## Week 0e: Generative AI
### The Creation Challenge

Machine Learning for Smarter Innovation

BSc-Level Course

October 7, 2025

# Outline

1. Act 1: The Challenge

2. Act 2: Variational Autoencoders

3. Act 3: Adversarial & Diffusion

4. Act 4: Synthesis

**Traditional ML:** "What is this?"

- Email spam detector: Classify existing emails
- Medical diagnosis: Analyze X-ray images
- Sentiment analysis: Judge customer reviews

**Limitation:** Only analyzes, never creates

**Generative AI:** "Create something new"

- Generate phishing emails for security training
- Synthesize medical images for rare diseases
- Write product descriptions automatically
- Compose music for video backgrounds

**Power:** Creation enables innovation

Fundamental shift: from pattern recognition to content generation

# Mathematical Foundation
Two Approaches to Learning

**Discriminative Models**
Learn: $P(y|x)$ - Conditional probability
**What it does:**

- Given $x$, predict label $y$
- Learns decision boundaries
- Divides input space

**Examples:** Logistic, RF, SVM
**Can sample new $x$?** NO - only classifies existing data

**Generative Models**
Learn: $P(x)$ - Joint or marginal distribution
**What it does:**

- Models entire data distribution
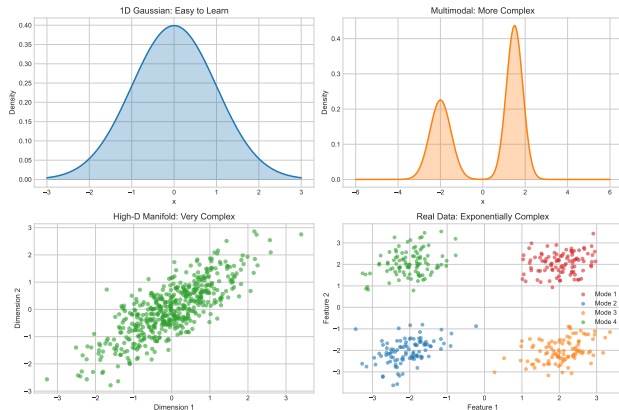- Samples new $x \sim P(x)$
- Creates novel instances

**Examples:** VAEs, GANs, Diffusion
**Can sample new $x$?** YES - generates from distribution

Key distinction: Discriminative draws boundaries, Generative learns distributions enabling sampling

**Challenges:**

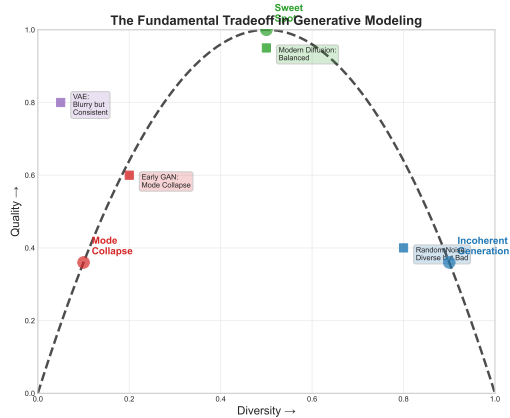- High-dimensional spaces

- Multimodal distributions

Real data lives on complex manifolds - learning full distribution is exponentially hard

**Requirements:**

- Capture all patterns

- Maintain realism

The Fundamental Tradeoff in Generative Modeling

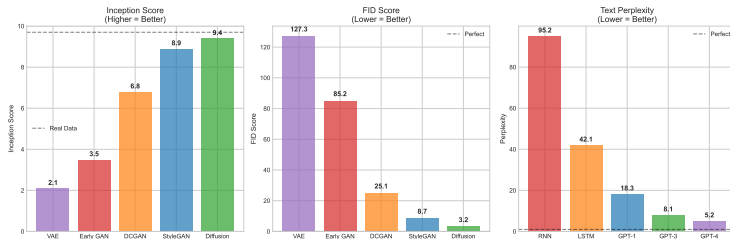**High Quality:** Mode collapse, repetitive
**Balanced:** Realistic variety
**High Diversity:** Unrealistic

Realistic AND diverse remains the central challenge

Inception Score (Higher = Better) — FID Score (Lower = Better) — Text Perplexity (Lower = Better)

**Inception Score (IS)**
- Range: 1-1000
- Higher = better
- Quality & diversity

Interpretation:
- >300: Excellent
- 100-300: Good
- <100: Poor

**FID Score**
- Range: 0-500
- Lower = better
- Feature distance

Interpretation:
- <10: Photorealistic
- 10-50: Good quality
- >50: Noticeable artifacts

**Perplexity (Text)**
- Range: 1-10,000
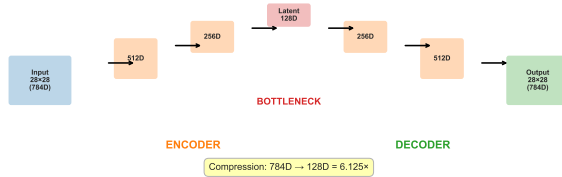- Lower = better
- Predictability

Interpretation:
- <20: Human-like
- 20-100: Coherent
- >100: Gibberish

Quantitative metrics enable objective quality assessment and model comparison

**Autoencoder Architecture: Compression Through Reconstruction**



BOTTLENECK

ENCODER          DECODER

Compression: 784D → 128D = 6.125×

**Encoder**
- 784D -¿ 128D
- Forces selective encoding
- Filters noise

**Latent**
- 128D bottleneck
- Key features only
- 6.1x compressed

**Decoder**
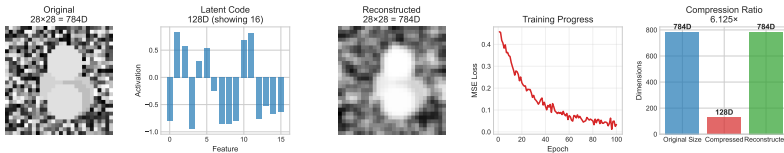- 128D -¿ 784D
- Lossy reconstruction
- Preserves essentials

Bottleneck forces meaningful compression

**Architecture:**
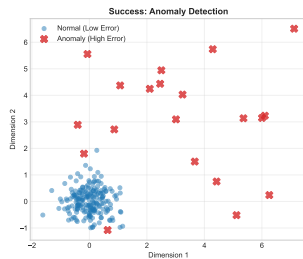
- Input: 784 pixels
- Encoder: 784 -¿ 128
- Decoder: 128 -¿ 784

**Training:**

- Loss: $L = ||x - \hat{x}||^2$
- Optimizer: Adam
- Compression: 6.125x

MSE drops 0.45 -¿ 0.03 over 100 epochs

Success: Efficient Dimensionality Reduction

Success: Learned Meaningful Features

Success: Noise Removal

Success: Anomaly Detection

[+] SUCCESSES:

Results:

Failure: Blurry Outputs

Failure: Poor Generation Quality

Failure: Holes in Latent Space

Failure: Mode Collapse / Limited Diversity

[-] FAILURES:

Metrics:

| IS | 2.1 |

**MSE Loss Forces Averaging**



Given two inputs $x_1$ and $x_2$

**MSE optimal reconstruction:** $\hat{x} = \frac{x_1 + x_2}{2}$

Result: Blurry average, not realistic sample

**Averaging in Distribution Space**



**MSE Loss: Convex (Forces Average)**



**Problem:**

**Math:**

# Variational Autoencoders (VAEs)
The Probabilistic Solution

**VAE Framework: Probabilistic Encoder-Decoder**

Reparameterization Trick



Probabilistic Latent Space

L = -E[log p(x|z)] + KL(q(z|x)||p(z))

**Key Innovation:**

- Encode to distribution: $q_\phi(z|x) = \mathcal{N}(\mu, \sigma^2)$
- Sample: $z = \mu + \sigma \odot \epsilon$

**Reparameterization:**

- Make $z$ deterministic
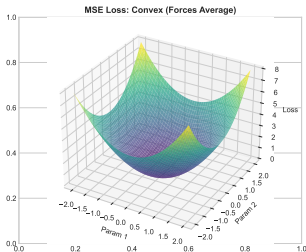- Gradient flows

Reparameterization enables gradient optimization

**VAE Loss:**

$$\mathcal{L} = -E[\log p(x|z)] + KL(q||p)$$

**Two terms:**

- Reconstruction
- KL regularization
- $\beta$-VAE balances

How Artists Improve Through Critique → GANs



**Adversarial Learning Cycle**

Student
(Generator)

Teacher
(Discriminator)

1. Creates

2. Evaluates

Artwork

4. Improves

3. Critiques

Critique

*Both Student and Teacher Improve Through Competition*

**Art Education:**

- Student creates
- Teacher critiques
- Student improves
  Adversarial learning inspired GANs.

**Insights:**

- Adversarial feedback drives improvement
- Both improve together

**Two Revolutionary Approaches to Generation**

**Adversarial Training**

Two Networks Compete

Generator ⟷ Discriminator
Compete

+ Sharp, realistic outputs

- Training instability

*Best for: Image generation*

**Diffusion Models**

Iterative Denoising

Noise → Clean (1000 steps)

+ Stable training

- Slow sampling

*Best for: Highest quality*

**Adversarial**
- Two networks compete
- Sharp, realistic

**Diffusion**
- Iterative denoising
- Stable, controllable

Both address VAE limitations

# GANs: The Forger vs Detective Game
Adversarial Training in Plain English

**GAN: The Forger vs Detective Game**



FORGER
(Generator Network)
Creates fake paintings

Fake Painting

DETECTIVE
(Discriminator Network)
Spots fakes vs real

*Feedback: "Too obvious!"*

Real Painting

Early Training: Detective wins easily

Late Training: Forger fools detective!

Equilibrium: 50% accuracy (perfect balance)

**Forger:**
- Creates fakes
- Fools detective

**Result:** Detective can't tell fake from real!

Competition drives both to excellence

**Detective:**
- Examines: real/fake?
- Gets better at detection

**Diffusion: The Reverse Corruption Process**

**Forward Process: Add Noise** | **Reverse Process: Remove Noise**



$q(x_t \mid x_{t-1}) = N(sqrt(1-beta_t) x_{t-1}, beta_t I)$

$p\_theta(x_{t-1} \mid x_t) - \text{Neural network predicts noise}$

Gradually corrupt clean data

*1000 tiny steps*

Learn to reverse corruption

*1000 denoising steps*

**Forward:**

- Clean -¿ noise
- 1000 steps

**Reverse:**

- Noise -¿ clean
- 1000 steps

**Key:** Learn to undo corruption

Like sculptor revealing statue

GAN Dynamics: Generator Learns to Match Real Distribution

**Generator:**

- Maps $z$ to $x$
- Loss: $-\log D(G(z))$

Equilibrium: Generator = Real, D accuracy = 50%

**Discriminator:**

- Separates real/fake
- Loss: $-\log D(x) - \log(1 - D(G))$

# GAN Training: Step-by-Step Example
Real Loss Values from MNIST Training



Loss Convergence Over Training



Discriminator Performance

**Training Progress Metrics**



Generation Quality Improvement

| Epoch | D_loss | G_loss | D_acc | FID |
|-------|--------|--------|-------|-----|
| 1 | 1.386 | 0.693 | 95% | 450 |
| 25 | 0.8 | 1.2 | 65% | 120 |
| 50 | 0.72 | 0.85 | 55% | 35 |
| 100 | 0.695 | 0.698 | 51% | 8.7 |

**Epoch 1:**

- D: 1.386, G: 0.693
- Images: noise

**Epoch 100:**

- D: 0.695, G: 0.698
- Images: realistic

# Diffusion Mathematical Framework
## Forward and Reverse Processes



Diffusion Mathematical Framework

**Forward:**

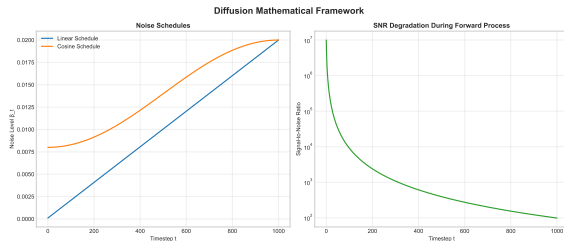$$q(x_t|x_{t-1}) = \mathcal{N}(\sqrt{1 - \beta_t}x_{t-1}, \beta_t I)$$

**Noise Schedule:**

- Linear: 0.0001 -¿ 0.02
- Cosine: Variable rate
- Matters: Smooth degradation

Linear noise schedule works for most cases

**Reverse:**

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(\mu_\theta, \Sigma_\theta)$$

**Training:**

$$L = E[||\epsilon - \epsilon_\theta(x_t, t)||^2]$$

**Intuition:** Predict noise, subtract it

# Latent Space Interpolation
Smooth Transitions in Generated Content



Latent Space Interpolation: Smooth Transitions

**Method:**

- Sample $z_1, z_2$
- Interpolate: $z_t = (1 - t)z_1 + tz_2$
- Generate: $x_t = G(z_t)$

**Applications:**

- Style transfer
- Face morphing
- Drug discovery

Meaningful latent spaces enable smooth interpolation

Diffusion Denoising: From Noise to Image in 1000 Steps

**Steps:**

- T=1000: Noise
- T=500: Structure
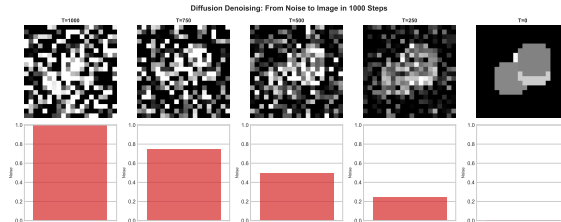- T=0: High quality

**Control:**

- Guidance scale
- Step count

Gradual refinement

Minimax Game Surface

Convergence to Nash Equilibrium

Jensen-Shannon Divergence Minimization

**Theoretical Guarantee**

At equilibrium:

p_generator = p_data

D(x) = 0.5  (50% accuracy)

*Discriminator cannot tell real from fake*

**Theory:**

- Minimax convergence
- Equilibrium: $p_g = p_{data}$

**Benefits:**

- Sharp, realistic
- Fine details

**Inception Score Over Training**

**FID Score Over Training**

**Time to Convergence**

**Quality-Speed Tradeoff**

**Observations:**

- Diffusion: Best
- GAN: 4x faster
- VAE: Fast, blurry

**Results (MNIST):**

| Method | IS | FID | Time |
|--------|-----|-----|-------|
| Random | 1.0 | 500 | - |
| VAE | 5.2 | 48 | 30min |
| GAN | 9.1 | 9 | 2hr |
| Diffusion | 9.3 | 3 | 8hr |

**Stable Diffusion API: Production-Ready Generation**



| User Prompt | → | API Request | → | Diffusion Model | → | Generated Image |

**Key Parameters:**

cfg_scale: 1-20 (prompt adherence)

steps: 10-150 (quality vs speed)

seed: reproducibility

**Production APIs:**

DALL-E 3: $0.04-0.12/image

Midjourney: Subscription

Stable Diffusion: $0.004/image

Example: "A futuristic city at sunset"

→ High-quality 1024x1024 image in 10-30 seconds

**Usage:**

```
response = requests.post(
    api_url,
    headers={"Auth": key},
    json={
        "text_prompts": [{"text": "city"}],
        "cfg_scale": 7,
        "steps": 30
    })
APIs: DALL-E 3, Midjourney, Stable Diffusion
```

**Parameters:**

- `cfg_scale`: 1-20
- `steps`: 10-150

**Cost:** $0.004/image

**The Generative AI Landscape**



**VAEs**
- Probabilistic
- Smooth latent
- Blurry outputs
- Stable training

**GANs**
- Adversarial
- Sharp outputs
- Mode collapse risk
- Training unstable

**Modern systems combine approaches**

**Diffusion**
- Iterative
- Highest quality
- Slow sampling
- Stable

**Transformers**
- Sequential
- Excellent for text
- Scalable
- Left-to-right

**VAEs:** Probabilistic, smooth latent, blurry
**GANs:** Adversarial, sharp outputs, unstable

**Diffusion:** Iterative denoising, high quality, slow
**Transformers:** Sequential, excellent text, scalable

**Decision Criteria:**

**1. What are you generating?**

- Images: Diffusion or GAN
- Text: Transformer (GPT family)
- Structured data: VAE
- Multimodal: Diffusion + Transformer

**2. Data size?**

- < 10k samples: VAE (stable)
- 10k-100k: GAN or VAE
- > 100k: Diffusion or Transformer

**3. Priority?**

- Quality: Diffusion (FID ¡ 5)
- Speed: GAN (single pass)
- Stability: VAE (always converges)
- Control: Diffusion (guidance)

**Recommendation Table:**

| Use Case | Best | Why |
|---|---|---|
| Photorealistic | Diffusion | Quality |
| Fast prototype | GAN | Speed |
| Data augment | VAE | Stable |
| Text gen | Transformer | Sequential |
| Style transfer | VAE | Interpolate |
| Research | VAE | Interpret |

**When NOT to Use:**

- VAE: Need sharp images
- GAN: Limited data, need stability
- Diffusion: Real-time inference required
- All: Insufficient compute resources

Model selection requires balancing quality, speed, stability against problem constraints

## Common Pitfalls: What Can Go Wrong
Failure Modes and Solutions

**VAE Pitfalls**
1. **Posterior Collapse**
   - KL -¿ 0
   - Fix: $\beta$-VAE, warm-up
2. **Blurry**
   - MSE averages
   - Fix: Perceptual loss

**GAN Pitfalls**
1. **Mode Collapse**
   - Limited variety
   - Fix: Minibatch disc
2. **Unstable**
   - Oscillates
   - Fix: Wasserstein, spectral norm

**Diffusion Pitfalls**
1. **Slow (1000 steps)**
   - Latency issue
   - Fix: DDIM (50 steps)
2. **Memory**
   - High-res costly
   - Fix: Latent diffusion

Each approach has characteristic failure modes with specific solutions

## Generative AI Best Practices
From Research to Production

**Training:**

1. **Start Simple**
   - Low res first (64×64 before 1024×1024)
   - Validate on toy datasets
2. **Monitor Obsessively**
   - Log every 100 steps
   - Visual sample inspection
   - Track FID/IS
3. **Use Pretrained**
   - Transfer learning saves weeks
   - Fine-tune Stable Diffusion
4. **Ablation Studies**
   - Test components independently

**Deployment:**

1. **Quality Control**
   - Human-in-the-loop review
   - Content filtering
   - Watermarking
2. **Performance**
   - Quantization (FP16, INT8)
   - Distillation for speed
   - Caching
3. **Safety**
   - Rate limiting
   - Content moderation
   - Prompt injection defenses
4. **Continuous Improvement**
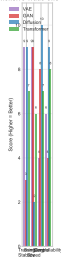   - User feedback
   - A/B testing

Production requires systematic validation and continuous monitoring

Comprehensive Trade-offs Comparison

**Stability:**

- VAEs, Diffusion: Stable
- GANs: Unstable

**Speed:**

- VAEs, GANs: Fast
- Diffusion: Slow

Choose based on requirements

**Quality:**

- Diffusion, GANs: Excellent
- VAEs: Blurry

**Control:**

- Diffusion, Transformers: High
- GANs: Limited

# State-of-the-Art Applications
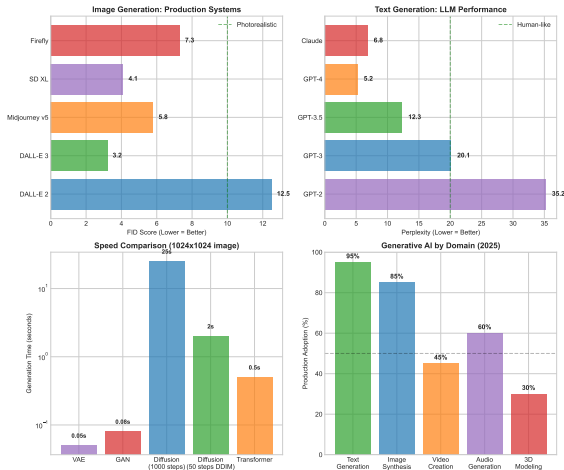
Production Generative AI Systems



Image Generation: Production Systems — FID Score (Lower = Better)
- Firefly: 7.3
- SD XL: 4.1
- Midjourney v5: 5.8
- DALL-E 3: 3.2
- DALL-E 2: 12.5
- Photorealistic (reference line)

Text Generation: LLM Performance — Perplexity (Lower = Better)
- Claude: 6.8
- GPT-4: 5.2
- GPT-3.5: 12.3
- GPT-3: 20.1
- GPT-2: 35.2
- Human-like (reference line)

Speed Comparison (1024x1024 image) — Generation Time (seconds)
- VAE: 0.05s
- GAN: 0.08s
- Diffusion (1000 steps): 20s
- Diffusion (50 steps DDIM): 2s
- Transformer: 0.5s

Generative AI by Domain (2025) — Production Adoption (%)
- Text Generation: 95%
- Image Synthesis: 85%
- Video Creation: 45%
- Audio Generation: 60%
- 3D Modeling: 30%

**Image:**

- DALL-E 3, Midjourney
- Stable Diffusion, Firefly

**Text:**

- GPT-4, Claude, Gemini
- Llama 2 (open)

Generative AI: Ethics and Future

**Learned:**
- VAEs: Probabilistic, blurry
- GANs: Adversarial, realistic
- Diffusion: Best quality
- Decision framework, pitfalls

**Future:**

**Ethics:**
- Deepfakes, copyright
- Bias, displacement

**Solutions:**
- Watermarking, auditing
- Governance