Announce: Fill out survey
Extra Credit for everyone (3%)
If 75% of class does the survey
Due: Fri of Thanksgiving Week.
Submit Draft Reports today
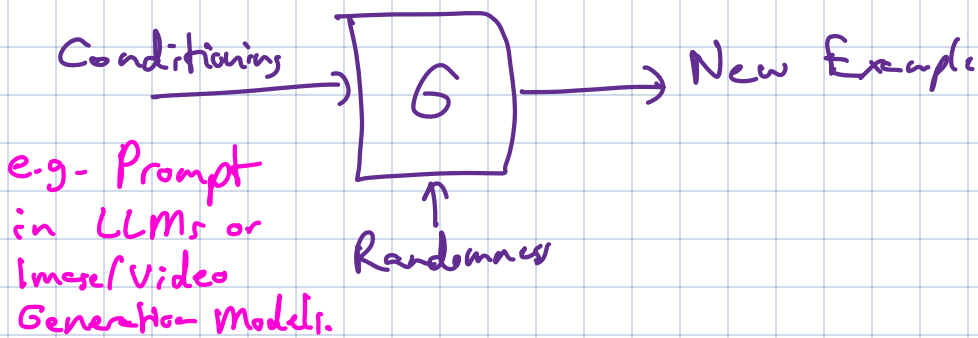   & do reviews this weekend

# Generative Models

← Understood/Conceptualized as sampling from an unknown distribution.

## Unconditional



→ New Example

Randomness

We'll use these for teaching the core ideas.

## Conditional

Conditioning → G → New Example

Randomness

This is the most practically useful setting.

e.g. Prompt in LLMs or Image/Video Generation Models.

# Ideas that don't work

## A) Use a classifier

Image $\mathbb{R}^{50 \times 50 \times 3}$ → Cat Classifier → $\mathbb{R}$ "cat score"
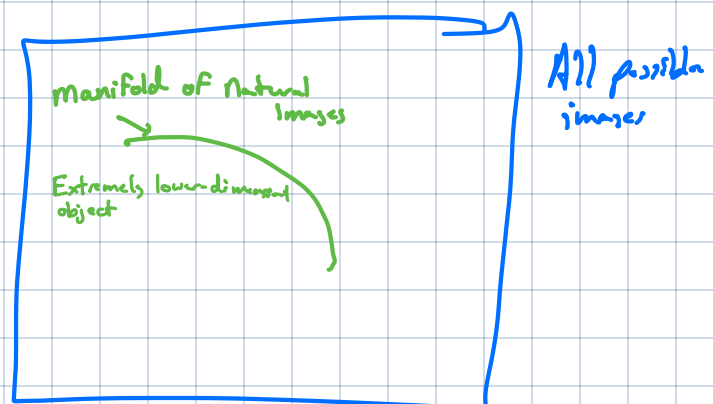
Try: Random Uniform Image
Followed by Gradient Ascent
on the Cat Score

Result: Noise-like image that classifier confidently classifies as a cat.

☹

manifold of Natural Images

Extremely low-dimensional object

All possible images

# B) Use an autoencoder

Core Ingredients: Labels are $\vec{x}_i$ itself.
Architecture has an encoder followed by a decoder.
Bottleneck in the middle.

$\vec{x} \rightarrow \boxed{E} \xrightarrow{\vec{z}} \boxed{D} \rightarrow \hat{x}$    Traditional Perspective: Decoder is scaffolding.

## Try using $D$ to generate samples.

Try: Draw $\vec{z}$ from a random distribution. — If $\vec{z}$ too small, get blurry junk
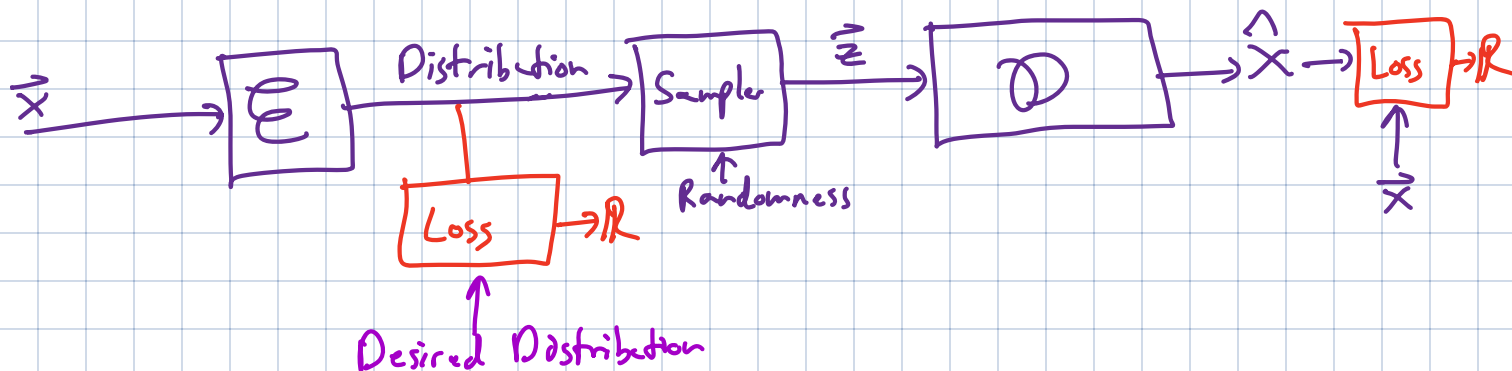If $\vec{z}$ bess, "Noise-like Imags."

What went wrong? The $\vec{z}$ that was random
was nothing like the $\vec{z}$ seen during training.

## VAE Approach

3 key ingredients:  1) Make $\vec{z}$ random during training too.

2) Add a loss on <u>distribution of $\vec{z}$</u>

3) Make this work with SGD

$\vec{x} \rightarrow \boxed{E} \xrightarrow{\text{Distribution}} \boxed{\text{Sample}} \xrightarrow{\vec{z}} \boxed{D} \rightarrow \hat{x} \rightarrow \boxed{\text{Loss}} \rightarrow \mathbb{R}$

Loss $\rightarrow \mathbb{R}$

Randomness

Desired Distribution

$\hat{x}$ ... $\vec{x}$

Desired Property of Distribution:  A) Continuous
B) Easy to sample
C) Easy to compute loss

Loss on Distributions...

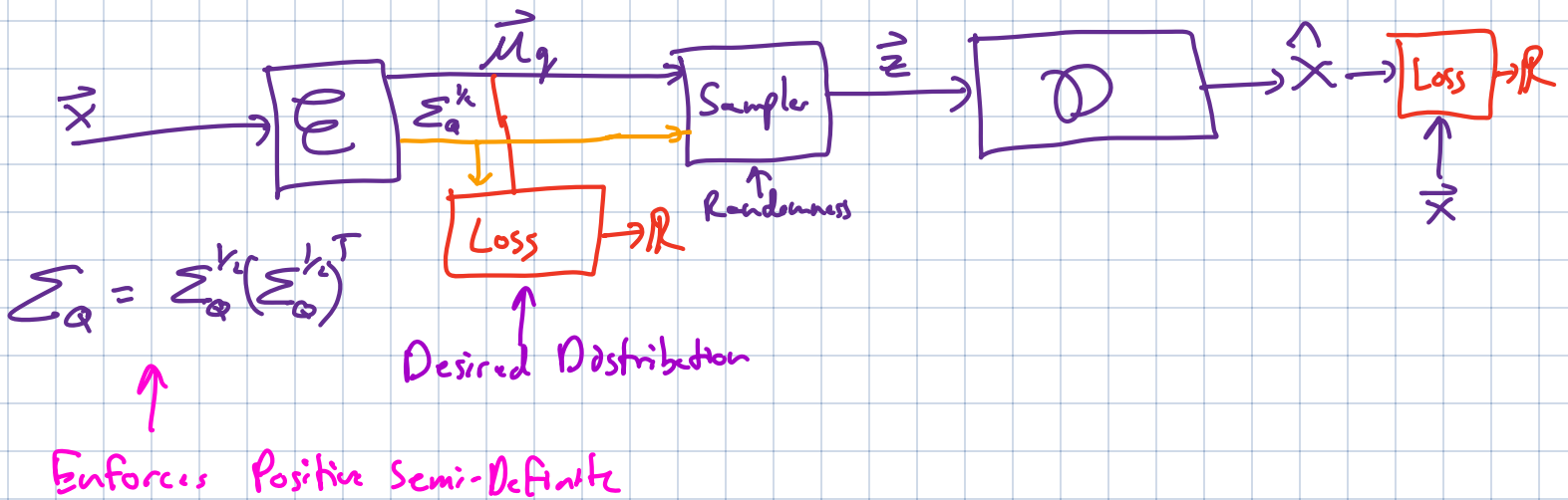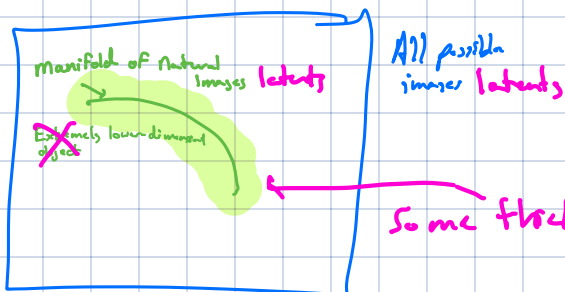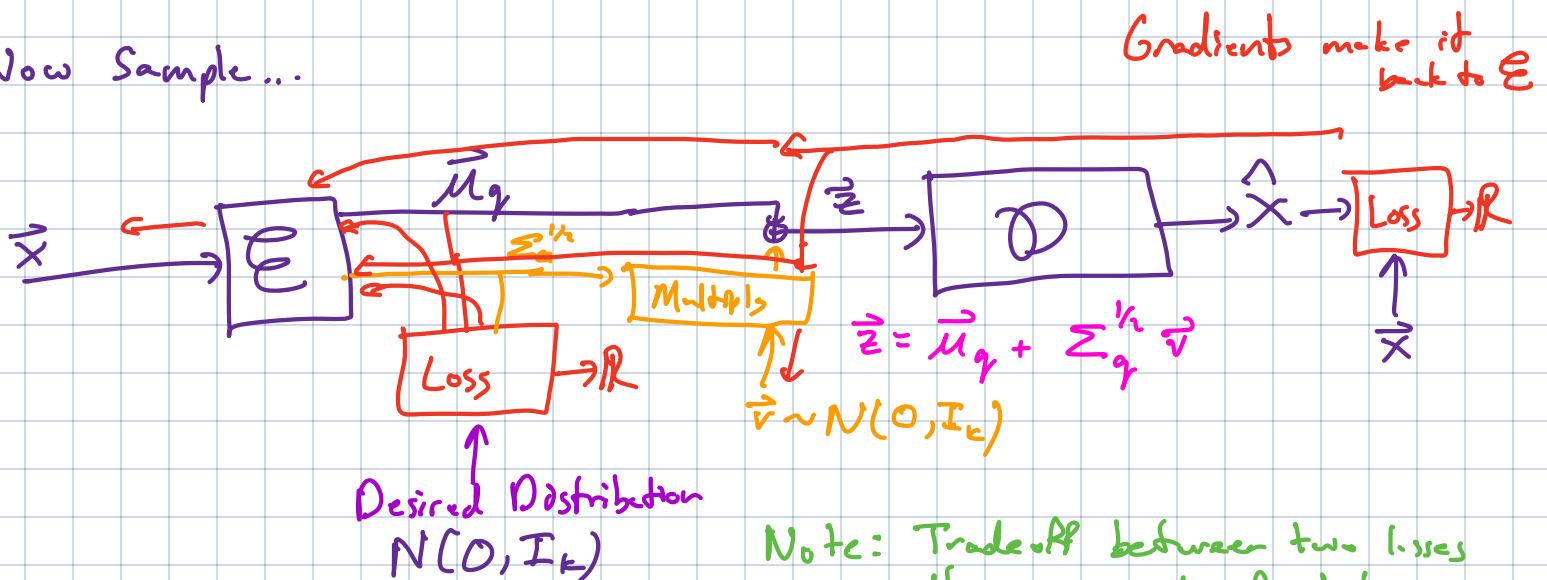KL Divergence  $KL(Q\|P) = \int Q(z) \ln \frac{Q(z)}{P(z)} \, dz$    Asymmetric, but for a good reason.

Our choice for distribution: $N(0, I_k)$

$$KL\left(N(\vec{\mu}_q, \Sigma_Q) \,\|\, N(0, I_k)\right) = \tfrac{1}{2} Tr(\Sigma_Q) + \vec{\mu}_q^T \vec{\mu}_q - k - \log \det \Sigma_Q$$



$$\Sigma_Q = \Sigma_Q^{1/2}(\Sigma_Q^{1/2})^T$$

Enforces Positive Semi-Definite

Desired Distribution

Randomness

Now Sample...

Gradients make it back to $\mathcal{E}$

$$\vec{z} = \vec{\mu}_q + \Sigma_q^{1/2}\vec{v}$$

$$\vec{v} \sim N(0, I_k)$$

Desired Distribution
$N(0, I_k)$

Note: Tradeoff between two losses
Hyperparameter for balance.

Distribution loss enforces
information bottleneck from $\vec{x}$
to $\vec{z}$

manifold of natural images latents

All possible images latents

Extremely lower-dimensional object

Some thickness due to VAE training.

Key New Tricks : 1) Using a KL Divergence Regularizer
on a distribution

2) Treating Sampling in a way that
allows gradients to just pass through it.

3) Accepting the stochasticity of random
noise in sampling as just more
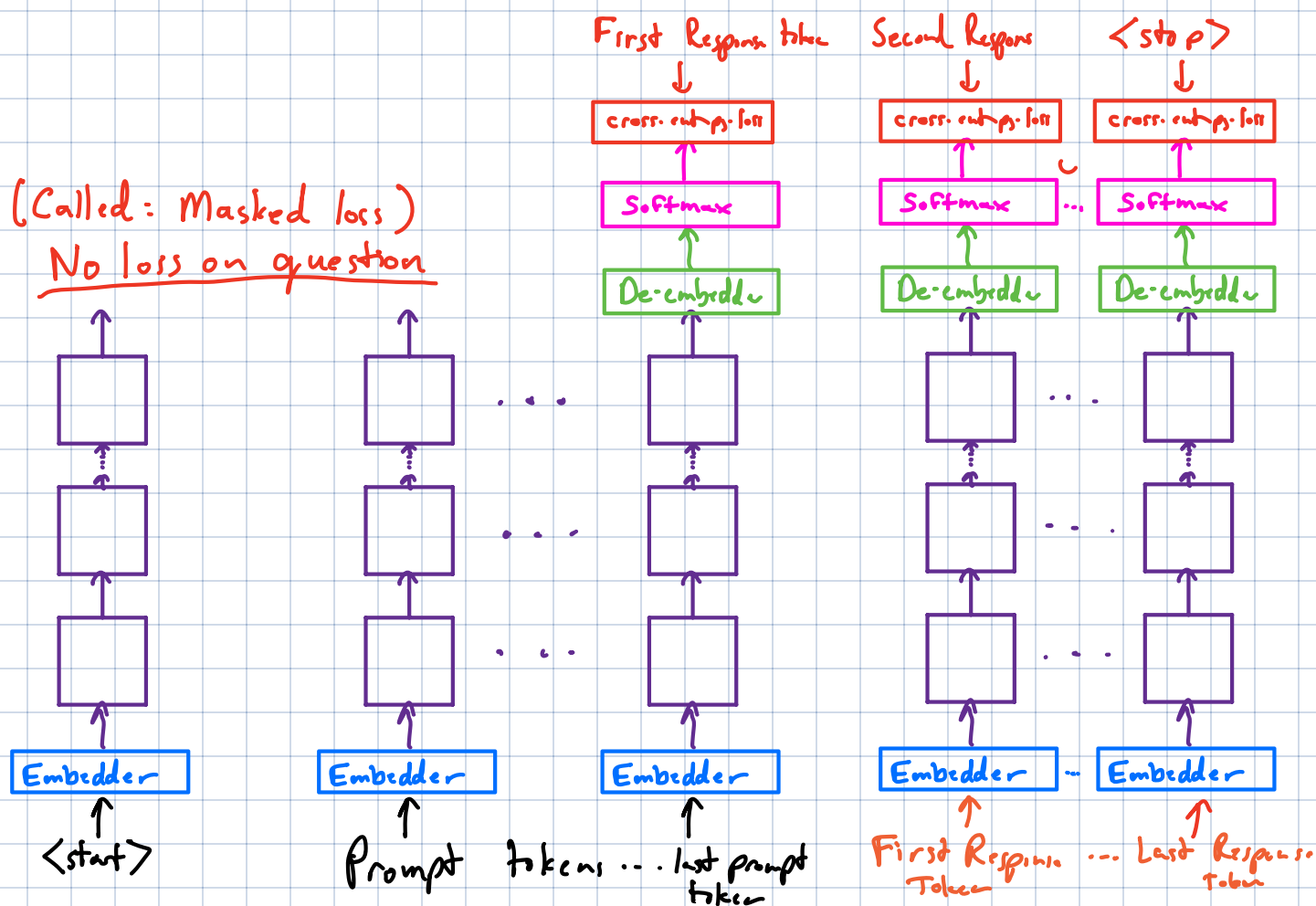stochasticity in SGD-style optimization.

These tricks are useful beyond this VAE setting.

e.g. Quantization-Aware-Training (QAT) uses (2)
during training to train a model whose performance
does not degrade much when quantized.
Keep weights in high-precision, but quantize in forward-pass.
Ordinarily, quantization is non-differentiable/gives zero-gradients.
But pretend it is added noise and pass gradient thru

You also saw in homework how introducing noise the
right way lets you get gradients for some otherwise
non-differentiable losses.

Post-training including RLVR.

Recall basic SFT for instruction following.



(Called: Masked loss)
No loss on question

First Response token    Second Response    <stop>

For LLMs, one key advance was "Chain-of-thought." So it is useful to have examples that show their work while solving a problem.

But how can we make a model better at solving problems?

Two parts to the answer:

A) Be willing to spend more compute while answering.
↖ test time compute

B) Train it to be better.

<u>Test-time compute:</u>  0) Oldest Approach: pure prompting.

"Think step by step. Be careful."

1) Repeated generation. Instead of one try, generate N responses.

At Inference/Test time : A) Take Majority/Plurality Vote
                                       B) Take highest prob

2) Sample better...

Observation (Empirical): RLVR-ed models improved performance and showed generalization/transfer of reasoning to new domains. But people noticed that their outputs were also higher-prob under the base models (without RLVR tuning) as well. So how much of this gain is distribution shaping?

Recent Paper: Karan & Du, "Reasoning with Sampling: Your base model is smarter than you think" Oct 16, 2025.

We want a good answer. That includes thinking. If we sample from original model, what can go wrong? Moving forward one token at a time, we can make a mistake... and then flounder.

Older Approach: Beam-Search with some k.

To do what? Sample from higher-likelihood <u>sequences.</u>

Alternative: sample not from $p(\vec{x})$
                  but from normalized $(p(\vec{x}))^{\alpha}$ when $\alpha > 1$