

# Machine Learning for Visual Computing Oral Exam重点

---

L.C.

2024.11

## 1. Core Machine Learning Foundations

---

### Question 1: "Can you explain the key differences between supervised and unsupervised learning?" (监督学习和无监督学习的区别)

Model Answer:

"Yes, let me explain the fundamental differences between supervised and unsupervised learning:

Supervised learning (监督学习) involves training a model using labeled data, where we have both input features and target outputs. The model learns to map inputs to outputs by minimizing the error between its predictions and the true labels. Common examples include image classification and regression tasks.

In contrast, unsupervised learning (无监督学习) works with unlabeled data. The model tries to discover hidden patterns or structures within the data without explicit guidance. Examples include clustering algorithms and dimensionality reduction techniques.

Let me give you a concrete example: In image classification (supervised), we might train a model with images labeled as 'cat' or 'dog'. However, in clustering (unsupervised), we might give the model unlabeled images and let it group similar images together without specifying what these groups should be."

### Question 2: "Why do we split our data into training, validation, and test sets? How do you typically choose the proportions?" (数据集划分的原因和比例)

Model Answer:

"The split of data into these three sets serves crucial purposes in machine learning:

The training set (训练集) is used to train the model and update its parameters. This is typically around 70% of the data.

The validation set (验证集), usually 15%, is used during training to:

- Tune hyperparameters
- Monitor for overfitting
- Make decisions about model architecture
- Select the best model

The test set (测试集), also typically 15%, is kept completely separate and only used once to evaluate the final model's performance. This gives us an unbiased estimate of the model's generalization ability.

The standard split is often 70:15:15, though this can vary depending on the total amount of data available and specific requirements of the project."

### Question 3: "Could you explain what MSE loss is and why we use it in linear regression?" (均方误差损失函数)

Model Answer:

"Mean Squared Error (MSE) loss (均方误差损失) is a fundamental loss function in machine learning, particularly in regression tasks. Let me break this down:

The mathematical formula is:

$$\text{MSE} = (1/n) \sum (y_i - \hat{y}_i)^2$$

Where:

- $y_i$  is the true value
- $\hat{y}_i$  is the predicted value
- $n$  is the number of samples

We use MSE in linear regression for several reasons:

1. It penalizes larger errors more heavily due to the squaring
2. It's differentiable, making it suitable for gradient-based optimization
3. It provides a single, positive number measuring model performance
4. The quadratic nature creates a convex optimization problem, ensuring a global minimum

In practice, minimizing MSE leads to finding the line of best fit that minimizes the squared vertical distances between the data points and the regression line."

### Question 4: "What is the Normal Equation in linear regression, and when would you use it versus gradient descent?" (正规方程)

Model Answer:

"The Normal Equation provides an analytical solution for finding the optimal parameters in linear regression. Let me explain:

The equation is:  $\theta = (X^T X)^{-1} X^T y$

Where:

- $\theta$  is the parameter vector we're solving for
- $X$  is the design matrix
- $y$  is the target vector

Advantages of the Normal Equation:

- Provides an exact solution in one step
- No need to choose learning rate
- No iteration required

However, we often prefer gradient descent when:

- The dataset is very large ( $n > 10000$ )
- The matrix is non-invertible
- We need to update the model online

The computational complexity is  $O(n^3)$  for the Normal Equation versus  $O(kn^2)$  for gradient descent, where  $k$  is the number of iterations."

## Question 5: "What is linear regression and how does it relate to optimization?" (线性回归与优化)

Model Answer:

"Linear regression is a fundamental supervised learning algorithm that models the relationship between input features and a continuous output variable. Let me explain its components and optimization process:

The model:  $y = wx + b$

Where:

- $y$  is the predicted output
- $w$  is the weight vector
- $x$  is the input features
- $b$  is the bias term

The optimization process involves:

1. Defining the loss function (typically MSE)
2. Finding the parameters ( $w, b$ ) that minimize this loss
3. Using either:
  - Gradient descent: iterative optimization
  - Normal equation: direct analytical solution

The beauty of linear regression lies in its simplicity and interpretability, making it a perfect starting point for understanding optimization in machine learning. The same principles extend to more complex models."

These questions and answers cover the core concepts while maintaining a conversational tone suitable for an oral exam. Remember to:

- Start with a clear, concise definition
- Provide relevant mathematical formulas when appropriate
- Include practical examples
- Explain why concepts are important
- Be prepared for follow-up questions

## Question 6: "How does polynomial fitting differ from linear fitting, and why might we use it?" (多项式拟合)

Model Answer:

"Polynomial fitting extends linear fitting by including higher-order terms. Let me explain:

In linear fitting, we have:  $y = wx + b$

In polynomial fitting, we have:  $y = w_1x^2 + w_2xy + w_3y^2 + w_4x + w_5y + w_6$

The key advantages of polynomial fitting include:

1. Capturing non-linear relationships in data
2. More flexible model capacity
3. Better fit for complex patterns

However, there are important considerations:

- Higher risk of overfitting
- Need for regularization
- More computationally expensive
- Requires careful degree selection

Interestingly, we can still use linear regression techniques by treating each polynomial term as a separate feature, effectively transforming our feature space while keeping our optimization methods the same."

## Question 7: "Can you explain what a feature matrix is and how we use it in machine learning?" (特征矩阵)

Model Answer:

"The feature matrix is a fundamental concept in machine learning. Let me break this down:

Structure:

- Each row represents one sample
- Each column represents one feature
- Dimensions are typically  $n \times d$  where:
  - $n$  is the number of samples
  - $d$  is the number of features

For example, in a housing price prediction task:

```
x = [  
    [size, bedrooms, age],  
    [1500,    3,    10],  
    [2000,    4,    15],  
    ...  
]
```

We use this matrix for:

1. Input to our models

2. Feature engineering
3. Computing correlations
4. Normalization/standardization

Understanding the feature matrix is crucial because it's the primary way we represent our data for machine learning algorithms."

## Question 8: "What role does matrix multiplication play in machine learning, and can you give some examples?" (矩阵乘法在机器学习中的应用)

Model Answer:

"Matrix multiplication is fundamental to machine learning operations. Let me explain its key roles:

1. Linear Transformations:
  - Computing layer outputs:  $y = Wx + b$
  - Feature transformations
  - Coordinate system changes
2. Common Operations:
  - Computing predictions:  $\hat{y} = Xw$
  - Calculating gradients:  $\nabla w = X^T(Xw - y)$
  - Feature interactions:  $X^TX$  for correlation matrix
3. Practical Applications:
  - Neural network layers
  - Principal Component Analysis (PCA)
  - Computing similarity metrics

The efficiency of matrix operations is crucial for modern machine learning, which is why frameworks like PyTorch and TensorFlow are optimized for matrix computations."

## Question 9: "What is regularization in the context of linear regression, and why do we use it?" (正则化)

Model Answer:

"Regularization is a technique to prevent overfitting in machine learning models. Let me explain its implementation and importance:

Types of Regularization:

1. L1 (Lasso): Add  $|w|$  term
2. L2 (Ridge): Add  $|w|^2$  term

The regularized loss function becomes:

$$L(w) = \text{MSE} + \lambda |w|^2$$

Where:

- $\lambda$  is the regularization strength

- $\|w\|^2$  is the L2 norm of weights

Benefits:

1. Prevents overfitting
2. Improves generalization
3. Handles multicollinearity
4. Makes solutions more stable

When using the Normal Equation with regularization:

$$w = (X^T X + \lambda I)^{-1} X^T y$$

This modification helps ensure the matrix is invertible and produces more robust solutions."

## Question 10: "How do we evaluate the performance of a regression model?" (回归模型评估)

Model Answer:

"There are several key metrics and methods for evaluating regression models:

### 1. Error Metrics:

- Mean Squared Error (MSE): Average squared differences
- Root Mean Squared Error (RMSE): Square root of MSE
- Mean Absolute Error (MAE): Average absolute differences
- R-squared ( $R^2$ ): Proportion of variance explained

### 2. Validation Techniques:

- Cross-validation
- Hold-out validation
- Time series splitting for temporal data

### 3. Diagnostic Plots:

- Residual plots
- Q-Q plots
- Prediction vs. Actual plots

### 4. Practical Considerations:

- Domain-specific requirements
- Computational efficiency
- Interpretability needs

We typically use multiple metrics because each provides different insights into model performance. For example, RMSE is in the same units as the target variable, making it interpretable, while  $R^2$  provides a normalized measure of fit."

## 2. Gradient-Based Optimization

### Question 1: "Could you explain what gradient descent is and how it works?" (梯度下降法的基本原理)

Model Answer:

"Let me explain gradient descent comprehensively:

Gradient descent is an iterative optimization algorithm used to find the minimum of a function. Think of it like finding the lowest point in a valley by taking steps downhill.

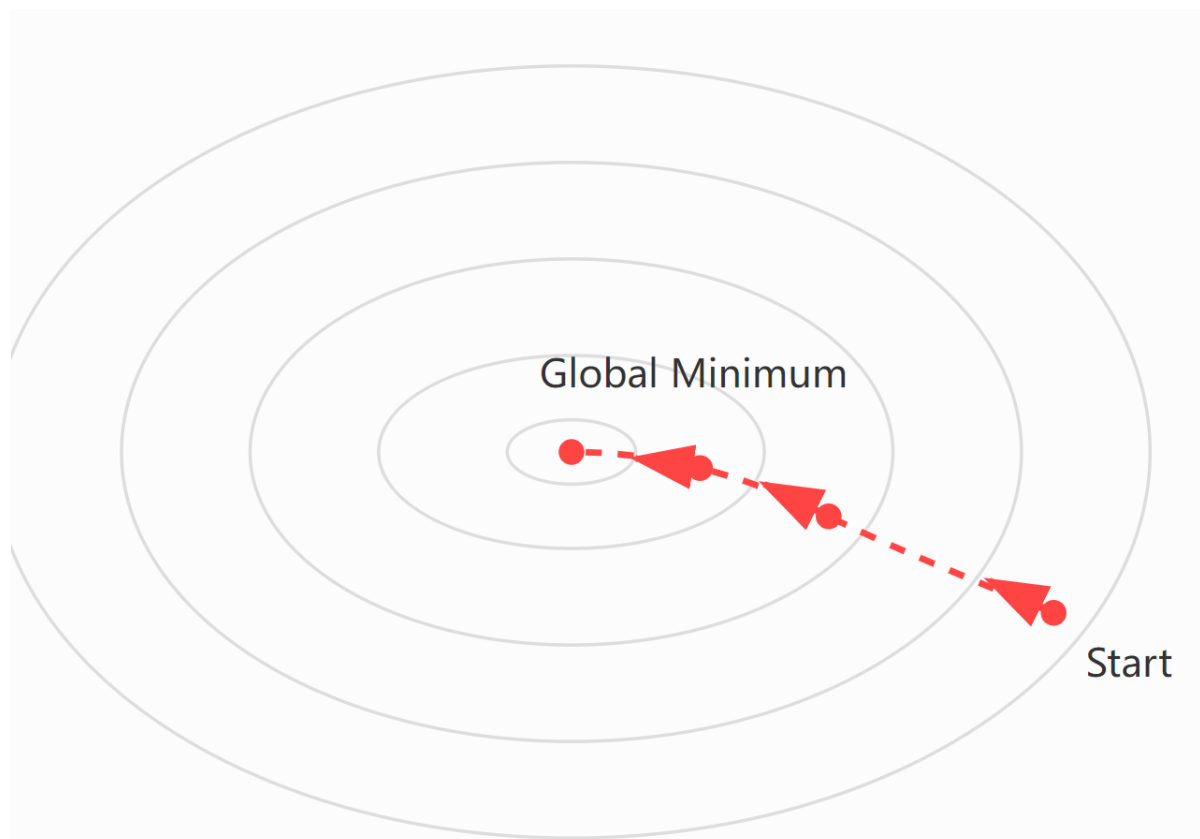
The basic update rule is:

$$\theta = \theta - \alpha \nabla f(\theta)$$

Where:

- $\theta$  is the parameter we're optimizing
- $\alpha$  is the learning rate (步长)
- $\nabla f(\theta)$  is the gradient of our objective function

Let me draw the process:



1. We start at some initial point
2. Calculate the gradient (slope) at that point
3. Take a step in the negative gradient direction
4. Repeat until convergence

Critical components:

1. Learning Rate Selection:

- Too large: might overshoot
- Too small: slow convergence

## 2. Convergence Criteria:

- Gradient magnitude near zero
- Change in parameters below threshold
- Maximum iterations reached

This forms the foundation for training most modern machine learning models."

## Question 2: "What's the difference between batch gradient descent and stochastic gradient descent?" (批量梯度下降vs随机梯度下降)

Model Answer:

"Let me explain the key differences between these approaches:

### 1. Batch Gradient Descent (批量梯度下降):

$$\text{Loss} = (1/N) \sum L_i(\theta) \quad \# \text{ Full dataset}$$

$$\theta = \theta - \alpha \nabla \text{Loss}$$

- Uses entire dataset for each update
- More stable updates
- Computationally expensive
- Better gradient estimate

### 2. Stochastic Gradient Descent (随机梯度下降):

$$\text{Loss} = L_i(\theta) \quad \# \text{ Single sample}$$

$$\theta = \theta - \alpha \nabla \text{Loss}$$

- Uses single sample per update
- Faster iterations
- Noisier updates
- More frequent parameter updates

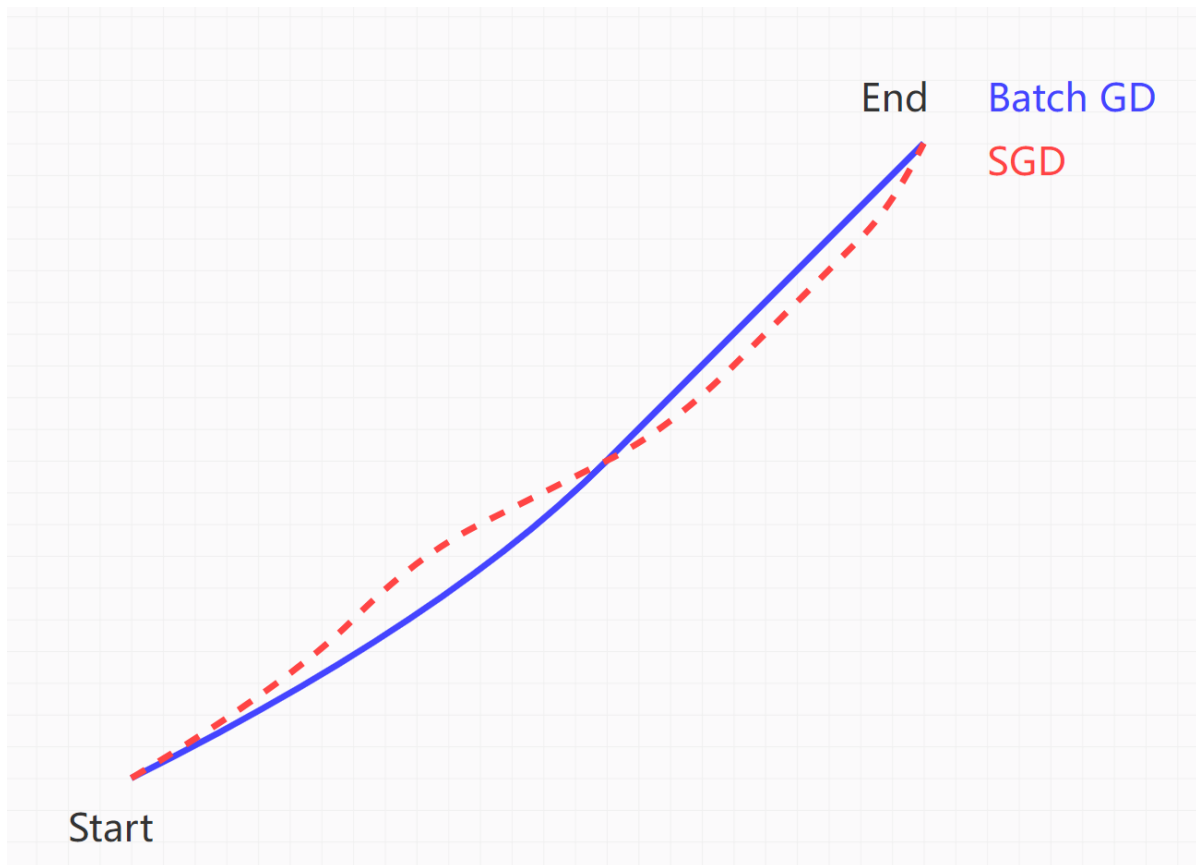
### 3. Mini-batch SGD (小批量梯度下降):

$$\text{Loss} = (1/m) \sum L_i(\theta) \quad \# m \text{ samples}$$

$$\theta = \theta - \alpha \nabla \text{Loss}$$

- Compromise between batch and SGD
- Typical batch sizes: 32, 64, 128
- Most commonly used in practice



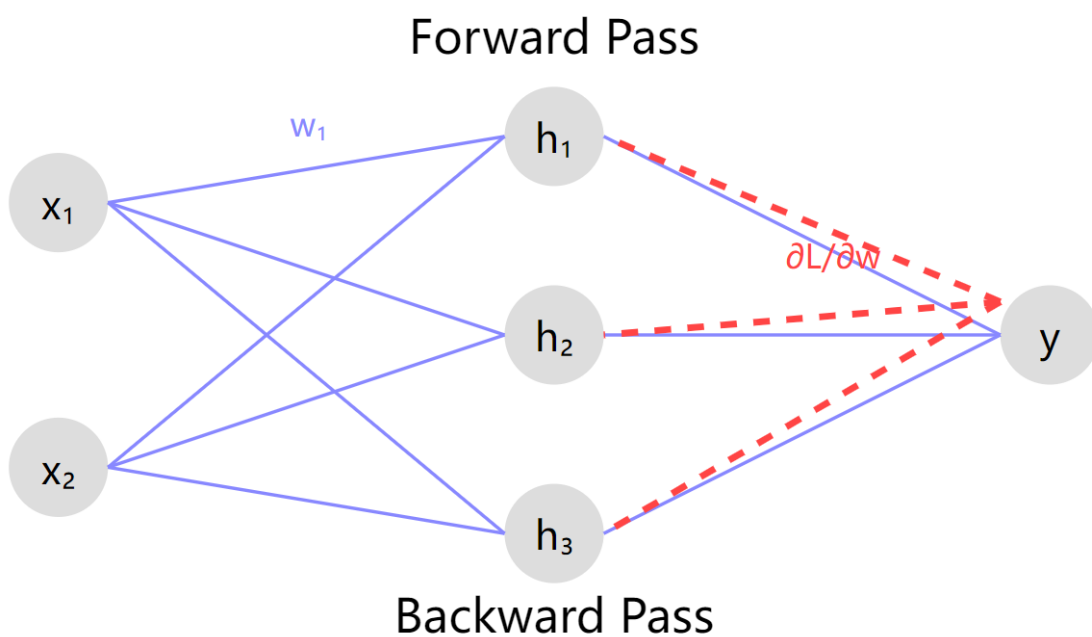


- Batch GD: smooth path to minimum
- SGD: noisy path but faster progress
- Mini-batch: balanced approach"

### Question 3: "Could you explain the chain rule and its role in backpropagation?" (链式法则与反向传播)

Model Answer:

"The chain rule is fundamental to backpropagation. Let me explain:



1. Chain Rule (链式法则):  
For composite functions:

$$y = f(g(x))$$
$$dy/dx = dy/dg \times dg/dx$$

2. In Neural Networks:

[Drawing a simple network]

$$x \rightarrow h1 \rightarrow h2 \rightarrow y$$

Computing  $\partial L / \partial w$  for each weight requires:

$$\partial L / \partial w = \partial L / \partial y \times \partial y / \partial h2 \times \partial h2 / \partial h1 \times \partial h1 / \partial w$$

3. Backpropagation Process:

a) Forward Pass:

- Compute activations layer by layer
- Store intermediate values

b) Backward Pass:

- Start from output layer
- Apply chain rule repeatedly
- Update weights using gradients

This efficient algorithm makes training deep neural networks possible."

## Question 4: "What are the common challenges in gradient descent optimization and how do we address them?" (梯度下降优化中的常见问题)

Model Answer:

"There are several key challenges in gradient descent optimization:

1. Learning Rate Issues:

- Too large: overshooting/divergence
  - Too small: slow convergence
- Solutions:
- Learning rate scheduling
  - Adaptive learning rates (Adam, RMSprop)

2. Local Minima and Saddle Points:

- Local minima: suboptimal solutions
  - Saddle points: zero gradient but not minimum
- Solutions:
- Momentum
  - Random initialization
  - Adding noise

3. Ill-Conditioning:

- Different parameters have different scales  
Solutions:
- Feature normalization
- Batch normalization
- Adaptive optimization methods

4. Vanishing/Exploding Gradients:  
Solutions:

- Proper initialization
- Gradient clipping
- Skip connections"

## Question 5: "Explain momentum in gradient descent. How and why does it help?" (动量梯度下降)

Model Answer:

"Momentum is a technique that accelerates gradient descent by adding a fraction of the previous update to the current one.

The formula is:

$$\begin{aligned}v &= \beta v - \alpha \nabla f(\theta) \\ \theta &= \theta + v\end{aligned}$$

Where:

- $\beta$  is the momentum coefficient (typically 0.9)
- $v$  is the velocity vector

Benefits:

1. Faster Convergence:
  - Accelerates in consistent directions
  - Reduces oscillations in narrow valleys
2. Avoiding Local Minima:
  - Maintains momentum through small bumps
  - Can escape shallow local minima
3. Practical Implementation:
  - Initialize  $v = 0$
  - Update  $v$  first, then  $\theta$
  - Typically improves training speed"

Let me continue with more advanced questions about optimization.

## Question 6: "Can you explain why backpropagation is more efficient than naive gradient computation in neural networks?" (反向传播的效率优势)

Model Answer:

"Let me explain this with a concrete example:

Consider a simple neural network:

[Drawing a 3-layer network]

Input → Hidden Layer → Output  
(x) → (h) → (y)

1. Naive Approach (直接计算):

- Would require computing gradients separately for each weight
- Computational complexity:  $O(n^2)$
- Redundant calculations

2. Backpropagation Approach (反向传播):

Let's walk through the computation:

```
# Forward Pass
h = σ(wx + b)      # Hidden layer activation
y = vh + c          # Output
L = (y - t)²        # Loss

# Backward Pass
∂L/∂y = 2(y - t)   # Output gradient
∂L/∂v = ∂L/∂y × h  # weight gradient
∂L/∂h = v^T × ∂L/∂y # Hidden layer gradient
```

Key Advantages:

1. Reuses intermediate computations
2. Computational complexity:  $O(n)$
3. Naturally parallelizable

The efficiency comes from storing intermediate results and using the chain rule systematically."

## Question 7: "What is the relationship between learning rate, batch size, and convergence in SGD?" (学习率、批量大小与收敛性的关系)

Model Answer:

"This is a crucial relationship in optimization. Let me break it down:

1. Learning Rate Effects:

[Drawing learning curves with different rates]

Too Large: zig-zag pattern, possible divergence  
Too Small: slow convergence  
Just Right: smooth convergence

## 2. Batch Size Impact:

### Large Batch:

- More stable gradients
- Better parallelization
- Might get stuck in sharp minima

### Small Batch:

- Noisy gradients
- Natural regularization
- Better generalization

## 3. The Interaction:

- Learning Rate  $\propto \sqrt{\text{Batch Size}}$
- When increasing batch size, increase learning rate
- Linear scaling rule:  $\text{lr} = \text{base\_lr} \times (\text{batch\_size}/\text{base\_batch})$

## 4. Practical Guidelines:

```
# Example learning rate schedule
base_lr = 0.1
batch_size = 128
lr = base_lr * sqrt(batch_size/32)

# Common batch sizes: 32, 64, 128, 256
# Common learning rates: 0.1, 0.01, 0.001
```

This relationship is crucial for efficient training."

## Question 8: "How do adaptive optimization methods differ from standard SGD?" (自适应优化方法)

Model Answer:

"Let me explain the key adaptive optimization methods:

### 1. Adam (Adaptive Moment Estimation):

```
# Adam update rule
m =  $\beta_1 m + (1 - \beta_1) \nabla \theta$  # First moment
v =  $\beta_2 v + (1 - \beta_2) \nabla \theta^2$  # Second moment
 $\theta = \theta - \alpha * m / \sqrt{v}$  # Parameter update
```

### 2. RMSprop:

[Drawing showing parameter-wise learning rates]

```
# RMSprop update rule
```

```
v =  $\beta v + (1-\beta)\nabla\theta^2$ 
```

```
 $\theta = \theta - \alpha * \nabla\theta / \sqrt{v}$ 
```

3. Comparison with SGD:

SGD:

- Single learning rate
- Manual schedule needed
- More interpretable

Adaptive Methods:

- Parameter-specific rates
- Automatic scaling
- Better with sparse gradients

4. When to Use Each:

- SGD: Better generalization
- Adam: Faster convergence
- RMSprop: Good for RNNs"

## Question 9: "What role does gradient normalization play in deep learning optimization?" (梯度归一化)

Model Answer:

"Gradient normalization is crucial for stable training:

1. Batch Normalization:

```
# During training
```

```
 $\mu = \text{mean}(x)$ 
```

```
 $\sigma = \text{std}(x)$ 
```

```
 $x_{\text{norm}} = (x - \mu) / \sigma$ 
```

```
 $y = \gamma x_{\text{norm}} + \beta$ 
```

2. Layer Normalization:

- Similar to batch norm but across features
- Important for transformers and RNNs

3. Gradient Clipping:

```
# Gradient clipping
```

```
if  $\|\nabla\theta\| > \text{threshold}$ :
```

```
 $\nabla\theta = \text{threshold} * \nabla\theta / \|\nabla\theta\|$ 
```

Benefits:

1. Controls gradient magnitude
2. Prevents exploding gradients
3. Stabilizes training

4. Enables larger learning rates"

Let me add these two crucial questions.

## Question 10: "Can you write out the pseudocode for gradient descent optimization?" (梯度下降伪代码)

Model Answer:

"Let me write the pseudocode for different variants of gradient descent:

1. Batch Gradient Descent:

```
def batch_gradient_descent(X, y, learning_rate=0.01, max_epochs=1000):
    # Initialize parameters
    w = initialize_weights()
    b = 0

    for epoch in range(max_epochs):
        # Compute predictions
        y_pred = X @ w + b

        # Compute gradients
        dw = (1/m) * X.T @ (y_pred - y) # m is sample size
        db = (1/m) * sum(y_pred - y)

        # Update parameters
        w = w - learning_rate * dw
        b = b - learning_rate * db

        # Check convergence
        if compute_loss() < threshold:
            break

    return w, b
```

2. Stochastic Gradient Descent:

```
def sgd(X, y, learning_rate=0.01, max_epochs=100):
    w = initialize_weights()
    b = 0

    for epoch in range(max_epochs):
        # Shuffle data
        indices = shuffle(range(len(X)))

        for i in indices:
            # Single sample update
            y_pred = X[i] @ w + b

            # Compute gradients
            dw = X[i].T * (y_pred - y[i])
            db = y_pred - y[i]

            # Update parameters
            w = w - learning_rate * dw
```

```
        b = b - learning_rate * db

    return w, b
```

### 3. Mini-batch Gradient Descent:

```
def minibatch_gd(X, y, batch_size=32, learning_rate=0.01):
    w = initialize_weights()
    b = 0

    for epoch in range(max_epochs):
        # Generate mini-batches
        batches = create_mini_batches(X, y, batch_size)

        for batch_X, batch_y in batches:
            # Forward pass
            y_pred = batch_X @ w + b

            # Compute gradients
            dw = (1/batch_size) * batch_X.T @ (y_pred - batch_y)
            db = (1/batch_size) * sum(y_pred - batch_y)

            # Update parameters
            w = w - learning_rate * dw
            b = b - learning_rate * db

    return w, b
```

Key Components to Explain:

1. Parameter initialization
2. Forward pass computation
3. Gradient calculation
4. Parameter updates
5. Convergence checking"

## Question 12: "Can you explain the backpropagation process in detail with a concrete example?" (反向传播详解)

Model Answer:

"Let me explain backpropagation using a simple neural network example:

[Drawing a 2-layer neural network]

```
Input (x) → Hidden (h) → Output (y)
2 neurons → 2 neurons → 1 neuron
```

Let's solve this step by step:

1. Network Setup:



```
# Forward equations
h1 =  $\sigma(w_{11}x_1 + w_{12}x_2 + b_1)$ 
h2 =  $\sigma(w_{21}x_1 + w_{22}x_2 + b_2)$ 
y =  $\sigma(v_1h_1 + v_2h_2 + b_3)$ 
Loss =  $(y - t)^2/2$  # t is target
```

2. Forward Pass Example:  
Given input  $x = [0.1, 0.2]$ :

```
Step 1: Hidden Layer
h1 =  $\sigma(0.15 \times 0.1 + 0.2 \times 0.2 + 0.1) = 0.57$ 
h2 =  $\sigma(0.25 \times 0.1 + 0.3 \times 0.2 + 0.1) = 0.61$ 

Step 2: Output Layer
y =  $\sigma(0.4 \times 0.57 + 0.45 \times 0.61 + 0.1) = 0.73$ 
```

3. Backward Pass:

```
# Chain rule application
 $\partial L / \partial y = y - t$ 
 $\partial L / \partial v_1 = \partial L / \partial y \times \partial y / \partial v_1 = (y - t) \times y(1-y) \times h_1$ 
 $\partial L / \partial v_2 = \partial L / \partial y \times \partial y / \partial v_2 = (y - t) \times y(1-y) \times h_2$ 

# Hidden layer gradients
 $\partial L / \partial h_1 = \partial L / \partial y \times \partial y / \partial h_1 = (y - t) \times y(1-y) \times v_1$ 
 $\partial L / \partial w_{11} = \partial L / \partial h_1 \times \partial h_1 / \partial w_{11} = \partial L / \partial h_1 \times h_1(1-h_1) \times x_1$ 
```

4. Update Process:

```
# Parameter updates
v1_new =  $v_1 - \text{learning\_rate} \times \partial L / \partial v_1$ 
w11_new =  $w_{11} - \text{learning\_rate} \times \partial L / \partial w_{11}$ 
# Similar updates for other weights
```

Let me provide a comprehensive example of backpropagation with answers to follow-up questions.

Detailed BP Example:

[Let's use a simple 2-2-1 neural network]

Input Layer (2 neurons) → Hidden Layer (2 neurons) → Output Layer (1 neuron)

Given:

```
Input: x = [0.1, 0.2]
Target: t = 0.8
Initial weights:
- w1 = [[0.15, 0.20], # Hidden layer weights
        [0.25, 0.30]]
- w2 = [0.40, 0.45] # Output layer weights
Biases: b1 = [0.1, 0.1], b2 = 0.1
Activation function: sigmoid(x) = 1/(1 + e^(-x))
Learning rate: α = 0.1
```

Step-by-Step Calculation:

1. Forward Pass:

```
# Hidden Layer
h1_in = 0.15×0.1 + 0.20×0.2 + 0.1 = 0.155
h1 = sigmoid(0.155) = 0.539

h2_in = 0.25×0.1 + 0.30×0.2 + 0.1 = 0.185
h2 = sigmoid(0.185) = 0.546

# Output Layer
y_in = 0.40×0.539 + 0.45×0.546 + 0.1 = 0.463
y = sigmoid(0.463) = 0.614
```

2. Backward Pass:

```
# Output Layer Error
E = (0.614 - 0.8)²/2 = 0.0344
δy = (y - t)×y×(1-y) = -0.0456

# Hidden Layer Error
δh1 = δy×w2[0]×h1×(1-h1) = -0.0042
δh2 = δy×w2[1]×h2×(1-h2) = -0.0045

# Weight Updates
Δw2[0] = -α×δy×h1 = 0.00246
Δw2[1] = -α×δy×h2 = 0.00249
```

Follow-up Questions That Might Be Asked:

Q1: "What happens if we change the activation function?"

A: "Changing the activation function affects both the forward and backward passes:

- ReLU: Simpler gradient (1 or 0), helps with vanishing gradients
- tanh: Similar to sigmoid but centered at 0, often better for deep networks
- Leaky ReLU: Prevents dead neurons by allowing small negative gradients

The choice impacts:

1. Training speed
2. Gradient flow
3. Network capacity"

Q2: "How does the network handle different input scales?"

A: "Input scaling is crucial for network performance:

1. Without scaling:

- Large inputs can saturate neurons
- Different features might dominate unfairly
- Slower convergence

2. Solutions:

- Standardization:  $(x - \mu)/\sigma$
- Min-Max scaling:  $(x - \min)/(\max - \min)$
- Batch Normalization in deeper networks"

Q3: "What modifications would you make for deeper networks?"

A: "For deeper networks, several modifications are essential:

1. Architecture changes:

- Add batch normalization
- Use skip connections (ResNet style)
- Implement layer normalization

2. Training strategies:

- Learning rate scheduling
- Gradient clipping
- Proper initialization (He/Xavier)"

Q4: "How do you handle the vanishing gradient problem?"

A: "The vanishing gradient problem can be addressed through:

1. Architectural solutions:

- Skip connections
- LSTM/GRU for sequential data
- Transformer architecture

2. Training techniques:

- Better activation functions (ReLU family)
- Proper initialization
- Gradient normalization"

Q5: "Why do we use non-linear activation functions?"

A: "Non-linear activation functions are crucial because:

1. They enable the network to learn complex patterns:

- Linear combinations alone can't approximate non-linear functions
- Each layer creates new feature combinations

2. Without non-linearities:

- Multiple layers would collapse into a single linear transformation
- Network loses its representation power
- Complex patterns couldn't be learned"

### 3. Neural Network Fundamentals

#### Question 1: "Can you explain what a perceptron is and how it evolved into the multi-layer perceptron?" (感知机到多层感知机的演变)

Model Answer:

"Let me explain this evolution step by step:

##### 1. Single Perceptron (单层感知机):

[Drawing of single perceptron]

Inputs → Weights → Sum → Step Function → Output

$x_1 \rightarrow w_1$

$x_2 \rightarrow w_2 \rightarrow \Sigma \rightarrow f \rightarrow y$

$x_3 \rightarrow w_3$

Mathematical representation:

$$y = f(\sum w_i x_i + b)$$

where  $f$  is step function:

$$f(x) = 1 \text{ if } x \geq 0$$

$$f(x) = 0 \text{ if } x < 0$$

##### 2. Limitations of Single Perceptron:

- Can only learn linearly separable patterns
- Cannot solve XOR problem
- Binary output only

##### 3. Multi-layer Perceptron (多层感知机):

Input Layer → Hidden Layer(s) → Output Layer

Key improvements:

- Multiple layers of neurons
- Non-linear activation functions
- Can approximate any continuous function
- Capable of learning complex patterns

##### 4. Modern Implementation:

```

class MLP:
    def __init__(self, layers=[3,4,1]):
        self.weights = initialize_weights(layers)
        self.biases = initialize_biases(layers)

    def forward(self, x):
        activation = x
        for w, b in zip(self.weights, self.biases):
            z = np.dot(w, activation) + b
            activation = sigmoid(z)
        return activation

```

This evolution marked the beginning of deep learning as we know it today."

## Question 2: "Could you explain the different activation functions used in neural networks and their pros and cons?" (激活函数)

Model Answer:

"Let me break down the main activation functions:

1. Sigmoid ( $\sigma$ ):

[Drawing sigmoid curve]

$$\sigma(x) = 1 / (1 + e^{(-x)})$$

Pros:

- Smooth gradient
- Output between [0,1]
- Clear probability interpretation

Cons:

- Vanishing gradient problem
- Not zero-centered
- Computationally expensive

2. ReLU (Rectified Linear Unit):

[Drawing ReLU function]

$$f(x) = \max(0, x)$$

Pros:

- Fast computation
- No vanishing gradient for  $x > 0$
- Sparse activation

Cons:

- "Dying ReLU" problem

- Not zero-centered
- Unbounded output

### 3. Leaky ReLU:

$$f(x) = \max(0.01x, x)$$

Pros:

- Prevents dying ReLU
- All benefits of ReLU
- Small negative gradient

### 4. tanh:

[Drawing tanh curve]

$$\tanh(x) = (e^x - e^{-x}) / (e^x + e^{-x})$$

Pros:

- Zero-centered
- Bounded output [-1,1]
- Better than sigmoid for deep networks

Cons:

- Still has vanishing gradient
- Computationally expensive

### 5. Usage Guidelines:

- Hidden layers: ReLU (default choice)
- Output layer:
  - Classification: Sigmoid/Softmax
  - Regression: Linear
  - Bounded output: tanh

## Question 3: "Explain the XOR problem and how MLPs solve it. Could you implement a solution?" (异或问题)

Model Answer:

"The XOR problem is a classic example showing the limitations of single-layer networks.

### 1. XOR Truth Table:

$x_1$	$x_2$		$y$
0	0		0
0	1		1
1	0		1
1	1		0

## 2. Why Single Layer Fails:

- XOR is not linearly separable
- Requires at least one hidden layer

## 3. MLP Solution:

```
# Minimal architecture
class XORNetwork:
    def __init__(self):
        # 2-2-1 architecture
        self.hidden_weights = np.array([[1.0, 1.0],
                                         [-1.0, -1.0]])
        self.output_weights = np.array([1.0, -1.0])
        self.hidden_bias = np.array([0.0, -1.0])
        self.output_bias = np.array([0.0])

    def sigmoid(self, x):
        return 1/(1 + np.exp(-x))

    def forward(self, x):
        # Hidden layer
        h = self.sigmoid(np.dot(x, self.hidden_weights.T) + self.hidden_bias)
        # Output layer
        y = self.sigmoid(np.dot(h, self.output_weights.T) + self.output_bias)
        return y
```

## 4. Geometric Interpretation:

- First hidden neuron: line separating (1,1)
- Second hidden neuron: line separating (0,0)
- Output combines these to form XOR function"

## Question 4: "Can you explain in detail how a Multi-Layer Perceptron works and its key components?" (多层感知机详解)

Model Answer:

"Let me break down the Multi-Layer Perceptron architecture and functionality:

### 1. Basic Architecture:

[Drawing MLP architecture]

Input Layer → Hidden Layer(s) → Output Layer  
(x) → (h) → (y)

### 2. Components and Mathematics:

#### a) Layer Types:

Input Layer:  $x = [x_1, x_2, \dots, x_n]$   
Hidden Layer:  $h = f(wx + b)$   
Output Layer:  $y = g(vh + c)$

Where:

- $W, V$  are weight matrices
- $b, c$  are bias vectors
- $f, g$  are activation functions

b) Detailed Computation Example:

```
# For a 3-4-2 network (3 inputs, 4 hidden, 2 outputs)
x = [x1, x2, x3]           # Input vector
W = [[w11, w12, w13],      # Hidden layer weights
      [w21, w22, w23],
      [w31, w32, w33],
      [w41, w42, w43]]
h = f(Wx + b)                # Hidden layer activation
y = g(Vh + c)                # Output layer
```

3. Key Features:

a) Non-linearity:

- Each hidden layer applies non-linear activation
- Enables learning complex patterns
- Hierarchical feature learning

b) Universal Approximation:

- Can approximate any continuous function
- Requires sufficient hidden units
- Depth vs width trade-off

4. Training Process:

a) Forward Pass:

```
z1 = Wx + b
h = f(z1)
z2 = Vh + c
y = g(z2)
```

b) Backward Pass:

```
# Compute gradients
δy = ∂L/∂y
δh = vT δy ⊙ f'(z1)
```

c) Weight Updates:

```
V = V - α∂L/∂V
W = W - α∂L/∂W
```

Example: MNIST Digit Classification (手写数字分类)



Let's build an MLP to classify handwritten digits (0-9):

### 1. Network Architecture:

Input Layer: 784 neurons (28x28 pixels)  
Hidden Layer 1: 128 neurons  
Hidden Layer 2: 64 neurons  
Output Layer: 10 neurons (0-9 digits)

Visualization:

[28x28 image] → [784] → [128] → [64] → [10]  
                  x      h<sub>1</sub>     h<sub>2</sub>     y

### 3. Example Computation:

For input image of digit "7":

```
# Input preprocessing
x = image/255.0 # Normalize to [0,1]

# Layer computations
h1 = ReLU(w1x + b1) # (128,)
h2 = ReLU(w2h1 + b2) # (64,)
y = softmax(w3h2 + b3) # (10,)

# Output: [0.01, 0.02, 0.01, 0.03, 0.02, 0.01, 0.02, 0.83, 0.03, 0.02]
#          0    1    2    3    4    5    6    7    8    9
```

Possible Oral Exam Questions:

Q1: "Why do we need two hidden layers in this architecture?"

A: "Two hidden layers allow the network to:

1. First layer: Extract basic features (edges, curves)
2. Second layer: Combine features into more complex patterns
3. Improve the network's capacity to learn complex relationships
4. Balance between model complexity and computational efficiency"

Q2: "Explain the choice of activation functions in this network."

A: "We use:

1. ReLU in hidden layers because:
  - Prevents vanishing gradient
  - Sparse activation
  - Fast computation
2. Softmax in output layer because:
  - Outputs sum to 1
  - Provides probability distribution
  - Suitable for multi-class classification"

Q3: "How would you handle overfitting in this network?"

A: "Several strategies could be applied:

1. Dropout between layers:

```
h1 = dropout(reLu(w1x + b1), p=0.5)
```

2. L2 regularization on weights

3. Early stopping based on validation accuracy

4. Data augmentation (rotation, scaling of digits)

5. Batch normalization"

Q4: "What are the computational considerations for this architecture?"

A: "Key considerations include:

1. Parameter count:

- Layer 1:  $784 \times 128 = 100,352$

- Layer 2:  $128 \times 64 = 8,192$

- Layer 3:  $64 \times 10 = 640$

Total: ~109K parameters

2. Memory requirements

3. Forward pass complexity

4. Batch size selection for training"

Q5: "How would you modify this architecture for different input sizes?"

A: "The architecture can be adapted by:

1. Adjusting input layer size to match new dimensions

2. Scaling hidden layer sizes proportionally

3. Keeping output layer fixed for same number of classes

4. Potentially adding/removing hidden layers based on complexity"

## Question 5: "How do we handle overfitting in neural networks, and what are the main techniques used?" (过拟合问题)

Model Answer:

"Let me explain overfitting and its solutions:

1. Identifying Overfitting:

[Drawing showing training vs validation curves]

- Training error decreases
- Validation error increases
- Model memorizes training data

2. Prevention Techniques:

a) Regularization:

```
# L2 regularization
loss = base_loss +  $\lambda \Sigma w^2$ 

# Dropout
def dropout_layer(x, p=0.5):
    mask = np.random.binomial(1, p, size=x.shape)
    return x * mask / p
```

b) Data Augmentation:

- Random rotations
- Flips
- Scaling
- Noise addition

c) Early Stopping:

```
if validation_error > best_error + patience:
    stop_training()
```

3. Architectural Solutions:

- Reduce model capacity
- Add batch normalization
- Use appropriate activation functions

4. Monitoring Tools:

- Learning curves
- Cross-validation
- Hold-out validation set
- Model complexity metrics"

## Question 6: "What is positional encoding and why is it important in neural networks?" (位置编码)

Model Answer:

"Position encoding is crucial for neural networks to understand spatial or sequential information:

1. Basic Concept:

```
PE(pos, 2i) = sin(pos/10000(2i/d))
PE(pos, 2i+1) = cos(pos/10000(2i/d))
```

2. Properties:

- Unique for each position
- Bounded between [-1,1]
- Allows extrapolation
- Preserves relative distances

### 3. Applications:

#### a) Transformers:

- Adds position information to tokens
- Enables parallel processing
- Maintains sequence order

#### b) NeRF:

- Encodes 3D coordinates
- Helps capture high-frequency details
- Improves rendering quality

### 4. Benefits:

- Better gradient flow
- Captures multi-scale information
- Enables position-aware learning"

## 4. Convolutional Neural Networks

---

### Question 1: "Can you explain what convolution is in CNNs and how it differs from traditional matrix multiplication?" (卷积运算)

Model Answer:

"Let me explain convolution in CNNs both conceptually and mathematically:

A convolution operation involves sliding a kernel (also called a filter) across an input image. At each position, we compute a weighted sum of the pixel values:

#### 1. The Process:

- Take a small filter (e.g.,  $3 \times 3$  matrix)
- Slide it across the image
- At each position: multiply element-wise and sum
- Store result in output feature map

For example, if we have a  $3 \times 3$  kernel for edge detection:

```
-1  -1  -1
-1   8  -1
-1  -1  -1
```

When we apply this to an image region:

1. Multiply corresponding elements
2. Sum all products
3. Store result in output position

Key differences from matrix multiplication:

- Parameter sharing: same kernel weights used across image
- Local connectivity: each output only depends on nearby inputs
- Translation invariance: can detect features regardless of position

This makes convolution particularly efficient for image processing tasks."

## Question 2: "What are the main components of a CNN architecture? Please explain each one's role." (CNN架构组件)

Model Answer:

"A typical CNN architecture consists of several key components:

### 1. Convolutional Layers (卷积层):

- Primary feature extractors
- Learn hierarchical features:
  - Early layers: edges, textures
  - Middle layers: patterns, parts
  - Deep layers: high-level concepts

### 2. Activation Functions (激活函数):

- Usually ReLU after convolutions
- Introduces non-linearity
- Helps with gradient flow
- Formula:  $f(x) = \max(0, x)$

### 3. Pooling Layers (池化层):

- Reduce spatial dimensions
- Types:
  - Max pooling: takes maximum value
  - Average pooling: takes average value
- Provides some translation invariance

### 4. Fully Connected Layers (全连接层):

- Usually at the end of network
- Combines features for final prediction
- Dense connections between neurons

The typical flow is:

Input  $\rightarrow$   $[[\text{Conv} \rightarrow \text{ReLU} \rightarrow \text{Pool}] \times N] \rightarrow [\text{FC} \rightarrow \text{ReLU}] \times M \rightarrow$  Output

This architecture progressively transforms raw pixels into meaningful features for classification or other tasks."

### Question 3: "Could you walk us through the innovations in some classic CNN architectures like AlexNet, ResNet, and DenseNet?" (经典CNN架构)

Model Answer:

"I'll explain the key innovations in these landmark architectures:

#### 1. AlexNet (2012):

- First major CNN success on ImageNet
- Key innovations:
  - ReLU activation
  - Local response normalization
  - Dropout regularization
  - Multiple GPUs training
- 5 convolutional layers, 3 fully connected layers
- Showed deep learning's potential for computer vision

#### 2. ResNet (2015):

- Introduced residual connections(残差连接):
  - $x_{out} = F(x) + x$
  - Helps with gradient flow
  - Enables training of very deep networks
- Key concept: learn residual function
- Solved vanishing gradient problem
- Can be hundreds of layers deep

#### 3. DenseNet (2017):

- Dense connectivity pattern:
  - Each layer connected to all previous layers
  - Better feature reuse
  - Improved gradient flow
- Advantages:
  - Fewer parameters
  - Better feature propagation
  - Natural feature reuse

The evolution from AlexNet → ResNet → DenseNet shows:

1. Increasing depth
2. Better gradient flow mechanisms
3. More efficient feature utilization"

## Question 4: "What are pooling operations in CNNs, and why do we use them?" (池化操作)

Model Answer:

"Let me explain pooling operations comprehensively:

### 1. Definition:

Pooling reduces spatial dimensions by summarizing regions into single values.

### 2. Common Types:

#### a) Max Pooling(最大池化):

```
2  5  |  1  3
6  8  |  4  2
-----
7  1  |  3  4
3  4  |  2  5
```

→ with 2×2 max pooling becomes:

```
8  4
7  5
```

#### b) Average Pooling(平均池化):

- Takes average of values in window
- Less common but useful for certain tasks

### 3. Benefits:

- Reduces computation cost
- Provides translation invariance
- Controls overfitting
- Reduces spatial dimensions

### 4. Typical Usage:

- Often 2×2 windows with stride 2
- After convolution layers
- Before fully connected layers"

## Question 5: "How does backpropagation work in CNNs, particularly through convolution and pooling layers?" (CNN中的反向传播)

Model Answer:

"Let me explain the backpropagation process in CNNs:

### 1. Convolution Layer Backprop:

- Compute gradient of loss w.r.t output
- Rotate kernel by 180°
- Perform convolution between gradient and layer input

- Update kernel weights using chain rule
2. Pooling Layer Backprop:
- For Max Pooling:
- Gradient flows only through maximum value
  - Other positions receive zero gradient
  - Records which input was maximum during forward pass

Example:

```
Forward Pass (2x2 max pool):
1  4  →  4
2  3

Backward Pass:
0  ∂L/∂y → propagates to position of max value
0  0
```

3. Chain Rule Application:

$$\partial L / \partial w = \partial L / \partial y * \partial y / \partial w$$

Where:

- L is loss
- y is layer output
- w is layer weights

This process enables end-to-end training of the entire network."

## Question 6: "Can you explain how CNN filters evolve through different layers of the network and what patterns they typically learn to detect?" (CNN过滤器和特征学习)

Model Answer:

"Let me explain the hierarchical feature learning in CNNs:

1. Early Layers (Input-Adjacent):

- Detect basic visual elements:
  - Edges in different orientations
  - Color transitions
  - Simple textures
- Filters typically show:
  - Horizontal/vertical lines
  - Diagonal edges
  - Color blobs

2. Middle Layers:

- Combine simple features into:
  - Corners and contours



- Textures and patterns
- Simple shapes
- More complex patterns emerge:
  - Repetitive textures
  - Basic object parts

### 3. Deep Layers:

- Learn highly abstract features:
  - Object parts (eyes, wheels)
  - Complex textures
  - Class-specific patterns

For example, in facial recognition:

- Layer 1: Edges and colors
- Layer 2: Facial contours
- Layer 3: Eyes, nose shapes
- Layer 4: Complete facial features

This hierarchical learning is what makes CNNs so powerful for visual tasks."

## **Question 7: "What are the trade-offs between different CNN architectures, and how do you choose the right one for a specific task?" (CNN架构选择)**

Model Answer:

"Let me discuss the key considerations when choosing CNN architectures:

### 1. Model Complexity vs. Dataset Size:

- Larger models (like DenseNet):
  - Need more training data
  - Higher computational cost
  - Better feature extraction
- Smaller models (like MobileNet):
  - Faster inference
  - Less prone to overfitting
  - Suitable for mobile devices

### 2. Architecture Characteristics:

- ResNet:
  - Good for very deep networks
  - Stable training
  - Gradient flow advantages
- DenseNet:
  - Better feature reuse

- Memory intensive
- Fewer parameters

### 3. Application Requirements:

- Real-time applications:
  - Choose efficient architectures
  - Consider inference speed
  - Focus on model compression
- High accuracy needed:
  - Deeper architectures
  - Ensemble methods
  - More complex models

### 4. Resource Constraints:

- Memory limitations
- Computational power
- Inference time requirements"

## **Question 8: "How do fully connected layers work in CNNs, and why are they typically placed at the end of the network?" (全连接层的作用)**

Model Answer:

"Fully connected layers play a crucial role in CNNs. Let me explain:

#### 1. Structure and Function:

- Every neuron connected to all neurons in previous layer
- Transforms spatial data into class scores
- Mathematical operation:  $y = Wx + b$

#### 2. Role in CNN:

- Combines high-level features
- Performs final classification
- Reduces spatial information to class probabilities

#### 3. Architecture Considerations:

- Usually multiple FC layers
- Often include dropout for regularization
- Typically decrease in size:
  - Example:  $4096 \rightarrow 4096 \rightarrow 1000$  (ImageNet)

#### 4. Modern Trends:

- Global Average Pooling
- Fewer FC layers

- Reduced parameters

Traditional vs Modern approach:

- Traditional: Multiple FC layers
- Modern: Often single FC or GAP
- Trade-off: Parameters vs Performance"

## Question 9: "How do you handle different input sizes in CNNs, and what are the implications for the architecture?" (输入尺寸处理)

Model Answer:

"This is an important practical consideration in CNN design:

### 1. Fixed-Size Approaches:

- Resize all inputs to fixed dimensions
- Advantages:
  - Simpler architecture
  - Batch processing
  - Consistent memory usage
- Disadvantages:
  - Potential information loss
  - Aspect ratio distortion

### 2. Variable-Size Solutions:

- Fully Convolutional Networks
- Global Pooling layers
- Spatial Pyramid Pooling
- Adaptive pooling

### 3. Implementation Strategies:

```
# Fixed size approach
resize_transform = transforms.Resize((224, 224))

# Adaptive approach
adaptive_pool = nn.AdaptiveAvgPool2d((1, 1))
```

### 4. Architectural Implications:

- Convolutional layers: Handle any size
- Pooling layers: Need careful sizing
- FC layers: Need fixed input size"

## Question 10: "Can you explain what a kernel is in CNNs and walk through a complete example of how convolution works with a specific input?" (卷积核和计算示例)

Model Answer:

"Let me explain kernels and demonstrate with a concrete example:

### 1. Kernel Definition:

A kernel is a small matrix of weights that slides over the input to detect features. Let me show a practical example:

Consider a 6×6 input image with a 3×3 kernel:

Input Image:

```
1  2  1  0  2  1
3  1  0  2  1  0
1  2  1  1  0  2
0  1  2  1  1  0
2  0  1  0  2  1
1  1  0  2  1  0
```

Edge Detection Kernel:

```
-1  -1  -1
-1   8  -1
-1  -1  -1
```

### 2. Convolution Process:

Let's compute one position step-by-step:

Input Region:	Kernel:	Calculation:
1 2 1	-1 -1 -1	$1 \times (-1) + 2 \times (-1) + 1 \times (-1)$
3 1 0	-1 8 -1	$= 3 \times (-1) + 1 \times (8) + 0 \times (-1)$
1 2 1	-1 -1 -1	$1 \times (-1) + 2 \times (-1) + 1 \times (-1)$
		$= -1 - 2 - 1 - 3 + 8 + 0 - 1 - 2 - 1 = -3$

### 3. Complete Output Calculation:

- Stride: 1
- No padding
- Output size:  $(N-F+1) \times (N-F+1)$ 
  - Where  $N=6$  (input size),  $F=3$  (kernel size)
  - Therefore output is 4×4

### 4. Real CNN Example:

Let's consider a typical first layer:

- Input: 224×224×3 RGB image
- Kernel: 3×3×3 (matches input depth)
- Number of kernels: 64

- Stride: 1
- Padding: 1

Output calculation:

- Output size:  $224 \times 224 \times 64$
- Each kernel produces one feature map
- Total parameters in layer:  $(3 \times 3 \times 3) \times 64 + 64$  biases = 1,792

5. Common Kernel Types:

- Edge detection (shown above)
- Blur kernel:

```
1/9  1/9  1/9
1/9  1/9  1/9
1/9  1/9  1/9
```

- Sharpening kernel:

```
0  -1  0
-1  5  -1
0  -1  0
```

The key points to understand are:

1. Kernels are learned during training
2. Different kernels detect different features
3. Earlier layers typically learn simpler patterns
4. The convolution operation maintains spatial relationships

This fundamental operation is what allows CNNs to effectively process visual information by detecting increasingly complex features through the network layers."

## 5. Feature Understanding and Visualization

**Question 1: "Can you explain what CAM (Class Activation Mapping) and Grad-CAM are, and how they help us understand CNNs?" (类激活图和梯度类激活图)**

Model Answer:

"Let me explain these important visualization techniques:

CAM (类激活映射) is a technique that helps us understand which regions of an image are important for classification decisions. Here's how it works:

1. Basic CAM:
  - Uses global average pooling layer just before the final classification
  - Weights of the final layer represent importance of each feature map
  - Mathematically:  $CAM = \sum_i w_i F_i$  where:

- $w_i$  are class-specific weights
- $F_i$  are feature maps

## 2. Grad-CAM's improvements:

- More general, can be applied to any CNN architecture
- Uses gradient information flowing to the last convolutional layer
- Formula:  $\text{Grad-CAM} = \text{ReLU}(\sum_i \alpha_i A_i)$  where:
  - $\alpha_i$  are feature map importance weights
  - $A_i$  are activation maps

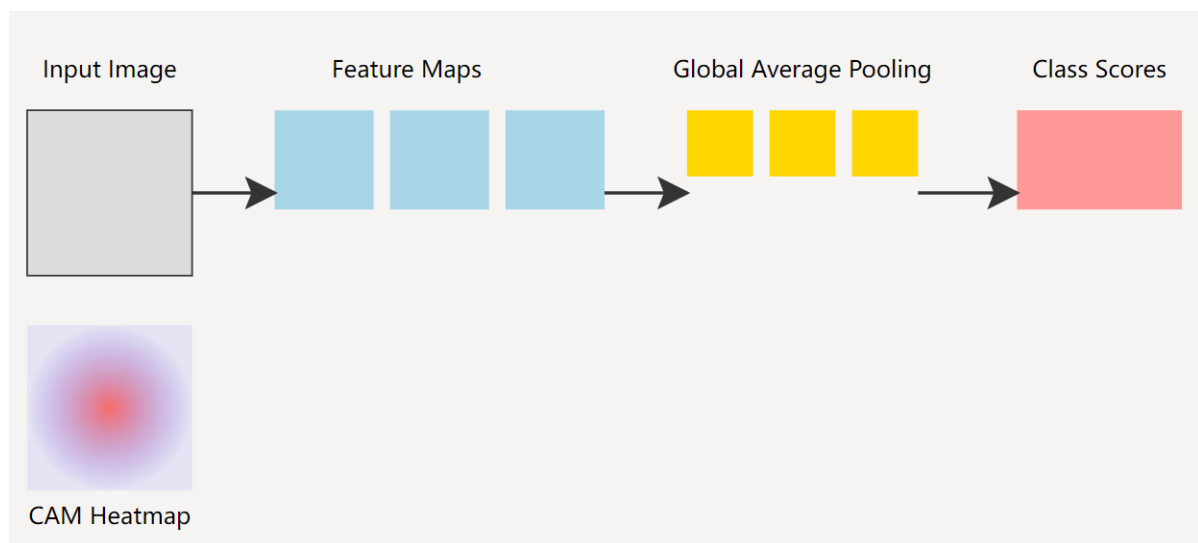
Practical example: If we're classifying a cat image:

- Red regions in the heatmap show areas most important for the decision
- Often highlights facial features for animal classification
- Can show if the model is focusing on relevant parts or background

Grad-CAM is particularly useful for:

- Model interpretation
- Debugging model decisions
- Identifying dataset bias
- Improving model trustworthiness"

Class Activation Mapping (CAM) Detailed Explanation:



Process:

### 1. Forward Pass:

- Input image goes through CNN
- Extract feature maps from final convolutional layer
- Apply global average pooling

### 2. Weight Combination:

- Multiply each feature map by its corresponding class weight
- Sum up weighted feature maps

### 3. Visualization:

- Upscale to input image size
- Apply colormap (usually red=high activation, blue=low)
- Overlay on original image

## Question 2: "What is Feature Inversion and how does it help us understand CNN representations?" (特征反演)

Model Answer:

"Feature Inversion is a powerful technique to visualize what information is preserved at different layers of a CNN. Let me explain:

The core concept:

- Given: Feature activations  $\Phi(x)$  at a particular layer
- Goal: Reconstruct an image  $x^*$  that produces similar features
- Optimization problem:  $x^* = \operatorname{argmin}_x ||\Phi(x) - \Phi(x_0)||^2 + \lambda R(x)$

Key components:

1. Content loss:  $||\Phi(x) - \Phi(x_0)||^2$ 
  - Ensures feature similarity
  - Usually uses L2 distance
2. Regularization  $R(x)$ :
  - Total variation regularization
  - Ensures visual plausibility
  - Controls noise and artifacts

The process reveals:

- Early layers preserve color and texture
- Middle layers maintain complex textures and patterns
- Deep layers keep high-level semantic information but lose spatial details

Practical applications:

1. Understanding layer representations
2. Analyzing network behavior
3. Neural style transfer foundations
4. Feature visualization research"

## Question 3: "Can you explain Gradient Ascent in the context of feature visualization?" (梯度上升可视化)

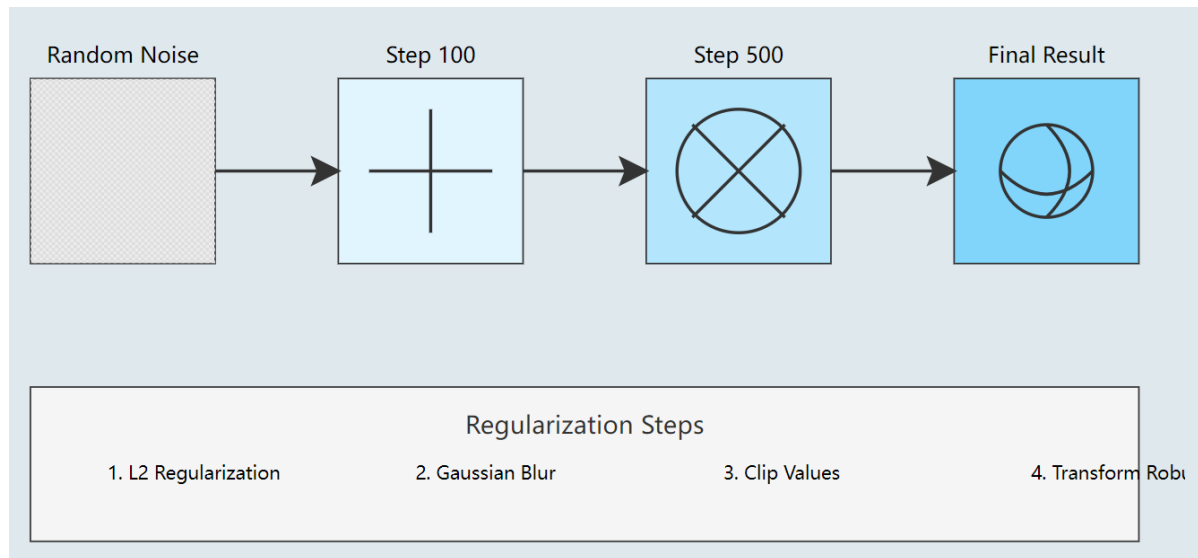
Model Answer:

"Gradient Ascent is a technique used to generate images that maximally activate specific neurons or channels in a CNN. Here's how it works:

Mathematical formulation:

$$I^* = \operatorname{argmax}_I \sum_i f_i(I)^2$$

Process:



1. Initialization:

- Start with random noise image
- Define target neuron/feature

2. Iterative Optimization:

- Forward pass to compute activation
- Compute gradients with respect to input
- Update image to maximize activation
- Apply regularization

3. Regularization Steps:

- L2 regularization to control magnitude
- Gaussian blur to reduce high-frequency noise
- Clip pixel values to valid range
- Apply transformations for robustness

The process continues until convergence or for a fixed number of iterations, revealing what patterns or features the chosen neuron is looking for.

Important regularization techniques:

1. Frequency penalization
2. Transformation robustness
3. Gaussian blur
4. Clipping pixel values

Examples of visualizations:

- Lower layers: edges, textures, simple patterns
- Middle layers: textures, repeated patterns
- Higher layers: object parts, complex structures

CAM & GA

- CAM is more about understanding model decisions



- Gradient Ascent is about understanding feature detectors
- Both are crucial for model interpretability
- Different layers reveal different levels of abstraction

## Question 4: "What is the feature space in CNNs and how do we analyze it?" (特征空间分析)

Model Answer:

"The feature space in CNNs represents the high-dimensional space where our data is embedded after passing through network layers. Let me elaborate:

Feature space characteristics:

1. Dimensionality
  - Early layers: high dimension, spatial information
  - Later layers: more compact, semantic information
  - Final layers: task-specific representations
2. Analysis methods:
  - a) Nearest Neighbor Search
    - Find similar images in feature space
    - Reveals semantic relationships
  - b) Dimensionality Reduction
    - PCA for linear projection
    - t-SNE for non-linear visualization
    - UMAP for preserving global structure
3. Visualization techniques:
  - Feature space clustering
  - Activation maximization
  - Interpolation between samples

Practical applications:

1. Understanding model decisions
2. Finding similar images
3. Detecting outliers
4. Evaluating feature quality"

## Question 5: "How do we analyze and visualize weights in CNNs?" (权重分析和可视化)

Model Answer:

"Weight analysis is crucial for understanding what our CNN has learned. Let me explain the key approaches:

1. First Layer Analysis:
  - Direct visualization of filters

- Often shows edge detectors, color detectors
- Patterns similar to Gabor filters

## 2. Weight Statistics:

- Distribution analysis
- Magnitude studies
- Sparsity patterns

## 3. Advanced Visualization:

- Weight evolution during training
- Connection patterns
- Impact on feature maps

Interpretation example:

- Low-level features: edges, colors
- Mid-level features: textures, patterns
- High-level features: semantic concepts

## **Question 6: "Could you explain the concept of feature embeddings and their applications?" (特征嵌入)**

Model Answer:

"Feature embeddings represent high-dimensional data in a lower-dimensional space while preserving important relationships. In CNNs:

### 1. Creation Process:

- Extract features from later layers (e.g., fc7 in VGG)
- Often 4096 or 2048 dimensions
- Can be reduced further using PCA/t-SNE

### 2. Key Applications:

- Image retrieval
- Similarity search
- Transfer learning
- Zero-shot learning

### 3. Properties:

- Semantic similarity preserved
- Enables efficient search
- Supports downstream tasks"

## Question 7: "What's the difference between feature visualization in shallow versus deep layers of a CNN?" (浅层vs深层特征可视化)

Model Answer:

"The differences are quite systematic:

Shallow Layers:

- Simple patterns (edges, textures)
- More localized information
- Higher spatial resolution
- Direct visual interpretation

Deep Layers:

- Complex semantic concepts
- Global information
- More abstract representations
- Requires sophisticated visualization techniques"

## 6. Style Transfer and Texture

---

### Question 1: "Can you explain what neural style transfer is and how it works?" (神经风格迁移)

Model Answer:

"Neural style transfer is a deep learning technique that combines the content of one image with the artistic style of another. Let me explain the key components:

We need three images:

1. Content image (内容图片) - The base image we want to stylize
2. Style image (风格图片) - The artwork whose style we want to capture
3. Generated image (生成图片) - Our output that we optimize

The process works through optimization of two main losses:

1. Content Loss (内容损失):
  - Measures how well we preserve the content structure
  - Usually computed using high-level features from a pretrained CNN
  - Typically uses features from later layers as they capture semantic information
2. Style Loss (风格损失):
  - Captures the artistic style using Gram matrices
  - Computed across multiple layers of the network
  - Focuses on texture and color patterns rather than spatial arrangement

The final loss function is:

$$L_{\text{total}} = \alpha L_{\text{content}} + \beta L_{\text{style}}$$

Where  $\alpha$  and  $\beta$  are weights controlling the content-style trade-off. This loss is then optimized using gradient descent to update the generated image."

## Question 2: "What is the Gram matrix and why is it crucial for style transfer?" (格拉姆矩阵)

Model Answer:

"The Gram matrix is a key mathematical tool in style transfer that captures the statistical correlations between features. Let me break this down:

Mathematical Definition:

- Given feature maps  $F$  of size  $(C \times H \times W)$
- Gram matrix  $G$  is computed as:  $G_{ij} = \sum_k F_{ik} F_{jk}$
- Where  $i, j$  are feature channels and  $k$  iterates over spatial positions

Key Properties:

1. Size Independence:
  - Results in  $C \times C$  matrix regardless of image size
  - Allows comparison between different-sized images
2. Style Representation:
  - Captures texture and style patterns
  - Loses spatial arrangement information
  - Preserves feature correlations
3. Implementation:
  - Usually computed at multiple layers
  - Different layers capture different aspects of style
  - Earlier layers capture fine details, later layers capture larger patterns

The Gram matrix is effective because it:

- Discards spatial information (which is content)
- Retains style information through feature correlations
- Is differentiable, allowing gradient-based optimization"

## Question 3: "Can you explain the different loss functions used in style transfer and their purposes?" (损失函数)

Model Answer:

"Let me provide a comprehensive explanation of the loss functions used in style transfer:

1. Content Loss Functions:

a) L2 Loss (MSE):

$$L_{\text{content}} = ||F_1(G) - F_1(C)||^2$$

Where:

- $F_l$  represents features at layer  $l$
- $G$  is generated image
- $C$  is content image

Use cases:

- Basic content preservation
- Feature matching in deep layers
- Simple to implement and optimize

b) MAE (Mean Absolute Error):

$$L_{MAE} = |F_l(G) - F_l(C)|$$

Benefits:

- More robust to outliers
- Less emphasis on large differences
- Better preservation of sharp features

2. Style Loss Functions:

a) Gram Matrix Loss:

$$L_{style} = \sum_l w_l ||G_l(S) - G_l(G)||^2_F$$

Where:

- $G_l$  is the Gram matrix at layer  $l$
- $w_l$  are layer weights
- $||\cdot||_F$  is Frobenius norm

b) SSIM (Structural Similarity):

Components:

- Luminance:  $l(x,y) = (2\mu_x\mu_y + c_1)/(\mu_x^2 + \mu_y^2 + c_1)$
- Contrast:  $c(x,y) = (2\sigma_x\sigma_y + c_2)/(\sigma_x^2 + \sigma_y^2 + c_2)$
- Structure:  $s(x,y) = (\sigma_{xy} + c_3)/(\sigma_x\sigma_y + c_3)$

Final SSIM:

$$SSIM(x,y) = [l(x,y)]^\alpha [c(x,y)]^\beta [s(x,y)]^\gamma$$

c) LPIPS:

$$L_{LPIPS} = \sum_l w_l ||F_l(x) - F_l(y)||_2$$

Features:

- Uses learned perceptual features
- Better matches human perception

- More sophisticated than pixel-based metrics

3. Total Loss Combination:

$$L_{total} = \alpha L_{content} + \beta L_{style} + \gamma L_{reg}$$

Where:

- $\alpha, \beta$  control content-style trade-off
- $\gamma L_{reg}$  is optional regularization term"

Each loss serves a specific purpose:

- Content preservation: L2, MAE
- Style matching: Gram matrix loss
- Perceptual quality: SSIM, LPIPS"

## Question 4: "How does texture synthesis work, and how does it relate to style transfer?" (纹理合成)

Model Answer:

"Texture synthesis is closely related to style transfer, focusing on generating new textures that match a given example. Let me explain the process:

Traditional Approach:

1. Analyze input texture
2. Generate new texture patch by patch
3. Ensure local consistency
4. Maintain global statistics

Neural Approach:

1. Extract features using CNN
2. Compute Gram matrices at multiple layers
3. Generate texture by optimizing:

$$L = \sum_i w_i |G_i(x) - G_i(s)|^2$$

Where:

- $G_i$  are Gram matrices
- $x$  is generated texture
- $s$  is source texture

Key Differences from Style Transfer:

1. No content loss needed
2. Focus on pure texture replication
3. Usually simpler optimization problem

Applications:

- Material synthesis
- Texture completion

- Background generation
- Game asset creation"

## Question 5: "What are feedforward networks in the context of style transfer?" (前馈网络风格迁移)

Model Answer:

"Feedforward networks offer a faster alternative to optimization-based style transfer. Let me explain:

Architecture:

1. Image Transformation Network:
  - Takes content image as input
  - Outputs stylized image directly
  - Usually uses encoder-decoder structure
2. Loss Network:
  - Pre-trained VGG network
  - Frozen weights
  - Computes content and style losses

Advantages:

1. Real-time style transfer
  - Single forward pass
  - No iteration needed
  - ~1000x faster than optimization
2. Consistent results
  - Deterministic output
  - More stable than optimization

Limitations:

1. One network per style
2. Less flexible than optimization
3. Potential quality trade-off

This approach is widely used in applications requiring real-time performance, such as mobile apps and video processing."

## Question 6: "How do we handle the multi-scale nature of textures in style transfer?" (多尺度纹理处理)

Model Answer:

"Multi-scale texture handling is crucial for effective style transfer. Let me explain the approach:

1. Feature Hierarchy:
  - Lower layers: fine details, brushstrokes

- Middle layers: textures, patterns
- Higher layers: larger structures

Implementation Strategy:

1. Multi-scale Feature Extraction:

- Use different VGG layers (conv1\_1, conv2\_1, conv3\_1, etc.)
- Compute Gram matrices at each scale
- Weight contributions differently

2. Progressive Upsampling:

- Start with lower resolution
- Gradually increase resolution
- Maintain consistency across scales

3. Attention to Different Scales:

```
style_layers = ['conv1_1', 'conv2_1', 'conv3_1', 'conv4_1']
style_weights = [1.0, 0.8, 0.5, 0.3]
```

This multi-scale approach helps capture both fine details and larger structural patterns in the style transfer process."

## Question 7: "What are the common challenges in style transfer and how do we address them?" (风格迁移的挑战)

Model Answer:

"There are several key challenges in style transfer that require specific solutions:

1. Content-Style Trade-off:

- Challenge: Balancing content preservation with style transfer
- Solution: Careful weight tuning ( $\alpha$  and  $\beta$  parameters)
- Example: Higher  $\alpha$  for portraits, lower for landscapes

2. Instability Issues:

- Challenge: Unstable optimization process
- Solutions:
  - Total Variation regularization
  - Learning rate scheduling
  - Multiple resolution passes

3. Style Inconsistency:

- Challenge: Inconsistent style application
- Solutions:
  - Instance Normalization
  - Adaptive Instance Normalization (AdaIN)
  - Style-attentional networks



#### 4. Computational Efficiency:

- Challenge: Slow optimization process
- Solutions:
  - Feedforward networks
  - Network pruning
  - Mobile-optimized architectures"

### Question 8: "Could you explain the role of attention mechanisms in modern style transfer?" (注意力机制在风格迁移中的应用)

Model Answer:

"Attention mechanisms have revolutionized style transfer by enabling more precise style-content matching:

#### 1. Self-Attention in Style:

- Computes attention maps within style features
- Captures long-range dependencies
- Formula:  $\text{Attention}(Q, K, V) = \text{softmax}(QK^T)V$

#### 2. Cross-Attention between Content and Style:

- Matches content regions with appropriate style patterns
- Enables local style transfer
- Better preservation of content structure

#### 3. Implementation:

```
# Simplified attention mechanism
Q = content_features
K, V = style_features
attention = torch.softmax(torch.matmul(Q, K.transpose(-2, -1)), dim=-1)
output = torch.matmul(attention, V)
```

#### 4. Benefits:

- More precise style transfer
- Better detail preservation
- Reduced artifacts
- Improved content-style alignment"

## 7. Sequence Learning and Modern Architectures

---

### Question 1: "Can you explain the different types of sequence modeling tasks and give examples in visual computing?" (序列建模任务的类型)

Model Answer:

"Yes, there are four main types of sequence modeling architectures:

1. One-to-One Mapping:

- Single input to single output
- Example: Image classification, where we map one image to one label

2. One-to-Many Mapping:

- Single input generating multiple sequential outputs
- Example: Image captioning, where one image generates a sequence of words
- Applications: Visual story generation, multi-label image annotation

3. Many-to-One Mapping:

- Sequence of inputs producing a single output
- Example: Video classification, where multiple frames determine one category
- Applications: Action recognition, motion analysis

4. Many-to-Many Mapping:

- Sequence to sequence transformation
- Example: Video-to-text description, machine translation
- Applications: Video prediction, animation synthesis

These architectures are fundamental to how we process sequential visual data. For instance, in image captioning, we first encode an image using a CNN, then generate text word by word using an RNN or Transformer."

### Question 2: "What are the main challenges with RNNs, and how does LSTM address them?" (RNN的问题与LSTM的解决方案)

Model Answer:

"RNNs face two major challenges in processing sequences: vanishing and exploding gradients. Let me explain these issues and LSTM's solutions:

1. RNN Issues:

- Vanishing Gradients: Long-term dependencies get lost as gradients become too small
- Exploding Gradients: Gradients become too large, making training unstable
- Information bottleneck: Simple state updates can't maintain complex information

2. LSTM Solution Architecture:

LSTM (Long Short-Term Memory) addresses these through a sophisticated gating mechanism:

a) Gate Structure:

- Forget Gate (f): Decides what information to discard
- Input Gate (i): Controls new information flow
- Output Gate (o): Determines what to output
- Cell State (c): Long-term memory storage

b) Mathematical Formulation:

```
ft = σ(wf[ht-1, xt] + bf)
it = σ(wi[ht-1, xt] + bi)
ot = σ(wo[ht-1, xt] + bo)
ct = ft ⊙ ct-1 + it ⊙ tanh(wc[ht-1, xt] + bc)
ht = ot ⊙ tanh(ct)
```

The effectiveness of LSTM comes from its ability to:

- Maintain long-term dependencies through cell state
- Control information flow through gates
- Learn what to remember and forget"

### Question 3: "Can you explain the basic architecture and working principle of RNN?" (RNN基本原理)

Model Answer:

"Let me explain RNNs comprehensively:

1. Basic Structure:

```
ht = tanh(whh * ht-1 + wxh * xt + bh)
yt = why * ht + by
```

Where:

- ht is current hidden state
- xt is current input
- yt is output
- W matrices are weights

2. Information Flow:

- Takes current input xt
- Combines with previous state ht-1
- Produces new state ht
- Generates output yt

3. Key Features:

- Parameter sharing across time steps
- Memory through hidden states
- Ability to process variable-length sequences

#### 4. Visual Computing Applications:

- Video frame prediction
- Action recognition
- Image sequence analysis"

### Question 4: "How does adding attention to RNN improve its performance? Can you explain the RNN+Attention architecture?" (RNN与注意力机制的结合)

Model Answer:

"RNN with attention combines sequential processing with focused feature selection:

#### 1. Architecture Components:

##### a) RNN Base:

- Processes sequence step by step
- Maintains hidden states

##### b) Attention Layer:

```
# At each decoder step t:  
score =  $v^T \tanh(w_1 \cdot \text{encoder\_states} + w_2 \cdot h_t)$   
attention = softmax(score)  
context =  $\sum(\text{attention} * \text{encoder\_states})$ 
```

#### 2. Working Process:

- Encoder RNN processes input sequence
- Stores all intermediate states
- Decoder uses attention to focus on relevant states
- Combines context with current state for prediction

#### 3. Advantages:

- Better handling of long sequences
- Direct access to all encoder states
- More interpretable through attention weights
- Improved gradient flow

#### 4. Example in Image Captioning:

```
Input Image → CNN → Feature Map  
                ↓  
            RNN with Attention  
                ↓  
            word Sequence
```

The attention mechanism allows the model to focus on different image regions when generating each word."

Additional Important Points:

1. Implementation Details:

- Attention can be applied at different layers
- Multiple attention heads possible
- Different attention score functions available

2. Training Considerations:

- Requires more memory than basic RNN
- Can be trained end-to-end
- May need gradient clipping

3. Visual Examples:

- Attention maps can be visualized
- Shows where model focuses
- Helps in model interpretation

## Question 5: "Could you explain how the attention mechanism works and why it's important?" (注意力机制)

Model Answer:

"The attention mechanism is a breakthrough in sequence modeling. Let me break it down:

1. Core Components:

- Query (Q): What we're looking for
- Key (K): What we match against
- Value (V): The actual information we want to extract

2. Attention Calculation:

$$\text{Attention}(Q, K, V) = \text{softmax}(QK^T / \sqrt{d_k})V$$

Where:

- $QK^T$  computes compatibility scores
- $\sqrt{d_k}$  is a scaling factor
- Softmax normalizes the scores
- Final multiplication with V produces weighted sum

3. Types of Attention:

a) Self-Attention:

- Relates different positions in a single sequence
- Used in transformers for contextual understanding
- Example: Understanding relationships between all parts of an image

b) Cross-Attention:

- Relates positions from two different sequences

- Example: Image captioning, relating image regions to words

#### 4. Visual Applications:

- In image recognition: Focusing on relevant image regions
- In video analysis: Tracking important temporal relationships
- In image generation: Controlling which parts to modify"

## Question 6: "Walk me through the architecture of a Transformer and explain its advantages." (Transformer架构)

Model Answer:

"The Transformer architecture revolutionized sequence modeling. Let me explain its key components:

#### 1. Overall Structure:

- Encoder: Processes input sequence
- Decoder: Generates output sequence
- Multiple layers of self-attention and feed-forward networks

#### 2. Key Components:

##### a) Multi-Head Attention:

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W_O$$

where  $\text{head}_i = \text{Attention}(Q W_i^Q, K W_i^K, V W_i^V)$

- Allows attention from different perspectives
- Captures different types of relationships

##### b) Position Encoding:

$$\text{PE}(\text{pos}, 2i) = \sin(\text{pos}/10000^{(2i/d)})$$

$$\text{PE}(\text{pos}, 2i+1) = \cos(\text{pos}/10000^{(2i/d)})$$

- Adds positional information
- Enables sequence order understanding

#### 3. Advantages:

- Parallel processing capability
- Better long-range dependency handling
- More stable training
- Scale to longer sequences

#### 4. Applications in Visual Computing:

- Vision Transformers (ViT) for image classification
- DETR for object detection
- Image GPT for image generation"

## Question 7: "What is the role of Softmax in sequence models, and how does it work?" (Softmax函数)

Model Answer:

"Softmax is crucial in sequence models for converting raw scores into probabilities. Let me explain:

1. Mathematical Definition:

$$\text{softmax}(x_i) = \frac{\exp(x_i)}{\sum \exp(x_j)}$$

2. Properties:

- Outputs sum to 1
- Preserves relative ordering
- Differentiable
- Emphasizes larger values

3. Applications in Sequence Models:

- Attention weight normalization
- Next token prediction
- Classification probabilities
- Sequence likelihood estimation

4. Implementation Considerations:

- Numerical stability (subtracting max)
- Temperature scaling
- Gradient computation
- Memory efficiency"

Follow-up Question 1: "How would you implement attention in a practical visual computing task?"

Model Answer:

"Let me explain using image captioning as an example:

1. Visual Feature Extraction:

- Use CNN to extract image features (e.g., ResNet)
- Get feature map of shape (H×W×C)

2. Attention Implementation:

```
# Image features as Keys and Values
K = transform_key(image_features)    # [H*W, d_k]
V = transform_value(image_features)  # [H*W, d_v]

# Text decoder state as Query
Q = transform_query(decoder_state)   # [1, d_k]

# Compute attention weights
scores = torch.matmul(Q, K.transpose(-2, -1)) # [1, H*W]
attention_weights = torch.softmax(scores/√d_k, dim=-1)

# Get context vector
context = torch.matmul(attention_weights, V)
```

3. Practical Considerations:

- Batch processing for efficiency
- Memory management for large images
- Gradient scaling for stability"

Follow-up Question 2: "What are the computational trade-offs between RNNs and Transformers?"

Model Answer:

"The key trade-offs are:

1. Computational Complexity:

- RNN:  $O(n)$  sequential operations
- Transformer:  $O(n^2)$  parallel operations

2. Memory Usage:

- RNN:  $O(1)$  memory for hidden states
- Transformer:  $O(n^2)$  for attention matrices

3. Training Characteristics:

- RNN: Harder to parallelize, stable training
- Transformer: Highly parallelizable, needs careful initialization

4. Sequence Length Handling:

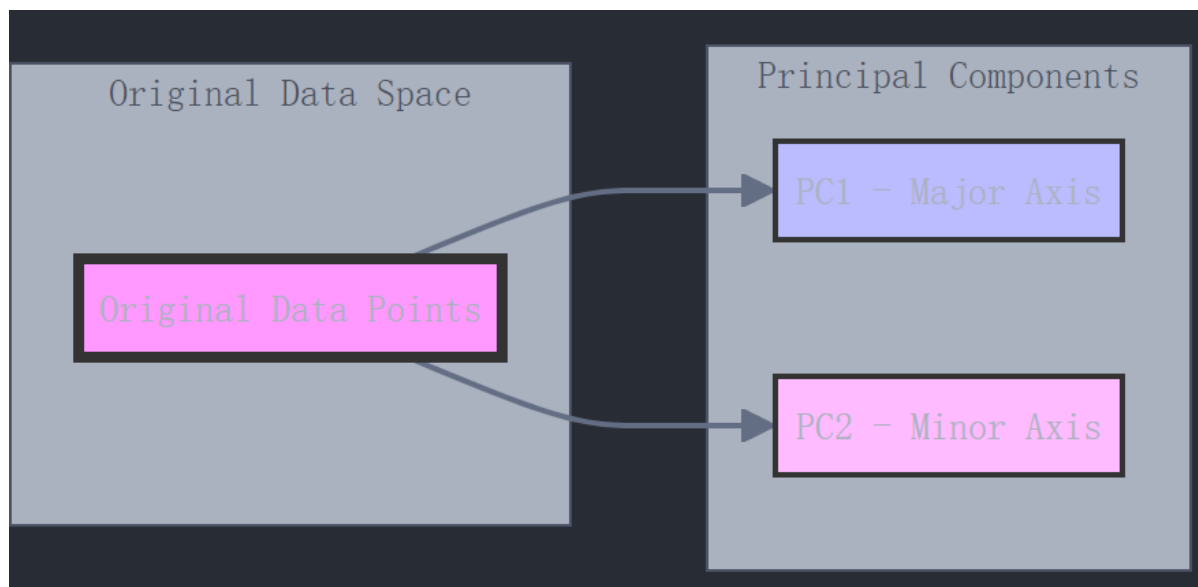
- RNN: Natural for variable lengths
- Transformer: Needs padding, position encoding"

## 8. Generative Models

**Question 1: "Can you explain what PCA is and its role in generative modeling?" (主成分分析)**

Model Answer:





"Principal Component Analysis (PCA) is a fundamental technique in dimensionality reduction and generative modeling. Let me explain its key concepts:

#### 1. Core Principle:

- PCA finds orthogonal directions of maximum variance in the data
- These directions are called principal components (主成分)
- They form a new coordinate system for the data

#### 2. Mathematical Process:

- Compute the covariance matrix of centered data
- Find eigenvectors and eigenvalues
- Sort components by explained variance
- Project data onto selected components

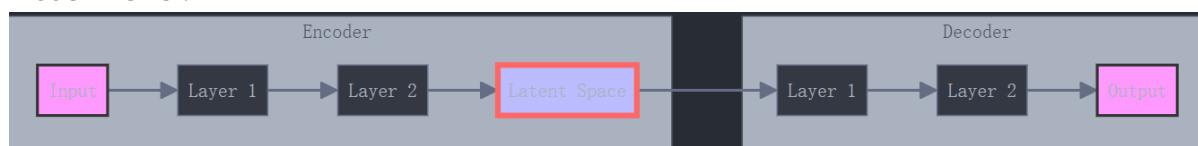
#### 3. Application in Generative Models:

- Dimensionality reduction before generation
- Understanding data distribution
- Feature decorrelation
- Basis for more complex generative techniques

PCA shows us how to represent complex data in simpler forms, which is a foundational concept for modern generative models."

## Question 2: "What are autoencoders, and how do they differ from PCA?" (自编码器)

Model Answer:



"An autoencoder is a neural network architecture that learns to compress and reconstruct data. Let me break down its components:

### 1. Structure:

- Encoder: Compresses input to latent space
- Latent space: Compressed representation
- Decoder: Reconstructs input from latent space

### 2. Key Differences from PCA:

- Non-linear transformations possible
- Can learn more complex patterns
- Requires training (unlike PCA's direct solution)
- May capture more meaningful features

### 3. Training Process:

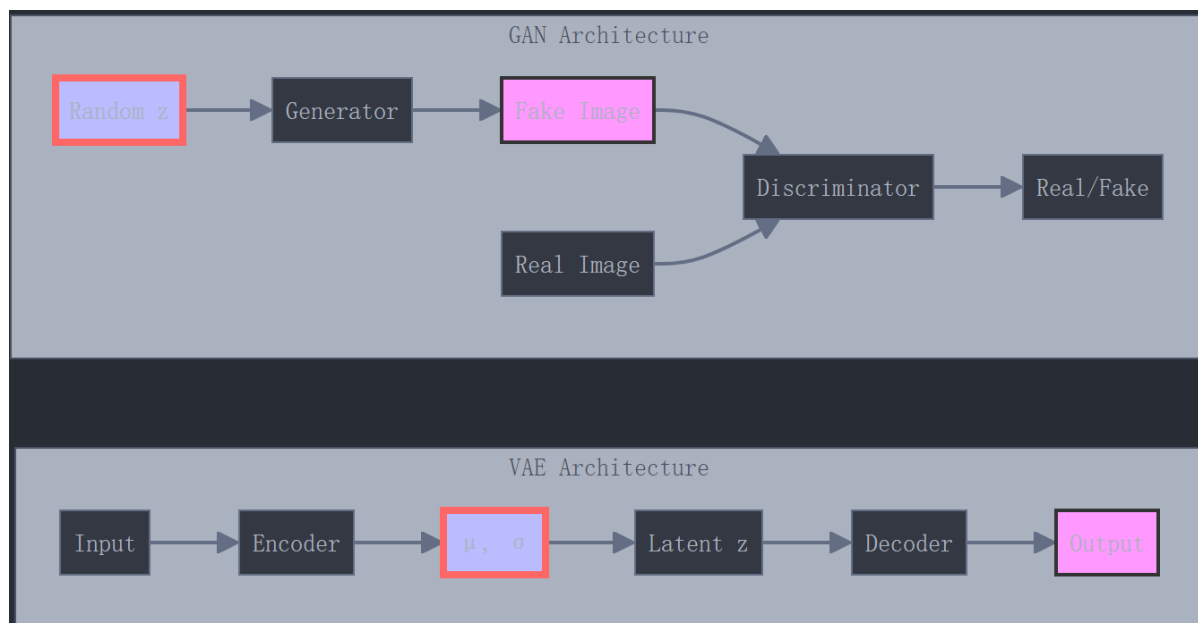
- Input = Target (self-supervised)
- Loss: reconstruction error
- Typically uses MSE or cross-entropy loss

### 4. Applications:

- Feature learning
- Denoising
- Anomaly detection
- Generative modeling foundation"

## Question 3: "Could you compare VAEs and GANs? What are their respective strengths and weaknesses?" (VAE和GAN的比较)

Model Answer:



"Let me compare these two fundamental generative approaches:

VAE (Variational Autoencoder):

### 1. Architecture:

- Encoder-decoder structure
- Probabilistic latent space
- Learns distribution parameters

#### 2. Training:

- Direct likelihood optimization
- Reconstruction + KL loss
- More stable training

#### 3. Advantages:

- Explicit latent space
- Better for representation learning
- More stable training
- Provides encoder for inference

GAN (Generative Adversarial Network):

#### 1. Architecture:

- Generator and Discriminator
- Adversarial training
- No explicit encoding

#### 2. Training:

- Minimax game
- Can be unstable
- Mode collapse issues

#### 3. Advantages:

- Usually sharper outputs
- Better visual quality
- Can learn complex distributions

Choice depends on application:

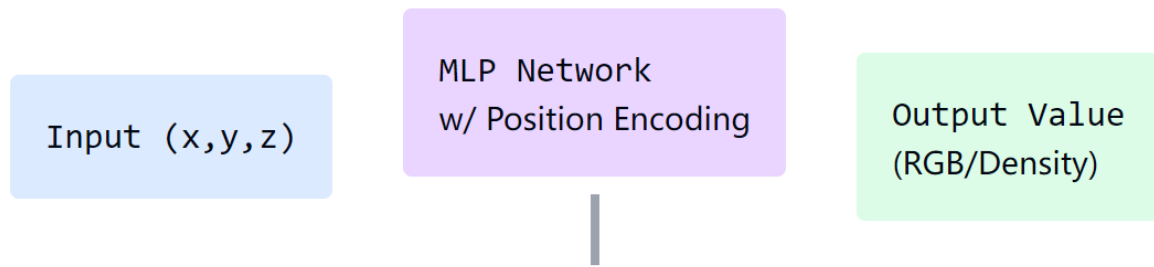
- Need latent space? → VAE
- Need highest quality? → GAN
- Need stability? → VAE
- Need realism? → GAN"

## **Question 4: "What is Neural Implicit Representation, and why is it important for generative models?" (神经隐式表示)**

Model Answer:

"Neural Implicit Representation is a modern approach to representing continuous signals with neural networks. Let me explain:

# Neural Implicit Representation



## 1. Core Concept:

- Instead of discrete values (像素/voxels)
- Learn a continuous function  $f(x) \rightarrow y$
- x: coordinates (spatial/temporal)
- y: signal value (color/density)

## 2. Key Components:

- MLP architecture
- Positional encoding
- Continuous optimization

## 3. Applications:

- NeRF (Neural Radiance Fields)
- Shape representation
- Image/video generation
- 3D scene reconstruction

## 4. Advantages:

- Resolution-independent
- Compact representation
- Continuous interpolation
- Natural regularization"

## Question 5: "Explain autoregressive models and how they work in generative tasks." (自回归模型)

Model Answer:

"Autoregressive models generate data sequentially, one element at a time. Let me elaborate:

### 1. Basic Principle:

$$p(x) = \prod_i p(x_i | x_1, \dots, x_{i-1})$$

### 2. Key Characteristics:

- Sequential generation
- Explicit probability modeling
- Each step depends on previous ones

- Natural ordering required

### 3. Applications:

- PixelCNN for images
- Language modeling
- Audio synthesis
- Time series generation

### 4. Advantages/Limitations:

- Tractable likelihood
- Stable training
- Slow generation
- Need for sequential ordering

This approach is particularly powerful when the data has natural sequential structure."

## Question 6: "How does the reparameterization trick work in VAEs, and why is it necessary?" (重参数化技巧)

Model Answer:

"The reparameterization trick is crucial for training VAEs as it enables backpropagation through random sampling. Let me explain:

### 1. Problem It Solves:

- Need to sample from learned distribution
- Can't backpropagate through random sampling
- Need differentiable operation

### 2. Implementation:

Instead of directly sampling  $z \sim N(\mu, \sigma^2)$ , we do:

- Sample  $\epsilon \sim N(0,1)$
- Compute  $z = \mu + \sigma * \epsilon$
- This makes the sampling process differentiable

### 3. Mathematical Expression:

Original:  $z \sim N(\mu, \sigma^2)$

Reparameterized:  $z = \mu + \sigma * \epsilon$ , where  $\epsilon \sim N(0,1)$

### 4. Benefits:

- Enables gradient flow
- Reduces variance in training
- Maintains randomness
- Allows end-to-end training"

## Question 7: "How do you handle mode collapse in GANs?" (模式崩塌)

Model Answer:

"Mode collapse is a common challenge in GAN training where the generator produces limited varieties of outputs. Let me explain the solutions:

1. Detection:

- Monitor diversity of generated samples
- Check discriminator loss patterns
- Visualize generated distributions

2. Solutions:

a) Architectural:

- Use minibatch discrimination
- Implement instance noise
- Add regularization terms

b) Training Strategies:

- Wasserstein loss
- Progressive growing
- Two-timescale update rule (TTUR)

3. Best Practices:

- Batch normalization
- Proper learning rate scheduling
- Regular diversity evaluation
- Balance generator/discriminator training

This is why modern GAN variants like StyleGAN incorporate specific mechanisms to prevent mode collapse."

## Question 8: "What is the role of the Evidence Lower Bound (ELBO) in VAEs?" (ELBO)

Model Answer:

[I would write the ELBO equation]

"The Evidence Lower Bound (ELBO) is the objective function for training VAEs. Let me break it down:

1. Mathematical Expression:

$$\text{ELBO} = \mathbb{E}[\log p(x|z)] - \text{KL}(q(z|x) || p(z))$$

Where:

- First term: reconstruction loss
- Second term: KL divergence regularization
- $q(z|x)$ : encoder (inference model)

- $p(z)$ : prior distribution
  - $p(x|z)$ : decoder (generative model)
2. Components:
- a) Reconstruction Term:
- Measures how well we can reconstruct input
  - Usually implemented as MSE or cross-entropy

b) KL Divergence Term:

- Regularizes latent space
- Pushes toward prior distribution
- Enables sampling from prior

3. Optimization:

- Maximize ELBO
- Balance between reconstruction and regularization
- Can be optimized using standard gradient descent"

## Question 9: "What are the different types of generative models and how do they compare?" (生成模型类型比较)

Model Answer:

[I would draw a taxonomy diagram]

"Let me categorize and compare the main types of generative models:

1. Explicit Density Models:

a) Tractable:

- Autoregressive models (PixelCNN)
- Normalizing flows
- Advantages: exact likelihood
- Disadvantages: sequential generation

b) Approximate:

- VAEs
- Advantages: fast generation, stable training
- Disadvantages: blurry outputs

2. Implicit Models:

- GANs
- Advantages: high quality samples
- Disadvantages: mode collapse, training instability

3. Hybrid Approaches:

- VAE-GAN
- Flow-based GANs

- Advantages: combines benefits
- Disadvantages: complex training

Each type has its own use cases and trade-offs."

## **Question 10: "What is the relationship between the latent space dimension and model performance in VAEs?" (潜在空间维度)**

Model Answer:

"The latent space dimension is a crucial hyperparameter in VAEs. Let me explain its effects:

### 1. Impact on Representation:

- Smaller dimension: more compression
- Larger dimension: more capacity
- Trade-off between compression and reconstruction

### 2. Effects:

#### a) Too Small:

- Poor reconstruction
- Over-regularization
- Limited generative capability

#### b) Too Large:

- Risk of posterior collapse
- Computational overhead
- Possible redundancy

### 3. Selection Criteria:

- Task complexity
- Data dimensionality
- Computational constraints
- Required generation quality

### 4. Best Practices:

- Start with small dimension
- Gradually increase while monitoring reconstruction
- Use regularization techniques
- Consider hierarchical structures"



## 9. Advanced Topics in Generation

### Question 1: "Could you explain what a diffusion process is and how it works in generative models?" (扩散过程)

Model Answer:

"A diffusion process is a technique for generative modeling that gradually converts data into noise and then learns to reverse this process. Let me explain the key components:

Forward Process (正向过程):

- Gradually adds Gaussian noise to data  $x_0$
- Through T steps:  $x_0 \rightarrow x_1 \rightarrow \dots \rightarrow x_t$
- Each step follows:  $x_t = \sqrt{\alpha_t}x_0 + \sqrt{(1-\alpha_t)}\epsilon_t$   
where  $\epsilon_t$  is random Gaussian noise

Reverse Process (反向过程):

- Learns to denoise:  $x_t \rightarrow x_{t-1} \rightarrow \dots \rightarrow x_0$
- Uses neural network to predict noise or clean data
- Each step estimates:  $p(x_{t-1} | x_t)$

This process is particularly effective because:

1. It's based on well-understood mathematics
2. The forward process is fixed and tractable
3. Only the reverse process needs to be learned
4. It produces high-quality samples"

### Question 2: "How does the training process work in diffusion models?" (训练过程)

Model Answer:

"The training process in diffusion models involves several key steps:

1. Training Objective:

- Minimize the difference between predicted and actual noise
- Loss function:  $L = E[||\epsilon - \epsilon_\theta(x_t, t)||^2]$
- Where  $\epsilon_\theta$  is the noise prediction network

2. Training Loop:

For each batch:

- Sample real images  $x_0$
- Sample timesteps  $t$
- Add noise to create  $x_t$
- Predict noise  $\epsilon_\theta(x_t, t)$
- Calculate loss and update parameters

3. Network Architecture:

- Usually U-Net based
- Conditions on timestep  $t$
- May include attention layers
- Takes noisy image as input

The key innovation is that we train on all noise levels simultaneously, making the model robust and stable."

### Question 3: "Could you explain the sampling process in diffusion models and why Gaussian noise is important?" (采样过程和高斯噪声)

Model Answer:

"Let me explain both the sampling process and the role of Gaussian noise:

Sampling Process:

1. Start with pure noise  $x_t \sim N(0, I)$
2. Iteratively apply:

$$x_{t-1} = 1/\sqrt{\alpha_t} (x_t - (1-\alpha_t)/\sqrt{(1-\bar{\alpha}_t)} \epsilon_{\theta}(x_t, t))$$

3. Continue until we reach  $x_0$

Gaussian Noise is important because:

1. Mathematical Properties:
  - Easily parameterizable
  - Well-understood statistics
  - Nice addition properties
2. Training Benefits:
  - Smooth transitions between steps
  - Stable gradients
  - Predictable behavior
3. Sampling Quality:
  - Natural image manifold coverage
  - Diverse generation
  - Controlled randomness"

### Question 4: "What are the key components of Latent Diffusion Models (LDM) and how do they differ from regular diffusion models?" (潜在扩散模型)

Model Answer:

"Latent Diffusion Models (LDM) introduce several important modifications to the basic diffusion framework:

Key Components:

### 1. Autoencoder:

- Compresses images to latent space
- Reduces computational cost
- Preserves semantic information

### 2. Diffusion in Latent Space:

- Operates on compressed representations
- More efficient than pixel-space diffusion
- Maintains generation quality

### 3. Conditioning Mechanisms:

- Text embeddings
- Cross-attention
- Additional control signals

The main advantages over regular diffusion:

1. Faster training and inference
2. Lower memory requirements
3. Better control over generation
4. More flexible conditioning options"

## **Question 5: "Could you walk me through a practical example of using diffusion models for image generation?" (实际应用示例)**

Model Answer:

"Let me explain with a concrete example of text-to-image generation:

Process Steps:

### 1. Input Phase:

```
Text prompt: "A red car in front of a sunset"  
- Encode text using transformer  
- Initialize random noise
```

### 2. Denoising Steps:

- Start with pure noise
- Apply model iteratively
- Each step becomes clearer
- Model conditions on text embedding

### 3. Final Generation:

- Get final denoised image
- Apply any post-processing
- Ensure quality control

## Key Considerations:

1. Number of sampling steps
2. Guidance scale for text conditioning
3. Seed for reproducibility
4. Batch processing for variations

This process can be extended to various applications like inpainting, image editing, and style transfer."

## Additional Follow-up Questions They Might Ask:

Follow-up 1: "How do we choose the number of diffusion steps?"

"The choice of diffusion steps involves a trade-off between quality and computation:

- Typical range is between 50-1000 steps
- More steps (like 1000) give higher quality but slower generation
- Fewer steps (like 50) are faster but may reduce quality
- Modern techniques like DDIM allow for fewer steps while maintaining quality
- We often validate different step counts empirically on our specific task"

Follow-up 2: "What are the trade-offs between speed and quality in diffusion models?"

"The main trade-offs are:

1. Sampling Steps:
  - More steps = better quality but slower generation
  - Fewer steps = faster but potential quality loss
2. Model Size:
  - Larger models = better quality but slower and more memory
  - Smaller models = faster but reduced capability
3. Latent Space:
  - Pixel space = higher quality but slower
  - Compressed latent space = faster but potential detail loss"

Follow-up 3: "How does conditioning work in diffusion models?"

"Conditioning in diffusion models can be implemented in several ways:

1. Classifier Guidance:
  - Uses gradients from a classifier
  - Helps steer generation towards desired attributes
2. Cross-Attention:
  - Incorporates condition (like text) through attention
  - Allows for fine-grained control
3. Direct Concatenation:
  - Append condition to input
  - Simple but effective for basic conditioning

The choice depends on the type of condition (text, class labels, images) and desired control level."

Follow-up 4: "What are the main challenges in training diffusion models?"

"The main challenges include:

1. Computational Resources:
  - Long training times
  - High memory requirements
  - Need for large datasets
2. Optimization Difficulties:
  - Balancing noise prediction accuracy
  - Managing training stability
  - Choosing correct hyperparameters
3. Quality Control:
  - Ensuring consistent generation quality
  - Handling mode collapse
  - Maintaining diversity in generation
4. Architectural Choices:
  - Selecting appropriate model architecture
  - Balancing model size vs performance
  - Implementing efficient attention mechanisms"

## 10. Visual Computing Specific Applications

---

**Question 1: "Can you explain the different ways we represent images and video data in visual computing?" (图像和视频表示)**

Model Answer:

"Visual data can be represented in several ways:

For Images:

- Pixel arrays (像素数组):  $H \times W \times C$  format
- Vector graphics (矢量图): Using control points and curves
- Feature representations: Using deep features or descriptors

For Videos:

- Sequence of frames (帧序列)
- Motion vectors (运动向量)
- Temporal features (时序特征)

Each representation has specific advantages:

- Pixel arrays are direct and simple
- Vector graphics are resolution-independent

- Feature representations support high-level understanding"

## **Question 2: "What are the main approaches for image retrieval and how do they differ?" (图像检索)**

Model Answer:

"There are several key approaches for image retrieval:

1. RGB Distance:

- Direct pixel comparison
- Simple but sensitive to position changes

2. Histogram Methods:

- Color histogram comparison
- Position-invariant but loses spatial information

3. VGG Features:

- Deep neural network features
- Better at semantic similarity
- More robust to variations

4. Gram Matrix:

- Captures style and texture information
- Used in style transfer applications

The trend has been moving from simple pixel-based methods to more sophisticated semantic approaches using deep learning."

## **Question 3: "How does image denoising work, and what are the key differences between traditional and CNN-based approaches?" (图像去噪)**

Model Answer:

"Image denoising typically involves two main approaches:

Traditional (Gaussian):

- Uses fixed kernels
- Local neighborhood operations
- Simple but can blur details

CNN-based:

- Learns optimal filters from data
- Can capture complex noise patterns
- Better at preserving details
- Requires training data

The main advantage of CNN-based methods is their ability to adapt to specific types of noise while maintaining image structure."

## Question 4: "What is NeRF and what are its key applications?" (神经辐射场)

Model Answer:

"NeRF (Neural Radiance Fields) is a neural rendering approach:

Core Concept:

- Represents 3D scenes as continuous functions
- Maps 5D coordinates (position + viewing direction) to color and density

Key Components:

- Positional encoding (位置编码)
- MLP network (多层感知机)
- Volume rendering (体渲染)

Applications:

1. Novel view synthesis
2. 3D scene reconstruction
3. Virtual/Augmented reality content creation"

## Question 5: "What is an image pyramid and why is it useful?" (图像金字塔)

Model Answer:

"An image pyramid is a multi-scale representation:

Structure:

- Series of images at different resolutions
- Each level is typically half the size of previous
- Created through successive downsampling

Applications:

1. Multi-scale feature detection
2. Efficient search across scales
3. Progressive image transmission
4. Texture synthesis

The main advantage is the ability to analyze images at multiple scales efficiently, capturing both fine details and global structure."

## Question 6: "Can you explain what semantic tasks are in visual computing and give some examples?" (语义任务)

Model Answer:

"Semantic tasks in visual computing deal with understanding the meaning and content of visual data. Let me explain the main types:

1. Semantic Segmentation (语义分割):

- Assigns class labels to each pixel
- Example: Identifying roads, buildings, people in scenes

#### 2. Object Detection (目标检测):

- Locates and classifies objects
- Outputs bounding boxes and class labels

#### 3. Scene Understanding (场景理解):

- Analyzes relationships between objects
- Interprets context and spatial layout

These tasks are fundamental in applications like autonomous driving and medical imaging."

## **Question 7: "What are different types of 3D surface representations and their applications?" (三维表面表示)**

Model Answer:

"3D surfaces can be represented in several ways:

#### 1. Point Clouds (点云):

- Collection of 3D points
- Simple but lacks connectivity
- Good for scanning and LIDAR data

#### 2. Mesh Representations (网格):

- Vertices, edges, and faces
- Common in graphics and animation
- Supports efficient rendering

#### 3. Volumetric Data (体积数据):

- Regular 3D grid
- Used in medical imaging
- Good for internal structures

Each representation has specific use cases depending on the application requirements and computational constraints."

## **Question 8: "How do we handle feature extraction at different scales in images?" (多尺度特征提取)**

Model Answer:

"Multi-scale feature extraction involves several key techniques:

#### 1. Image Pyramid Methods:

- Gaussian pyramid for smoothing
- Laplacian pyramid for feature detection
- DoG (Difference of Gaussian) for scale-invariant features



## 2. Multi-resolution Processing:

- Analyzes images at different resolutions
- Captures both fine details and global structure
- Useful for hierarchical feature extraction

## 3. Scale-space Theory:

- Provides theoretical framework
- Handles scale-invariant feature detection
- Important for applications like SIFT

These approaches are crucial for robust computer vision systems that need to handle objects at different scales."