

Lecture 7:

Deep Neural Networks and Architecture Patterns

Ho-min Park

homin.park@ghent.ac.kr

powersimmani@gmail.com

Lecture Contents

Part 1: The Need for Deep Neural Networks

Part 2: Modern Activation Functions

Part 3: Advanced Architecture Patterns

Part 1/3:

The Need for Deep Neural Networks

1. Limitations of Shallow Networks
2. The Power of Depth - Hierarchical Representations
3. Feature Reuse and Composition
4. Parameter Efficiency
5. Challenges of Deep Networks
6. Vanishing Gradient Problem
7. Exploding Gradient Problem
8. Overview of Solutions

Limitations of Shallow Networks

UNIVERSAL APPROXIMATION THEOREM

Single hidden layer networks can theoretically approximate any function
BUT exponentially many hidden units may be required



No Hierarchical Patterns 1

Cannot efficiently capture hierarchical structure in images and text

Compositional Struggle 2

Must learn all combinations from scratch without reusing patterns

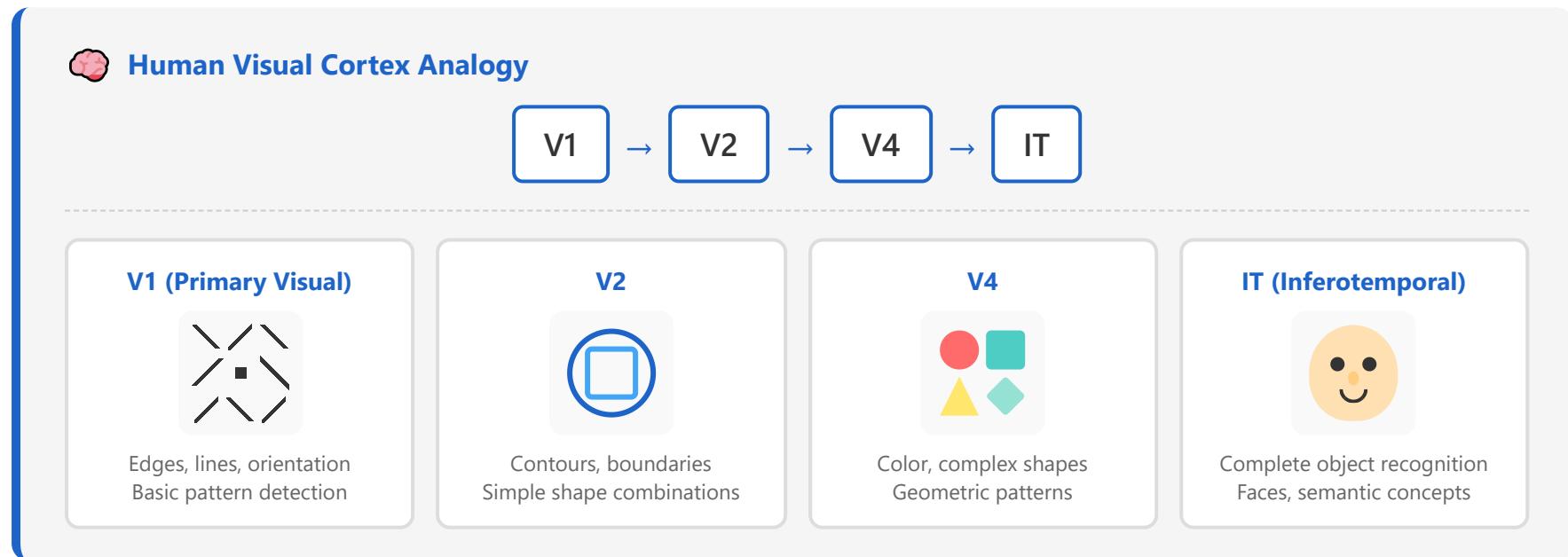
Poor Generalization 3

Weak performance on high-dimensional data with limited samples

Limited Abstraction 4

Manual feature engineering still required for complex tasks

The Power of Depth: Hierarchical Representations

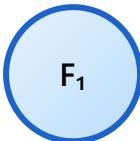


✓ Transforms representation space to be more linearly separable

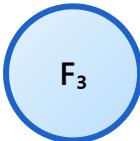
✓ Enables transfer learning: lower layers generalize across tasks

Feature Reuse and Composition

Layer I



F_2



k features



Layer I+1



P_2



P_4

...

k^2 patterns

COMPOSITIONAL EFFICIENCY

Combine k features → Create k^2 patterns



Reusable Components



Exponential Expressiveness



Modular Learning

EXAMPLE

"Eye" feature reused for detecting different faces, animals → Reduces redundancy

Parameter Efficiency: Deep vs. Shallow Networks

Shallow Network

1000 units

Layers:

Total Units:

Parameters:

1
1,000
1M+

Deep Network

100 units

100 units

100 units

Layers:

Total Units:

Parameters:

3
300
30K

Computational Complexity (n inputs)

Shallow

$O(2^n)$

vs

Deep

$O(n^2)$

✓ Better generalization

✓ Faster training

✓ Reduced redundancy

✓ Faster inference

Challenges of Deep Networks

Training Complexity



- ▶ Non-convex optimization
- ▶ Longer training time
- ▶ High computational cost

Gradient Problems



- ▶ Vanishing gradients in early layers
- ▶ Exploding gradients
- ▶ Unstable parameter updates

Generalization Issues



- ▶ Overfitting risk
- ▶ Memorizing training data
- ▶ Poor test performance

Configuration Sensitivity



- ▶ Critical initialization requirements
- ▶ Hyperparameter sensitivity
- ▶ Architecture design choices

⚠ HISTORICAL BARRIER (Pre-2006)

Deep networks often performed worse than shallow ones due to these challenges

These challenges motivated research into better training techniques and architectures

Vanishing Gradient Problem



MATHEMATICAL CAUSE

gradient $\propto \prod (\partial a_i / \partial z_i) \rightarrow \text{Product of many terms} < 1$

⚠ Consequences

- ▶ Early layers learn extremely slowly
- ▶ Weights barely change
- ▶ Training loss plateaus
- ▶ Network behaves like shallow



ReLU Activation



Skip Connections



Careful Init

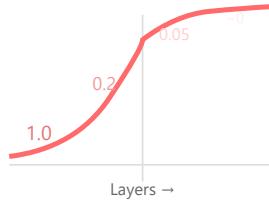


Batch Norm

Activation Function Comparison

Sigmoid (σ)

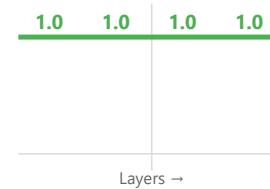
✗ Gradient Vanishing



$$\sigma'(x) = \sigma(x)(1-\sigma(x)) \leq 0.25$$

ReLU

✓ Gradient Preserved



$$ReLU'(x) = 1 \text{ (if } x > 0)$$



Activation Functions & Gradient Vanishing Risk

⚠ High Risk

Sigmoid (σ)

$$\sigma'(z) = \sigma(z)(1-\sigma(z))$$

Max gradient: 0.25

$$0.25^5 = 0.00098 \times$$

⚠ Medium Risk

Tanh

$$\tanh'(z) = 1 - \tanh^2(z)$$

Max gradient: 1.0

Saturates at extremes → vanish

✓ Lower Risk

ReLU

$$ReLU'(z) = 1 \text{ if } z > 0 \text{ else } 0$$

Gradient: 0 or 1

No saturation for $z > 0$ ✓



Chain Rule Multiplication: How Gradients Vanish

Backpropagation through layers:

$$\frac{\partial L}{\partial W_1} = \frac{\partial L}{\partial a_4} \times \frac{\partial a_4}{\partial z_4} \times \frac{\partial z_4}{\partial a_3} \times \frac{\partial a_3}{\partial z_3} \times \frac{\partial z_3}{\partial a_2} \times \frac{\partial a_2}{\partial z_2} \times \frac{\partial z_2}{\partial a_1} \times \frac{\partial a_1}{\partial z_1} \times \frac{\partial z_1}{\partial W_1}$$

Blue terms: activation derivatives (σ' , \tanh' , $ReLU'$)

Example with Sigmoid activation:

✗ Bad Case (Sigmoid)

Layer 1: grad × $W_1 \times \sigma'$

$$= 1.0 \times 1.0 \times 0.25 = 0.25$$

$$\text{Layer 2: } 0.25 \times 1.0 \times 0.25 = 0.0625$$

$$\text{Layer 3: } 0.0625 \times 1.0 \times 0.25 = 0.0156$$

$$\text{Layer 4: } 0.0156 \times 1.0 \times 0.25 \approx 0.004$$

✓ Good Case (ReLU)

Layer 1: $1.0 \times 1.0 \times 1.0 = 1.0$

Layer 2: $1.0 \times 1.0 \times 1.0 = 1.0$

Layer 3: $1.0 \times 1.0 \times 1.0 = 1.0$

Layer 4: $1.0 \times 1.0 \times 1.0 = 1.0$

10 layers: $0.25^{10} \approx 0.00000095$
Practically zero! 

Gradient preserved!
Can train very deep networks ✓

Real Vanishing Scenario

Sigmoid Max Gradient

$$\sigma' = 0.25$$



10 Layers

$$n = 10$$



Gradient Scale

$$0.25^{10} \approx 0.000001$$

With 20 layers: $0.25^{20} \approx 9.09 \times 10^{-13}$ → Effectively zero! 

Exploding Gradient Problem



MATHEMATICAL CAUSE

gradient $\propto \prod (\partial a_i / \partial z_i)$ → Product of many terms > 1

Consequences

Unstable parameter updates

Loss becomes NaN or infinity

Training loss spikes dramatically

Weights become very large



Gradient Clipping



Xavier/He Init



Batch Normalization



Activation Functions & Gradient Explosion Risk

⚠ High Risk

Sigmoid (σ)

$$\sigma'(z) = \sigma(z)(1-\sigma(z))$$

Max gradient: 0.25

⚠ Medium Risk

Tanh

$$\tanh'(z) = 1 - \tanh^2(z)$$

Max gradient: 1.0

✓ Lower Risk

ReLU

$$\text{ReLU}'(z) = 1 \text{ if } z>0 \text{ else } 0$$

Gradient: 0 or 1

0.25⁵ = 0.00098 ❌

Poor weight init → explode

But large W → still explode

1	2
3	4

Chain Rule Multiplication: How Gradients Explode

Backpropagation through layers:

$$\frac{\partial L}{\partial W_1} = \frac{\partial L}{\partial a_4} \times \frac{\partial a_4}{\partial z_4} \times \frac{\partial z_4}{\partial a_3} \times \frac{\partial a_3}{\partial z_3} \times \frac{\partial z_3}{\partial a_2} \times \frac{\partial a_2}{\partial z_2} \times \frac{\partial z_2}{\partial a_1} \times \frac{\partial a_1}{\partial z_1} \times \frac{\partial z_1}{\partial W_1}$$

Red terms: activation derivatives (σ' , \tanh' , ReLU')

Example with poorly initialized weights:

❌ Bad Case (W ~ 2.0)

$$\begin{aligned} \text{Layer 1: } & \text{grad} \times W_1 \times \sigma' \\ &= 1.0 \times 2.0 \times 0.25 = 0.5 \\ \text{Layer 2: } & 0.5 \times 2.0 \times 0.25 = 0.25 \\ \text{Layer 3: } & 0.25 \times 2.0 \times 0.25 = 0.125 \\ \text{Layer 4: } & 0.125 \times 2.0 \times 0.25 = 0.0625 \\ \text{If } W=3: & 3 \times 0.25 = 0.75 \rightarrow \text{vanish!} \\ \text{If } W=5: & 5 \times 0.25 = 1.25 \rightarrow \text{explode!} \end{aligned}$$

✓ Good Case (Xavier Init)

$W \sim N(0, 1/\sqrt{n})$
Keeps gradient variance stable
 $E[\partial L/\partial W] \approx \text{constant}$
No explosion or vanishing

Initial Weight

W = 2.5



💥 Real Explosion Scenario

5 Layers

n = 5



Gradient Scale

$2.5^5 = 97.7$

With 20 layers: $2.5^{20} = 9.09 \times 10^8 \rightarrow \text{NaN!}$ 🔥

Solutions to Deep Learning Challenges

Better Activation Functions



Prevents vanishing gradients during backpropagation

ReLU, Leaky ReLU, ELU, GELU

→ Vanishing Gradient

Normalization Techniques



Stabilizes training and speeds up convergence

Batch Norm, Layer Norm

→ Training Stability

Careful Optimization



Adaptive learning rates for stable updates

Adam, RMSprop, AdaGrad

→ Training Efficiency

Smart Initialization



Maintains gradient scale across layers

Xavier, He Initialization

→ Gradient Issues

Skip Connections



Creates gradient highways through the network

ResNet, Highway Networks

→ Vanishing Gradient

Regularization & Clipping



Prevents overfitting and gradient explosion

Dropout, Weight Decay, Grad Clipping

→ Overfitting & Explosion

Modern Deep Learning

Robust training is achieved by combining multiple techniques together

Part 2/3:

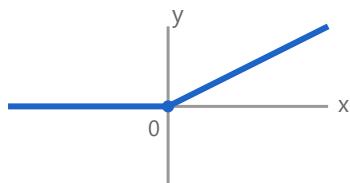
Modern Activation Functions

- 9.** The ReLU Revolution
- 10.** Leaky ReLU, PReLU
- 11.** ELU, SELU
- 12.** Swish, GELU
- 13.** Activation Function Selection Guide
- 14.** Dead ReLU Problem
- 15.** Gradient Flow Analysis
- 16.** Layer-wise Activation Patterns

The ReLU Revolution

DEFINITION

$$\text{ReLU}(x) = \max(0, x)$$



HISTORICAL IMPACT

Enabled training of **AlexNet (2012)**

First deep CNN to win **ImageNet**

→ Sparked the deep learning revolution



Prevents vanishing gradients (derivative = 1 for $x > 0$)



Computational efficiency (no exponential calculations)



Sparse activation (~50% neurons output zero)



Unchanged gradient flow through active neurons

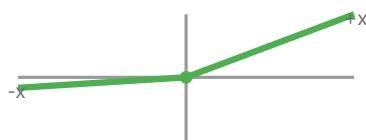
LIMITATIONS: Dead ReLU problem · Non-zero centered · Unbounded output

Leaky ReLU & PReLU: Addressing Dead Neurons

Leaky ReLU

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

$\alpha = 0.01$ (fixed)

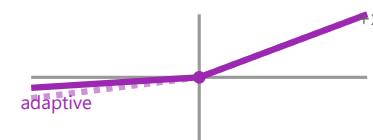


- ✓ Prevents dead neurons
- ✓ Small negative slope (0.01)
- ✓ Minimal computational cost

PReLU

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

α learned during training



- ✓ Learnable parameter α
- ✓ Adapts per channel/layer
- ✓ Better in deep networks

Performance Comparison

Aspect

Leaky ReLU

PReLU

Flexibility

Fixed α

Adaptive α

Complexity

Minimal

Slight increase

Use Case

General purpose

Deep CNNs

ELU & SELU: Advanced Activation Functions

ELU

Exponential Linear Unit

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

$$f(x) = \alpha(e^x - 1) \text{ if } x \leq 0$$



- ✓ Smooth exponential curve
- ✓ Negative values → zero mean
- ✓ Reduces bias shift
- ✓ More robust to noise

SELU

Scaled Exponential Linear Unit

$$f(x) = \lambda \times \text{ELU}(x)$$

$$\lambda \approx 1.05, \alpha \approx 1.67$$



- ✓ Self-normalizing properties
- ✓ Maintains mean/variance
- ✓ No Batch Norm needed

SPECIAL PROPERTY

Auto-normalization without Batch Normalization

Usage Comparison

Aspect	ELU	SELU
Computation	Exponential (moderate)	Exponential (moderate)
Best Use	General improvements	Fully-connected networks
Key Benefit	Faster learning vs ReLU	Self-normalization

Swish & GELU: Modern State-of-the-Art Activations

Swish

Self-Gating Activation

$$f(x) = x \times \sigma(\beta x)$$

$\beta = 1$ (Swish-1)



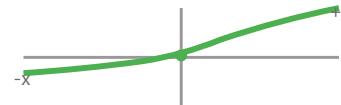
- ✓ Self-gating mechanism
- ✓ Smooth, non-monotonic
- ✓ Neural architecture search

GELU

Gaussian Error Linear Unit

$$f(x) = x \times \Phi(x)$$

Φ = Gaussian CDF



- ✓ Probabilistic interpretation
- ✓ Smooth, non-monotonic
- ✓ State-of-the-art in NLP

ADOPTION

BERT, GPT (Default)

Key Properties

- ▶ Smooth everywhere
- ▶ Outperform ReLU in Transformers
- ▶ Higher computational cost

Best Use Cases

- 🎯 Transformer models
- 🎯 NLP tasks (BERT, GPT)
- 🎯 State-of-the-art performance

Activation Function Selection Guide

DEFAULT CHOICE

Start with ReLU

Fast, works well for most cases

Deep CNNs

Leaky ReLU / PReLU

Prevents dead neurons in deep networks



Faster Convergence

ELU

Slightly slower but better accuracy



Fully-Connected Deep Networks

SELU

Eliminates need for Batch Normalization



Transformers / NLP

GELU

Standard in BERT, GPT models



Limited Compute

ReLU / Leaky ReLU

Fastest computational speed



Experimental / Research

Swish / GELU

Potential accuracy gains



General rule: Modern smooth activations (**GELU**, **Swish**) often worth the extra computational cost

Dead ReLU Problem

DEFINITION

Neurons that always output zero for all inputs (permanently inactive)

CAUSE

Large negative gradients push weights into negative region

CONSEQUENCE

Neuron stops learning,
effectively removed

IMPACT

Reduced capacity, wasted
parameters

⚠ Symptoms & Detection

Significant portion of neurons inactive (**sometimes >40%**) · Monitor percentage of always-zero activations during training

✓ Solutions



Leaky ReLU / PReLU



Lower Learning Rate



Better Initialization



Gradient Clipping

Gradient Flow Analysis

GOAL

Understand how gradients propagate through different activations

ReLU

⚠ Risky

Binary gradient (0 or 1)

Can block gradients completely when inactive

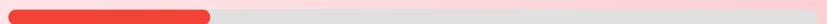


Sigmoid / Tanh

✗ Poor

Max derivative: 0.25

Gradients shrink exponentially



Leaky ReLU

✓ Good

Always allows gradient

Small slope in negative region



Smooth (ELU, GELU, Swish)

★ Best

Continuous gradients

Better for optimization



Ideal Property

Gradients neither vanish nor explode across depth

Modern Approach

Good Activations + Normalization + Skip Connections

Layer-wise Activation Patterns

KEY OBSERVATION

Different layers may benefit from different activations



Empirical Finding

Uniform activation often works well, mixed can be better

Research Direction

Neural Architecture Search explores per-layer choices

Practical Tip

If experimenting, start from output and work backward

Implementation

Start simple (uniform), optimize if needed

Transformer Pattern Example

GELU in feedforward layers · Linear (no activation) in attention mechanisms

Part 3/3:

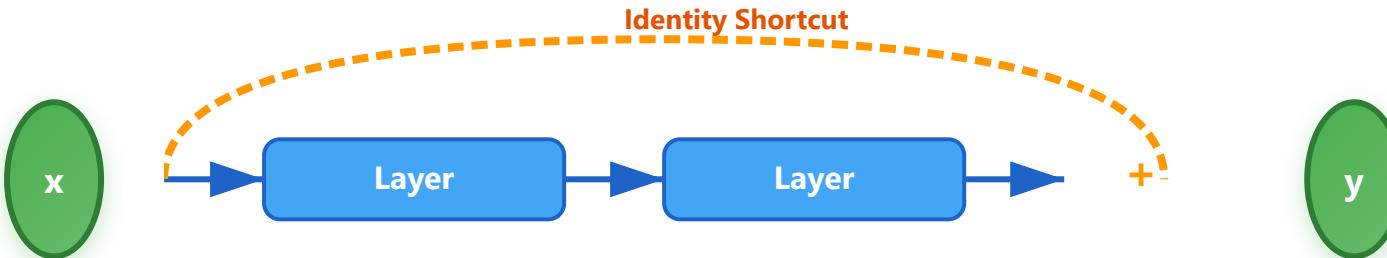
Advanced Architecture Patterns

- 17.** Skip Connection (ResNet)
- 18.** Dense Connection (DenseNet)
- 19.** Bottleneck Architecture
- 20.** The Role of 1x1 Convolution
- 21.** Inception Module
- 22.** Depthwise Separable Convolution
- 23.** Neural Architecture Search
- 24.** Model Compression Techniques
- 25.** Practical Design Guidelines

Skip Connection (ResNet): Breakthrough Architecture

CORE FORMULA

$$y = F(x) + x$$



Extreme Depth

Enabled training 152+ layer networks



Gradient Highway

Direct path for gradients to early layers



Residual Learning

Learns $F(x)$ instead of full $H(x)$



Math Benefit

Gradient always has component of 1



Impact & Adoption

ImageNet Champions

ResNet-50/101/152

Performance

Deeper = Better

Adoption

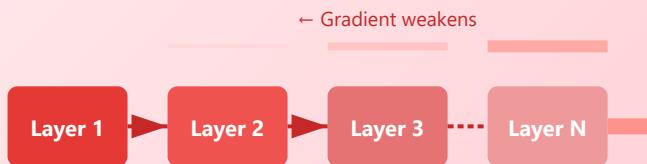
Universal Standard



Deep Dive: Understanding Skip Connections

ResNet (Residual Network), introduced by Kaiming He et al. in 2015, revolutionized deep learning by solving a fundamental problem: as neural networks get deeper, they become increasingly difficult to train. Surprisingly, adding more layers to a network would often make it perform *worse*, not better. Skip connections elegantly solved this paradox.

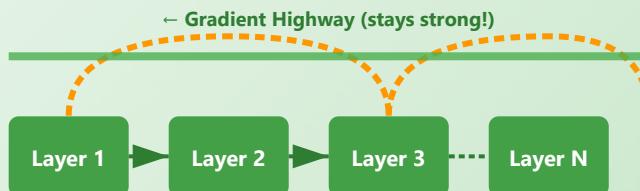
✗ The Problem: Vanishing Gradients



Gradients vanish in deep networks!

Early layers barely learn

✓ The Solution: Skip Connections



Skip connections preserve gradients!

All layers learn effectively



The Mathematics Behind ResNet

Traditional Network vs. Residual Network

Traditional: $H(x) = \text{desired output} \rightarrow$ Network must learn the entire mapping

Residual: $H(x) = F(x) + x \rightarrow$ Network only learns the "residual" $F(x) = H(x) - x$

The key insight is that learning a residual mapping $F(x)$ is easier than learning the full mapping $H(x)$. If the optimal function is close to identity, it's easier to push $F(x)$ toward zero than to fit an identity mapping with a stack of nonlinear layers.

Why Gradients Flow Better

Gradient computation: $\partial y / \partial x = \partial F(x) / \partial x + 1$

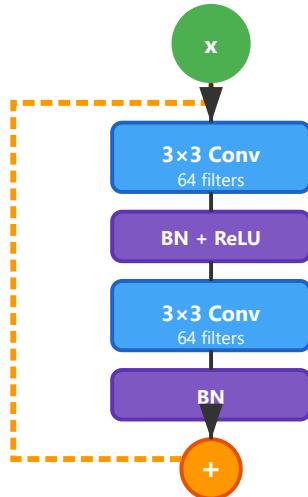
The gradient always has a component of **1** from the identity shortcut! This means even if $\partial F(x) / \partial x$ becomes very small, the gradient never completely vanishes. This "gradient highway" ensures that early layers continue to receive meaningful learning signals.

💡 Key Insight: Identity Mapping

If additional layers are not needed, the network can simply push $F(x) \rightarrow 0$, making the block an identity function ($y \approx x$). This makes it easy to add layers without risk of degradation – at worst, extra layers become identity mappings and don't hurt performance.

🏗 Residual Block Architecture

Basic Block (ResNet-18/34)



Bottleneck Block (ResNet-50+)



Why Bottleneck?

- 1×1 conv reduces dimensions
- 3×3 conv processes features
- 1×1 conv expands back

Benefits:

- ✓ Fewer parameters
- ✓ Faster computation
- ✓ Enables deeper networks



ResNet Family Variants

ResNet-18

18 layers

11.7M params

ResNet-34

34 layers

21.8M params

ResNet-50

50 layers

25.6M params

ResNet-101

101 layers

44.5M params

ResNet-152

152 layers

60.2M params

ResNeXt

Grouped convolutions

Improved efficiency

Interestingly, ResNet-50 uses bottleneck blocks and has *fewer* parameters than ResNet-34 (which uses basic blocks), yet achieves better accuracy. The bottleneck design is more parameter-efficient while allowing for greater depth.



Historical Context & Evolution

2012 - AlexNet

8 layers. Sparked the deep learning revolution by winning ImageNet with a large margin.

2014 - VGGNet

19 layers. Showed that deeper networks with smaller filters work better, but faced training difficulties.

2014 - GoogLeNet/Inception

22 layers. Introduced inception modules with multiple parallel pathways.

2015 - ResNet 🏆

152 layers! Won ImageNet 2015 with 3.57% top-5 error rate. Skip connections made very deep networks trainable.

2016+ - Descendants

DenseNet (dense connections), ResNeXt (grouped convolutions), SENet (squeeze-and-excitation), EfficientNet, and modern transformers all build upon residual connection principles.

Real-World Applications

Medical Imaging

ResNet is widely used for disease detection in X-rays, CT scans, and MRI images. Its ability to capture fine-grained features makes it ideal for identifying tumors and abnormalities.

Autonomous Vehicles

Object detection and scene understanding in self-driving cars rely on ResNet backbones for real-time processing of road conditions.

Face Recognition

Most modern facial recognition systems use ResNet-based architectures for accurate identification across varying conditions.

Transfer Learning

Pre-trained ResNet models serve as powerful feature extractors for countless downstream tasks, saving weeks of training time.



Summary: Why ResNet Changed Everything

Before ResNet, deeper networks often performed worse than shallower ones due to vanishing gradients and degradation problems. ResNet's simple yet profound insight – adding the input directly to the output – created a "gradient highway" that enables training of extremely deep networks. This breakthrough has influenced virtually every modern deep learning architecture, from computer vision to natural language processing. The principle of residual connections is now considered a fundamental building block of deep learning.

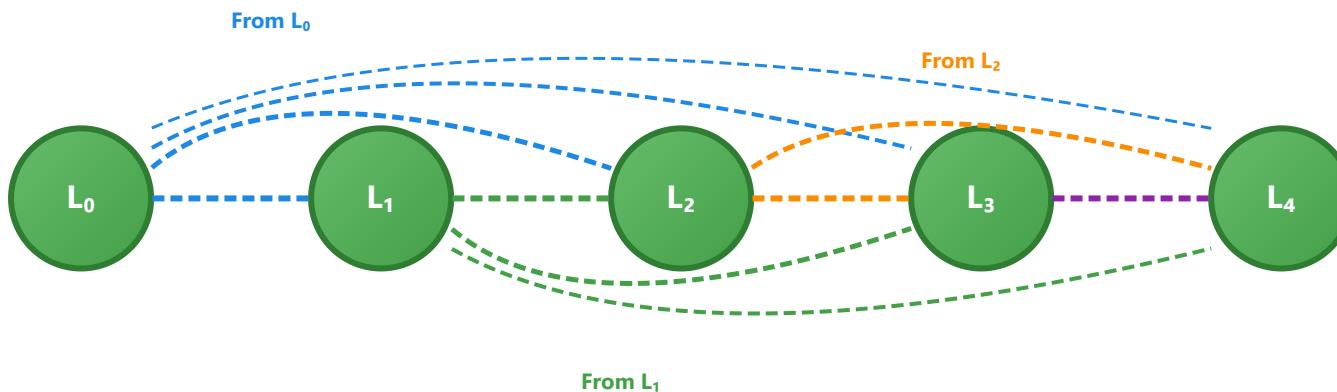
Dense Connection (DenseNet): Maximum Connectivity

CORE IDEA

Each layer connects to ALL subsequent layers

$$\mathbf{x}_l = H_l([\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_{l-1}])$$

Dense Connectivity Pattern



✓ Advantages

Maximum information flow · Feature reuse · Parameter efficiency

Efficiency

Fewer parameters than ResNet for same accuracy



Memory Consideration

Requires storing all intermediate features (higher memory)



Use Cases

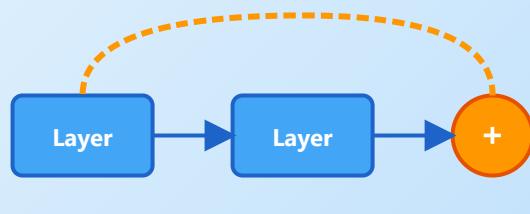
Image classification · Feature extraction · Limited parameters

GROWTH RATE (k)

Deep Dive: Understanding DenseNet

DenseNet (Densely Connected Convolutional Networks), introduced by Gao Huang et al. in 2017, takes the concept of skip connections from ResNet to its logical extreme. Instead of adding skip connections between some layers, DenseNet connects *every layer to every other layer* in a feed-forward fashion. This seemingly simple change leads to dramatic improvements in gradient flow, feature reuse, and parameter efficiency.

ResNet: Additive Skip



$$y = F(x) + x$$

Element-wise addition

DenseNet: Concatenation



$$x_l = H_l([x_0, x_1, \dots, x_{l-1}])$$

Channel-wise concatenation



Key Difference: Addition vs. Concatenation

ResNet uses element-wise addition ($y = F(x) + x$), which requires input and output to have the same dimensions.

DenseNet uses channel-wise concatenation ($[x_0, x_1, \dots]$), which preserves all features from previous layers. This means each layer has direct access to the gradients from the loss function and the original input, providing implicit deep supervision.



The Mathematics: Growth Rate & Feature Maps

Growth Rate (k)

Number of input channels to layer ℓ : $k_0 + k \times (\ell - 1)$

Where k_0 is the number of channels in the input layer, and k is the growth rate. Each layer produces k feature maps that are concatenated to all subsequent layers.

Total Number of Connections in a Dense Block

$L(L+1)/2$ connections for L layers

Example: 12-layer dense block = $12 \times 13 / 2 = 78$ connections

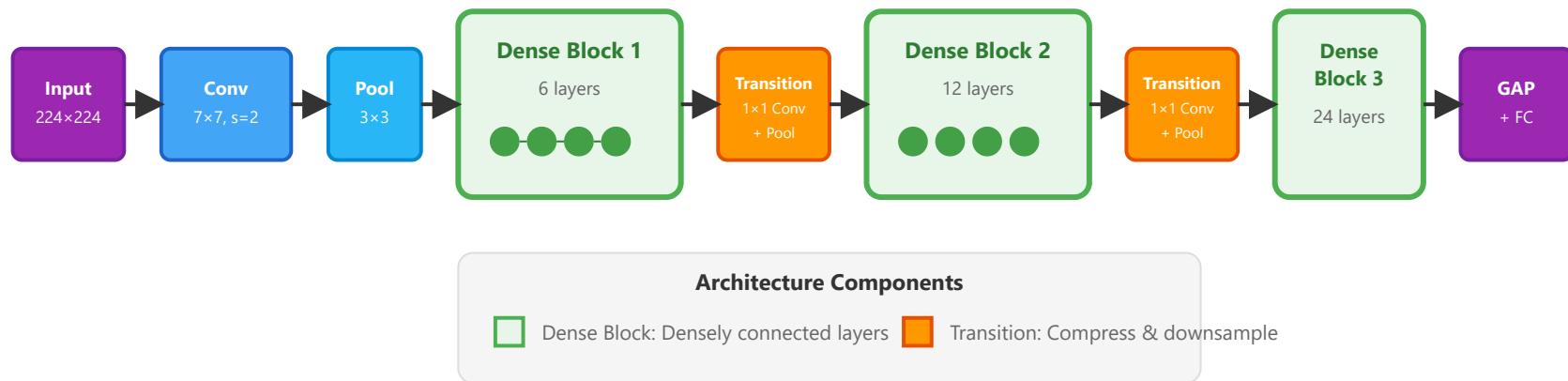
Why Small k Works

Collective knowledge: Each layer accesses ALL preceding feature maps

Unlike traditional networks where each layer only sees the previous layer's output, DenseNet layers can access the "collective knowledge" of all preceding layers. This means even a small k (like 12) is sufficient because the network doesn't need to re-learn or replicate features – it can directly reuse them.



DenseNet Architecture: Dense Blocks & Transitions



Dense Block (BN-ReLU-Conv)

Each layer: BN → ReLU → 3x3 Conv (k filters)

Uses "pre-activation" design (BN-ReLU before Conv). Each layer produces exactly k feature maps that are concatenated to all subsequent layers within the block.

Transition Layer (Compression)

Structure: BN → 1x1 Conv → 2x2 AvgPool

Reduces feature map dimensions between dense blocks. Compression factor θ (typically 0.5) reduces channels by half, making the model more compact.



DenseNet Family Variants

DenseNet-121

[6, 12, 24, 16] layers

DenseNet-169

[6, 12, 32, 32] layers

DenseNet-201

[6, 12, 48, 32] layers

DenseNet-264

[6, 12, 64, 48] layers

8.0M params

14.1M params

20.0M params

33.3M params

Parameter Efficiency Comparison

DenseNet-201 achieves similar accuracy to **ResNet-101** but with only **20M parameters** compared to ResNet's **44.5M** – less than half! This remarkable efficiency comes from feature reuse: since all previous features are directly accessible, the network doesn't need to re-learn redundant features.

DenseNet vs ResNet: When to Use Which?

Choose ResNet When:

- ✓ Memory is limited
- ✓ Need fast inference speed
- ✓ Larger batch sizes required
- ✓ Standard transfer learning tasks

vs

Choose DenseNet When:

- ✓ Parameter efficiency is critical
- ✓ Small datasets (better regularization)
- ✓ Need strong feature extraction
- ✓ Medical imaging / fine-grained tasks

Real-World Applications

Medical Image Analysis

DenseNet excels in medical imaging due to its strong feature propagation. Used extensively for chest X-ray diagnosis, retinal disease detection, and cancer identification where subtle features matter.

Remote Sensing

Satellite image classification and land-use mapping benefit from DenseNet's ability to capture multi-scale features and its parameter efficiency for edge deployment.

Semantic Segmentation

Dense connections provide rich feature representations at multiple scales, making DenseNet backbones popular in segmentation architectures like FC-DenseNet.

Mobile Applications

The parameter efficiency of DenseNet makes it suitable for resource-constrained environments where model size directly impacts deployment feasibility.

Summary: Why DenseNet Matters

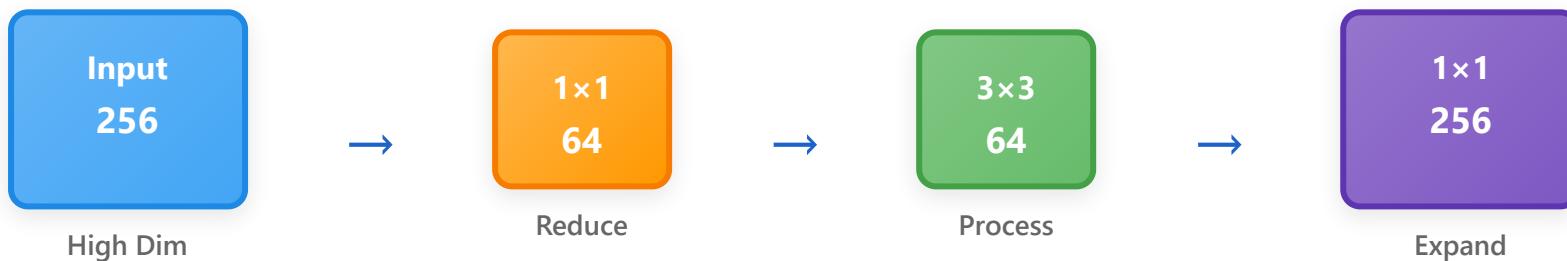
DenseNet's dense connectivity pattern offers three key advantages: **(1) Strong gradient flow** – every layer has direct access to gradients from the loss, providing implicit deep supervision; **(2) Feature reuse** – concatenation preserves all features, eliminating redundant learning; **(3) Parameter efficiency** – small growth rate k is sufficient since features are reused, not re-learned. While requiring more memory during training, DenseNet's compact models and strong feature extraction make it a powerful choice for many applications, especially when parameters are limited or fine-grained feature detection is crucial.

Bottleneck Architecture: Efficient Design

PURPOSE

Reduce computational cost while maintaining expressiveness

Bottleneck Structure



DIMENSIONAL FLOW EXAMPLE

256 channels → **64** channels → **64** channels → **256** channels



Computational Savings
~70% fewer FLOPs



Trade-off
More layers, fewer ops



Used In
ResNet, Inception, etc.

Modern Standard: Almost all efficient architectures use bottlenecks

The Role of 1×1 Convolution

PRIMARY FUNCTIONS

Dimensionality control · Channel mixing · Adding non-linearity



Channel Reduction

Projects high-dimensional features to lower dimensions



Channel Expansion

Increases feature map channels without spatial processing



Cross-Channel Mixing

Learns relationships between different channels

Example: Dimensionality Reduction

256

Input Channels

→
 1×1 Conv

64

Output Channels



Efficient

Much cheaper than 3×3 or 5×5



Computational efficiency vs larger kernels



Network-in-Network: Adds depth without overhead

Universal Use: Inception · ResNet · MobileNet · Transformers (feedforward)

Inception Module: Multi-Scale Architecture

CORE IDEA

Multiple parallel conv paths with different receptive fields

Inception Structure

Input Feature Map

1×1 reduce

1×1 reduce

1×1 reduce

Max pool

1×1 conv

3×3 conv

5×5 conv

1×1 conv

Concatenation (Filter Concat)



Multi-Scale

Captures features at different scales



Efficiency

1×1 convs reduce dimensions first



Diversity

Rich feature diversity

GoogLeNet (v1)

ImageNet 2014

Inception v2/v3

Factorized Convs

Inception-ResNet

Skip Connections

Modern Influence

EfficientNet, NAS

Depthwise Separable Convolution

KEY INSIGHT

Spatial and channel-wise processing can be separated

Standard Convolution

Single Operation
All Channels Together

Params: $D^2 \times M \times N$
(expensive)

Depthwise Separable

Depthwise Conv
Each Channel Separate

Pointwise (1×1)
Combine Channels

Params: $D^2 \times M + M \times N$
(efficient)

Efficiency Gains

Computational Savings
8-9x fewer ops

Parameter Reduction
Massive reduction



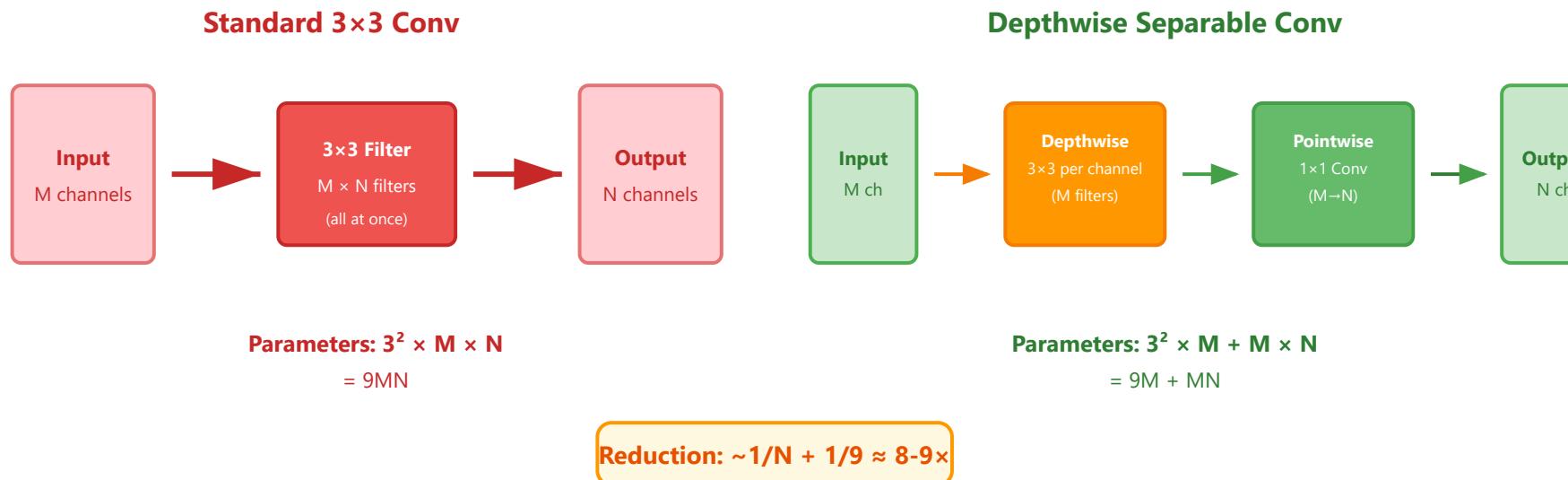
Used In
MobileNet, EfficientNet, Xception



Trade-off
Slightly lower accuracy, massive efficiency

How It Works: Step by Step

Standard convolution processes **spatial patterns** and **channel relationships** simultaneously. Depthwise Separable Convolution breaks this into two simpler operations, dramatically reducing computation while maintaining most of the representational power.



Concrete Example

3x3 Conv: 64 channels → 128 channels

Standard Conv parameters:

$$3 \times 3 \times 64 \times 128 = 73,728$$

Depthwise (3×3 per channel):

$$3 \times 3 \times 64 = 576$$

Pointwise (1×1 to expand):

$$1 \times 1 \times 64 \times 128 = 8,192$$

Depthwise Separable Total:

$$576 + 8,192 = 8,768$$

Reduction ratio:

$$73,728 / 8,768 = 8.4\times \text{ smaller!}$$

Why Does This Work?

The key insight is that **spatial correlations** (edges, textures) and **cross-channel correlations** (combining color/feature channels) are largely independent. By processing them separately, we lose minimal representational power while gaining massive computational efficiency. This makes deep learning practical on mobile devices and edge computing.

Neural Architecture Search (NAS)

GOAL

Automatically discover optimal network architectures using ML

NAS Process



Search Space



Search Strategy



Optimal Architecture

Search Space

Operations · Connections · Layer counts · Kernel sizes · Channel counts

Search Strategies

Reinforcement learning · Evolutionary algorithms · Gradient-based

Success Stories

NASNet · EfficientNet · AmoebaNet - State-of-the-art results

Computational Cost

Originally thousands of GPU-hours · Now more efficient methods



Future Direction: Architecture search becoming standard practice for deployment optimization

Model Compression Techniques



Pruning

Remove unimportant weights/neurons

Structured or unstructured



Quantization

Reduce precision (FP32 → INT8)

~4x memory/speed improvement



Knowledge Distillation

Train small student to mimic large teacher

Transfer knowledge efficiently



Low-Rank Factorization

Decompose weight matrices into smaller components

Reduce parameter count



Weight Sharing

Group weights to reduce unique parameters

Efficient representation

Typical Compression Results

Compression Ratio

5-10x

Accuracy Loss

<1%



Mobile Deployment

Essential for phones, IoT devices



Tools

TensorFlow Lite · PyTorch Mobile · ONNX Runtime

Practical Design Guidelines

🚀 START SIMPLE

Begin with proven architectures (ResNet, EfficientNet), modify as needed

Architecture Choices



- ▶ Depth vs. width: Deeper generally better
- ▶ Residual connections for >10 layers

Key Components



- ▶ Batch norm: After conv, before activation
- ▶ Activation: ReLU/GELU (avoid sigmoid/tanh)

Regularization Stack



- ▶ Data augmentation
- ▶ Dropout
- ▶ Weight decay

Architecture Principles



- ▶ Gradually reduce spatial size
- ▶ Increase channels through network



Validate on Real Hardware

Consider **inference speed** and **memory**, not just accuracy

Thank you

Ho-min Park

homin.park@ghent.ac.kr

powersimmani@gmail.com