusion where you train a clip model on noisy images
- Can do text-guided inpainting by finetuning model to inpaint with text conditioning
 - Masked regions don't necessarily correspond to text captions during training – suboptimal, but in practice still seems to work
 - The model is able to see a fair amount of examples where the masked region does correspond, and over time learns to provide desired behavior



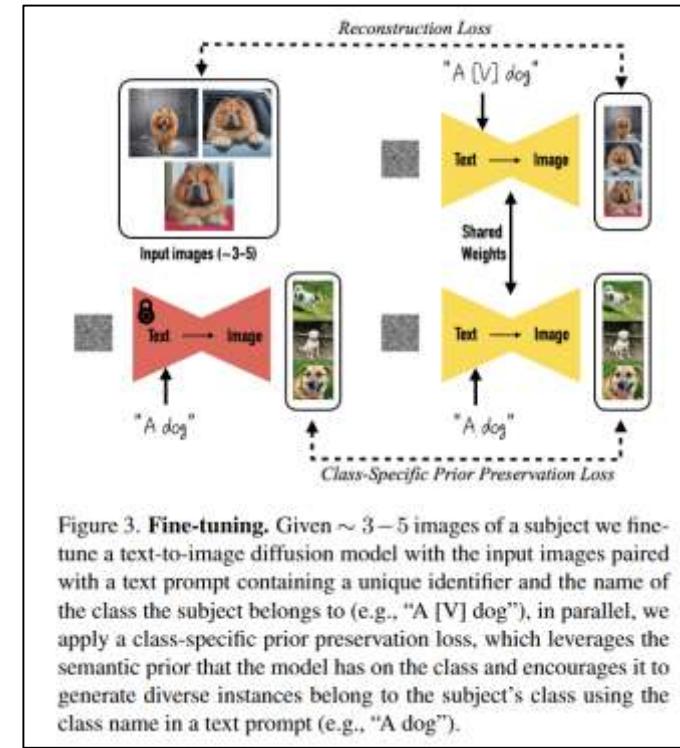
What is cascaded diffusion?

- Framework for generating high-resolution images using a cascade of separate diffusion models
 - Consists of a base diffusion model that generates a low res image (eg 32 x 32), then several superresolution models that upsample to higher resolutions
 - Models not trained end-to-end; they're separate models
- Superresolution model turns high-res gaussian noise into a high-res image, conditioned on a low-resolution image
- Key trick is augmentation to superresolution model to minimize train-test mismatch (train is natural images, test is high-res images; different distributions)
 - Variants include truncated diffusion results (adds noise), or gaussian noise/blur to bridge the gap



How does DreamBooth work?

- Tackles subject-driven generation: method to personalize large text-to-image diffusion models so they can generate novel, high-fidelity images of a specific subject (like a pet or object) using only 3–5 reference images
- Fine-tunes a pretrained diffusion model to bind a **unique textual token** to a specific subject instance (eg a particular dog). Enables model to recognize and regenerate the object, placed in various contexts
- Method:
 - Finetune model using few images, paired with prompts like “a [V] dog”, where [V] is a rare token. The prompt should include the class so the model can leverage the class’s prior
 - It’s also important to have a “class-specific prior preservation loss”, where we have the model continue training on general images that involve the class of the specific subject. This significantly reduces overfitting; notably, authors found that training in generated images (ie, psuedolabels) works best. This encourages model to continue to generate diverse samples in the subject’s class
- Early stopping also used to mitigate overfitting; training is short, with a small learning rate



How does MAGVIT work?

First, trains 3D VQ-GAN using 3D conv based architecture

- Encodes videos into codebook, enabling space-time compression (quantizes 16 frame, 128x218 video into 4x16x16 tokens, where each token is a spatiotemporal supervoxel)
- Results in a 3D-VQ encoder to obtain video tokens, and a 3D-VQ decoder to convert tokens back to pixel space

Second, trains a BERT-like bidirectional transformer to do masked token modeling

- Learns to predict masked tokens based on context, class, and task specific prompts
- Input to BERT composed of three parts, all with the same reconstruction loss.

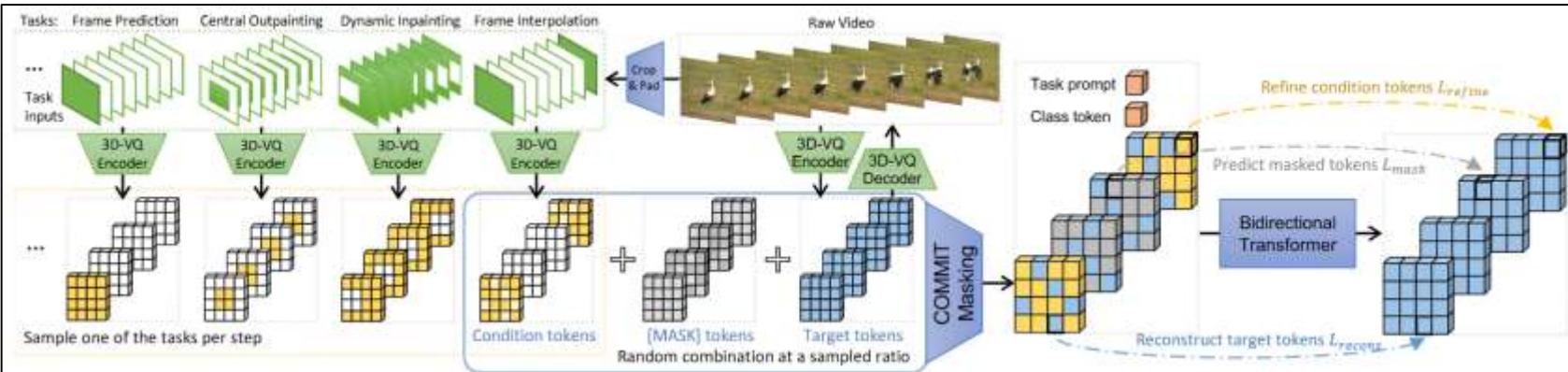
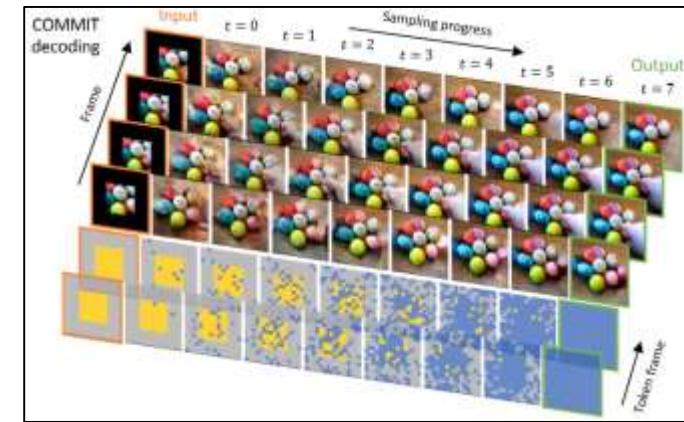
- **Condition Tokens:** Partially masked video is fed into 3D-VQ encoder to obtain tokens, and these are propagated to BERT. Allows us to model the train-test mismatch, since the encoder now has to deal with sparse video inputs as conditioning. Intuitively, refining known condition tokens ensures that it can process noisy/incomplete inputs well.
- **Mask Tokens:** Want BERT to predict these regions.
- **Target Tokens:** We feed some target tokens; by feeding in GT this helps ensure all tokens are well modeled.

- Tasks includes frame prediction, interpolation, inpainting, etc

Decoding/inference is done non-autoregressively

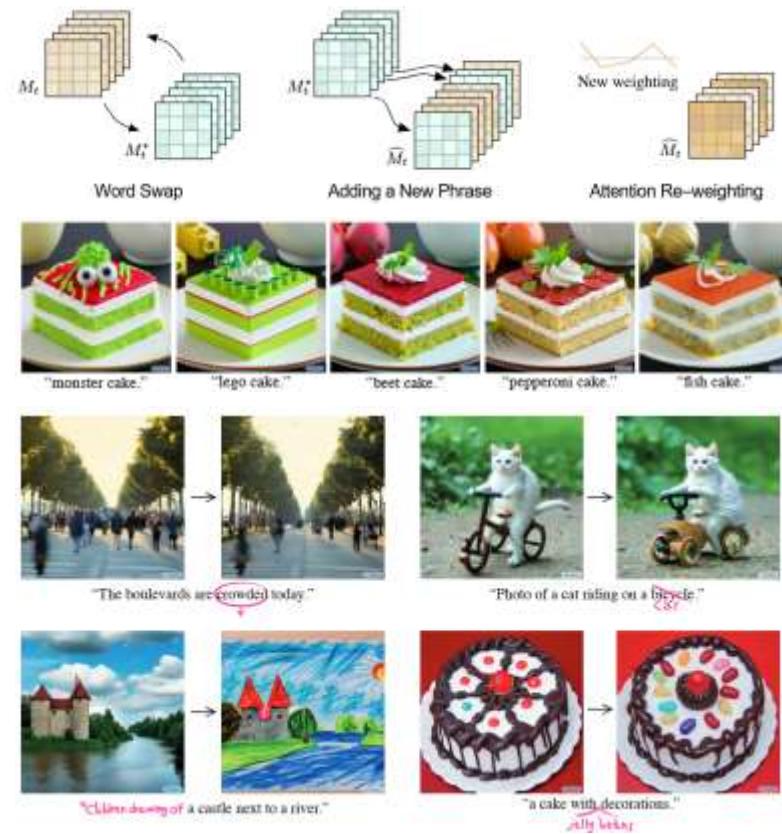
- we start with token sequence with condition tokens embedded and others masked. Then, iteratively: predict masked tokens, replace a fraction of them based on high confidence.
- Entire video generated in the first iteration, and then in subsequent iterations we continue to refine every frame in the video simultaneously in parallel.

Note: no guarantee that conditioned pixels will remain the same after generation, but this can be a feature; allows liberty for consistency and flexibility



How does prompt-to-prompt editing work for diffusion models?

- Allows fine-grained, targeted editing of images generated by text-to-image diffusion models by using only textual prompts; no masks
 - Crucially, properly maintains the image structure and content of unchanged words
- Based on the Imagen diffusion architecture, which is U-Shaped and conditions on entire sequences of text embeddings via cross attention
- Does not require any retraining; an inference-only modification that modifies/injects the cross-attention map. First seed is fixed, then procedure is to at each diffusion timestep:
 - Run denoising through original prompt, and save cross-attention map
 - Run denoising through target prompt, and save cross-attention map
 - Composite a new cross attention map that preserves the structure of the original prompt, but selectively changes it using the target cross-attention map
 - Run denoising with the composite cross-attention map, and use this for the next iteration
- Tasks enabled:
 - Word Swap:** Replaces a word in the prompt
 - Adding New Phrase:** Append tokens to the prompt
 - Attention Re-Weighing:** Scale influence of a token



How does DALL-E 2 work?

- First, we train a CLIP model on noisy (diffusion gaussian noise) images with their text pairs
 - This helps distill images into their base semantic meaning
 - This multimodal embedding space also allows use to do interpolation & vector arithmetic
- Second, we train a “diffusion prior model”: predict CLIP image embedding given CLIP text embeddings
 - This is a decoder-only transformer, that operates with a causal attention mask on the tokenized text, CLIP text embedding, diffusion timestep, and noised image CLIP embedding
 - Intuitively, the image embedding can be more appropriate compared to directly using the text embedding, because the text/image “shared” embedding space is not perfectly aligned
- Third, train a classifier-free conditional U-Net diffusion model using the diffusion prior
- At inference, given a text prompt we embed it, predict the image embedding, and then decode using diffusion
- Can do many things with this setup:
 - Generate images from text
 - Encode an image, then decode with diffusion to generate variations; 2 image embeddings can be interpolated as well
 - Encode 2 text samples and diffuse their interpolation



Figure 3: Variations of an input image by encoding with CLIP and then decoding with a diffusion model. The variations preserve both semantic information like presence of a clock in the painting and the overlapping strokes in the logo, as well as synthetic elements like the saturation in the painting and the color gradients in the logo, while varying the non-essential details.



a photo of a landscape in winter → a photo of a landscape in fall

How does Lumiere work?

- Lumiere is a space-time aware diffusion model that generates entire video clips in a single pass (and then iteratively refines it)
- Uses a Space-Time U-Net for diffusion
- Loads weights from a pretrained Imagen network (U-net based diffusion model trained for images; conditions on text using cross-attention & does classifier-free guidance)
 - We “inflate” the imagen weights by adapting them to video & incorporating temporal components
 - We add 1D convolutions applied across the time dimension; spatial attention becomes spatiotemporal attention by attending temporally
 - The weights of the pretrained network are fixed; only added temporal components are trained.
- Additionally can do image-to-video (by conditioning on the first frame) or text-to-video (with text conditioning)

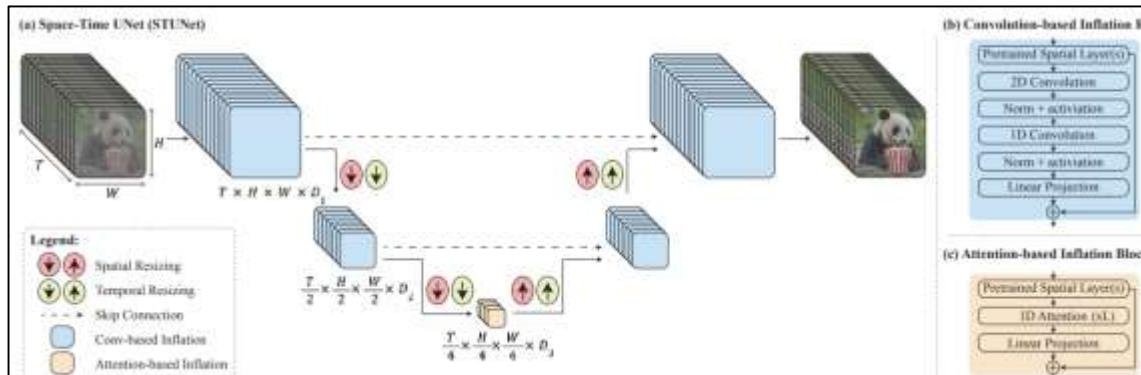


Figure 4: **STUNet architecture.** We “inflate” a pre-trained T2I U-Net architecture (Ho et al., 2022a) into a Space-Time UNet (STUNet) that down- and up-sample the video in both space and time. (a) An illustration of STUNet’s activation maps; color indicates features resulting from different temporal modules; (b) *Convolution-based* blocks which consist of pre-trained T2I layers followed by a factorized space-time convolution, and (c) *Attention-based* blocks at the coarsest U-Net level in which the pre-trained T2I layers are followed by temporal attention. Since the video representation is compressed at the coarsest level, we stack several temporal attention layers with limited computational overhead. See Sec. 3.1 for details.

How does Zero123 work?

- Novel view synthesis by finetuning a pretrained Stable Diffusion model to take in camera pose conditioning
 - Stable diffusion is just a latent-space diffusion method with a U-Net architecture
- Trained on synthetic objaverse dataset (10M 3D objects)
- Able to generalize well on real images, due to the internet-scale trained stable diffusion training, as well as high-quality & diverse synthetic data
- Zero123++ improves method further via tiling and other tricks
 - Tiling involves generating 6 predefined views of an object simultaneously, all within one image to improve 3D consistency

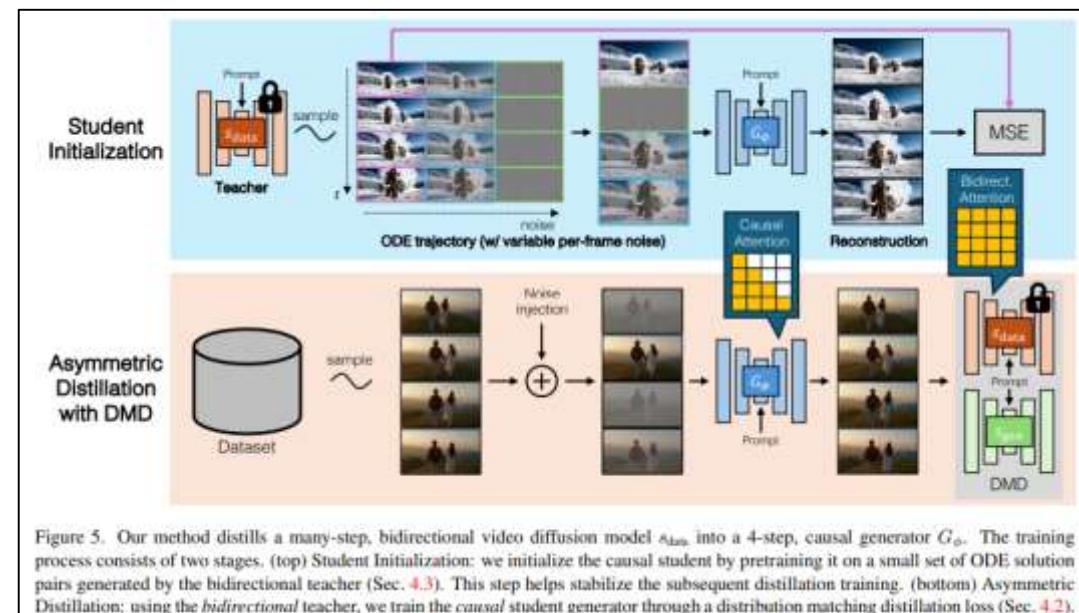


At a high level, what are the trade-offs between diffusion vs VQ for generation?

- Diffusion models (eg Stable Diffusion, Imagen):
 - are generally more detailed/photorealistic
 - slower & need to take many steps
 - Very good editability & controllability
- VQ based models (eg DALLE-1, VQ-VQE, VQGAN, MAGVIT):
 - Faster
 - Very naturally fits LLM-framework / transformers

How does CausVid turn a many-step bidirectional diffusion model into a few-step causal one?

- In this setup, small chunks are generated sequentially, using causal attention masking across chunks and bidirectional attention within chunks, for 10 fps streaming
- Student architecture based on DiT with latent space from a 3D VAE
- Distills a 50 step bidirectional teacher to a 4-step autoregressive student (asymmetric; bidirectional to causal)
 - First, student model is initialized by training it to denoise teacher-generated videos. Student operates in one step. supervised by generating teacher videos -> adding noise (student input) -> having teacher denoise (student gt). This is supposed to be an easier task that converges faster.
 - In second stage, student denoises image frames causally. Supervised with “Distribution Matching Distillation” (DMD) loss, which is basically a reverse KL divergence to promote realism. Computation involves the score function, which a DMD network “ s_{gen} ” predicts.
- KV cache used at runtime for efficiency

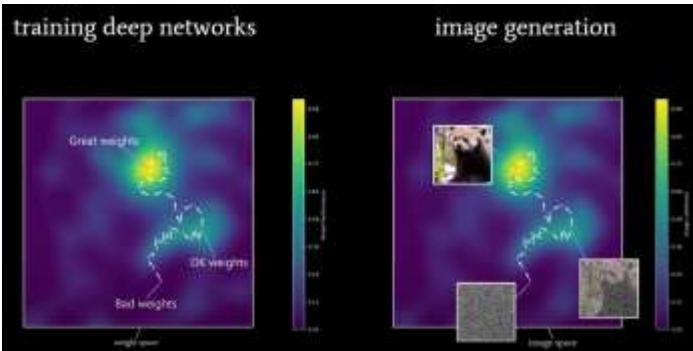


How can RL be used to improve diffusion?

- Can adapt text-to-image diffusion models to objectives that are difficult to express via prompting, such as aesthetic quality (from human feedback)
- Goal is to train models to satisfy objectives directly, instead of matching a data distribution
 - In principle, you could also reshape your data distribution to achieve the same effect in a supervised learning regime. However, this is usually limited, inefficient, and inflexible
 - RL is goal-driven, rather than gt-data-driven – gt data can have its limits. Hard to source data that matches preferences
- Treats each denoising step as a Markov decision process, with actions optimizing a downstream reward
- No need for explicit labels (eg images); just need rewards from generated samples
- Can handle sparse feedback for methods that are based on sequences of action over time
 - Supervised models by contrast focus on immediate correctness per example

What's the connection between Diffusion models and Langevin sampling?

- Langevin Sampling is a method to sample from potentially complex, high dimensional probability distributions by using only the **score function** – the gradient of the log-likelihood
 - Intuitively, the score represents the direction towards greater likelihood; the log scales its magnitude, so low probability regions are amplified.
 - In the algorithm, we add noise to the score; intuitively, this adds randomness & prevents us from getting stuck in a local max
- It can be shown that diffusion is doing the same thing & simulates Langevin dynamics; the diffusion objective is equivalent to predicting the score
 - In this interpretation, diffusion is not just denoising; this makes sense, since we're actually adding noise & simulating Langevin dynamics. The noise adds some high-frequency details, and the diffusion model selectively removes some components to create structures
 - Another way to view this is that diffusion models are kind of doing test-time SGD, where the gradient is approximated by the NN. Just like how minibatches have randomness to take us out of local optima, diffusion models use the normal distribution to do so.



Algorithm: Langevin sampling

$$x_0 \sim \mathcal{N}(0, 1)$$

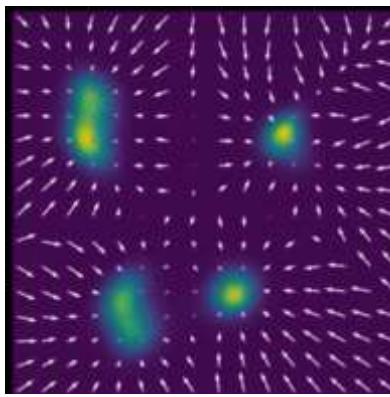
for $t \leftarrow 0$ to 1000

$$z_t \sim \mathcal{N}(0, 1)$$

$$x_{t+1} = x_t + \frac{\epsilon}{2} F(x_t) + \sqrt{2\epsilon} z_t$$

return x_{1000}

$$F(x_t) = \nabla_x \log p(x_t)$$



Example score values for 2D distribution

Algorithm: Minibatch Gradient Descent

$$\mathbf{w}_0 = \mathcal{N}(0, \mathbf{I}); t = 0$$

until $\mathcal{L}(\mathbf{w})$ is low enough

sample batch of examples $\mathbf{B} = \{\mathbf{x}_i\}_{i=1}^b \sim p_{\text{data}}$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \epsilon \nabla_{\mathbf{w}} \mathcal{L}_{\mathbf{B}}(\mathbf{w}_t)$$

return last \mathbf{w}_t

Algorithm: Generate Image

$$\mathbf{x}_0 = \mathcal{N}(0, \mathbf{I}); t = 0$$

until $p_{\text{images}}(\mathbf{x}_t)$ is high enough

$$\mathbf{z}_t \sim \mathcal{N}(0, 1)$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \frac{\epsilon}{2} F_{\text{images}}(\mathbf{x}_t) + \sqrt{2\epsilon} \mathbf{z}_t$$

return last \mathbf{x}_t

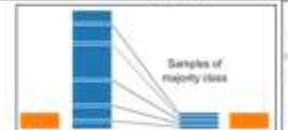
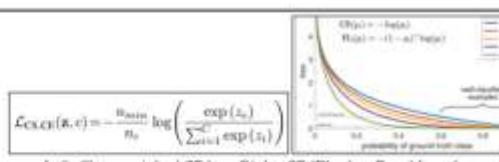
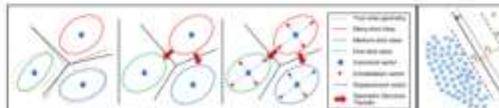
Question Goes Here

- Answer goes here

Data Imbalance

What are the methods to deal with data imbalance?

Summary of Methods to Improve Learning with Data Imbalance

Name	Description	Remarks	Figure
Data Collection	<ul style="list-style-type: none"> Collect more data to improve class balance 	Best performance, but costs the most.	
Oversampling	<ul style="list-style-type: none"> Duplicate data copies of the minority class(es), to improve balance. 	Can cause overfitting and poor generalization.	
Undersampling	<ul style="list-style-type: none"> Remove some data from the majority class(es), to improve balance. For example, this can be random, or you can choose the “furthest neighbor points”. One technique is NearMiss, which tries to only keep points from majority class which are close to the decision boundary (i.e. avg distance to the minority class are smallest) 	A good choice when you have a ton of data from majority class. Obviously, also improves training time. May discard valuable data. One option is to use an ensemble (similar to bagging).	 Left: Undersampling. Right: NearMiss.
Data Augmentation	<ul style="list-style-type: none"> Create synthetic or augmented data for minority class(es). This can be through data augmentation or even GANs. In SMOTE, you randomly select a close neighbor of a minority class point, and create a random point within their line segment. This requires a well-defined semantic space. In Mixup, new images are created by interpolating pairs images & their labels. 		 Left: SMOTE. Right: Mixup.
Reweighting	<ul style="list-style-type: none"> Provide a weight for each class in the cost function (e.g. cross entropy), to provide more emphasis on minority classes. Focal Loss modifies the cross-entropy function curvature so easy, confident predictions have less loss weight. Widely used for detection, where there are many “easy” BG negative bboxes and only a few “hard” FG positive to pay attention to. 	A ICLR 21 paper states that although mathematically equivalent, resampling is generally better than reweighting because of factors from SGD optimization.	 Left: Class-weighted CE loss. Right: CE (Blue) vs Focal Loss (green)
Two-Stage Training	<ul style="list-style-type: none"> Consists of an imbalanced training stage first (to learn good representations) and a balanced fine-tuning stage (using another method, to re-balance). 		
Metric Learning	<ul style="list-style-type: none"> Want to enforce good margins in an embedding space (especially for few-shot examples), in a class-balanced sensitive way. GistNet tries to transfer the decision boundary shape of many-shot classes to the medium/few-shot classes. Large-margin softmax tries to make margins larger, useful for generalizability of few/medium shot. Label Distribution Aware Margin (LDAM) is an improvement on this to consider class data size. 		 Left: GistNet. Right: LDAM.

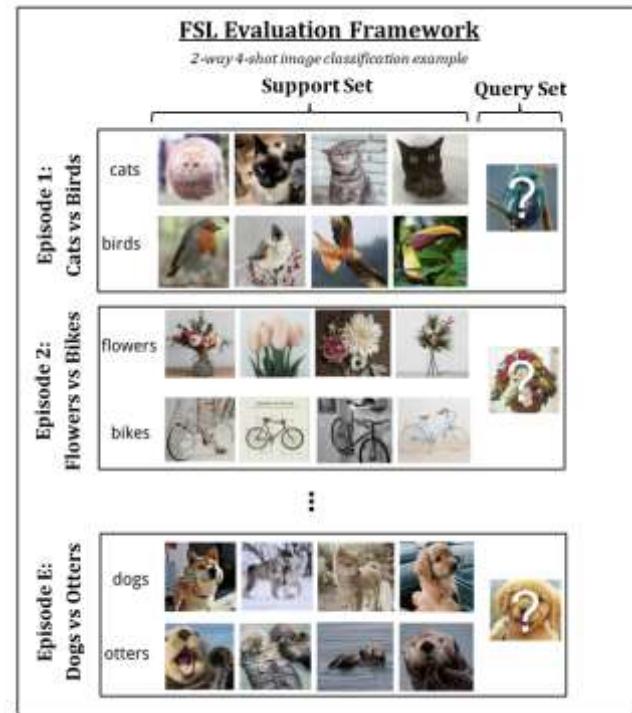
Question Goes Here

- Answer goes here

Few-Shot Learning

At a high level, what is the few-shot learning problem?

- In **few-shot learning (FSL)**, the goal is to learn tasks with only a few training examples.
- We usually evaluate a FSL algorithm's capability for data-efficient learning with multiple **N-way, K-shot episodes**.
- Here, “**way**” = “**classes**” and “**shot**” = “**training examples per class**”. An **episode** can be thought of as an individually defined “classification task”. In detail, each episode will have:
 - A **support set** containing K training samples & labels per class (so, samples total), which the FSL algorithm trains on
 - Commonly, K is 1 or 5; N is usually 5 but can be over 100
 - A **query set**, which we test the FSL algorithm on after it trains on the support set
 - Generally at least NK in size
 - For testing purposes we have labels, but the FSL algorithm doesn't train on them; it's purely a forward pass
- In general, FSL episodes are easier with fewer classes (lower way) and more examples per class (higher shot).
- FSL episodes are harder with more classes and less examples per class.
- A variant of few-shot learning is **zero-shot learning**, which learns solely on the category name, textual descriptions, or attributes.
- Note that many few-shot methods are more tailored to that problem, and may be insufficient to tackle the more general **long-tailed problem**.



What is Meta Learning?

- In meta learning, the goal is to optimize a model to perform well across a distribution of tasks, instead of just one task. It should quickly adapt to new tasks using only a small amount of data – goal is to “learn to learn”

Optimization-Based

- Here, the goal is to learn an initialization that allows for fast & effective adaptation via fine-tuning on a few examples
- Involves training a model so it can adapt quickly with only a few gradient steps
- Examples include MAML (Model Agnostic Meta Learning) and REPTILE
 - In MAML, we compute preliminary updated parameters on multiple tasks, then make a meta-update by evaluating performance on new/held-out tasks (requiring second order derivatives)
 - Intuitively, we gather what a small update looks like for multiple tasks, then adjust the parameters so that future updates work even better.
 - You want to “learn how to learn” efficiently, rather than learning for a specific downstream task; the goal is to learn something of a higher (meta) order, not necessarily the task directly
 - For example, humans learn certain universally applicable mental tricks/techniques that allow us to learn novel tasks/concepts rapidly

Model-Based

- HyperNetworks directly regress/predict the parameters of another network
- Have been shown to work for very large networks (hundreds of millions)
- Note, HyperNetworks are conceptually similar to LoRA

Applications:

- Few-Shot Learning
- Quick generation of NeRFs
- Personalized/dynamic networks
- Neural architecture search

Limitations:

- Can be prone to overfitting and bad generalization
- Can be hard to train/optimize

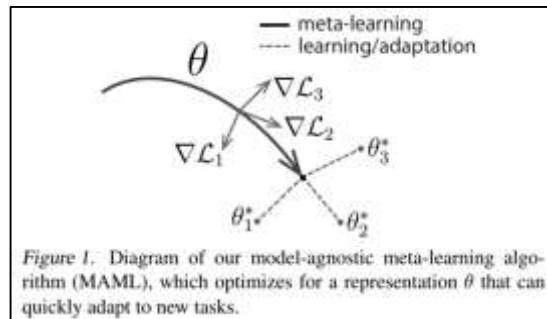


Figure 1. Diagram of our model-agnostic meta-learning algorithm (MAML), which optimizes for a representation θ that can quickly adapt to new tasks.

Algorithm 1 Model-Agnostic Meta-Learning

```
Require:  $p(\mathcal{T})$ : distribution over tasks
Require:  $\alpha, \beta$ : step size hyperparameters
1: randomly initialize  $\theta$ 
2: while not done do
3:   Sample batch of tasks  $\mathcal{T}_i \sim p(\mathcal{T})$ 
4:   for all  $\mathcal{T}_i$  do
5:     Evaluate  $\nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$  with respect to  $K$  examples
6:     Compute adapted parameters with gradient descent:  $\theta_i^r = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$ 
7:   end for
8:   Update  $\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta_i^r})$ 
9: end while
```

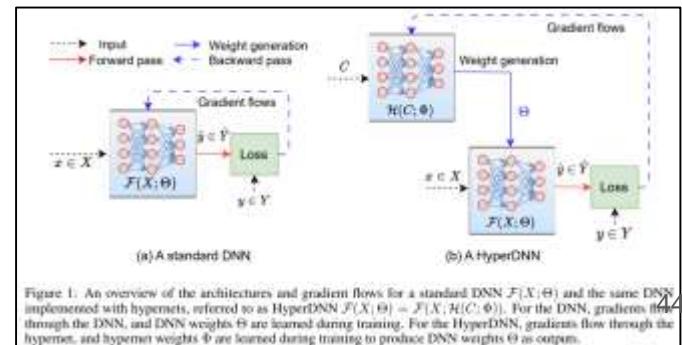
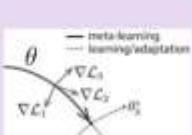
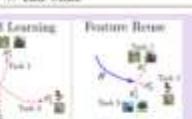


Figure 1. An overview of the architectures and gradient flows for a standard DNN $F(X; \Theta)$ and the same DNN implemented with hypernets, referred to as HyperDNN $F(X; \Theta) = F(X; H(C; \Phi))$. For a standard DNN, gradients flow through the DNN, and DNN weights Θ are learned during training. For the HyperDNN, gradients flow through the hypernet, and hypernet weights Φ are learned during training to produce DNN weights Θ as components.

Summarize the following optimization meta learning approaches for few-shot: MAML, Reptile, LEO.

Summary of Methods for Few-Shot Learning			
Metric Meta Learning Optimization Meta Learning Transfer Learning			
Name	Category	Description	Figure
MAML: Model Agnostic Meta Learning (ICML 17, > 5000 cits)	Optimization Meta Learning	<ul style="list-style-type: none"> Goal: Learn good NN initialization that's easy to fine-tune Parameters of NN explicitly trained so a few gradient steps with a few training data from a new task will produce good generalization performance Very literal approach. Repeatedly 1) samples tasks 2) for each task, takes a gradient step on individual copy of weights, and 3) meta-updates weights, backpropagating through step 2 so that weights would minimize loss after step 2's gradient step: $\min_{\theta} \sum_{T_i \sim p(T)} \mathcal{L}_{T_i}(f_{\theta} - \alpha \nabla_{\theta} \mathcal{L}_{T_i}(f_{\theta}))$ The 3rd step involves taking the gradient of a gradient (2nd order derivative), unrolling through the computational graph. However, it's shown that First-Order MAML is almost as good, by treating the weights from step 2 as constants. There's evidence MAML mostly reuses/learns high-quality features, rather than "learning how to learn rapidly". 	 <pre> Algorithm 1 Model-Agnostic Meta-Learning Require: $p(T)$: distribution over tasks Require: α, β: step size hyperparameters 1: randomly initialize θ 2: while not done do 3: Sample batch of tasks $T_i \sim p(T)$ 4: for all T_i do 5: Evaluate $\nabla_{\theta} \mathcal{L}_{T_i}(f_{\theta})$ with respect to K examples 6: Compute adapted parameters with gradient descent: $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{T_i}(f_{\theta})$ 7: end for <i>Note: the meta-update is using different set of data</i> 8: Update $\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_{T_i \sim p(T)} \mathcal{L}_{T_i}(f_{\theta'})$ 9: end while </pre> 
Reptile (OpenAI, 18, >1000 citations)	Optimization Meta Learning	<ul style="list-style-type: none"> Like MAML, goal is to find a parameter configuration which is close to the optimal parameter manifolds of all tasks Repeatedly 1) samples tasks, 2) separately trains parameters on each task for multiple gradient steps, and 3) meta-updates the original parameters towards the averaged new parameters Multiple SGD steps makes it different than joint, mixture training They theoretically show that this generalizes MAML and also implicitly involves 2nd derivatives, but is more efficient 	 <pre> Algorithm 2 Reptile, batched version Initialize θ for iteration = 1, 2, ... do Sample tasks $\tau_1, \tau_2, \dots, \tau_n$ for $i = 1, 2, \dots, n$ do Compute $W_i = \text{SGD}(L_{\tau_i}, \theta, k)$ end for Update $\theta \leftarrow \theta + \beta \frac{1}{n} \sum_{i=1}^n (W_i - \theta)$ end for </pre>
LEO: Meta-Learning w/ Latent Embedding Optimization (DeepMind, ICLR 19, >700 cits)	Optimization Meta Learning	<ul style="list-style-type: none"> While MAML operates directly in a high dimensional parameter space, LEO performs meta-learning within a low-dimensional latent space Thus, parameters are generated/predicted 	

At a high level, what is metric learning? How can it be applied to few-shot?

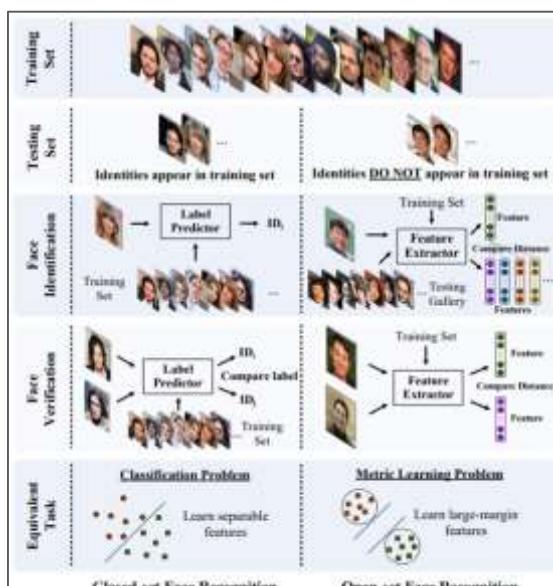
- Metric Learning (i.e. similarity learning) aims to measure how related two objects are, in a lower-dimensional, semantic manifold
- It can be learned with a regression or classification based framework, using labeled pairs of similar or dissimilar images
- Alternatively, triplets of can be used where is known to be more similar to than .
 - This is weaker supervision than regression or classification, since instead of providing an exact measure of similarity (usually one-hot), one only has to provide the relative order of similarity. This also allows the use of data augmentation, for self supervised learning (e.g. simCLR)
 - However, mining for the “hard” examples for to train on is important, and nontrivial
- The embeddings can readily be used for one-shot learning tasks (e.g. face verification databases), and support easy additions at test-time without retraining. You just need to embed the support elements, and use K nearest neighbors to classify.
- Besides FSL, other applications for metric learning include ranking (for retrieval or recommender systems), large-scale classification, or clustering. It is also common to use the learned representation from self-supervised metric learning for a wide range of downstream tasks via transfer learning.



Explain how few-shot can tackled with transfer learning, optimization meta learning, and metric meta learning.

FSL via Transfer Learning

- Pre-train on a large, separate dataset with many classes and instances per class (where classes distinct from those seen during evaluation), e.g. ImageNet
- Training can be fully supervised or self-supervised
- During evaluation, **transfer the learned representations**, for example by
 - Fine-tuning on the NK episode training examples, or
 - Using it as a **fixed feature extractor** and learning a linear classifier (e.g. softmax) on top
- Intuitively, at the transfer learning step, you're not “teaching it to reason”, just “**recontextualizing**” with the few examples that you have
- Most basic approach, but surprisingly, **very competitive** with more complicated meta learning based methods



FSL via Optimization-Based Meta Learning

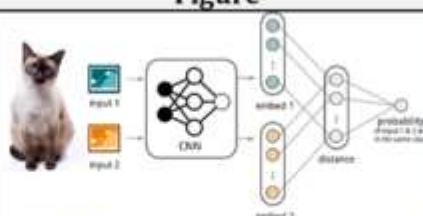
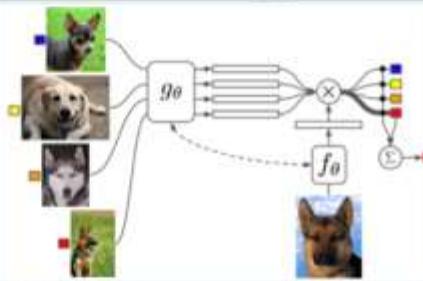
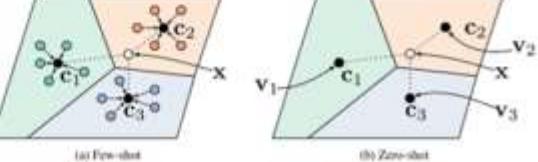
- Insight: **Replicate the FSL episode-based evaluation procedure during training**
- Intuitively, optimize for “fast-weights/parameters” in network which can be adapted in a data-efficient manner to episodes without overfitting, using only a few gradient steps
- However, some believe that this is more less due to **reusing/learning high-quality features**, rather than “rapid learning of each task”.

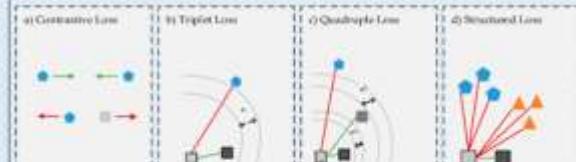
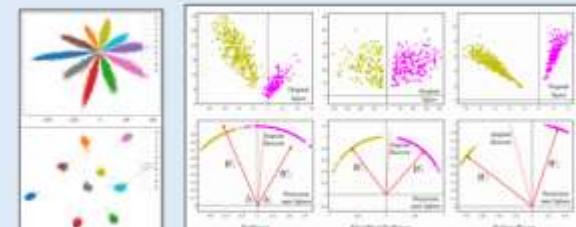
FSL via Metric-Based Meta Learning

- Based on metric learning; want to learn a mapping from images to an **embedding space**, and use a (possibly learned) **similarity function** between points
 - Goal: images from the same category are closed together and from different categories are far apart.
 - Ideally, the mapping will hold true for unseen categories; during evaluation, we embed the support set and use nearest-neighbors to classify
 - Essentially assumes that can be modeled by the form , where k is a similarity function (e.g. embedding Euclidian distance)
- “Older”, more classical approach to FSL
- Especially useful for tasks like **face verification**, where you have a **frequently changing database** of objects (pictures of people) that you want to recognize, but only a few data samples for each object. Unlike Meta-learning approaches, you would **not need to retrain** when you add or remove people from the database.
- Can involve contrastive methods (where you have positive and negative sampling), or methods based on softmax (which is a generally easier to apply)

Summary of Methods for Few-Shot Learning

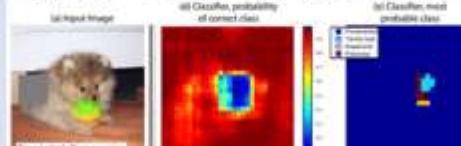
[Metric Meta Learning](#) | [Optimization Meta Learning](#) | [Transfer Learning](#)

Name	Category	Description	Figure
Siamese NNs (ICML 15, >2500 cits)	Metric Meta Learning	<ul style="list-style-type: none"> Feeds pairs of images into a CNN to get two embeddings The L1 distance between embeddings fed into a linear layer to predict similarity probability score of if they are same class Trained with binary cross entropy loss At test time, a query image is classified as the class of the support image with highest class-similarity probability 	
Matching Networks (Google, NIPS 16, >3700 cits)	Metric Meta Learning	<ul style="list-style-type: none"> Given support $S = \{x_i, y_i\}_{i=1}^k$ where y_i are one-hot vectors, learns embedding functions f, g such that for query x we have $P(y x, S) = \sum_{i=1}^k a(x, x_i)y_i$ a is an "attention kernel" weighing one-hot vectors, based on cosine similarity softmax: $a(x, x_i) = \frac{\exp(\text{cosine}(f(x), g(x_i)))}{\sum_{j=1}^k \exp(\text{cosine}(f(x), g(x_j)))}$ Training mirrors evaluation: during training, loss explicitly minimized over small episodic samples of support & queries Potentially, f, g can be the same feature extractor, but they instead propose a more complicated bidirectional LSTM 	
Relation Networks (CVPR 18, >2000 cits)	Metric Meta Learning	<ul style="list-style-type: none"> Similar to Siamese NNs gets pairs of embeddings, but concatenates them and feeds into a NN to predict relation score instead of using L1 distance Also, instead of cross entropy, takes a regression approach and uses MSE with a one-hot vector 	
Prototypical Networks (NIPS 17, >3500 cits)	Metric Meta Learning	<ul style="list-style-type: none"> Embeds all support images, and defines a prototype embedding for each class mean Given query x, distribution over classes given as softmax over negative squared Euclidean distance It's argued that this inductive bias reduces overfitting Training mirrors evaluation: episodes sampled during training Trained with cross entropy 	

Summary of Methods for Few-Shot Learning			
Metric Meta Learning Optimization Meta Learning Transfer Learning			
Name	Category	Description	Figure
Contrastive Metric Learning: <u>Triplet Loss</u> (FaceNet, CVPR 15, >10K cits), Quadruplet Loss (CVPR 17, 1K), Structured Loss (CVPR 16, 1.2K), N-Pair Loss, Magnet Loss, Clustering loss, MoCo (CVPR20, 21 K)	Metric Meta Learning	<ul style="list-style-type: none"> Goal: pull together embeddings with same label, and push away those with dissimilar labels, with explicit pos/neg samples Classical contrastive loss uses pairs of samples and margin α: $\mathcal{L} = \sum_{y_1=y_2} D^2(f(x_1), f(x_2)) + \sum_{y_1 \neq y_2} \max(0, D^2(\alpha - f(x_1), f(x_2)))$ Triplet loss uses A and P with same class, and N of different class: $\mathcal{L}(A, P, N) = \max (\ f(A) - f(P)\ ^2 - \ f(A) - f(N)\ ^2 + \alpha, 0)$ Structured and N-Pair Losses tries to improve sampling with all samples in a batch, instead of ignoring them and only using one Magnet and Clustering loss are based on the dataset as a whole, but is complex and hard to scale. Often paired with heavy data aug, large batch size, and/or memory bank to efficiently reuse embeddings from the prev. iteration Main difficulty of contrastive learning is sampling comprehensively, mining the ‘hard negatives’. Also, batch-sample based learning enforces local margins, but a global margin for all labels is hard to achieve in practice. Thus, non-contrastive losses are often easier to use. Useful for unsupervised representation learning (e.g. simCLR) 	 
Margin Non-Contrastive Metric Learning: <u>Center Loss</u> (ECCV16, 27K cit.), <u>LargeMargin Softmax</u> (IJCAI16), <u>SphereFace</u> (CVPR 17, 2K cit.), <u>CosFace</u> (CVPR 18, 1300 cit.), <u>ArcFace</u> (CVPR 19, 2300 cit.)	Metric Meta Learning	<ul style="list-style-type: none"> Goal: pull together embeddings with same label, and push away those with dissimilar labels, without explicit pos/neg samples Center loss adds an additional term to softmax loss, to pull each batch of N features towards its (jointly trained) class center Large-Margin Softmax notes that for sample (x_i, y_i), the softmax loss for a FC layer exponentiates $W_{y_i}^T x_i = \ W_{y_i}\ \ x_i\ \cos(\theta_{y_i})$ in the numerator (selects the g.t. class vector in W and takes dot prod). Then, a larger margin is explicitly enforced by enforcing larger θ. SphereFace improves margin by enforcing the centers to be on a hypersphere w/ normalization, improving the decision boundary Generally only applicable when supervised. For unsupervised settings or when you have a lot of out-of-distribution samples, contrastive learning may be better. General tradeoff with within-domain accuracy (smaller margins, standard classification) vs out-of-distribution generalization (larger margins, longtail/fewshot/openset verif.) 	 <p>Left: Softmax vs Center Loss. Right: Softmax vs SphereFace</p>

Explainable Artificial AI (XAI)

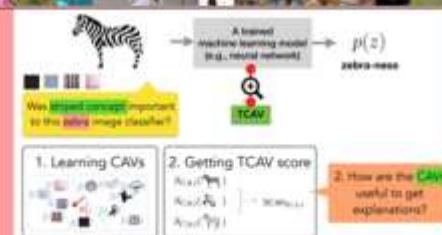
At a high level, describe the following basic XAI methods: Layer weight/activation visualization, maximally activating retrievals, and occlusion maps.

Show Layer Activations	Show forward pass activations (as an image) for an input	Check for localized regions in later layers. If an activation map is zero for many different inputs, this can indicate dead filters	
Show Layer Weights	Show convolutional filter weights (as an image) at layer of trained network	Usually, well-trained networks should display nice and smooth filters without noisy patterns.	
Maximally Activating Retrievals	Take large dataset of images, feed them through the network. Keep images maximally activate some neuron	Gives an understanding of what specific neurons are looking for in its receptive field.	
Occlusion Maps	Plot the probability of the correct class as a function of the position of an occluder object (e.g. a gray square), as a 2D heatmap.	Helps understand where a CNN is looking for its prediction, e.g., checks if it's looking a context clues/backgrounds rather than the main object	

At a high level, describe the following backprop-based XAI methods: Guided backprop, DeepDream, CAM, Grad-CAM.

Guided Backpropagation	Backpropagate prediction all way back to the input while zeroing-out negative gradients (only accounts for positive contributions, reducing noise)	Fine-grained visualization, in input space to help understand the important parts of the input in making prediction.		
DeepDream	Apply gradient ascent on an image or noise, freezing network weights, to optimize the probability classification of a certain class or (classes).	Shows the salient features used for the task, through generation. Can reveal biases/correlations in network.		Towers & Pagodas Buildings Birds & Insects Dumbbells
Class Activation Mappings (CAM)	In architecture, number of feature maps in the last conv layer equals the number of classes. Maps are global averaged pooled and fed directly into softmax. This forces a heatmap-like representation. CAM is the weighted average of those feature maps.	Shows the most important areas in input image used for the prediction.		
Grad-CAM	Weighted sum of feature maps activations; weights are from the gradients with respect to a target class, averaged per feature map.	Shows the most important areas in input image used for the prediction.		(a) Original Image (b) Guided Backprop 'Cat' (c) Grad-CAM 'Cat' (d) DGrad-CAM 'Cat'

At a high level, describe the following numerical-based XAI methods: T-SNE, TCAV.

T-SNE	A method for non-linear dimensionality reduction, based on neighbor distances. Can be applied to feature embeddings for network inputs.	Can check for clusters within data.		
Testing with Concept Activation Vectors (TCAV)	Given images of concepts and random images, trains linear classifier for vector orthogonal to hyperplane decision boundary.	Flexible way to quantitatively check how much abstract concepts (e.g. gender, race) are related to a trained-CNN prediction. Concept need not be seen during training.		

Security / Adversarial Attacks

Explain what white box and black box attacks are, and how they're typically implemented.

- **White box** attacks have access to neural network weights
 - You can use **gradient ascent** to backpropagate back to the original input image
- **Black box** attacks do not have access to the weights of the neural network.
 - A **substitute model** can be used to emulate the original model, using a transferable attack strategy
 - A **query feedback** mechanism can be used, along with random/grid/local search. You just keep modifying the input until you break the CNN and produce incorrect results. Modifications can be things like affine transform, changing hue, etc.
- Most attacks that are designed to fool a particular CNN also fool many other CNNs
- There are two theories on why adversarial attacks exist:
 - Consequence of SGD-based approximate optimization
 - Reliance on human-invisible non-robust feature patterns (i.e. overfitting to the training dataset)

Explain what real-world adversarial attacks are, and give some examples.

- Real world attacks face some additional challenges:
 - Needs to have viewpoint invariance and robustness to lighting, occlusion, etc
 - If it's on something like a street sign, cannot be on the whole image, only on the foreground obj
 - Cannot just be imperceptible; needs to be perceptible by camera
 - Needs to be robust to some fabrication error (eg by printer)
- The paper “Synthesizing Robust Adversarial Examples” shows that attacks can be 3d printed and work in the real world.
 - A 3D turtle texture is optimized to be an adversarial attack from many viewpoints and lighting conditions with differentiable rendering
- In the paper “Robust Physical-World Attacks on Deep Learning Visual Classification”, real-world attacks are performed on stop signs
 - The adversarial optimization procedure is performed over a distribution (dataset) of “realistic” stop signs, by using a hybrid real & synthetic augmentation approach.
 - A mask is used to ensure the attack pattern is on the foreground stop sign
 - There is also a loss that encourages colors/patterns to be easily printable

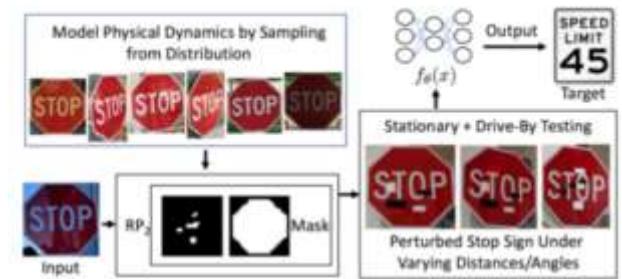


Figure 2: RP₂ pipeline overview. The input is the target Stop sign. RP₂ samples from a distribution that models physical dynamics (in this case, varying distances and angles), and uses a mask to project computed perturbations to a shape that resembles graffiti. The adversary prints out the resulting perturbations and sticks them to the target Stop sign.

What are some typical defenses for adversarial attacks?

- Adversarial training (eg training from FGSM or iterative FGSM attacks)
 - In general, any attack can be transformed into a defense in this way
- Ensemble voting
- Adversarial classifier

$$\mathbf{X}^{adv} = \mathbf{X} + \epsilon \operatorname{sign}(\nabla_{\mathbf{X}} J(\mathbf{X}, y_{true}))$$

FGSM procedure, where J is the loss function

Efficient Deep Learning

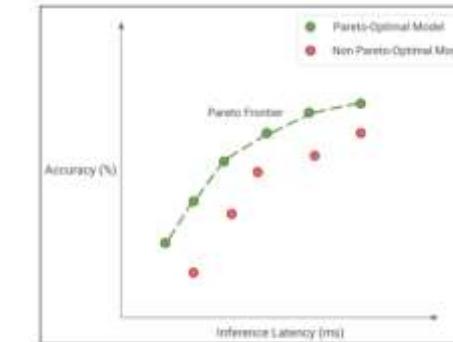
At a high level, what is the goal of efficient DL?

- The goal of **efficient deep learning** is to improve the **quality/footprint ratio** of models:

$$\text{Model Efficiency} = \frac{\text{Quality}}{\text{Footprint}}$$

Measured by metric, e.g. Accuracy, precision, recall, f1, MSE, etc.
In terms of runtime, computational hardware needed, or data.

- While usually not a priority for general ML research, it's important for:
 - Very large models (which cost substantial amounts of money)
 - Deployment at the edge (e.g. inference of NNs on smartphones)
- In general, there is a **Pareto optimal frontier** that describes optimal efficiency for intrinsic trade-offs between quality and footprint. The goal of efficient DL is to approach the frontier.
- Methods can be broken down into 5 categories:
 - Compression:** Can we compress parts of the given model graph?
 - Learning Techniques:** Can we train the model better? Techniques may involve the data and/or loss functions.
 - Architectural:** Can we use/design layers and architectures which are fundamentally more efficient?
 - Automation:** Can we search for more efficient models?
 - CS-Based:**
 - Caching:** if some inputs are more common than others. Before calling model, check cache
 - Batching:** Gather prediction requests until you have a batch, then perform them in parallel as a batch. Batch size tradeoffs between throughput and latency



Explain some compression and architectural based methods for efficient DL.

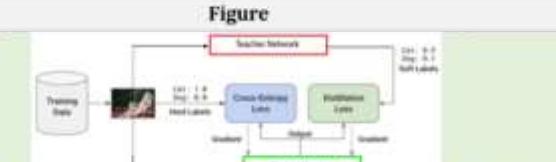
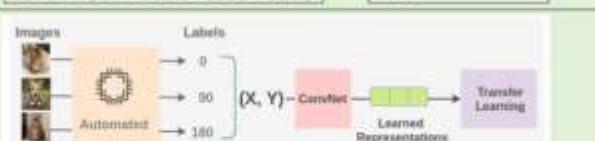
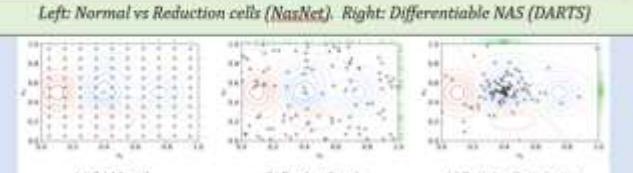
Overview of Methods for Efficient Deep Learning

Footprint Reduction | Quality Improvement | Both

Method Name	Category	Description	Figure
Numerical Quantization	Compression	<ul style="list-style-type: none"> Replaces high-bit signed floats with lower-bit signed ints, weights and/or activations Conversions to/from ints are performed using scale mapping (range-based normalization) <ul style="list-style-type: none"> E.g. suppose float range is [0,1] and we want to convert 0.5 to int8 which can represent ints in [-128,128]; the result is 64 At the extreme end, can even be binarized (BNNs) Generally, multiplications reduced precision; additions lighter and can be full precision Simplest method is to apply quantization post-training, but has more accuracy loss Can also make training quantization-aware (e.g. by "simulating" it, allowing the network to adapt and achieve better performance) Can lead to significant (~4x) reductions in model size and inference latency 	<p>(a) Inverse arithmetic only inference</p> <p>(b) Training with correlated quantization</p> <p>The graph plots Accuracy (Y-axis) against Floating Point Op count (X-axis). It shows a linear decrease in accuracy as the number of operations increases, with a shaded region indicating the variance or uncertainty in the data points.</p>
Model Pruning	Compression	<ul style="list-style-type: none"> Removes weights (i.e. set to 0, making applicable for sparsity optimization) or remove nodes/filters from a trained model (more easily supported, but can hurt effectiveness more) Goal: reduce memory and make calculations faster Several heuristics to remove weights or nodes, e.g. based on weight magnitude (replace those close to 0 with 0), activation outputs (replace "dead" neurons which always output small values), or derivatives wrt loss Network can be retrained after pruning to restore accuracy (iteratively) 	<p>before pruning</p> <p>after pruning</p> <p>pruning synapses</p> <p>pruning neurons</p> <p>Structured / Block Sparsity (4x4) in a Weight Matrix (Grey = Weight is Pruned)</p>
Depthwise Separable Convolutions	Architectural	<ul style="list-style-type: none"> Decomposes the regular convolution, for increased efficiency: <ul style="list-style-type: none"> Depthwise convolution operation preserves depth C of input. C many "groups" of $n \times n \times 1$ kernels are used, each separately applied to each channel of the input to get C channels 1×1 convolutions provide desired depth While fewer in parameters, performance lags behind large ResNet variants Originally from MobileNet paper 	<p>Normal (top) vs depthwise (bottom) conv</p>

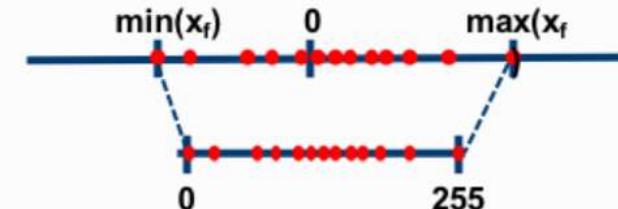


Explain some learning and automation based methods for efficient DL.

Overview of Methods for Efficient Deep Learning			
		Footprint Reduction Quality Improvement Both	
Method Name	Category	Description	
Student-Teacher Distillation	Learning Techniques	<ul style="list-style-type: none"> First train a larger "teacher" model, then train a more lightweight "student" model with both the original training data and the softmax classification scores of the teacher (with KL divergence) The availability of "dark knowledge" non-one-hot information from the teacher can improve performance for the student 	
Data Augmentation	Learning Techniques	<ul style="list-style-type: none"> Useful to improve performance and data efficiency Can use techniques like RandAugment and Mixup, or even GANs to synthesize new examples Synthetic sampling like SMOTE can allow for re-balancing to make up for dataset skewness 	
Self-Supervised Learning Transfer	Learning Techniques	<ul style="list-style-type: none"> Goal: utilize auxiliary tasks in a way where we can generate virtually unlimited labels from our existing images, to learn useful representations Resulting network can then be used for transfer learning on a domain-specific downstream task Useful to improve performance and data efficiency Examples: colorization, superresolution, inpainting, jigsaw, rotation, clustering, data-aug based contrastive learning 	
Neural Architecture Search (NAS)	Automation	<ul style="list-style-type: none"> Automates NN design by searching space of filter sizes, layer types, depth, etc for best performance and/or efficiency <ul style="list-style-type: none"> Have normal cells (in/out have same spatial size) & reduction cells Several black-box optimization approaches can be taken including grid search, random search, RL, and evolutionary algorithms Differentiable NAS (e.g. DARTS) is also possible by a continuous relaxation of the search space allowing gradient-based optimization <ul style="list-style-type: none"> Discrete operation choices replaced by mixture operations in DAG graph Bilevel optimization problem: need to optimize the parameters & architecture weights. Can approximately solve by iteratively optimizing 	 <p>Left: Normal vs Reduction cells (NasNet). Right: Differentiable NAS (DARTS)</p>
Hyperparameter Search	Automation	<ul style="list-style-type: none"> Automates tuning NN hyperparameters for better performance, without changing footprint Several black-box optimization approaches can be taken including grid search, random search, Bayesian optimization, RL, and evolutionary algorithms 	 <p>(a) Grid Search (b) Random Search (c) Bayesian Optimization</p>

How does QAT work? What types of errors can be incurred?

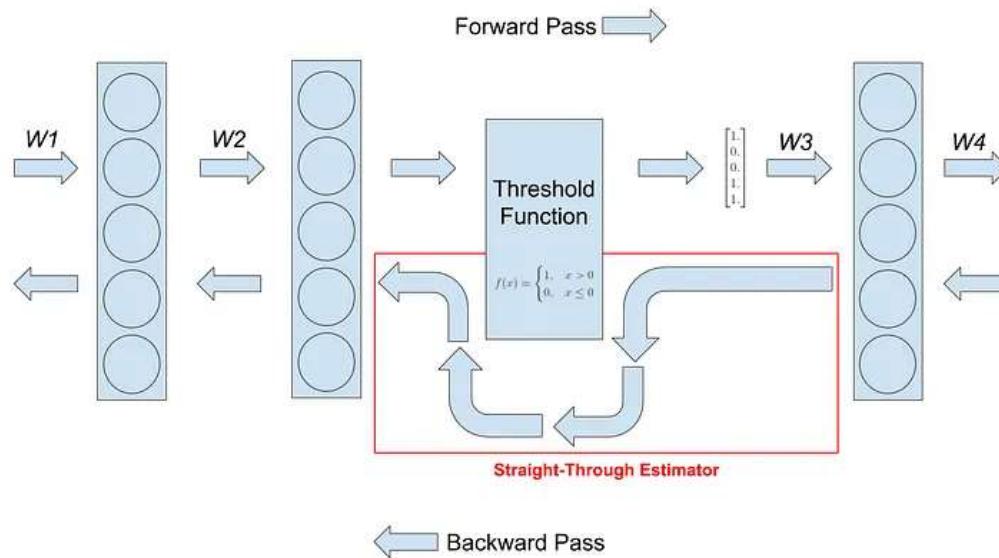
2. **Insertion of Quantization Nodes:** Modify the model by inserting "fake" quantization nodes in the network. These nodes simulate the effect of quantization (like rounding and clipping) in the forward pass during training, but the backward pass is still done using high-precision calculations. This step is crucial as it helps the model to 'learn' the quantization noise and adapt to it.
 3. **Calibration:** Run a few batches of data through the modified model to calibrate the quantization parameters. This step involves determining the range (minimum and maximum values) for activations and setting the scale and zero-point for quantization. The calibration can be done using either a subset of the training data or a separate calibration dataset.
 4. **Fine-tuning:** Retrain (fine-tune) the network with these quantization nodes in place. During this process, the model adjusts its weights to compensate for the quantization errors introduced. The goal is to minimize the impact of quantization on the model's accuracy. The training can be done for a few epochs using the same hyperparameters as the original training.
-
1. **Rounding Error:** This is the most common type of error in quantization. When converting from a high-precision format (like 32-bit floating-point) to a lower precision (like 8-bit integers), values are rounded to the nearest representable value in the lower precision format. This rounding can introduce small errors in the weights and activations of the network.
 2. **Clipping Error:** When quantizing, the range of values that can be represented in the lower precision format is often smaller than in the original format. If a value falls outside this representable range, it is clipped to the nearest bound (either the minimum or maximum representable value). This clipping can lead to significant errors, especially if the original value distribution has a long tail.



$$x_q = \text{round} \left((x_f - \underbrace{\min_{x_f}}_{zp_x}) \underbrace{\frac{2^n - 1}{\max_{x_f} - \min_{x_f}}}_{q_x} \right) = \text{round}(q_x x_f - \underbrace{\min_{x_f} q_x}_{zp_x}) = \text{round}(q_x x_f - zp_x)$$

How does the backwards pass work during QAT?

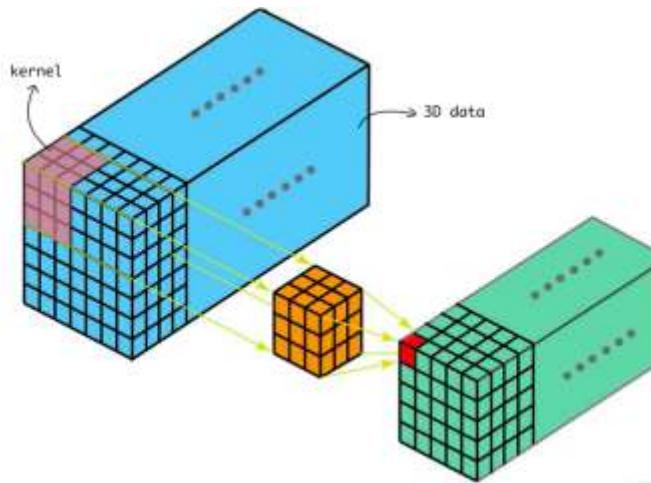
- Since a QAT model is a piecewise constant function where gradients vanishes almost everywhere (as quantization is non-differentiable), we cannot perform vanilla backprop
- Instead, we use a “straight-through estimator” (STE), ignoring the fakequant node and bypassing it as if it was an identity function.
- This “gradient” is not the actual gradient of the loss function, yet empirically/heuristically seems to work (this is still not fully understood)



3D Deep Learning

How do 3D convolutions, and deconv/transpose convs work?

- Instead of 3D $n \times n \times C$ kernels/filters, we use 4D $n \times n \times n \times C$ kernels/filters to slide over the 4D input of $W \times H \times D \times C$
 - Sliding dot product aggregates local features in 3D; width, height, and depth
- Deconvolution is similar, but uses padding in between to increase spatial size of output feature map



How do graph convolutions work?

At a high level, the idea is as follows. Suppose we have graph $G = (V, E)$, containing:

- N many D -dimensional vertices $V = \{v_i | v_i \in \mathbb{R}^D\}$, $|V| = N$
 - This can be summarized as a feature matrix $X_{N \times D}$
- Edges E
 - This can be summarized as an adjacency matrix $A_{N \times N}$

We wish to have a neural network with layers:

$$H^{l+1} = f(H^l, A)$$

The first input is $H^0 = X$, and $H^L = Z_{N \times F}$ are the output features (of dimension F) for each vertex.

A very simple way to have propagation of features from connected nodes on the graph is to have $f(H^l, A) = \sigma(AH^lW^l)$.

This is kind of like a FC layer (HW), where edges are “selected” using the adjacency matrix.

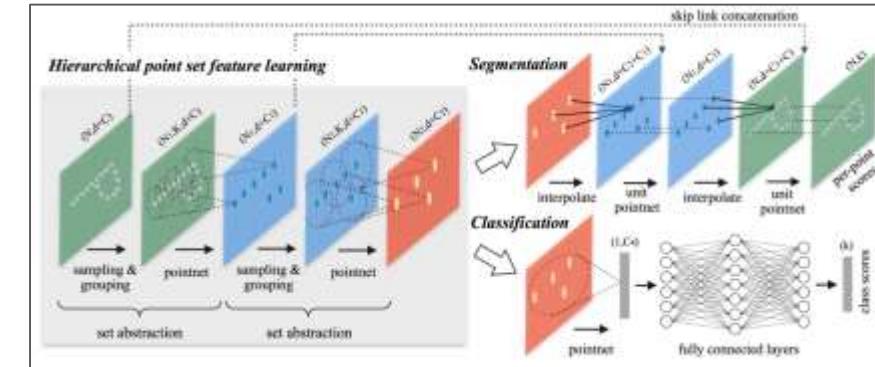
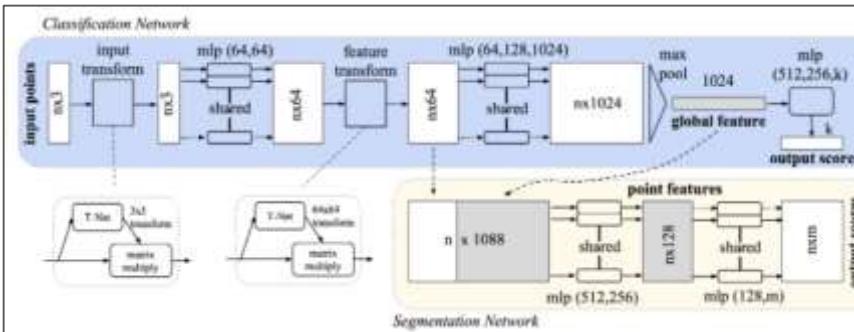
In practice, we actually use the following, which incorporates symmetric normalization and self-loops on A .

$$f(H^{(l)}, A) = \sigma\left(\hat{D}^{-\frac{1}{2}} \hat{A} \hat{D}^{-\frac{1}{2}} H^{(l)} W^{(l)}\right)$$

There are some more complicated aspects to this, involving spectral graph theory and Fourier transforms.

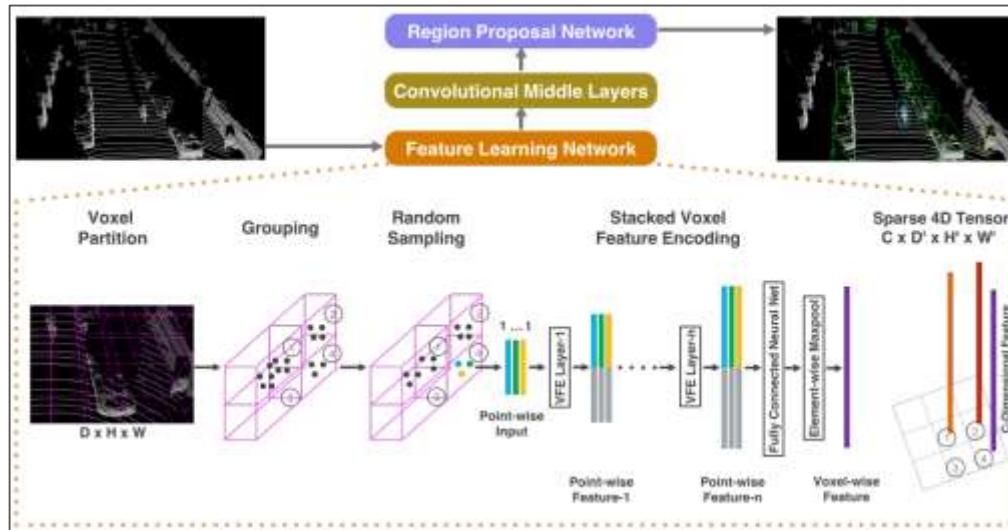
How do PointNets and PointNet++ work?

- PointNet and PointNet++ consume raw point clouds; no need to voxelize or render
- PointNet:
 - Learns both local and global point features
 - Per-point local feature vectors by shared MLPs on n individual points
 - T-Nets try to make it invariant to rotation by learning a 3×3 transformation matrix
 - Global feature vector obtained by max pooling $n \times 1024 \rightarrow 1024$
 - Applications:
 - Classify point cloud: global feature vector fed into MLP + softmax + CE loss
 - Segment point cloud: global feature vector concat with local feature vectors, per point classification
 - Useful for object part segmentation or scene semantic segmentation
 - Notes:
 - Max pooling is an example of a **symmetric function, which is order invariant**
- PointNet++:
 - Repeatedly uses PointNets to aggregate sets of nearby points
 - For segmentation, need to interpolate to obtain the original input dimensions



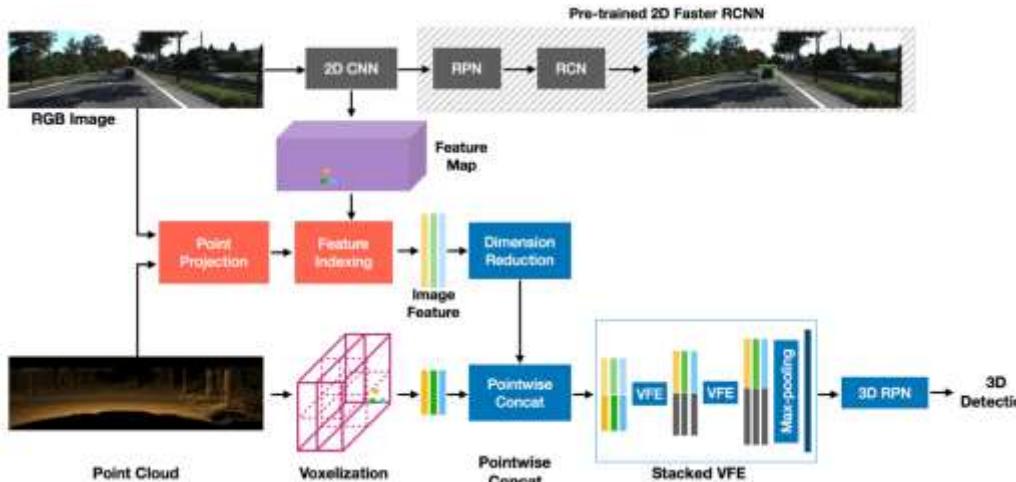
How do voxelnets work?

- Performs 3D object detection on point clouds
- Steps:
 - Partitions pointcloud into 3D voxel grid, and then applies a PointNet w/ global concat to each cell to get a sparse 3D feature map of dimensions $H \times W \times D \times C$
 - Applies 3D convolutions onto 3D feature map (still in form $H \times W \times D \times C$)
 - Merge depth and channels; flatten to $H \times W \times C$
 - RPN finds bboxes in 2D feature map, which also regress bbox
 - 7 Bbox params: xyz center, length, width, height, degree yaw about z axis



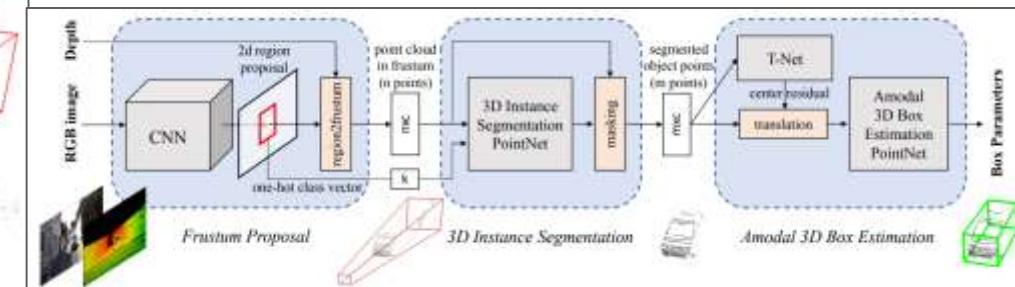
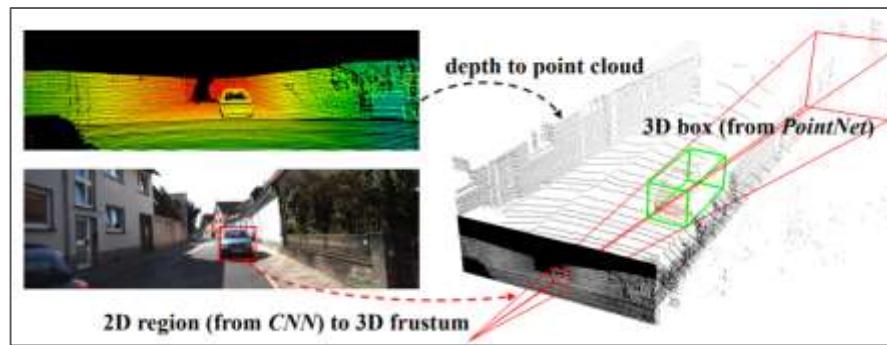
How does the MVX-Net extend voxelnets?

- Instead of 3D detection only using pointcloud data, also incorporates image data
- 3D points are projected to the image using the calibration information and learned corresponding image feature map features are appended to the 3D points using interpolation



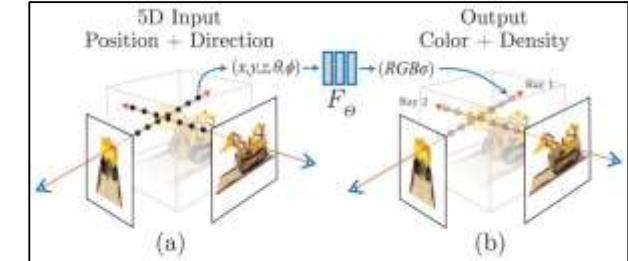
How do Frustum PointNets work?

- Instead of just using a 3D-lidar based pointcloud (e.g. the voxelnet approach), uses RGBD + image setup
- Steps:
 - Make 2D detection on image
 - Project 2D BBox to a frustum of points, similar to Psuedo-LiDAR
 - Segment points w/ pointnet to remove background clutter
 - Feed 3D instance segmentation of points to regress amodal (shape-completed) 3D BBox parameters
- For maximum effectiveness, combine with BEV methods to still detect 2D occluded objects



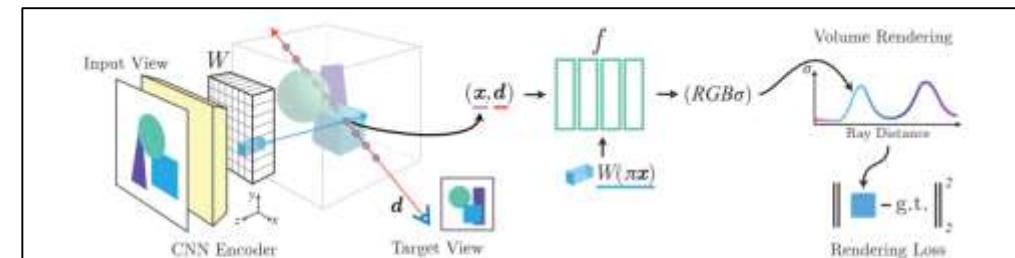
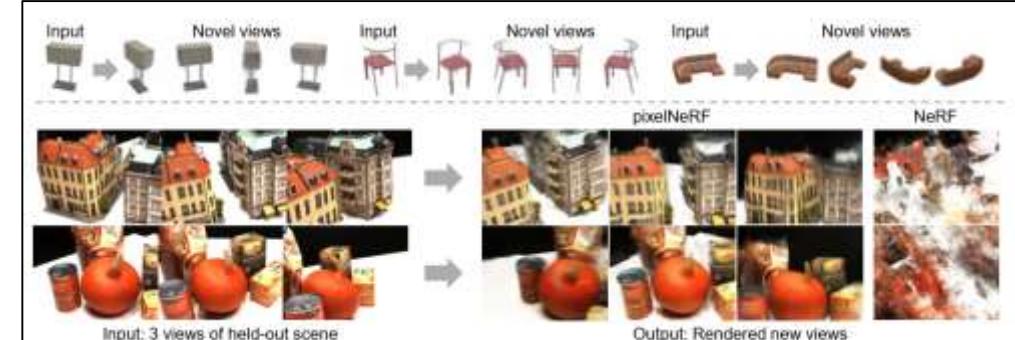
How do NeRFs work?

- Method for novel view synthesis & 3D reconstruction based on an MLP-based implicit volumetric scene
 - Given set of **images from a static scene with known camera poses**, render the scene from novel viewpoints
 - Requires good coverage eg for small object (20-100 imgs), indoor room (100-1000+ imgs)
- First points x_i are obtained by ray casting; for a pixel in a novel view, $x_i = r(t_i) = o + t_i d$
 - o is camera origin, d is ray direction, t is depth along ray
- MLP f_θ takes in 3D points x_i and predicts:
 - **Radiance (color)** $c_i = f_\theta(\gamma(x_i), \gamma(d)) \in \mathbb{R}^3$
 - Feeding in view direction allows us to learn specular reflections
 - **Volume Density** $\sigma_i = f_\theta(\gamma(x_i)) \in \mathbb{R}_+$
 - $\sigma_i \approx 0$ mostly empty air, transparent
 - $\sigma_i \rightarrow \infty$ opaque, eg solid objects or heavy fog
 - Intermediate values of σ_i represent semi-transparent things (eg dirty glass, light fog)
- **Positional Encoding** function $\gamma(x) = [\gamma(x), \gamma(y), \gamma(z)]$ is used which empirically allows the MLP to learn better sharp color transitions and fine details
 - $\gamma(x) = [\sin(2^0 \pi x), \cos(2^0 \pi x), \dots, \sin(2^{L-1} \pi x), \cos(2^{L-1} \pi x)]$
- **Volume Rendering Equation** is $\hat{C}(r) = \sum_{i=1}^N T_i * \alpha_i * c_i$; accumulates light along a camera ray through a semi-transparent 3D scene through weighted avg
 - $\alpha_i = 1 - \frac{1}{e^{\sigma_i \delta_i}} \in [0,1]$: **Opacity** (contribution strength) of point i
 - δ_i is the step size
 - For air $\alpha_i \approx 0$; for solid surfaces $\alpha_i \approx 1$
 - $T_i = \exp(-\sum_j^i \sigma_j \delta_j) \in [0,1]$: **Transmittance**, telling us how much light makes it to point i , without being blocked by earlier points
 - If earlier parts of ray are dense $T_i \approx 0$, if nothing in the way $T_i \approx 1$
- Training objective: minimize pixel error between predicted and GT RGB: $L = \sum_{rays} \|\hat{C}(r) - C_{gt}(r)\|^2$
- Can extract explicit 3D representation by volume grid sampling + marching cubes, or rendering depthmaps + TSDF fusion
- Advantages over traditional MVS:
 - Great at modeling soft shadows, fine details (eg hair), translucency, view-dependent specularities (eg glossy reflections)
 - Avoids artifacts of meshes, textures, etc
- Limitations:
 - Can be slow to render at inference (this has been largely mitigated in recent works)
 - Works best for static scenes; difficulties with dynamic scenes
 - Sensitive to coverage
 - Limited controllability in implicit form
 - Vanilla formulation optimized per-scene, does not generalize to new scenes



How does PixelNeRF work?

- Learns a NeRF from one or view images in a generalized feed-forward way, eliminating per-scene test-time optimization
- Projects 3D query points onto image feature map using intrinsics and uses bilinear interpolation to sample the image feature
 - Sample features then passed through linear layer and added to NeRF network at multiple intermediate activations
 - Local feature may be ambiguous/incorrect due to occlusions; model needs to learn priors to hallucinate in those cases
- Trained on only 2D multi-view image supervision; no GT 3D data or masks



How does DreamFusion work?

- Text-to-3D generation using pretrained & frozen 2D text-to-image diffusion model (Imagen); no 3D training data
- Training:
 - First, instantiates a NeRF per prompt, differentiably rendering random viewpoints
 - Then, perturb render with noise
 - Imagen predicts added noise
 - Noise defines a gradient direction to update NeRF parameters; backprop into NeRF weights
- Intuitively, we are updating the weights of the NeRF so that the rendered images are more likely under the prompt-conditional distribution defined by Imagen

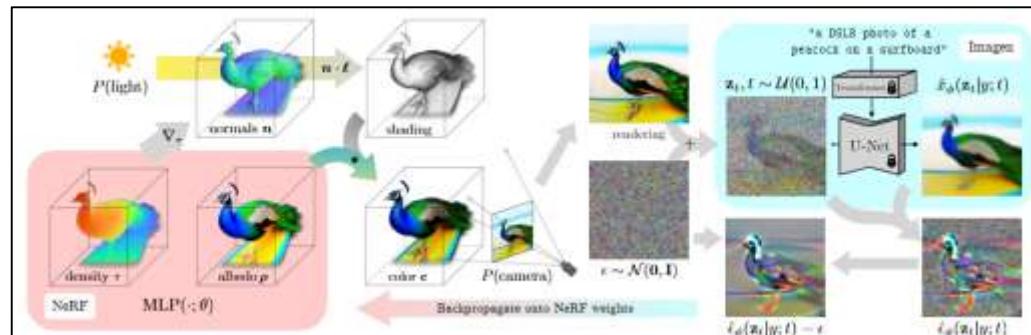
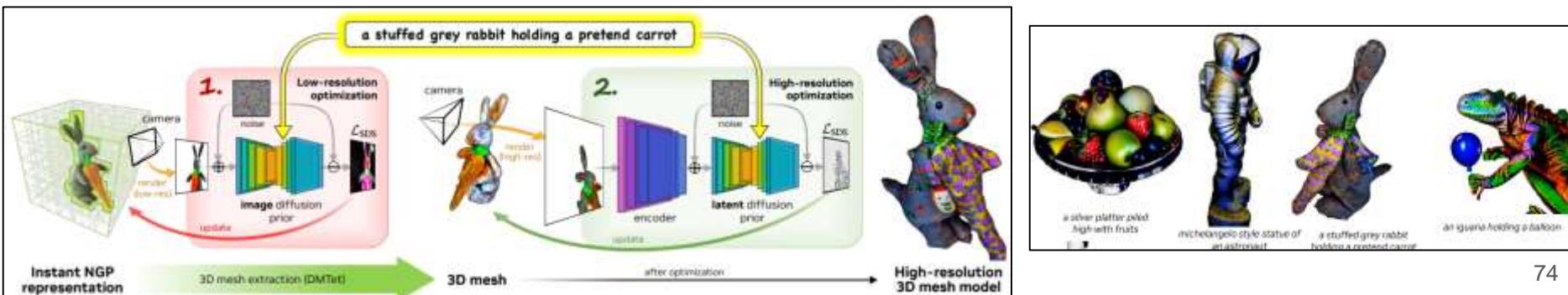


Figure 3: DreamFusion generates 3D objects from a natural language caption such as “*a DSLR photo of a peacock on a surfboard.*” The scene is represented by a Neural Radiance Field that is randomly initialized and trained from scratch for each caption. Our NeRF parameterizes volumetric density and albedo (color) with an MLP. We render the NeRF from a random camera, using normals computed from gradients of the density to shade the scene with a random lighting direction. Shading reveals geometric details that are ambiguous from a single viewpoint. To compute parameter updates, DreamFusion diffuses the rendering and reconstructs it with a (frozen) conditional Imagen model to predict the injected noise $\hat{\epsilon}_\phi(z_t|y; t)$. This contains structure that should improve fidelity, but is high variance. Subtracting the injected noise produces a low variance update direction $\text{stopgrad}[\hat{\epsilon}_\phi - \epsilon]$ that is backpropagated through the rendering process to update the NeRF MLP parameters.



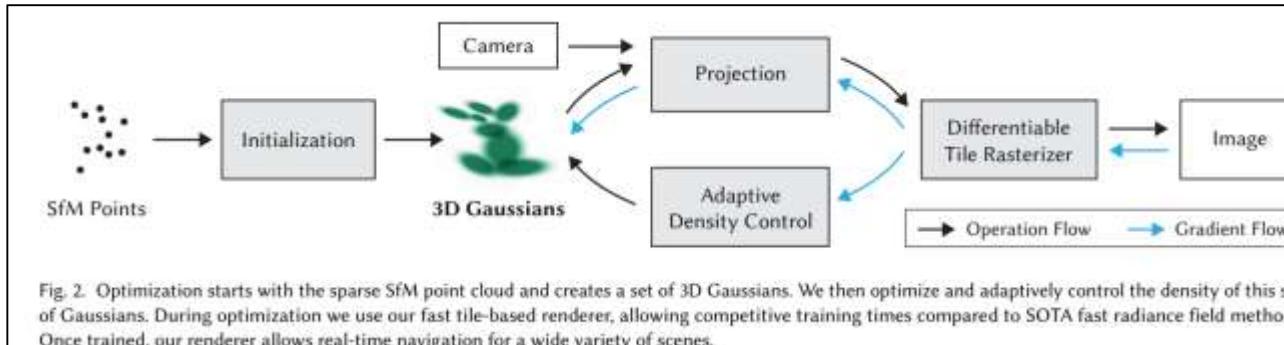
How does Magic3D improve DreamFusion?

- DreamFusion is slow & low-res; Magic3D improves via 2 stages
- Stage One: Low-Res 3D Implicit Generation
 - Use grid-based NeRF (Instant-NGP), which is a popular option for more efficiency
 - Optimize the same way DeepFusion does, with a low-res (64x64) diffusion network
- Stage Two: High Res 3D Mesh Generation
 - First convert Instant-NGP into an SDF
 - Differentiably turn SDF into textured mesh using differentiable marching cubes
 - Differentiably render mesh at random viewpoints
 - Low res network no longer used; instead, optimize using a high-res (512x512) diffusion network
 - By doing differentiable mesh extraction instead of non-differentiable, there's more flexibility with the mesh topology



What is Gaussian Splatting?

- Explicit 3D representation that models objects as a set of Anisotropic 3D Gaussians
 - Each gaussian has: a **Mean** $\mu_i \in \mathbb{R}^3$, **Covariance Matrix** $\Sigma_i \in \mathbb{R}^{3 \times 3}$, **Color** $c_i \in \mathbb{R}^3$, **Opacity** $o_i \in [0,1]$
 - Anisotropic: can have different variances in different directions (can be ellipsoidal)
- Optionally, also add **spherical harmonics** to allow for view-depending effects like specular highlights
 - Common choice is to use first 9 basis functions $Y_0(v), \dots, Y_8(v)$ where $v = (x, y, z)$ is viewing direction
 - Each gaussian would have 27 new parameters representing color; $c^R \in \mathbb{R}^9, c^G \in \mathbb{R}^9, c^B \in \mathbb{R}^9$
 - Then, calculate RGB as $[\sum_{i=0}^8 c_i^R * Y_i(v), \sum_{i=0}^8 c_i^G * Y_i(v), \sum_{i=0}^8 c_i^B * Y_i(v)]$
 - Intuitively, spherical harmonics are orthogonal functions defined on a sphere that can approximate any function of a sphere. So we're learning the weights of these basis functions to provide view-dependent color.
- Can be rendered / "splatted" onto 2D pixel grid using alpha blending. First we sort gaussians by depth, and then contributions are accumulated with $C_{final} = \sum_i T_i \alpha_i c_i$
 - $T_i = \prod_{j=1}^{i-1} (1 - \alpha_j)$ is the transmittance up to the i -th Gaussian
 - $\alpha_i(p) = o_i \cdot \exp\left(-\frac{1}{2}(p - x_i)^T (\Sigma_i)^{-1} (p - x_i)\right)$ is the alpha value
- Original, vanilla formulation doesn't use a neural network; parameters optimized directly using a set of images with known poses, with pixel-space loss
- Generally, point-cloud via some other algorithm (eg COLMAP) is used to initialize position
- With some heuristic scheduling, **prune** invisible gaussians (eg low alpha, or outside bbox) and **densify** for new gaussians (if gaussian has large gradient)
- Common regularizations:
 - Opacity regularization to prevent too many opaque layers / over saturation
 - Covariance regularization to incentivize compact & sharp gaussians



Left: Gaussian Splat. Right: Same splat with all gaussians semi-transparent

What's the connection between Gaussian Splatting & NeRFs?

- **Similarities**

- Both are radiance fields; Function that map 3D locations & directions to color.
- optimized using photometric loss
- Require substantial coverage of scene
- Focus is photorealistic rendering from novel viewpoints

- **Differences**

- NeRF is an implicit representation, while GS is explicit
- NeRFs are generally slower to train & render; GS allows for real-time rendering
- NeRFs are harder to interpret

- **Converting NeRF → GS**

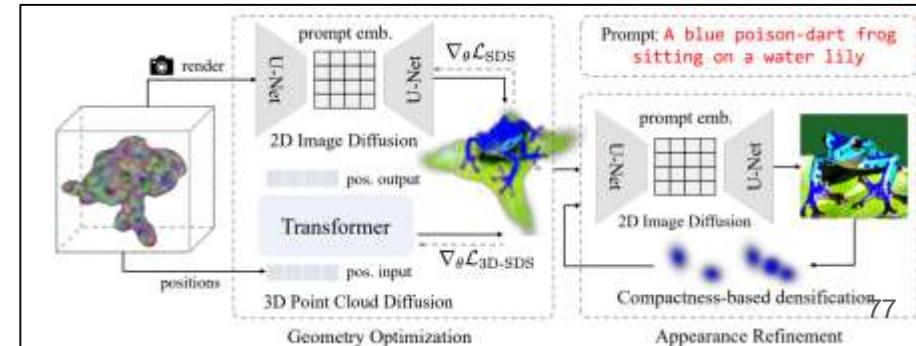
- Initialize gaussians by querying NeRF MLP
- Optimize/estimate shape of gaussian by analyzing local point structure.
- Further optimize by using rendered NeRF images

- **GS → NeRF**

- Render images from GS and train NeRF with that supervision

How does GSGEN work?

- Text-to-3D method by optimizing a per-instance Gaussian Splat
- First stage: Geometry Optimization
 - Given prompt, first use Point-E (a text-to-pointcloud diffusion model) to initialize gaussian centers
 - Optimize by randomly differentiably rendering views, adding noise, and updating gaussians with diffusion gradient
 - Additional loss term comes from the same Point-E 3D network, to continue updating positions of gaussians
- Second Stage: Appearance Refinement
 - Continue optimizing using 2D image diffusion network
 - Add gaussian densification algorithm based on close neighbors to improve smoothness
 - Prune gaussians with low alpha



How does LGM work?

- Generates 3D gaussian splats from text or single-view images, in a feed-forward manner
- First, obtain 4 images from a text prompt (Mvdream, text to MV image diffusion) or image prompt (ImageDream, image to MV image diffusion)
- Given the 4 images as input, regresses per-image pixel-aligned gaussian splats ($128 \times 128 \times 4 = 65536$ gaussians)
 - Input images are always in a canonical pose (0,90,180,270 at fixed elevation 0)
- After fusing the 4 gaussians, render novel views and supervise in rgb space
- Contends from domain gap – training using 3D asset renders, but inference using diffusion generated multi-view images

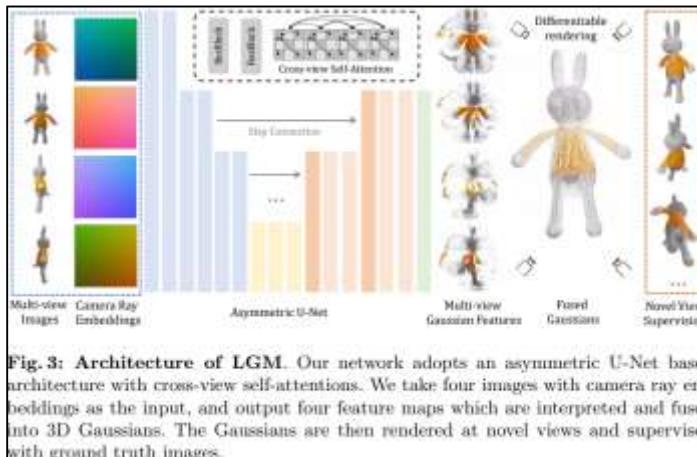


Fig. 3: Architecture of LGM. Our network adopts an asymmetric U-Net based architecture with cross-view self-attentions. We take four images with camera ray embeddings as the input, and output four feature maps which are interpreted and fused into 3D Gaussians. The Gaussians are then rendered at novel views and supervised with ground truth images.

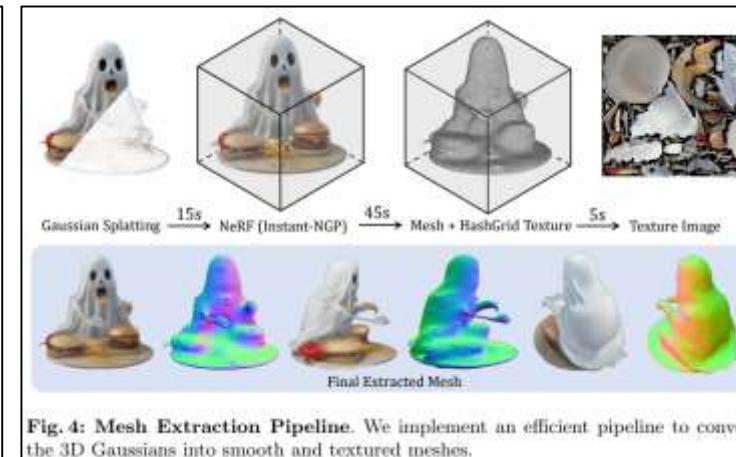


Fig. 4: Mesh Extraction Pipeline. We implement an efficient pipeline to convert the 3D Gaussians into smooth and textured meshes.

How does Point-E work?

- Text-to-3dPointCloud method with two stages
- Text-to-Image Stage:
 - First, use a fine-tuned GLIDE model to generate a synthetic 2D view from a text prompt
- Image-to-3D Stage:
 - A point cloud diffusion model generates a 3D RGB point cloud
 - Model is transformer-based, invariant to point order
 - Points represented as tensor of shape $n \times 6$, where n is the number of points and we have 3D position + RGB values, normalized to [-1,1]
 - Points passed through linear layer to get embedding; in total, input sequence has point embeddings, timestep, and image embeddings
 - Encoder processes sequence and the last n tokens are projected as ϵ and Σ for the diffusion step
- Trained on millions of 3D models

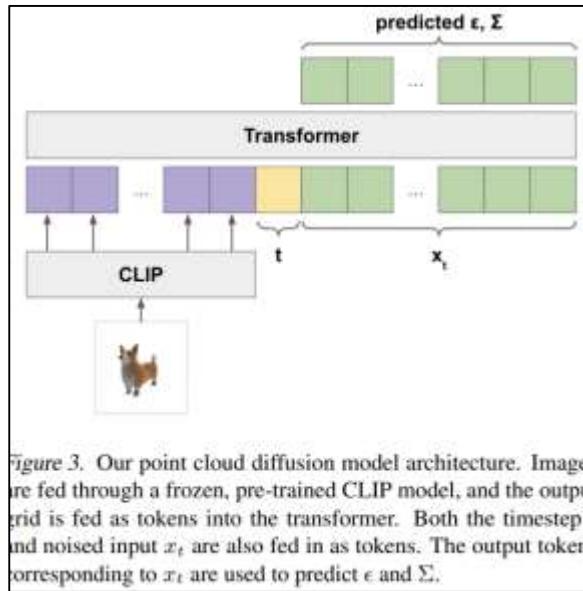
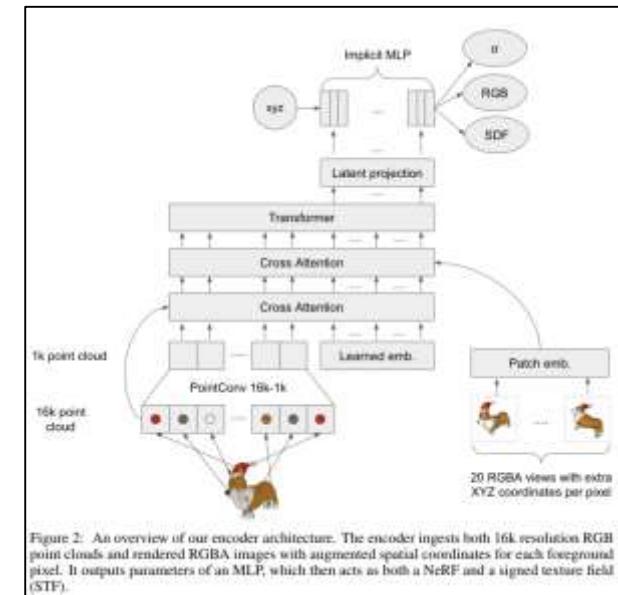
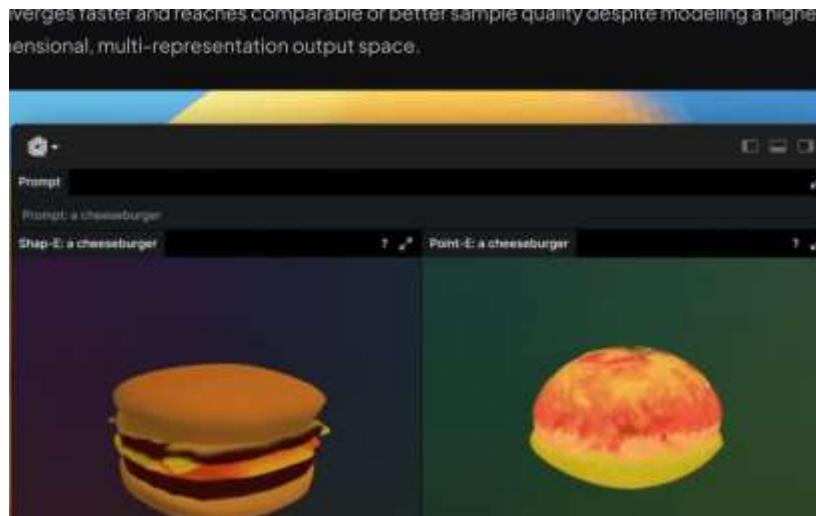


Figure 3. Our point cloud diffusion model architecture. Images are fed through a frozen, pre-trained CLIP model, and the output grid is fed as tokens into the transformer. Both the timestep t and noised input x_t are also fed in as tokens. The output tokens corresponding to x_t are used to predict ϵ and Σ .



How does Shape-E work?

- Text-to-NeRF model, with two stages
- Encoder Stage
 - Train encoder to predict parameters of a NeRF given a point cloud and its RGB views. NeRF supervised with rendering loss
- Diffusion Training
 - Learn to generate latent MLP parameters conditioned on text/image prompts



What are triplanes & k-planes?

- 3D representation that provides a more efficient alternative to 3D voxel grids
- Comprised of 3 orthogonal 2D feature planes aligned with principal axes:
 - F_{xy} aligned with XY plane, F_{xz} aligned with XZ plane, F_{yz} aligned with YZ plane
 - Each plane is a 2D grid of high dimensional feature vectors, shape $N \times N \times C$
- To sample feature at point $p = (x, y, z)$
 - Project onto F_{xy} (x, y), F_{xz} (x, z), & F_{yz} (y, z)
 - For each projection, apply bilinear interpolation to get the corresponding feature vectors f_{xy}, f_{xz}, f_{yz}
 - Aggregate features by concatenation/summation
- Final feature can be used for downstream applications, eg into a MLP to decode NeRF parameters
- Space usage is $O(N^2C)$ instead of $O(N^3C)$ for voxel grids
- Can be further generalized to K-Planes, where a d-dimensional space is turned into “d choose 2” 2D feature planes
- One key application is for 4D NeRFs, with 6 planes: xy, yz, xz, xt, yt, zt
 - NeRF predicts rgb & density given x, y, z, t, d (viewing direction)
 - Needs a multiview dynamic video dataset as supervision (Multiview video, multiple cameras capturing a scene over time)
 - A 4D point (x, y, z, t) is projected onto all 6 planes; features bilinearly interpolated and multiplied (Hadamard product)
 - This is done repeatedly through volumetric ray rendering & then supervised with photometric loss
 - K-plane formulation allows for regularizing in certain dimensions (eg Laplacian smoothness in time)

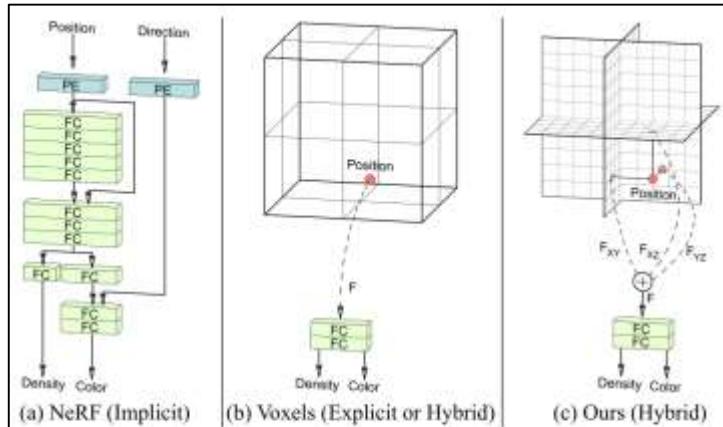


Figure 2. Neural implicit representations use fully connected layers (FC) with positional encoding (PE) to represent a scene, which can be slow to query (a). Explicit voxel grids or hybrid variants using small implicit decoders are fast to query, but scale poorly with resolution (b). Our hybrid explicit–implicit tri-plane representation (c) is fast and scales efficiently with resolution, enabling greater detail for equal capacity.



How are 4D Gaussian Splats formulated & optimized?

- Introduces a deformation field that warps a canonical, initial set of 3D gaussians over time
- First stage: a fixed set of initial gaussians are learned (eg over the first few frames)
- Second stage: Learn gaussian deformation field network
 - Hex-plane based, where 6 multi-resolution voxel planes are learned and fed into MLP to regress time-dependent changes on the canonical gaussians
 - Changes include position, rotation, scale, and optionally alpha & color
 - Deformed gaussians differentiably splatted & supervised using photometric loss
- Two kinds of data setup, requiring camera poses
 - Multiple cameras capturing the same scene simultaneously from different viewpoints -- preferred
 - Monocular video (single camera) – more challenging

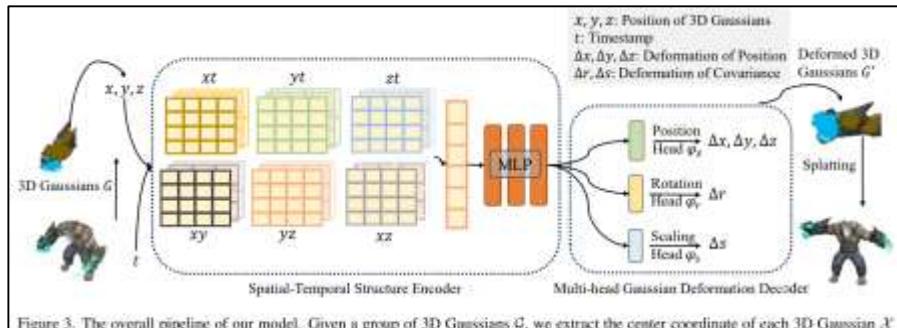
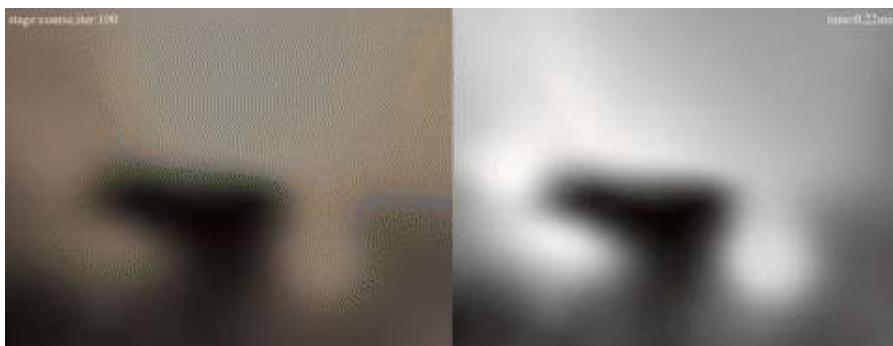
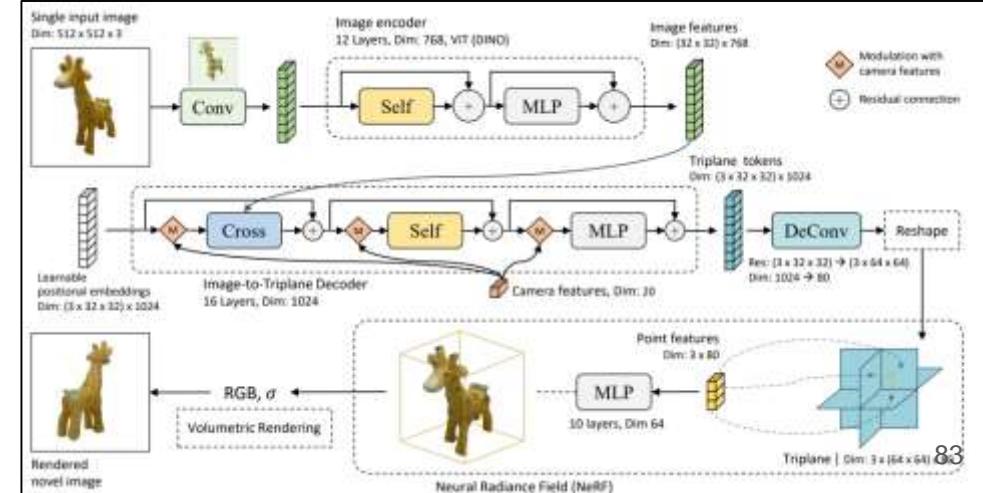


Figure 3. The overall pipeline of our model. Given a group of 3D Gaussians \mathcal{G} , we extract the center coordinate of each 3D Gaussian \mathcal{X} and timestamp t to compute the voxel feature by querying multi-resolution voxel planes. Then a tiny multi-head Gaussian deformation decoder is used to decode the feature and get the deformed 3D Gaussians \mathcal{G}' at timestamp t . The deformed Gaussians are then splatted to get the rendered images.



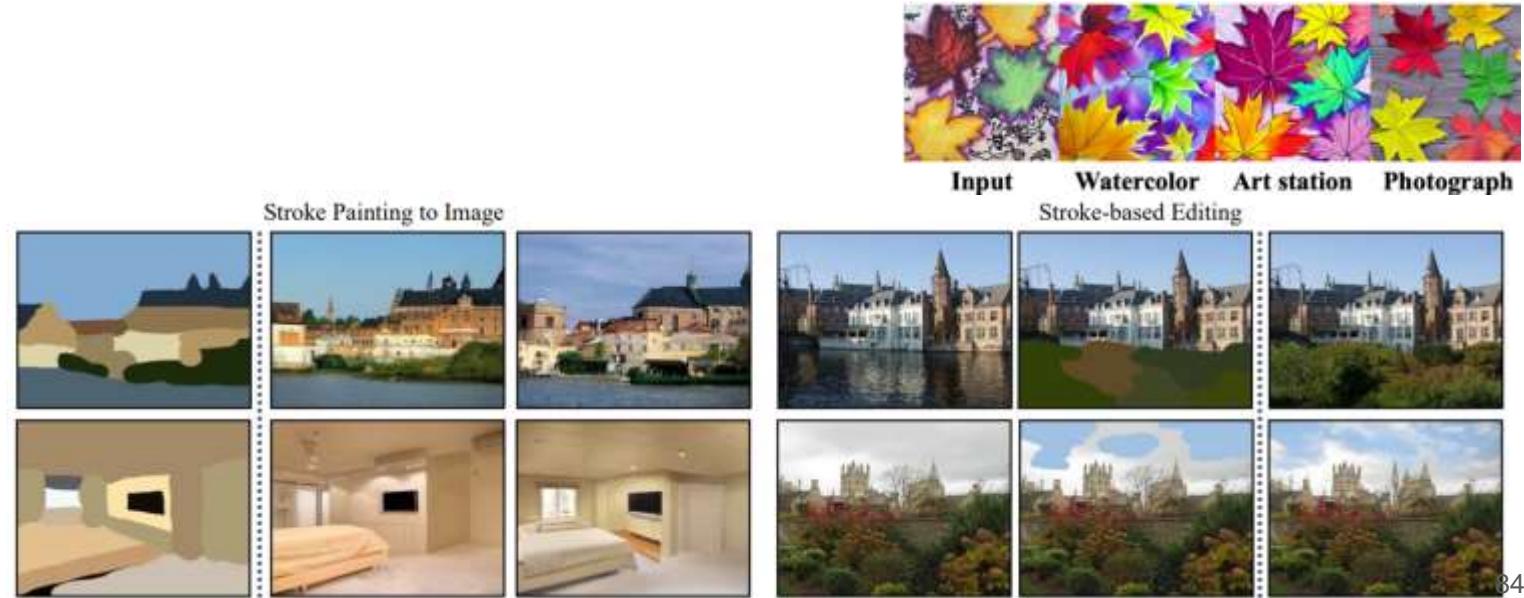
How does LRM work?

- Single-view 3D reconstruction at scale
 - Trained end-to-end on a million 3D models, mix of real & synthetic data
- Three main steps:
 - **Pre-trained DINO ViT** processes images into patches
 - **Image-To-Triplane decoder** converts 2D features into a triplane representation using cross-attention & camera-conditioned modulation using extrinsics & intrinsics
 - **Triplane NerF**: MLP predicts RGB & density from queried triplane features. Random novel views rendered from this and supervised using RGB loss
- Training feeds in camera embedded features (from intrinsics/extrinsics); at test time, we feed a fixed vector
 - Feeding this in during training led to better convergence
 - A trade-off; sometimes leads to degradation in performance



How does SDEdit work?

- Framework for guided image generation/editing, derived using the score-based continuous-time SDE formulation of diffusion
- Can do stroke-based editing as well as style transfer
- Method: by “hijacking”/exploiting the fact that diffusion models can be run from arbitrary noise levels (not just from pure noise), we can intentionally corrupt the guide and denoise towards a sample that lies both near the guide and on the natural image manifold
 - Eg, the colors used in the stroke painting matters – influences types of objects generated
 - Similarly, for style transfer, we destroy low-frequency noise (style) but keep high-frequency noise (structure). Similar intuition to how diffusion models incrementally add details.
- Uses a standard pretrained model & does not change weights



How does CAT4D work?

Given a monocular video (real or generated), generates a 4D gaussian using 3 steps:

Step 1: Multi-view Video Diffusion Model Training (monocular video to multi-view videos)

- Takes in M input conditional set of images I^{cond} , poses P^{cond} , & times T^{cond} , trained to learn distribution of N target images I^{tgt} given their camera parameters P^{tgt} & times T^{tgt} : $p(I^{tgt} | I^{cond}, P^{cond}, T^{cond}, p^{tgt}, T^{tgt})$
- Based off of CAT3D: multi-view latent diffusion model, U-Net based with 3D self-attention to connect multi-view latents
- Uses classifier-free guidance with two components to disentangle camera & time

Step 2: Generating Consistent Multi-View Videos

- Uses an **Alternating Sampling Strategy** to generate large consistent multiview video grids (many cameras x timesteps)
- Uses a sliding-window approach to gradually build full grid to enforce viewpoint & time consistency
- First, Multi-View sampling is done (fix timestep, generate multiple consistent views w/ sliding window)
- Then, Temporal Sampling is done (fix camera, generate timestamps w/ sliding window)
- Above procedure can be iterative refined using SDEdit – partial refinement & consistency corrections (instead of regenerating from scratch every time)

Step 3: 4DG generation

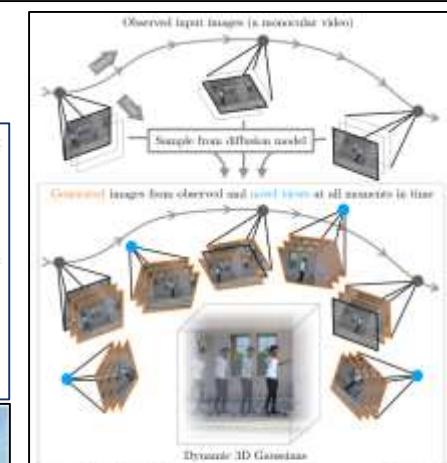
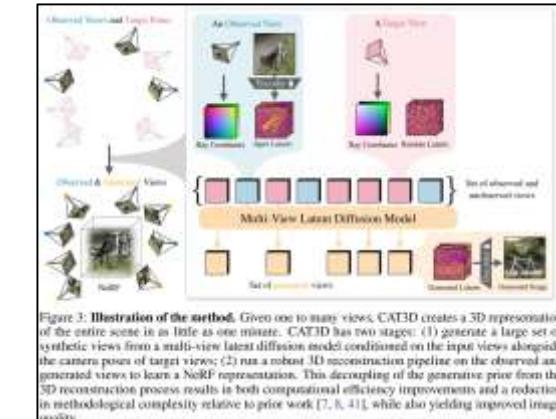
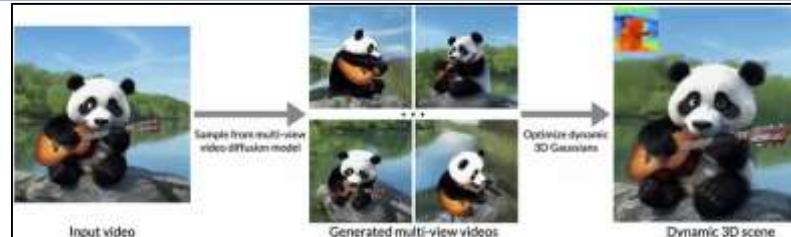
- Use the multi-view videos and learn a gaussian splat along with deformation field

Data-driven 4D prior is difficult to build due to lack of diverse data; approach here is to use combinations of synthetic 4D data, fixed camera videos, and object-centric video captures. The latter two is augmented through CAT3D (to provide camera motion) and Lumiere (to provide object motion) respectively

$$\hat{\epsilon}_{\theta}(z, P_{tgt}, c_I, c_T) = \underbrace{\epsilon_{\theta}(z, P_{tgt}, \emptyset, \emptyset)}_{\text{unconditioned}} + s_I \cdot \underbrace{\left[\epsilon_{\theta}(z, P_{tgt}, c_I, \emptyset) - \epsilon_{\theta}(z, P_{tgt}, \emptyset, \emptyset) \right]}_{\text{view-guided}} + s_T \cdot \underbrace{\left[\epsilon_{\theta}(z, P_{tgt}, c_I, c_T) - \epsilon_{\theta}(z, P_{tgt}, c_I, \emptyset) \right]}_{\text{full-guided}}$$

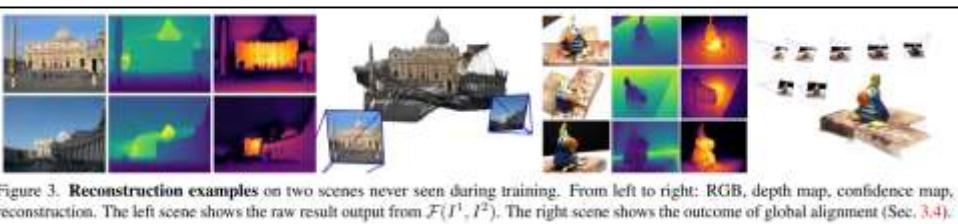
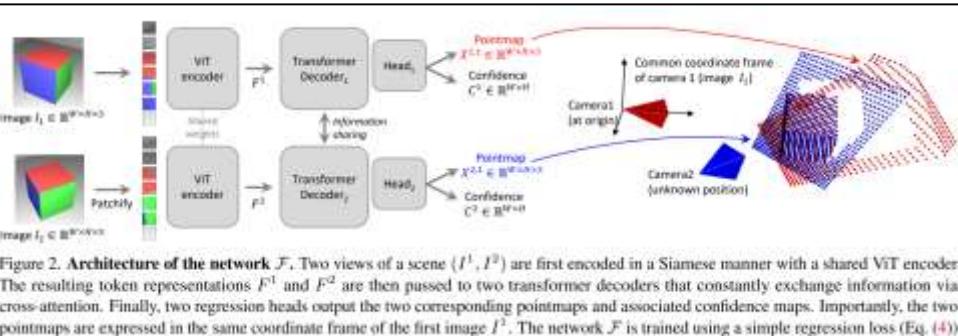
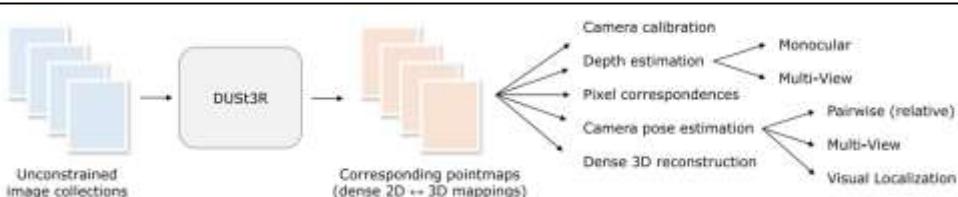
$c_I = \{I_{cond}, P_{cond}\}$
 $c_T = \{T_{cond}, T_{tgt}\}$

- First term: unconditional prediction
- Second term: encourages the sample to align with the camera/view input
- Third term: encourages the sample to align with the temporal input



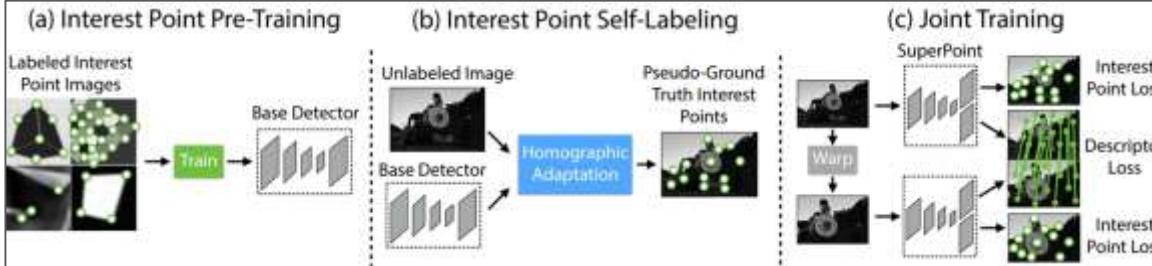
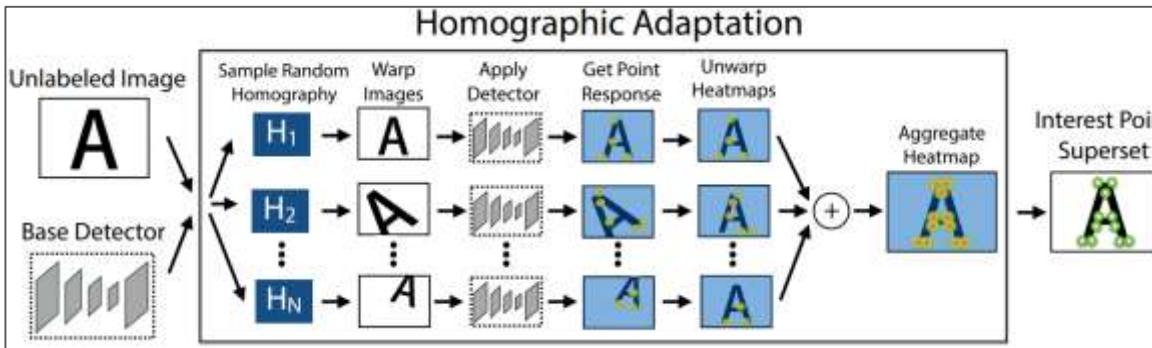
How does Dust3r work?

- Unified, general, simple model that tackles a lot of 3D vision tasks which were historically fragmented
 - Predicts dense pointmaps (2D grid where each pixel maps to a 3D point) directly from image pairs, w/o intrinsics/poses
 - Pointmap $\in \mathbb{R}^{W \times H \times 3}$
 - Crucially, both pointmaps from a stereo pair are expressed/aligned in the same coordinate frame, for easy geometric operations
 - Similar to a dense 3D pointcloud, but with one point per pixel; point corresponding to pixel approximately lies along that pixel's viewing ray
 - Also predicts confidence map (for where point is not well defined, like the sky)
- Downstream tasks using pointmaps:
- Camera calibration**
 - If we assume principal point is centered & pixels are squares, we can easily solve for focal length
 - Camera pose estimation**
 - Compare pointmaps $X^{1,1}$ & $X^{1,2}$ and solve optimization problem to find extrinsics that align them (pointmap for the same image, but in different coordinate frames by swapping the input order)
 - Alternatively, use PnP + RANSAC
 - Depth estimation**
 - Pixel correspondences**
 - Dense 3D reconstruction**
- Trained using a combination of real & synthetic RGBD data



How does SuperPoint work?

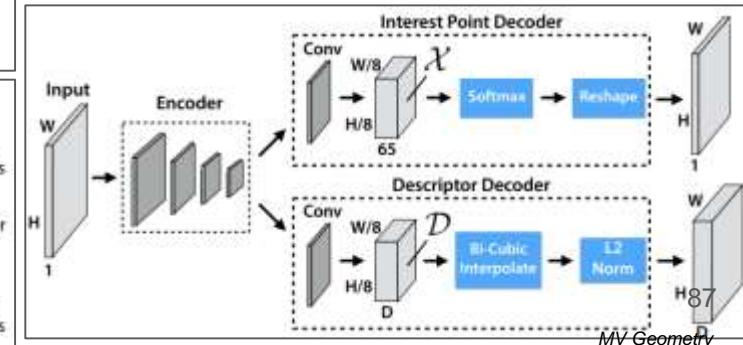
- Unlike tasks like human keypoint estimation which can be tackled with strong annotated supervision, general interest point prediction is semantically ill-defined
- SuperPoint instead takes an approach using synthetic data + pseudo GT
 - a) Train a **base detector** on a large **synthetic dataset** of simple geometric shapes
 - b) Warp target domain, unlabeled input images multiple times to help the base model see the image from many different viewpoints & scales, and save those **pseudo-GT** results. The warping is called “Homographic Adaptation” & helps bridge the domain gap. In the paper, this is iterative done twice to bootstrap.
 - c) Using pseudo-GT, train 2 head model: **per-pixel interest point binary classification** (CE loss) & **feature description** (hinge loss)
 - Training occurs with a pair of synthetically warped image with pseudo-GT & known correspondence from a randomly sampled homography
 - Hinge loss is applied to all pairs of descriptor cells, where the dot product of matching pairs are incentivized to be larger than a margin, and mismatching pairs smaller than a margin



$$\mathcal{L}_d(\mathcal{D}, \mathcal{D}', S) = \frac{1}{(H_e W_e)^2} \sum_{h=1}^{H_e} \sum_{w=1}^{W_e} l_d(\mathbf{d}_{hw}, \mathbf{d}'_{h'w'}; s_{hw h'w'}),$$

where

$$l_d(\mathbf{d}, \mathbf{d}'; s) = \lambda_d * s * \max(0, m_p - \mathbf{d}^T \mathbf{d}') + (1 - s) * \max(0, \mathbf{d}^T \mathbf{d}' - m_n),$$



Full Stack Deep Learning

<https://fullstackdeeplearning.com/>

[Training & Debugging Phase] Compare data vs model parallelism when doing distributed training.

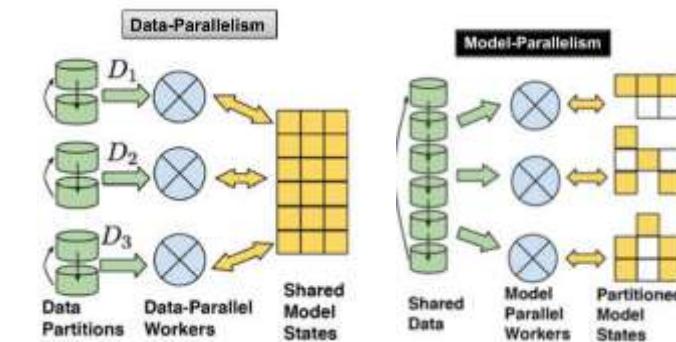
Distributed training is a must for large datasets and large models, and involves multiple GPUs and/or machines to train a single model. Note that forward computations and backprop are inherently serial in nature, so one can think of these methods as “workarounds”.

- **Data parallelism**

- Helpful for large datasets
- Splits batch into multiple GPUs (each gpu must fit a copy of the model on it), evaluate/backprop on each, and combine to get overall gradient to update across all GPUs
- Generally simple to do
- [Hogwild paper](#) (NIPS 11, 2000+ cit) showed that you don't even need locks; SGD works even if the shared model sometimes receives older gradients.

- **Model Parallelism**

- Needed when a model does not fit on a GPU
- Splits the model parameters on separate GPUs
- Introduces a lot of complexity, so see if this can be avoided by getting larger GPUs
- One way this is implemented (AWS Sagemaker) is by splitting consecutive layers into different GPUs, and having an execution pipeline where different devices can work on forward and backward passes for different data samples at the same time



[Testing Phase] Draw a diagram of an ML system, and how you can test it.

Training Infrastructure Tests

- Avoid bugs in training pipeline
- Train on a single batch or epoch, to make sure it can overfit

Evaluation Testing

- Fully understand performance characteristics; see if guardrails are necessary
- Eval on all metrics, datasets, and slices
- Invariance robustness testing
- Compare to previous models, using thresholds. Avoid performance "degradation leaks".
- Simulation testing, if possible (needs rendering)
- Test for privacy, fairness, bias

Basic Pred. Tests

- Checks basic prediction functionality
- Load pretrained model, test on few key examples

Shadow Tests

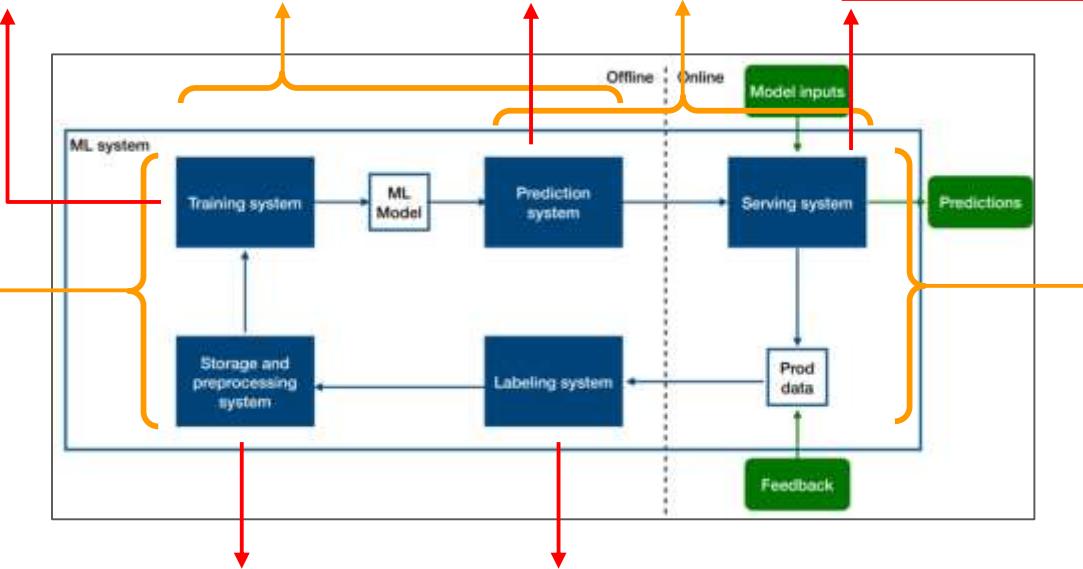
- Checks prod. bugs and compare performance w/ old model
- Run new model in parallel with old model in production, but in "shadow mode"
- Unlike dev/sandbox env, more realistic

A/B Tests

- Understand how new model affects user experience and business metrics
- Canary model on a tiny fraction of incoming requests
- Compare metrics on control and variation cohort

Training Tests

- Makes sure training is reproducible with current architecture
- Pull a fixed dataset and train it
- Check to make sure performance is consistent or better than before



Data Expectation Tests

- Catch data quality issues before training
- Check for null values, numerical ranges, and formatting using rules and thresholds
- Example package: "Great Expectations"

Labeling QA

- Catch poor quality labels
- Train labelers well
- Assign them a "trust score"; promote best performers to be "QA experts"
- Aggregate labels by voting

Production Monitoring

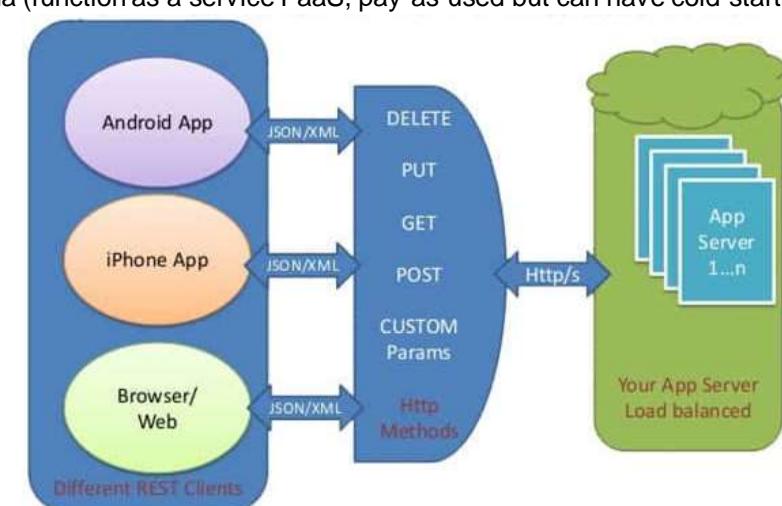
- Want to keep model healthy in deployment
- Extremely important; what's used to produce next model and decide what new data to train on
- Often hard to do directly if you don't have labels; can instead check **business metrics** and **model inputs/predictions** for irregularities
- Need to check for bugs, outliers, class imbalance, **concept drift** (dependent var. changes)
- **Data drift** (independent var. changes) can be detected by comparing a sample of previous data and current prod. data using distance metric. The data can be features, or something more basic like HSV.

- Don't use KL-Divergence, too sensitive and range dependent (div by zero)
- Use **Kolmogorov-Smirnov Stat** (max dist between histogram CDF) or **D1 distance** (sum of histogram CDF distances)
- Can experiment with more complex metrics like wasserstein dist or earth mover's dist.
- Or, use **rule-based metrics** (w/ great expectations, e.g. min, max, mean in acceptable range, Col A > col B, etc)

[Deployment Phase] Explain what a REST APIs is, and its advantages.

- A widely used, standardized, simple format for the access of resources and communicating over the internet, when it comes to client-server models
- Works over HTTPs and uses GET (read), POST (create), PUT (update), and DELETE, a server destination URI/URL, along with simple data formats like JSON and XML. Can also contain headers with auth data.
 - Client sends request, server sends response
- Crucially, is **stateless** on the server side; each request processed independently
 - If notions of a session are required, then this is managed on the client-side, and passed onto other services as necessary by the client.
 - This creates a clear distinction between front-end (client) and back-end (server)
- The stateless feature allows easy **horizontal scaling** (no state syncing required) and **caching**.
- Can have an actual server running 24/7, or use something like AWS Lambda (function as a service FaaS, pay-as-used but can have cold-starts sometimes. Considered “serverless” since you don’t have to manage them)
 - Typically uses a dependency management system like docker
- Can be implemented in python using packages like Flask or FastAPI

TLDR: API format over HTTP that is stateless & thus easily scalable/cacheable



[Deployment Phase] What is Federated Learning, at a high level?

- Federated Learning: training a global model from data on local devices, without ever having access to the data
- approach for collaborative machine learning without compromising data privacy and security

2. Local Training: Local devices or servers, such as smartphones or edge devices, perform training on their own data locally. This training typically consists of updating the model based on the data available on each device.
3. Model Updates: Instead of sending raw data to the central server, local devices or servers send model updates (such as gradients or weight updates) to the central server. These updates represent the changes needed to improve the global model.
4. Aggregation: The central server collects these model updates from various devices and aggregates them to update the global model. This can be done using techniques like federated averaging, where the central server averages the updates from multiple devices to update the global model.

