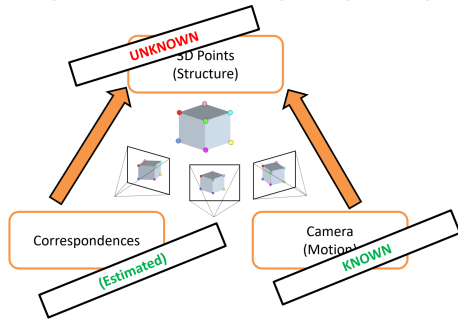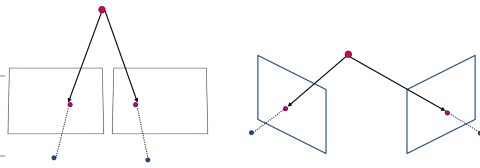# Final
# Lec 13   Epipolar Geometry + Calibration

## Simple Stereo;
Corresp + Camera = Disparity = depth$^{-1}$



- The two cameras need not have parallel optical axes.
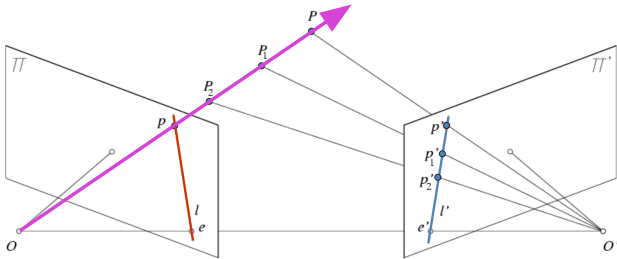- Assume camera intrinsics are calibrated



**Same hammer:**
Find the correspondences, then solve for structure

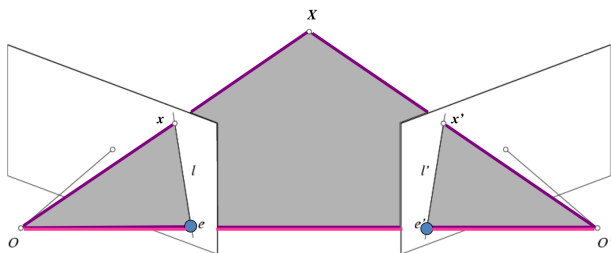### General case, known camera, find depth:
### Option 2

1. **Find correspondences**
2. Triangulate

## Camera helps Correspondence:
### Epipolar Geometry



## Correspondence gives camera:
### Epipolar Geometry



Camera (Motion) 与 Correspondence 关系串联:
Epipolar Geometry!

Given Camera: How to
find Epipolar? Option 1:
Homography: two image plane
to one plane! ⇒ parallel

Option 1: Rectify via homography



Original stereo pair

Then find correspondences on the horizontal scan line

After rectification

Option2 · Use math (Epipolar Geometry) to find epipolar

Given O · O', p. then since
Intrinsics are calibrated. so we can
confirm O·O' equivalent position.
连线 OP. OP 上每一点与 O' 连线,
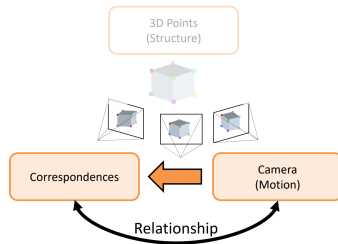与 Π' 上的交点形成的连线便是
epipolar

## Epipolar constraint



- Potential matches for *p* have to lie on the corresponding epipolar line *l'*.
- Potential matches for *p'* have to lie on the corresponding epipolar line *l*.

因为 OP. OO' 形成了 Epipolar Plane!
△: 极点、Epipole: OO' (baseline) 与
两个 Image Plane 的交点、
如图可见:
Epipolars must pass through epipoles!

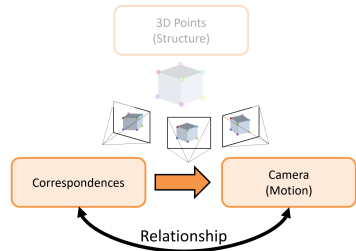## Parts of Epipolar geometry



- **Baseline** – line connecting the two camera centers
- **Epipolar Plane** – plane containing baseline (1D family)
- **Epipoles**
  = intersections of baseline with image planes
  = projections of the other camera center
  = vanishing points of the baseline
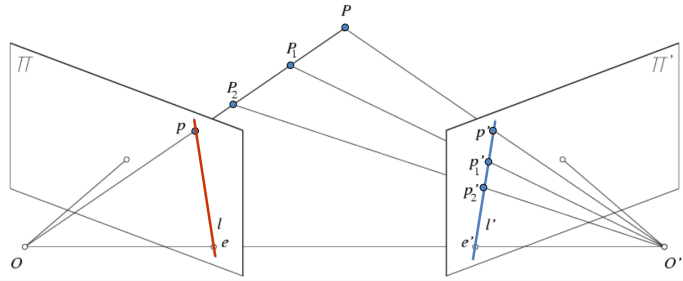
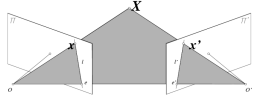Epipoles *infinitely* far away, epipolar lines parallel

## Ok so where were we?

- Setup: Calibrated Camera (both extrinsic & intrinsic)
- Goal: 3D reconstruction of corresponding points in the image
- We need to find correspondences!
- ➜ 1D search along the epipolar line!
- ➜ Need: Compute the epipolar line from camera

## Ok so what exactly are l and l'?
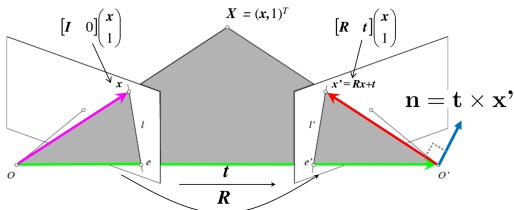


### Step 0: Factor out intrinsics



$$x = K[R\ t]X$$
$$K^{-1}x = [R\ t]X$$

- Let's factor out the effect of K (do everything in 3D)
- Make it into a ray with $K^{-1}$ and use depth = 1
- This is called the *normalized* image coordinates. It may be thought of as a set of points with identity K

$$x_{\text{norm}} = K^{-1}x_{\text{pixel}} = [I\ 0]X, \qquad x'_{\text{norm}} = K'^{-1}x'_{\text{pixel}} = [R\ t]X$$

- Assume that the points are normalized from here on

### Epipolar constraint: Calibrated case



$$\mathbf{n} = \mathbf{t} \times \mathbf{x}'$$

The vectors $x$, $t$, and $x'$ are coplanar

What can you say about their relationships, given $\mathbf{n} = \mathbf{t} \times \mathbf{x}'$?

$$\boxed{\mathbf{x}' \cdot (\mathbf{t} \times \mathbf{x}') = 0}$$

$$\mathbf{x}' \cdot (\mathbf{t} \times (R\mathbf{x} + \mathbf{t})) = 0$$
$$\mathbf{x}' \cdot (\mathbf{t} \times R\mathbf{x} + \mathbf{t} \times \mathbf{t})) = 0$$

$$\boxed{\mathbf{x}' \cdot (\mathbf{t} \times R\mathbf{x}) = 0}$$

Given $K$ (intrinsic) and $o'$
(以 $o$ 为参考系) 则:

$$X_{\pi} = [I\ 0]X = x \qquad X_{\pi'} = [R\ t]X = x'$$
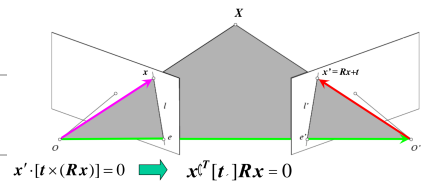
其中 $R$, $t$ 是 $o'$ 相对于 $o$ 的 rotation & translation $X_{\pi'}$

则: $t \times (Rx + t) \perp x'$

$\therefore x' \cdot (t \times Rx) = 0$

$\therefore x'^T E x = 0$,

$$E = [t_\times] R$$



$$x' \cdot [t \times (Rx)] = 0 \implies x'^T [t_\times] Rx = 0$$

Recall: $\mathbf{a} \times \mathbf{b} = \begin{bmatrix} 0 & -a_z & a_y \\ a_z & 0 & -a_x \\ -a_y & a_x & 0 \end{bmatrix} \begin{bmatrix} b_x \\ b_y \\ b_z \end{bmatrix} = [\mathbf{a}_\times]\mathbf{b}$

The vectors $x$, $t$, and $x'$ are coplanar

$$x' \cdot [t \times (Rx)] = 0 \implies x'^T \underbrace{[t_\times] R}_{E} x = 0 \implies x'^T E x = 0$$

Recall: $\mathbf{a} \times \mathbf{b} = \begin{bmatrix} 0 & -a_z & a_y \\ a_z & 0 & -a_x \\ -a_y & a_x & 0 \end{bmatrix} \begin{bmatrix} b_x \\ b_y \\ b_z \end{bmatrix} = [\mathbf{a}_\times]\mathbf{b}$

$$\boxed{\text{Essential Matrix} \\ \text{(Longuet-Higgins, 1981)}}$$

The vectors $x$, $t$, and $x'$ are coplanar

$$\boxed{x'^T E x = 0}$$

$E\,x$ is the epipolar line associated with $x$ ($l' = E\,x$)

- Recall: a line is given by $ax + by + c = 0$ or

$$\mathbf{l}^T \mathbf{x} = 0 \quad \text{where} \quad \mathbf{l} = \begin{bmatrix} a \\ b \\ c \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

Recall, knowing the camera gives you the essential matrix (i.e. the plane per point)

So the DoF has to match up

Essential matrix: 3 x 3, 9 numbers, but rank2 means 2 columns fully define = 6 parameters -1 for scale = 5 DoF

Extrinsic Camera (R, T): 3 for rotation, 3 for translation, but -1 for scale = 5 DoF!

$$\boxed{x'^T E x = 0}$$

$E\,x$ is the epipolar line associated with $x$ ($l' = E\,x$)
$E^T x'$ is the epipolar line associated with $x'$ ($l = E^T x'$)
$E\,e = 0$ and $E^T e' = 0$
$E$ is singular (rank two)
$E$ has five degrees of freedom

### Some property of E. If K is not the same?

### Epipolar constraint: Uncalibrated case

- Recall that we normalized the coordinates
$$x = K^{-1}\hat{x} \quad x' = K'^{-1}\hat{x}' \qquad \hat{x} = \begin{bmatrix} u \\ v \\ 1 \end{bmatrix}$$
where $\hat{x}$ is the image coordinates
- But in the *uncalibrated* case, $K$ and $K'$ are unknown!
- We can write the epipolar constraint in terms of *unknown* normalized coordinates:
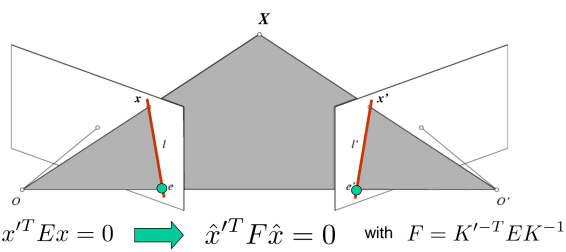
$$x'^T E x = 0$$
$$(K'^{-1}\hat{x}')'^T E (K^{-1}\hat{x}) = 0$$
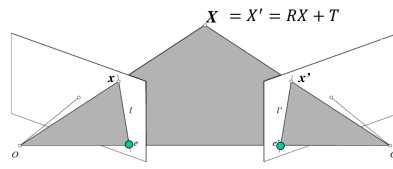$$\hat{x}'^T K'^{-T} E (K^{-1}\hat{x}) = 0$$
$$\hat{x}'^T F \hat{x} = 0$$

$$F = K'^{-T} E K^{-1}$$

$$\boxed{\text{Fundamental Matrix} \\ \text{(Faugeras and Luong, 1992)}}$$

We know about the camera, $K_1$, $K_2$ and $[R \ t]$:

$$X = X' = RX + T$$

$$x'^T E x = 0 \quad \Rightarrow \quad \hat{x}'^T F \hat{x} = 0 \quad \text{with} \quad F = K'^{-T} E K^{-1}$$

and found the corresponding points: $x \leftrightarrow x'$

$$x = K X \qquad x' = K'X'$$
$$= K'(RX + T)$$

- $F\hat{x}$ is the epipolar line associated with $\hat{x}$ ($l' = F\hat{x}$)
- $F^T\hat{x}'$ is the epipolar line associated with $\hat{x}'$ ($l = F^T \hat{x}'$)
- $Fe = 0$ and $F^Te' = 0$
- $F$ is singular (rank two)
- $F$ has *seven* degrees of freedom

How many unknowns + how many equations do we have?

only unknowns!

Solve by formulating Ax=0, see H&Z ch.12

但并非任何时候 OX 与 O'X' 两个 ray 会相交:

Find 3D point

Ray's don't always intersect because of noise!!!

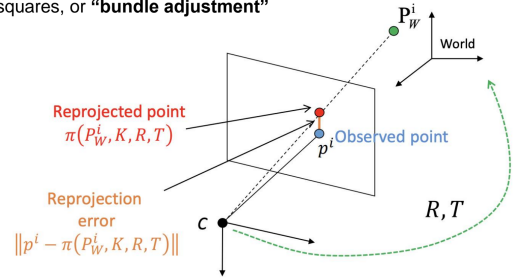因此: $\Longrightarrow$

Solve with non-linear least squares

$\Leftarrow$ Summary

Even if you do everything right, you will still be off because of noise, this is called the Reprojection Error

In practice with noise, want to directly minimize this with non-linear least squares, or **"bundle adjustment"**

$P_W^i$    World

Reprojected point $\pi(P_W^i, K, R, T)$

Observed point $p^i$

Reprojection error $\|p^i - \pi(P_W^i, K, R, T)\|$

$c$    $R, T$

Solve with non-linear least squares, iteratively

## Summary: Two-view, known camera

0. Assuming known camera intrinsics + extrinsics

1. Find correspondences:
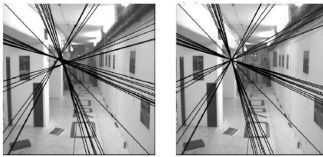   - Reduce this to 1D search with Epipolar Geometry!

2. Get depth:
   - If simple stereo, disparity (difference of corresponding points) is inversely proportional to depth
   - In the general case, triangulate.

Now consider if we have correspondence, can we estimate F ?

3D Points (Structure)

Correspondences $\Rightarrow$ Camera (Motion)

$$x = (u,v,1)^T, \quad x' = (u',v',1)$$

$$[u' \ v' \ 1]\begin{bmatrix} f_{11} & f_{12} & f_{13} \\ f_{21} & f_{22} & f_{23} \\ f_{31} & f_{32} & f_{33} \end{bmatrix}\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = 0$$

$$[u'u \ u'v \ u' \ v'u \ v'v \ v' \ u \ v \ 1]\begin{bmatrix} f_{11} \\ f_{12} \\ f_{13} \\ f_{21} \\ f_{22} \\ f_{23} \\ f_{31} \\ f_{32} \\ f_{33} \end{bmatrix} = 0$$

Solve homogeneous linear system using eight or more matches

$x'^T F x = 0$

$\Leftarrow$ Eight points cause if F' s.t. $x^T F' x = 0$ then: $x^T (cF')x = 0 \Rightarrow$ DOF : 8

Enforce rank-2 constraint (take SVD of $F$ and throw out the smallest singular value)

$E = T_x R$

If we know E, we can recover t and R

$$\begin{bmatrix} e_{11} & e_{12} & e_{13} \\ e_{21} & e_{22} & e_{23} \\ e_{31} & e_{32} & e_{33} \end{bmatrix} = \begin{bmatrix} 0 & -t_z & t_y \\ t_z & 0 & -t_x \\ -t_y & t_x & 0 \end{bmatrix}\begin{bmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{bmatrix}$$

Given that $T_x$ is a Skew-Symmetric matrix ($a_{ij} = -a_{ji}$) and $R$ is an Orthonormal matrix, it is possible to "decouple" $T_x$ and $R$ from their product using "Singular Value Decomposition".

The geometry of three views is described by a 3 x 3 x 3 tensor called the *trifocal tensor*

The geometry of four views is described by a 3 x 3 x 3 x 3 tensor called the *quadrifocal tensor*

After this it starts to get complicated…

Get the essential matrix with K (or some estimates of K) …

in practice you calibrate your cameras so you know K or have a very good estimate

$$E = K'^T F K.$$

Now if we have correspondence, how to calibrate!

What are the camerea parameters?
- Extrinsics (R, T)
- Intrinsics (K)

How am I situated in the world + what is the shape of the ray

Approach two: Solve linear System!

World origin

Focal length

Camera Rotation

Camera Translation

## How to estimate the camera?

1. Estimate the fundamental/essential matrix!

2. Another method: Calibration

$$x = K[R \ \ t]X$$

$$\begin{bmatrix} su \\ sv \\ s \end{bmatrix} = \begin{bmatrix} * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{bmatrix}\begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

If we know the points in 3D we can estimate the camera!!

## Can we factorize M back to K [R | T]?

Yes.

Why? because K and R have a very special form:

$$\begin{bmatrix} su \\ sv \\ s \end{bmatrix} = \begin{bmatrix} m_{11} & m_{12} & m_{13} & m_{14} \\ m_{21} & m_{22} & m_{23} & m_{24} \\ m_{31} & m_{32} & m_{33} & m_{34} \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

*intrinsic*

$$\begin{bmatrix} f_x & s & o_x \\ 0 & f_y & o_y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{bmatrix}$$

QR decomposition

Practically, use camera calibration packages (there is a good one in OpenCV)

### Solve for m's entries using linear least squares

**Ax=0** form

$$\begin{bmatrix} X_1 & Y_1 & Z_1 & 1 & 0 & 0 & 0 & 0 & -u_1X_1 & -u_1Y_1 & -u_1Z_1 & -u_1 \\ 0 & 0 & 0 & 0 & X_1 & Y_1 & Z_1 & 1 & -v_1X_1 & -v_1Y_1 & -v_1Z_1 & -v_1 \\ & & & & & & & & & & & \\ X_n & Y_n & Z_n & 1 & 0 & 0 & 0 & 0 & -u_nX_n & -u_nY_n & -u_nZ_n & -u_n \\ 0 & 0 & 0 & 0 & X_n & Y_n & Z_n & 1 & -v_nX_n & -v_nY_n & -v_nZ_n & -v_n \end{bmatrix} \begin{bmatrix} m_{11} \\ m_{12} \\ m_{13} \\ m_{14} \\ m_{21} \\ m_{22} \\ m_{23} \\ m_{24} \\ m_{31} \\ m_{32} \\ m_{33} \\ m_{34} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$
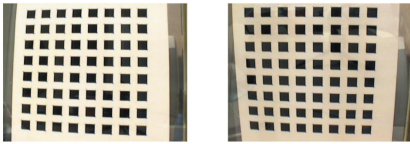
Similar to how you solved for homography!

---

## Need at least 6 pairs of (3D coord, 2D image coord)

### Inserting a 3D known object…

Also called "Tsai'scalibration" requires non-coplanar 3D points, is not very practical…

Modern day calibration uses a planar calibration target



Developed in 2000 by Zhang at Microsoft research

### Doesn't plane give you homography?

Yes! If it's a plane, it's only a homography, so instead of recovering 3x4 matrix, you will recover 3x3 in Zhang's method

The 3x3 gives first two columns of R and T

$$\begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{bmatrix} = \begin{bmatrix} \alpha_u & 0 & u_0 \\ 0 & \alpha_v & v_0 \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} r_{11} & r_{12} & t_1 \\ r_{21} & r_{22} & t_2 \\ r_{31} & r_{32} & t_3 \end{bmatrix}$$
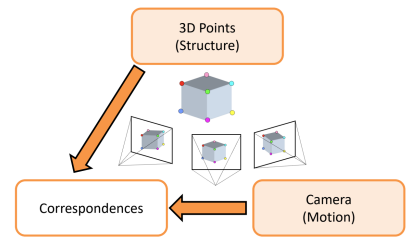
*Moreover. if 3D points are in a plane, than it is :)*

---

# Lec 14: SfM (Structure from Motion)

### Camera Calibration; aka Perspective-n-Point



### Stereo (w/2 cameras); aka Triangulation



### You can easily get correspondence via projection from 3D points + Camera



---

*In the operations introduced before, they showed on to use triangle relationship. But what if none of these are known?*

### Structure from motion



Reconstruction (side)

(top)

- Input: images with points in correspondence $p_{i,j} = (u_{i,j}, v_{i,j})$

- Output
  - structure: 3D location $x_i$ for each point $p_i$
  - motion: camera parameters $R_j$, $t_j$ possibly $K_j$

- Objective function: minimize *reprojection error*

### Ultimate: Structure-from-Motion



Start from nothing known (except maybe intrinsics), exploit the relationship to slowly get the right answer

*Given $p_{i,j}$: 第i个点在第j个相机中的像素生标*

*⇒ Structure & Motion*

*\* Correspondence : Unknown !*

## Feature matching

Match features between each pair of images

Refine matching using RANSAC to estimate fundamental matrix between each pair

- Point: 3D position in space ($\mathbf{X}_j$)

- Camera ($C_i$):
  - A 3D position ($\mathbf{c}_i$)
  - A 3D orientation ($\mathbf{R}_i$)
  - Intrinsic parameters (**focal length**, aspect ratio, …)
  - 7 parameters (3+3+1) in total

*Get optimized*

*$\chi_i, R_j, t_j$ in one optimization problem!*



- Minimize sum of squared reprojection errors:

$$g(\mathbf{X}, \mathbf{R}, \mathbf{T}) = \sum_{i=1}^{m} \sum_{j=1}^{n} w_{ij} \cdot \left\| \mathbf{P}(\mathbf{x}_i, \mathbf{R}_j, \mathbf{t}_j) - \begin{bmatrix} u_{i,j} \\ v_{i,j} \end{bmatrix} \right\|^2$$

*indicator variable*: is point $i$ visible in image $j$?

*predicted image location*  *observed image location*

- Minimizing this function is called *bundle adjustment*
  - Optimized using non-linear least squares, e.g. Levenberg-Marquardt
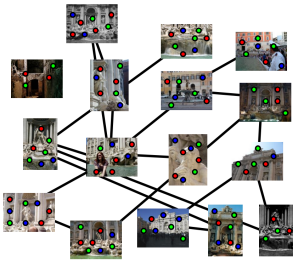
### Challenges:
- Large number of parameters (1000's of cameras, millions of points)
- Very non-linear objective function

- Important tool: Bundle Adjustment [Triggs *et al.* '00]
  - Joint non-linear optimization of both cameras and points
  - Very powerful, elegant tool
- The bad news:
  - Starting from a random initialization is very likely to give the wrong answer
  - Difficult to initialize all the cameras at once

### The good news:
- Structure from motion with two cameras is (relatively) easy
- Once we have an initial model, it's easy to add new cameras

### Idea:
- Start with a small seed reconstruction, and grow

*⇒ Incremental SfM:*

## Incremental SfM: Algorithm

1. Pick a strong initial pair of images
2. Initialize the model using two-frame SfM
3. While there are connected images remaining:
   a. Pick the image which sees the most existing 3D points
   b. Estimate the pose of that camera
   c. Triangulate any new points
   d. Run bundle adjustment

- We want a pair with many matches, but which has as large a baseline as possible

*Strong:*

*Want many matches but want baseline as large as possible.*



✓ lots of matches
✗ small baseline

✓ large baseline
✗ very few matches

✓ large baseline
✓ lots of matches

### Multi-view Stereo (Lots of calibrated images)
- Input: calibrated images from several viewpoints (known camera: intrinsics and extrinsics)
- Output: 3D Model



Figures by Carlos Hernandez

In general, conducted in a controlled environment with multi-camera setup that are all calibrated

*The problem of SfM is that its output is sparse point cloud. With SfM's output of calibration information, can we form dense point cloud? ⇒ Multi-View stereo*

### Multi-view stereo: Basic idea



Photometric error across different depths

error
depth

In this manner, solve for a depth map over the whole reference view

*For an image pixel patch, consider ray through it with different depth, and see which depth most fit in other images best ⇒ Depth map*

- Can match windows using more than 1 other image, giving a **stronger match signal**

- If you have lots of potential images, can **choose the best subset** of images to match per reference image

- Can reconstruct a depth map for each reference frame, and the merge into a **complete 3D model**

Source: Y. Furukawa

width of a pixel

all of these points project to the same pair of pixels

Large Baseline          Small Baseline

- What's the optimal baseline?
  - Too small: large depth error
  - Too large: difficult search problem

Discretized Scene Volume

Input Images (Calibrated)

Goal: Assign RGB values to voxels in V *photo-consistent* with images

---

For 3D reconstruction. another approach: volumetric stereo

For every voxel, if projected on these cameras have minor error to gt, then this voxel remains, while others who don't satisfy this will be removed.

---

Lec 15 & 16 & 17 : NeRF.

### Problem Statement

⇐ What problem NeRF want to solve?

**Input:**
A set of calibrated Images

**Output:**
A 3D scene representation that renders novel views

- Need to know the camera parameters: **extrinsic** (viewpoint) & **intrinsics** (focal length, distortion, etc)

## Structure from Motion! (last lecture)

Input: set of images.
Output: extrinsics, intrinsics, 3D points, pixel correspondences

### What was before NeRF? "Photogrammetry"

- Problem: Given calibrated cameras, recover highly detailed 3D **surface** model
- Often the output is textured meshes

Before: how to: image + calibration ⟹

⇐ Photogrammetry : ⟹ Complicated!

Advantage of NeRF

High quality reconstruction with view-dependent effects

Can represent non-opaque objects

## NeRF's Three Key Components

Original 3D representation - - - - - -

### Volumetric representations

Model the *entire* space
Can be explicit (voxels) or implicit (NeRF)

### Polygonal Meshes

- A mesh is a set of vertices with faces that defines the topology
- Mesh = {Vertices, Faces}
  - Vertices: N x 3
  - Faces: F x {3, 4, ...} specifying the edges of a polygon
  - Triangle faces most common but tetrahedrons (tets) are also.
- **Surface** is explicitly modeled by the faces
- Most common modeling representation

$(x, y, z, \theta, \phi) \rightarrow F_{\Omega} \rightarrow (r, g, b, \sigma)$

Neural Volumetric 3D Scene Representation

Ray

3D volume

Camera

Differentiable Volumetric Rendering Function

Optimization via Analysis-by-Synthesis

# Lightfield / Lumigraph

Levoy and Hanrahan, SIGGRAPH 1996
Gortler et al. SIGGRAPH 1996

- Previous approaches for modeling the Plenoptic Function
- Take a lot of pictures from many views
- Interpolate the rays to render a novel view

**Stanford Gantry 128 cameras**

**Lytro camera**

Figure from Marc Levoy

光场 / 光栅图：有海量光线数据，通过计算和内插 这些已记录的光线，来合成任意新视角的图像

- These methods are called Image Based Rendering, because they literally interpolate the ray colors to make a new image

- i.e. no 3D information is recovered (you have to know the camera)

**Plenoptic Function**

**NeRF**

NeRF requires *integration* along the viewing ray to compute the Plenoptic Function

Bottom line: it models a 5D plenoptic function!

## Density: Second key difference from lightfields, plenoptic function ✗

- Continuous probability density function (PDF) over "stuff"

- Connected to opacity: high density == very opaque, solid

$$(x, y, z, \theta, \phi) \rightarrow F_\Omega \rightarrow (r, g, b, \sigma)$$

Spatial location · Viewing direction · Output color · Output density

Two difference between lightfield (NeRF) and plenoptic function.

## Where NeRF stands

**Appearance Based Reconstruction (Image Based Rendering)**

- can do Image Based Rendering well, while also being a 3D representation
- Does not suffer from limitations of surface models
- Easy to optimize from images

**NeRFs**

**Physics based Reconstruction (3D Surface Modeling)**

Lightfield/Lumigraph (No 3D representation)

Layered Depth Images (LDIs) · Multi-Plane Images (MPIs)

One 3D Surface, View-Dependent Texture Mapping

One 3D Surface, Single Albedo Texture

Conventional Graphics Pipeline

## ✰ "Analysis-by-Synthesis"

- Search for a world state from which you can explain many observations through synthesis

- In English: "If you understand (analyze) something, you can create (synthesize) it"

- (For NeRF): "If you really know what a scene looks like, you can render it from any view"

- (For Chemistry): "If you know how a molecule is structured, you can synthesize it from other molecules"
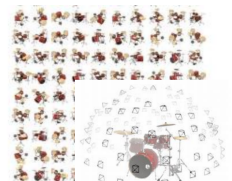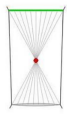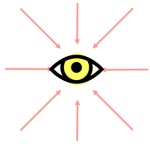
- Commonly used paradigm across CV!

\* Core Function NeRF want to learn: For a point (x, y, z), if look through it via direction (θ, φ), what's its observered rgb value and its volumn density?

Training Strategy: Find a way to generate rgb for one ray with rgb & σ info. as pred. \*

### "Neural Radiance Fields"

How an image is made ("Inference")

$(x, y, z, \theta, \phi) \rightarrow (r, g, b, \sigma)$
Spatial · Viewing · Output · Output

**Volumetric 3D Representation θ**
MLP 9 layers, 256 channels

**Differentiable Volumetric Rendering**

Rendered Image: I'

"Training" Objective (aka Analysis-by-Synthesis):

$$\min_\theta \| \text{Rendered Image: } I' - \text{Observed Image: } I \|_2$$

Toy setting: 2D. no φ & φ.

## How to get MLPs to represent higher frequency functions?

Challenge observed: 介

### Let's simplify, do this in 2D:

$(x, y) \rightarrow \theta \rightarrow (r, g, b) \xrightarrow{\text{Retrieve color from this network for every pixel}} \text{Rendered Image: I'}$

$F_\Omega$ MLP

Optimize with "Training" Objective (aka Analysis-by-Synthesis):

$$\frac{\partial L}{\partial \theta} = \frac{\partial (rgb - rgb')}{\partial \theta}$$

$$\min_\theta \| \text{Rendered Image: } I' - \text{Observed Image: } I \|_2$$

Iteration 1000

MLP output — Supervision image

**Why so?**
**From example, we can see that with slight diff in $x$, rgb can differ greatly!** But if with PE, then 'slight diff in $x$' $\Rightarrow$ distinct PE input
So it can respond high-freq change in rgb!

## Positional Encoding

Standard input — Positionally Encoded input

Fourier Features $\gamma(p) = \left( \sin(2^0 \pi p), \cos(2^0 \pi p), \cdots, \sin(2^{L-1}\pi p), \cos(2^{L-1}\pi p) \right)$

Input | Target
x  y
A
B

With Positional Encoding
x  y
A
B

Target Image

## NeRF Network Architecture

**Now: the architecture:**

**① Prediction of $\sigma$ doesn't depend on direction $\phi \& \psi$! (Intuitively Yes)**

**② 'Coord & Direction' inputs both undergo PE.**  **③ ReLU or Sigmoid? Check desired output range!**

Next section you will implement this:

3D coord

$\sigma > 0!$

x (3D) — PE — Linear (256) ReLU — Linear (256) ReLU — Linear (256) ReLU — Linear (256) ReLU — Concat — Linear (256) ReLU — Linear (256) ReLU — Linear (256) ReLU — Linear (256) ReLU

Linear (1) — ReLU — density (1D)

rgb $\in (-1, 1)$

Linear (256) — Linear (128) ReLU — Linear (3) — Sigmoid — rgb (3D)

$r_d$ (3D) — PE — Concat

$\psi \& \phi$

## Volumetric formulation for NeRF



$\mathbf{c}(t), \sigma(t)$

Ray $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$

Camera

at a point on the ray $\mathbf{r}(t)$, we can query color $\mathbf{c}(t)$ and density $\sigma(t)$

How to integrate all the info along the ray to get a color per ray?

Remember, expected color is equal to

$$\int T(t)\sigma(t)\mathbf{c}(t)dt \approx \sum_i T_i \alpha_i \mathbf{c}_i = \sum_i w_i \mathbf{c}_i$$

$T(t)\sigma(t)$ and $T_i\alpha_i$ are "rendering weights" — probability distribution along the ray (continuous and discrete, respectively)

You can also render entities other than color in 3D, for example it's depth, or any other N-D vector $v_i$

Volume rendered "feature" $= \sum_i w_i v_i$

## Alpha Blending

for two image case, A and B, both partially transparent:

$(C_b, \alpha_b)$

$(C_a, \alpha_a)$

$I = C_a \alpha_a + C_b \alpha_b (1 - \alpha_a)$

How much light is the previous layer letting through?

General D layer case:
$I = \sum_{i=1}^{D} C_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j)$

layer D
layer 2
layer 1

## Summary

for a ray $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$:

differentiable w.r.t. $\mathbf{c}, \sigma$

$$\mathbf{c} \approx \sum_{i=1}^{n} w_i \mathbf{c}_i = \sum_{i=1}^{n} T_i \alpha_i \mathbf{c}_i$$
weights — colors

How much light is blocked earlier along ray:
$$T_i = \prod_{j=1}^{i-1} (1 - \alpha_j)$$

$\delta_i$: stepsize

How much light is contributed by ray segment $i$:
$$\alpha_i = 1 - \exp(-\sigma_i \delta_i)$$


Ray — 3D volume — Camera

$\delta_i$: **stepsize.**

**\*Finally, how to pred rgb of one ray (one pixel) with rgb & $\sigma$ of points along the way?**
**$T_i$: The prob that ray first hit $c_i$**

**$\sigma_i$: logit $\Rightarrow \alpha_i = 1 - \exp(-\sigma_i \delta_i)$**

**The formula of Volumetric Rendering actually originate from alpha blending:**

$\Leftarrow$ **$I = C_a \alpha_a + C_b (1 - \alpha_a) \alpha_b$**

**how much light left to pass through.**

$$I = \sum_{i=1}^{D} C_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j)$$

# When transfered to Volumetric Rendering:
## for differentiable : $\sigma_i \to \alpha_i$ :

$$\alpha_i = 1 - \exp(-\sigma_i \delta_i)$$

So: $C = \sum_{i=1}^{n} c_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j)$

## Summary

for a ray $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$:

differentiable w.r.t. $\mathbf{c}, \sigma$

$$\mathbf{c} \approx \sum_{i=1}^{n} w_i \mathbf{c}_i = \sum_{i=1}^{n} T_i \alpha_i \mathbf{c}_i$$

weights    colors

How much light is blocked earlier along ray:

$$T_i = \prod_{j=1}^{i-1} (1 - \alpha_j)$$

How much light is contributed by ray segment $i$:

$$\alpha_i = 1 - \exp(-\sigma_i \delta_i)$$

# Now : Some bells & whistles:
## How to sample reasonable ? Intuitively,
## most of the space can be vacuum!

## Hierarchical Sampling vs. Acceleration Structures

**Hierarchical Sampling**

Iteratively use samples from NeRF to more efficiently sample visible scene content
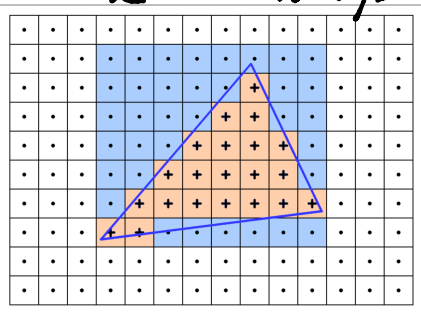
**Acceleration Structures**

Distill/cache properties of NeRF into a structure that helps generate samples: e.g. Occupancy Grids

Straightforward compute —> storage tradeoff



**Key Idea: sample points proportionally to expected effect on final rendering**

$$C \approx \sum_{i=1}^{N} T_i \alpha_i c_i$$

treat weights as probability distribution for new samples



**What about aliasing during coarse sampling?**

Solution: train two NeRFs! —> lower resolution for first "coarse" level



# Solution :
# Hierarchical
# Search

**What about aliasing during coarse sampling?**

Solution: train two NeRFs! —> higher resolution for second "fine" level



# Coarse → Fine, and use two NeRFs!

Rasterization = conversion of primitives to pixels (details in CS184)

**Rasteriza[t]**

**Blending (+ S[)**



# Last topic : Gaussian Splatting
## Preliminary: About Rasterization (光栅化): Instead of casting
## rays, we Map Object to pixels.

## Differentiable Gaussian Rendering

**What is the representation of a 3D Gaussian?**

Position $\mathbf{p}$

$$\mathcal{G}_V(\mathbf{x} - \mathbf{p}) = \frac{1}{2\pi |\mathbf{V}|^{\frac{1}{2}}} e^{-\frac{1}{2}(\mathbf{x}-\mathbf{p})^T \mathbf{V}^{-1} (\mathbf{x}-\mathbf{p})}$$

$$S = \begin{bmatrix} s_x & 0 & 0 \\ 0 & s_y & 0 \\ 0 & 0 & s_z \end{bmatrix}$$   Factorize as scale and rotation: $\mathbf{V} = RSS^T R^T$

$R \in SO(3)$   Each Gaussian also has an opacity and view-dependent color (via SH coefficients): $\alpha, \mathbf{c}$

**How to project to 2D and rasterize?**

$\mathbf{p}', R', S$

**How to model/aggregate appearance?**

Camera coordinate system

Q: What is the image-space projection of a 3D Gaussian?

**A: Can approximate as a 2D Gaussian!**

$$\pi(\mathbf{x}) = \mathbf{u} \qquad z \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = K \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

$\pi$: Projection function for mapping 3D points to pixels

**2D mean:** $\mu_{2D} = \pi(\mu_{3D})$

**2D covariance:**

$$J = \frac{\partial \pi}{\partial \mathbf{x}}(\mu_{3D})$$

$$\Sigma_{2D} = J \Sigma_{3D} J^T$$

1. Sort Gaussians from closest to furthest from the camera

2. For each pixel $\mathbf{u}$, compute opacity for each gaussian $\mathcal{G}_k$:

$$\mu_{2D} = \pi(\mu_{3D})$$

$$\Sigma_{2D} = J \Sigma_{3D} J^T$$

$$\bar{\alpha}_k = \alpha_k \frac{e^{-(\mathbf{u}-\mu_{2D}^k)^T (\Sigma_{2D}^k)^{-1} (\mathbf{u}-\mu_{2D}^k)}}{2\pi |\Sigma_{2D}^k|^{0.5}}$$

+ Lots of efficient GPU optimization strategies

Initialize with sparse point cloud from SfM

SfM Points   Initialization   3D Gaussians

Split/clone Gaussians based on heuristics

Under-Reconstruction   Clone   Optimization Continues

Over-Reconstruction   Split   Optimization Continues

# A glance at Gaussian Splatting.
# Finally, holy 'grail': Dynamic Scene & Using
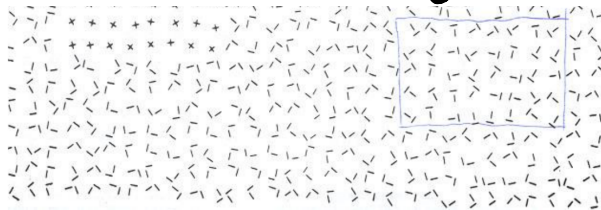## prior knowledge

# Lec 18: Texture

## Instance (实例) v.s. Category (门类)

⇒ Texture (纹理)

⭐ Texture depicts spatially repeating patterns
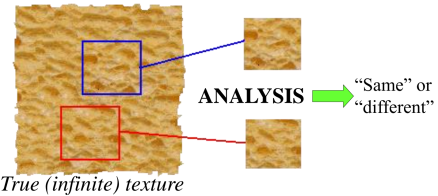- Many natural phenomena are textures

← Def

E.g. For two places on the right : Should identify as 'same'



radishes    rocks    yogurt



Human vision is sensitive to the difference of some types of elements and appears to be "numb" on other types of differences.

Human vision operates in two distinct modes:

**1. Preattentive vision**

parallel, instantaneous (~100--200ms), without scrutiny, independent of the number of patterns, covering a large visual field.

**2. Attentive vision**

serial search by focal attention in 50ms steps limited to small aperture.

- Instance recognition:
  - "Find me this particular chair again"
  - Often simple template matching works OK
    - Even better with many small templates, a.k.a. feature descriptors
- Category recognition:
  - "find me all chairs"
  - Templates don't work. Why?
  - Focus on things that might be <u>invariant</u> across the category
  - Relates to concept of "texture"

### Texture Analysis



ANALYSIS → "Same" or "different"

*True (infinite) texture*

Compare textures and decide if they're made of the same "stuff".
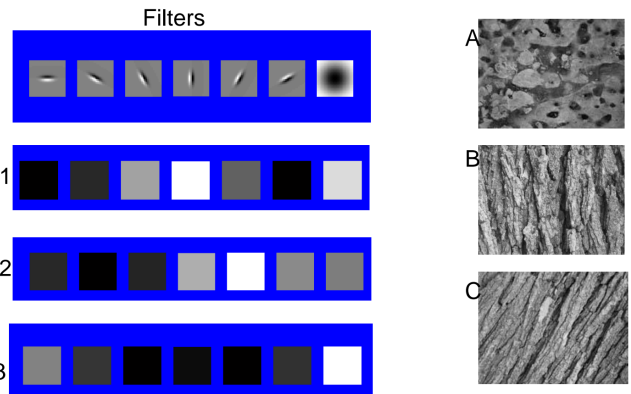
## ⚡ Julesz Conjecture:

*Textures cannot be spontaneously discriminated if they have the **same first-order and second-order statistics** of texture features (textons) and differ only in their third-order or higher-order statistics.*

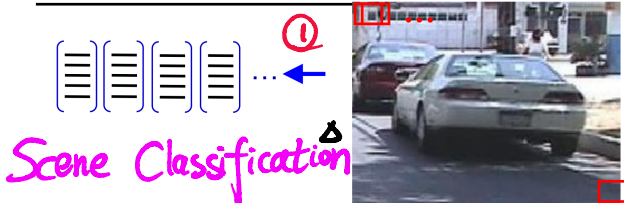## Two questions of texture modeling

- What are the texture features (textons)?
  - Pixels
  - Pixel patches
  - Outputs of V1-like filters
  - Clusters of patches / filter outputs
  - CNN features
  - Etc.

  纹理基元 ⭜

- How do we aggregate statistics
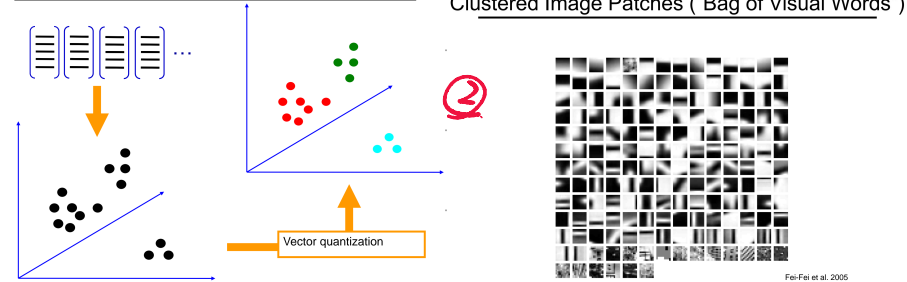  - Various types of histograms
  - Implicit or explicit

We can view texton as filter *

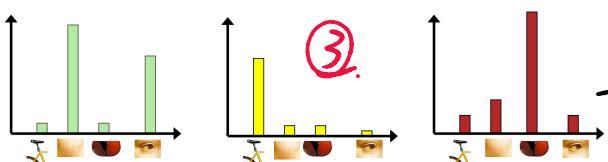Conv will indicate how image respond to this texton

### Filters



1
2
3



A
B
C

## Patch Features



①

**Scene Classification** △



### Clustering (usually k-means)



②

Vector quantization

### Clustered Image Patches ("Bag of Visual Words")



Fei-Fei et al. 2005

⚡ Brain Structure has similarities with this.



③

△: Method : Patch Feature → K-means to ① form '视觉字典' ② → 图像表示为视觉单词频率直方图 ③

# Lec 19: Image to image translation



Convolutional Neural Networks

e.g. Depth Prediction
Instead: give label of depthmap, train network to do regression (e.g., $\|z_i - \hat{z}_i\|$ where $z_i$ is the ground-truth and $\hat{z}_i$ the prediction of the network at pixel i).

Input HxWx3 RGB Image | Output HxWx1 Depth Image | True HxWx1 Depth Image

Surface Normals
Instead: train normal network to minimize $\|n_i - \hat{n}_i\|$ where $n_i$ is ground-truth and $\hat{n}_i$ prediction at pixel i.
Input: HxWx3 RGB Image | Output: HxWx3 Normals

Neural Network 现广泛应用于 CV. 其中一类任务, task 中:

$$\mathbb{R}^{H \times W \times 3/1} \longrightarrow NN \longrightarrow \mathbb{R}^{H \times W \times ?}$$ , i.e., 给每个 pixel 分配一个 feature.

## Generic Task: Image to image translation

如 Depth Prediction / Surface Normal Task / Denoise /

Semantic Segmentation / ......, 分别给 $\mathbb{R}^1, \mathbb{R}^3, \mathbb{R}^3, \mathbb{R}^{num\ of\ class}$

**"Semantic Segmentation"**
Each pixel has label, inc. **background**, and unknown
Usually visualized by colors.
Note: don't distinguish between object *instances*
Input | Label | Input | Label

Denoising neural network

Generic: Image-to-Image Translation
Labels to Street Scene | Labels to Facade | BW to Color
Aerial to Map | Day to Night | Edges to Photo

We need to:
1. Have large receptive fields to figure out what we're looking at
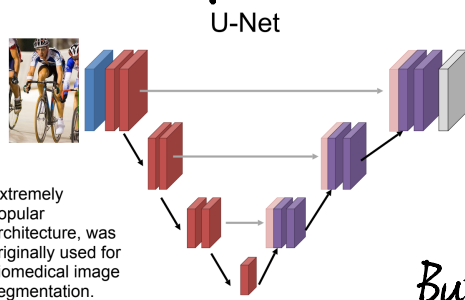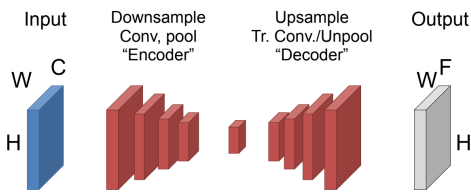2. Not waste a ton of time or memory while doing so

These two objectives are in total conflict

设计网络时, 使用卷积为主要操作。但

感受野随深度↑才↑, 而深度↑ compute↑

⇒ Conflict! How to solve it?

## Putting it Together
Convolutions + pooling downsample/compress/encode
Transpose convs./unpoolings upsample/uncompress/decode

Input | Downsample Conv, pool "Encoder" | Upsample Tr. Conv./Unpool "Decoder" | Output

W C | | | W F
H | | | H

## U-Net
Extremely popular architecture, was originally used for biomedical image segmentation

Better way : Novel

Architecture : Unet !

But, e.g., in image colorization.

## Image Colorization

Grayscale image: L channel
$\mathbf{X} \in \mathbb{R}^{H \times W \times 1}$
$L \rightarrow \mathcal{F} \rightarrow ab$

Color information: ab channels
$\hat{\mathbf{Y}} \in \mathbb{R}^{H \times W \times 2}$

## Better Loss Function
$$\theta^* = \arg\min_\theta \ell(\mathcal{F}_\theta(\mathbf{X}), \mathbf{Y})$$

- Regression with L2 loss inadequate
$$L_2(\hat{\mathbf{Y}}, \mathbf{Y}) = \frac{1}{2}\sum_{h,w}\|\mathbf{Y}_{h,w} - \hat{\mathbf{Y}}_{h,w}\|_2^2$$

- Use per-pixel multinomial classification
$$L(\hat{\mathbf{Z}}, \mathbf{Z}) = -\frac{1}{HW}\sum_{h,w}\sum_q \mathbf{Z}_{h,w,q}\log(\hat{\mathbf{Z}}_{h,w,q})$$

Colors in *ab* space
(discrete)

Loss design is critical.

## Question: Any universal loss?

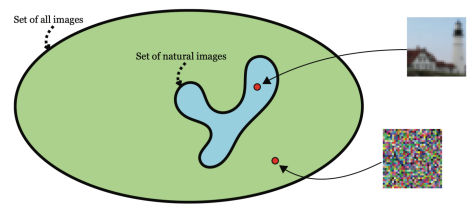有的。 Core idea : teacher-student. `teacher` component

as scaffolding.

**Generated images**

**Generative Adversarial Network (GANs)**



Generated vs Real (classifier)

Real photos

[Goodfellow, Pouget-Abadie, Mirza, Xu, Warde-Farley, Ozair, Courville, Bengio 2014]

$$\mathbf{x} \quad G \quad G(\mathbf{x}) \quad D \quad \textbf{fake } (0.9)$$

$$\mathbf{y} \quad D \quad \textbf{real } (0.1)$$

$$\arg\max_D \mathbb{E}_{\mathbf{x},\mathbf{y}} \left[ \boxed{\log D(G(\mathbf{x}))} + \boxed{\log(1 - D(\mathbf{y}))} \right]$$

**G** tries to synthesize fake images that *fool* the *best* **D**:

**G**'s perspective: **D** is a loss function.

$$\arg\min_G \max_D \mathbb{E}_{\mathbf{x},\mathbf{y}} \left[ \log D(G(\mathbf{x})) + \log(1 - D(\mathbf{y})) \right]$$

Rather than being hand-designed, it is *learned*.

$$\mathbf{x} \quad G \quad G(\mathbf{x}) \quad D \quad \textbf{fake } pair$$

$$\mathbf{x} \quad G \quad G(\mathbf{x}) \quad D \quad \textbf{real } pair$$

$$\arg\min_G \max_D \mathbb{E}_{\mathbf{x},\mathbf{y}} \left[ \log D(\mathbf{x}, G(\mathbf{x})) + \log(1 - D(\mathbf{x}, \mathbf{y})) \right]$$

$$\arg\min_G \max_D \mathbb{E}_{\mathbf{x},\mathbf{y}} \left[ \log D(\mathbf{x}, G(\mathbf{x})) + \log(1 - D(\mathbf{x}, \mathbf{y})) \right]$$

Also, we need 'pair'. Why? 这种机制称为 Conditional GAN,
判别器不仅要看输出图像是否逼真, 还要看 Y 是否与输入 X 匹配.

Lec 20 : Generative Models of Images
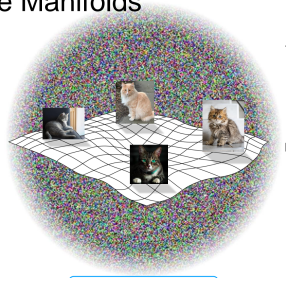
Set of all images
Set of natural images

Image中 pixel 组合与 value 取值所造成的
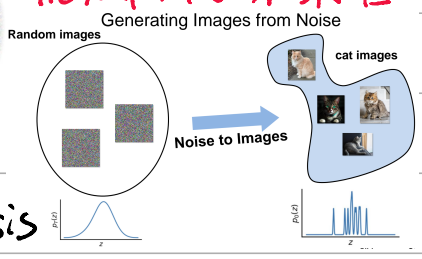permutations 繁多. 但其中很小一部分图片,
我们才认为是 '有意义的' 图像.
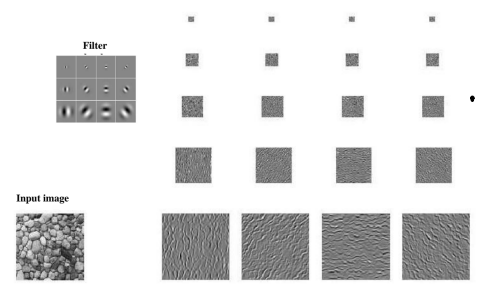
**Natural Image Manifolds**

Most images are "noise"

"Meaningful" images tend to form some manifold within the space of all images

认为有意义图像在巨大图像空间中并非随
机分布. 而聚集在一个特定的低维'流形'上

**Generating Images from Noise**

Random images

Noise to Images

cat images

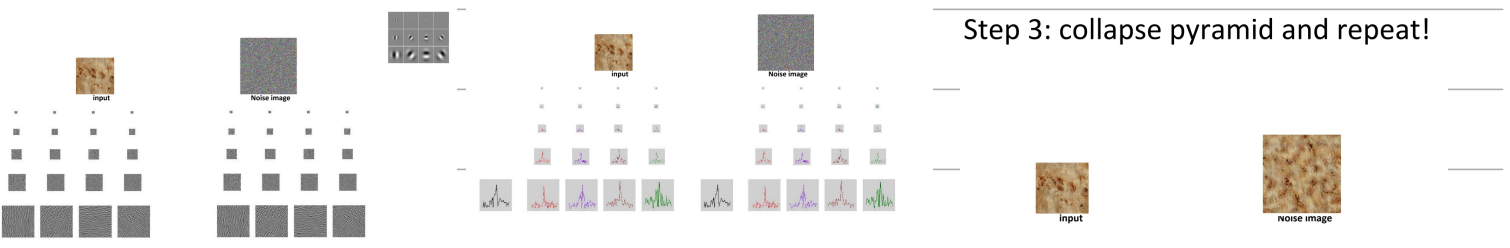**Multi-scale filter decomposition (steerable pyramid)**

Filter

Input image

从 Parametric Texture Synthesis
开始: 最早用多尺度滤波器(金字塔)来合成

Step 1: Convolve with filterbank

Step 2: match per--channel histograms

Step 3: collapse pyramid and repeat!
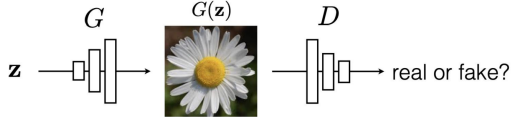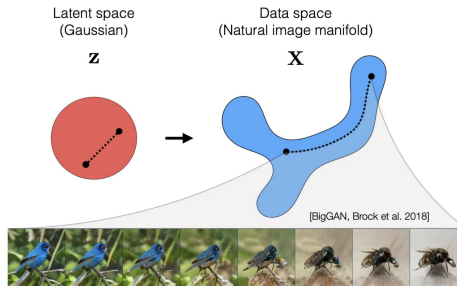
input

Noise image

input

Noise image

## GANs as generative models

- **G** tries to synthesize fake images that *fool* the *best* **D**
- **D** tries to identify the fakes



$$\arg\min_G \max_D \; \mathbb{E}_{\mathbf{z},\mathbf{x}}\big[\; \log D(G(\mathbf{z})) \;+\; \log\left(1 - D(\mathbf{x})\right) \;\big]$$

## GANs can "walk" on the manifold



Latent space (Gaussian) $\mathbf{z}$     Data space (Natural image manifold) $\mathbf{x}$

[BigGAN, Brock et al. 2018]

之后有 CNN-based
方式，甚至可实现 style
transfer 类做。
之后 GAN 盛行。

---

# Diffusion: (Recently Popular)

noise→image : hard; image→noise : *easy.*

**Key insight**

- Globally, creation is much harder than destruction

- But locally, they are almost reversable!

但局部它们加噪过
程又可逆！用 NN 预
测噪声，然后去噪



raspberry images

random images

**Denoising diffusion neural network**



**Diffusion neural network**

This network can be a U-Net or other suitable image-to-image network

raspberry images

random images

## Curious property of Diffusion

+ We are training the model to reconstruct the training set
  + But it fails!
  + Instead, it generate novel images
  + Which is what makes it great

+ Perhaps it models images as textures
  + Keeping important correlations and throwing away the rest
  + But we don't know the "model space" of these textures

Denoising: train a U-net, learn to 一步步
去噪。虽然训练目标是"重建"训练数据
(去噪)，但模型最终学会的并非死记硬背。
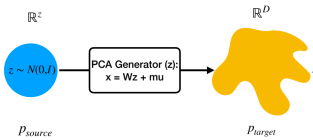Diffusion 可与 LLM 结合，以文字控制生成过程。

---

# Lec21: Flow Matching

## Generative Story

通常生成范式: latent space $\xrightarrow{\text{sample}}$ Images

- Any Generative Model has a process of sampling an image

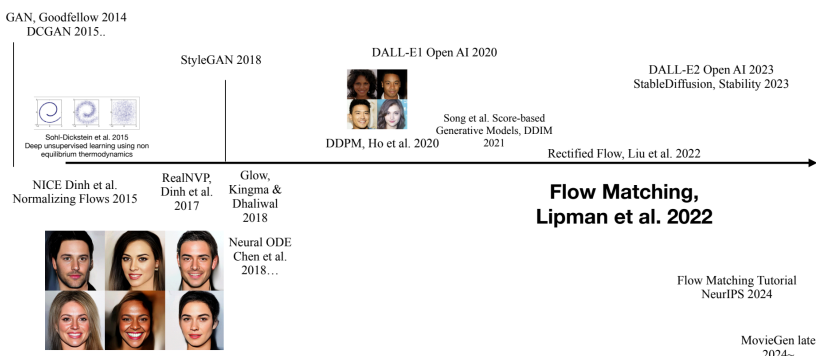- For ex, here's the generative story for PCA in its probabilistic interpretation:

1. Sample from a Gaussian Distribution
   $z \sim N(0, I)$
2. Project to Images (W = Eigenvectors, Mu = avg datapoint)

$$x = Wz + \mu$$



$z \sim N(0,I)$   PCA Generator (z): $x = Wz + mu$
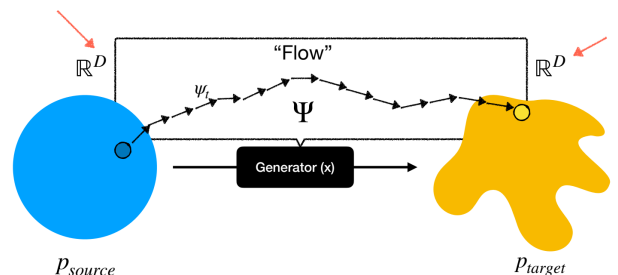
$p_{source}$     $p_{target}$

- GANs really opened up the possibility of image generation

- But people didn't like it for many reasons

  - Severe mode collapse

  - Unstable training mechanics

- Flow/Diffusion is a reactionary movement against GANs, next natural evolution

No GAN
Movement"

---

## History



GAN, Goodfellow 2014
DCGAN 2015..

StyleGAN 2018

DALL-E1 Open AI 2020

DALL-E2 Open AI 2023
StableDiffusion, Stability 2023

Sohl-Dickstein et al. 2015
Deep unsupervised learning using non equilibrium thermodynamics

Song et al. Score-based Generative Models, DDIM 2021

DDPM, Ho et al. 2020

Rectified Flow, Liu et al. 2022

NICE Dinh et al. Normalizing Flows 2015

RealNVP, Dinh et al. 2017

Glow, Kingma & Dhaliwal 2018

**Flow Matching, Lipman et al. 2022**

Neural ODE Chen et al. 2018...

Flow Matching Tutorial NeurIPS 2024

MovieGen late 2024~

## Flow based Generative Models



$\mathbb{R}^D$    "Flow"    $\mathbb{R}^D$

$\psi_t$    $\Psi$
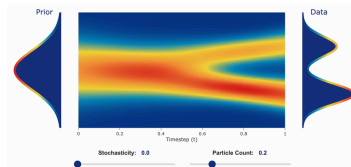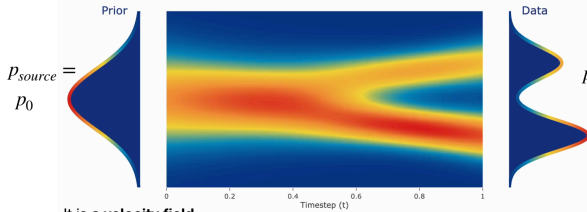
Generator (x)

$p_{source}$     $p_{target}$

1. Latent space dim is same as the target!

2. Takes T steps to go from src to tgt

但 Flow base model 不追求 sample directly 而是 flow.

## What is Flow?



$p_{source} = p_0$

$p_{target} = p_1$

- It is a **velocity field.**
- It's like a river with some currents, every point defines how fast you move (velocity)
- You ride this river to go from one distribution to next

## Riding the river = Integration



Simplest "Euler Integration":
$$x_{t+\Delta t} = x_t + v_\theta(x_t, t)\Delta t$$

- Riding this rive means you add little bits of velocity defined at each location
- This is called "Integration", also called solving the Ordinary Differential Equation (ODE) with initial state $x_0$, through some differential parametrized by a network: $\frac{dx}{dt} = v_\theta(x, t)$
- You can add stochasticity when riding it, then it becomes SDE (more next lecture)
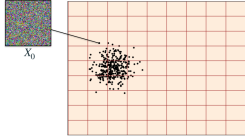
希望学习一个速度场 $u_\theta(x, t)$，通过在这个速度场中"漂流"积分.

我们可以将噪声分布（源）平滑地变换为复杂图像分布.

## The flow can be though about learning a warping function

like, gaussian noice

$X_t = \psi_t(X_0)$, $t \in [0,1]$

Warping    Source $X_0 \sim p$

## Initial approach trained flow with Maximum Likelihood

$$D_{KL}(q \| p_1) = -\mathbb{E}_{x \sim q}\log p_1(x) + c$$

$X_t = \psi_t(X_0)$, $t \in [0,1]$

Warping    Source $X_0 \sim p$

- Normalizing Flow, Continuous Normalizing Flow
- Chaining $\psi_t$ needs to satisfy the continuity equation!!!!!
- This requires ODE integration DURING training with invertible neural networks
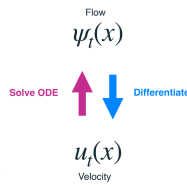
早期传统 Normalizing Flow

（归一化流）用最大化似然训练.

让 P1 分布尽可能贴合真实分布 q.

训练极其缓慢

## Previous Normalizing Flow works
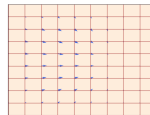
Caveat

- Tries to directly deal with this continuity equation constraint
- Very slow to train (need to integrate while training)
- Other constraints like invertibility of $\psi_t$
- Nice idea with promising results but limited capability + not practical to train

## Instead, model Flow with Velocity

Flow
$\psi_t(x)$

$\frac{d}{dt}\psi_t(x) = u_t(\psi_t(x))$

Solve ODE ↑    ↓ Differentiate

$u_t(x)$
Velocity

- **Pros**: velocities are **_linear_**
- **Cons**: simulate to sample
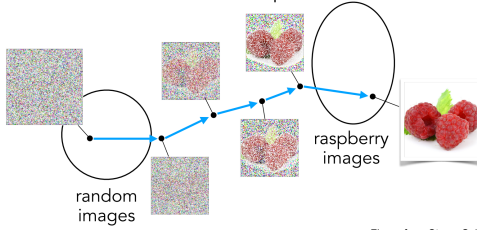
不学位置映射 $\psi_t(x)$，而学速度场 $u_t(x)$：不问"下一刻我在哪"
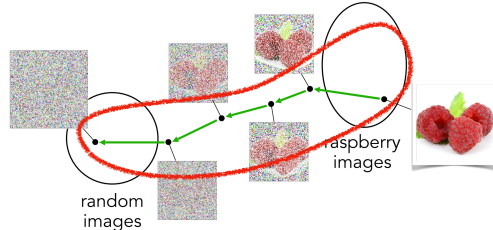
而问"我现在往哪个方向走，且速度多少"

In Flow approach:

## Training

1. Take real data, corrupt it to left distribution somehow
2. Learn to undo the process!



random images    raspberry images

## $$$ question, how to pick the intermediate path?

How to generate this Green path?



random images    raspberry images

## How to construct $x_t$

**TLDR:** Sample noise, add it, then reconstruct the data

Flow matching says you can **pick any combination**, as long as it starts from a sample in the source (e.g. gaussian) and ends with a sample in the target distribution (image)

$$x_t = \alpha_t x_0 + \sigma_t x_1$$

$x_0 \sim p_0(x)$        $x_1 \sim p_1(x)$

## Flow training



- For each data $x_1$
- Sample some noise $x_0$
- Combine it however you want to get $x_t$
- Now learn to predict the velocity at $x_t$
  - What is the velocity? Depends on how you got $x_t$

## What is the velocity supervision?

$x_t = \alpha_t x_0 + \sigma_t x_1$

$x_t = (1-t)x_0 + tx_1$

$\frac{dx_t}{dt} = -x_0 + x_1$

$= x_1 - x_0$

$$\mathbb{E}_{t, x_0, x_1}\left\|u_t^\theta(X_t) - (X_1 - X_0)\right\|^2$$



$X_t = (1-t)X_0 + tX_1$

理想速度就是 $x_1 - x_0$. 因此

loss: $\mathbb{E}_{t, x_0, x_1} \|u_t^\theta(X_t) - (X_1 - X_0)\|^2$

# Inside a Training Loop
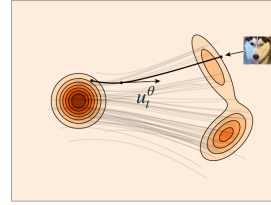## Flow Matching

```python
x = next(dataset)
t = torch.rand(1) # Sample timestep (0,1)
noise = torch.randn_like(x) # Sample noise
x_t = (1-t) * noise + (t) * x # Get noisy x_t

flow_pred = model(x_t, t) # Predict noise in x_t
flow_gt = x - noise # ground truth flow (w/ linear sched)
loss = F.mse_loss(flow_pred, flow_gt) # Update model
loss.backward()
optimizer.step()
```

## During inference:

- Just take a small step in the velocity

- Use any ODE Solver, i.e. integration you like, like Euler integration:

$$x_{t+\Delta t} = x_t + \Delta t \cdot \frac{dx}{dt}\Big|_{x_t,t}$$



**Sample**
from $X_0 \sim p$

## Training: Model parameterization

- You can make your network output undo the noise in many different ways, predicting x, v, noise, or flow

$$v_t = \alpha_t x_1 - \sigma_t x_0 \qquad \begin{aligned} u_t &= x_1 - x_0 = \epsilon - x_0 \\ u_t &= x_t - x_0 \end{aligned}$$

- These are all equivalent because of the linear relationship with $x_t$. You can derive all of these as long as you know one of them

$$x_t = \alpha_t x_0 + \sigma_t x_1$$

- For example

$$\mathbb{E}[(\hat{\mathbf{x}}_0 - \mathbf{x}_0)^2] = \mathbb{E}\left[\left(\frac{\mathbf{x}_t - \sigma(t)\hat{\epsilon}}{\alpha(t)} - \frac{\mathbf{x}_t - \sigma(t)\epsilon}{\alpha(t)}\right)^2\right] = \mathbb{E}\left[\frac{\sigma(t)^2}{\alpha(t)^2}(\hat{\epsilon} - \epsilon)^2\right].$$

## Training: Flow Matching vs. Diffusion

**Algorithm 1: Flow Matching training.**
**Input** : dataset $q$, noise $p$
Initialize $v^\theta$
**while** *not converged* **do**
  $t \sim \mathcal{U}([0,1])$    ▷ sample time
  $x_1 \sim q(x_1)$    ▷ sample data
  $x_0 \sim p(x_0)$    ▷ sample noise
  $x_t = \Psi_t(x_0|x_1)$    ▷ conditional flow
  Gradient step with $\nabla_\theta \|v_t^\theta(x_t) - \dot{x}_t\|^2$
**Output:** $v^\theta$

$p_t(x_t|x_1)$ general
$p(x_0)$ is general

**Algorithm 2: Diffusion training.**
**Input** : dataset $q$, noise $p$
Initialize $s^\theta$
**while** *not converged* **do**
  $t \sim \mathcal{U}([0,1])$    ▷ sample time
  $x_1 \sim q(x_1)$    ▷ sample data
  $x_t = p_t(x_t|x_1)$    ▷ sample conditional prob
  Gradient step with
  $\nabla_\theta \|s_t^\theta(x_t) - \nabla_{x_t} \log p_t(x_t|x_1)\|^2$
**Output:** $v^\theta$

$p_t(x_t|x_1)$ closed-form from of SDE $dx_t = f_t dt + g_t dw$
- *Variance Exploding*: $p_t(x|x_1) = \mathcal{N}(x|x_1, \sigma_{1-t}^2 I)$
- *Variance Preserving*: $p_t(x|x_1) = \mathcal{N}(x|\alpha_{1-t} x_1, (1-\alpha_{1-t}^2)I)$
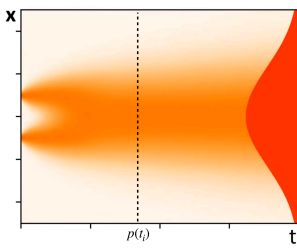    $\alpha_t = e^{-\frac{1}{2}T(t)}$

$p(x_0)$ is Gaussian
$p_0(\cdot|x_1) \approx p$

Algorithm: Flow Matching: 采样一个 noise $x_0$ 与 data $x_1$，强制
中间状态在两点连线上，而训练目标：$v_\theta(x_t) \to x_1 - x_0$

Diffusion：采样$x_1$，依 $p_t(x_t|x_1)$ 随机加噪得$x_t$
训练目标：预测噪声，即 $\nabla \log p_t$

Lec22: Diffusion Sampling

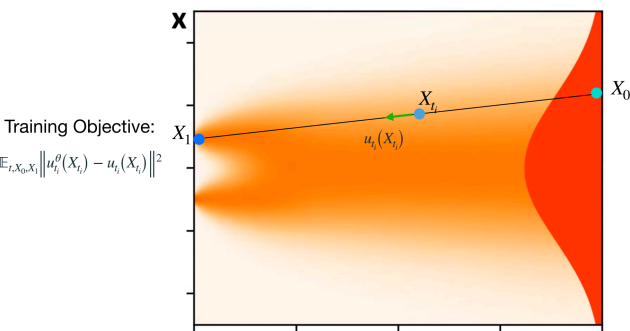## Revisit diffusion models with a 1D example



t: Miika Aittala

$p(t_i)$

Recall in 1D example: how to: noise → distribution?
We need to know the 'speed vector'. i.e., 'flow',
of $X_{t_i} \Rightarrow u_{t_i}(X_{t_i})$    We may use a model to
represent this: $u_{t_i}^\theta(X_{t_i})$
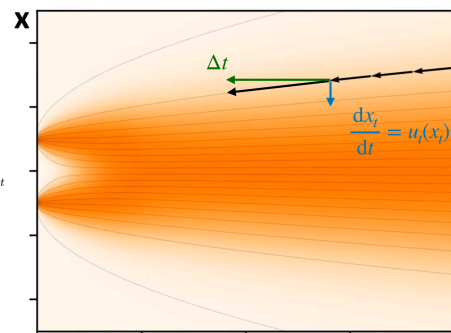So objective: $\mathbb{E}_{t, x_0, x_1} \| u_{t_i}^\theta(X_{t_i}) - u_{t_i}(X_{t_i})\|_2$



**Training Objective:**
$\mathbb{E}_{t, x_0, x_1} \|u_{t_i}^\theta(X_{t_i}) - u_{t_i}(X_{t_i})\|^2$

$X_1$   $X_t$   $u_t(X_t)$   $X_0$

## Solving the flow ODE with discretization



Euler step:
$x_{t+\Delta t} = x_t + \Delta t \cdot \frac{dx}{dt}\Big|_{x_t,t}$

$\Delta t$

$\frac{dx_t}{dt} = u_t(x_t)$

With this model who can give $\frac{dx}{dt}\Big|_{x_t,t}$
We can flow in a way:
$X_{t+\Delta t} = X_t + \Delta t \cdot \frac{dx}{dt}\Big|_{x_t,t}$

Truncation error: fails to approximate ideal trajectory by finite steps.

This is not a perfect approach
⇐ 离散化 can be problematic
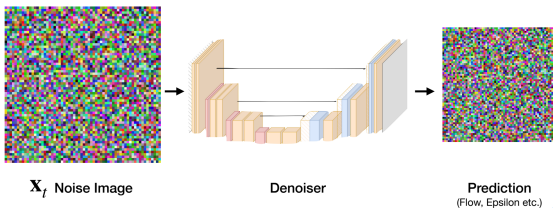Model may give $\frac{dx}{dt}\big|_{x_t, t}$ ⇒ which is a little bit inaccurate

Model fails to approximate the marginal flow.

1. Naive solution: sampling with more steps.

2. Time steps are long at high noise levels and short at low noise levels

3. Higher-order ODE solver

该 sol 很巧妙！

1. Stochastic sampler (SDE) injects fresh noise throughout the evolution in addition to reducing the noise.

"给模型一些迂回的余地！" $x_{t+1}$
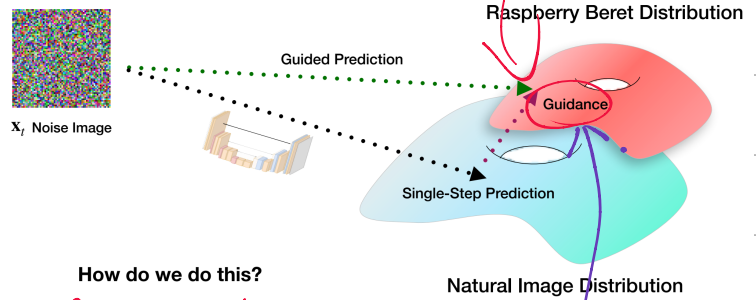
Diffusion！

$\hat{x}_t$

A specific implementation $x_t$ of 'flow' idea

Without condition injection and we want conditioned generated image, 生成image distribute over entire natural image space！

⇒ We want to inject a guidance ......
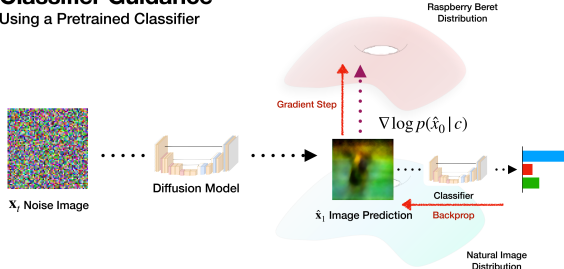
**Diffusion Guidance**
Motivation



$\mathbf{x}_t$ Noise Image     Denoiser     Prediction (Flow, Epsilon etc.)

Raspberry Beret Distribution



$\mathbf{x}_t$ Noise Image     Denoiser     Prediction

Unguided Diffusion Isn't Too Useful!     Natural Image Distribution

**Diffusion Guidance**
Push Toward a Conditional Mode



Raspberry Beret Distribution
Guided Prediction
$\mathbf{x}_t$ Noise Image
Guidance
Single-Step Prediction
Natural Image Distribution

How? Two approaches

**Original**
Classifier Guidance
Guide with a pretrained classifier.

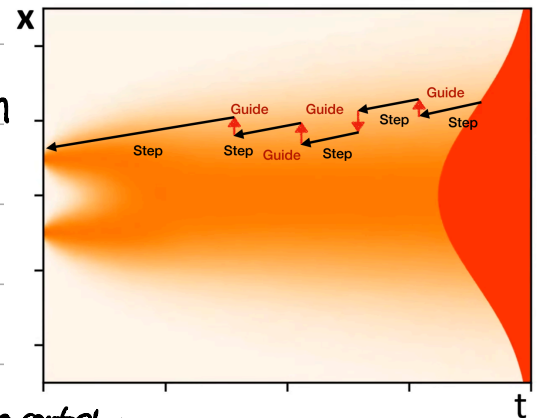**Current**
Classifier-Free Guidance
Guide a diffusion model with itself.

**How do we do this?**

Classifier Guidance
For a generated step image, use a classifier to generate guidance vector.

**Classifier Guidance**
Using a Pretrained Classifier



Raspberry Beret Distribution
Gradient Step
$\nabla \log p(\hat{x}_0 | c)$
$\mathbf{x}_t$ Noise Image     Diffusion Model     $\hat{x}_1$ Image Prediction     Classifier     Backprop
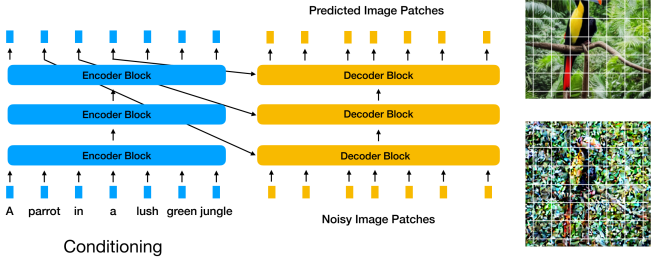Natural Image Distribution

But this approach is not favored: ① Images in the middle: OOD for classifier, so guidance vector may not be reliable！

② Train these two models jointly or separately are both trouble-some！ (Hard to train)
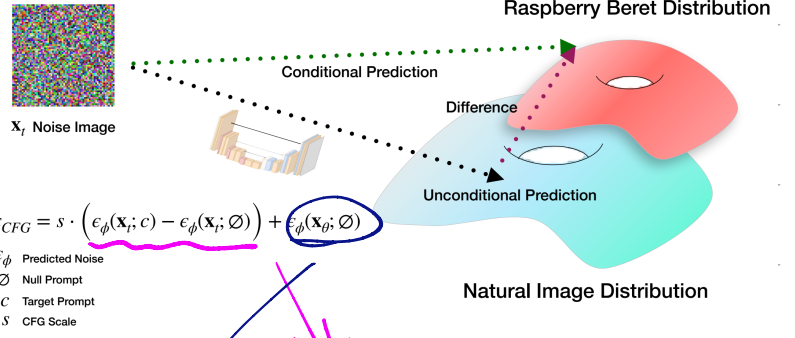


X
Guide     Guide     Guide
Step     Step     Step
Step     Guide     Step
Step
t

## Diffusion Transformer Architecture
DiT, PixArt Alpha, MMDiT, etc.



Predicted Image Patches

| Encoder Block | Decoder Block |
| Encoder Block | Decoder Block |
| Encoder Block | Decoder Block |

A parrot in a lush green jungle

Noisy Image Patches

Conditioning

Approach 1: With Transformer's encoder-decoder paradigm.

Model Knows Unconditional, Text-Conditioned Distributions



$X_t$ Noise Image

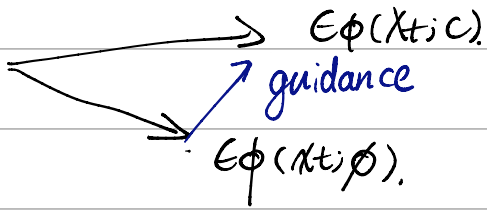Raspberry Beret Distribution

Conditional Prediction

Difference

Unconditional Prediction

Natural Image Distribution

$$\epsilon_{CFG} = s \cdot \left( \epsilon_\phi(\mathbf{x}_i; c) - \epsilon_\phi(\mathbf{x}_i; \varnothing) \right) + \epsilon_\theta(\mathbf{x}_\theta; \varnothing)$$

$\epsilon_\phi$ Predicted Noise
$\varnothing$ Null Prompt
$c$ Target Prompt
$s$ CFG Scale

Approach 2: Use UNet

$\epsilon_\phi(\;\underset{\sim}{;}\;\underset{\sim}{;}\;)$    image. condition

At one step, generate

both : $\epsilon_\phi(x_t; c)$ & $\epsilon_\phi(x_t; \varnothing)$

Their gap can be a guidance!

$\rightarrow \epsilon_\phi(x_t; c)$    And amplify $\epsilon_\phi(x_t; c) - \epsilon_\phi(x_t; \varnothing)$

guidance    让生成图片强烈地被推向 guidance 方向.

$\rightarrow \epsilon_\phi(x_t; \varnothing)$.