

Deep Learning and Transformers: Theory, Mathematics, and Implementation

[Author Names]

2026

Contents

I	Mathematical Foundations	1
1	Linear Algebra for Deep Learning	2
1.1	Vector Spaces and Transformations	2
1.1.1	Vectors as Data Representations	2
1.1.2	Linear Transformations	3
1.1.3	Matrices as Linear Transformations	4
1.2	Matrix Operations	4
1.2.1	Matrix Multiplication	4
1.2.2	Computational Complexity of Matrix Operations	5
1.2.3	Batch Matrix Multiplication	6
1.2.4	Transpose	6
1.2.5	Hardware Context for Matrix Operations	7
1.3	Dot Products and Similarity	8
1.4	Matrix Decompositions	9
1.4.1	Eigenvalues and Eigenvectors	9
1.4.2	Singular Value Decomposition	10
1.5	Norms and Distance Metrics	12
1.6	Practical Deep Learning Examples	12
1.6.1	Embedding Layers and Memory Requirements	12
1.6.2	Complete Transformer Layer Analysis	13
1.6.3	Common Dimension Errors and Debugging	13
1.7	Exercises	14
2	Calculus and Optimization	17
2.1	Multivariable Calculus	17
2.1.1	Partial Derivatives	17
2.1.2	Gradients	18
2.1.3	The Chain Rule	18
2.1.4	Jacobian and Hessian Matrices	19
2.2	Gradient Descent	19
2.2.1	The Gradient Descent Algorithm	19
2.2.2	Stochastic Gradient Descent (SGD)	20
2.2.3	Momentum	20
2.2.4	Adam Optimizer	21
2.3	Gradient Computation Complexity	21
2.3.1	FLOPs for Gradient Computation	21
2.3.2	Memory Requirements for Activations	22
2.3.3	Automatic Differentiation: Forward vs Reverse Mode	22
2.3.4	Gradient Checkpointing	23
2.4	Backpropagation	23
2.4.1	Computational Graphs	23
2.4.2	Backpropagation Algorithm	24
2.4.3	Why Backpropagation is $O(n)$ Not $O(n^2)$	24

2.5	Optimizer Memory Requirements	25
2.5.1	Memory Comparison by Optimizer	25
2.5.2	Impact on GPU Memory Budget	26
2.6	Learning Rate Schedules	27
2.6.1	Learning Rate Impact on Convergence and GPU Utilization	27
2.6.2	Learning Rate Scaling with Batch Size	27
2.6.3	Practical Learning Rates for Transformers	27
2.6.4	Common Schedules	28
2.7	Hardware Considerations for Gradient Computation	28
2.7.1	Gradient Computation on GPUs	28
2.7.2	Mixed Precision Training	29
2.7.3	Gradient Accumulation	29
2.7.4	Distributed Gradient Synchronization	30
2.8	Exercises	31
3	Probability and Information Theory	33
3.1	Probability Fundamentals	33
3.1.1	Random Variables and Distributions	33
3.1.2	Conditional Probability and Bayes' Theorem	34
3.2	Information Theory	34
3.2.1	Entropy	34
3.2.2	Cross-Entropy	35
3.2.3	Kullback-Leibler Divergence	35
3.3	Cross-Entropy Loss: Computational and Memory Analysis	36
3.3.1	Memory Requirements for Logits	36
3.3.2	Computational Cost of Softmax	36
3.3.3	Optimizations for Large Vocabularies	37
3.4	KL Divergence in Practice	37
3.4.1	Applications in Modern Deep Learning	38
3.4.2	Numerical Stability Considerations	38
3.5	Hardware Implications of Softmax and Large Vocabularies	39
3.5.1	Softmax Computation on GPUs	39
3.5.2	Memory Bandwidth and Large Vocabularies	40
3.5.3	Why Vocabulary Size Impacts Training Speed	40
3.5.4	Optimization Techniques	41
3.6	Exercises	41
II	Neural Network Fundamentals	43
4	Feed-Forward Neural Networks	44
4.1	From Linear Models to Neural Networks	44
4.1.1	The Perceptron	44
4.1.2	Multi-Class Classification: Softmax Regression	44
4.2	Multi-Layer Perceptrons	45
4.2.1	Why Depth Matters	45
4.3	Memory and Computation Analysis	45
4.3.1	Parameter Count vs FLOPs	45
4.3.2	Memory Requirements for Activations	46
4.3.3	GPU Utilization for Different Layer Sizes	46
4.3.4	Batch Size Impact on Efficiency	47
4.3.5	Transformer Feed-Forward Networks	48
4.4	Activation Functions	48
4.4.1	Computational Cost of Activation Functions	49

4.4.2	Hardware Support and Fused Kernels	50
4.4.3	Why GELU is Preferred in Transformers	50
4.5	Universal Approximation Theorem	50
4.6	Weight Initialization	51
4.6.1	Variance Preservation Through Layers	51
4.6.2	Impact on Training Speed	52
4.6.3	GPU Memory During Initialization	52
4.6.4	Example: BERT-base Initialization	52
4.7	Regularization	53
4.7.1	L2 Regularization	53
4.7.2	Dropout	53
4.7.3	Computational Overhead of Dropout	54
4.7.4	Memory Requirements for Dropout	54
4.7.5	Inference Mode Differences	54
4.7.6	Dropout in Transformer Models	55
4.8	Exercises	55
5	Convolutional Neural Networks	57
5.1	Convolution Operation	57
5.1.1	Output Dimensions	57
5.2	Multi-Channel Convolutions	58
5.3	Computational Analysis of Convolutions	58
5.3.1	FLOPs for Convolution Operations	58
5.3.2	Memory Requirements for Feature Maps	59
5.3.3	GPU Optimization: im2col and Winograd	59
5.3.4	Comparison with Transformer Attention	60
5.4	Parameter Efficiency: CNNs vs Transformers	61
5.4.1	Why CNNs are More Parameter-Efficient for Images	61
5.4.2	Vision Transformer Comparison	62
5.4.3	Hybrid Architectures	63
5.5	Hardware Optimization for Convolutions	64
5.5.1	Tensor Core Utilization for Convolutions	64
5.5.2	cuDNN Optimizations	65
5.5.3	Memory Bandwidth vs Compute	65
5.5.4	Batch Size Impact on Convolution Performance	66
5.6	Pooling Layers	66
5.7	Classic Architectures	67
5.7.1	VGG-16 (2014)	67
5.7.2	ResNet (2015)	67
5.8	Batch Normalization	67
5.9	Exercises	67
6	Recurrent Neural Networks	69
6.1	Vanilla RNNs	69
6.1.1	Backpropagation Through Time (BPTT)	70
6.1.2	Vanishing and Exploding Gradients	70
6.1.3	Quantitative Analysis of Gradient Decay	71
6.2	Long Short-Term Memory (LSTM)	71
6.2.1	LSTM Computational Analysis	72
6.3	Gated Recurrent Unit (GRU)	73
6.4	Bidirectional RNNs	73
6.5	RNN Applications	74
6.6	RNNs vs Transformers: A Computational Comparison	74
6.6.1	Sequential vs Parallel Computation	74

6.6.2	Memory Complexity: $O(nd)$ vs $O(n^2)$	75
6.6.3	Training Time Comparison	75
6.6.4	Hardware Limitations of RNNs	76
6.6.5	Why Transformers Replaced RNNs	76
6.7	Exercises	77

III Attention Mechanisms 79

7 Attention Mechanisms: Fundamentals 80

7.1	Motivation: The Seq2Seq Bottleneck	80
7.1.1	RNN Encoder-Decoder Architecture	80
7.1.2	Attention Solution	81
7.2	Additive Attention (Bahdanau)	82
7.3	Scaled Dot-Product Attention	85
7.3.1	Why Scaling Matters: Variance Analysis	85
7.3.2	Computational Efficiency	86
7.4	Attention Score Computation Methods	89
7.4.1	Taxonomy of Attention Mechanisms	89
7.4.2	Comparative Analysis	90
7.5	Query-Key-Value Paradigm	90
7.5.1	Intuition	90
7.5.2	Projecting to QKV	91
7.6	Hardware Implications of Attention	93
7.6.1	Parallelization: RNNs vs Attention	93
7.6.2	Memory Bandwidth vs Compute	93
7.6.3	Batch Processing Efficiency	94
7.7	Attention Variants	94
7.7.1	Self-Attention vs Cross-Attention	94
7.7.2	Masked Attention	95
7.8	Exercises	95

8 Self-Attention and Multi-Head Attention 96

8.1	Self-Attention Mechanism	96
8.1.1	Hardware Considerations and Memory Layout	97
8.2	Multi-Head Attention	98
8.2.1	Parallel Computation and Memory Layout	99
8.2.2	Tensor Core Utilization	100
8.3	Positional Encoding	100
8.3.1	Positional Encoding Variants	101
8.4	Computational Complexity	102
8.4.1	Memory Complexity Analysis	102
8.4.2	Time Complexity Breakdown	103
8.4.3	Scaling Experiments	103
8.5	Causal (Masked) Self-Attention	104
8.5.1	Efficient Causal Mask Implementation	104
8.6	Attention Patterns and Interpretability	105
8.7	Hardware-Specific Optimizations	106
8.7.1	Flash Attention	106
8.7.2	Fused Kernels	106
8.7.3	Tensor Core Optimization	107
8.8	Memory-Efficient Attention Variants	107
8.8.1	Sparse Attention	107
8.8.2	Linear Attention	108

8.8.3	Low-Rank Attention	108
8.9	Exercises	109
9	Attention Variants and Mechanisms	110
9.1	Cross-Attention	110
9.1.1	Transformer Decoder Attention Layers	111
9.2	Relative Position Representations	111
9.2.1	Shaw et al. Relative Attention	111
9.2.2	T5 Relative Position Bias	112
9.3	Alternative Attention Scoring Functions	112
9.3.1	Additive (Bahdanau)	112
9.3.2	Multiplicative (Luong)	112
9.3.3	Scaled Dot-Product (Transformers)	112
9.3.4	General	112
9.4	Attention Masking	112
9.4.1	Padding Mask	112
9.4.2	Combined Masks	113
9.5	Attention Dropout	113
9.6	Layer Normalization with Attention	113
9.6.1	Post-Norm (Original Transformer)	113
9.6.2	Pre-Norm (More Common Now)	113
9.7	Visualizing Attention	113
9.7.1	Attention Heatmaps	114
9.7.2	Interpreting Multiple Heads	114
9.8	Practical Implementation Considerations	114
9.8.1	Memory-Efficient Attention	114
9.8.2	Fused Attention Kernels	114
9.9	Efficient Attention Variants	115
9.9.1	Local Attention	115
9.9.2	Sparse Attention	115
9.9.3	Linear Attention	116
9.9.4	Low-Rank Attention	117
9.9.5	Comprehensive Complexity Comparison	117
9.9.6	Implementation Considerations	118
9.10	Exercises	119
IV	Transformer Architecture	121
10	The Transformer Model	122
10.1	Transformer Architecture Overview	122
10.1.1	High-Level Structure	122
10.2	Transformer Encoder	123
10.2.1	Single Encoder Layer	123
10.2.2	Complete Encoder Stack	126
10.3	Position-wise Feed-Forward Networks	127
10.4	Transformer Decoder	129
10.4.1	Single Decoder Layer	129
10.4.2	Complete Decoder Stack	131
10.5	Computational Complexity and Hardware Analysis	133
10.5.1	FLOPs Analysis	133
10.5.2	Memory Bandwidth Considerations	134
10.5.3	Scaling to Large Models	135
10.6	Complete Transformer Architecture	136

10.6.1	Full Encoder-Decoder Model	136
10.6.2	Original Transformer Configuration	136
10.7	Residual Connections and Layer Normalization	137
10.7.1	Residual Connections	137
10.7.2	Layer Normalization	137
10.7.3	Pre-Norm vs Post-Norm	138
10.8	Training Objectives	139
10.8.1	Sequence-to-Sequence Training	139
10.8.2	Autoregressive Generation	140
10.9	Transformer Variants: Architectural Patterns	140
10.9.1	Encoder-Only Architecture (BERT)	140
10.9.2	Decoder-Only Architecture (GPT)	140
10.9.3	Encoder-Decoder Architecture (T5, BART)	141
10.9.4	Choosing the Right Architecture	142
10.10	Exercises	143
11	Training Transformers	144
11.1	Training Objectives and Loss Functions	144
11.1.1	Masked Language Modeling	145
11.1.2	Causal Language Modeling	145
11.1.3	Sequence-to-Sequence Training	146
11.2	Backpropagation Through Transformers	146
11.2.1	Gradient Flow Analysis	147
11.2.2	Gradients Through Residual Connections	147
11.2.3	Gradients Through Layer Normalization	147
11.2.4	Gradients Through Attention	148
11.2.5	Gradients Through Feed-Forward Networks	149
11.2.6	Computational Cost of Backpropagation	149
11.3	Optimization Algorithms	150
11.3.1	Adam Optimizer	150
11.3.2	AdamW: Decoupled Weight Decay	151
11.3.3	LAMB: Large Batch Training	151
11.3.4	Optimizer Memory Comparison	152
11.4	Learning Rate Schedules	152
11.4.1	The Necessity of Warmup	152
11.4.2	Warmup Plus Linear Decay	153
11.4.3	Inverse Square Root Decay	153
11.4.4	Cosine Annealing	154
11.5	Mixed Precision Training	154
11.5.1	FP16 Training Algorithm	154
11.5.2	Memory Savings	155
11.5.3	Hardware Acceleration	155
11.5.4	BF16: An Alternative to FP16	156
11.6	Gradient Accumulation	156
11.6.1	Algorithm and Implementation	156
11.6.2	Trade-offs and Considerations	157
11.6.3	Practical Example	157
11.7	Gradient Checkpointing	158
11.7.1	The Memory-Computation Trade-off	158
11.7.2	Implementation Strategies	159
11.7.3	Practical Impact	159
11.7.4	When to Use Gradient Checkpointing	160
11.8	Distributed Training Strategies	160
11.8.1	Data Parallelism	160

11.8.2	Model Parallelism	161
11.8.3	Pipeline Parallelism	161
11.8.4	Tensor Parallelism	162
11.8.5	ZeRO: Zero Redundancy Optimizer	162
11.8.6	Comparison of Strategies	163
11.9	Batch Size and Sequence Length Selection	163
11.9.1	Batch Size Considerations	163
11.9.2	Memory Scaling with Batch Size	164
11.9.3	Sequence Length Considerations	164
11.9.4	Dynamic Batching	165
11.9.5	Practical Guidelines	166
11.10	Regularization Techniques	166
11.10.1	Dropout	167
11.10.2	Weight Decay	167
11.10.3	Label Smoothing	168
11.10.4	Gradient Clipping	168
11.11	Training Time and Cost Estimates	169
11.11.1	BERT-base Training	169
11.11.2	GPT-2 Training	170
11.11.3	GPT-3 Training	170
11.11.4	Scaling Laws	171
11.11.5	Cost-Performance Trade-offs	172
11.12	Practical Training Recipe	172
11.12.1	Data Preparation	172
11.12.2	Model Initialization	172
11.12.3	Hyperparameter Selection	173
11.12.4	Training Loop	173
11.12.5	Monitoring and Debugging	174
11.13	Exercises	174
12	Computational Analysis of Transformers	176
12.1	Computational Complexity	176
12.1.1	Self-Attention Complexity	176
12.1.2	Feed-Forward Network Complexity	178
12.1.3	Per-Layer Total Complexity	179
12.1.4	Complexity Analysis	180
12.2	Memory Requirements	181
12.2.1	Model Parameters	181
12.2.2	Activation Memory	182
12.2.3	Training Memory Budget	183
12.2.4	Hardware Selection Guide	185
12.3	Inference Optimization	186
12.3.1	KV Caching for Autoregressive Decoding	186
12.3.2	Batched Inference	187
12.3.3	Quantization for Inference	188
12.3.4	Model Distillation	188
12.3.5	Inference Optimization Summary	188
12.4	Scaling Laws	189
12.4.1	Kaplan et al. Scaling Laws	189
12.4.2	Chinchilla Scaling Laws	189
12.5	Exercises	189

V	Modern Transformer Variants	191
13	BERT: Bidirectional Encoder Representations	192
13.1	BERT Architecture	192
13.1.1	Model Specification	192
13.1.2	Parameter Breakdown and Memory Requirements	193
13.1.3	Input Representation	194
13.2	Pre-Training Objectives	194
13.2.1	Masked Language Modeling (MLM)	194
13.2.2	Next Sentence Prediction (NSP)	196
13.3	Training Details and Computational Cost	197
13.3.1	Hardware and Training Configuration	197
13.3.2	Computational Cost Analysis	197
13.4	Fine-Tuning BERT	198
13.4.1	Fine-Tuning Procedure and Memory Requirements	198
13.4.2	Classification Tasks	198
13.4.3	Token-Level Tasks	198
13.4.4	Question Answering (SQuAD)	199
13.5	BERT Variants	199
13.5.1	RoBERTa (Robustly Optimized BERT)	199
13.5.2	ALBERT (A Lite BERT)	200
13.5.3	DistilBERT: Knowledge Distillation for Compression	201
13.5.4	Memory and Speed Comparisons	201
13.6	Hardware Requirements and Deployment	203
13.6.1	GPU Memory Requirements	203
13.6.2	Batch Size Limits by GPU Type	203
13.6.3	Inference Speed Analysis	204
13.7	Analysis and Interpretability	205
13.7.1	What BERT Learns	205
13.7.2	Probing Tasks	205
13.8	Exercises	205
14	GPT: Generative Pre-Training	207
14.1	GPT Architecture	207
14.1.1	Decoder-Only Transformers	207
14.1.2	GPT Model Sizes	208
14.2	Pre-Training: Autoregressive Language Modeling	210
14.2.1	Training Objective	210
14.2.2	Pre-Training Data	211
14.2.3	Training Infrastructure and Costs	212
14.3	In-Context Learning and Few-Shot Prompting	214
14.3.1	Autoregressive Generation with KV Caching	214
14.3.2	Zero-Shot, One-Shot, Few-Shot	215
14.3.3	Zero-Shot, One-Shot, Few-Shot	215
14.3.4	Emergent Abilities	216
14.4	Scaling Laws	217
14.4.1	Parameter Scaling	217
14.4.2	Compute-Optimal Training	218
14.4.3	Hardware Requirements for Inference	219
14.5	Instruction Tuning and RLHF	221
14.5.1	Instruction Tuning	221
14.5.2	RLHF (Reinforcement Learning from Human Feedback)	221
14.6	GPT Capabilities and Limitations	221
14.6.1	Capabilities	221

14.6.2	Limitations	222
14.7	Exercises	222
15	T5 and BART: Encoder-Decoder Architectures	223
15.1	T5: Text-to-Text Transfer Transformer	223
15.1.1	Unified Text-to-Text Framework	223
15.1.2	T5 Architecture	224
15.1.3	Pre-Training Objective: Span Corruption	226
15.1.4	T5 Model Sizes and Scaling	227
15.1.5	T5 Training Details	228
15.2	BART: Denoising Autoencoder	230
15.2.1	BART Architecture and Design Philosophy	230
15.2.2	BART Parameter Breakdown and Memory Requirements	231
15.2.3	Denoising Objectives and Corruption Strategies	232
15.2.4	BART Training Details	234
15.3	Encoder-Decoder Efficiency Analysis	234
15.3.1	Computational Cost of Cross-Attention	234
15.3.2	Comparison: Encoder-Decoder vs Decoder-Only	235
15.4	Comparing T5 and BART	237
15.5	Prefix Language Models	237
15.5.1	Prefix LM Objective	237
15.6	Applications and Fine-tuning	238
15.6.1	Summarization	238
15.6.2	Translation	238
15.6.3	Question Answering	238
15.7	Mixture of Denoisers (UL2)	238
15.8	Exercises	239
16	Efficient Transformers	240
16.1	The Quadratic Bottleneck	240
16.1.1	Complexity Analysis	240
16.2	Sparse Attention Patterns	241
16.2.1	Efficiency Taxonomy	241
16.2.2	Fixed Sparse Patterns	241
16.2.3	BigBird: Random + Window + Global	243
16.3	Linear Attention Methods	243
16.3.1	Linformer	244
16.3.2	Performer (Kernel-based)	245
16.4	Memory-Efficient Attention	246
16.4.1	Flash Attention	246
16.4.2	Memory-Efficient Transformers	247
16.5	Comparison of Efficient Methods	248
16.5.1	Comprehensive Benchmarks	248
16.5.2	Memory Scaling Analysis	248
16.5.3	Speed Benchmarks on A100 GPU	248
16.5.4	Quality Trade-offs	249
16.5.5	When to Use Each Method	249
16.6	Long-Context Models	250
16.6.1	Longformer	250
16.6.2	Reformer	250
16.7	Exercises	251

VI	Advanced Topics	253
17	Vision Transformers	254
17.1	From Images to Sequences	254
17.1.1	The Patch Embedding Approach	254
17.1.2	Position Encodings for 2D	255
17.2	Vision Transformer (ViT) Architecture	255
17.2.1	Complete ViT Model	255
17.2.2	ViT Model Variants	256
17.2.3	Memory Requirements and Computational Analysis	256
17.3	Training Vision Transformers	257
17.3.1	Pre-training Strategies	257
17.3.2	Data Augmentation and Regularization	258
17.3.3	DeiT: Data-efficient Image Transformers	258
17.4	Masked Autoencoders (MAE)	258
17.4.1	Self-Supervised Pre-training for Vision	258
17.5	Hierarchical Vision Transformers	259
17.5.1	Motivation for Hierarchical Architectures	259
17.5.2	Swin Transformer	259
17.5.3	Pyramid Vision Transformer (PVT)	261
17.5.4	Hybrid Architectures: CoAtNet	261
17.6	ViT vs CNN Comparison	262
17.6.1	Parameter Efficiency	262
17.6.2	Computational Complexity Analysis	262
17.6.3	Data Requirements and Inductive Bias	263
17.6.4	When to Use Each Architecture	264
17.7	Exercises	264
18	Multimodal Transformers	266
18.1	Multimodal Learning Fundamentals	266
18.1.1	Fusion Strategies	266
18.1.2	Alignment Objectives	267
18.2	CLIP: Contrastive Language-Image Pre-training	267
18.2.1	CLIP Architecture	267
18.2.2	Computational Analysis of CLIP Training	268
18.2.3	Zero-Shot Classification with CLIP	269
18.2.4	CLIP Variants and Training Requirements	270
18.3	DALL-E and Stable Diffusion	270
18.3.1	DALL-E: Text-to-Image Generation	270
18.3.2	Stable Diffusion	271
18.4	Vision-Language Understanding	271
18.4.1	BLIP: Bootstrapped Language-Image Pre-training	271
18.4.2	Flamingo: Few-Shot Learning	272
18.5	Computational Analysis of Multimodal Transformers	273
18.5.1	Image Encoding Cost	273
18.5.2	Text Encoding Cost	273
18.5.3	Cross-Modal Attention Cost	274
18.5.4	Total Memory Requirements	274
18.6	Training Challenges for Multimodal Transformers	274
18.6.1	Batch Size Requirements for Contrastive Learning	274
18.6.2	Memory Optimization Techniques	275
18.6.3	Distributed Training Infrastructure	276
18.6.4	Training Time and Cost Estimates	276
18.7	Audio Transformers	277

18.7.1	Whisper: Speech Recognition	277
18.7.2	Audio-Text Pre-training	277
18.8	Unified Multimodal Models	278
18.8.1	Perceiver and Perceiver IO	278
18.8.2	GPT-4V and LLaVA	278
18.9	Exercises	278
19	Long Context Transformers	280
19.1	Context Length Limitations	280
19.1.1	The Quadratic Memory Bottleneck	280
19.2	Position Encoding for Long Context	281
19.2.1	The Extrapolation Challenge	281
19.2.2	Position Interpolation	282
19.2.3	Rotary Position Embedding (RoPE)	283
19.2.4	ALiBi: Attention with Linear Biases	284
19.3	Efficient Attention for Long Context	285
19.3.1	Sparse Attention Patterns	285
19.3.2	Longformer: Local and Global Attention	285
19.3.3	BigBird: Random, Window, and Global Attention	286
19.3.4	Comparison of Sparse Attention Methods	287
19.4	Recurrent Transformers	287
19.4.1	Transformer-XL	287
19.5	Retrieval-Augmented Generation	288
19.5.1	RAG Architecture	288
19.5.2	RETRO: Retrieval-Enhanced Transformer	289
19.6	Memory-Augmented Transformers	289
19.6.1	Compressive Transformer	289
19.6.2	Memorizing Transformers	289
19.7	Long Context Models in Practice	290
19.7.1	LongT5: Efficient Encoder-Decoder	290
19.7.2	GPT-4 and Claude: Production Long Context Systems	290
19.7.3	Practical Considerations for Long Context	291
19.8	Comparison and Trade-offs	291
19.8.1	Method Comparison	291
19.8.2	Hardware and Memory Considerations	292
19.8.3	Recommendations by Use Case	293
19.9	Exercises	294
20	Pre-training Strategies and Transfer Learning	297
20.1	Pre-training Objectives	297
20.1.1	Language Modeling Objectives	297
20.1.2	Denoising Objectives	298
20.1.3	Computational Cost Comparison	299
20.1.4	Contrastive Objectives	300
20.2	Data Curation and Processing	300
20.2.1	Data Scale and Requirements	300
20.2.2	Data Quality versus Quantity	301
20.2.3	Data Filtering and Cleaning	301
20.2.4	Data Deduplication	302
20.3	Training Compute Requirements	302
20.3.1	FLOPs Analysis	302
20.3.2	GPU-Hours and Cost Estimates	303
20.3.3	Scaling Laws	304
20.4	Curriculum Learning	304

20.4.1	Progressive Training Strategies	304
20.4.2	Progressive Training	305
20.4.3	Domain-Adaptive Pre-training	306
20.5	Hardware Requirements and Infrastructure	306
20.5.1	BERT-base Training Infrastructure	306
20.5.2	GPT-3 Training Infrastructure	307
20.5.3	LLaMA-65B Training Infrastructure	307
20.6	Efficient Pre-training Techniques	308
20.6.1	Mixed Precision Training	308
20.6.2	Gradient Checkpointing	308
20.6.3	ZeRO Optimizer	309
20.6.4	Pipeline Parallelism	309
20.7	Parameter-Efficient Fine-tuning	309
20.7.1	Motivation	309
20.7.2	LoRA: Low-Rank Adaptation	310
20.7.3	Adapter Layers	310
20.7.4	Prompt Tuning	310
20.8	Multi-Task and Multi-Stage Pre-training	311
20.8.1	Multi-Task Pre-training	311
20.8.2	Multi-Stage Pre-training	311
20.9	Transfer Learning Analysis	312
20.9.1	Measuring Transfer	312
20.9.2	What Makes Good Pre-training?	312
20.10	Exercises	313

VII Practical Implementation 314

21 Implementing Transformers in PyTorch 315

21.1	Multi-Head Attention Implementation	315
21.1.1	Core Components	315
21.1.2	Memory-Efficient Attention	316
21.1.3	Dimension Tracking Example	317
21.2	Position Encodings	317
21.2.1	Sinusoidal Encoding	317
21.2.2	Learned Positional Embeddings	317
21.3	Masking Strategies	317
21.3.1	Causal Mask for GPT	317
21.3.2	Padding Mask	317
21.4	Training Optimizations	318
21.4.1	Fused Kernels for Layer Normalization	318
21.4.2	Mixed Precision Training	318
21.4.3	Gradient Accumulation	319
21.4.4	Gradient Checkpointing	320
21.5	Model Initialization	321
21.5.1	Best Practices	321
21.6	Memory Profiling and Optimization	321
21.6.1	Understanding Memory Usage	321
21.6.2	Identifying Memory Bottlenecks	322
21.6.3	Optimization Strategies	322
21.6.4	Case Study: Optimizing BERT-base	323
21.7	Debugging Transformers	324
21.7.1	Common Issues	324
21.7.2	Validation Checks	324

21.8	Inference Optimization	325
21.8.1	TorchScript Compilation	325
21.8.2	KV Cache for Autoregressive Generation	325
21.8.3	ONNX Export	327
21.8.4	TensorRT Optimization	328
21.8.5	Quantization	329
21.8.6	Inference Benchmarking	331
21.9	Complete Training Pipeline	332
21.9.1	Training Script Structure	332
21.10	Complete Implementation Examples	333
21.10.1	Full Transformer from Scratch	333
21.10.2	Optimized Training Loop	337
21.10.3	Comprehensive Benchmarks	339
21.11	Distributed Training	340
21.11.1	Data Parallel	340
21.11.2	Distributed Data Parallel (DDP)	340
21.12	Performance Optimization	341
21.12.1	DataLoader Optimization	341
21.12.2	Asynchronous Data Transfer	341
21.12.3	Profiling with torch.profiler	341
21.12.4	Batch Size Tuning	342
21.12.5	Compilation with torch.compile	343
21.13	Distributed Training Implementation	344
21.13.1	Understanding Distributed Strategies	344
21.13.2	DistributedDataParallel Setup	344
21.13.3	Gradient Synchronization	345
21.13.4	Scaling Efficiency	346
21.14	Exercises	346
22	Hardware Optimization and Deployment	349
22.1	Hardware Architectures	349
22.1.1	GPU Architecture for Deep Learning	349
22.1.2	Computational Intensity	350
22.1.3	Tensor Core Optimization	350
22.1.4	TPU Architecture	351
22.2	Memory Optimization Techniques	351
22.2.1	Coalesced Memory Access	352
22.2.2	Shared Memory and Bank Conflicts	352
22.2.3	Memory Bandwidth Utilization	352
22.3	Kernel Fusion and Operation Optimization	353
22.3.1	Fusion Opportunities in Transformers	353
22.3.2	Flash Attention: Fused Attention Implementation	354
22.3.3	Implementing Fused Kernels	354
22.4	Model Quantization	354
22.4.1	Quantization Fundamentals	354
22.4.2	Post-Training Quantization (PTQ)	355
22.4.3	Quantization-Aware Training (QAT)	355
22.5	Model Pruning	356
22.5.1	Pruning Strategies	356
22.5.2	Iterative Pruning	356
22.6	Knowledge Distillation	357
22.6.1	Distillation Loss	357
22.6.2	DistilBERT Approach	357
22.7	Multi-GPU Training and Optimization	357

22.7.1	Interconnect Technologies	357
22.7.2	Data Parallelism and Gradient Synchronization	358
22.7.3	Pipeline and Tensor Parallelism	358
22.7.4	Overlapping Communication and Computation	359
22.8	Inference Optimization	359
22.8.1	ONNX Runtime	359
22.8.2	TensorRT	359
22.8.3	Batching Strategies	359
22.9	Production Deployment	360
22.9.1	Serving Frameworks	360
22.9.2	Deployment Checklist	360
22.10	Hardware Selection and Cost Analysis	361
22.10.1	CPU vs GPU Trade-offs	361
22.10.2	Training Hardware Selection	362
22.10.3	Inference Hardware Selection	362
22.10.4	Hardware Selection Decision Tree	363
22.11	Exercises	363
23	Best Practices and Production Case Studies	367
23.1	Model Selection	367
23.1.1	Architecture Selection Framework	367
23.1.2	Model Size Selection	368
23.1.3	Pre-trained versus Training from Scratch	368
23.1.4	Cost-Benefit Analysis	368
23.2	Training Best Practices	369
23.2.1	Learning Rate Selection and Tuning	369
23.2.2	Batch Size Selection	369
23.2.3	Checkpointing and Monitoring Strategy	369
23.3	Memory Management	370
23.3.1	Estimating Memory Requirements	370
23.3.2	Choosing Batch Size and Sequence Length	370
23.3.3	Gradient Checkpointing	371
23.3.4	Mixed Precision Training	371
23.3.5	Multi-GPU Strategies	371
23.4	Debugging Transformers	372
23.4.1	Systematic Debugging Workflow	372
23.4.2	Gradient Analysis	373
23.4.3	Common Training Issues and Solutions	373
23.5	Inference Optimization	373
23.5.1	Batch Size for Throughput versus Latency	373
23.5.2	KV Caching for Generation	374
23.5.3	Quantization Strategies	374
23.5.4	Model Distillation	374
23.5.5	Hardware Selection	375
23.6	Cost Optimization	375
23.6.1	Training Cost Analysis	375
23.6.2	Training Time Estimation	375
23.6.3	Inference Cost Analysis	376
23.6.4	Cost Optimization Strategies	376
23.7	Production Deployment	376
23.7.1	Model Serving Frameworks	376
23.7.2	Scaling Strategies	377
23.7.3	Monitoring and Logging	377
23.7.4	A/B Testing and Deployment	377

23.8	Practical Checklists	378
23.8.1	Before Training Checklist	378
23.8.2	During Training Checklist	378
23.8.3	Before Deployment Checklist	378
23.9	Hyperparameter Tuning	379
23.9.1	Critical Hyperparameters (Ordered by Impact)	379
23.9.2	Tuning Strategy	380
23.10	Common Pitfalls and Solutions	380
23.10.1	Architecture Pitfalls	380
23.10.2	Training Pitfalls	381
23.10.3	Deployment Pitfalls	381
23.11	Case Study: BERT for Search Ranking	381
23.11.1	Problem Setup	381
23.11.2	Architecture Decisions	381
23.11.3	Training Strategy	382
23.11.4	Production Deployment	382
23.12	Case Study: GPT for Code Generation	382
23.12.1	Problem Setup	382
23.12.2	Model and Data	382
23.12.3	Training	383
23.12.4	Evaluation	383
23.13	Future Directions	383
23.13.1	Architectural Innovations	383
23.13.2	Training Innovations	384
23.14	Conclusion	384
23.14.1	Key Takeaways	384
23.14.2	Final Advice	385
23.15	Exercises	385

Part I

Mathematical Foundations

Chapter 1

Linear Algebra for Deep Learning

Chapter Overview

Linear algebra forms the mathematical foundation of deep learning. Neural networks perform sequences of linear transformations interspersed with nonlinear operations, making matrices and vectors the fundamental objects of study. This chapter develops the linear algebra concepts essential for understanding how deep learning models transform data, how information flows through neural architectures, and how we can interpret the geometric operations these models perform.

Unlike a pure mathematics course, our treatment emphasizes the specific linear algebra operations that appear repeatedly in deep learning: matrix multiplication for transforming representations, dot products for measuring similarity, and matrix decompositions for understanding structure. We pay particular attention to dimensions and shapes, as tracking how tensor dimensions transform through operations is crucial for implementing and debugging deep learning systems.

Learning Objectives

After completing this chapter, you will be able to:

1. Represent data as vectors and transformations as matrices with clear understanding of dimensions
2. Perform matrix operations and understand their geometric interpretations
3. Calculate and interpret dot products as similarity measures
4. Understand eigendecompositions and singular value decompositions and their applications
5. Apply matrix norms and use them in regularization
6. Recognize how linear algebra operations map to neural network computations

1.1 Vector Spaces and Transformations

1.1.1 Vectors as Data Representations

In deep learning, we represent data as vectors in high-dimensional spaces. A vector $\mathbf{x} \in \mathbb{R}^n$ is an ordered collection of n real numbers, which we can interpret geometrically as a point in n -dimensional space or as an arrow from the origin to that point.

Definition 1.1 (Vector). A vector $\mathbf{x} \in \mathbb{R}^n$ is an n -tuple of real numbers:

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad (1.1)$$

where each $x_i \in \mathbb{R}$ is called a component or element of the vector.

The dimension n is the number of components in the vector. We write vectors as column vectors by default.

Example 1.1 (Image as Vector). Consider a grayscale image of size 28×28 pixels, such as an image from the MNIST handwritten digit dataset. Each pixel has an intensity value between 0 (black) and 255 (white). We can represent this image as a vector $\mathbf{x} \in \mathbb{R}^{784}$ by concatenating all pixel values:

$$\mathbf{x} = \begin{bmatrix} x_{1,1} \\ x_{1,2} \\ \vdots \\ x_{28,28} \end{bmatrix} \in \mathbb{R}^{784} \quad (1.2)$$

For color images with three channels (red, green, blue), a 224×224 RGB image becomes a vector in \mathbb{R}^{150528} ($224 \times 224 \times 3 = 150,528$). The enormous dimensionality of image data motivates the need for powerful models that can find meaningful patterns in such high-dimensional spaces.

Example 1.2 (Text as Vector). In natural language processing, we represent words as vectors called *word embeddings*. A common choice is to represent each word as a vector in \mathbb{R}^{300} or \mathbb{R}^{768} . For instance, the word “king” might be represented as:

$$\mathbf{w}_{\text{king}} = \begin{bmatrix} 0.23 \\ -0.45 \\ 0.87 \\ \vdots \\ 0.12 \end{bmatrix} \in \mathbb{R}^{300} \quad (1.3)$$

These embeddings are learned such that semantically similar words have similar vector representations. The famous example is that $\mathbf{w}_{\text{king}} - \mathbf{w}_{\text{man}} + \mathbf{w}_{\text{woman}} \approx \mathbf{w}_{\text{queen}}$, suggesting that vector arithmetic can capture semantic relationships.

1.1.2 Linear Transformations

Definition 1.2 (Linear Transformation). A function $T : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a **linear transformation** if for all vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ and all scalars $a, b \in \mathbb{R}$:

$$T(a\mathbf{x} + b\mathbf{y}) = aT(\mathbf{x}) + bT(\mathbf{y}) \quad (1.4)$$

Linear transformations preserve vector space structure: they map lines to lines and preserve the origin ($T(\mathbf{0}) = \mathbf{0}$).

1.1.3 Matrices as Linear Transformations

Every linear transformation from \mathbb{R}^n to \mathbb{R}^m can be represented by an $m \times n$ matrix.

Definition 1.3 (Matrix). An $m \times n$ matrix \mathbf{A} is a rectangular array of numbers with m rows and n columns:

$$\mathbf{A} = \begin{bmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,n} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m,1} & a_{m,2} & \cdots & a_{m,n} \end{bmatrix} \in \mathbb{R}^{m \times n} \quad (1.5)$$

The notation $\mathbf{A} \in \mathbb{R}^{m \times n}$ specifies the dimensions explicitly: m rows and n columns.

Key Point 1.1. Dimension Tracking: For matrix-vector multiplication $\mathbf{Ax} = \mathbf{y}$:

$$\underbrace{\mathbf{A}}_{\mathbb{R}^{m \times n}} \underbrace{\mathbf{x}}_{\mathbb{R}^n} = \underbrace{\mathbf{y}}_{\mathbb{R}^m} \quad (1.6)$$

The inner dimensions must match (n), and the result has the outer dimensions (m).

Example 1.3 (Neural Network Layer). A single fully-connected neural network layer performs:

$$\mathbf{h} = \mathbf{W}\mathbf{x} + \mathbf{b} \quad (1.7)$$

where $\mathbf{x} \in \mathbb{R}^{n_{\text{in}}}$, $\mathbf{W} \in \mathbb{R}^{n_{\text{out}} \times n_{\text{in}}}$, $\mathbf{b} \in \mathbb{R}^{n_{\text{out}}}$, $\mathbf{h} \in \mathbb{R}^{n_{\text{out}}}$.

For transforming a 784-dimensional input to 256-dimensional hidden representation:

$$\underbrace{\mathbf{h}}_{\mathbb{R}^{256}} = \underbrace{\mathbf{W}}_{\mathbb{R}^{256 \times 784}} \underbrace{\mathbf{x}}_{\mathbb{R}^{784}} + \underbrace{\mathbf{b}}_{\mathbb{R}^{256}} \quad (1.8)$$

This layer has $256 \times 784 = 200,704$ weights plus 256 biases, totaling **200,960 trainable parameters**.

Concrete Numerical Example: With $n_{\text{in}} = 3$, $n_{\text{out}} = 2$:

$$\mathbf{W} = \begin{bmatrix} 0.5 & -0.3 & 0.8 \\ 0.2 & 0.6 & -0.4 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 0.1 \\ -0.2 \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} 1.0 \\ 2.0 \\ -0.5 \end{bmatrix} \quad (1.9)$$

Computing:

$$\mathbf{W}\mathbf{x} = \begin{bmatrix} 0.5(1.0) - 0.3(2.0) + 0.8(-0.5) \\ 0.2(1.0) + 0.6(2.0) - 0.4(-0.5) \end{bmatrix} = \begin{bmatrix} -0.5 \\ 1.6 \end{bmatrix} \quad (1.10)$$

$$\mathbf{h} = \begin{bmatrix} -0.5 \\ 1.6 \end{bmatrix} + \begin{bmatrix} 0.1 \\ -0.2 \end{bmatrix} = \begin{bmatrix} -0.4 \\ 1.4 \end{bmatrix} \quad (1.11)$$

1.2 Matrix Operations

1.2.1 Matrix Multiplication

Definition 1.4 (Matrix Multiplication). For $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{B} \in \mathbb{R}^{n \times p}$, their product $\mathbf{C} = \mathbf{AB} \in \mathbb{R}^{m \times p}$ is:

$$c_{i,k} = \sum_{j=1}^n a_{i,j} b_{j,k} \quad (1.12)$$

Example 1.4 (Matrix Multiplication Computation). Compute $\mathbf{C} = \mathbf{AB}$ where:

$$\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \in \mathbb{R}^{2 \times 2}, \quad \mathbf{B} = \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} \in \mathbb{R}^{2 \times 2} \quad (1.13)$$

Computing each entry:

$$c_{1,1} = 1(5) + 2(7) = 19 \quad (1.14)$$

$$c_{1,2} = 1(6) + 2(8) = 22 \quad (1.15)$$

$$c_{2,1} = 3(5) + 4(7) = 43 \quad (1.16)$$

$$c_{2,2} = 3(6) + 4(8) = 50 \quad (1.17)$$

Therefore: $\mathbf{C} = \begin{bmatrix} 19 & 22 \\ 43 & 50 \end{bmatrix}$

1.2.2 Computational Complexity of Matrix Operations

Understanding the computational cost of matrix operations is essential for designing efficient deep learning systems.

Theorem 1.1 (Matrix Multiplication Complexity). Computing $\mathbf{C} = \mathbf{AB}$ where $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{B} \in \mathbb{R}^{n \times p}$ requires:

$$\text{FLOPs} = 2mnp \quad (1.18)$$

floating-point operations (multiply-accumulate operations count as 2 FLOPs each).

Example 1.5 (Transformer Attention Complexity). In transformer self-attention, we compute $\mathbf{A} = \mathbf{QK}^\top$ where $\mathbf{Q}, \mathbf{K} \in \mathbb{R}^{n \times d_k}$ (sequence length n , key dimension d_k).

Dimensions: $\underbrace{\mathbf{Q}}_{\mathbb{R}^{n \times d_k}} \underbrace{\mathbf{K}^\top}_{\mathbb{R}^{d_k \times n}} = \underbrace{\mathbf{A}}_{\mathbb{R}^{n \times n}}$

Computational cost: $2n \cdot d_k \cdot n = 2n^2 d_k$ FLOPs

For GPT-3 with $n = 2048$ tokens and $d_k = 128$:

$$\text{FLOPs} = 2 \times (2048)^2 \times 128 = 1,073,741,824 \approx 1.07 \text{ GFLOPs} \quad (1.19)$$

This quadratic scaling in sequence length ($O(n^2)$) is why long-context transformers are computationally expensive.

Example 1.6 (Feed-Forward Network Cost). A transformer feed-forward network applies two linear transformations:

$$\mathbf{h} = \mathbf{W}_1 \mathbf{x} + \mathbf{b}_1 \quad \text{where } \mathbf{W}_1 \in \mathbb{R}^{d_{ff} \times d_{\text{model}}} \quad (1.20)$$

$$\mathbf{y} = \mathbf{W}_2 \mathbf{h} + \mathbf{b}_2 \quad \text{where } \mathbf{W}_2 \in \mathbb{R}^{d_{\text{model}} \times d_{ff}} \quad (1.21)$$

For a batch of B sequences of length n , input is $\mathbf{X} \in \mathbb{R}^{B \times n \times d_{\text{model}}}$.

First transformation: $2 \cdot (Bn) \cdot d_{\text{model}} \cdot d_{ff}$ FLOPs

Second transformation: $2 \cdot (Bn) \cdot d_{ff} \cdot d_{\text{model}}$ FLOPs

Total: $4Bn \cdot d_{\text{model}} \cdot d_{ff}$ FLOPs

For BERT-base ($d_{\text{model}} = 768$, $d_{ff} = 3072$, $n = 512$, $B = 32$):

$$\text{FLOPs} = 4 \times 32 \times 512 \times 768 \times 3072 = 154,618,822,656 \approx 154.6 \text{ GFLOPs} \quad (1.22)$$

1.2.3 Batch Matrix Multiplication

Modern deep learning frameworks process multiple examples simultaneously using batched operations.

Definition 1.5 (Batch Matrix Multiplication). For tensors $\mathbf{A} \in \mathbb{R}^{B \times m \times n}$ and $\mathbf{B} \in \mathbb{R}^{B \times n \times p}$, batch matrix multiplication produces $\mathbf{C} \in \mathbb{R}^{B \times m \times p}$ where:

$$\mathbf{C}[b] = \mathbf{A}[b] \mathbf{B}[b] \quad \text{for } b = 1, \dots, B \quad (1.23)$$

Example 1.7 (Multi-Head Attention Dimensions). In multi-head attention with $h = 12$ heads, batch size $B = 32$, sequence length $n = 512$, and head dimension $d_k = 64$:

Query tensor: $\mathbf{Q} \in \mathbb{R}^{B \times h \times n \times d_k} = \mathbb{R}^{32 \times 12 \times 512 \times 64}$

Key tensor: $\mathbf{K} \in \mathbb{R}^{B \times h \times n \times d_k} = \mathbb{R}^{32 \times 12 \times 512 \times 64}$

Attention scores: $\mathbf{A} = \mathbf{Q} \mathbf{K}^\top \in \mathbb{R}^{32 \times 12 \times 512 \times 512}$

This requires $B \times h \times 2n^2 d_k = 32 \times 12 \times 2 \times 512^2 \times 64 = 12,884,901,888 \approx 12.9$ GFLOPs.

Key Point 1.2. Broadcasting in PyTorch/NumPy: When dimensions don't match, broadcasting rules automatically expand dimensions by aligning them from the right, stretching size-1 dimensions to match, and adding missing dimensions as size-1. For example, adding a bias vector \mathbb{R}^{768} to a tensor $\mathbb{R}^{32 \times 512 \times 768}$ broadcasts the bias across batch and sequence dimensions, effectively treating it as $\mathbb{R}^{1 \times 1 \times 768}$ and expanding it to match the full shape.

1.2.4 Transpose

Definition 1.6 (Transpose). The **transpose** of $\mathbf{A} \in \mathbb{R}^{m \times n}$, denoted $\mathbf{A}^\top \in \mathbb{R}^{n \times m}$, swaps rows and columns:

$$[\mathbf{A}^\top]_{i,j} = a_{j,i} \quad (1.24)$$

Important properties:

$$(\mathbf{A}^\top)^\top = \mathbf{A} \quad (1.25)$$

$$(\mathbf{AB})^\top = \mathbf{B}^\top \mathbf{A}^\top \quad (1.26)$$

1.2.5 Hardware Context for Matrix Operations

Understanding how matrix operations map to hardware is crucial for writing efficient deep learning code.

Memory Layout: Row-Major vs Column-Major

Matrices are stored in memory as one-dimensional arrays, and the layout significantly affects performance. In row-major order, used by C and PyTorch, rows are stored consecutively in memory. In column-major order, used by Fortran and MATLAB, columns are stored consecutively. For a matrix $\mathbf{A} = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$, row-major storage produces the sequence $[a, b, c, d]$ while column-major storage produces $[a, c, b, d]$.

Key Point 1.3. Cache Efficiency: Accessing memory sequentially is 10-100× faster than random access due to CPU cache lines, which typically hold 64 bytes of consecutive memory. This means you should always iterate in the storage order. For row-major matrices, iterate rows in the outer loop to access consecutive memory locations, avoiding strided access patterns that jump across rows and cause cache misses.

For row-major matrices, iterate rows in the outer loop:

```
# Good: Sequential memory access
for i in range(m):
    for j in range(n):
        result += A[i, j] # Accesses consecutive memory

# Bad: Strided memory access
for j in range(n):
    for i in range(m):
        result += A[i, j] # Jumps across rows
```

GPU Acceleration and BLAS Libraries

Modern deep learning relies on highly optimized linear algebra libraries that provide standardized interfaces for common operations. The Basic Linear Algebra Subprograms (BLAS) standard defines three levels of operations: Level 1 for vector operations like dot products and norms with $O(n)$ complexity, Level 2 for matrix-vector operations like \mathbf{Ax} with $O(n^2)$ complexity, and Level 3 for matrix-matrix operations like \mathbf{AB} with $O(n^3)$ complexity. Common CPU implementations include Intel MKL, OpenBLAS, and Apple Accelerate, while GPU implementations include NVIDIA cuBLAS and AMD rocBLAS. These libraries achieve near-peak hardware performance through careful optimization of memory access patterns, instruction scheduling, and hardware-specific features.

Example 1.8 (GPU Matrix Multiplication). NVIDIA GPUs use specialized Tensor Cores for accelerated matrix multiplication, achieving dramatically higher throughput than standard CUDA cores. The A100 GPU delivers 312 TFLOPS peak performance for FP16 operations with Tensor Cores, has 1.6 TB/s memory bandwidth, and includes 40 MB of L2 cache to reduce memory access latency.

For matrix multiplication $\mathbf{C} = \mathbf{AB}$ with $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{4096 \times 4096}$, we need $2 \times 4096^3 = 137,438,953,472 \approx 137.4$ GFLOPs of computation and must transfer $3 \times 4096^2 \times 4 = 201,326,592 \approx 192$ MB of data (three matrices at 4 bytes per float). On an A100, this takes approximately $\frac{137.4 \text{ GFLOPS}}{312,000 \text{ GFLOPS}} \approx 0.44$ ms, making it compute-bound since the computation time exceeds the memory transfer time of

$$\frac{192 \text{ MB}}{1,600,000 \text{ MB/s}} \approx 0.12 \text{ ms.}$$

Compute-Bound vs Memory-Bound Operations

Definition 1.7 (Arithmetic Intensity). **Arithmetic intensity** measures the ratio of computation to memory access:

$$\text{Arithmetic Intensity} = \frac{\text{FLOPs}}{\text{Bytes Transferred}} \quad (1.27)$$

Operations with high arithmetic intensity are **compute-bound**, meaning they are limited by computational throughput, while operations with low arithmetic intensity are **memory-bound**, meaning they are limited by memory bandwidth.

Example 1.9 (Arithmetic Intensity Analysis). Element-wise operations like ReLU, which computes $\mathbf{y} = \max(0, \mathbf{x})$, perform n comparisons while transferring $2n$ elements at 4 bytes each for a total of $8n$ bytes, yielding an arithmetic intensity of only $\frac{n}{8n} = 0.125$ FLOP/byte. This makes element-wise operations memory-bound, as the GPU spends more time waiting for data than computing.

In contrast, matrix multiplication $\mathbf{C} = \mathbf{AB}$ with $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$ performs $2n^3$ FLOPs while transferring $3n^2 \times 4 = 12n^2$ bytes, yielding an arithmetic intensity of $\frac{2n^3}{12n^2} = \frac{n}{6}$ FLOP/byte. For $n = 1024$, this gives 170.7 FLOP/byte, making the operation compute-bound and well-suited for GPU acceleration. For smaller matrices with $n = 64$, the arithmetic intensity drops to 10.7 FLOP/byte, placing it in a transitional regime where both compute and memory bandwidth matter.

Key Point 1.4. Matrix Blocking for Cache Efficiency: Large matrix multiplications are broken into smaller blocks that fit in cache, computing $\mathbf{C}_{ij} = \sum_k \mathbf{A}_{ik} \mathbf{B}_{kj}$ where each block is typically 32×32 or 64×64 elements. This blocking strategy reduces cache misses from $O(n^3)$ to $O(n^3/\sqrt{M})$ where M is the cache size, dramatically improving performance by ensuring that frequently accessed data remains in fast cache memory rather than requiring slow main memory accesses.

1.3 Dot Products and Similarity

Definition 1.8 (Dot Product). For vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, the **dot product** is:

$$\mathbf{x}^\top \mathbf{y} = \sum_{i=1}^n x_i y_i \quad (1.28)$$

Theorem 1.2 (Geometric Dot Product). For non-zero vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$:

$$\mathbf{x}^\top \mathbf{y} = \|\mathbf{x}\|_2 \|\mathbf{y}\|_2 \cos(\theta) \quad (1.29)$$

where θ is the angle between vectors and $\|\mathbf{x}\|_2 = \sqrt{\mathbf{x}^\top \mathbf{x}}$ is the Euclidean norm.

[Cosine Similarity] The **cosine similarity** between two non-zero vectors is:

$$\text{sim}(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^\top \mathbf{y}}{\|\mathbf{x}\|_2 \|\mathbf{y}\|_2} = \cos(\theta) \in [-1, 1] \quad (1.30)$$

Example 1.10 (Attention Similarity Scores). In transformer attention, we compute similarity between query and key vectors using dot products:

$$\mathbf{q} = \begin{bmatrix} 0.5 \\ 0.8 \\ 0.3 \end{bmatrix}, \quad \mathbf{k}_1 = \begin{bmatrix} 0.6 \\ 0.7 \\ 0.2 \end{bmatrix}, \quad \mathbf{k}_2 = \begin{bmatrix} -0.3 \\ 0.1 \\ 0.9 \end{bmatrix} \quad (1.31)$$

Computing similarities:

$$\mathbf{q}^\top \mathbf{k}_1 = 0.5(0.6) + 0.8(0.7) + 0.3(0.2) = 0.92 \quad (1.32)$$

$$\mathbf{q}^\top \mathbf{k}_2 = 0.5(-0.3) + 0.8(0.1) + 0.3(0.9) = 0.20 \quad (1.33)$$

The query \mathbf{q} is more similar to \mathbf{k}_1 (score 0.92) than to \mathbf{k}_2 (score 0.20). These scores determine attention weights.

1.4 Matrix Decompositions

1.4.1 Eigenvalues and Eigenvectors

Definition 1.9 (Eigenvalues and Eigenvectors). For a square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$, a non-zero vector $\mathbf{v} \in \mathbb{R}^n$ is an **eigenvector** with corresponding **eigenvalue** $\lambda \in \mathbb{R}$ if:

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v} \quad (1.34)$$

Geometrically, an eigenvector is only scaled (not rotated) when \mathbf{A} is applied. The eigenvalue λ is the scaling factor.

Example 1.11 (Computing Eigenvalues). Find eigenvalues of:

$$\mathbf{A} = \begin{bmatrix} 3 & 1 \\ 1 & 3 \end{bmatrix} \quad (1.35)$$

Solving $\det(\mathbf{A} - \lambda\mathbf{I}) = 0$:

$$\det \begin{bmatrix} 3 - \lambda & 1 \\ 1 & 3 - \lambda \end{bmatrix} = (3 - \lambda)^2 - 1 = \lambda^2 - 6\lambda + 8 = 0 \quad (1.36)$$

$$= (\lambda - 4)(\lambda - 2) = 0 \quad (1.37)$$

Eigenvalues: $\lambda_1 = 4, \lambda_2 = 2$

For $\lambda_1 = 4$, eigenvector: $\mathbf{v}_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

For $\lambda_2 = 2$, eigenvector: $\mathbf{v}_2 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ -1 \end{bmatrix}$

1.4.2 Singular Value Decomposition

Theorem 1.3 (Singular Value Decomposition). *Any matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ can be decomposed as:*

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top \quad (1.38)$$

where $\mathbf{U} \in \mathbb{R}^{m \times m}$ is an orthogonal matrix of left singular vectors, $\mathbf{\Sigma} \in \mathbb{R}^{m \times n}$ is a diagonal matrix with singular values $\sigma_1 \geq \sigma_2 \geq \dots \geq 0$, and $\mathbf{V} \in \mathbb{R}^{n \times n}$ is an orthogonal matrix of right singular vectors.

Key Point 1.5. *SVD always exists for any matrix, unlike eigendecomposition which requires special conditions.*

Low-Rank Approximation and Model Compression

SVD enables efficient model compression by approximating matrices with lower-rank factorizations.

Theorem 1.4 (Eckart-Young Theorem). *The best rank- k approximation to \mathbf{A} in Frobenius norm is:*

$$\mathbf{A}_k = \sum_{i=1}^k \sigma_{ii} \mathbf{v}_i^\top \quad (1.39)$$

where σ_i are the k largest singular values with corresponding singular vectors \mathbf{v}_i .

The approximation error is:

$$\|\mathbf{A} - \mathbf{A}_k\|_F = \sqrt{\sum_{i=k+1}^r \sigma_i^2} \quad (1.40)$$

where $r = \min(m, n)$ is the rank of \mathbf{A} .

Example 1.12 (SVD for Model Compression - Detailed Analysis). Consider weight matrix $\mathbf{W} \in \mathbb{R}^{512 \times 2048}$ containing 1,048,576 parameters. The full SVD gives $\mathbf{W} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$ where $\mathbf{U} \in \mathbb{R}^{512 \times 512}$, $\mathbf{\Sigma} \in \mathbb{R}^{512 \times 2048}$, and $\mathbf{V} \in \mathbb{R}^{2048 \times 2048}$.

For a rank- k approximation, we keep only the top k singular values to obtain $\mathbf{W} \approx \mathbf{W}_1 \mathbf{W}_2 = \mathbf{U}_k \mathbf{\Sigma}_k \mathbf{V}_k^\top$ where $\mathbf{U}_k \in \mathbb{R}^{512 \times k}$, $\mathbf{\Sigma}_k \in \mathbb{R}^{k \times k}$, and $\mathbf{V}_k \in \mathbb{R}^{k \times 2048}$. We can absorb the diagonal matrix $\mathbf{\Sigma}_k$ into either factor, giving $\mathbf{W}_1 = \mathbf{U}_k \mathbf{\Sigma}_k \in \mathbb{R}^{512 \times k}$ and $\mathbf{W}_2 = \mathbf{V}_k^\top \in \mathbb{R}^{k \times 2048}$.

The original matrix has $512 \times 2048 = 1,048,576$ parameters, while the compressed form has $512k + 2048k = 2560k$ parameters, yielding a compression ratio of $\frac{2560k}{1,048,576} = \frac{k}{409.6}$. For $k = 64$, we have $2560 \times 64 = 163,840$ parameters, achieving 84.4% compression. For $k = 128$, we have 327,680 parameters, achieving 68.8% compression.

In terms of memory savings with 32-bit floats, the original matrix requires $1,048,576 \times 4 = 4,194,304$ bytes or approximately 4.0 MB. The compressed version with $k = 64$ requires only $163,840 \times 4 = 655,360$ bytes or approximately 0.625 MB, saving 3.375 MB per layer. For a model with 100 such layers, this yields a total savings of 337.5 MB, significantly reducing memory

footprint and enabling deployment on resource-constrained devices.

Example 1.13 (Accuracy vs Compression Trade-off). Consider a weight matrix with singular values that decay exponentially:

$$\sigma_i = \sigma_1 \cdot e^{-\alpha i} \quad (1.41)$$

The relative approximation error for rank- k approximation is:

$$\frac{\|\mathbf{W} - \mathbf{W}_k\|_F}{\|\mathbf{W}\|_F} = \sqrt{\frac{\sum_{i=k+1}^r \sigma_i^2}{\sum_{i=1}^r \sigma_i^2}} \quad (1.42)$$

Typical results for transformer feed-forward layers:

Rank k	Compression	Relative Error	Accuracy Drop
256	50%	0.05	<0.1%
128	75%	0.12	0.3%
64	87.5%	0.25	1.2%
32	93.75%	0.45	3.5%

Sweet spot: 50-75% compression with minimal accuracy loss.

SVD in Modern Architectures

Key Point 1.6. LoRA (Low-Rank Adaptation): Instead of fine-tuning all parameters in a pre-trained model, LoRA adds low-rank updates through the decomposition $\mathbf{W}' = \mathbf{W} + \Delta\mathbf{W} = \mathbf{W} + \mathbf{B}\mathbf{A}$, where the original weights $\mathbf{W} \in \mathbb{R}^{d \times k}$ remain frozen, and only the low-rank factors $\mathbf{B} \in \mathbb{R}^{d \times r}$ and $\mathbf{A} \in \mathbb{R}^{r \times k}$ with $r \ll \min(d, k)$ are trained. For a layer with $d = 4096$ and $k = 4096$, full fine-tuning requires updating 16,777,216 parameters, while LoRA with rank $r = 8$ requires only $8 \times (4096 + 4096) = 65,536$ trainable parameters, achieving a 99.6% parameter reduction while maintaining comparable performance.

Implementation:

Computing SVD and low-rank approximation in PyTorch:

```
import torch

# Original weight matrix
W = torch.randn(512, 2048)

# Compute SVD
U, S, Vt = torch.linalg.svd(W, full_matrices=False)

# Rank-k approximation
k = 64
W_compressed = U[:, :k] @ torch.diag(S[:k]) @ Vt[:k, :]

# Factored form for efficient computation
W1 = U[:, :k] @ torch.diag(S[:k]) # 512 x 64
W2 = Vt[:k, :] # 64 x 2048

# Verify approximation
error = torch.norm(W - W_compressed, p='fro')
```

```

relative_error = error / torch.norm(W, p='fro')
print(f"Relative error: {relative_error:.4f}")

# Memory comparison
original_params = W.numel()
compressed_params = W1.numel() + W2.numel()
compression_ratio = compressed_params / original_params
print(f"Compression: {(1-compression_ratio)*100:.1f}%")

```

1.5 Norms and Distance Metrics

Definition 1.10 (Vector Norms). For vector $\mathbf{x} \in \mathbb{R}^n$:

$$\text{L1 norm (Manhattan): } \|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i| \quad (1.43)$$

$$\text{L2 norm (Euclidean): } \|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^n x_i^2} \quad (1.44)$$

$$\text{L}\infty \text{ norm (Max): } \|\mathbf{x}\|_\infty = \max_i |x_i| \quad (1.45)$$

Definition 1.11 (Matrix Norms). For matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$:

$$\text{Frobenius norm: } \|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n a_{i,j}^2} = \sqrt{\text{tr}(\mathbf{A}^\top \mathbf{A})} \quad (1.46)$$

Norms are used in regularization to prevent overfitting by penalizing large weights.

Implementation:

In PyTorch:

```

import torch

# Vector norms
x = torch.tensor([3.0, 4.0])
l2_norm = torch.norm(x, p=2) # 5.0
l1_norm = torch.norm(x, p=1) # 7.0

# Matrix Frobenius norm
W = torch.randn(256, 784)
frob_norm = torch.norm(W, p='fro')

```

1.6 Practical Deep Learning Examples

1.6.1 Embedding Layers and Memory Requirements

Example 1.14 (Vocabulary Embeddings). Large language models use embedding layers to map discrete tokens to continuous vector representations. For GPT-3, the vocabulary contains $V = 50,257$ tokens, each mapped to a vector of dimension $d_{\text{model}} = 12,288$, requiring an embedding matrix $\mathbf{E} \in \mathbb{R}^{50,257 \times 12,288}$ with $50,257 \times 12,288 = 617,558,016$ parameters. Storing these embeddings in 32-bit floating-point format requires $617,558,016 \times 4 = 2,470,232,064$ bytes or approximately 2.3 GB of memory, while 16-bit format reduces this to approximately 1.15 GB.

For a batch of $B = 32$ sequences of length $n = 2048$, the input consists of integer token IDs in $\mathbb{R}^{32 \times 2048}$, which the embedding layer transforms into dense representations in $\mathbb{R}^{32 \times 2048 \times 12,288}$, requiring $32 \times 2048 \times 12,288 \times 4 = 3,221,225,472$ bytes or approximately 3.0 GB of memory. This demonstrates why large batch sizes and long sequences quickly exhaust GPU memory, necessitating techniques like gradient checkpointing and mixed-precision training.

1.6.2 Complete Transformer Layer Analysis

Example 1.15 (Full Transformer Layer Breakdown). Consider a single transformer layer with $d_{\text{model}} = 768$, $h = 12$ attention heads, $d_k = d_v = 64$ per head, and feed-forward dimension $d_{\text{ff}} = 3072$. The multi-head attention mechanism requires four weight matrices: $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{768 \times 768}$ for projecting to queries, keys, and values, plus $\mathbf{W}_O \in \mathbb{R}^{768 \times 768}$ for the output projection, totaling $4 \times 768^2 = 2,359,296$ parameters. The feed-forward network uses $\mathbf{W}_1 \in \mathbb{R}^{3072 \times 768}$ with bias $\mathbf{b}_1 \in \mathbb{R}^{3072}$ for the expansion, and $\mathbf{W}_2 \in \mathbb{R}^{768 \times 3072}$ with bias $\mathbf{b}_2 \in \mathbb{R}^{768}$ for the projection back, totaling $2 \times (768 \times 3072) + 3072 + 768 = 4,722,432$ parameters. Two layer normalization operations add $2 \times 2 \times 768 = 3,072$ parameters (scale and shift for each), bringing the total per layer to $2,359,296 + 4,722,432 + 3,072 = 7,084,800$ parameters.

For BERT-base with 12 such layers, the transformer stack contains $12 \times 7,084,800 = 85,017,600$ parameters (excluding embeddings), requiring $85,017,600 \times 4 = 340,070,400$ bytes or approximately 324 MB of memory for 32-bit weights. The computational cost for processing a batch of $B = 32$ sequences of length $n = 512$ includes approximately 154.6 GFLOPs for attention (from $4Bnd_{\text{model}}^2 + 2Bn^2d_{\text{model}}$) and 154.6 GFLOPs for the feed-forward network (from $4Bnd_{\text{model}}d_{\text{ff}}$), totaling approximately 309 GFLOPs per layer or 3.7 TFLOPs for all 12 layers. On an NVIDIA A100 GPU with 312 TFLOPS peak performance, a forward pass takes approximately 12 ms, though actual performance depends on memory bandwidth, kernel launch overhead, and other factors.

1.6.3 Common Dimension Errors and Debugging

Key Point 1.7. Dimension Mismatch Errors: *The most common bugs in deep learning involve incompatible tensor dimensions. When debugging dimension errors, start by printing tensor shapes using `print(x.shape)` to verify actual dimensions against expected values. Check whether the batch dimension is present (is it $\mathbb{R}^{B \times \dots}$ or just \mathbb{R}^{\dots} ?), verify the sequence length dimension (is it $\mathbb{R}^{B \times n \times d}$ or $\mathbb{R}^{B \times d \times n}$?), confirm that matrix multiplication has compatible inner dimensions, and watch for unintentional broadcasting that may hide shape mismatches. These systematic checks quickly identify the source of dimension errors and guide appropriate fixes.*

Implementation:

Common dimension fixes in PyTorch:

```
import torch
```

```

# Problem: Shape mismatch in matrix multiplication
Q = torch.randn(32, 512, 768) # [batch, seq_len, d_model]
K = torch.randn(32, 512, 768)

# Wrong: Q @ K gives error (768 != 512)
# scores = Q @ K # Error!

# Correct: Transpose last two dimensions of K
scores = Q @ K.transpose(-2, -1) # [32, 512, 512]

# Problem: Missing batch dimension
x = torch.randn(512, 768) # Missing batch dimension
W = torch.randn(768, 3072)

# Add batch dimension
x = x.unsqueeze(0) # [1, 512, 768]
output = x @ W      # [1, 512, 3072]

# Problem: Broadcasting confusion
x = torch.randn(32, 512, 768)
bias = torch.randn(768)

# This works due to broadcasting
output = x + bias # bias broadcasts to [32, 512, 768]

# Explicit broadcasting (clearer)
output = x + bias.view(1, 1, 768)

```

1.7 Exercises

Exercise 1.1. Given $\mathbf{x} = [2, -1, 3]^\top$ and $\mathbf{y} = [1, 4, -2]^\top$, compute:

1. The dot product $\mathbf{x}^\top \mathbf{y}$
2. The L2 norms $\|\mathbf{x}\|_2$ and $\|\mathbf{y}\|_2$
3. The cosine similarity between \mathbf{x} and \mathbf{y}

Exercise 1.2. For a transformer layer with $d_{\text{model}} = 768$ and feed-forward dimension $d_{\text{ff}} = 3072$:

1. Calculate the number of parameters in the two linear transformations
2. If processing a batch of $B = 32$ sequences of length $n = 512$, what are the dimensions of the input tensor?
3. How many floating-point operations (FLOPs) are required for one forward pass through this layer?

Exercise 1.3. Prove that for symmetric matrix $\mathbf{A} = \mathbf{A}^\top$, eigenvectors corresponding to distinct eigenvalues are orthogonal.

Exercise 1.4. A weight matrix $\mathbf{W} \in \mathbb{R}^{1024 \times 4096}$ is approximated using SVD with rank r .

1. Express the number of parameters as a function of r
2. What value of r achieves 75% compression?
3. What is the memory savings in MB (assuming 32-bit floats)?

Exercise 1.5. Consider computing attention scores $\mathbf{A} = \mathbf{Q}\mathbf{K}^\top$ where $\mathbf{Q}, \mathbf{K} \in \mathbb{R}^{B \times n \times d_k}$ with $B = 16, n = 1024, d_k = 64$.

1. What are the dimensions of the output \mathbf{A} ?
2. Calculate the total FLOPs required
3. Compute the arithmetic intensity (FLOPs per byte transferred, assuming 32-bit floats)
4. Is this operation compute-bound or memory-bound on a GPU with 312 TFLOPS and 1.6 TB/s bandwidth?

Exercise 1.6. An embedding layer has vocabulary size $V = 32,000$ and embedding dimension $d = 512$.

1. How many parameters does the embedding matrix contain?
2. What is the memory requirement in MB for 32-bit floats?
3. For a batch of $B = 64$ sequences of length $n = 256$, what is the memory required for the embedded representations?
4. If we use LoRA with rank $r = 16$ to adapt the embeddings, how many trainable parameters are needed?

Exercise 1.7. Compare the computational cost of two equivalent operations:

1. Computing $(\mathbf{A}\mathbf{B})\mathbf{x}$ where $\mathbf{A} \in \mathbb{R}^{m \times n}, \mathbf{B} \in \mathbb{R}^{n \times p}, \mathbf{x} \in \mathbb{R}^p$
2. Computing $\mathbf{A}(\mathbf{B}\mathbf{x})$

For $m = 512, n = 2048, p = 512$, which order is more efficient and by what factor?

Exercise 1.8. A matrix multiplication $\mathbf{C} = \mathbf{A}\mathbf{B}$ with $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{2048 \times 2048}$ is performed on a GPU.

1. Calculate the total FLOPs
2. Calculate the memory transferred (assuming matrices are read once and result written once)
3. Compute the arithmetic intensity

4. If the GPU has 100 TFLOPS compute and 900 GB/s memory bandwidth, what is the theoretical execution time assuming perfect utilization?
5. Which resource (compute or memory) is the bottleneck?

Chapter 2

Calculus and Optimization

Chapter Overview

Training deep learning models requires optimizing complex, high-dimensional functions. This chapter develops the calculus and optimization theory necessary to understand how neural networks learn from data. We cover multivariable calculus, gradient computation, and the optimization algorithms that power modern deep learning.

The centerpiece of this chapter is backpropagation, the algorithm that efficiently computes gradients in neural networks. We derive backpropagation from first principles, showing how the chain rule enables gradient computation through arbitrarily deep computational graphs. We then explore gradient descent and its variants, which use these gradients to iteratively improve model parameters.

Learning Objectives

After completing this chapter, you will be able to:

1. Compute gradients and Jacobians for multivariable functions
2. Apply the chain rule to composite functions
3. Understand and implement the backpropagation algorithm
4. Implement gradient descent and its variants (SGD, momentum, Adam)
5. Analyze convergence properties of optimization algorithms
6. Apply learning rate schedules and regularization techniques

2.1 Multivariable Calculus

2.1.1 Partial Derivatives

Definition 2.1 (Partial Derivative). For function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, the **partial derivative** with respect to x_i is:

$$\frac{\partial f}{\partial x_i} = \lim_{h \rightarrow 0} \frac{f(x_1, \dots, x_i + h, \dots, x_n) - f(x_1, \dots, x_i, \dots, x_n)}{h} \quad (2.1)$$

Example 2.1 (Computing Partial Derivatives). For $f(x_1, x_2) = x_1^2 + 3x_1x_2 + x_2^2$:

$$\frac{\partial f}{\partial x_1} = 2x_1 + 3x_2 \quad (2.2)$$

$$\frac{\partial f}{\partial x_2} = 3x_1 + 2x_2 \quad (2.3)$$

At point $(x_1, x_2) = (1, 2)$:

$$\left. \frac{\partial f}{\partial x_1} \right|_{(1,2)} = 2(1) + 3(2) = 8 \quad (2.4)$$

$$\left. \frac{\partial f}{\partial x_2} \right|_{(1,2)} = 3(1) + 2(2) = 7 \quad (2.5)$$

2.1.2 Gradients

Definition 2.2 (Gradient). For function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, the **gradient** is the vector of partial derivatives:

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix} \in \mathbb{R}^n \quad (2.6)$$

The gradient points in the direction of steepest ascent of the function.

Example 2.2 (Gradient of Loss Function). For mean squared error loss:

$$L(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N (y_i - \mathbf{w}^\top \mathbf{x}^{(i)})^2 \quad (2.7)$$

The gradient with respect to \mathbf{w} is:

$$\nabla_{\mathbf{w}} L = -\frac{2}{N} \sum_{i=1}^N (y_i - \mathbf{w}^\top \mathbf{x}^{(i)}) \mathbf{x}^{(i)} \quad (2.8)$$

For $N = 1$, $\mathbf{w} = [w_1, w_2]^\top$, $\mathbf{x} = [1, 2]^\top$, $y = 5$, and current prediction $\hat{y} = \mathbf{w}^\top \mathbf{x} = 3$:

$$\nabla_{\mathbf{w}} L = -2(5 - 3) \begin{bmatrix} 1 \\ 2 \end{bmatrix} = \begin{bmatrix} -4 \\ -8 \end{bmatrix} \quad (2.9)$$

The negative gradient $-\nabla_{\mathbf{w}} L = [4, 8]^\top$ points toward better parameters.

2.1.3 The Chain Rule

Theorem 2.1 (Chain Rule for Functions). For composite function $h(\mathbf{x}) = f(g(\mathbf{x}))$ where $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $f : \mathbb{R}^m \rightarrow \mathbb{R}$:

$$\frac{\partial h}{\partial x_i} = \sum_{j=1}^m \frac{\partial f}{\partial g_j} \frac{\partial g_j}{\partial x_i} \quad (2.10)$$

In vector form:

$$\nabla_{\mathbf{x}} h = \mathbf{J}_g^\top \nabla_{\mathbf{z}} f \quad (2.11)$$

where $\mathbf{z} = g(\mathbf{x})$ and $\mathbf{J}_g \in \mathbb{R}^{m \times n}$ is the Jacobian of g .

Example 2.3 (Chain Rule Application). For neural network layer: $\mathbf{y} = \sigma(\mathbf{W}\mathbf{x} + \mathbf{b})$ where σ is applied element-wise.

Let $\mathbf{z} = \mathbf{W}\mathbf{x} + \mathbf{b}$ (pre-activation). Then:

$$\frac{\partial L}{\partial \mathbf{x}} = \mathbf{W}^\top \left(\frac{\partial L}{\partial \mathbf{y}} \odot \sigma'(\mathbf{z}) \right) \quad (2.12)$$

where \odot denotes element-wise multiplication.

Concrete example: For ReLU activation $\sigma(z) = \max(0, z)$:

$$\sigma'(z) = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{if } z \leq 0 \end{cases} \quad (2.13)$$

If $\mathbf{z} = [2.0, -1.0, 0.5]^\top$, then $\sigma'(\mathbf{z}) = [1, 0, 1]^\top$.

2.1.4 Jacobian and Hessian Matrices

Definition 2.3 (Jacobian Matrix). For function $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$, the **Jacobian matrix** is:

$$\mathbf{J}_f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \cdots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} \in \mathbb{R}^{m \times n} \quad (2.14)$$

Definition 2.4 (Hessian Matrix). For function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, the **Hessian matrix** contains second derivatives:

$$\mathbf{H}_f(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix} \in \mathbb{R}^{n \times n} \quad (2.15)$$

The Hessian describes the local curvature of the function. For smooth functions, \mathbf{H} is symmetric.

2.2 Gradient Descent

2.2.1 The Gradient Descent Algorithm

Gradient descent iteratively moves parameters in the direction opposite to the gradient:

[H] Gradient Descent Objective function $f(\mathbf{w})$, initial parameters $\mathbf{w}^{(0)}$, learning rate η , iterations T
 Optimized parameters $\mathbf{w}^{(T)}$ $t = 0 \rightarrow T - 1$ Compute gradient: $\mathbf{g}^{(t)} = \nabla f(\mathbf{w}^{(t)})$
 Update parameters: $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \mathbf{g}^{(t)}$

Key Point 2.1. *The learning rate η controls the step size. Too large: divergence. Too small: slow convergence.*

Example 2.4 (Gradient Descent on Quadratic). Minimize $f(w) = w^2$ starting from $w^{(0)} = 3$ with $\eta = 0.1$:

$$t = 0 : \quad w^{(0)} = 3, \quad g^{(0)} = 2w^{(0)} = 6, \quad w^{(1)} = 3 - 0.1(6) = 2.4 \quad (2.16)$$

$$t = 1 : \quad w^{(1)} = 2.4, \quad g^{(1)} = 4.8, \quad w^{(2)} = 2.4 - 0.1(4.8) = 1.92 \quad (2.17)$$

$$t = 2 : \quad w^{(2)} = 1.92, \quad g^{(2)} = 3.84, \quad w^{(3)} = 1.92 - 0.1(3.84) = 1.536 \quad (2.18)$$

The parameters converge to $w^* = 0$ (the minimum).

2.2.2 Stochastic Gradient Descent (SGD)

For large datasets, computing the full gradient is expensive. SGD approximates the gradient using mini-batches.

[H] Stochastic Gradient Descent (SGD) Dataset $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$, batch size B , learning rate η , epochs E Optimize parameters \mathbf{w} Initialize \mathbf{w} randomly
epoch $e = 1$ E Shuffle dataset \mathcal{D}
each mini-batch $\mathcal{B} \subset \mathcal{D}$ of size B Compute mini-batch gradient: $\mathbf{g} = \frac{1}{B} \sum_{(\mathbf{x}, y) \in \mathcal{B}} \nabla_{\mathbf{w}} L(\mathbf{w}; \mathbf{x}, y)$
Update: $\mathbf{w} \leftarrow \mathbf{w} - \eta \mathbf{g}$ \mathbf{w}

Implementation:

PyTorch SGD implementation:

```
import torch
import torch.nn as nn

# Model and loss
model = nn.Linear(10, 1)
criterion = nn.MSELoss()

# SGD optimizer
optimizer = torch.optim.SGD(model.parameters(), lr=0.01)

# Training loop
for epoch in range(100):
    for x_batch, y_batch in dataloader:
        # Forward pass
        y_pred = model(x_batch)
        loss = criterion(y_pred, y_batch)

        # Backward pass
        optimizer.zero_grad() # Clear previous gradients
        loss.backward()        # Compute gradients
        optimizer.step()        # Update parameters
```

2.2.3 Momentum

Momentum accelerates SGD by accumulating a velocity vector:

[H] SGD with Momentum Learning rate η , momentum coefficient β (typically 0.9) Initialize velocity $\mathbf{v} = \mathbf{0}$
 each iteration Compute gradient $\mathbf{g} = \nabla_{\mathbf{w}}L(\mathbf{w})$
 Update velocity: $\mathbf{v} \leftarrow \beta\mathbf{v} + \mathbf{g}$
 Update parameters: $\mathbf{w} \leftarrow \mathbf{w} - \eta\mathbf{v}$
 Momentum helps navigate ravines and accelerates convergence in relevant directions.

2.2.4 Adam Optimizer

Adam (Adaptive Moment Estimation) combines momentum with adaptive learning rates:

[H] Adam Optimizer Learning rate α (default 0.001), $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-8}$ Initialize $\mathbf{m}_0 = \mathbf{0}$ (first moment), $\mathbf{v}_0 = \mathbf{0}$ (second moment), $t = 0$
 not converged $t \leftarrow t + 1$
 Compute gradient: $\mathbf{g}_t = \nabla_{\mathbf{w}}L(\mathbf{w}_{t-1})$
 Update biased first moment: $\mathbf{m}_t = \beta_1\mathbf{m}_{t-1} + (1 - \beta_1)\mathbf{g}_t$
 Update biased second moment: $\mathbf{v}_t = \beta_2\mathbf{v}_{t-1} + (1 - \beta_2)\mathbf{g}_t^2$
 Bias-corrected first moment: $\hat{\mathbf{m}}_t = \mathbf{m}_t / (1 - \beta_1^t)$
 Bias-corrected second moment: $\hat{\mathbf{v}}_t = \mathbf{v}_t / (1 - \beta_2^t)$
 Update parameters: $\mathbf{w}_t = \mathbf{w}_{t-1} - \alpha\hat{\mathbf{m}}_t / (\sqrt{\hat{\mathbf{v}}_t} + \epsilon)$

Adam is the most commonly used optimizer for training transformers and large language models.

2.3 Gradient Computation Complexity

Understanding the computational and memory costs of gradient computation is essential for training large models efficiently.

2.3.1 FLOPs for Gradient Computation

Key Point 2.2. *Computing gradients via backpropagation requires approximately $2\times$ the FLOPs of the forward pass: $1\times$ for the backward pass itself, plus the original $1\times$ forward pass.*

Example 2.5 (BERT-base Gradient Computation). For BERT-base (110M parameters, 12 layers, $d_{\text{model}} = 768$) processing sequence length $n = 512$:

Forward pass:

- Self-attention: $12 \times 4n^2d = 12 \times 4(512)^2(768) \approx 48$ GFLOPs
- Feed-forward: $12 \times 2nd(4d) = 12 \times 2(512)(768)(3072) \approx 36$ GFLOPs
- Other operations: ≈ 12 GFLOPs
- **Total forward: ≈ 96 GFLOPs**

Backward pass:

- Gradient computation through each layer: ≈ 96 GFLOPs
- Gradient accumulation for weight updates: ≈ 97 GFLOPs
- **Total backward: ≈ 193 GFLOPs**

Total per training step: ≈ 289 GFLOPs

For batch size $B = 32$: $289 \times 32 \approx 9.2$ TFLOPs per batch.

2.3.2 Memory Requirements for Activations

During backpropagation, intermediate activations must be stored for gradient computation.

Definition 2.5 (Activation Memory). For a network with L layers processing batch size B , activation memory is:

$$M_{\text{act}} = B \sum_{\ell=1}^L d_{\ell} \quad (2.19)$$

where d_{ℓ} is the dimension of layer ℓ 's output.

Example 2.6 (BERT-base Activation Memory). For BERT-base with batch size $B = 32$, sequence length $n = 512$, $d = 768$:

Per transformer layer:

- Query, Key, Value projections: $3 \times Bnd = 3 \times 32 \times 512 \times 768 \times 4 \text{ bytes} \approx 113 \text{ MB}$
- Attention scores: $B \times h \times n \times n = 32 \times 12 \times 512 \times 512 \times 4 \text{ bytes} \approx 402 \text{ MB}$
- Attention output: $Bnd \approx 38 \text{ MB}$
- Feed-forward intermediate: $B \times n \times 4d \approx 151 \text{ MB}$
- **Per layer total:** $\approx 704 \text{ MB}$

For 12 layers: $704 \times 12 \approx 8.4 \text{ GB}$

This excludes gradients and optimizer states!

2.3.3 Automatic Differentiation: Forward vs Reverse Mode

Definition 2.6 (Forward Mode AD). Forward mode computes derivatives by propagating tangent vectors forward through the computational graph. For function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$, computing ∇f requires n forward passes.

Definition 2.7 (Reverse Mode AD). Reverse mode (backpropagation) computes derivatives by propagating adjoints backward. Computing ∇f requires 1 forward pass + 1 backward pass, regardless of n .

Key Point 2.3. For neural networks where $n \gg m$ (millions of parameters, one loss), reverse mode is vastly more efficient: $O(1)$ passes vs $O(n)$ passes.

Example 2.7 (Forward vs Reverse Mode Comparison). For a network with $n = 10^8$ parameters (100M) and scalar loss ($m = 1$):

Forward mode:

- Requires 10^8 forward passes
- Each pass: ≈ 100 GFLOPs
- Total: 10^{10} GFLOPs ≈ 10 PFLOPs
- Time on A100 GPU (312 TFLOPS): $\approx 32,000$ seconds ≈ 9 hours

Reverse mode (backpropagation):

- Requires 1 forward + 1 backward pass
- Total: ≈ 300 GFLOPs
- Time on A100 GPU: ≈ 0.001 seconds
- **Speedup: ≈ 32 million \times**

2.3.4 Gradient Checkpointing

Gradient checkpointing trades computation for memory by recomputing activations during the backward pass.

[H] Gradient Checkpointing Network with L layers, checkpoint every k layers Forward Pass $\ell = 1$ to L Compute $\mathbf{h}^{(\ell)} = f^{(\ell)}(\mathbf{h}^{(\ell-1)})$
 $\ell \bmod k = 0$ Save $\mathbf{h}^{(\ell)}$ to memory (checkpoint)

Backward Pass $\ell = L$ to 1 $\mathbf{h}^{(\ell)}$ not in memory Recompute forward from last checkpoint to layer ℓ
 Compute gradient $\nabla_{\mathbf{h}^{(\ell-1)}} L$ using $\mathbf{h}^{(\ell)}$

Example 2.8 (Checkpointing Trade-off). For BERT-base (12 layers) with checkpointing every 3 layers:

Without checkpointing:

- Memory: 8.4 GB (all activations)
- Computation: 289 GFLOPs (1 forward + 1 backward)

With checkpointing (every 3 layers):

- Memory: $8.4/3 \approx 2.8$ GB (only checkpoints)
- Computation: $96 + 193 + 72 = 361$ GFLOPs (1 forward + 1 backward + 0.75 forward recompute)
- **Memory reduction: $3\times$, Computation increase: $1.25\times$**

For GPT-3 (175B parameters), checkpointing is essential to fit in GPU memory.

2.4 Backpropagation

Backpropagation efficiently computes gradients in neural networks using the chain rule.

2.4.1 Computational Graphs

A computational graph represents the sequence of operations in a neural network. Each node is an operation, and edges carry values/gradients.

Example 2.9 (Simple Computational Graph). For $L = (y - \hat{y})^2$ where $\hat{y} = w_2\sigma(w_1x + b_1) + b_2$:

Forward pass:

$$z_1 = w_1x + b_1 = 2.0(1.0) + 0.5 = 2.5 \quad (2.20)$$

$$a_1 = \sigma(z_1) = \sigma(2.5) = 0.924 \quad (\text{sigmoid}) \quad (2.21)$$

$$z_2 = w_2a_1 + b_2 = 1.5(0.924) + 0.3 = 1.686 \quad (2.22)$$

$$L = (y - z_2)^2 = (3.0 - 1.686)^2 = 1.726 \quad (2.23)$$

Backward pass:

$$\frac{\partial L}{\partial z_2} = 2(z_2 - y) = 2(1.686 - 3.0) = -2.628 \quad (2.24)$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial z_2} \cdot a_1 = -2.628(0.924) = -2.428 \quad (2.25)$$

$$\frac{\partial L}{\partial a_1} = \frac{\partial L}{\partial z_2} \cdot w_2 = -2.628(1.5) = -3.942 \quad (2.26)$$

$$\frac{\partial L}{\partial z_1} = \frac{\partial L}{\partial a_1} \cdot \sigma'(z_1) = -3.942(0.070) = -0.276 \quad (2.27)$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial z_1} \cdot x = -0.276(1.0) = -0.276 \quad (2.28)$$

2.4.2 Backpropagation Algorithm

[H] Backpropagation Training example (\mathbf{x}, y) , network with L layers Gradients $\{\nabla_{\mathbf{W}^{(\ell)}} L, \nabla_{\mathbf{b}^{(\ell)}} L\}_{\ell=1}^L$

Forward Pass $\mathbf{h}^{(0)} = \mathbf{x}$

$$\ell = 1 \quad L \quad \mathbf{z}^{(\ell)} = \mathbf{W}^{(\ell)} \mathbf{h}^{(\ell-1)} + \mathbf{b}^{(\ell)}$$

$$\mathbf{h}^{(\ell)} = \sigma^{(\ell)}(\mathbf{z}^{(\ell)}) \quad \hat{y} = \mathbf{h}^{(L)}$$

Compute loss: $L = \text{Loss}(y, \hat{y})$

$$\text{Backward Pass } \boldsymbol{\delta}^{(L)} = \nabla_{\mathbf{h}^{(L)}} L \odot \sigma'^{(L)}(\mathbf{z}^{(L)})$$

$$\ell = L \quad 1 \quad \nabla_{\mathbf{W}^{(\ell)}} L = \boldsymbol{\delta}^{(\ell)} (\mathbf{h}^{(\ell-1)})^\top$$

$$\nabla_{\mathbf{b}^{(\ell)}} L = \boldsymbol{\delta}^{(\ell)}$$

$$\ell > 1 \quad \boldsymbol{\delta}^{(\ell-1)} = (\mathbf{W}^{(\ell)})^\top \boldsymbol{\delta}^{(\ell)} \odot \sigma'^{(\ell-1)}(\mathbf{z}^{(\ell-1)}) \quad \text{All gradients}$$

Key Point 2.4. Backpropagation computes gradients in $O(n)$ time where n is the number of parameters, compared to $O(n^2)$ for naive methods. This efficiency enables training of billion-parameter models.

2.4.3 Why Backpropagation is $O(n)$ Not $O(n^2)$

Theorem 2.2 (Backpropagation Complexity). For a neural network with n parameters and m operations, backpropagation computes all gradients in $O(m)$ time, where typically $m = O(n)$.

Intuition. Each operation in the forward pass corresponds to one gradient computation in the backward pass. The chain rule allows us to reuse intermediate gradients:

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial z_j} \cdot \frac{\partial z_j}{\partial w_i} \quad (2.29)$$

We compute $\frac{\partial L}{\partial z_j}$ once and reuse it for all parameters that affect z_j . This sharing prevents the $O(n^2)$ cost of computing each gradient independently. \square

Example 2.10 (Complexity Comparison). For a network with $n = 10^8$ parameters:

Naive finite differences:

$$\frac{\partial L}{\partial w_i} \approx \frac{L(w_i + \epsilon) - L(w_i)}{\epsilon} \quad (2.30)$$

Requires n forward passes: $O(n \cdot m) = O(n^2)$ operations.

Backpropagation:

- Forward pass: $O(m) = O(n)$ operations
- Backward pass: $O(m) = O(n)$ operations
- Total: $O(n)$ operations

Speedup: $O(n) = 10^8 \times$

2.5 Optimizer Memory Requirements

Different optimizers have vastly different memory requirements, which becomes critical for large models.

2.5.1 Memory Comparison by Optimizer

Optimizer	Memory per Parameter	Total Memory Factor
SGD (no momentum)	4 bytes (fp32)	$1\times$
SGD with momentum	8 bytes (param + velocity)	$2\times$
Adam	16 bytes (param + 2 moments)	$4\times$
Adam (mixed precision)	10 bytes (fp16 param + fp32 master + 2 moments)	$2.5\times$

Table 2.1: Memory requirements per parameter for different optimizers

Example 2.11 (BERT-base Memory Requirements). For BERT-base with 110M parameters:

Model parameters:

- FP32: $110 \times 10^6 \times 4 \text{ bytes} = 440 \text{ MB}$
- FP16: $110 \times 10^6 \times 2 \text{ bytes} = 220 \text{ MB}$

SGD with momentum:

- Parameters: 440 MB
- Momentum buffer: 440 MB
- **Total: 880 MB**

Adam optimizer:

- Parameters: 440 MB
- First moment (**m**): 440 MB
- Second moment (**v**): 440 MB
- Gradients: 440 MB

- **Total: 1,760 MB \approx 1.7 GB**

Adam with mixed precision:

- FP16 parameters: 220 MB
- FP32 master copy: 440 MB
- FP32 first moment: 440 MB
- FP32 second moment: 440 MB
- FP16 gradients: 220 MB
- **Total: 1,760 MB \approx 1.7 GB**

Note: Mixed precision doesn't reduce optimizer memory, but enables larger batch sizes.

Example 2.12 (GPT-3 Memory Requirements). For GPT-3 (175B parameters) with Adam optimizer:

Model + optimizer states:

- Parameters (FP16): $175 \times 10^9 \times 2 = 350$ GB
- Master copy (FP32): $175 \times 10^9 \times 4 = 700$ GB
- First moment (FP32): 700 GB
- Second moment (FP32): 700 GB
- Gradients (FP16): 350 GB
- **Total: 2,800 GB \approx 2.8 TB**

This requires distributed training across multiple GPUs. With $8 \times$ A100 GPUs (80 GB each = 640 GB total), we need model parallelism and optimizer state sharding (e.g., ZeRO optimizer).

2.5.2 Impact on GPU Memory Budget

Key Point 2.5. *For large models, optimizer states often consume more memory than the model itself. Adam uses $4 \times$ parameter memory, leaving less room for batch size and activations.*

Example 2.13 (Memory Budget Breakdown). Training BERT-base on A100 GPU (80 GB memory):

Memory allocation:

- Model parameters: 0.44 GB
- Optimizer states (Adam): 1.32 GB
- Activations (batch size 32): 8.4 GB
- Gradients: 0.44 GB

- Framework overhead: ≈ 2 GB
- **Total:** ≈ 12.6 GB

Remaining: 67.4 GB available for larger batch sizes or longer sequences.
With batch size 256: Activations ≈ 67 GB, total ≈ 71 GB (fits comfortably).

2.6 Learning Rate Schedules

Learning rate schedules adjust η during training to improve convergence.

2.6.1 Learning Rate Impact on Convergence and GPU Utilization

Key Point 2.6. *Learning rate affects both convergence speed and hardware efficiency. Larger learning rates enable larger batch sizes, improving GPU utilization.*

Example 2.14 (Learning Rate vs Convergence Speed). Training BERT-base on 1M examples: Optimal learning rate (5×10^{-4}) achieves $6.7\times$ faster convergence than conservative rate.

2.6.2 Learning Rate Scaling with Batch Size

Theorem 2.3 (Linear Scaling Rule). *When increasing batch size by factor k , scale learning rate by k to maintain convergence behavior:*

$$\eta_{\text{new}} = k \cdot \eta_{\text{base}} \quad (2.31)$$

This holds approximately for $k \leq 8$. For larger k , use gradual warmup.

Example 2.15 (Batch Size and Learning Rate Scaling). Training BERT-base with different batch sizes:

Larger batches improve GPU utilization but require proportionally larger learning rates. Throughput increases $4.6\times$ from batch 32 to 512.

2.6.3 Practical Learning Rates for Transformers

Model	Batch Size	Peak Learning Rate
BERT-base	256	1×10^{-4}
BERT-large	256	5×10^{-5}
GPT-2 (117M)	512	2.5×10^{-4}
GPT-2 (1.5B)	512	1.5×10^{-4}
GPT-3 (175B)	3.2M	6×10^{-5}
T5-base	128	1×10^{-4}
T5-11B	2048	1×10^{-4}

Table 2.4: Typical learning rates for transformer models

Key Point 2.7. *Larger models generally require smaller learning rates for stability. GPT-3 uses 6×10^{-5} despite massive batch size of 3.2M tokens.*

2.6.4 Common Schedules

Step Decay:

$$\eta_t = \eta_0 \gamma^{\lfloor t/s \rfloor} \quad (2.32)$$

where $\gamma < 1$ (e.g., 0.1) and s is step size (e.g., every 10 epochs).

Exponential Decay:

$$\eta_t = \eta_0 e^{-\lambda t} \quad (2.33)$$

Cosine Annealing:

$$\eta_t = \eta_{\min} + \frac{1}{2}(\eta_{\max} - \eta_{\min}) \left(1 + \cos \left(\frac{t\pi}{T} \right) \right) \quad (2.34)$$

Warmup + Decay (Transformers):

$$\eta_t = \frac{d_{\text{model}}^{-0.5}}{\max(t, \text{warmup_steps}^{-0.5})} \cdot \min(t^{-0.5}, t \cdot \text{warmup_steps}^{-1.5}) \quad (2.35)$$

The warmup phase prevents instability in early training of transformers.

2.7 Hardware Considerations for Gradient Computation

Modern deep learning relies on specialized hardware for efficient gradient computation.

2.7.1 Gradient Computation on GPUs

GPUs excel at gradient computation due to massive parallelism in matrix operations.

Example 2.16 (GPU vs CPU Gradient Computation). Computing gradients for BERT-base (110M parameters) on one training batch:

NVIDIA A100 GPU:

- Forward pass: 0.31 ms (96 GFLOPs ÷ 312 TFLOPS)
- Backward pass: 0.62 ms (193 GFLOPs ÷ 312 TFLOPS)
- **Total: 0.93 ms per batch**
- Throughput: 1,075 batches/second

Intel Xeon CPU (32 cores):

- Forward pass: 45 ms (96 GFLOPs ÷ 2.1 TFLOPS)
- Backward pass: 90 ms (193 GFLOPs ÷ 2.1 TFLOPS)
- **Total: 135 ms per batch**
- Throughput: 7.4 batches/second

GPU speedup: 145×

2.7.2 Mixed Precision Training

Mixed precision uses FP16 for computation and FP32 for accumulation, reducing memory and increasing speed.

[H] Mixed Precision Training Model parameters \mathbf{w} (FP32 master copy) each training step Convert \mathbf{w} to FP16: $\mathbf{w}_{16} = \text{FP16}(\mathbf{w})$
 Forward pass in FP16: $\hat{y} = f(\mathbf{x}; \mathbf{w}_{16})$
 Compute loss: $L = \text{Loss}(y, \hat{y})$
 Scale loss: $L_{\text{scaled}} = s \cdot L$ (prevent underflow)
 Backward pass in FP16: $\mathbf{g}_{16} = \nabla_{\mathbf{w}_{16}} L_{\text{scaled}}$
 Unscale gradients: $\mathbf{g}_{16} = \mathbf{g}_{16}/s$
 Convert to FP32: $\mathbf{g} = \text{FP32}(\mathbf{g}_{16})$
 Update FP32 master: $\mathbf{w} \leftarrow \mathbf{w} - \eta \mathbf{g}$

Example 2.17 (Mixed Precision Impact). Training BERT-base on A100 GPU:

FP32 training:

- Forward + backward: 0.93 ms
- Memory: 12.6 GB
- Max batch size: 32
- Throughput: 34,400 samples/sec

Mixed precision (FP16) training:

- Forward + backward: 0.48 ms ($1.94\times$ faster)
- Memory: 8.2 GB (35% reduction)
- Max batch size: 64
- Throughput: 133,300 samples/sec ($3.87\times$ faster)

Mixed precision provides $1.94\times$ computational speedup and enables $2\times$ larger batches, yielding $3.87\times$ total throughput improvement.

2.7.3 Gradient Accumulation

Gradient accumulation simulates large batch sizes by accumulating gradients over multiple forward-backward passes.

[H] Gradient Accumulation Desired batch size B , physical batch size b , accumulation steps $k = B/b$
 Initialize gradients: $\mathbf{g}_{\text{acc}} = \mathbf{0}$
 $i = 1$ to k Sample mini-batch \mathcal{B}_i of size b
 Forward pass: $L_i = \text{Loss}(\mathcal{B}_i)$
 Backward pass: $\mathbf{g}_i = \nabla L_i$
 Accumulate: $\mathbf{g}_{\text{acc}} \leftarrow \mathbf{g}_{\text{acc}} + \mathbf{g}_i$ Average: $\mathbf{g}_{\text{acc}} \leftarrow \mathbf{g}_{\text{acc}}/k$
 Update parameters: $\mathbf{w} \leftarrow \mathbf{w} - \eta \mathbf{g}_{\text{acc}}$
 Clear gradients: $\mathbf{g}_{\text{acc}} = \mathbf{0}$

Example 2.18 (Gradient Accumulation for Large Batches). Training GPT-2 (1.5B parameters) on single A100 GPU (80 GB):

Without accumulation:

- Max batch size: 4 (memory limit)
- Update frequency: every 4 samples
- Training unstable (batch too small)

With gradient accumulation (32 steps):

- Physical batch size: 4
- Effective batch size: $4 \times 32 = 128$
- Update frequency: every 128 samples
- Memory: same as batch size 4
- Training stable and efficient

Trade-off: $32\times$ more forward-backward passes per update, but enables training large models on limited hardware.

2.7.4 Distributed Gradient Synchronization

For multi-GPU training, gradients must be synchronized across devices.

[H] Data Parallel Training with Gradient Synchronization N GPUs, global batch size B , local batch size $b = B/N$ each GPU $i = 1, \dots, N$ in parallel Sample local mini-batch \mathcal{B}_i of size b

Forward pass: $L_i = \text{Loss}(\mathcal{B}_i)$

Backward pass: $\mathbf{g}_i = \nabla L_i$ All-reduce gradients: $\mathbf{g} = \frac{1}{N} \sum_{i=1}^N \mathbf{g}_i$

each GPU $i = 1, \dots, N$ in parallel Update local parameters: $\mathbf{w}_i \leftarrow \mathbf{w}_i - \eta \mathbf{g}$

Example 2.19 (Distributed Training Efficiency). Training BERT-base on $8\times$ A100 GPUs with NVLink:

Single GPU baseline:

- Batch size: 32
- Time per step: 0.93 ms
- Throughput: 34,400 samples/sec

8 GPUs (data parallel):

- Global batch size: 256
- Time per step: 0.93 ms (computation) + 0.12 ms (communication)
- Total: 1.05 ms
- Throughput: 243,800 samples/sec
- **Scaling efficiency:** $243,800 / (8 \times 34,400) = 88.6\%$

Communication overhead is 11.4% due to gradient all-reduce. NVLink (600 GB/s) enables efficient synchronization.

Key Point 2.8. For large models, gradient synchronization can become a bottleneck. Techniques like gradient compression, ZeRO optimizer, and pipeline parallelism reduce communication over-

head.

2.8 Exercises

Exercise 2.1. Compute the gradient of $f(\mathbf{w}) = \mathbf{w}^\top \mathbf{A} \mathbf{w} + \mathbf{b}^\top \mathbf{w} + c$ where $\mathbf{A} \in \mathbb{R}^{n \times n}$ is symmetric, $\mathbf{w}, \mathbf{b} \in \mathbb{R}^n$, and $c \in \mathbb{R}$.

Exercise 2.2. Implement backpropagation for a 2-layer network with ReLU activation. Given input $\mathbf{x} = [1.0, 0.5]^\top$, weights $\mathbf{W}^{(1)} \in \mathbb{R}^{3 \times 2}$, $\mathbf{W}^{(2)} \in \mathbb{R}^{1 \times 3}$, and target $y = 2.0$, compute all gradients.

Exercise 2.3. For Adam optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\alpha = 0.001$:

1. Why is bias correction necessary?
2. What are the effective learning rates after steps $t = 1, 10, 100, 1000$?
3. How does Adam handle sparse gradients compared to SGD?

Exercise 2.4. A transformer is trained with learning rate warmup over 4000 steps, then inverse square root decay. If $d_{\text{model}} = 512$:

1. Plot the learning rate schedule for 100,000 steps
2. What is the learning rate at step 1, 4000, and 10,000?
3. Why is warmup beneficial for transformer training?

Exercise 2.5. Calculate the memory requirements for training GPT-2 (1.5B parameters) with Adam optimizer:

1. Model parameters in FP16
2. Optimizer states (FP32 master copy + 2 moments)
3. Gradients in FP16
4. Total memory for model + optimizer
5. How many A100 GPUs (80 GB each) are needed?

Exercise 2.6. For BERT-base processing sequence length 512 with batch size 64:

1. Calculate total FLOPs for one training step (forward + backward)

2. Estimate time per step on A100 GPU (312 TFLOPS)
3. How does mixed precision (FP16) affect throughput?
4. What is the maximum batch size that fits in 80 GB memory?

Exercise 2.7. Compare gradient computation methods for a network with 10^7 parameters:

1. How many forward passes does finite differences require?
2. How many passes does backpropagation require?
3. If one forward pass takes 10 ms, compare total time
4. Why is reverse mode AD preferred over forward mode?

Exercise 2.8. Implement gradient checkpointing for a 24-layer transformer:

1. Without checkpointing, how much activation memory is needed?
2. With checkpointing every 6 layers, what is the memory reduction?
3. What is the computational overhead (extra forward passes)?
4. At what model size does checkpointing become necessary?

Exercise 2.9. Analyze distributed training efficiency for 8 GPUs:

1. If gradient all-reduce takes 15 ms and computation takes 100 ms, what is the scaling efficiency?
2. How does batch size affect communication overhead?
3. Compare ring all-reduce vs tree all-reduce for 64 GPUs
4. When does gradient compression become beneficial?

Chapter 3

Probability and Information Theory

Chapter Overview

Deep learning is fundamentally a probabilistic framework. Neural networks learn probability distributions over data, make predictions with uncertainty, and are trained using probabilistic objectives. This chapter develops the probability theory and information theory necessary to understand these probabilistic aspects of deep learning.

We cover probability distributions, conditional probability, expectation, and variance—the building blocks for understanding neural network outputs as probabilistic models. We then introduce information theory concepts like entropy, cross-entropy, and KL divergence, which form the basis for loss functions used in training.

Learning Objectives

After completing this chapter, you will be able to:

1. Work with probability distributions and compute expectations
2. Apply Bayes' theorem to understand conditional probabilities
3. Understand entropy as a measure of uncertainty
4. Derive and apply cross-entropy loss for classification
5. Use KL divergence to measure distribution differences
6. Interpret neural network outputs as probability distributions

3.1 Probability Fundamentals

3.1.1 Random Variables and Distributions

Definition 3.1 (Random Variable). A **random variable** X is a function that maps outcomes from a sample space to real numbers. We distinguish between:

- **Discrete random variables:** Take countable values (e.g., class labels)
- **Continuous random variables:** Take values in continuous ranges

Definition 3.2 (Probability Mass Function (PMF)). For discrete random variable X , the **probability mass function** is:

$$P(X = x) = p(x) \tag{3.1}$$

satisfying: (1) $0 \leq p(x) \leq 1$ for all x , and (2) $\sum_x p(x) = 1$

Example 3.1 (Classification as Discrete Distribution). In image classification with 10 classes (digits 0-9), a neural network outputs a probability distribution using softmax:

$$P(Y = k|\mathbf{x}) = \frac{\exp(z_k)}{\sum_{j=1}^{10} \exp(z_j)} \quad (3.2)$$

For logits $\mathbf{z} = [2.1, 0.5, -1.2, 3.4, 0.8, -0.5, 1.1, -2.0, 0.3, 1.8]$, the model predicts class 3 with highest probability $\approx 68.9\%$.

3.1.2 Conditional Probability and Bayes' Theorem

Definition 3.3 (Conditional Probability). The probability of event A given event B :

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad \text{if } P(B) > 0 \quad (3.3)$$

Theorem 3.1 (Bayes' Theorem). For events A and B with $P(B) > 0$:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \quad (3.4)$$

where $P(A|B)$ is the posterior, $P(B|A)$ is the likelihood, $P(A)$ is the prior, and $P(B)$ is the evidence.

3.2 Information Theory

3.2.1 Entropy

Definition 3.4 (Shannon Entropy). For discrete random variable X with PMF $p(x)$:

$$H(X) = - \sum_x p(x) \ln p(x) = \mathbb{E}[-\ln P(X)] \quad (3.5)$$

Entropy measures average uncertainty. Higher entropy means more uncertainty.

Example 3.2 (Computing Entropy). **Fair coin:** $P(\text{heads}) = P(\text{tails}) = 0.5$

$$H = -[0.5 \log_2(0.5) + 0.5 \log_2(0.5)] = 1 \text{ bit (maximum)} \quad (3.6)$$

Biased coin: $P(\text{heads}) = 0.9, P(\text{tails}) = 0.1$

$$H \approx 0.469 \text{ bits (lower, more predictable)} \quad (3.7)$$

3.2.2 Cross-Entropy

Definition 3.5 (Cross-Entropy). For true distribution p and predicted distribution q :

$$H(p, q) = - \sum_x p(x) \log q(x) = \mathbb{E}_{x \sim p}[-\log q(x)] \quad (3.8)$$

Theorem 3.2 (Cross-Entropy Loss for Classification). For true label y and predicted probabilities $\hat{\mathbf{p}}$:

$$L = -\log \hat{p}_y \quad (3.9)$$

Example 3.3 (Cross-Entropy Loss Calculation). For 3-class classification with true label $y = 2$:

- Predicted: $\hat{\mathbf{p}} = [0.2, 0.6, 0.2] \Rightarrow L = -\log(0.6) \approx 0.511$
- More confident: $\hat{\mathbf{p}} = [0.1, 0.8, 0.1] \Rightarrow L = -\log(0.8) \approx 0.223$ (better)
- Wrong prediction: $\hat{\mathbf{p}} = [0.7, 0.2, 0.1] \Rightarrow L = -\log(0.2) \approx 1.609$ (bad)

Implementation:

PyTorch cross-entropy loss:

```
import torch
import torch.nn as nn

# Logits: shape (batch_size, num_classes)
logits = torch.tensor([[2.0, 1.0, 0.1],
                       [0.5, 2.5, 1.0]])
labels = torch.tensor([0, 1])

# CrossEntropyLoss applies softmax internally
criterion = nn.CrossEntropyLoss()
loss = criterion(logits, labels)
print(f"Loss: {loss.item():.4f}")
```

3.2.3 Kullback-Leibler Divergence

Definition 3.6 (KL Divergence). The KL divergence from distribution q to p :

$$D_{\text{KL}}(p||q) = \sum_x p(x) \log \frac{p(x)}{q(x)} = H(p, q) - H(p) \quad (3.10)$$

Properties: (1) $D_{\text{KL}}(p||q) \geq 0$ with equality iff $p = q$, (2) Not symmetric: $D_{\text{KL}}(p||q) \neq D_{\text{KL}}(q||p)$

Key Point 3.1. *Minimizing KL divergence is equivalent to minimizing cross-entropy when p is fixed. Training neural networks with cross-entropy loss is maximum likelihood estimation.*

3.3 Cross-Entropy Loss: Computational and Memory Analysis

While cross-entropy loss appears simple mathematically, its implementation in large-scale language models presents significant computational and memory challenges. Understanding these practical considerations is essential for training modern transformers efficiently.

3.3.1 Memory Requirements for Logits

The memory footprint of cross-entropy computation is dominated by the logits tensor before softmax normalization. For a language model with batch size B , sequence length n , and vocabulary size V , the logits tensor has shape $B \times n \times V$ and requires $4BnV$ bytes in FP32 (or $2BnV$ bytes in FP16). This memory requirement becomes the primary bottleneck for models with large vocabularies.

Consider BERT-base with vocabulary size $V = 30,000$. For a batch of $B = 32$ sequences with length $n = 512$, the logits tensor requires $32 \times 512 \times 30,000 \times 4 = 1,966,080,000$ bytes, or approximately 1.83 GB of GPU memory just for the unnormalized scores. This is before computing the softmax normalization, which requires storing both the exponentials and the normalization constants. The total memory for the forward pass of cross-entropy loss approaches 3.7 GB for this single operation, consuming a substantial fraction of an NVIDIA V100's 16 GB memory or nearly half of a consumer RTX 3090's 24 GB.

The situation becomes more severe for larger models. GPT-2 with vocabulary size $V = 50,257$ requires 3.1 GB for logits alone with the same batch configuration. GPT-3 and modern large language models often use vocabularies exceeding 50,000 tokens, making the logits tensor one of the largest single memory allocations during training. This explains why vocabulary size is a critical hyperparameter that directly impacts the maximum batch size and sequence length that can fit in GPU memory.

The memory pressure is particularly acute during the backward pass. The gradient of the cross-entropy loss with respect to the logits has the same shape $B \times n \times V$ as the forward logits, effectively doubling the memory requirement. Additionally, the softmax operation requires storing its output for the backward pass, adding another $4BnV$ bytes. The total memory for cross-entropy loss computation (forward and backward) approaches $12BnV$ bytes, or approximately 11 GB for BERT-base with the configuration above.

3.3.2 Computational Cost of Softmax

The softmax operation itself has computational complexity $O(BnV)$ for the forward pass, requiring one exponential operation per logit and a sum reduction over the vocabulary dimension. Modern GPUs can compute exponentials efficiently, but the memory bandwidth required to read and write the large logits tensor often becomes the bottleneck rather than the arithmetic operations themselves.

For BERT-base processing a batch of 32 sequences with 512 tokens each, the softmax operation must compute $32 \times 512 \times 30,000 = 491,520,000$ exponentials. On an NVIDIA A100 GPU with peak FP32 throughput of 19.5 TFLOPS, the arithmetic alone would take approximately 25 microseconds if perfectly parallelized. However, the actual runtime is dominated by memory bandwidth: reading 1.83 GB of logits and writing 1.83 GB of probabilities requires approximately 2.4 milliseconds at the A100's memory bandwidth of 1.5 TB/s, making the operation $100\times$ slower than the arithmetic-limited case. This memory bandwidth bottleneck is characteristic of softmax and explains why vocabulary size has such a direct impact on training speed.

The backward pass of softmax requires computing the Jacobian matrix, but due to the structure of the softmax function, this can be done efficiently without materializing the full $V \times V$ Jacobian. The gradient computation has the same $O(BnV)$ complexity as the forward pass and similar memory bandwidth requirements. The total time for forward and backward softmax over large vocabularies

typically accounts for 15-25% of the total training time per batch, making it a significant optimization target.

3.3.3 Optimizations for Large Vocabularies

Several techniques have been developed to address the memory and computational costs of cross-entropy loss with large vocabularies. These optimizations are essential for training modern language models efficiently.

Sampled Softmax reduces computational cost by approximating the full softmax using only a subset of the vocabulary. During training, instead of computing the normalization constant over all V tokens, we sample a small set of K negative examples (typically $K \approx 1,000$ to $10,000$) and compute softmax over only the true token plus the K samples. This reduces the per-token computational cost from $O(V)$ to $O(K)$, providing speedups of $5\text{-}50\times$ for large vocabularies. The sampling is typically done using a proposal distribution that approximates the unigram frequency of tokens, ensuring that common words are sampled more frequently than rare words.

The memory savings are equally dramatic. With sampled softmax, the logits tensor has shape $B \times n \times (K + 1)$ instead of $B \times n \times V$, reducing memory from 1.83 GB to approximately 61 MB for BERT-base with $K = 1,000$ samples. This $30\times$ memory reduction enables much larger batch sizes or longer sequences, directly improving training efficiency. However, sampled softmax introduces bias in the gradient estimates, which can slow convergence. In practice, it is most effective during the early stages of training and is often replaced with full softmax for final fine-tuning.

Adaptive Softmax exploits the Zipfian distribution of natural language, where a small number of frequent words account for the majority of token occurrences. The vocabulary is partitioned into clusters based on frequency: frequent words in a small head cluster, and rare words in larger tail clusters. The model first predicts which cluster the token belongs to, then predicts the specific token within that cluster. This hierarchical approach reduces the effective vocabulary size for most predictions from V to approximately \sqrt{V} , providing substantial speedups for very large vocabularies.

For a vocabulary of $V = 100,000$ tokens partitioned into clusters of sizes $[1,000, 9,000, 90,000]$, the average computational cost per token is approximately $3 + 1,000 = 1,003$ operations (cluster prediction plus within-cluster prediction for the head cluster) compared to 100,000 for full softmax. This represents a $100\times$ speedup for the most common tokens. The memory requirements are similarly reduced, as the model only needs to store logits for the predicted cluster rather than the full vocabulary. Adaptive softmax is particularly effective for language modeling tasks where the token distribution is highly skewed, and it has been successfully used in models like Transformer-XL and adaptive input representations.

Vocabulary Pruning and Subword Tokenization address the problem at its source by reducing the vocabulary size V itself. Subword tokenization methods like Byte-Pair Encoding (BPE) and WordPiece decompose rare words into common subword units, allowing models to use vocabularies of 30,000-50,000 tokens instead of 100,000+ word-level tokens while maintaining good coverage. This directly reduces both memory and computation by a factor of $2\text{-}3\times$. Modern models like GPT-3 and BERT use vocabularies of approximately 50,000 tokens, carefully chosen to balance vocabulary size against the increased sequence length from subword splitting.

The choice of vocabulary size involves a fundamental trade-off: smaller vocabularies reduce memory and computation per token but increase sequence length (more tokens per sentence), while larger vocabularies reduce sequence length but increase per-token costs. For transformer models where attention has $O(n^2)$ complexity, reducing sequence length by using larger vocabularies can actually improve overall efficiency despite the higher per-token costs. Empirical studies suggest that vocabularies of 32,000-50,000 tokens provide a good balance for most languages, though this varies with the specific language and domain.

3.4 KL Divergence in Practice

KL divergence appears throughout modern deep learning as a measure of distribution similarity. Understanding its computational properties and applications is essential for implementing techniques like

variational autoencoders, knowledge distillation, and reinforcement learning from human feedback.

3.4.1 Applications in Modern Deep Learning

Variational Autoencoders (VAEs) use KL divergence as a regularization term to ensure that the learned latent distribution $q(z|x)$ remains close to a prior distribution $p(z)$, typically a standard Gaussian. The VAE loss function combines reconstruction loss with a KL divergence term:

$$\mathcal{L}_{\text{VAE}} = \mathbb{E}_{q(z|x)}[\log p(x|z)] - D_{\text{KL}}(q(z|x)||p(z)) \quad (3.11)$$

For a Gaussian encoder with mean μ and variance σ^2 , the KL divergence to a standard Gaussian has a closed form:

$$D_{\text{KL}}(q||p) = \frac{1}{2} \sum_{i=1}^d (\mu_i^2 + \sigma_i^2 - \log \sigma_i^2 - 1) \quad (3.12)$$

This closed form makes VAEs computationally efficient, as the KL term requires only $O(d)$ operations for a d -dimensional latent space, typically much smaller than the reconstruction loss computation. For a VAE with 512-dimensional latent space, the KL divergence computation takes less than 1 microsecond on a modern GPU, making it negligible compared to the encoder and decoder networks.

Knowledge Distillation transfers knowledge from a large teacher model to a smaller student model by minimizing the KL divergence between their output distributions. The student is trained to match not just the hard labels but the full probability distribution produced by the teacher:

$$\mathcal{L}_{\text{distill}} = \alpha \mathcal{L}_{\text{CE}}(y, \hat{y}_{\text{student}}) + (1 - \alpha) T^2 D_{\text{KL}}(\hat{y}_{\text{teacher}}||\hat{y}_{\text{student}}) \quad (3.13)$$

where T is a temperature parameter that softens the distributions. The KL divergence term encourages the student to learn the relative confidences between classes that the teacher has learned, not just the most likely class. This is particularly valuable when the teacher assigns non-negligible probability to multiple classes, indicating genuine ambiguity or similarity between categories.

The computational cost of knowledge distillation is dominated by running both teacher and student models, with the KL divergence computation itself being relatively cheap at $O(BnV)$ for batch size B , sequence length n , and vocabulary size V . For BERT-base distillation with vocabulary size 30,000, computing the KL divergence over a batch of 32 sequences with 512 tokens requires approximately 1.5 milliseconds on an A100 GPU, compared to 50-100 milliseconds for the forward passes of the teacher and student models. The memory overhead is also modest, requiring storage of both teacher and student logits but no additional activations for backpropagation through the teacher.

Reinforcement Learning from Human Feedback (RLHF) uses KL divergence to constrain the policy learned through reinforcement learning to remain close to the original supervised fine-tuned model. This prevents the model from exploiting the reward model by generating adversarial outputs that score highly but are nonsensical. The RLHF objective includes a KL penalty term:

$$\mathcal{L}_{\text{RLHF}} = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}}[r(x, y)] - \beta D_{\text{KL}}(\pi_{\theta}||\pi_{\text{ref}}) \quad (3.14)$$

where π_{θ} is the policy being optimized, π_{ref} is the reference model, and β controls the strength of the KL constraint. Computing this KL divergence requires running both the policy and reference models on the same inputs and computing the divergence over the vocabulary at each token position. For large language models with vocabularies of 50,000+ tokens and sequences of 1,000+ tokens, this KL computation can consume 10-20% of the total training time, making it a non-negligible cost in RLHF training pipelines.

3.4.2 Numerical Stability Considerations

Computing KL divergence naively can lead to numerical instability due to the logarithm of very small probabilities. When $q(x)$ is close to zero, $\log q(x)$ approaches negative infinity, and the product $p(x) \log q(x)$ can produce NaN values or catastrophic cancellation. Similarly, when computing $\log(p(x)/q(x))$, direct division can lose precision for very small probabilities.

The numerically stable approach computes KL divergence in log-space using the log-sum-exp trick. Instead of computing probabilities via softmax and then taking logarithms, we work directly with log-probabilities:

$$D_{\text{KL}}(p||q) = \sum_x p(x)(\log p(x) - \log q(x)) = \sum_x \exp(\log p(x)) \cdot (\log p(x) - \log q(x)) \quad (3.15)$$

This formulation avoids computing very small probabilities explicitly. Modern deep learning frameworks like PyTorch provide `F.kl_div` that operates on log-probabilities directly, ensuring numerical stability even when probabilities span many orders of magnitude.

Another source of instability arises when $p(x) > 0$ but $q(x) = 0$, which makes the KL divergence infinite. In practice, this occurs when the model assigns zero probability to an event that actually occurs in the data. To prevent this, implementations typically add a small epsilon ($\epsilon \approx 10^{-8}$) to probabilities before computing logarithms, or use label smoothing to ensure that the target distribution p never assigns exactly zero probability to any class. Label smoothing replaces hard targets with a mixture of the true label and a uniform distribution:

$$p_{\text{smooth}}(x) = (1 - \epsilon)p_{\text{true}}(x) + \epsilon/V \quad (3.16)$$

where $\epsilon \approx 0.1$ is typical. This not only improves numerical stability but also acts as a regularizer that prevents overconfident predictions and often improves generalization.

3.5 Hardware Implications of Softmax and Large Vocabularies

The softmax operation and large vocabulary sizes have profound implications for hardware utilization and training efficiency. Understanding these hardware-level considerations is essential for optimizing transformer training and deployment.

3.5.1 Softmax Computation on GPUs

Softmax is a memory-bandwidth-bound operation rather than compute-bound on modern GPUs. The operation requires reading the input logits, computing exponentials, summing the exponentials, and writing the normalized probabilities. For a vocabulary of size V , this involves reading and writing $2V$ values while performing only $2V$ arithmetic operations (one exponential and one division per element). On an NVIDIA A100 GPU with 312 TFLOPS of FP16 compute but only 1.5 TB/s of memory bandwidth, the arithmetic could theoretically complete in nanoseconds, but the memory transfers take microseconds.

Consider computing softmax over a vocabulary of $V = 50,000$ tokens for a single position. Reading 50,000 FP32 values (200 KB) and writing 50,000 FP32 values (200 KB) requires 400 KB of memory bandwidth. At the A100's bandwidth of 1.5 TB/s, this takes approximately 0.27 microseconds. The arithmetic operations (50,000 exponentials and 50,000 divisions) would take approximately 0.003 microseconds at peak throughput, making the operation $90\times$ memory-bandwidth-limited. This ratio worsens for smaller vocabularies and improves slightly for larger ones, but softmax remains fundamentally bandwidth-bound across all practical vocabulary sizes.

The memory-bandwidth-bound nature of softmax has several implications. First, reducing precision from FP32 to FP16 provides nearly $2\times$ speedup by halving the memory traffic, with minimal impact on accuracy for most applications. Second, fusing the softmax operation with subsequent operations (like the cross-entropy loss computation) can eliminate intermediate memory traffic by keeping values in registers, providing additional speedups of $1.5\text{--}2\times$. Modern deep learning frameworks implement fused softmax-cross-entropy kernels that compute both operations in a single GPU kernel, reducing memory traffic from $4V$ to $2V$ values per position.

Third, the memory-bandwidth bottleneck means that vocabulary size has a nearly linear impact on softmax runtime. Doubling the vocabulary from 25,000 to 50,000 tokens approximately doubles the softmax computation time, as the memory traffic doubles while the arithmetic remains negligible. This linear scaling makes vocabulary size one of the most direct levers for controlling training speed, and it explains why careful vocabulary selection is critical for efficient training.

3.5.2 Memory Bandwidth and Large Vocabularies

The memory bandwidth requirements of large vocabularies extend beyond just the softmax operation to the entire forward and backward pass of the language model head. The final linear layer that projects from the model dimension d to the vocabulary size V has weight matrix $\mathbf{W} \in \mathbb{R}^{d \times V}$, which must be read from memory for every forward pass and written during every backward pass.

For BERT-base with $d = 768$ and $V = 30,000$, this weight matrix contains $768 \times 30,000 = 23,040,000$ parameters, requiring 92 MB in FP32 or 46 MB in FP16. For a batch of $B = 32$ sequences with $n = 512$ tokens, the forward pass must read this 92 MB matrix once and perform $32 \times 512 \times 768 \times 30,000 = 377$ billion multiply-accumulate operations. On an A100 GPU, reading 92 MB takes approximately 61 microseconds at 1.5 TB/s bandwidth, while the arithmetic takes approximately 600 microseconds at 312 TFLOPS FP16 throughput. In this case, the operation is compute-bound rather than bandwidth-bound, but the memory bandwidth still accounts for approximately 10% of the total time.

The situation changes dramatically for smaller batch sizes or shorter sequences. With batch size $B = 1$ and sequence length $n = 1$ (as in autoregressive generation), the forward pass performs only $768 \times 30,000 = 23$ million operations, taking 0.07 microseconds at peak throughput. The memory bandwidth to read the weight matrix remains 61 microseconds, making the operation $870\times$ bandwidth-bound. This explains why autoregressive generation is so much slower than parallel training: the small batch size prevents amortizing the memory bandwidth cost over many operations, and the model spends most of its time waiting for memory rather than computing.

The backward pass has similar memory bandwidth requirements but must also write the gradient of the weight matrix, doubling the memory traffic. For the BERT-base example, the backward pass requires reading 92 MB (weights) and writing 92 MB (gradients), totaling 184 MB of memory traffic. This takes approximately 122 microseconds, compared to approximately 600 microseconds for the arithmetic, making the backward pass approximately 17% bandwidth-bound. The total memory bandwidth for the forward and backward pass of the language model head is 276 MB per batch, which accumulates to significant overhead over thousands of training steps.

3.5.3 Why Vocabulary Size Impacts Training Speed

Vocabulary size impacts training speed through three primary mechanisms: memory capacity, memory bandwidth, and arithmetic operations. Understanding the relative contribution of each mechanism helps guide optimization strategies.

Memory Capacity: As discussed in Section 3.3, the logits tensor requires $4BnV$ bytes, which directly limits the maximum batch size and sequence length that fit in GPU memory. For BERT-base with $V = 30,000$, reducing the vocabulary to $V = 15,000$ would halve the logits memory from 1.83 GB to 915 MB, allowing a $2\times$ larger batch size or $1.4\times$ longer sequences (since attention memory scales as $O(n^2)$). Larger batch sizes improve GPU utilization and reduce the number of training steps required for convergence, directly improving training efficiency.

The memory capacity constraint is particularly severe for large language models. GPT-3 with 175 billion parameters requires approximately 700 GB of memory just for the model weights and optimizer states in FP32 (or 350 GB in FP16 with mixed precision). Adding 1.83 GB for logits might seem negligible, but when training across multiple GPUs with model parallelism, the logits must be replicated on each GPU that computes the language model head, multiplying the memory cost. For a model parallelized across 8 GPUs, the logits consume 14.6 GB of total memory, which becomes significant relative to the per-GPU memory budget.

Memory Bandwidth: The softmax operation and language model head are bandwidth-bound, as discussed above. Reducing vocabulary size from 50,000 to 25,000 tokens reduces the memory traffic for softmax by $2\times$, directly improving runtime by approximately $2\times$ for this operation. Since softmax and the language model head together account for 20-30% of total training time, this translates to an overall speedup of approximately $1.15\text{-}1.2\times$ for the entire training pipeline. This speedup is achieved without any loss in model quality, making vocabulary reduction through subword tokenization one of the most effective optimizations for language model training.

Arithmetic Operations: The language model head performs $BndV$ multiply-accumulate operations, which scales linearly with vocabulary size. However, as discussed above, this operation is typically

compute-bound only for large batch sizes, and the arithmetic cost is often dominated by memory bandwidth. For very large vocabularies ($V > 100,000$) and large batch sizes, the arithmetic can become significant, but for typical configurations with $V \approx 30,000$ -50,000, memory bandwidth is the primary bottleneck.

The combined effect of these three mechanisms means that vocabulary size has a superlinear impact on training speed. Doubling the vocabulary size reduces the maximum batch size (due to memory capacity), increases the per-batch runtime (due to memory bandwidth and arithmetic), and may require more training steps (due to the smaller batch size). Empirically, doubling the vocabulary from 25,000 to 50,000 tokens typically increases total training time by 1.5 - $2\times$, making vocabulary selection a critical hyperparameter for efficient training.

3.5.4 Optimization Techniques

Several techniques mitigate the hardware costs of large vocabularies. **Vocabulary pruning** removes rare tokens that appear fewer than a threshold number of times in the training data, reducing vocabulary size without significantly impacting coverage. For example, removing tokens that appear fewer than 100 times in a large corpus might reduce vocabulary from 50,000 to 35,000 tokens while affecting less than 0.1% of tokens in the data. This 30% reduction in vocabulary size provides approximately $1.2\times$ training speedup with negligible impact on model quality.

Subword tokenization methods like BPE and WordPiece achieve smaller vocabularies by decomposing rare words into common subword units. This allows models to handle unlimited vocabulary with a fixed-size token set, typically 30,000-50,000 tokens. The trade-off is increased sequence length, as rare words are split into multiple tokens. However, for transformer models where attention has $O(n^2)$ complexity, the increased sequence length is often offset by the reduced per-token costs, resulting in net speedups of 1.3 - $1.5\times$ compared to word-level tokenization.

Adaptive softmax and **sampled softmax**, discussed in Section 3.3, provide algorithmic approaches to reducing the computational cost of large vocabularies. These techniques are particularly effective for very large vocabularies ($V > 100,000$) where the vocabulary size dominates training time. For typical transformer models with vocabularies of 30,000-50,000 tokens, the simpler approaches of vocabulary pruning and subword tokenization are often sufficient and easier to implement.

Model parallelism distributes the vocabulary across multiple GPUs, allowing each GPU to compute softmax over only a subset of the vocabulary. For a vocabulary of 50,000 tokens distributed across 4 GPUs, each GPU computes softmax over 12,500 tokens, reducing the per-GPU memory and computation by $4\times$. However, this requires an all-reduce operation to compute the global normalization constant, which adds communication overhead. Model parallelism for the vocabulary is most effective for very large vocabularies and when training on high-bandwidth interconnects like NVLink or InfiniBand.

The choice of optimization technique depends on the specific model and training configuration. For most applications, subword tokenization with vocabularies of 30,000-50,000 tokens provides a good balance between vocabulary size and sequence length. For very large models or very large vocabularies, adaptive softmax or model parallelism may be necessary to achieve acceptable training speeds.

3.6 Exercises

Exercise 3.1. A neural network outputs $\hat{\mathbf{p}} = [0.15, 0.60, 0.20, 0.05]$ for 4 classes. Compute: (1) entropy $H(\hat{\mathbf{p}})$, (2) cross-entropy loss if true label is class 2, (3) optimal output distribution.

Exercise 3.2. Show that $H(p, q) = H(p) + D_{\text{KL}}(p||q)$, proving cross-entropy minimization equals KL divergence minimization when p is fixed.

Exercise 3.3. For binary classifier with $\hat{p} = 0.8$ and true label class 1: (1) Compute binary cross-entropy loss, (2) Find $\frac{\partial L}{\partial \hat{p}}$, (3) Compare loss for $\hat{p} \in \{0.99, 0.2\}$.

Exercise 3.4. Calculate the memory requirements for storing logits in a GPT-2 model with vocabulary size $V = 50,257$, batch size $B = 16$, and sequence length $n = 1024$. How much memory is saved by using FP16 instead of FP32? If you have an NVIDIA A100 with 40 GB of memory, what is the maximum batch size you can use if logits consume at most 25% of available memory?

Exercise 3.5. For sampled softmax with $K = 5,000$ negative samples and vocabulary size $V = 100,000$: (1) Calculate the speedup factor for the forward pass compared to full softmax, (2) Compute the memory reduction for a batch of 32 sequences with 512 tokens each, (3) Discuss why sampled softmax introduces bias in gradient estimates.

Exercise 3.6. An NVIDIA A100 GPU has memory bandwidth of 1.5 TB/s and FP16 compute throughput of 312 TFLOPS. For softmax over a vocabulary of $V = 30,000$ tokens: (1) Calculate the time to read and write the logits and probabilities (400 KB total), (2) Calculate the time to compute 30,000 exponentials and divisions at peak throughput, (3) Determine whether the operation is compute-bound or bandwidth-bound and by what factor.

Exercise 3.7. In knowledge distillation, the KL divergence loss is scaled by T^2 where T is the temperature parameter. Explain why this scaling is necessary by: (1) Showing how temperature affects the magnitude of gradients, (2) Deriving the gradient of $D_{\text{KL}}(\text{softmax}(\mathbf{z}/T) \parallel \text{softmax}(\mathbf{z}'/T))$ with respect to \mathbf{z}' , (3) Demonstrating that without T^2 scaling, the distillation loss would vanish as $T \rightarrow \infty$.

Part II

Neural Network Fundamentals

Chapter 4

Feed-Forward Neural Networks

Chapter Overview

Feed-forward neural networks are the foundation of deep learning. These networks transform inputs through sequences of linear and nonlinear operations to produce outputs. This chapter develops the architecture, training, and theory of feed-forward networks, establishing concepts that extend to all modern deep learning models including transformers.

Learning Objectives

After completing this chapter, you will be able to:

1. Understand the architecture of feed-forward neural networks
2. Implement forward and backward passes through MLPs
3. Apply appropriate activation functions and understand their properties
4. Initialize network weights properly to enable training
5. Apply regularization techniques to prevent overfitting
6. Understand the universal approximation theorem

4.1 From Linear Models to Neural Networks

4.1.1 The Perceptron

Definition 4.1 (Perceptron). The perceptron is a binary classifier:

$$\hat{y} = \text{sign}(\mathbf{w}^\top \mathbf{x} + b) = \begin{cases} +1 & \text{if } \mathbf{w}^\top \mathbf{x} + b > 0 \\ -1 & \text{otherwise} \end{cases} \quad (4.1)$$

where $\mathbf{w} \in \mathbb{R}^n$ are weights, $b \in \mathbb{R}$ is bias, $\mathbf{x} \in \mathbb{R}^n$ is input.

4.1.2 Multi-Class Classification: Softmax Regression

Definition 4.2 (Softmax Function). For logits $\mathbf{z} = [z_1, \dots, z_C]^\top \in \mathbb{R}^C$:

$$\text{softmax}(\mathbf{z})_k = \frac{\exp(z_k)}{\sum_{j=1}^C \exp(z_j)} \quad (4.2)$$

Example 4.1 (Softmax Computation). For logits $\mathbf{z} = [2.0, 1.0, 0.1]$: Sum of exponentials = 11.212, giving probabilities $[0.659, 0.242, 0.099]$. The model predicts class 1 with 65.9 percent confidence.

4.2 Multi-Layer Perceptrons

Definition 4.3 (Multi-Layer Perceptron). An L-layer MLP transforms input through layers:

$$\mathbf{z}^{(\ell)} = \mathbf{W}^{(\ell)} \mathbf{h}^{(\ell-1)} + \mathbf{b}^{(\ell)} \quad (4.3)$$

$$\mathbf{h}^{(\ell)} = \sigma^{(\ell)}(\mathbf{z}^{(\ell)}) \quad (4.4)$$

where $\mathbf{W}^{(\ell)} \in \mathbb{R}^{n_\ell \times n_{\ell-1}}$ is the weight matrix and $\sigma^{(\ell)}$ is the activation function.

Example 4.2 (3-Layer MLP for MNIST). Architecture for MNIST digit classification:

- Input: $\mathbf{x} \in \mathbb{R}^{784}$ (flattened 28×28 image)
- Hidden 1: $\mathbf{h}^{(1)} \in \mathbb{R}^{256}$ with ReLU
- Hidden 2: $\mathbf{h}^{(2)} \in \mathbb{R}^{128}$ with ReLU
- Output: $\mathbf{z}^{(3)} \in \mathbb{R}^{10}$ with softmax

Parameter count: $200,960 + 32,896 + 1,290 = 235,146$ parameters.

4.2.1 Why Depth Matters

Without nonlinear activations, multiple layers collapse to single linear transformation. With nonlinearities, deep networks learn complex functions efficiently.

4.3 Memory and Computation Analysis

Understanding the memory and computational requirements of feed-forward networks is essential for training large models efficiently. The relationship between parameter count, floating-point operations (FLOPs), and memory usage determines the practical limits of model size and batch size on available hardware.

4.3.1 Parameter Count vs FLOPs

The parameter count of a neural network determines its memory footprint for storing weights, while the FLOPs (floating-point operations) determine the computational cost of forward and backward passes. These two quantities scale differently with network architecture, leading to important trade-offs in model design.

For a single fully-connected layer computing $\mathbf{y} = \mathbf{W}\mathbf{x} + \mathbf{b}$ where $\mathbf{W} \in \mathbb{R}^{m \times n}$, the parameter count is $mn + m$ (weights plus biases). The forward pass requires mn multiply-accumulate operations for the matrix-vector product plus m additions for the bias, totaling approximately $2mn$ FLOPs. The backward pass requires computing gradients with respect to inputs ($\nabla_{\mathbf{x}}L = \mathbf{W}^\top \nabla_{\mathbf{y}}L$, requiring $2mn$ FLOPs), gradients with respect to weights ($\nabla_{\mathbf{W}}L = \nabla_{\mathbf{y}}L\mathbf{x}^\top$, requiring $2mn$ FLOPs), and gradients with respect to biases ($\nabla_{\mathbf{b}}L = \nabla_{\mathbf{y}}L$, requiring m FLOPs). The total computational cost for forward and backward passes is approximately $6mn$ FLOPs, or $3\times$ the parameter count.

This $3\times$ ratio between FLOPs and parameters holds approximately for fully-connected layers and provides a useful rule of thumb: training a model for one step requires approximately $6\times$ as many FLOPs as the model has parameters ($2\times$ for forward pass, $4\times$ for backward pass including gradient computation). For a model with 100 million parameters, one training step requires approximately 600 million FLOPs, or 0.6 GFLOPs. At 1,000 training steps, this totals 600 GFLOPs of computation.

However, this ratio varies significantly with architecture. Convolutional layers have much higher FLOPs per parameter due to weight sharing: a 3×3 convolutional filter with C_{in} input channels and C_{out} output channels has $9C_{\text{in}}C_{\text{out}}$ parameters but requires $9C_{\text{in}}C_{\text{out}}HW$ FLOPs for an $H \times W$ feature map, giving a FLOPs-to-parameter ratio of HW . For a 224×224 image, this ratio is 50,176, making convolutional layers far more compute-intensive per parameter than fully-connected layers. Conversely, embedding layers have zero FLOPs (they perform table lookups rather than arithmetic) despite having many parameters, making them memory-intensive but computationally cheap.

4.3.2 Memory Requirements for Activations

During training, neural networks must store intermediate activations for use in the backward pass, and these activations often consume more memory than the model parameters themselves. Understanding activation memory is critical for determining maximum batch size and sequence length.

For a feed-forward layer computing $\mathbf{h} = \sigma(\mathbf{W}\mathbf{x} + \mathbf{b})$ with batch size B , the network must store the input activations $\mathbf{x} \in \mathbb{R}^{B \times n}$, the pre-activation values $\mathbf{z} = \mathbf{W}\mathbf{x} + \mathbf{b} \in \mathbb{R}^{B \times m}$, and the post-activation values $\mathbf{h} \in \mathbb{R}^{B \times m}$. In FP32, this requires $4B(n + 2m)$ bytes of memory. For a typical transformer feed-forward layer with $n = 768$ (model dimension) and $m = 3072$ (intermediate dimension), processing batch size $B = 32$ requires $4 \times 32 \times (768 + 2 \times 3072) = 901,120$ bytes, or approximately 0.86 MB per layer. For a 12-layer BERT-base model, activation memory totals approximately 10.3 MB per batch, which is modest compared to the 440 MB required for model parameters.

However, activation memory scales linearly with batch size while parameter memory remains constant. Increasing batch size from 32 to 256 increases activation memory by $8\times$, from 10.3 MB to 82.4 MB, while parameter memory remains 440 MB. For very large batch sizes, activation memory can exceed parameter memory. At batch size 1024, activation memory for BERT-base reaches 329.6 MB, approaching the parameter memory. This scaling explains why large batch sizes eventually become memory-limited: the activations grow without bound while parameters remain fixed.

The situation is more severe for transformer models due to attention mechanisms. Self-attention requires storing attention score matrices of size $B \times h \times n \times n$ where h is the number of attention heads and n is the sequence length. For BERT-base with $h = 12$ heads, batch size $B = 32$, and sequence length $n = 512$, the attention scores require $4 \times 32 \times 12 \times 512 \times 512 = 402,653,184$ bytes, or approximately 384 MB per layer. Across 12 layers, attention scores alone consume 4.6 GB of memory, dwarfing both the parameter memory (440 MB) and the feed-forward activation memory (10.3 MB). This explains why sequence length has such a dramatic impact on memory usage: doubling the sequence length quadruples the attention memory due to the $O(n^2)$ scaling.

4.3.3 GPU Utilization for Different Layer Sizes

GPU utilization—the fraction of peak computational throughput actually achieved—varies dramatically with layer dimensions and batch size. Understanding these utilization patterns is essential for designing efficient architectures and selecting appropriate hyperparameters.

Modern GPUs achieve peak performance on large matrix multiplications where dimensions are multiples of the GPU's tile size (typically 16 or 32 for FP16 operations). For an NVIDIA A100 GPU with peak FP16 throughput of 312 TFLOPS, a matrix multiplication $\mathbf{C} = \mathbf{A}\mathbf{B}$ where $\mathbf{A} \in \mathbb{R}^{m \times k}$ and

$\mathbf{B} \in \mathbb{R}^{k \times n}$ achieves near-peak performance when m , k , and n are all large (greater than 1024) and multiples of 16. Under these conditions, the GPU can achieve 280-300 TFLOPS, or 90-95% of peak throughput.

However, for smaller dimensions, utilization drops dramatically. A matrix multiplication with $m = 32$, $k = 768$, $n = 768$ (corresponding to batch size 32 and BERT-base dimensions) requires $2 \times 32 \times 768 \times 768 = 37,748,736$ FLOPs. At peak throughput, this would take 0.12 microseconds, but the actual runtime is approximately 15 microseconds, indicating only 0.8% utilization. The poor utilization arises because the small batch dimension ($m = 32$) provides insufficient parallelism to saturate the GPU's 6,912 CUDA cores. Each CUDA core can process one operation per clock cycle, so saturating the GPU requires at least 6,912 concurrent operations. With $m = 32$, only 32 rows can be processed in parallel, leaving 99.5% of the GPU idle.

Increasing batch size directly improves GPU utilization. With batch size 256, the same operation requires $2 \times 256 \times 768 \times 768 = 301,989,888$ FLOPs, taking approximately 50 microseconds for actual runtime. This corresponds to 6.0 TFLOPS, or 1.9% of peak throughput—still poor, but $2.4\times$ better than batch size 32. At batch size 2048, the operation achieves approximately 45 TFLOPS, or 14.4% of peak throughput. Full utilization (90%+) requires batch sizes exceeding 8192 for these dimensions, which is impractical for most training scenarios due to memory constraints and optimization difficulties with very large batches.

The feed-forward layers in transformers achieve better utilization than attention layers due to their larger intermediate dimension. For BERT-base, the first feed-forward layer computes $\mathbf{W}_1 \mathbf{h}$ where $\mathbf{W}_1 \in \mathbb{R}^{3072 \times 768}$ and $\mathbf{h} \in \mathbb{R}^{B \times 768}$. With batch size 32, this requires $2 \times 32 \times 768 \times 3072 = 150,994,944$ FLOPs, taking approximately 25 microseconds for 6.0 TFLOPS throughput (1.9% utilization). The larger output dimension (3072 vs 768) provides more parallelism, but utilization remains poor due to the small batch size. At batch size 256, the feed-forward layer achieves approximately 60 TFLOPS (19.2% utilization), and at batch size 2048, it reaches approximately 180 TFLOPS (57.7% utilization). These higher utilization rates explain why feed-forward layers account for a larger fraction of training time than their FLOPs would suggest: they achieve better hardware efficiency than attention layers.

4.3.4 Batch Size Impact on Efficiency

Batch size is the primary lever for controlling GPU utilization and training efficiency. Larger batches amortize the fixed costs of launching GPU kernels, loading weights from memory, and synchronizing across devices, leading to higher throughput measured in samples per second. However, larger batches also require more memory and may necessitate adjustments to learning rate and training schedule.

For BERT-base training on an NVIDIA A100 GPU, the relationship between batch size and throughput is approximately logarithmic: doubling the batch size increases throughput by $1.5\text{-}1.7\times$ rather than $2\times$. With batch size 8, BERT-base achieves approximately 120 samples per second. At batch size 16, throughput increases to 200 samples per second ($1.67\times$ improvement). At batch size 32, throughput reaches 320 samples per second ($1.6\times$ improvement). At batch size 64, throughput reaches 480 samples per second ($1.5\times$ improvement). The diminishing returns arise because larger batches improve GPU utilization but eventually become limited by memory bandwidth rather than compute throughput.

The memory cost of larger batches is substantial. Batch size 8 requires approximately 4.2 GB of GPU memory for BERT-base (including model parameters, optimizer states, and activations). Batch size 16 requires 6.8 GB ($1.62\times$ increase). Batch size 32 requires 12.0 GB ($1.76\times$ increase). Batch size 64 requires 22.6 GB ($1.88\times$ increase). The super-linear scaling of memory with batch size arises because activation memory scales linearly with batch size while parameter and optimizer memory remain constant, and the activation memory eventually dominates. An A100 GPU with 80 GB of memory can accommodate batch size 256 for BERT-base, but larger batches require gradient accumulation or distributed training.

The optimal batch size balances throughput, memory usage, and optimization dynamics. From a hardware efficiency perspective, larger batches are always better, as they improve GPU utilization and samples-per-second throughput. However, from an optimization perspective, very large batches can slow convergence by reducing the number of parameter updates per epoch. Empirically, batch sizes of 256-2048 work well for BERT-base, providing good hardware efficiency (40-60% GPU utilization) while maintaining reasonable convergence speed. Larger batches require careful tuning of learning rate and

warmup schedule to maintain training stability and final model quality.

4.3.5 Transformer Feed-Forward Networks

The feed-forward networks in transformer models follow a specific architecture that differs from traditional MLPs. Each transformer layer contains a two-layer feed-forward network with an expansion factor of 4: the first layer projects from model dimension d to intermediate dimension $4d$, applies an activation function (typically GELU), and the second layer projects back to dimension d . This architecture is used universally in BERT, GPT, T5, and other transformer models.

For BERT-base with $d = 768$, the feed-forward network has dimensions $768 \rightarrow 3072 \rightarrow 768$. The first layer has weight matrix $\mathbf{W}_1 \in \mathbb{R}^{3072 \times 768}$ with 2,359,296 parameters, and the second layer has weight matrix $\mathbf{W}_2 \in \mathbb{R}^{768 \times 3072}$ with 2,359,296 parameters, totaling 4,718,592 parameters per transformer layer. Across 12 layers, the feed-forward networks contain 56,623,104 parameters, or 51.5% of BERT-base's 110 million total parameters. This makes the feed-forward networks the largest component of the model by parameter count, exceeding the attention layers (38.6% of parameters) and embeddings (9.9% of parameters).

The computational cost of the feed-forward network is similarly dominant. For batch size B and sequence length n , the first layer requires $2Bn \times 768 \times 3072$ FLOPs, and the second layer requires $2Bn \times 3072 \times 768$ FLOPs, totaling $4Bn \times 768 \times 3072 = 9,437,184Bn$ FLOPs per transformer layer. For $B = 32$ and $n = 512$, this totals 154,140,098,048 FLOPs per layer, or approximately 154 GFLOPs. Across 12 layers, the feed-forward networks require 1.85 TFLOPs per forward pass, compared to 1.57 TFLOPs for attention layers. The feed-forward networks account for 54.1% of the total computational cost, slightly more than their share of parameters due to the large intermediate dimension.

The memory requirements for feed-forward activations are modest compared to attention. For batch size $B = 32$ and sequence length $n = 512$, the intermediate activations after the first layer have shape $32 \times 512 \times 3072$, requiring $4 \times 32 \times 512 \times 3072 = 201,326,592$ bytes, or approximately 192 MB per layer. Across 12 layers, feed-forward activations total 2.3 GB, which is substantial but less than the 4.6 GB required for attention score matrices. The feed-forward activations scale linearly with sequence length ($O(n)$) rather than quadratically ($O(n^2)$), making them less problematic for long sequences.

The $4\times$ expansion factor used in transformer feed-forward networks is a design choice that balances model capacity, computational cost, and memory usage. Larger expansion factors (e.g., $8\times$ or $16\times$) increase model capacity and can improve performance on some tasks, but they also increase parameter count, FLOPs, and memory proportionally. Smaller expansion factors (e.g., $2\times$) reduce computational cost but may limit model expressiveness. The $4\times$ factor has proven effective across a wide range of tasks and model sizes, from BERT-base ($768 \rightarrow 3072$) to GPT-3 ($12288 \rightarrow 49152$), and has become a standard architectural choice.

4.4 Activation Functions

Definition 4.4 (ReLU).

$$\text{ReLU}(z) = \max(0, z) \quad (4.5)$$

Derivative: $\text{ReLU}'(z) = \mathbb{I}[z > 0]$

Definition 4.5 (GELU). Gaussian Error Linear Unit (default in transformers):

$$\text{GELU}(z) = z \cdot \Phi(z) \quad (4.6)$$

where Φ is standard normal CDF. Approximation:

$$\text{GELU}(z) \approx 0.5z \left(1 + \tanh \left[\sqrt{\frac{2}{\pi}} (z + 0.044715z^3) \right] \right) \quad (4.7)$$

Key Point 4.1. *Transformer models use GELU (BERT, GPT) or variants like Swish for feed-forward networks.*

4.4.1 Computational Cost of Activation Functions

The choice of activation function has direct implications for both computational cost and memory bandwidth utilization. While activation functions appear simple mathematically, their performance characteristics on modern hardware vary significantly, making activation selection an important consideration for efficient neural network training.

ReLU is the most computationally efficient activation function, requiring only a single comparison and conditional assignment per element. On modern GPUs, ReLU can be implemented as a single instruction using the maximum operation: $\text{ReLU}(z) = \max(0, z)$. For a layer with n activations, ReLU requires n comparisons and n conditional moves, totaling approximately $2n$ operations. On an NVIDIA A100 GPU with 312 TFLOPS of FP16 throughput, computing ReLU for a batch of $B = 32$ sequences with $n = 512$ tokens and $d = 768$ dimensions requires $32 \times 512 \times 768 = 12,582,912$ operations, completing in approximately 0.04 microseconds at peak throughput. However, the actual runtime is dominated by memory bandwidth: reading and writing the activation tensor requires $2 \times 32 \times 512 \times 768 \times 2 = 50$ MB of memory traffic, taking approximately 33 microseconds at the A100's 1.5 TB/s bandwidth. This makes ReLU approximately $825\times$ memory-bandwidth-bound rather than compute-bound.

GELU is significantly more expensive computationally than ReLU due to the Gaussian error function $\Phi(z)$, which requires computing the cumulative distribution function of the standard normal distribution. The exact GELU implementation requires evaluating the error function, which typically involves polynomial approximations with 10-15 arithmetic operations per element. The tanh-based approximation shown in Definition 4.5 reduces this to approximately 8 operations per element: one cube, two multiplications, one addition, one square root, one tanh evaluation (itself requiring 5-6 operations), and two final multiplications. For the same BERT-base configuration with $32 \times 512 \times 768$ activations, GELU requires approximately $8 \times 12,582,912 = 100,663,296$ operations, taking approximately 0.32 microseconds at peak throughput. The memory bandwidth remains 50 MB, taking 33 microseconds, so GELU is still approximately $100\times$ memory-bandwidth-bound but significantly less so than ReLU.

The computational overhead of GELU compared to ReLU is approximately $4\times$ in terms of arithmetic operations, but the actual runtime difference is much smaller due to memory bandwidth limitations. In practice, GELU adds approximately 10-15% to the total activation computation time compared to ReLU, as both operations spend most of their time waiting for memory transfers rather than computing. For a full BERT-base forward pass taking approximately 50 milliseconds, replacing ReLU with GELU in all 12 layers adds approximately 0.5-1 milliseconds, or 1-2% of total training time. This modest overhead explains why modern transformers universally adopt GELU despite its higher computational cost: the improved training dynamics and final model quality outweigh the small performance penalty.

Swish, defined as $\text{Swish}(z) = z \cdot \sigma(z)$ where σ is the sigmoid function, has computational cost similar to GELU. The sigmoid function requires computing an exponential and a division, totaling approximately 6-8 operations per element including the final multiplication. Swish therefore has comparable performance to GELU, typically within 5-10% in runtime. The choice between GELU and Swish is usually based on empirical performance on specific tasks rather than computational considerations, as their efficiency is nearly identical.

4.4.2 Hardware Support and Fused Kernels

Modern deep learning frameworks provide fused kernels that combine activation functions with preceding operations to reduce memory traffic. A fused linear-GELU kernel computes $\text{GELU}(\mathbf{W}\mathbf{x} + \mathbf{b})$ in a single GPU kernel, eliminating the need to write the intermediate result $\mathbf{z} = \mathbf{W}\mathbf{x} + \mathbf{b}$ to memory and then read it back for the GELU computation. This fusion reduces memory traffic from $3V$ to $2V$ values (where V is the number of activations), providing speedups of $1.3\text{--}1.5\times$ for the combined operation.

For BERT-base with hidden dimension $d = 768$ and feed-forward intermediate dimension $d_{\text{ff}} = 3072$, the first feed-forward layer computes $\text{GELU}(\mathbf{W}_1\mathbf{h} + \mathbf{b}_1)$ where $\mathbf{W}_1 \in \mathbb{R}^{3072 \times 768}$. Without fusion, this requires writing $32 \times 512 \times 3072 = 50,331,648$ FP16 values (100 MB) to memory after the linear layer, then reading them back for GELU, totaling 200 MB of memory traffic. With fusion, only the final GELU output is written to memory (100 MB), reducing traffic by 50% and improving runtime from approximately 100 microseconds to 67 microseconds on an A100 GPU. Across 12 transformer layers with 2 feed-forward layers each, this fusion saves approximately 0.8 milliseconds per forward pass, or 1.6% of total training time.

NVIDIA’s cuDNN library and PyTorch’s JIT compiler automatically apply these fusions when possible, but they require that the activation function be known at compile time. Custom activation functions or dynamically selected activations may not benefit from fusion, resulting in 30-50% slower performance. This hardware consideration provides another reason to prefer standard activations like ReLU, GELU, and Swish over custom alternatives: the extensive optimization effort invested in these common operations by hardware vendors and framework developers translates directly to faster training.

4.4.3 Why GELU is Preferred in Transformers

Despite its higher computational cost, GELU has become the standard activation function for transformer models, used in BERT, GPT-2, GPT-3, T5, and most modern language models. This preference is driven by empirical performance rather than computational efficiency: models trained with GELU consistently achieve better final accuracy than those trained with ReLU, particularly on language understanding tasks.

The theoretical motivation for GELU is that it provides a smoother approximation to the ReLU function, with non-zero gradients for negative inputs. While ReLU has gradient zero for all $z < 0$, GELU has small but non-zero gradients in this region, allowing the network to recover from neurons that have been pushed into the negative regime. This property is particularly valuable in deep networks where gradient flow through many layers can be fragile. For a 24-layer BERT-large model, the probability that a gradient signal survives through all layers is significantly higher with GELU than with ReLU, as GELU never completely blocks gradient flow.

Empirically, BERT-base trained with GELU achieves 84.6% accuracy on the MNLI natural language inference task, compared to 83.9% with ReLU—a 0.7 percentage point improvement that is statistically significant and practically meaningful. For GPT-2, the perplexity on the WebText validation set is 18.3 with GELU compared to 19.1 with ReLU, indicating better language modeling performance. These improvements justify the 1-2% computational overhead of GELU, as the improved model quality translates to better downstream task performance and potentially reduced training time to reach a target accuracy.

The success of GELU has inspired variants like Swish and Mish that share the property of smooth, non-zero gradients everywhere. Swish, defined as $\text{Swish}(z) = z \cdot \sigma(z)$, has similar performance to GELU on most tasks and is used in some efficient transformer architectures like EfficientNet. Mish, defined as $\text{Mish}(z) = z \cdot \tanh(\text{softplus}(z))$, provides slightly better performance than GELU on some vision tasks but has higher computational cost. The landscape of activation functions continues to evolve, but GELU remains the standard for language models due to its strong empirical performance and reasonable computational cost.

4.5 Universal Approximation Theorem

Theorem 4.1 (Universal Approximation). *A single-hidden-layer neural network with nonlinear activation can approximate any continuous function on compact domain to arbitrary precision, given sufficient hidden units.*

Caveat: The theorem says nothing about how many units needed, how to find weights, or generalization. Deep networks often more efficient than wide networks.

4.6 Weight Initialization

Definition 4.6 (Xavier Initialization). For layer with n_{in} inputs and n_{out} outputs:

$$w_{ij} \sim \mathcal{N}\left(0, \frac{2}{n_{\text{in}} + n_{\text{out}}}\right) \quad (4.8)$$

Best for tanh and sigmoid activations.

Definition 4.7 (He Initialization). For ReLU networks:

$$w_{ij} \sim \mathcal{N}\left(0, \frac{2}{n_{\text{in}}}\right) \quad (4.9)$$

Accounts for ReLU zeroing half the activations.

4.6.1 Variance Preservation Through Layers

Proper weight initialization ensures that activations and gradients maintain reasonable magnitudes as they propagate through deep networks. Without careful initialization, activations can explode (growing exponentially with depth) or vanish (shrinking to zero), making training impossible. The initialization schemes above are designed to preserve variance through forward and backward passes.

Consider a linear layer $\mathbf{y} = \mathbf{W}\mathbf{x}$ where $\mathbf{x} \in \mathbb{R}^{n_{\text{in}}}$ has zero mean and unit variance, and weights w_{ij} are independent with zero mean and variance σ_w^2 . The variance of each output element is:

$$\text{Var}(y_i) = \text{Var}\left(\sum_{j=1}^{n_{\text{in}}} w_{ij}x_j\right) = \sum_{j=1}^{n_{\text{in}}} \text{Var}(w_{ij})\text{Var}(x_j) = n_{\text{in}}\sigma_w^2 \quad (4.10)$$

To preserve variance ($\text{Var}(y_i) = 1$), we need $\sigma_w^2 = 1/n_{\text{in}}$. This is the basis for Xavier initialization, which uses $\sigma_w^2 = 2/(n_{\text{in}} + n_{\text{out}})$ to balance forward and backward pass variance preservation. The factor of 2 in the numerator accounts for the fact that gradients flow backward through the transpose of the weight matrix, which has dimensions $n_{\text{out}} \times n_{\text{in}}$.

For ReLU activations, the analysis is modified because ReLU zeros out half the activations on average. If the input has variance 1, the output of ReLU has variance approximately 0.5 (since half the values become zero). To compensate, He initialization uses $\sigma_w^2 = 2/n_{\text{in}}$, doubling the variance compared to the linear case. This ensures that after the ReLU activation, the variance returns to approximately 1, maintaining signal strength through deep networks.

The importance of proper initialization becomes apparent in deep networks. For a 100-layer network with Xavier initialization, activations maintain roughly constant variance through all layers. With naive initialization using $\sigma_w^2 = 1$ (too large), activations grow exponentially: after 10 layers, the variance is approximately 10^{10} , causing numerical overflow. With $\sigma_w^2 = 0.01$ (too small), activations shrink exponentially: after 10 layers, the variance is approximately 10^{-20} , causing numerical underflow. Both scenarios make training impossible, as gradients either explode or vanish.

4.6.2 Impact on Training Speed

Proper initialization not only enables training but also significantly affects convergence speed. Networks initialized with appropriate schemes reach target accuracy in fewer training steps, reducing total training time and computational cost.

For BERT-base trained on the MNLI natural language inference task, the impact of initialization is dramatic. With He initialization (appropriate for the GELU activations used in BERT), the model reaches 84% validation accuracy after approximately 15,000 training steps, requiring 3.5 hours on an NVIDIA A100 GPU. With Xavier initialization (suboptimal for GELU), the model reaches the same accuracy after approximately 22,000 steps, requiring 5.1 hours—a 46% increase in training time. With naive initialization using $\sigma_w^2 = 0.01$, the model fails to converge even after 50,000 steps, as the gradients vanish in the deep network.

The mechanism behind this speedup is that proper initialization places the network in a region of parameter space where gradients have appropriate magnitude for learning. With He initialization, the average gradient norm for BERT-base is approximately 1.0 in early training, allowing the Adam optimizer with learning rate 10^{-4} to make meaningful parameter updates. With Xavier initialization, the average gradient norm is approximately 0.3, requiring either a higher learning rate (which risks instability) or more training steps to achieve the same parameter changes. With naive initialization, the gradient norm is approximately 0.001, making learning extremely slow regardless of learning rate.

The computational cost of initialization itself is negligible. Generating random numbers for 110 million parameters in BERT-base requires approximately 50 milliseconds on a CPU, compared to hours or days of training time. Modern deep learning frameworks like PyTorch provide efficient initialization functions that run on the GPU, reducing initialization time to less than 10 milliseconds. This one-time cost is amortized over thousands of training steps, making proper initialization essentially free from a computational perspective while providing substantial benefits for training speed and stability.

4.6.3 GPU Memory During Initialization

Initialization requires temporarily allocating memory for random number generation, which can be significant for very large models. For a model with P parameters, initialization requires $4P$ bytes to store the parameters in FP32, plus additional memory for the random number generator state. For BERT-base with 110 million parameters, this totals 440 MB plus approximately 10 MB for RNG state, totaling 450 MB. This is modest and fits comfortably in any modern GPU.

However, for very large models like GPT-3 with 175 billion parameters, initialization requires $4 \times 175 \times 10^9 = 700$ GB of memory just for the parameters in FP32. This exceeds the memory of any single GPU, requiring distributed initialization across multiple devices. The typical approach is to initialize parameters on CPU in chunks, transfer each chunk to the appropriate GPU, and convert to FP16 to reduce memory. This process can take several minutes for GPT-3, but it remains a one-time cost that is negligible compared to the weeks of training time required.

Modern frameworks provide memory-efficient initialization strategies for large models. PyTorch’s `torch.nn.init` module supports in-place initialization, which avoids allocating temporary tensors. For models using mixed precision training, parameters can be initialized directly in FP16, halving the memory requirement. For models using model parallelism, each GPU initializes only its shard of the parameters, distributing the memory cost across devices. These optimizations make initialization practical even for models with hundreds of billions of parameters.

4.6.4 Example: BERT-base Initialization

BERT-base uses a variant of He initialization adapted for GELU activations. The initialization scheme is:

- Embedding layers: $\mathcal{N}(0, 0.02^2)$ (fixed small variance)
- Linear layers: $\mathcal{N}(0, \sigma^2)$ where $\sigma = \sqrt{2/n_{\text{in}}}$
- Layer norm parameters: $\gamma = 1, \beta = 0$

- Biases: $b = 0$

For the feed-forward layers in BERT-base, the first layer has $n_{\text{in}} = 768$, giving $\sigma = \sqrt{2/768} \approx 0.051$. The second layer has $n_{\text{in}} = 3072$, giving $\sigma = \sqrt{2/3072} \approx 0.026$. These initialization variances ensure that activations maintain unit variance through the network, enabling stable training from the first iteration.

The impact of this initialization can be measured empirically. At initialization (before any training), BERT-base with proper He initialization has average activation magnitude approximately 1.0 in all layers, and gradient magnitude approximately 1.0 for all parameters. With naive initialization using $\sigma = 0.01$ for all layers, the activation magnitude in the final layer is approximately 0.001, and gradients for early layers are approximately 10^{-6} , making learning extremely slow. With too-large initialization using $\sigma = 0.1$, the activation magnitude in the final layer is approximately 100, and gradients are approximately 1000, causing training instability and divergence.

The lesson is clear: proper initialization is not optional but essential for training deep networks efficiently. The specific initialization scheme (Xavier vs He vs other variants) matters less than ensuring that variance is preserved through the network. For transformer models with GELU activations, He initialization or slight variants thereof work well and are used universally in BERT, GPT, T5, and other modern architectures.

4.7 Regularization

4.7.1 L2 Regularization

Add penalty to loss:

$$L_{\text{total}} = L_{\text{data}} + \frac{\lambda}{2} \sum_{\ell} \left\| \mathbf{W}^{(\ell)} \right\|_F^2 \quad (4.11)$$

L2 regularization, also known as weight decay, penalizes large parameter values to prevent overfitting. The regularization term adds the squared Frobenius norm of all weight matrices to the loss function, encouraging the optimizer to keep weights small. The hyperparameter λ controls the strength of regularization: larger λ produces smaller weights and stronger regularization.

The computational cost of L2 regularization is modest. Computing the squared norm $\|\mathbf{W}\|_F^2 = \sum_{ij} w_{ij}^2$ requires one multiplication and one addition per parameter, totaling $2P$ operations for a model with P parameters. For BERT-base with 110 million parameters, this requires 220 million operations, or 0.22 GFLOPs. Compared to the 96 GFLOPs required for a forward pass, the regularization computation adds only 0.23% overhead. On an NVIDIA A100 GPU, computing the regularization term takes approximately 0.7 microseconds, which is negligible compared to the 50 milliseconds for a full forward-backward pass.

The gradient of the L2 regularization term is even simpler: $\nabla_{\mathbf{W}} \left(\frac{\lambda}{2} \|\mathbf{W}\|_F^2 \right) = \lambda \mathbf{W}$. This adds a term proportional to the current weights to the gradient, which can be implemented as a simple scaling operation during the optimizer step. Most optimizers, including PyTorch's Adam and SGD, support weight decay as a built-in parameter that applies this scaling automatically without requiring explicit computation of the regularization term. This makes L2 regularization essentially free from a computational perspective.

The memory overhead of L2 regularization is zero, as it requires no additional storage beyond the parameters themselves. The regularization term is computed on-the-fly during the backward pass and does not need to be stored. This makes L2 regularization an attractive regularization technique for large models where memory is at a premium.

4.7.2 Dropout

Definition 4.8 (Dropout). During training, randomly set activations to zero with probability p . During inference, scale by $(1 - p)$.

Dropout is a powerful regularization technique that randomly drops (sets to zero) a fraction of activations during training. This prevents the network from relying too heavily on any single neuron and encourages learning robust features. The dropout probability p is typically 0.1 to 0.5, with higher values providing stronger regularization at the cost of slower convergence.

4.7.3 Computational Overhead of Dropout

The computational cost of dropout consists of random number generation and masking operations. For each activation tensor with N elements, dropout requires generating N random numbers, comparing each to the threshold p , and multiplying the activations by the resulting binary mask. Additionally, the surviving activations must be scaled by $1/(1 - p)$ to maintain expected values.

For a BERT-base layer with batch size $B = 32$, sequence length $n = 512$, and hidden dimension $d = 768$, the activation tensor has $32 \times 512 \times 768 = 12,582,912$ elements. Generating 12.6 million random numbers on a GPU takes approximately 50 microseconds using CUDA’s cuRAND library. The masking operation (element-wise multiplication) requires 12.6 million operations, taking approximately 0.04 microseconds at peak throughput but actually taking approximately 20 microseconds due to memory bandwidth limitations (reading activations, reading mask, writing masked activations). The scaling operation requires another 12.6 million operations, taking approximately 20 microseconds. The total dropout overhead is approximately 90 microseconds per layer.

For a 12-layer BERT-base model with dropout applied after attention and feed-forward layers (2 dropout operations per layer), the total dropout overhead is $12 \times 2 \times 90 = 2,160$ microseconds, or approximately 2.2 milliseconds per forward pass. Compared to the 50 milliseconds for the full forward pass, dropout adds approximately 4.4% overhead. The backward pass has similar overhead, as dropout must be applied to gradients as well, bringing the total dropout overhead to approximately 4.4 milliseconds per training step, or 4.4% of total training time.

This overhead is non-negligible but acceptable given the regularization benefits. Dropout typically improves final model accuracy by 0.5-2 percentage points on downstream tasks, which justifies the 4-5% increase in training time. For models where training time is critical, dropout can be reduced or eliminated, but this often requires other forms of regularization (like L2 regularization or data augmentation) to maintain model quality.

4.7.4 Memory Requirements for Dropout

Dropout requires storing the binary dropout mask for use in the backward pass. For an activation tensor with N elements, the mask requires N bits, or $N/8$ bytes. For BERT-base with $32 \times 512 \times 768$ activations per layer, the mask requires $12,582,912/8 = 1,572,864$ bytes, or approximately 1.5 MB per dropout operation. With 2 dropout operations per layer and 12 layers, the total mask memory is $12 \times 2 \times 1.5 = 36$ MB.

This memory overhead is modest compared to the activation memory itself (approximately 10 GB for BERT-base with batch size 32), adding only 0.36% overhead. However, for very large batch sizes or long sequences, the mask memory can become significant. At batch size 256 and sequence length 2048, the mask memory for BERT-base would be $12 \times 2 \times 256 \times 2048 \times 768/8 = 1,207,959,552$ bytes, or approximately 1.15 GB. This is still manageable on modern GPUs with 40-80 GB of memory, but it represents a non-trivial fraction of the memory budget.

Modern deep learning frameworks optimize dropout memory by using compact representations. PyTorch stores dropout masks as boolean tensors (1 byte per element) rather than float tensors (4 bytes per element), reducing memory by $4\times$. Some implementations use bit-packed representations (1 bit per element) to reduce memory by $32\times$, though this requires custom CUDA kernels and is not standard in most frameworks. For most applications, the memory overhead of dropout is acceptable and does not limit batch size or sequence length.

4.7.5 Inference Mode Differences

During inference, dropout is disabled: all activations are kept, and no scaling is applied (assuming the training-time scaling approach where activations are divided by $1 - p$). This means inference is faster

than training, as it avoids the random number generation and masking operations. For BERT-base, disabling dropout reduces inference time from approximately 50 milliseconds to 48 milliseconds per batch, a 4% speedup. This speedup is modest but can be significant for latency-sensitive applications where every millisecond counts.

The alternative approach, called inverted dropout, scales activations during training by $1/(1 - p)$ and does nothing during inference. This is the approach used in most modern frameworks, as it makes inference code simpler (no scaling required) and slightly faster. The computational cost is identical to standard dropout, but the implementation is cleaner and less error-prone.

4.7.6 Dropout in Transformer Models

Transformer models apply dropout at multiple points in the architecture:

- Attention dropout: Applied to attention weights after softmax
- Residual dropout: Applied to the output of attention and feed-forward layers before adding to the residual connection
- Embedding dropout: Applied to input embeddings

BERT-base uses dropout probability $p = 0.1$ at all these locations, totaling 4 dropout operations per transformer layer (attention dropout, attention residual dropout, feed-forward residual dropout, and embedding dropout for the first layer). With 12 layers, this totals approximately 50 dropout operations per forward pass, consuming approximately 4.5 milliseconds or 9% of total training time. This overhead is higher than for simple feed-forward networks due to the multiple dropout locations, but it provides strong regularization that is essential for good generalization on downstream tasks.

For GPT-3, dropout is applied more sparingly: only residual dropout with $p = 0.1$ is used, and attention dropout is disabled. This reduces the dropout overhead to approximately 2 dropout operations per layer, or 192 operations for the 96-layer model. The total dropout overhead is approximately 17 milliseconds per forward pass, or approximately 5% of total training time. The reduced dropout is compensated by the massive scale of the training data (300 billion tokens), which provides implicit regularization through data diversity.

The lesson is that dropout overhead scales with the number of dropout operations and the size of the activation tensors. For models with many layers and large hidden dimensions, dropout can consume 5-10% of training time. This overhead is generally acceptable given the regularization benefits, but for models where training time is critical, reducing the number of dropout operations or using lower dropout probabilities can provide speedups with minimal impact on final model quality.

4.8 Exercises

Exercise 4.1. Design 3-layer MLP for binary classification of 100-dimensional inputs. Specify layer dimensions, activations, and parameter count.

Exercise 4.2. Compute forward pass through 2-layer network with given weights and ReLU activation.

Exercise 4.3. For layer with 512 inputs and 256 outputs using ReLU: (1) What is He initialization variance? (2) Why different from Xavier? (3) What happens with zero initialization?

Exercise 4.4. Prove that without nonlinear activations, L-layer network equivalent to single layer.

Chapter 5

Convolutional Neural Networks

Chapter Overview

Convolutional Neural Networks (CNNs) revolutionized computer vision by exploiting spatial structure. This chapter develops convolution operations, pooling, and modern CNN architectures including ResNet.

Learning Objectives

1. Understand convolution operations and compute output dimensions
2. Design CNN architectures with appropriate pooling and stride
3. Understand translation equivariance
4. Implement modern CNN architectures (ResNet, VGG)

5.1 Convolution Operation

Definition 5.1 (2D Convolution). For input $\mathbf{X} \in \mathbb{R}^{H \times W}$ and kernel $\mathbf{K} \in \mathbb{R}^{k_h \times k_w}$:

$$(\mathbf{X} \star \mathbf{K})_{i,j} = \sum_{m=0}^{k_h-1} \sum_{n=0}^{k_w-1} \mathbf{X}_{i+m,j+n} \cdot \mathbf{K}_{m,n} \quad (5.1)$$

Example 5.1 (3x3 Convolution). Input 4×4 , kernel 3×3 (edge detector), output 2×2 . Computing first position: sum of element-wise products gives edge response.

5.1.1 Output Dimensions

Theorem 5.1 (Output Size). For input size $H \times W$, kernel $k_h \times k_w$, padding p , stride s :

$$H_{out} = \left\lfloor \frac{H + 2p - k_h}{s} \right\rfloor + 1 \quad (5.2)$$

5.2 Multi-Channel Convolutions

Definition 5.2 (Convolutional Layer). For input $\mathbf{X} \in \mathbb{R}^{C_{\text{in}} \times H \times W}$ with C_{out} output channels:

$$\mathbf{Y}^{(i)} = \sum_{c=1}^{C_{\text{in}}} \mathbf{X}^{(c)} \star \mathbf{K}^{(i,c)} + b^{(i)} \quad (5.3)$$

Example 5.2 (RGB Convolution). Input: $\mathbf{X} \in \mathbb{R}^{3 \times 224 \times 224}$. Conv layer: 64 filters 3×3 , stride 1, padding 1.

Parameters: $64 \times 3 \times 3 \times 3 + 64 = 1,792$

Output: $\mathbf{Y} \in \mathbb{R}^{64 \times 224 \times 224}$

Compare to fully-connected: ≈ 483 billion parameters!

Key Point 5.1. *Convolution provides: (1) Parameter sharing, (2) Local connectivity, (3) Translation equivariance. Massive parameter reduction compared to fully-connected layers.*

5.3 Computational Analysis of Convolutions

Understanding the computational cost and memory requirements of convolutional layers is essential for designing efficient architectures and comparing CNNs with alternative approaches like transformers. The relationship between parameters, FLOPs, and memory usage in convolutions differs fundamentally from fully-connected layers, leading to distinct performance characteristics on modern hardware.

5.3.1 FLOPs for Convolution Operations

The computational cost of a convolutional layer is determined by the number of multiply-accumulate operations required to compute all output feature maps. For a convolutional layer with input shape $C_{\text{in}} \times H \times W$, kernel size $k \times k$, and C_{out} output channels, each output position requires $C_{\text{in}} \times k \times k$ multiply-accumulate operations. With output spatial dimensions $H_{\text{out}} \times W_{\text{out}}$ and C_{out} output channels, the total FLOPs is:

$$\text{FLOPs}_{\text{conv}} = 2 \times C_{\text{out}} \times C_{\text{in}} \times k^2 \times H_{\text{out}} \times W_{\text{out}} \quad (5.4)$$

The factor of 2 accounts for the multiply-accumulate operation (one multiplication and one addition per operation). This formula reveals that convolution FLOPs scale linearly with both input and output channels, quadratically with kernel size, and linearly with output spatial dimensions.

For the RGB convolution example in Example 5.2 with input $3 \times 224 \times 224$, kernel size 3×3 , and 64 output channels with stride 1 and padding 1, the output dimensions are $64 \times 224 \times 224$. The FLOPs calculation is $2 \times 64 \times 3 \times 9 \times 224 \times 224 = 173,408,192$ FLOPs, or approximately 173 MFLOPs. Despite having only 1,792 parameters, this layer requires 173 million floating-point operations, giving a FLOPs-to-parameter ratio of approximately 96,768. This ratio is dramatically higher than fully-connected layers, which have a FLOPs-to-parameter ratio of approximately 2-3.

The high FLOPs-to-parameter ratio of convolutions has important implications for model design. Convolutional layers are compute-intensive relative to their memory footprint, making them well-suited for modern GPUs that have abundant compute throughput but limited memory bandwidth. A ResNet-50 model with 25.6 million parameters requires approximately 4.1 billion FLOPs for a single forward pass on a 224×224 image, giving an overall FLOPs-to-parameter ratio of 160. This means that during

training, the computational cost dominates over the memory cost of loading parameters, and GPU utilization is primarily limited by compute throughput rather than memory bandwidth.

The scaling behavior of convolution FLOPs explains several architectural design choices in modern CNNs. Early layers operating on high-resolution feature maps (224×224 or larger) consume the majority of FLOPs despite having relatively few parameters. For ResNet-50, the first convolutional layer with kernel size 7×7 and 64 output channels accounts for only 0.4% of parameters but 5.8% of total FLOPs. Conversely, later layers operating on low-resolution feature maps (7×7 or smaller) have many parameters but relatively few FLOPs. The final fully-connected layer in ResNet-50 accounts for 7.8% of parameters but only 0.1% of FLOPs. This distribution motivates the use of larger kernels and more channels in early layers (where spatial dimensions are large) and smaller kernels with many channels in later layers (where spatial dimensions are small).

5.3.2 Memory Requirements for Feature Maps

During training, convolutional networks must store intermediate feature maps for use in the backward pass, and these activations typically consume far more memory than the model parameters. Understanding activation memory is critical for determining maximum batch size and input resolution.

For a convolutional layer with input shape $B \times C_{\text{in}} \times H \times W$ (where B is batch size) and output shape $B \times C_{\text{out}} \times H_{\text{out}} \times W_{\text{out}}$, the network must store both the input feature map and the output feature map. In FP32, this requires $4B(C_{\text{in}}HW + C_{\text{out}}H_{\text{out}}W_{\text{out}})$ bytes of memory. For the RGB convolution example with batch size $B = 32$, input $3 \times 224 \times 224$, and output $64 \times 224 \times 224$, the activation memory is $4 \times 32 \times (3 \times 224 \times 224 + 64 \times 224 \times 224) = 130,809,792$ bytes, or approximately 125 MB. This is $70,000\times$ larger than the parameter memory (1,792 parameters = 7,168 bytes), demonstrating that activation memory dominates for convolutional layers.

The memory consumption of a full CNN scales with the number of layers and the spatial dimensions of feature maps. For ResNet-50 processing batch size 32 with input $3 \times 224 \times 224$, the total activation memory is approximately 8.2 GB in FP32. This includes the input image (6.4 MB), early high-resolution feature maps (hundreds of MB), and later low-resolution feature maps (tens of MB). The parameter memory for ResNet-50 is only 102 MB (25.6 million parameters \times 4 bytes), making activations $80\times$ larger than parameters. This ratio increases with batch size: at batch size 256, activations consume 65.6 GB while parameters remain 102 MB, a ratio of $643\times$.

The quadratic scaling of activation memory with spatial resolution has profound implications for input image size. Doubling the input resolution from 224×224 to 448×448 increases the number of pixels by $4\times$, and since early feature maps maintain similar spatial dimensions to the input, activation memory increases by approximately $4\times$. For ResNet-50 with batch size 32, increasing resolution from 224×224 to 448×448 increases activation memory from 8.2 GB to approximately 32.8 GB, exceeding the capacity of most GPUs. This explains why high-resolution image processing typically requires smaller batch sizes or gradient accumulation: the activation memory grows faster than available GPU memory.

Modern techniques for reducing activation memory include gradient checkpointing, which recomputes activations during the backward pass rather than storing them, trading computation for memory. For ResNet-50, gradient checkpointing can reduce activation memory by $5\text{--}10\times$ at the cost of increasing training time by 20-30%. This trade-off is often worthwhile for training with larger batch sizes or higher resolutions, as the improved convergence from larger batches can offset the increased computation time.

5.3.3 GPU Optimization: im2col and Winograd

Efficient implementation of convolution on GPUs requires specialized algorithms that transform the convolution operation into a form amenable to highly optimized matrix multiplication routines. The two primary approaches are im2col (image-to-column) and Winograd convolution, each with distinct performance characteristics.

The im2col algorithm transforms convolution into matrix multiplication by unrolling the input feature map into a large matrix where each column contains the input values for one output position. For a convolutional layer with input $C_{\text{in}} \times H \times W$, kernel size $k \times k$, and output $C_{\text{out}} \times H_{\text{out}} \times W_{\text{out}}$, im2col creates a matrix of shape $(C_{\text{in}}k^2) \times (H_{\text{out}}W_{\text{out}})$ by extracting all $k \times k$ patches from the input.

The convolution kernels are reshaped into a matrix of shape $C_{\text{out}} \times (C_{\text{in}}k^2)$. The convolution is then computed as a single matrix multiplication: $\text{output} = \text{kernels} \times \text{im2col}(\text{input})$, producing a matrix of shape $C_{\text{out}} \times (H_{\text{out}}W_{\text{out}})$ that is reshaped to the final output dimensions.

For the RGB convolution example with input $3 \times 224 \times 224$, kernel size 3×3 , and 64 output channels, im2col creates a matrix of shape $27 \times 50,176$ (since $C_{\text{in}}k^2 = 3 \times 9 = 27$ and $H_{\text{out}}W_{\text{out}} = 224 \times 224 = 50,176$). The kernel matrix has shape 64×27 . The matrix multiplication 64×27 times $27 \times 50,176$ requires $2 \times 64 \times 27 \times 50,176 = 173,408,192$ FLOPs, matching the direct convolution calculation. However, the im2col matrix requires $27 \times 50,176 \times 4 = 5,419,008$ bytes (5.2 MB) of temporary storage, which is $757\times$ larger than the original input (7,168 bytes for $3 \times 224 \times 224$ in FP32).

The advantage of im2col is that it leverages highly optimized BLAS (Basic Linear Algebra Subprograms) libraries like cuBLAS on NVIDIA GPUs, which achieve 80-95% of peak hardware throughput for large matrix multiplications. For the 64×27 times $27 \times 50,176$ multiplication on an NVIDIA A100 GPU with 312 TFLOPS FP16 throughput, the operation completes in approximately 0.6 microseconds at 90% efficiency, achieving 280 TFLOPS. Direct convolution implementations without im2col typically achieve only 40-60% efficiency due to irregular memory access patterns and difficulty saturating the GPU's parallel execution units.

The disadvantage of im2col is the memory overhead. For batch size $B = 32$, the im2col matrix grows to $32 \times 27 \times 50,176 = 43,352,064$ elements, requiring 167 MB of temporary storage. This memory must be allocated and deallocated for each convolutional layer, adding memory pressure and potentially causing out-of-memory errors for large batch sizes or high-resolution inputs. Modern implementations mitigate this by processing the batch in chunks or fusing the im2col transformation with the matrix multiplication to avoid materializing the full im2col matrix.

Winograd convolution is an alternative algorithm that reduces the number of multiplications required for small convolutions (typically 3×3 or 5×5 kernels) by using a mathematical transformation that trades multiplications for additions. For 3×3 convolutions, Winograd reduces the number of multiplications by $2.25\times$ compared to direct convolution, from 9 multiplications per output to 4 multiplications per output. This reduction translates directly to FLOPs savings: the RGB convolution example requires only $173,408,192/2.25 = 77,070,752$ FLOPs with Winograd, a 56% reduction.

However, Winograd convolution has several limitations. First, it requires additional memory for intermediate transformations, typically $2\text{-}3\times$ the input size. Second, it is numerically less stable than direct convolution, particularly in FP16, due to the transformation matrices having large condition numbers. Third, it is only applicable to small kernel sizes (3×3 and 5×5) and becomes inefficient for larger kernels. Fourth, the transformation overhead becomes significant for small spatial dimensions, making Winograd most effective for early layers with large feature maps.

In practice, modern deep learning frameworks like PyTorch and TensorFlow automatically select between im2col , Winograd, and direct convolution based on layer dimensions, batch size, and hardware characteristics. For 3×3 convolutions on high-resolution feature maps ($\geq 56 \times 56$) with batch size ≥ 16 , Winograd typically provides $1.5\text{-}2\times$ speedup over im2col . For larger kernels (5×5 or 7×7) or smaller feature maps, im2col is preferred. For very small batch sizes (≤ 4), direct convolution may be fastest due to lower overhead. NVIDIA's cuDNN library implements all three algorithms and includes heuristics to select the optimal approach for each layer configuration.

5.3.4 Comparison with Transformer Attention

Comparing the computational characteristics of convolutional layers with transformer self-attention reveals fundamental trade-offs between local and global receptive fields, parameter efficiency, and computational scaling.

A convolutional layer with kernel size $k \times k$ has a local receptive field: each output position depends only on a $k \times k$ neighborhood of the input. To achieve a global receptive field spanning the entire input, multiple convolutional layers must be stacked. For an input of size $H \times W$, achieving a receptive field covering the full input requires approximately $\log_k(\max(H, W))$ layers. For a 224×224 image with 3×3 convolutions, this requires approximately $\log_3(224) \approx 5$ layers. Each layer adds computational cost, but the cost per layer remains $O(C_{\text{out}}C_{\text{in}}k^2HW)$, scaling linearly with spatial dimensions.

In contrast, self-attention in transformers has a global receptive field: each output position attends to

all input positions in a single layer. For an input sequence of length $n = HW$ (treating the 2D image as a 1D sequence) with model dimension d , self-attention requires computing query-key products for all pairs of positions, resulting in $O(n^2d)$ FLOPs. For a 224×224 image with $n = 50,176$ positions and $d = 768$ (typical for Vision Transformers), self-attention requires approximately $2 \times 50,176^2 \times 768 = 3.86 \times 10^{12}$ FLOPs, or 3.86 TFLOPs per layer. This is $22,000\times$ more expensive than the RGB convolution example (173 MFLOPs), despite both operating on the same input resolution.

The quadratic scaling of attention with spatial resolution makes it prohibitively expensive for high-resolution images. Doubling the resolution from 224×224 to 448×448 increases attention FLOPs by $16\times$ (since n increases by $4\times$ and attention scales as n^2), while convolution FLOPs increase by only $4\times$ (linear scaling with spatial dimensions). For a 448×448 image, self-attention requires 61.8 TFLOPs per layer, making it impractical without modifications like hierarchical attention or local attention windows.

Vision Transformers (ViTs) address this computational challenge by dividing the image into patches and treating each patch as a token. For a 224×224 image with patch size 16×16 , the sequence length is $n = (224/16)^2 = 196$ patches. Self-attention on 196 patches with $d = 768$ requires $2 \times 196^2 \times 768 = 59,015,168$ FLOPs, or approximately 59 MFLOPs per layer. This is $65\times$ less expensive than attention on individual pixels and comparable to the RGB convolution example (173 MFLOPs). However, the patch-based approach sacrifices fine-grained spatial resolution: each patch is treated as a single token, and the model cannot attend to individual pixels within a patch.

The parameter efficiency of convolutions versus attention also differs significantly. A convolutional layer with C_{in} input channels, C_{out} output channels, and kernel size $k \times k$ has $C_{\text{out}}C_{\text{in}}k^2$ parameters. For the RGB convolution example, this is $64 \times 3 \times 9 = 1,728$ parameters. A self-attention layer with model dimension d has query, key, and value projection matrices, each of size $d \times d$, totaling $3d^2$ parameters (ignoring the output projection). For $d = 768$, this is $3 \times 768^2 = 1,769,472$ parameters, which is $1,024\times$ more than the convolutional layer. However, the attention parameters are independent of spatial resolution, while convolution parameters are independent of spatial resolution as well. The key difference is that attention parameters scale with d^2 while convolution parameters scale with $C_{\text{in}}C_{\text{out}}k^2$, and typically $d \gg C_{\text{in}}$ for early layers.

For complete models, ResNet-50 has 25.6 million parameters and requires 4.1 GFLOPs per image, while ViT-Base has 86 million parameters and requires 17.6 GFLOPs per image. The ViT has $3.4\times$ more parameters and $4.3\times$ more FLOPs, but achieves comparable or better accuracy on ImageNet classification. The higher computational cost of ViT is offset by its ability to leverage large-scale pretraining on datasets like ImageNet-21k or JFT-300M, where the global receptive field and flexibility of attention provide advantages over the inductive biases of convolution.

5.4 Parameter Efficiency: CNNs vs Transformers

The parameter efficiency of convolutional networks compared to transformers is a critical consideration for model design, particularly for vision tasks where input dimensions are large. Understanding why CNNs achieve strong performance with fewer parameters than transformers reveals fundamental differences in their architectural inductive biases.

5.4.1 Why CNNs are More Parameter-Efficient for Images

Convolutional networks achieve parameter efficiency through three key mechanisms: weight sharing, local connectivity, and hierarchical feature learning. These properties are particularly well-suited to natural images, which exhibit strong spatial locality and translation invariance.

Weight sharing in convolutions means that the same kernel is applied to all spatial positions in the input. A 3×3 convolutional kernel with 64 input channels and 64 output channels has $64 \times 64 \times 9 = 36,864$ parameters, regardless of whether the input is 32×32 or 224×224 . This kernel is applied $H \times W$ times (once per output position), effectively sharing the same 36,864 parameters across all spatial locations. In contrast, a fully-connected layer connecting a $224 \times 224 \times 64$ input to a $224 \times 224 \times 64$ output would require $224^2 \times 64 \times 224^2 \times 64 = 1.3 \times 10^{11}$ parameters, which is 3.5 million times larger.

Weight sharing reduces parameters by a factor equal to the spatial dimensions, which is enormous for images.

Local connectivity means that each output position depends only on a small neighborhood of the input, rather than the entire input. For a 3×3 convolution, each output depends on only 9 input positions (plus all input channels). This locality assumption is well-matched to natural images, where nearby pixels are highly correlated and distant pixels are largely independent. By restricting connectivity to local neighborhoods, convolutions avoid the quadratic parameter growth of fully-connected layers while still capturing the relevant spatial structure.

Hierarchical feature learning in CNNs builds global receptive fields through stacking local operations. Early layers learn low-level features like edges and textures with small receptive fields, middle layers learn mid-level features like object parts with medium receptive fields, and late layers learn high-level features like whole objects with large receptive fields. This hierarchy is achieved by stacking convolutional layers with pooling or strided convolutions to progressively reduce spatial dimensions. For ResNet-50, the receptive field grows from 7×7 in the first layer to 427×427 in the final layer, covering the entire 224×224 input multiple times over. This hierarchical approach requires far fewer parameters than directly modeling global dependencies, as each layer only needs to model local relationships.

The parameter efficiency of CNNs is evident in model comparisons. ResNet-50 achieves 76.1% top-1 accuracy on ImageNet with 25.6 million parameters, while ViT-Base achieves 77.9% accuracy with 86 million parameters— $3.4\times$ more parameters for a 1.8 percentage point improvement. EfficientNet-B0 achieves 77.1% accuracy with only 5.3 million parameters, demonstrating that carefully designed CNNs can match or exceed transformer performance with $16\times$ fewer parameters. The parameter efficiency of CNNs makes them particularly attractive for deployment on resource-constrained devices like mobile phones or embedded systems, where model size directly impacts memory usage and inference latency.

However, the parameter efficiency of CNNs comes with trade-offs. The strong inductive biases of weight sharing and local connectivity make CNNs sample-efficient for small datasets but potentially limit their capacity to learn from very large datasets. Vision Transformers, with their weaker inductive biases and higher parameter counts, can leverage massive datasets like ImageNet-21k (14 million images) or JFT-300M (300 million images) to achieve superior performance. When pretrained on JFT-300M, ViT-Large (307 million parameters) achieves 87.8% accuracy on ImageNet, significantly outperforming any CNN. The optimal architecture depends on the available data: CNNs excel with limited data, while transformers excel with abundant data.

5.4.2 Vision Transformer Comparison

Vision Transformers (ViTs) adapt the transformer architecture from NLP to computer vision by treating images as sequences of patches. Understanding the architectural differences and performance trade-offs between ViTs and CNNs is essential for selecting the appropriate model for a given task and dataset.

A Vision Transformer divides an input image into non-overlapping patches, linearly embeds each patch, and processes the sequence of patch embeddings with standard transformer layers. For a 224×224 image with patch size 16×16 , the image is divided into $(224/16)^2 = 196$ patches. Each patch is flattened to a vector of size $16 \times 16 \times 3 = 768$ and linearly projected to the model dimension d (typically 768 for ViT-Base). The resulting sequence of 196 tokens is processed by 12 transformer layers with multi-head self-attention and feed-forward networks, identical to BERT.

The parameter breakdown for ViT-Base reveals where parameters are allocated. The patch embedding layer has $768 \times 768 = 589,824$ parameters (projecting flattened patches to model dimension). Each of the 12 transformer layers has approximately 7.1 million parameters: 2.4 million for self-attention (query, key, value, and output projections) and 4.7 million for the feed-forward network (two linear layers with $4\times$ expansion). The classification head has $768 \times 1000 = 768,000$ parameters for ImageNet's 1000 classes. The total is approximately 86 million parameters, with 85% in the transformer layers and 15% in embeddings and classification head.

Comparing ViT-Base to ResNet-50 reveals fundamental differences in parameter allocation. ResNet-50 has 25.6 million parameters distributed across convolutional layers (23.5 million, 92%), batch normalization (1.1 million, 4%), and the classification head (1.0 million, 4%). The convolutional parameters are concentrated in later layers operating on low-resolution feature maps: the final residual block has

7.1 million parameters despite operating on 7×7 feature maps, while the first convolutional layer has only 9,408 parameters despite operating on 224×224 inputs. This distribution reflects the CNN’s hierarchical design, where early layers extract simple features with few parameters and late layers combine features with many parameters.

In contrast, ViT-Base distributes parameters uniformly across layers: each of the 12 transformer layers has approximately 7.1 million parameters, regardless of the stage of processing. This uniform distribution reflects the transformer’s lack of hierarchical structure: all layers operate on the same sequence length (196 patches) and model dimension (768), performing the same operations. The absence of hierarchy means that ViT must learn hierarchical features implicitly through the attention mechanism, rather than having them built into the architecture as in CNNs.

The computational cost comparison is similarly revealing. ResNet-50 requires 4.1 GFLOPs per image, with 3.8 GFLOPs (93%) in convolutional layers and 0.3 GFLOPs (7%) in other operations. The FLOPs are concentrated in early layers: the first residual block accounts for 1.2 GFLOPs (29%) despite having only 0.2 million parameters (0.8%), while the final residual block accounts for 0.1 GFLOPs (2.4%) despite having 7.1 million parameters (28%). This distribution reflects the $O(C_{\text{in}}C_{\text{out}}k^2HW)$ scaling of convolution FLOPs: early layers have large HW but small C , while late layers have small HW but large C .

ViT-Base requires 17.6 GFLOPs per image, with 16.8 GFLOPs (95%) in transformer layers and 0.8 GFLOPs (5%) in patch embedding and classification. The FLOPs are distributed uniformly across layers: each transformer layer accounts for approximately 1.4 GFLOPs (8%). The self-attention in each layer requires 0.6 GFLOPs (computed as $2 \times 196^2 \times 768 \times 12/10^9$ for query-key products and attention-value products across 12 heads), while the feed-forward network requires 0.8 GFLOPs (computed as $2 \times 196 \times 768 \times 3072 \times 2/10^9$ for two linear layers). The uniform distribution reflects the constant sequence length and model dimension throughout the network.

The accuracy comparison on ImageNet reveals the impact of pretraining scale. When trained from scratch on ImageNet-1k (1.3 million images), ResNet-50 achieves 76.1% top-1 accuracy while ViT-Base achieves only 72.3% accuracy—3.8 percentage points worse despite having $3.4\times$ more parameters and $4.3\times$ more FLOPs. This performance gap demonstrates that ViT’s weak inductive biases require more data to learn effectively. However, when pretrained on ImageNet-21k (14 million images) and fine-tuned on ImageNet-1k, ViT-Base achieves 81.8% accuracy, surpassing ResNet-50 by 5.7 percentage points. With even larger pretraining on JFT-300M (300 million images), ViT-Large achieves 87.8% accuracy, establishing a new state-of-the-art.

The lesson is clear: CNNs are more parameter-efficient and sample-efficient for small to medium datasets, making them the preferred choice when data is limited or computational resources are constrained. Vision Transformers excel when large-scale pretraining data is available, leveraging their flexibility and capacity to achieve superior performance. The optimal choice depends on the specific use case: CNNs for resource-constrained deployment or limited data, transformers for maximum accuracy with abundant data and compute.

5.4.3 Hybrid Architectures

Hybrid architectures combine convolutional and transformer components to leverage the strengths of both approaches. These models use convolutional layers for early feature extraction, exploiting the parameter efficiency and translation equivariance of convolutions, then apply transformer layers for global reasoning, exploiting the flexibility and long-range modeling of attention.

The Convolutional Vision Transformer (CvT) replaces the patch embedding in ViT with a convolutional stem consisting of several convolutional layers with stride. For a 224×224 input, the convolutional stem progressively reduces spatial dimensions to 56×56 , 28×28 , and 14×14 while increasing channels to 64, 192, and 384. At each stage, a transformer layer processes the feature map (treating spatial positions as tokens), then a strided convolution reduces dimensions for the next stage. This hierarchical design combines the parameter efficiency of convolutions with the global modeling of transformers, achieving 81.6% ImageNet accuracy with only 20 million parameters— $4.3\times$ fewer than ViT-Base while matching its accuracy.

The Swin Transformer introduces hierarchical transformers with shifted windows, creating a pyra-

mid structure similar to CNNs. The input is divided into 4×4 patches (rather than 16×16 in ViT), creating a sequence of $56 \times 56 = 3,136$ tokens. Transformer layers process this sequence using local attention within 7×7 windows (rather than global attention), reducing computational cost from $O(n^2)$ to $O(n)$. After several layers, adjacent patches are merged to create a 28×28 sequence with doubled channel dimension, and the process repeats. This hierarchical design achieves 83.3% ImageNet accuracy with 29 million parameters, outperforming both ResNet-50 and ViT-Base while using fewer parameters than ViT.

The success of hybrid architectures demonstrates that the dichotomy between CNNs and transformers is not absolute. By combining convolutional inductive biases for early processing with transformer flexibility for late processing, hybrid models achieve better parameter efficiency and accuracy than either pure CNNs or pure transformers. This trend suggests that future vision models will increasingly blend architectural components rather than adhering strictly to one paradigm or the other.

5.5 Hardware Optimization for Convolutions

Efficient execution of convolutional networks on modern hardware requires understanding the interaction between algorithm design, memory hierarchy, and specialized compute units. This section examines how convolutions map to GPU architectures and how to maximize hardware utilization.

5.5.1 Tensor Core Utilization for Convolutions

Modern NVIDIA GPUs include Tensor Cores, specialized hardware units that accelerate matrix multiplication for specific data types and dimensions. Understanding how convolutions map to Tensor Cores is essential for achieving peak performance.

Tensor Cores on NVIDIA A100 GPUs perform matrix multiplication on 16×16 tiles in FP16, producing FP32 accumulation. Each Tensor Core can execute one $16 \times 16 \times 16$ matrix multiplication per clock cycle, computing $C = A \times B$ where A is 16×16 , B is 16×16 , and C is 16×16 . The A100 has 432 Tensor Cores running at 1.41 GHz, providing peak throughput of $432 \times 2 \times 16^3 \times 1.41 \times 10^9 = 312$ TFLOPS in FP16. To achieve this peak throughput, matrix dimensions must be multiples of 16, and the matrices must be large enough to saturate all Tensor Cores.

Convolutions map to Tensor Cores through the im2col transformation described in Section 5.3. For a convolutional layer with C_{out} output channels, C_{in} input channels, kernel size $k \times k$, and output spatial dimensions $H_{\text{out}} \times W_{\text{out}}$, im2col creates a matrix multiplication of shape $(C_{\text{out}}) \times (C_{\text{in}}k^2) \times (H_{\text{out}}W_{\text{out}})$. To achieve high Tensor Core utilization, all three dimensions should be multiples of 16 and sufficiently large.

For the RGB convolution example with $C_{\text{out}} = 64$, $C_{\text{in}} = 3$, $k = 3$, $H_{\text{out}} = W_{\text{out}} = 224$, the matrix multiplication has dimensions $64 \times 27 \times 50,176$. The output dimension (64) is a multiple of 16, which is good. The inner dimension (27) is not a multiple of 16, which reduces efficiency: Tensor Cores will pad to 32, wasting $5/32 = 15.6\%$ of compute. The batch dimension (50,176) is large and a multiple of 16, which is good. Overall, this configuration achieves approximately 80-85% of peak Tensor Core throughput, limited primarily by the non-multiple-of-16 inner dimension.

To improve Tensor Core utilization, modern CNN architectures use channel counts that are multiples of 16 or 32. ResNet-50 uses channel counts of 64, 128, 256, 512, and 1024, all of which are multiples of 16. EfficientNet uses channel counts like 32, 40, 80, 112, 192, 320, all chosen to be multiples of 8 or 16. These choices ensure that matrix dimensions align with Tensor Core tile sizes, maximizing hardware efficiency. The performance impact is substantial: a convolutional layer with 63 output channels achieves only 75% of the throughput of a layer with 64 output channels, despite having 98.4% as many parameters.

Batch size also affects Tensor Core utilization. For the RGB convolution example, the spatial dimension $H_{\text{out}}W_{\text{out}} = 50,176$ is large enough to saturate Tensor Cores even with batch size 1. However, for later layers with smaller spatial dimensions, batch size becomes critical. A layer with output dimensions $7 \times 7 = 49$ requires batch size at least 16 to provide sufficient parallelism ($49 \times 16 = 784$ output positions). With batch size 1, this layer achieves only 10-15% of peak throughput, as most Tensor Cores remain idle. With batch size 32, utilization increases to 60-70%, and with batch size

128, it reaches 85-90%. This scaling explains why larger batch sizes improve training throughput: they provide more parallelism to saturate the hardware.

5.5.2 cuDNN Optimizations

NVIDIA’s cuDNN library provides highly optimized implementations of convolutional operations, incorporating years of engineering effort to maximize performance on NVIDIA GPUs. Understanding cuDNN’s optimization strategies provides insight into how to design efficient CNN architectures.

cuDNN implements multiple convolution algorithms and automatically selects the fastest for each layer configuration. The primary algorithms are: (1) implicit GEMM (im2col-based matrix multiplication), (2) Winograd convolution for 3×3 and 5×5 kernels, (3) direct convolution for small batch sizes, and (4) FFT-based convolution for large kernels. For each forward pass, cuDNN benchmarks all applicable algorithms and caches the fastest choice, amortizing the benchmarking cost over many iterations.

For the RGB convolution example with 3×3 kernel, batch size 32, and 224×224 spatial dimensions, cuDNN typically selects Winograd convolution, which provides $1.5\text{-}2\times$ speedup over implicit GEMM. The Winograd algorithm reduces FLOPs from 173 MFLOPs to 77 MFLOPs ($2.25\times$ reduction) and achieves approximately 200 TFLOPS on an A100 GPU, or 64% of peak throughput. The lower-than-expected utilization (compared to 80-85% for implicit GEMM) arises because Winograd has higher memory bandwidth requirements and less regular computation patterns, making it harder to saturate Tensor Cores.

For later layers with smaller spatial dimensions, cuDNN typically selects implicit GEMM. A layer with output dimensions 7×7 , 512 input channels, 512 output channels, and 3×3 kernel has matrix multiplication dimensions $512 \times 4,608 \times 49$ (where $4,608 = 512 \times 9$). With batch size 32, the batch dimension becomes $49 \times 32 = 1,568$, giving dimensions $512 \times 4,608 \times 1,568$. This configuration achieves approximately 250 TFLOPS on an A100 GPU, or 80% of peak throughput. The high utilization arises because all dimensions are large and multiples of 16, providing excellent Tensor Core efficiency.

cuDNN also provides fused operations that combine multiple layers into a single GPU kernel, reducing memory traffic. A fused convolution-bias-ReLU kernel computes $\text{ReLU}(\text{Conv}(\mathbf{X}) + \mathbf{b})$ in a single pass, eliminating the need to write intermediate results to memory. For the RGB convolution example, this fusion reduces memory traffic from 375 MB (write conv output, read for bias add, write bias result, read for ReLU, write ReLU output) to 250 MB (write final output only), providing a $1.3\text{-}1.5\times$ speedup. cuDNN automatically applies these fusions when possible, but they require that the operation sequence be known at compile time.

The performance impact of cuDNN optimizations is dramatic. A naive convolution implementation in pure CUDA typically achieves 20-40 TFLOPS on an A100 GPU, or 6-13% of peak throughput. cuDNN’s optimized implementations achieve 200-280 TFLOPS, or 64-90% of peak throughput—a $5\text{-}10\times$ speedup. This performance gap explains why all modern deep learning frameworks (PyTorch, TensorFlow, JAX) use cuDNN as their backend for convolutional operations rather than implementing convolutions from scratch.

5.5.3 Memory Bandwidth vs Compute

Understanding whether a convolutional layer is compute-bound or memory-bandwidth-bound is essential for optimization. Compute-bound layers are limited by arithmetic throughput and benefit from algorithmic improvements like Winograd, while memory-bound layers are limited by data transfer rates and benefit from memory optimizations like fusion.

The arithmetic intensity of a convolutional layer is the ratio of FLOPs to bytes transferred: $\text{AI} = \text{FLOPs}/\text{bytes}$. For a layer to be compute-bound on an A100 GPU with 312 TFLOPS FP16 compute and 1.5 TB/s memory bandwidth, the arithmetic intensity must exceed $312 \times 10^{12} / (1.5 \times 10^{12}) = 208$ FLOPs per byte. Layers with lower arithmetic intensity are memory-bound.

For the RGB convolution example, the FLOPs are 173 MFLOPs. The memory transfers include reading the input ($3 \times 224 \times 224 \times 2 = 301,056$ bytes in FP16), reading the kernel ($64 \times 3 \times 9 \times 2 = 3,456$ bytes), and writing the output ($64 \times 224 \times 224 \times 2 = 6,422,528$ bytes), totaling 6.7 MB. The arithmetic

intensity is $173 \times 10^6 / (6.7 \times 10^6) = 25.8$ FLOPs per byte, which is far below the 208 threshold. This layer is memory-bound: it spends most of its time waiting for data transfers rather than computing.

The memory-bound nature of this layer explains why Winograd provides less than the theoretical $2.25\times$ speedup: reducing FLOPs from 173 MFLOPs to 77 MFLOPs ($2.25\times$ reduction) does not proportionally reduce runtime because the layer is limited by memory bandwidth, not compute. The actual speedup is approximately $1.5\times$, as Winograd reduces some memory traffic through better cache utilization but cannot eliminate the fundamental memory bottleneck.

In contrast, later layers with larger channel counts and smaller spatial dimensions are typically compute-bound. A layer with 512 input channels, 512 output channels, 7×7 spatial dimensions, and 3×3 kernel has $2 \times 512 \times 512 \times 9 \times 7 \times 7 = 230,686,720$ FLOPs. The memory transfers include reading the input ($512 \times 7 \times 7 \times 2 = 50,176$ bytes), reading the kernel ($512 \times 512 \times 9 \times 2 = 4,718,592$ bytes), and writing the output ($512 \times 7 \times 7 \times 2 = 50,176$ bytes), totaling 4.8 MB. The arithmetic intensity is $230,686,720 / (4.8 \times 10^6) = 48$ FLOPs per byte, which is still below the 208 threshold but much higher than the early layer. With batch size 32, the spatial dimension becomes $7 \times 7 \times 32 = 1,568$, and the arithmetic intensity increases to $48 \times 32 = 1,536$ FLOPs per byte, making the layer strongly compute-bound.

The transition from memory-bound to compute-bound as networks deepen has important implications for optimization. Early layers benefit from memory optimizations like fused operations and efficient data layouts, while late layers benefit from compute optimizations like Tensor Core utilization and algorithmic improvements. Profiling tools like NVIDIA Nsight Systems can identify whether each layer is memory-bound or compute-bound, guiding optimization efforts.

5.5.4 Batch Size Impact on Convolution Performance

Batch size is the primary lever for controlling GPU utilization in convolutional networks. Larger batches provide more parallelism, improving hardware efficiency and throughput measured in images per second. However, larger batches also require more memory and may affect convergence.

For ResNet-50 on an A100 GPU, the relationship between batch size and throughput is approximately logarithmic. With batch size 1, ResNet-50 achieves approximately 140 images per second. At batch size 8, throughput increases to 680 images per second ($4.9\times$ improvement). At batch size 32, throughput reaches 1,920 images per second ($2.8\times$ improvement). At batch size 128, throughput reaches 3,840 images per second ($2.0\times$ improvement). The diminishing returns arise because larger batches improve GPU utilization but eventually become limited by memory bandwidth and kernel launch overhead.

The memory cost of larger batches scales linearly with batch size for parameters and optimizer states (which are independent of batch size) but linearly for activations. For ResNet-50, batch size 1 requires approximately 1.2 GB of GPU memory (100 MB for parameters, 1.1 GB for activations). Batch size 8 requires 3.8 GB (100 MB parameters, 3.7 GB activations). Batch size 32 requires 12.4 GB (100 MB parameters, 12.3 GB activations). Batch size 128 requires 47.6 GB (100 MB parameters, 47.5 GB activations). An A100 GPU with 80 GB of memory can accommodate batch size 128 for ResNet-50, but larger batches require gradient accumulation or distributed training.

The optimal batch size balances throughput, memory usage, and convergence. From a hardware efficiency perspective, larger batches are always better, as they improve GPU utilization and images-per-second throughput. However, from an optimization perspective, very large batches can slow convergence by reducing the number of parameter updates per epoch. Empirically, batch sizes of 256-1024 work well for ResNet-50 on ImageNet, providing good hardware efficiency (70-85% GPU utilization) while maintaining reasonable convergence speed. Larger batches require careful tuning of learning rate and warmup schedule to maintain training stability and final model accuracy.

5.6 Pooling Layers

Definition 5.3 (Max Pooling). For window $k \times k$ and stride s :

$$\text{MaxPool}(\mathbf{X})_{i,j} = \max_{m,n \in \text{window}} \mathbf{X}_{si+m, sj+n} \quad (5.5)$$

Pooling reduces spatial dimensions, increases receptive field, and provides translation invariance.

5.7 Classic Architectures

5.7.1 VGG-16 (2014)

Deep network with small 3×3 filters. Pattern: $[\text{Conv}3 \times 3]^n \rightarrow \text{MaxPool} \rightarrow \text{Double channels}$
Total: 138 million parameters

5.7.2 ResNet (2015)

Definition 5.4 (Residual Block). Learn residual:

$$\mathbf{y} = \mathcal{F}(\mathbf{x}) + \mathbf{x} \quad (5.6)$$

ResNet-50: 25.6M parameters, enables training 100+ layer networks.

Key Point 5.2. *Residual connections enable extremely deep networks by allowing gradients to flow through skip connections. Analogous to skip connections in transformers.*

5.8 Batch Normalization

Definition 5.5 (Batch Normalization). For mini-batch, normalize each feature:

$$\hat{\mathbf{x}}_i = \frac{\mathbf{x}_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} \quad (5.7)$$

$$\mathbf{y}_i = \gamma \hat{\mathbf{x}}_i + \beta \quad (5.8)$$

where γ, β are learnable.

Benefits: Reduces covariate shift, allows higher learning rates, acts as regularization.

5.9 Exercises

Exercise 5.1. For $32 \times 32 \times 3$ input, compute dimensions after: $\text{Conv}(64, 5 \times 5, s=1, p=2)$, $\text{MaxPool}(2 \times 2, s=2)$, $\text{Conv}(128, 3 \times 3, s=1, p=1)$, $\text{MaxPool}(2 \times 2, s=2)$. Count parameters.

Exercise 5.2. Show two 3×3 convolutions equal one 5×5 receptive field. Compare parameter counts.

Exercise 5.3. Design CNN for CIFAR-10 with 3 blocks, channels [64, 128, 256]. Calculate total parameters.

Chapter 6

Recurrent Neural Networks

Chapter Overview

Recurrent Neural Networks (RNNs) process sequential data by maintaining hidden states that capture information from previous time steps. This chapter develops RNNs from basic recurrence to modern architectures like LSTMs and GRUs, establishing foundations for understanding transformers.

Learning Objectives

1. Understand recurrent architectures for sequential data
2. Implement vanilla RNNs, LSTMs, and GRUs
3. Understand vanishing/exploding gradient problems
4. Apply RNNs to sequence modeling tasks
5. Understand bidirectional and multi-layer RNNs

6.1 Vanilla RNNs

Definition 6.1 (Recurrent Neural Network). An RNN processes sequence $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T$ by maintaining hidden state $\mathbf{h}_t \in \mathbb{R}^h$:

$$\mathbf{h}_t = \tanh(\mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{W}_{xh}\mathbf{x}_t + \mathbf{b}_h) \quad (6.1)$$

$$\mathbf{y}_t = \mathbf{W}_{hy}\mathbf{h}_t + \mathbf{b}_y \quad (6.2)$$

where:

- $\mathbf{W}_{hh} \in \mathbb{R}^{h \times h}$: hidden-to-hidden weights
- $\mathbf{W}_{xh} \in \mathbb{R}^{h \times d}$: input-to-hidden weights
- $\mathbf{W}_{hy} \in \mathbb{R}^{k \times h}$: hidden-to-output weights
- \mathbf{h}_0 initialized (often zeros)

Example 6.1 (RNN Forward Pass). Character-level language model with vocabulary size $V = 5$, hidden size $h = 3$.

Input sequence: "hello" encoded as one-hot vectors $\mathbf{x}_1, \dots, \mathbf{x}_5 \in \mathbb{R}^5$

Initialize: $\mathbf{h}_0 = [0, 0, 0]^\top$

Time step 1: Process 'h'

$$\mathbf{h}_1 = \tanh(\mathbf{W}_{hh}\mathbf{h}_0 + \mathbf{W}_{xh}\mathbf{x}_1 + \mathbf{b}_h) \in \mathbb{R}^3 \quad (6.3)$$

$$\mathbf{y}_1 = \mathbf{W}_{hy}\mathbf{h}_1 + \mathbf{b}_y \in \mathbb{R}^5 \quad (6.4)$$

$$\hat{\mathbf{p}}_1 = \text{softmax}(\mathbf{y}_1) \quad (\text{predict next character}) \quad (6.5)$$

Time step 2: Process 'e' using \mathbf{h}_1

$$\mathbf{h}_2 = \tanh(\mathbf{W}_{hh}\mathbf{h}_1 + \mathbf{W}_{xh}\mathbf{x}_2 + \mathbf{b}_h) \quad (6.6)$$

Hidden state \mathbf{h}_t carries information from all previous time steps.

6.1.1 Backpropagation Through Time (BPTT)

[H] Backpropagation Through Time Sequence $\{\mathbf{x}_1, \dots, \mathbf{x}_T\}$, targets $\{\mathbf{y}_1, \dots, \mathbf{y}_T\}$ Gradients for all parameters

Forward Pass $t = 1 \dots T$ $\mathbf{h}_t = \tanh(\mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{W}_{xh}\mathbf{x}_t + \mathbf{b}_h)$

$\mathbf{y}_t = \mathbf{W}_{hy}\mathbf{h}_t + \mathbf{b}_y$

$L_t = \text{Loss}(\mathbf{y}_t, \text{target}_t)$

Backward Pass Initialize $\frac{\partial L}{\partial \mathbf{h}_{T+1}} = \mathbf{0}$

$t = T \dots 1$ Compute $\frac{\partial L}{\partial \mathbf{h}_t}$ (includes gradient from $t + 1$)

Accumulate $\frac{\partial L}{\partial \mathbf{W}_{hh}}, \frac{\partial L}{\partial \mathbf{W}_{xh}}, \frac{\partial L}{\partial \mathbf{W}_{hy}}$

6.1.2 Vanishing and Exploding Gradients

The fundamental challenge in training RNNs on long sequences arises from the multiplicative nature of gradient backpropagation through time. When computing gradients with respect to early hidden states, the chain rule requires multiplying Jacobian matrices across all intermediate time steps, leading to exponential growth or decay of gradient magnitudes.

The gradient of the loss with respect to an early hidden state \mathbf{h}_0 involves the product of Jacobians across all time steps:

$$\frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} = \prod_{t=1}^T \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} = \prod_{t=1}^T \mathbf{W}_{hh}^\top \text{diag}(\tanh'(\mathbf{z}_t)) \quad (6.7)$$

where \mathbf{z}_t is the pre-activation at time t . Each Jacobian $\frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}}$ has spectral norm bounded by $\|\mathbf{W}_{hh}\| \cdot \|\text{diag}(\tanh'(\mathbf{z}_t))\|$. Since $\tanh'(z) \in (0, 1]$ with maximum value 1 at $z = 0$, the derivative term is at most 1 and typically much smaller for saturated activations. This means the Jacobian norm is approximately $\|\mathbf{W}_{hh}\|$ in the best case.

For a sequence of length $T = 100$, if $\|\mathbf{W}_{hh}\| = 0.95$ (slightly less than 1), the gradient magnitude decays as $0.95^{100} \approx 0.006$, reducing gradients by a factor of 167. If $\|\mathbf{W}_{hh}\| = 0.9$, the decay is $0.9^{100} \approx 2.7 \times 10^{-5}$, reducing gradients by a factor of 37,000. This exponential decay makes it nearly impossible for the network to learn long-range dependencies: the gradient signal from time step 100 is effectively zero by the time it reaches time step 0. In practice, vanilla RNNs struggle to learn dependencies longer than 10-20 time steps due to vanishing gradients.

Conversely, if $\|\mathbf{W}_{hh}\| = 1.05$, the gradient magnitude grows as $1.05^{100} \approx 131.5$, amplifying gradients by a factor of 131. If $\|\mathbf{W}_{hh}\| = 1.1$, the growth is $1.1^{100} \approx 13,781$, causing gradients to explode. Exploding gradients lead to numerical overflow (NaN values) and training instability, where loss suddenly spikes to infinity. While gradient clipping (capping gradient norms at a threshold like 1.0) provides a practical solution for exploding gradients, it does not address the fundamental problem of vanishing gradients.

The vanishing gradient problem is particularly severe because the spectral norm of \mathbf{W}_{hh} must be precisely 1.0 to avoid both vanishing and exploding gradients, and maintaining this property during training is extremely difficult. Initialization schemes like orthogonal initialization set \mathbf{W}_{hh} to have

spectral norm 1.0 initially, but gradient descent updates quickly perturb this property. Even with careful initialization, vanilla RNNs rarely learn dependencies beyond 20-30 time steps in practice.

6.1.3 Quantitative Analysis of Gradient Decay

To understand the severity of vanishing gradients, consider a concrete example with BERT-base dimensions. Suppose we have a vanilla RNN with hidden dimension $h = 768$ (matching BERT-base) and sequence length $n = 512$ (BERT's maximum sequence length). The recurrence matrix $\mathbf{W}_{hh} \in \mathbb{R}^{768 \times 768}$ has 589,824 parameters. If we initialize \mathbf{W}_{hh} orthogonally (spectral norm exactly 1.0) and the tanh derivatives average 0.5 (typical for non-saturated activations), the effective Jacobian norm per time step is approximately $1.0 \times 0.5 = 0.5$.

Over 512 time steps, the gradient magnitude decays as $0.5^{512} \approx 10^{-154}$, which is far below machine precision for FP32 (approximately 10^{-38}) or even FP64 (approximately 10^{-308}). The gradient effectively becomes exactly zero after about 130 time steps in FP32 or 1,000 time steps in FP64. This means a vanilla RNN cannot learn any dependencies spanning more than 130 tokens when using FP32 arithmetic, regardless of optimization algorithm or learning rate. The mathematical structure of the recurrence fundamentally limits the learnable dependency length.

For comparison, consider the gradient flow in a transformer with the same dimensions. The self-attention mechanism computes attention scores $\mathbf{A} = \text{softmax}(\frac{\mathbf{QK}^\top}{\sqrt{d_k}})$ and outputs $\mathbf{O} = \mathbf{AV}$. The gradient $\frac{\partial L}{\partial \mathbf{V}}$ flows directly from the output through the attention weights, without any multiplicative accumulation across time steps. The gradient magnitude remains approximately constant regardless of sequence length, enabling transformers to learn dependencies spanning thousands of tokens. This fundamental difference in gradient flow explains why transformers replaced RNNs for nearly all sequence modeling tasks: they solve the vanishing gradient problem by design.

The LSTM architecture addresses vanishing gradients through its cell state mechanism, which provides an additive path for gradient flow. The cell state update $\mathbf{c}_t = \mathbf{c}_{t-1} \odot \mathbf{t} + \mathbf{t} \odot \tilde{\mathbf{c}}_t$ includes an additive term rather than purely multiplicative updates. The gradient with respect to \mathbf{c}_{t-1} is:

$$\frac{\partial \mathbf{c}_t}{\partial \mathbf{c}_{t-1}} = \text{diag}(\mathbf{t}) \quad (6.8)$$

which is a diagonal matrix with entries in $(0, 1)$ determined by the forget gate. If the forget gate learns to output values close to 1 for important information, the gradient can flow backward through many time steps without vanishing. However, this requires the network to learn appropriate forget gate values, and in practice, LSTMs still struggle with dependencies beyond 100-200 time steps. The cell state provides a highway for gradients, but it does not eliminate the vanishing gradient problem entirely.

6.2 Long Short-Term Memory (LSTM)

Definition 6.2 (LSTM Cell). LSTM uses gating mechanisms to control information flow:

$$\mathbf{t} = \sigma(\mathbf{W}_f[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_f) \quad (\text{forget gate}) \quad (6.9)$$

$$\mathbf{t} = \sigma(\mathbf{W}_i[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_i) \quad (\text{input gate}) \quad (6.10)$$

$$\tilde{\mathbf{c}}_t = \tanh(\mathbf{W}_c[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_c) \quad (\text{candidate cell}) \quad (6.11)$$

$$\mathbf{c}_t = \mathbf{c}_{t-1} \odot \mathbf{t} + \mathbf{t} \odot \tilde{\mathbf{c}}_t \quad (\text{cell state}) \quad (6.12)$$

$$\mathbf{t} = \sigma(\mathbf{W}_o[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_o) \quad (\text{output gate}) \quad (6.13)$$

$$\mathbf{h}_t = \mathbf{c}_t \odot \tanh(\mathbf{c}_t) \quad (\text{hidden state}) \quad (6.14)$$

where σ is sigmoid, \odot is element-wise multiplication, and $[\cdot, \cdot]$ is concatenation.

Key components:

- **Cell state \mathbf{c}_t :** Long-term memory, flows with minimal modification
- **Forget gate \mathbf{i}_t :** What to remove from cell state
- **Input gate \mathbf{i}_t :** What new information to store
- **Output gate \mathbf{o}_t :** What to output from cell state

Example 6.2 (LSTM Parameter Count). For input dimension $d = 512$ and hidden dimension $h = 1024$:

Each gate has weight matrix for $[\mathbf{h}_{t-1}, \mathbf{x}_t] \in \mathbb{R}^{h+d}$:

$$\text{Single gate: } (h + d) \times h + h = (1024 + 512) \times 1024 + 1024 \quad (6.15)$$

$$= 1,572,864 + 1,024 = 1,573,888 \quad (6.16)$$

LSTM has 4 gates (forget, input, cell, output):

$$\text{Total: } 4 \times 1,573,888 = 6,295,552 \text{ parameters} \quad (6.17)$$

Compare to transformer attention with same dimensions: often fewer parameters and better parallelization!

6.2.1 LSTM Computational Analysis

Understanding the computational cost of LSTMs is essential for comparing them to transformers and explaining why transformers have become dominant despite LSTMs' theoretical advantages for sequential processing. The LSTM's gating mechanisms provide powerful modeling capabilities but come with significant computational overhead that limits their efficiency on modern hardware.

For an LSTM with input dimension d and hidden dimension h , each time step requires computing four gates (forget, input, candidate, output), each involving a matrix multiplication with the concatenated input $[\mathbf{h}_{t-1}, \mathbf{x}_t] \in \mathbb{R}^{h+d}$. The computational cost per time step is:

$$\text{FLOPs per step} = 4 \times 2h(h + d) = 8h(h + d) \quad (6.18)$$

where the factor of 2 accounts for multiply-accumulate operations, and the factor of 4 accounts for the four gates. For BERT-base dimensions with $d = h = 768$, this gives $8 \times 768 \times (768 + 768) = 9,437,184$ FLOPs per time step. For a sequence of length $n = 512$, the total cost is $512 \times 9,437,184 = 4,831,838,208$ FLOPs, or approximately 4.8 GFLOPs.

This computational cost is deceptively modest compared to transformers. A single transformer layer with the same dimensions requires approximately 12.9 GFLOPs for self-attention (with $n = 512$) plus 9.4 GFLOPs for the feed-forward network, totaling 22.3 GFLOPs—about $4.6\times$ more than the LSTM. However, this comparison is misleading because it ignores the critical difference in parallelization: the transformer can process all 512 positions simultaneously, while the LSTM must process them sequentially.

The sequential nature of LSTMs means that the 4.8 GFLOPs cannot be parallelized across time steps. On an NVIDIA A100 GPU with peak throughput of 312 TFLOPS (FP16), the theoretical minimum time to process a sequence of length 512 is $\frac{4.8 \times 10^9}{312 \times 10^{12}} = 15.4$ microseconds if we could achieve perfect parallelization. However, the sequential dependency forces us to process one time step at a time, with each step taking approximately $\frac{9.4 \times 10^6}{312 \times 10^{12}} = 0.03$ microseconds at peak throughput. In practice, small matrix multiplications achieve only 1-5% of peak throughput due to insufficient parallelism, so each time step actually takes approximately 1-3 microseconds, giving a total sequence processing time of 512-1,536 microseconds (0.5-1.5 milliseconds).

For comparison, a transformer layer can process the entire sequence in parallel. The self-attention computation requires three matrix multiplications ($\mathbf{Q} = \mathbf{XW}_Q$, $\mathbf{K} = \mathbf{XW}_K$, $\mathbf{V} = \mathbf{XW}_V$) with dimensions $512 \times 768 \times 768$, followed by the attention score computation $\mathbf{A} = \text{softmax}(\frac{\mathbf{QK}^\top}{\sqrt{d_k}})$ and output

computation $\mathbf{O} = \mathbf{A}\mathbf{V}$. These operations can be batched into large matrix multiplications that achieve 40-60% of peak GPU throughput, completing in approximately 50-100 microseconds total. The transformer is 5-30 \times faster than the LSTM despite having more FLOPs, purely due to better parallelization.

The memory requirements for LSTM hidden states are modest compared to transformer attention matrices. For batch size B and sequence length n , the LSTM must store hidden states $\mathbf{h}_t \in \mathbb{R}^{B \times h}$ and cell states $\mathbf{c}_t \in \mathbb{R}^{B \times h}$ for each time step, requiring $2Bnh \times 4 = 8Bnh$ bytes in FP32. For BERT-base dimensions with $B = 32$, $n = 512$, $h = 768$, this totals $8 \times 32 \times 512 \times 768 = 100,663,296$ bytes, or approximately 96 MB. This is substantially less than the 384 MB required for transformer attention scores in a single layer, making LSTMs more memory-efficient for long sequences.

However, this memory advantage is offset by the sequential processing requirement. While transformers can trade memory for speed by using gradient checkpointing (recomputing activations during the backward pass rather than storing them), LSTMs cannot benefit from this technique as effectively because the sequential dependency prevents parallelization of the recomputation. Gradient checkpointing reduces transformer memory by 3-5 \times with only 20-30% slowdown, but for LSTMs, the slowdown is 2-3 \times because the recomputation cannot be parallelized. This makes gradient checkpointing less attractive for LSTMs, limiting their ability to scale to very long sequences.

6.3 Gated Recurrent Unit (GRU)

Definition 6.3 (GRU Cell). GRU simplifies LSTM by merging cell and hidden states:

$$\mathbf{z}_t = \sigma(\mathbf{W}_z[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_z) \quad (\text{update gate}) \quad (6.19)$$

$$r_t = \sigma(\mathbf{W}_r[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_r) \quad (\text{reset gate}) \quad (6.20)$$

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_h[r_t \odot \mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_h) \quad (\text{candidate}) \quad (6.21)$$

$$\mathbf{h}_t = (1 - \mathbf{z}_t) \odot \mathbf{h}_{t-1} + \mathbf{z}_t \odot \tilde{\mathbf{h}}_t \quad (\text{hidden state}) \quad (6.22)$$

Advantages over LSTM:

- Fewer parameters (3 gates vs 4)
- Simpler architecture
- Often similar performance
- Faster training

6.4 Bidirectional RNNs

Definition 6.4 (Bidirectional RNN). Process sequence in both directions:

$$\vec{\mathbf{h}}_t = \text{RNN}_{\text{forward}}(\mathbf{x}_t, \vec{\mathbf{h}}_{t-1}) \quad (6.23)$$

$$\overleftarrow{\mathbf{h}}_t = \text{RNN}_{\text{backward}}(\mathbf{x}_t, \overleftarrow{\mathbf{h}}_{t+1}) \quad (6.24)$$

$$\mathbf{h}_t = [\vec{\mathbf{h}}_t; \overleftarrow{\mathbf{h}}_t] \quad (6.25)$$

Bidirectional RNNs capture context from both past and future, useful when entire sequence is available (not for online/causal tasks).

Example: BERT uses bidirectional transformers (attention, not RNN), capturing full context.

6.5 RNN Applications

Sequence-to-Sequence:

- Machine translation: Encoder RNN \rightarrow Decoder RNN
- Text summarization
- Speech recognition

Sequence Labeling:

- Part-of-speech tagging
- Named entity recognition
- Output at each time step

Sequence Generation:

- Language modeling
- Music generation
- Sample from output distribution

6.6 RNNs vs Transformers: A Computational Comparison

The transition from RNNs to transformers represents one of the most significant architectural shifts in deep learning history. While RNNs were the dominant architecture for sequence modeling from the 1990s through 2017, transformers have almost entirely replaced them for natural language processing tasks. Understanding the computational, memory, and hardware efficiency differences between these architectures explains this dramatic shift and provides insight into modern deep learning system design.

6.6.1 Sequential vs Parallel Computation

The fundamental difference between RNNs and transformers lies in their computational structure. RNNs process sequences sequentially, computing \mathbf{h}_t from \mathbf{h}_{t-1} and \mathbf{x}_t at each time step. This sequential dependency means that position t cannot be computed until position $t - 1$ is complete, preventing parallelization across time steps. For a sequence of length n , the RNN requires n sequential operations, each taking time T_{step} , for a total time of nT_{step} . Even with infinite computational resources, this sequential bottleneck cannot be overcome.

Transformers, by contrast, compute all positions simultaneously using self-attention. The query, key, and value projections $\mathbf{Q} = \mathbf{XW}_Q$, $\mathbf{K} = \mathbf{XW}_K$, $\mathbf{V} = \mathbf{XW}_V$ are batched matrix multiplications that process all n positions in parallel. The attention scores $\mathbf{A} = \text{softmax}(\frac{\mathbf{QK}^\top}{\sqrt{d_k}})$ and outputs $\mathbf{O} = \mathbf{AV}$ similarly operate on all positions simultaneously. The total computation time is independent of sequence length in terms of sequential depth—all positions are processed in a constant number of parallel operations.

For concrete comparison, consider processing a sequence of length $n = 512$ with model dimension $d = 768$ (BERT-base dimensions) on an NVIDIA A100 GPU. An LSTM requires 512 sequential steps, each taking approximately 1-3 microseconds (as computed in Section 6.2), for a total time of 512-1,536 microseconds (0.5-1.5 milliseconds). A transformer layer processes the entire sequence in approximately 50-100 microseconds—5-30 \times faster despite having more total FLOPs. The speedup comes entirely from parallelization: the transformer exploits the GPU's 6,912 CUDA cores to process all 512 positions simultaneously, while the LSTM can only process one position at a time.

This parallelization advantage scales with sequence length. For $n = 2048$ (GPT-2's maximum length), the LSTM time increases to 2-6 milliseconds, while the transformer time increases to approximately 200-400 microseconds—still 5-30 \times faster. For $n = 4096$, the LSTM requires 4-12 milliseconds,

while the transformer requires 800-1,600 microseconds. The transformer’s advantage is consistent across sequence lengths because both architectures scale linearly with n in terms of FLOPs, but the transformer can parallelize while the LSTM cannot.

6.6.2 Memory Complexity: $O(nd)$ vs $O(n^2)$

The memory requirements of RNNs and transformers scale differently with sequence length, leading to different bottlenecks for long sequences. RNNs store hidden states $\mathbf{h}_t \in \mathbb{R}^{B \times d}$ for each time step, requiring $O(Bnd)$ memory for a batch of size B with sequence length n and hidden dimension d . This linear scaling with sequence length makes RNNs memory-efficient for long sequences: doubling the sequence length doubles the memory requirement.

Transformers store attention score matrices $\mathbf{A} \in \mathbb{R}^{B \times h \times n \times n}$ for each layer, requiring $O(Bhn^2)$ memory where h is the number of attention heads. This quadratic scaling with sequence length makes transformers memory-intensive for long sequences: doubling the sequence length quadruples the memory requirement. For BERT-base with $B = 32$, $h = 12$, $n = 512$, attention scores require 384 MB per layer, or 4.6 GB across 12 layers. For $n = 2048$, this increases to 6.1 GB per layer, or 73.7 GB across 12 layers—exceeding the memory of most GPUs.

However, this comparison is incomplete because it ignores the different memory access patterns. RNNs must load the recurrence matrix $\mathbf{W}_{hh} \in \mathbb{R}^{d \times d}$ from memory at each time step, requiring n memory loads of size d^2 . For BERT-base dimensions with $d = 768$ and $n = 512$, this totals $512 \times 768^2 \times 4 = 1,207,959,552$ bytes, or approximately 1.15 GB of memory bandwidth. Transformers load the attention weight matrices $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{d \times d}$ once per layer, requiring $3d^2$ memory loads. For the same dimensions, this totals $3 \times 768^2 \times 4 = 7,077,888$ bytes, or approximately 6.75 MB—170× less memory bandwidth than the LSTM.

The memory bandwidth difference is critical for understanding hardware efficiency. Modern GPUs have memory bandwidth of 1-2 TB/s (e.g., A100 has 1.6 TB/s), which limits the rate at which data can be loaded from GPU memory to compute units. For the LSTM, loading 1.15 GB of weights requires $\frac{1.15 \text{ GB}}{1.6 \text{ TB/s}} = 0.72$ milliseconds, which is comparable to the 0.5-1.5 milliseconds of compute time. This means the LSTM is memory-bandwidth-bound: the GPU spends as much time loading weights as performing computation, achieving only 50% efficiency. For the transformer, loading 6.75 MB requires $\frac{6.75 \text{ MB}}{1.6 \text{ TB/s}} = 4.2$ microseconds, which is negligible compared to the 50-100 microseconds of compute time. The transformer is compute-bound, achieving 90-95% efficiency.

6.6.3 Training Time Comparison

The combined effects of parallelization and memory bandwidth lead to dramatic differences in training time between RNNs and transformers. For BERT-base training on the BookCorpus and English Wikipedia datasets (approximately 3.3 billion words, or 16 billion tokens with WordPiece tokenization), the original BERT paper reports training time of 4 days on 16 Cloud TPU chips (64 TPU cores total). Each TPU core has peak throughput of 45 TFLOPS (FP16), giving a total of 2,880 TFLOPS across 64 cores.

To train an LSTM with equivalent capacity (110 million parameters) on the same dataset, we can estimate the training time based on the computational and memory bandwidth analysis above. An LSTM with hidden dimension $d = 768$ has approximately $4 \times 768 \times (768 + 768) = 4,718,592$ parameters per layer (4 gates, each with weight matrix for concatenated input). To reach 110 million parameters, we need approximately $\frac{110,000,000}{4,718,592} \approx 23$ layers. Each layer requires 4.8 GFLOPs per sequence of length 512, so 23 layers require 110.4 GFLOPs per sequence.

For 16 billion tokens with sequence length 512, we have $\frac{16 \times 10^9}{512} = 31,250,000$ sequences. At 110.4 GFLOPs per sequence, the total computation is $31,250,000 \times 110.4 \times 10^9 = 3.45 \times 10^{18}$ FLOPs, or 3.45 exaFLOPs. At 2,880 TFLOPS peak throughput, this would take $\frac{3.45 \times 10^{18}}{2,880 \times 10^{12}} = 1,197,917$ seconds, or approximately 14 days, if we could achieve 100% hardware utilization.

However, as discussed above, LSTMs achieve only 1-5% of peak throughput due to sequential processing and memory bandwidth limitations. At 3% utilization, the training time increases to $\frac{14 \text{ days}}{0.03} = 467$ days—more than a year. Even with aggressive optimizations, LSTM training would

likely require 100-200 days on the same hardware that trains BERT in 4 days. This $25\text{-}50\times$ slowdown makes LSTMs impractical for large-scale pretraining, explaining why transformers have completely replaced RNNs for modern language models.

The training time difference becomes even more extreme for larger models. GPT-3 with 175 billion parameters was trained on 300 billion tokens in approximately 34 days on 10,000 NVIDIA V100 GPUs. An LSTM with equivalent capacity would require an estimated 3-5 years on the same hardware, making it economically infeasible. The transformer’s parallelization advantage is the primary enabler of large-scale language model pretraining.

6.6.4 Hardware Limitations of RNNs

The poor hardware efficiency of RNNs stems from three fundamental limitations: insufficient parallelism for GPU saturation, memory bandwidth bottlenecks, and poor cache utilization. Modern GPUs are designed for massively parallel workloads with thousands of concurrent operations, but RNNs provide only batch-level parallelism. For batch size $B = 32$ and hidden dimension $d = 768$, the LSTM matrix multiplication $\mathbf{W}[\mathbf{h}_{t-1}, \mathbf{x}_t]$ has dimensions $768 \times (32 \times 1536)$, which can process only 32 rows in parallel. An A100 GPU with 6,912 CUDA cores can theoretically process 6,912 operations simultaneously, but the LSTM provides only 32 concurrent operations—leaving 99.5% of the GPU idle.

Increasing batch size improves GPU utilization but has diminishing returns. At batch size 256, the LSTM achieves approximately 5-10% of peak throughput. At batch size 2048, it reaches 15-20% of peak throughput. However, batch sizes beyond 256 are impractical for most training scenarios due to memory constraints and optimization difficulties. Even at batch size 2048, the LSTM achieves only 15-20% utilization, compared to 40-60% for transformers at the same batch size. The sequential dependency fundamentally limits the LSTM’s ability to exploit GPU parallelism.

The memory bandwidth bottleneck exacerbates the parallelism problem. As computed above, the LSTM requires 1.15 GB of memory bandwidth per sequence, compared to 6.75 MB for the transformer—a $170\times$ difference. For batch size 32, the LSTM requires $32 \times 1.15 = 36.8$ GB of memory bandwidth, which takes $\frac{36.8 \text{ GB}}{1.6 \text{ TB/s}} = 23$ milliseconds to load. The actual compute time is approximately 16-48 milliseconds ($0.5\text{-}1.5$ milliseconds per sequence \times 32 sequences), so memory bandwidth accounts for 30-60% of the total time. The transformer requires $32 \times 6.75 = 216$ MB of memory bandwidth, taking only 0.135 milliseconds to load—negligible compared to the 1.6-3.2 milliseconds of compute time.

Cache utilization is similarly poor for RNNs. Modern GPUs have L1 cache (128 KB per streaming multiprocessor) and L2 cache (40 MB for A100) that can dramatically accelerate memory access for data that fits in cache. The transformer’s attention computation reuses the query, key, and value matrices across all positions, enabling effective cache utilization. For BERT-base, the $\mathbf{Q}, \mathbf{K}, \mathbf{V}$ matrices have size $32 \times 512 \times 768 \times 4 = 50.3$ MB, which fits in L2 cache. The LSTM’s recurrence matrix \mathbf{W}_{hh} has size $768 \times 768 \times 4 = 2.36$ MB, which also fits in cache, but it must be loaded at each time step, and the sequential dependency prevents batching these loads. The cache hit rate for LSTMs is typically 20-40%, compared to 60-80% for transformers, further reducing hardware efficiency.

6.6.5 Why Transformers Replaced RNNs

The dominance of transformers over RNNs can be summarized in three key advantages: parallelization, gradient flow, and hardware efficiency. The parallel computation structure of transformers enables $5\text{-}30\times$ faster training than RNNs on modern GPUs, making large-scale pretraining feasible. The direct attention connections enable gradient flow across arbitrary distances without vanishing, allowing transformers to learn dependencies spanning thousands of tokens. The high GPU utilization (40-60% vs 1-5% for RNNs) makes transformers $10\text{-}60\times$ more hardware-efficient, reducing training costs proportionally.

These advantages compound: the $5\text{-}30\times$ speedup from parallelization, combined with the $10\text{-}60\times$ improvement in hardware efficiency, yields an overall $50\text{-}1,800\times$ advantage for transformers. In practice, transformers train $100\text{-}500\times$ faster than RNNs for equivalent model capacity and dataset size. This speedup difference is the primary reason why transformers have completely replaced RNNs for natural language processing: the economic cost of training an LSTM-based language model is prohibitive compared to a transformer.

However, RNNs retain advantages for specific use cases. For online or streaming applications where inputs arrive sequentially and outputs must be produced in real-time, RNNs can process each input immediately without waiting for the full sequence. Transformers require the entire sequence to compute attention, making them unsuitable for true streaming applications. For very long sequences exceeding 10,000 tokens, the $O(n^2)$ memory complexity of transformers becomes prohibitive, while RNNs' $O(n)$ memory scaling remains manageable. For edge deployment on devices with limited memory and compute, RNNs' smaller memory footprint can be advantageous.

Despite these niche advantages, transformers dominate modern deep learning due to their superior training efficiency and scalability. The development of efficient attention mechanisms (Chapter 16) addresses the $O(n^2)$ memory bottleneck for long sequences, and techniques like streaming transformers enable online processing. The architectural innovations that made transformers successful—parallel computation, direct gradient flow, and hardware efficiency—represent fundamental advances that are unlikely to be superseded by sequential architectures.

Key Point 6.1. *While RNNs were dominant for sequences, transformers now excel in most NLP tasks due to: (1) Better parallelization enabling 5-30 \times faster training, (2) Direct long-range dependencies via attention avoiding vanishing gradients, (3) Superior hardware efficiency achieving 40-60% GPU utilization vs 1-5% for RNNs. The combined effect is 100-500 \times faster training, making transformers economically superior for large-scale pretraining. RNNs remain useful only for online/streaming tasks and extremely long sequences where $O(n^2)$ memory is prohibitive.*

6.7 Exercises

Exercise 6.1. For vanilla RNN with input dim $d = 128$, hidden dim $h = 256$, and sequence length $T = 50$: (1) Count total parameters in \mathbf{W}_{hh} , \mathbf{W}_{xh} , and \mathbf{W}_{hy} , (2) Compute total FLOPs for forward pass through all time steps, (3) Estimate GPU utilization on an A100 (312 TFLOPS peak) with batch size 32, assuming each time step achieves 2% of peak throughput. Why is utilization so low?

Exercise 6.2. Derive the gradient $\frac{\partial L}{\partial \mathbf{W}_{hh}}$ for a 3-step sequence. Show how the gradient involves products of Jacobians $\frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}}$. If $\|\mathbf{W}_{hh}\| = 0.9$ and \tanh' averages 0.5, compute the gradient magnitude decay factor from time step 3 to time step 0. At what sequence length would gradients vanish below FP32 precision (10^{-38})?

Exercise 6.3. Compare parameter counts and FLOPs per sequence for: (1) LSTM with $d = 512$, $h = 512$, $n = 512$, (2) GRU with $d = 512$, $h = 512$, $n = 512$, (3) Transformer attention layer with $d_{\text{model}} = 512$, $d_k = 64$, $h = 8$ heads, $n = 512$. Which architecture has the most parameters? Which has the most FLOPs? Which achieves the highest GPU utilization and why?

Exercise 6.4. Implement bidirectional LSTM in PyTorch for sequence "The cat sat on the mat" with vocabulary size 10, embedding dim 16, hidden dim 32. Process the sequence and show output dimensions. Compute the total memory required for hidden states and cell states in FP32. How does this compare to the memory required for attention scores in a transformer with the same

dimensions?

Exercise 6.5. For BERT-base dimensions ($d = 768$, $n = 512$), compute: (1) Memory required for LSTM hidden and cell states across 12 layers with batch size 32, (2) Memory required for transformer attention scores across 12 layers with batch size 32 and 12 attention heads, (3) The sequence length at which LSTM memory equals transformer memory. Explain why transformers are memory-limited for long sequences while LSTMs are compute-limited.

Exercise 6.6. Estimate the training time for a 110M parameter LSTM on 16 billion tokens (sequence length 512) using 64 TPU cores with 2,880 TFLOPS total peak throughput. Assume the LSTM achieves 3% of peak throughput due to sequential processing. Compare this to BERT-base training time of 4 days on the same hardware. What is the speedup factor? Explain the three main reasons for the difference: parallelization, memory bandwidth, and GPU utilization.

Part III

Attention Mechanisms

Chapter 7

Attention Mechanisms: Fundamentals

Chapter Overview

Attention mechanisms revolutionized sequence modeling by allowing models to focus on relevant parts of the input when producing each output. This chapter introduces attention from first principles, developing the query-key-value paradigm that underpins modern transformers.

Attention solves a fundamental limitation of RNN encoder-decoder models: compressing entire input sequence into single fixed-size vector. Instead, attention computes dynamic, context-dependent representations by weighted combination of all input positions.

Learning Objectives

1. Understand the motivation for attention in sequence-to-sequence models
2. Master the query-key-value attention paradigm
3. Implement additive (Bahdanau) and multiplicative (Luong) attention
4. Understand scaled dot-product attention
5. Compute attention weights and apply to values
6. Visualize and interpret attention distributions

7.1 Motivation: The Seq2Seq Bottleneck

7.1.1 RNN Encoder-Decoder Architecture

The sequence-to-sequence (seq2seq) problem requires mapping an input sequence $\mathbf{x}_1, \dots, \mathbf{x}_n$ to an output sequence $\mathbf{y}_1, \dots, \mathbf{y}_m$ of potentially different length. This formulation encompasses machine translation, text summarization, question answering, and many other natural language processing tasks. Before attention mechanisms, the standard approach used recurrent neural networks in an encoder-decoder architecture that suffered from a fundamental information bottleneck.

The encoder RNN processes the input sequence sequentially, updating its hidden state at each time step:

$$\mathbf{h}_t^{\text{enc}} = \text{RNN}(\mathbf{x}_t, \mathbf{h}_{t-1}^{\text{enc}}) \quad (7.1)$$

After processing all n input tokens, the final hidden state $\mathbf{c} = \mathbf{h}_n^{\text{enc}}$ serves as the context vector—a fixed-size representation intended to capture the entire input sequence. This context vector, typically 512 or 1024 dimensions for LSTM-based systems, must encode all relevant information from the source sequence regardless of its length.

The decoder RNN then generates the output sequence conditioned on this context vector:

$$\mathbf{h}_t^{\text{dec}} = \text{RNN}([\mathbf{y}_{t-1}, \mathbf{c}], \mathbf{h}_{t-1}^{\text{dec}}) \quad (7.2)$$

where $[\mathbf{y}_{t-1}, \mathbf{c}]$ denotes concatenation of the previous output token embedding and the context vector. The decoder must rely on this single fixed-size vector throughout the entire generation process, accessing the same \mathbf{c} when producing the first output word and the last.

The Information Bottleneck: Compressing an entire input sequence into a single fixed-size vector creates severe information loss, particularly for long sequences. Consider translating a 50-word English sentence to French. The encoder must compress 50 words of semantic content, syntactic structure, and contextual relationships into a 512-dimensional vector. This is fundamentally insufficient—the context vector becomes an information bottleneck that limits the model’s capacity to handle complex or lengthy inputs.

Empirical evidence from early neural machine translation systems demonstrated this limitation quantitatively. For English-French translation using LSTM encoder-decoders with 1000-dimensional hidden states, translation quality (measured by BLEU score) remained stable for source sentences up to 20-25 words but degraded significantly beyond this length. Sentences of 30-40 words showed BLEU score drops of 5-10 points compared to shorter sentences, and sentences exceeding 50 words often produced nearly incomprehensible translations. The fixed-size context vector simply could not retain sufficient information about long, complex source sentences.

Memory and Computational Characteristics: The RNN encoder-decoder architecture requires $O(n + m)$ memory for storing hidden states during the forward pass, where n is the source length and m is the target length. For a typical translation task with $n = 50$, $m = 50$, and hidden dimension $d_h = 1024$, this amounts to $(50 + 50) \times 1024 \times 4 = 400$ KB per sequence in FP32. However, the sequential nature of RNN processing prevents parallelization across time steps. Each hidden state \mathbf{h}_t depends on \mathbf{h}_{t-1} , forcing strictly sequential computation. On a GPU capable of processing thousands of operations in parallel, this sequential constraint severely limits throughput.

For a batch of 32 sequences, the encoder processes $32 \times 50 = 1600$ time steps sequentially, even though the GPU could theoretically process all 1600 in parallel if the operations were independent. This sequential bottleneck means RNN encoder-decoders achieve only 5-10% of peak GPU utilization during training, wasting the majority of available compute capacity.

7.1.2 Attention Solution

Attention mechanisms solve the information bottleneck by allowing the decoder to access all encoder hidden states directly, rather than relying on a single compressed representation. The key insight is that when generating each output word \mathbf{y}_t , different input words have different relevance. When translating “The cat sat on the mat” to French, generating “chat” (cat) should focus primarily on the input word “cat,” while generating “assis” (sat) should focus on “sat.” The decoder’s information needs change dynamically throughout generation.

Rather than computing a single context vector \mathbf{c} for the entire sequence, attention computes a different context vector \mathbf{c}_t for each output position t . This context vector is a weighted sum of all encoder hidden states:

$$\mathbf{c}_t = \sum_{i=1}^n \alpha_{t,i} \mathbf{h}_i^{\text{enc}} \quad (7.3)$$

where the attention weights $\alpha_{t,i}$ indicate how much the decoder should focus on input position i when generating output position t . These weights form a probability distribution: $\alpha_{t,i} \geq 0$ and $\sum_{i=1}^n \alpha_{t,i} = 1$.

The attention weights are computed dynamically based on the current decoder state \mathbf{s}_t and each encoder hidden state \mathbf{h}_i . This allows the model to learn which input positions are relevant for each output position, adapting the context vector to the decoder’s current needs. When generating the first word of a translation, the attention might focus on the beginning of the source sentence. When generating the last word, attention shifts to the end of the source.

Memory Trade-off: Attention increases memory requirements from $O(n + m)$ to $O(nm)$ because we must store attention weights $\alpha_{t,i}$ for all pairs of input and output positions. For translation with $n = 50$ and $m = 50$, this requires storing a $50 \times 50 = 2500$ element attention matrix. At 4 bytes per element (FP32), this is 10 KB per sequence—modest compared to the benefits. However, this quadratic scaling becomes significant for very long sequences. For document-level translation with $n = 1000$ and $m = 1000$, the attention matrix requires $1000^2 \times 4 = 4$ MB per sequence, or 128 MB for batch size 32.

Parallelization Benefit: The crucial advantage is that attention enables parallelization. Unlike RNN hidden states that must be computed sequentially, attention weights for all output positions can be computed simultaneously during training when the target sequence is known. This transforms the sequential $O(m)$ decoder steps into a single parallel operation, dramatically improving GPU utilization from 5-10% to 60-80% in practice.

Example 7.1 (Translation with Attention). Consider translating the English sentence "The cat sat on the mat" to French: "Le chat était assis sur le tapis." Without attention, the encoder compresses all six English words into a single 512-dimensional context vector, which the decoder uses to generate all seven French words. The context vector must simultaneously encode that "cat" translates to "chat," "sat" translates to "était assis," and "mat" translates to "tapis"—a challenging compression task.

With attention, when generating "chat" (cat), the attention mechanism computes weights that heavily favor the input position containing "cat." The attention distribution might be $[0.05, 0.82, 0.03, 0.02, 0.03, 0.05]$, placing 82% of the weight on position 2 (the word "cat"). The context vector \mathbf{c}_2 is then dominated by the encoder hidden state for "cat," providing the decoder with direct access to the relevant input information.

When generating "assis" (sat), the attention distribution shifts to $[0.03, 0.08, 0.75, 0.04, 0.05, 0.05]$, now focusing 75% on position 3 (the word "sat"). The context vector \mathbf{c}_4 adapts to provide information about "sat" rather than "cat." This dynamic reweighting allows the decoder to access different parts of the input as needed, eliminating the information bottleneck of the fixed context vector.

Empirically, attention-based translation systems improved BLEU scores by 3-5 points on standard benchmarks and maintained consistent quality even for sentences exceeding 50 words—a regime where RNN encoder-decoders failed catastrophically.

7.2 Additive Attention (Bahdanau)

Bahdanau attention, introduced in 2015 for neural machine translation, was the first widely successful attention mechanism. It computes attention weights using an additive scoring function that combines the decoder state and encoder hidden states through learned transformations. While later superseded by more efficient mechanisms, understanding Bahdanau attention provides crucial insights into attention design and the evolution toward modern transformers.

Definition 7.1 (Bahdanau Attention). Given encoder hidden states $\mathbf{h}_1, \dots, \mathbf{h}_n \in \mathbb{R}^{d_h}$ and decoder hidden state $\mathbf{s}_t \in \mathbb{R}^{d_s}$ at time t , Bahdanau attention computes a context vector through four steps:

Step 1: Compute alignment scores

$$e_{t,i} = \mathbf{v}^\top \tanh(\mathbf{W}_1 \mathbf{s}_t + \mathbf{W}_2 \mathbf{h}_i) \quad (7.4)$$

where $\mathbf{W}_1 \in \mathbb{R}^{d_a \times d_s}$, $\mathbf{W}_2 \in \mathbb{R}^{d_a \times d_h}$, $\mathbf{v} \in \mathbb{R}^{d_a}$, and d_a is the attention dimension (typically 256-512).

Step 2: Compute attention weights (softmax)

$$\alpha_{t,i} = \frac{\exp(e_{t,i})}{\sum_{j=1}^n \exp(e_{t,j})} \quad (7.5)$$

Step 3: Compute context vector

$$\mathbf{c}_t = \sum_{i=1}^n \alpha_{t,i} \mathbf{h}_i \quad (7.6)$$

Step 4: Use in decoder

$$\mathbf{s}_t = \text{RNN}([\mathbf{y}_{t-1}, \mathbf{c}_t], \mathbf{s}_{t-1}) \quad (7.7)$$

Computational Cost Analysis: The additive scoring function in Step 1 requires substantial computation for each query-key pair. For a single alignment score $e_{t,i}$, we must:

1. Compute $\mathbf{W}_1 \mathbf{s}_t$: $2d_a d_s$ FLOPs (matrix-vector multiplication)
2. Compute $\mathbf{W}_2 \mathbf{h}_i$: $2d_a d_h$ FLOPs
3. Add the results: d_a FLOPs
4. Apply tanh: $\approx 3d_a$ FLOPs (exponentials and divisions)
5. Compute $\mathbf{v}^\top(\cdot)$: $2d_a$ FLOPs

Total per alignment score: approximately $2d_a(d_s + d_h + 3)$ FLOPs. For a translation task with source length n and target length m , we compute nm alignment scores, requiring:

$$\text{Bahdanau alignment FLOPs} \approx 2nmd_a(d_s + d_h + 3) \quad (7.8)$$

For typical dimensions $n = 50$, $m = 50$, $d_a = 256$, $d_s = d_h = 512$:

$$2 \times 50 \times 50 \times 256 \times (512 + 512 + 3) \approx 1.3 \text{ billion FLOPs} \quad (7.9)$$

This is substantial, but the more critical issue is that these operations do not map efficiently to GPU hardware. The additive scoring function involves element-wise operations (tanh), vector additions, and small matrix-vector products that achieve poor utilization on GPUs optimized for large matrix multiplications. In practice, Bahdanau attention achieves only 15-25% of peak GPU throughput.

Memory Requirements: The attention mechanism requires storing:

- Encoder hidden states: $n \times d_h$ elements
- Alignment scores: $m \times n$ elements (for all decoder positions during training)
- Attention weights: $m \times n$ elements
- Intermediate activations: $m \times n \times d_a$ elements for the tanh layer

For $n = 50$, $m = 50$, $d_h = 512$, $d_a = 256$ in FP32:

$$\text{Encoder states: } 50 \times 512 \times 4 = 102 \text{ KB} \quad (7.10)$$

$$\text{Alignment scores: } 50 \times 50 \times 4 = 10 \text{ KB} \quad (7.11)$$

$$\text{Attention weights: } 50 \times 50 \times 4 = 10 \text{ KB} \quad (7.12)$$

$$\text{Intermediate: } 50 \times 50 \times 256 \times 4 = 2.5 \text{ MB} \quad (7.13)$$

The intermediate activations dominate memory usage, requiring 2.5 MB per sequence or 80 MB for batch size 32. This is manageable for short sequences but scales poorly to longer contexts.

Parameter Count: Bahdanau attention introduces $O(d_a(d_s + d_h))$ parameters:

$$\mathbf{W}_1 \in \mathbb{R}^{d_a \times d_s} : d_a d_s \text{ parameters} \quad (7.14)$$

$$\mathbf{W}_2 \in \mathbb{R}^{d_a \times d_h} : d_a d_h \text{ parameters} \quad (7.15)$$

$$\mathbf{v} \in \mathbb{R}^{d_a} : d_a \text{ parameters} \quad (7.16)$$

For $d_a = 256$, $d_s = d_h = 512$: $(256 \times 512) + (256 \times 512) + 256 = 262,400$ parameters. While not enormous, these parameters must be learned specifically for the attention mechanism, adding to the model's overall capacity requirements.

Key Point 7.1. Attention weights $\alpha_{t,i}$ form a probability distribution: $\alpha_{t,i} \geq 0$ and $\sum_{i=1}^n \alpha_{t,i} = 1$. This ensures the context vector \mathbf{c}_t is a convex combination of encoder states, interpolating between them rather than extrapolating. The softmax normalization is crucial for training stability—without it, attention weights could grow unbounded, causing gradient explosion.

Example 7.2 (Bahdanau Attention Computation). Consider a small example with encoder hidden states $\mathbf{h}_1, \mathbf{h}_2, \mathbf{h}_3 \in \mathbb{R}^4$, decoder state $\mathbf{s}_2 \in \mathbb{R}^4$, and attention dimension $d_a = 3$. We compute attention for the second decoder position.

Step 1: Compute alignment scores for each encoder position. Suppose after applying $\mathbf{W}_1 \mathbf{s}_2 + \mathbf{W}_2 \mathbf{h}_i$ and passing through \tanh and \mathbf{v}^\top , we obtain:

$$e_{2,1} = 0.8 \quad (7.17)$$

$$e_{2,2} = 2.1 \quad (7.18)$$

$$e_{2,3} = 0.5 \quad (7.19)$$

These raw scores indicate that encoder position 2 has the highest compatibility with the current decoder state, but the scores are not yet normalized.

Step 2: Apply softmax to convert scores to a probability distribution:

$$\sum_j \exp(e_{2,j}) = \exp(0.8) + \exp(2.1) + \exp(0.5) \quad (7.20)$$

$$\approx 2.23 + 8.17 + 1.65 = 12.05 \quad (7.21)$$

Computing each attention weight:

$$\alpha_{2,1} = \frac{\exp(0.8)}{12.05} = \frac{2.23}{12.05} \approx 0.185 \quad (7.22)$$

$$\alpha_{2,2} = \frac{\exp(2.1)}{12.05} = \frac{8.17}{12.05} \approx 0.678 \quad (7.23)$$

$$\alpha_{2,3} = \frac{\exp(0.5)}{12.05} = \frac{1.65}{12.05} \approx 0.137 \quad (7.24)$$

The decoder places 67.8% of its attention on encoder position 2, with the remaining attention distributed between positions 1 and 3. This sharp distribution indicates high confidence about which input position is relevant.

Step 3: Compute the context vector as a weighted sum:

$$\mathbf{c}_2 = 0.185\mathbf{h}_1 + 0.678\mathbf{h}_2 + 0.137\mathbf{h}_3 \in \mathbb{R}^4 \quad (7.25)$$

If $\mathbf{h}_1 = [1.0, 0.5, -0.3, 0.8]^\top$, $\mathbf{h}_2 = [0.3, 0.9, 0.6, -0.2]^\top$, $\mathbf{h}_3 = [-0.4, 0.2, 0.7, 0.5]^\top$:

$$\mathbf{c}_2 = 0.185 \begin{bmatrix} 1.0 \\ 0.5 \\ -0.3 \\ 0.8 \end{bmatrix} + 0.678 \begin{bmatrix} 0.3 \\ 0.9 \\ 0.6 \\ -0.2 \end{bmatrix} + 0.137 \begin{bmatrix} -0.4 \\ 0.2 \\ 0.7 \\ 0.5 \end{bmatrix} \quad (7.26)$$

$$= \begin{bmatrix} 0.185 + 0.203 - 0.055 \\ 0.093 + 0.610 + 0.027 \\ -0.056 + 0.407 + 0.096 \\ 0.148 - 0.136 + 0.069 \end{bmatrix} = \begin{bmatrix} 0.333 \\ 0.730 \\ 0.447 \\ 0.081 \end{bmatrix} \quad (7.27)$$

The context vector is dominated by \mathbf{h}_2 due to the high attention weight $\alpha_{2,2} = 0.678$, but includes contributions from the other encoder states proportional to their attention weights.

7.3 Scaled Dot-Product Attention

Scaled dot-product attention, introduced in the "Attention is All You Need" paper, represents a fundamental simplification and improvement over additive attention. By replacing the learned additive scoring function with a simple scaled dot product, this mechanism achieves superior computational efficiency while maintaining or improving model performance. This design choice enabled the transformer architecture to scale to billions of parameters and become the foundation of modern large language models.

Definition 7.2 (Scaled Dot-Product Attention). Given queries $\mathbf{Q} \in \mathbb{R}^{m \times d_k}$, keys $\mathbf{K} \in \mathbb{R}^{n \times d_k}$, and values $\mathbf{V} \in \mathbb{R}^{n \times d_v}$, scaled dot-product attention computes:

Step 1: Compute attention scores

$$\mathbf{E} = \mathbf{Q}\mathbf{K}^\top \in \mathbb{R}^{m \times n} \quad (7.28)$$

where entry $e_{i,j} = \mathbf{q}_i^\top \mathbf{k}_j$ measures the compatibility of query i with key j .

Step 2: Scale by $\sqrt{d_k}$

$$\mathbf{E}_{\text{scaled}} = \frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}} \quad (7.29)$$

Step 3: Softmax over keys (row-wise)

$$\mathbf{A} = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}\right) \in \mathbb{R}^{m \times n} \quad (7.30)$$

Step 4: Apply attention to values

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \mathbf{A}\mathbf{V} \in \mathbb{R}^{m \times d_v} \quad (7.31)$$

The complete formula in one line:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}\right) \mathbf{V} \quad (7.32)$$

7.3.1 Why Scaling Matters: Variance Analysis

The scaling factor $1/\sqrt{d_k}$ is not merely a normalization convenience—it is essential for maintaining stable gradients during training. To understand why, we analyze the variance of dot products between queries and keys.

Assume query and key vectors have independent elements with zero mean and unit variance: $\mathbb{E}[\mathbf{q}_i] = \mathbb{E}[\mathbf{k}_i] = 0$ and $\text{Var}(\mathbf{q}_i) = \text{Var}(\mathbf{k}_i) = 1$. The dot product between a query and key is:

$$\mathbf{q}^\top \mathbf{k} = \sum_{i=1}^{d_k} q_i k_i \quad (7.33)$$

Since the elements are independent, the variance of the sum equals the sum of variances:

$$\text{Var}(\mathbf{q}^\top \mathbf{k}) = \sum_{i=1}^{d_k} \text{Var}(q_i k_i) = \sum_{i=1}^{d_k} \mathbb{E}[q_i^2] \mathbb{E}[k_i^2] = \sum_{i=1}^{d_k} 1 \cdot 1 = d_k \quad (7.34)$$

Without scaling, the dot product has variance d_k , which grows linearly with the key dimension. For $d_k = 64$, typical dot products have standard deviation $\sqrt{64} = 8$. For $d_k = 512$, the standard deviation grows to $\sqrt{512} \approx 22.6$. These large magnitudes cause severe problems for the softmax function.

Softmax Saturation Problem: The softmax function is defined as:

$$\text{softmax}(\mathbf{z})_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)} \quad (7.35)$$

When input magnitudes are large, softmax saturates—one element dominates and receives nearly all the probability mass, while others receive exponentially small probabilities. Consider a simple example with two elements:

$$\text{softmax}([z, 0]) = \left[\frac{\exp(z)}{\exp(z) + 1}, \frac{1}{\exp(z) + 1} \right] \quad (7.36)$$

For $z = 10$: $\text{softmax}([10, 0]) \approx [0.9999, 0.0001]$. For $z = 20$: $\text{softmax}([20, 0]) \approx [1.0, 2 \times 10^{-9}]$. The distribution becomes a hard selection rather than a soft weighting.

Gradient Flow Analysis: The gradient of softmax with respect to its input is:

$$\frac{\partial \text{softmax}(\mathbf{z})_i}{\partial z_j} = \text{softmax}(\mathbf{z})_i (\delta_{ij} - \text{softmax}(\mathbf{z})_j) \quad (7.37)$$

When softmax saturates with one element near 1 and others near 0, these gradients become tiny. For the dominant element i where $\text{softmax}(\mathbf{z})_i \approx 1$:

$$\frac{\partial \text{softmax}(\mathbf{z})_i}{\partial z_i} \approx 1 \cdot (1 - 1) = 0 \quad (7.38)$$

For non-dominant elements where $\text{softmax}(\mathbf{z})_j \approx 0$:

$$\frac{\partial \text{softmax}(\mathbf{z})_i}{\partial z_j} \approx 1 \cdot (0 - 0) = 0 \quad (7.39)$$

All gradients vanish, preventing the model from learning. This is analogous to the vanishing gradient problem in deep networks, but occurring within a single attention layer.

Scaling Solution: Dividing by $\sqrt{d_k}$ normalizes the variance:

$$\text{Var} \left(\frac{\mathbf{q}^\top \mathbf{k}}{\sqrt{d_k}} \right) = \frac{1}{d_k} \text{Var}(\mathbf{q}^\top \mathbf{k}) = \frac{1}{d_k} \cdot d_k = 1 \quad (7.40)$$

With unit variance, dot products typically range from -3 to $+3$ (within three standard deviations), keeping softmax in its sensitive region where gradients are substantial. This maintains effective gradient flow throughout training.

Numerical Example: Consider $d_k = 64$ versus $d_k = 512$ with random unit-variance queries and keys. Without scaling, for $d_k = 64$, a typical attention score might be $\mathbf{q}^\top \mathbf{k} = 12.3$. After softmax over 10 keys with similar magnitudes, the distribution might be $[0.45, 0.18, 0.12, 0.08, 0.06, 0.04, 0.03, 0.02, 0.01, 0.01]$ —reasonably distributed. The gradient of the top element is approximately $0.45 \times (1 - 0.45) = 0.248$, which is healthy.

For $d_k = 512$ without scaling, the same query-key pair might produce $\mathbf{q}^\top \mathbf{k} = 35.2$. After softmax, the distribution becomes $[0.9997, 0.0001, 0.0001, 0.0001, \dots]$ —completely saturated. The gradient is approximately $0.9997 \times (1 - 0.9997) = 0.0003$, which is 800 times smaller. Over many layers, these tiny gradients compound, making training extremely difficult or impossible.

With scaling by $\sqrt{512} \approx 22.6$, the score becomes $35.2/22.6 \approx 1.56$, producing a softmax distribution like $[0.38, 0.15, 0.12, 0.10, \dots]$ with gradient $0.38 \times (1 - 0.38) = 0.236$ —similar to the $d_k = 64$ case. The scaling makes attention behavior independent of the key dimension, enabling stable training across different model sizes.

7.3.2 Computational Efficiency

Scaled dot-product attention achieves dramatically better computational efficiency than additive attention, both in raw FLOP count and in hardware utilization. This efficiency difference is why transformers can scale to billions of parameters while additive attention models remained limited to hundreds of millions.

FLOP Count Comparison: For m queries and n keys with dimension d_k :

Scaled dot-product attention:

$$\mathbf{QK}^\top : 2mnd_k \text{ FLOPs} \quad (7.41)$$

$$\text{Scaling} : mn \text{ FLOPs (division)} \quad (7.42)$$

$$\text{Softmax} : \approx 5mn \text{ FLOPs (exp, sum, divide)} \quad (7.43)$$

$$\mathbf{AV} : 2mnd_v \text{ FLOPs} \quad (7.44)$$

$$\text{Total} : 2mn(d_k + d_v) + 6mn \approx 2mn(d_k + d_v) \quad (7.45)$$

For $d_k = d_v = 64$, $m = n = 512$:

$$2 \times 512 \times 512 \times (64 + 64) = 67,108,864 \text{ FLOPs} \approx 67 \text{ MFLOPs} \quad (7.46)$$

Bahdanau attention: As computed earlier, for $d_a = 256$, $d_s = d_h = 512$, $m = n = 512$:

$$2 \times 512 \times 512 \times 256 \times (512 + 512 + 3) \approx 69 \text{ billion FLOPs} \quad (7.47)$$

Scaled dot-product attention requires approximately $1000\times$ fewer FLOPs than Bahdanau attention for this configuration! The difference grows with sequence length since Bahdanau's cost scales with the attention dimension d_a while scaled dot-product depends only on d_k .

Hardware Efficiency: Beyond raw FLOP count, scaled dot-product attention maps naturally to highly optimized GPU operations. The core computation \mathbf{QK}^\top is a dense matrix multiplication (GEMM), which is the most optimized operation on modern GPUs. NVIDIA's cuBLAS library and Tensor Cores are specifically designed for GEMM, achieving 80-90% of theoretical peak performance.

In contrast, Bahdanau attention requires element-wise operations (tanh), vector additions, and many small matrix-vector products. These operations achieve only 15-25% of peak GPU performance due to memory bandwidth limitations and poor parallelization. The tanh activation requires computing exponentials for each element, which is slow compared to the fused multiply-add operations in GEMM.

Memory Bandwidth Considerations: Modern GPUs are often memory-bandwidth limited rather than compute-limited. The NVIDIA A100 has 312 TFLOPS of FP16 compute but only 1.5 TB/s memory bandwidth. For operations to be compute-bound, they must perform many FLOPs per byte loaded from memory.

Matrix multiplication \mathbf{QK}^\top for $\mathbf{Q}, \mathbf{K} \in \mathbb{R}^{512 \times 64}$ loads $2 \times 512 \times 64 \times 2 = 131 \text{ KB}$ (FP16) and performs $2 \times 512 \times 512 \times 64 = 67 \text{ MFLOPs}$, achieving $67,000,000/131,072 \approx 511 \text{ FLOPs per byte}$. This high arithmetic intensity keeps the GPU compute units busy.

Bahdanau's element-wise operations load data, perform a few operations, and store results—achieving only 1-5 FLOPs per byte. The GPU spends most of its time waiting for memory rather than computing, wasting the available compute capacity.

Practical Performance: On an NVIDIA A100 GPU, computing attention for a batch of 32 sequences with $n = 512$ and $d_k = 64$:

- Scaled dot-product attention: $\approx 0.8 \text{ ms}$ (achieving 250 TFLOPS, 80% of peak)
- Bahdanau attention: $\approx 15 \text{ ms}$ (achieving 15 TFLOPS, 5% of peak)

The $19\times$ speedup from scaled dot-product attention is what enables training GPT-3 scale models (175B parameters) in reasonable time. With Bahdanau attention, training would take $19\times$ longer, making such models economically infeasible.

Example 7.3 (Scaled Dot-Product Computation). Consider a single query attending to 3 keys with $d_k = 4$ and $d_v = 5$:

$$\mathbf{q} = \begin{bmatrix} 1.0 \\ 0.5 \\ -0.3 \\ 0.8 \end{bmatrix}, \quad \mathbf{K} = \begin{bmatrix} 0.8 & 0.2 & -0.1 & 0.5 \\ 0.3 & 0.7 & 0.4 & -0.2 \\ -0.5 & 0.1 & 0.9 & 0.6 \end{bmatrix} \quad (7.48)$$

Step 1: Compute dot products between the query and each key:

$$\mathbf{q}^\top \mathbf{k}_1 = 1.0(0.8) + 0.5(0.2) + (-0.3)(-0.1) + 0.8(0.5) \quad (7.49)$$

$$= 0.8 + 0.1 + 0.03 + 0.4 = 1.33 \quad (7.50)$$

$$\mathbf{q}^\top \mathbf{k}_2 = 1.0(0.3) + 0.5(0.7) + (-0.3)(0.4) + 0.8(-0.2) \quad (7.51)$$

$$= 0.3 + 0.35 - 0.12 - 0.16 = 0.37 \quad (7.52)$$

$$\mathbf{q}^\top \mathbf{k}_3 = 1.0(-0.5) + 0.5(0.1) + (-0.3)(0.9) + 0.8(0.6) \quad (7.53)$$

$$= -0.5 + 0.05 - 0.27 + 0.48 = -0.24 \quad (7.54)$$

Step 2: Scale by $\sqrt{d_k} = \sqrt{4} = 2$:

$$\text{scaled scores} = \left[\frac{1.33}{2}, \frac{0.37}{2}, \frac{-0.24}{2} \right] = [0.665, 0.185, -0.120] \quad (7.55)$$

Without scaling, the scores would be $[1.33, 0.37, -0.24]$. For this small example with $d_k = 4$, the difference is modest. But for $d_k = 64$, unscaled scores would be $\sqrt{64/4} = 4$ times larger, and for $d_k = 512$, they would be $\sqrt{512/4} \approx 11.3$ times larger, causing severe softmax saturation.

Step 3: Apply softmax to obtain attention weights:

$$\sum_j \exp(\text{score}_j) = \exp(0.665) + \exp(0.185) + \exp(-0.120) \quad (7.56)$$

$$\approx 1.945 + 1.203 + 0.887 = 4.035 \quad (7.57)$$

Computing each weight:

$$\alpha_1 = \frac{1.945}{4.035} \approx 0.482 \quad (7.58)$$

$$\alpha_2 = \frac{1.203}{4.035} \approx 0.298 \quad (7.59)$$

$$\alpha_3 = \frac{0.887}{4.035} \approx 0.220 \quad (7.60)$$

The attention is distributed across all three keys, with the highest weight on key 1 (48.2%) but substantial attention to keys 2 and 3 as well. This soft distribution allows the model to incorporate information from multiple positions.

Step 4: Apply attention weights to values. Suppose:

$$\mathbf{V} = \begin{bmatrix} 0.5 & 0.8 & -0.2 & 0.6 & 0.3 \\ 0.2 & -0.4 & 0.7 & 0.1 & 0.9 \\ -0.3 & 0.5 & 0.4 & -0.6 & 0.2 \end{bmatrix} \in \mathbb{R}^{3 \times 5} \quad (7.61)$$

The output is:

$$\text{output} = 0.482\mathbf{v}_1 + 0.298\mathbf{v}_2 + 0.220\mathbf{v}_3 \quad (7.62)$$

$$= 0.482 \begin{bmatrix} 0.5 \\ 0.8 \\ -0.2 \\ 0.6 \\ 0.3 \end{bmatrix} + 0.298 \begin{bmatrix} 0.2 \\ -0.4 \\ 0.7 \\ 0.1 \\ 0.9 \end{bmatrix} + 0.220 \begin{bmatrix} -0.3 \\ 0.5 \\ 0.4 \\ -0.6 \\ 0.2 \end{bmatrix} \quad (7.63)$$

$$= \begin{bmatrix} 0.241 + 0.060 - 0.066 \\ 0.386 - 0.119 + 0.110 \\ -0.096 + 0.209 + 0.088 \\ 0.289 + 0.030 - 0.132 \\ 0.145 + 0.268 + 0.044 \end{bmatrix} = \begin{bmatrix} 0.235 \\ 0.377 \\ 0.201 \\ 0.187 \\ 0.457 \end{bmatrix} \in \mathbb{R}^5 \quad (7.64)$$

The output vector is a weighted combination of the value vectors, with weights determined by the query-key similarities. This output can then be used by subsequent layers in the transformer.

7.4 Attention Score Computation Methods

The evolution of attention mechanisms reflects a progression toward greater computational efficiency and hardware compatibility. Understanding the trade-offs between different scoring functions illuminates why scaled dot-product attention became the standard for transformers, despite the apparent simplicity of alternatives.

7.4.1 Taxonomy of Attention Mechanisms

Attention mechanisms differ primarily in how they compute the compatibility score between a query and a key. This scoring function determines both the computational cost and the expressiveness of the attention mechanism. We examine four major variants that represent key points in the design space.

Additive Attention (Bahdanau): Introduced in 2015 for neural machine translation, additive attention computes scores through a learned feedforward network:

$$\text{score}(\mathbf{q}, \mathbf{k}) = \mathbf{v}^\top \tanh(\mathbf{W}_1 \mathbf{q} + \mathbf{W}_2 \mathbf{k}) \quad (7.65)$$

This approach projects the query and key into a shared space of dimension d_a , applies a nonlinearity, and projects to a scalar. The nonlinearity allows the model to learn complex, non-linear compatibility functions. However, this flexibility comes at substantial computational cost: each score requires $O(d_a(d_q + d_k))$ operations and introduces $O(d_a(d_q + d_k))$ learnable parameters.

For typical dimensions $d_q = d_k = 512$ and $d_a = 256$, computing one score requires approximately $2 \times 256 \times (512 + 512) = 524,288$ FLOPs. For a sequence of length 512 attending to itself, this amounts to $512^2 \times 524,288 \approx 137$ billion FLOPs just for score computation. The tanh nonlinearity and multiple matrix-vector products prevent efficient GPU utilization, achieving only 15-25% of peak performance.

Dot-Product Attention: The simplest scoring function is the unscaled dot product:

$$\text{score}(\mathbf{q}, \mathbf{k}) = \mathbf{q}^\top \mathbf{k} \quad (7.66)$$

This requires only $2d_k$ FLOPs (one multiply-add per dimension) and introduces zero parameters. The computation maps perfectly to GPU hardware through matrix multiplication: computing all scores \mathbf{QK}^\top is a single GEMM operation achieving 80-90% of peak performance.

However, unscaled dot-product attention suffers from the variance problem discussed earlier. For $d_k = 512$, dot products have standard deviation $\sqrt{512} \approx 22.6$, causing softmax saturation and vanishing gradients. This instability makes unscaled dot-product attention impractical for training deep networks, despite its computational advantages.

Scaled Dot-Product Attention: Adding the scaling factor $1/\sqrt{d_k}$ solves the variance problem while preserving computational efficiency:

$$\text{score}(\mathbf{q}, \mathbf{k}) = \frac{\mathbf{q}^\top \mathbf{k}}{\sqrt{d_k}} \quad (7.67)$$

The scaling adds negligible computational cost (one division per score, or n^2 operations for an $n \times n$ attention matrix) but stabilizes training by maintaining unit variance. This mechanism achieves the best of both worlds: the computational efficiency of dot-product attention with the training stability of additive attention. For $n = 512$ and $d_k = 64$, computing all scores requires $2 \times 512^2 \times 64 = 33.5$ million FLOPs—250× fewer than additive attention.

General Attention (Luong): Luong attention introduces a learned transformation matrix:

$$\text{score}(\mathbf{q}, \mathbf{k}) = \mathbf{q}^\top \mathbf{Wk} \quad (7.68)$$

where $\mathbf{W} \in \mathbb{R}^{d_q \times d_k}$ is a learned parameter matrix. This allows the model to learn a task-specific similarity metric rather than using raw dot products. The transformation can project queries and keys into a shared space where their dot product is more meaningful.

Computing one score requires $2d_qd_k$ FLOPs (matrix-vector product $\mathbf{W}\mathbf{k}$, then dot product with \mathbf{q}), and the mechanism introduces d_qd_k parameters. For $d_q = d_k = 512$, this is 262,144 parameters—substantial but manageable. The computational cost is $2 \times 512 \times 512 = 524,288$ FLOPs per score, similar to additive attention but without the nonlinearity.

General attention achieves better GPU utilization than additive attention because the core operation $\mathbf{Q}\mathbf{W}\mathbf{K}^\top$ can be computed as two matrix multiplications. However, it still requires more computation than scaled dot-product attention and introduces parameters that must be learned.

7.4.2 Comparative Analysis

The following table summarizes the key characteristics of each attention mechanism:

Method	Computation	Parameters	GPU Util.	Used In
Additive (Bahdanau)	$\mathbf{v}^\top \tanh(\mathbf{W}_1\mathbf{q} + \mathbf{W}_2\mathbf{k})$	$O(d_a(d_q + d_k))$	15-25%	Early seq2seq
Dot-product	$\mathbf{q}^\top \mathbf{k}$	0	80-90%	Not used
Scaled dot-product	$\mathbf{q}^\top \mathbf{k} / \sqrt{d_k}$	0	80-90%	Transformers
General (Luong)	$\mathbf{q}^\top \mathbf{W}\mathbf{k}$	$O(d_qd_k)$	50-70%	Some seq2seq

Why Scaled Dot-Product Won: The dominance of scaled dot-product attention in modern transformers reflects several factors beyond raw computational efficiency. First, the parameter-free nature means the model’s capacity is concentrated in the QKV projection matrices and feed-forward networks rather than the scoring function. This architectural choice scales better to very large models—GPT-3 with 175 billion parameters would require even more parameters if using additive or general attention.

Second, the simplicity of scaled dot-product attention makes it easier to optimize and implement efficiently. Hardware vendors can optimize specifically for the $\mathbf{Q}\mathbf{K}^\top$ operation, and software frameworks can apply specialized kernels. FlashAttention and other efficient attention implementations focus on scaled dot-product attention because its regular structure enables aggressive optimization.

Third, empirical results show that scaled dot-product attention performs as well as or better than more complex alternatives on most tasks. The learned transformations in additive or general attention provide little benefit in practice, suggesting that the model can learn appropriate representations through the QKV projections rather than the scoring function itself.

Memory Bandwidth Analysis: Modern GPUs are often memory-bandwidth limited, meaning performance depends on how efficiently data is moved between memory and compute units. Scaled dot-product attention achieves high arithmetic intensity—many FLOPs per byte loaded—because the matrix multiplication $\mathbf{Q}\mathbf{K}^\top$ reuses each element of \mathbf{Q} and \mathbf{K} multiple times.

For $\mathbf{Q}, \mathbf{K} \in \mathbb{R}^{512 \times 64}$ in FP16, we load $2 \times 512 \times 64 \times 2 = 131$ KB and perform $2 \times 512^2 \times 64 = 33.5$ million FLOPs, achieving $33,500,000 / 131,072 \approx 256$ FLOPs per byte. This high intensity keeps compute units busy rather than waiting for memory.

Additive attention achieves only 5-20 FLOPs per byte because element-wise operations and small matrix-vector products provide little data reuse. The GPU spends most of its time waiting for memory transfers, wasting compute capacity. This fundamental difference in memory access patterns explains much of the performance gap between attention mechanisms.

7.5 Query-Key-Value Paradigm

7.5.1 Intuition

The query-key-value (QKV) framework provides an elegant abstraction for understanding attention mechanisms through the lens of information retrieval. This paradigm, borrowed from database systems and search engines, offers intuitive explanations for attention’s behavior while precisely defining its mathematical operations.

Consider a database system where you want to retrieve relevant information. You provide a query describing what you’re looking for, the system compares your query against keys (indexed descriptions of stored content), and returns the values (actual content) associated with the most relevant keys.

Attention mechanisms operate identically: queries represent "what I'm looking for," keys represent "what information is available," and values represent "the actual information to retrieve."

In the context of neural networks, these three components serve distinct roles. The query \mathbf{q} encodes the current position's information needs—what aspects of the input are relevant for processing this position. The keys \mathbf{k}_i encode what information each input position offers—what content is available at that position. The values \mathbf{v}_i encode the actual information to be retrieved—the representations that will be combined to form the output.

This separation of concerns is crucial. By decoupling "what to look for" (queries) from "what is available" (keys) and "what to retrieve" (values), the attention mechanism gains flexibility. The same input can be queried in different ways by different positions, and the retrieved information can differ from the indexing representation. This three-way separation enables the model to learn rich, task-specific attention patterns.

Concrete Example: In machine translation, when generating the French word "chat" (cat) from the English sentence "The cat sat on the mat," the decoder's query encodes "I need information about the subject noun." The keys encode what each English word represents: "the" offers determiner information, "cat" offers subject noun information, "sat" offers verb information, etc. The attention mechanism computes high similarity between the query and the key for "cat," then retrieves the value associated with "cat"—a rich representation encoding its meaning, grammatical role, and context.

Importantly, the key and value for "cat" can differ. The key might emphasize grammatical features (noun, singular, animate) that help match queries, while the value emphasizes semantic features (animal, feline, pet) that are useful for generation. This separation allows the attention mechanism to index on one set of features while retrieving another.

7.5.2 Projecting to QKV

In transformers, queries, keys, and values are not provided directly but are computed from the input through learned linear projections. This design choice allows the model to learn task-specific representations for each role rather than using the raw input embeddings.

Given input $\mathbf{X} \in \mathbb{R}^{n \times d_{\text{model}}}$ where n is the sequence length and d_{model} is the model dimension, we compute:

$$\mathbf{Q} = \mathbf{XW}^Q \qquad \mathbf{W}^Q \in \mathbb{R}^{d_{\text{model}} \times d_k} \qquad (7.69)$$

$$\mathbf{K} = \mathbf{XW}^K \qquad \mathbf{W}^K \in \mathbb{R}^{d_{\text{model}} \times d_k} \qquad (7.70)$$

$$\mathbf{V} = \mathbf{XW}^V \qquad \mathbf{W}^V \in \mathbb{R}^{d_{\text{model}} \times d_v} \qquad (7.71)$$

Each projection matrix is a learned parameter that transforms the input into the appropriate representation space. The query and key projections map to the same dimension d_k (typically d_{model}/h where h is the number of attention heads) because they must be compatible for dot products. The value projection maps to dimension d_v , which is often equal to d_k but can differ.

Why Learn Separate Projections? One might ask: why not use the input \mathbf{X} directly as queries, keys, and values? The answer lies in representation learning. The raw input embeddings encode general semantic and syntactic information, but attention requires specialized representations. The query projection learns to emphasize features relevant for determining what to attend to. The key projection learns to emphasize features relevant for being attended to. The value projection learns to emphasize features relevant for the output representation.

These three projections can learn different aspects of the input. For example, in a language model, the query projection might emphasize the current word's part of speech and semantic category to determine what context is needed. The key projection might emphasize each word's grammatical role and position to help queries find relevant context. The value projection might emphasize semantic content and relationships to provide useful information for prediction.

Computational Cost: Each projection is a matrix multiplication requiring $2nd_{\text{model}}d_k$ FLOPs (for queries and keys) or $2nd_{\text{model}}d_v$ FLOPs (for values). With three projections and $d_k = d_v$:

$$\text{QKV projection FLOPs} = 3 \times 2nd_{\text{model}}d_k = 6nd_{\text{model}}d_k \qquad (7.72)$$

For typical transformer configurations where $d_k = d_{\text{model}}/h$ and we consider all h heads together (so $hd_k = d_{\text{model}}$):

$$\text{QKV projection FLOPs} = 6nd_{\text{model}}^2 \quad (7.73)$$

For BERT-base with $n = 512$ and $d_{\text{model}} = 768$:

$$6 \times 512 \times 768^2 = 1,811,939,328 \text{ FLOPs} \approx 1.8 \text{ GFLOPs} \quad (7.74)$$

This is substantial but represents only about 20% of the total attention computation for typical sequence lengths. The attention score computation (\mathbf{QK}^\top) and output computation (\mathbf{AV}) dominate for longer sequences.

Parameter Count: The three projection matrices introduce $d_{\text{model}}(2d_k + d_v)$ parameters per attention head. For h heads with $d_k = d_v = d_{\text{model}}/h$:

$$\text{QKV parameters} = h \times d_{\text{model}} \times 3 \times \frac{d_{\text{model}}}{h} = 3d_{\text{model}}^2 \quad (7.75)$$

For BERT-base with $d_{\text{model}} = 768$: $3 \times 768^2 = 1,769,472$ parameters per attention layer. With 12 layers, the QKV projections account for $12 \times 1.77 = 21.2$ million parameters out of BERT's total 110 million—about 19% of the model.

Example 7.4 (QKV Projection). Consider a sequence of 5 tokens, each represented by a $d_{\text{model}} = 512$ dimensional vector:

$$\mathbf{X} \in \mathbb{R}^{5 \times 512} \quad (7.76)$$

We project to $d_k = d_v = 64$ (as in a single attention head of a model with $h = 8$ heads):

$$\mathbf{Q} = \mathbf{XW}^Q \in \mathbb{R}^{5 \times 64} \quad (\mathbf{W}^Q \in \mathbb{R}^{512 \times 64}) \quad (7.77)$$

$$\mathbf{K} = \mathbf{XW}^K \in \mathbb{R}^{5 \times 64} \quad (\mathbf{W}^K \in \mathbb{R}^{512 \times 64}) \quad (7.78)$$

$$\mathbf{V} = \mathbf{XW}^V \in \mathbb{R}^{5 \times 64} \quad (\mathbf{W}^V \in \mathbb{R}^{512 \times 64}) \quad (7.79)$$

Each projection matrix has $512 \times 64 = 32,768$ parameters. Computing each projection requires $2 \times 5 \times 512 \times 64 = 327,680$ FLOPs, for a total of $3 \times 327,680 = 983,040$ FLOPs across all three projections.

Attention computation: After projection, we compute attention:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{QK}^\top}{\sqrt{64}}\right) \mathbf{V} \quad (7.80)$$

The attention matrix $\mathbf{A} = \text{softmax}(\mathbf{QK}^\top/\sqrt{64}) \in \mathbb{R}^{5 \times 5}$ has entry a_{ij} representing how much position i attends to position j . For example:

$$\mathbf{A} = \begin{bmatrix} 0.45 & 0.25 & 0.15 & 0.10 & 0.05 \\ 0.10 & 0.50 & 0.25 & 0.10 & 0.05 \\ 0.05 & 0.15 & 0.40 & 0.30 & 0.10 \\ 0.05 & 0.10 & 0.20 & 0.50 & 0.15 \\ 0.05 & 0.05 & 0.10 & 0.25 & 0.55 \end{bmatrix} \quad (7.81)$$

Position 1 attends most strongly to itself (45%) and position 2 (25%). Position 5 attends most strongly to itself (55%) and position 4 (25%). This pattern might emerge in a language model where each position attends to nearby context, with stronger attention to the current position and recent tokens.

The output $\mathbf{AV} \in \mathbb{R}^{5 \times 64}$ provides an attended representation for each position, combining information from all positions according to the attention weights. This output can then be processed by subsequent layers.

7.6 Hardware Implications of Attention

The shift from RNN-based sequence models to attention-based transformers represents not just an algorithmic change but a fundamental realignment with modern hardware capabilities. Understanding why attention enables effective GPU utilization illuminates both the success of transformers and the design principles for future architectures.

7.6.1 Parallelization: RNNs vs Attention

Recurrent neural networks process sequences sequentially by design. Each hidden state \mathbf{h}_t depends on the previous hidden state \mathbf{h}_{t-1} :

$$\mathbf{h}_t = \text{RNN}(\mathbf{x}_t, \mathbf{h}_{t-1}) \quad (7.82)$$

This recurrence creates a data dependency chain: we cannot compute \mathbf{h}_t until \mathbf{h}_{t-1} is available, which requires \mathbf{h}_{t-2} , and so on back to \mathbf{h}_0 . For a sequence of length n , we must perform n sequential operations, even if we have thousands of parallel compute units available.

Modern GPUs contain thousands of CUDA cores (NVIDIA A100 has 6912 cores) capable of executing operations simultaneously. When processing a batch of 32 sequences of length 512, we have $32 \times 512 = 16,384$ positions to process. An RNN must process these sequentially in 512 time steps, using only 32 parallel threads (one per sequence in the batch). This achieves only $32/6912 \approx 0.5\%$ of the GPU's parallel capacity.

Attention mechanisms eliminate this sequential bottleneck. The attention computation for position i depends only on the input sequence, not on previous attention computations:

$$\text{output}_i = \sum_{j=1}^n \alpha_{ij} \mathbf{v}_j \quad (7.83)$$

All attention weights α_{ij} can be computed simultaneously because they depend only on the queries and keys, which are computed from the input via matrix multiplication. During training, when the entire target sequence is known, we can compute attention for all positions in parallel. For 32 sequences of length 512, we can utilize all 16,384 positions simultaneously, achieving near-full GPU utilization.

Practical Impact: On an NVIDIA A100 GPU, processing a batch of 32 sequences of length 512 with hidden dimension 768:

- LSTM encoder: ≈ 45 ms (sequential processing, 5% GPU utilization)
- Transformer encoder: ≈ 3 ms (parallel processing, 75% GPU utilization)

The $15\times$ speedup from parallelization is what makes training large transformers feasible. GPT-3 with 175 billion parameters was trained on 300 billion tokens—a task that would be impractical with sequential RNN processing.

7.6.2 Memory Bandwidth vs Compute

Modern GPUs have enormous compute capacity but limited memory bandwidth. The NVIDIA A100 provides 312 TFLOPS (FP16 with Tensor Cores) but only 1.5 TB/s memory bandwidth. This means the GPU can perform $312 \times 10^{12} / (1.5 \times 10^{12}) \approx 208$ FLOPs for every byte loaded from memory. Operations must achieve high arithmetic intensity (FLOPs per byte) to be compute-bound rather than memory-bound.

Attention's core operation \mathbf{QK}^\top is a matrix multiplication, which achieves high arithmetic intensity through data reuse. For $\mathbf{Q}, \mathbf{K} \in \mathbb{R}^{512 \times 64}$:

- Data loaded: $2 \times 512 \times 64 \times 2 = 131$ KB (FP16)
- FLOPs performed: $2 \times 512^2 \times 64 = 33.5$ million
- Arithmetic intensity: $33,500,000 / 131,072 \approx 256$ FLOPs/byte

This high intensity means the GPU’s compute units stay busy rather than waiting for memory. The operation is compute-bound, achieving 80-90% of peak FLOPS.

In contrast, RNN operations like element-wise activations and small matrix-vector products achieve only 1-10 FLOPs per byte. The GPU spends most of its time waiting for memory transfers, achieving only 5-15% of peak FLOPS. This fundamental difference in memory access patterns explains much of the performance gap between RNNs and transformers.

Memory Requirements: Attention does require more memory than RNNs due to storing the attention matrix. For a sequence of length n with batch size b and h heads:

$$\text{Attention matrix memory} = b \times h \times n^2 \times 4 \text{ bytes (FP32)} \quad (7.84)$$

For BERT-base with $b = 32$, $h = 12$, $n = 512$:

$$32 \times 12 \times 512^2 \times 4 = 402,653,184 \text{ bytes} \approx 403 \text{ MB} \quad (7.85)$$

This is substantial but manageable on modern GPUs with 40-80 GB of memory. The memory cost is the price paid for parallelization—we trade memory for speed.

7.6.3 Batch Processing Efficiency

Attention mechanisms benefit dramatically from batching because matrix multiplications become more efficient with larger matrices. For small matrices, memory transfer overhead dominates. For large matrices, the compute units stay busy and achieve high utilization.

Consider computing \mathbf{QK}^\top for a single sequence ($b = 1$) versus a batch ($b = 32$):

- Single sequence: $\mathbf{Q}, \mathbf{K} \in \mathbb{R}^{512 \times 64}$, output $\in \mathbb{R}^{512 \times 512}$
 - FLOPs: $2 \times 512^2 \times 64 = 33.5$ million
 - Time on A100: ≈ 0.15 ms (achieves 30% peak, memory-bound)
- Batch of 32: $\mathbf{Q}, \mathbf{K} \in \mathbb{R}^{16384 \times 64}$, output $\in \mathbb{R}^{16384 \times 16384}$
 - FLOPs: $2 \times 16384^2 \times 64 = 34.4$ billion
 - Time on A100: ≈ 120 ms (achieves 85% peak, compute-bound)

The batched computation achieves $34.4/33.5 = 1024$ times more FLOPs in $120/0.15 = 800$ times more time, demonstrating the efficiency gain from batching. Per-sequence processing time drops from 0.15 ms to $120/32 = 3.75$ ms, but the throughput increases from $1/0.15 = 6.7$ sequences/ms to $32/120 = 0.27$ sequences/ms... wait, that’s wrong. Let me recalculate: throughput increases from $1/0.15 \approx 6.7$ sequences/second to $32/120 \approx 267$ sequences/second—a $40\times$ improvement in throughput.

This batching efficiency is crucial for both training (where large batches improve gradient estimates and enable higher learning rates) and inference (where serving multiple requests simultaneously improves throughput). Attention’s regular structure makes it particularly amenable to batching, unlike RNNs where variable-length sequences complicate batch processing.

7.7 Attention Variants

7.7.1 Self-Attention vs Cross-Attention

Self-Attention: $\mathbf{Q}, \mathbf{K}, \mathbf{V}$ all from same source

$$\mathbf{Q} = \mathbf{K} = \mathbf{V} = \mathbf{XW} \quad (7.86)$$

Used in: Transformer encoder, BERT

Cross-Attention: Queries from one source, keys and values from another

$$\mathbf{Q} = \mathbf{X}_{\text{dec}} \mathbf{W}^Q, \quad \mathbf{K} = \mathbf{V} = \mathbf{X}_{\text{enc}} \mathbf{W}^{K/V} \quad (7.87)$$

Used in: Transformer decoder (attending to encoder output)

7.7.2 Masked Attention

For autoregressive models (GPT), prevent attending to future positions:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top + \mathbf{M}}{\sqrt{d_k}}\right) \mathbf{V} \quad (7.88)$$

where mask $\mathbf{M}_{ij} = -\infty$ if $j > i$, else $\mathbf{M}_{ij} = 0$.

After softmax, $\exp(-\infty) = 0$, so no attention to future!

7.8 Exercises

Exercise 7.1. Compute Bahdanau attention for sequence length 4, decoder state dim 3, attention dim 2. Given specific \mathbf{W}_1 , \mathbf{W}_2 , \mathbf{v} , encoder states, and decoder state, calculate all attention weights.

Exercise 7.2. For scaled dot-product attention with $\mathbf{Q} \in \mathbb{R}^{10 \times 64}$, $\mathbf{K} \in \mathbb{R}^{20 \times 64}$, $\mathbf{V} \in \mathbb{R}^{20 \times 128}$: (1) What is output dimension? (2) What is attention matrix shape? (3) How many FLOPs for computing $\mathbf{Q}\mathbf{K}^\top$?

Exercise 7.3. Show that without scaling, for $d_k = 64$ and unit variance elements, dot products have variance 64. Demonstrate numerically how this affects softmax gradients.

Exercise 7.4. Implement scaled dot-product attention in PyTorch. Test with sequences of length 5 and 10, dimensions $d_k = 32$, $d_v = 48$. Visualize attention weights as heatmap.

Chapter 8

Self-Attention and Multi-Head Attention

Chapter Overview

Self-attention is the core innovation enabling transformers. This chapter develops self-attention from first principles, then introduces multi-head attention—the mechanism that allows transformers to attend to multiple types of relationships simultaneously.

Learning Objectives

1. Understand self-attention and its advantages over RNNs
2. Implement multi-head attention from scratch
3. Compute output dimensions and parameter counts
4. Understand positional encodings for sequence order
5. Analyze computational complexity of attention
6. Apply masking for causal (autoregressive) attention

8.1 Self-Attention Mechanism

Definition 8.1 (Self-Attention). For input sequence $\mathbf{X} \in \mathbb{R}^{n \times d}$, self-attention computes output where each position attends to all positions:

$$\mathbf{Q} = \mathbf{X}\mathbf{W}^Q, \quad \mathbf{K} = \mathbf{X}\mathbf{W}^K, \quad \mathbf{V} = \mathbf{X}\mathbf{W}^V \quad (8.1)$$

$$\text{SelfAttn}(\mathbf{X}) = \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}} \right) \mathbf{V} \quad (8.2)$$

where $\mathbf{W}^Q, \mathbf{W}^K \in \mathbb{R}^{d \times d_k}$ and $\mathbf{W}^V \in \mathbb{R}^{d \times d_v}$.

Self-attention exhibits several fundamental properties that distinguish it from recurrent architectures. The mechanism is permutation equivariant, meaning that if the input sequence order changes, the output changes correspondingly—there is no inherent notion of sequence order without positional encodings. Every position in the sequence attends to every other position through all-to-all connections, creating direct paths between any pair of tokens regardless of their distance in the sequence. This contrasts sharply with RNNs, where information must propagate sequentially through intermediate hidden states, potentially degrading over long distances.

The parallel computation property is perhaps the most significant advantage for modern hardware. Unlike RNNs which process sequences sequentially due to the recurrence relation $\mathbf{h}_t = f(\mathbf{h}_{t-1}, \mathbf{x}_t)$, self-attention computes all output positions simultaneously. This enables full utilization of GPU parallelism, where thousands of cores can work concurrently on different positions and attention heads. The long-range dependency modeling is direct rather than transitive: position 0 can attend to position 1000 with a single attention operation, whereas an RNN requires 1000 sequential steps, each potentially losing information through the recurrent bottleneck.

Example 8.1 (Self-Attention Computation). Input: 3 word embeddings, each $d = 4$ dimensional

$$\mathbf{X} = \begin{bmatrix} 1.0 & 0.5 & 0.2 & 0.8 \\ 0.3 & 1.2 & 0.7 & 0.4 \\ 0.6 & 0.9 & 1.1 & 0.3 \end{bmatrix} \in \mathbb{R}^{3 \times 4} \quad (8.3)$$

Projection matrices with $d_k = d_v = 3$:

$$\mathbf{W}^Q, \mathbf{W}^K, \mathbf{W}^V \in \mathbb{R}^{4 \times 3} \quad (8.4)$$

Step 1: Project to QKV

$$\mathbf{Q} = \mathbf{XW}^Q \in \mathbb{R}^{3 \times 3} \quad (8.5)$$

$$\mathbf{K} = \mathbf{XW}^K \in \mathbb{R}^{3 \times 3} \quad (8.6)$$

$$\mathbf{V} = \mathbf{XW}^V \in \mathbb{R}^{3 \times 3} \quad (8.7)$$

Step 2: Compute attention scores

$$\mathbf{QK}^\top \in \mathbb{R}^{3 \times 3} \quad (8.8)$$

Entry (i, j) measures how much position i attends to position j .

Step 3: Scale and softmax

$$\mathbf{A} = \text{softmax} \left(\frac{\mathbf{QK}^\top}{\sqrt{3}} \right) \in \mathbb{R}^{3 \times 3} \quad (8.9)$$

Each row sums to 1 (probability distribution over positions to attend to).

Step 4: Apply to values

$$\text{Output} = \mathbf{AV} \in \mathbb{R}^{3 \times 3} \quad (8.10)$$

Each output position is weighted combination of all input value vectors.

8.1.1 Hardware Considerations and Memory Layout

The memory layout of attention matrices in GPU memory significantly impacts performance. When computing self-attention for a batch of sequences, the attention matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ for each head must be materialized in GPU global memory. For BERT-base with 12 attention heads and maximum sequence length 512, each attention matrix contains $512 \times 512 = 262,144$ elements. Storing these in FP32 format requires $262,144 \times 4 = 1,048,576$ bytes, or approximately 1 MB per head per sequence. With 12 heads, this amounts to 12 MB per sequence just for the attention weights themselves, not including the query, key, and value matrices.

The memory requirements scale dramatically with batch size. For a batch of 32 sequences—a typical training batch size—the attention matrices alone consume $12 \times 32 = 384$ MB of GPU memory. This explains why training transformers on long sequences quickly exhausts available GPU memory. For sequence length 2048, the attention matrices grow to $2048^2 \times 4 = 16,777,216$ bytes per head, or approximately 16 MB. With 12 heads and batch size 32, this becomes $16 \times 12 \times 32 = 6,144$ MB, or

roughly 6 GB just for attention weights. An NVIDIA A100 with 40 GB of memory can accommodate this, but longer sequences of 4096 tokens would require $4096^2 \times 4 \times 12 \times 32 / (1024^3) \approx 24$ GB for attention matrices alone, leaving little room for activations, gradients, and model parameters.

The memory access patterns during attention computation determine whether the operation is compute-bound or memory-bound. Computing the attention scores \mathbf{QK}^\top involves reading the query matrix $\mathbf{Q} \in \mathbb{R}^{n \times d_k}$ and key matrix $\mathbf{K} \in \mathbb{R}^{n \times d_k}$ from global memory, performing $O(n^2 d_k)$ floating-point operations, and writing the result $\mathbf{S} \in \mathbb{R}^{n \times n}$ back to memory. For small batch sizes, the computation is fast but the memory transfers dominate. An NVIDIA A100 has memory bandwidth of approximately 1.5 TB/s and peak FP16 compute throughput of 312 TFLOPS. For BERT-base with $n = 512$ and $d_k = 64$, computing \mathbf{QK}^\top requires reading $2 \times 512 \times 64 \times 2 = 131,072$ bytes (in FP16) and performing $512^2 \times 64 \times 2 = 33,554,432$ FLOPs. The arithmetic intensity is $33,554,432 / 131,072 \approx 256$ FLOPs per byte, which is reasonably high. However, the subsequent softmax operation and multiplication by \mathbf{V} have lower arithmetic intensity, making attention memory-bound for small batches.

Cache locality plays a crucial role in attention performance. Modern GPUs have a memory hierarchy with small but fast on-chip SRAM (shared memory) and large but slower off-chip DRAM (global memory). The attention computation as typically implemented requires multiple passes over the data: first computing \mathbf{QK}^\top , then applying softmax, then multiplying by \mathbf{V} . Each pass reads data from global memory, processes it, and writes results back. For long sequences, the attention matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is too large to fit in SRAM, forcing repeated global memory accesses. Flash Attention (Section 8.7) addresses this by tiling the computation to fit in SRAM, dramatically reducing memory traffic and improving performance by 2-4 \times for long sequences.

8.2 Multi-Head Attention

Single-head attention with a single set of query, key, and value projections may capture only one type of relationship between tokens. In natural language, tokens relate to each other in multiple ways simultaneously: syntactically (subject-verb agreement, dependency structure), semantically (synonymy, antonymy, topic coherence), and positionally (proximity, relative ordering). A single attention head must compress all these relationship types into a single attention distribution, potentially losing important information. Multi-head attention addresses this limitation by computing multiple attention functions in parallel, each with its own learned projection matrices, allowing the model to attend to different types of relationships simultaneously.

Definition 8.2 (Multi-Head Attention). With h attention heads, each with dimension $d_k = d_v = d_{\text{model}}/h$:

For each head $i = 1, \dots, h$:

$$\mathbf{Q}^{(i)} = \mathbf{XW}^{Q(i)}, \quad \mathbf{K}^{(i)} = \mathbf{XW}^{K(i)}, \quad \mathbf{V}^{(i)} = \mathbf{XW}^{V(i)} \quad (8.11)$$

$$\text{head}_i = \text{Attention}(\mathbf{Q}^{(i)}, \mathbf{K}^{(i)}, \mathbf{V}^{(i)}) \quad (8.12)$$

Concatenate and project:

$$\text{MultiHead}(\mathbf{X}) = [\text{head}_1; \dots; \text{head}_h] \mathbf{W}^O \quad (8.13)$$

where $\mathbf{W}^{Q(i)}, \mathbf{W}^{K(i)}, \mathbf{W}^{V(i)} \in \mathbb{R}^{d_{\text{model}} \times d_k}$ and $\mathbf{W}^O \in \mathbb{R}^{hd_k \times d_{\text{model}}}$.

Example 8.2 (BERT-base Multi-Head Attention). BERT-base parameters:

- Model dimension: $d_{\text{model}} = 768$
- Number of heads: $h = 12$

- Dimension per head: $d_k = d_v = 768/12 = 64$
- Sequence length: $n = 512$ (maximum)

For single head:

$$\mathbf{Q}^{(i)} = \mathbf{X}\mathbf{W}^{Q(i)} \in \mathbb{R}^{512 \times 64} \quad (\mathbf{W}^{Q(i)} \in \mathbb{R}^{768 \times 64}) \quad (8.14)$$

$$\mathbf{K}^{(i)} = \mathbf{X}\mathbf{W}^{K(i)} \in \mathbb{R}^{512 \times 64} \quad (8.15)$$

$$\mathbf{V}^{(i)} = \mathbf{X}\mathbf{W}^{V(i)} \in \mathbb{R}^{512 \times 64} \quad (8.16)$$

Attention matrix: $\mathbf{A}^{(i)} \in \mathbb{R}^{512 \times 512}$ (huge!)

Concatenate all 12 heads:

$$[\text{head}_1; \dots; \text{head}_{12}] \in \mathbb{R}^{512 \times 768} \quad (8.17)$$

Output projection:

$$\text{Output} = [\text{head}_1; \dots; \text{head}_{12}] \mathbf{W}^O \in \mathbb{R}^{512 \times 768} \quad (8.18)$$

where $\mathbf{W}^O \in \mathbb{R}^{768 \times 768}$.

Parameter count:

$$\text{QKV projections: } 3h \cdot d_{\text{model}} \cdot d_k = 3 \times 12 \times 768 \times 64 = 1,769,472 \quad (8.19)$$

$$\text{Output projection: } d_{\text{model}}^2 = 768^2 = 589,824 \quad (8.20)$$

$$\text{Total: } 2,359,296 \text{ parameters per attention layer} \quad (8.21)$$

8.2.1 Parallel Computation and Memory Layout

Multiple attention heads can be computed in parallel on modern GPUs, with each head assigned to different streaming multiprocessors or computed concurrently through batched matrix operations. The key design decision is the memory layout: should the heads be stored in an interleaved fashion where all heads for a given position are contiguous, or should each head's data be stored separately? The interleaved layout $[\text{head}_1(\text{pos}_1), \text{head}_2(\text{pos}_1), \dots, \text{head}_h(\text{pos}_1), \text{head}_1(\text{pos}_2), \dots]$ provides better cache locality when concatenating heads for the output projection, since all data for a position is contiguous. The separated layout $[\text{head}_1(\text{pos}_1), \text{head}_1(\text{pos}_2), \dots, \text{head}_2(\text{pos}_1), \dots]$ allows each head to be processed independently with better memory coalescing within a head. Most implementations use the separated layout during attention computation and transpose to interleaved layout before the output projection.

The standard choice of $d_k = d_{\text{model}}/h$ ensures that the total number of parameters remains constant regardless of the number of heads. With h heads each of dimension d_k , the total dimension after concatenation is $h \cdot d_k = d_{\text{model}}$, matching the input dimension. This design choice means that using more heads does not increase the parameter count—it simply partitions the representation space into more subspaces. For BERT-base with $d_{\text{model}} = 768$ and $h = 12$, each head has dimension $d_k = 64$. The QKV projection matrices have shape 768×64 per head, for a total of $3 \times 12 \times 768 \times 64 = 1,769,472$ parameters. If instead a single head with $d_k = 768$ were used, the QKV projections would have shape 768×768 each, for a total of $3 \times 768^2 = 1,769,472$ parameters—exactly the same. The difference lies not in parameter count but in representational capacity: multiple heads can learn diverse attention patterns, while a single large head must compress all patterns into one.

Load balancing across heads is generally not a concern during training, as all heads are computed in parallel through batched matrix operations. However, during inference with dynamic batching or when pruning less important heads, load imbalance can occur. Some heads may be more important than others for the task at hand, and recent work has shown that many heads can be pruned without significant performance degradation. For example, in BERT-base, pruning 40% of attention heads (keeping only 7-8 heads per layer) typically reduces accuracy by less than 1% on downstream tasks, while reducing inference time by approximately 20%. This suggests that the 12 heads provide redundancy and

that the model could function with fewer heads, though training with more heads may help optimization by providing multiple gradient pathways.

8.2.2 Tensor Core Utilization

Modern NVIDIA GPUs include specialized Tensor Cores that accelerate matrix multiplication for reduced-precision formats. Tensor Cores on A100 GPUs can perform FP16 matrix multiplication at 312 TFLOPS, compared to 19.5 TFLOPS for standard FP32 CUDA cores—a $16\times$ difference. However, Tensor Cores have alignment requirements: matrix dimensions should be multiples of 8 for FP16 or multiples of 16 for INT8 to achieve peak throughput. This hardware constraint influences architecture design choices.

For BERT-base with $d_k = 64$, the dimension is a multiple of 8, enabling efficient Tensor Core utilization. The query-key multiplication \mathbf{QK}^\top has dimensions $(n \times 64) \times (64 \times n)$, where $n = 512$ is also a multiple of 8. The attention-value multiplication \mathbf{AV} has dimensions $(n \times n) \times (n \times 64)$, again with aligned dimensions. In practice, implementations pad dimensions to the nearest multiple of 8 if necessary. For example, if $d_k = 63$, it would be padded to 64, wasting 1.6% of computation but gaining the $16\times$ Tensor Core speedup—a worthwhile trade-off.

The memory bandwidth requirements for multi-head attention depend on the batch size and sequence length. For BERT-base with batch size 32 and sequence length 512, the QKV projections read $32 \times 512 \times 768 \times 2 = 25,165,824$ bytes (in FP16) and write $3 \times 32 \times 512 \times 64 \times 12 \times 2 = 75,497,472$ bytes for all heads. The attention computation reads these QKV matrices and writes attention outputs, totaling approximately 100 MB of memory traffic per layer. With 12 layers in BERT-base, this amounts to 1.2 GB of memory traffic per forward pass, which takes approximately $1.2/1.5 \approx 0.8$ ms on an A100 with 1.5 TB/s bandwidth. The actual time is higher due to kernel launch overhead, non-coalesced accesses, and compute time, typically around 2-3 ms per forward pass for BERT-base on an A100.

Comparing one head with $d_k = 768$ versus 12 heads with $d_k = 64$ reveals why multiple heads are better for hardware. The single large head would compute attention scores \mathbf{QK}^\top with dimensions $(512 \times 768) \times (768 \times 512)$, requiring $512^2 \times 768 \times 2 = 402,653,184$ FLOPs. The 12 smaller heads each compute $(512 \times 64) \times (64 \times 512)$, requiring $512^2 \times 64 \times 2 = 33,554,432$ FLOPs per head, or $12 \times 33,554,432 = 402,653,184$ FLOPs total—exactly the same. However, the 12 heads can be computed in parallel across different streaming multiprocessors, achieving better GPU utilization. Additionally, the smaller matrices fit better in cache, reducing memory traffic. The single large head would produce an attention matrix of size $512 \times 512 \times 4 = 1,048,576$ bytes, while the 12 smaller heads produce 12 matrices of the same size, totaling 12 MB. The memory usage is higher for multiple heads, but the parallelism and cache benefits outweigh this cost.

8.3 Positional Encoding

Self-attention is inherently permutation equivariant, meaning it treats the input as an unordered set rather than a sequence. If we shuffle the input tokens, the attention mechanism produces correspondingly shuffled outputs, with no awareness that the order has changed. For sequence modeling tasks like language understanding and generation, word order is crucial—"dog bites man" has a very different meaning from "man bites dog." To inject positional information into the model, we add positional encodings to the input embeddings before the first attention layer.

Definition 8.3 (Sinusoidal Positional Encoding). For position pos and dimension i :

$$\text{PE}_{(\text{pos}, 2i)} = \sin\left(\frac{\text{pos}}{10000^{2i/d_{\text{model}}}}\right) \quad (8.22)$$

$$\text{PE}_{(\text{pos}, 2i+1)} = \cos\left(\frac{\text{pos}}{10000^{2i/d_{\text{model}}}}\right) \quad (8.23)$$

The sinusoidal positional encoding has several desirable properties. Each position receives a unique encoding, ensuring that the model can distinguish between different positions. The use of periodic

functions with different frequencies allows the model to potentially extrapolate to longer sequences than seen during training—if the model learns to interpret the sinusoidal patterns, it can apply this understanding to positions beyond the training maximum. Different dimensions use different frequencies, with lower dimensions oscillating rapidly (high frequency) and higher dimensions oscillating slowly (low frequency). This multi-scale representation allows the model to capture both fine-grained local position information and coarse-grained global position information. Finally, the relative position between any two positions can be expressed as a linear transformation of their absolute positional encodings, which may help the model learn relative position relationships.

The usage is straightforward: the positional encoding matrix $\text{PE} \in \mathbb{R}^{n_{\max} \times d_{\text{model}}}$ is precomputed for the maximum sequence length n_{\max} , and for each input sequence of length $n \leq n_{\max}$, we add the first n rows of PE to the token embeddings: $\mathbf{X}_{\text{input}} = \mathbf{X}_{\text{embed}} + \text{PE}_{1:n}$. This addition happens before the first transformer layer, and the positional information propagates through the network via the residual connections.

Example 8.3 (Positional Encoding Values). For $d_{\text{model}} = 512$:

Position 0:

$$\text{PE}_{(0,0)} = \sin(0) = 0 \quad (8.24)$$

$$\text{PE}_{(0,1)} = \cos(0) = 1 \quad (8.25)$$

$$\vdots \quad (8.26)$$

$$\text{PE}_{(0,510)} = \sin(0) = 0 \quad (8.27)$$

$$\text{PE}_{(0,511)} = \cos(0) = 1 \quad (8.28)$$

Position 1:

$$\text{PE}_{(1,0)} = \sin\left(\frac{1}{10000^{0/512}}\right) = \sin(1) \approx 0.841 \quad (8.29)$$

$$\text{PE}_{(1,1)} = \cos\left(\frac{1}{10000^{0/512}}\right) = \cos(1) \approx 0.540 \quad (8.30)$$

Higher dimension indices have lower frequencies (longer periods).

8.3.1 Positional Encoding Variants

While sinusoidal positional encoding was used in the original Transformer, several alternative approaches have been developed, each with different trade-offs in terms of memory usage, extrapolation capability, and performance.

Learned positional embeddings treat position encodings as trainable parameters rather than fixed functions. A learnable embedding matrix $\mathbf{E}_{\text{pos}} \in \mathbb{R}^{n_{\max} \times d_{\text{model}}}$ is initialized randomly and optimized during training alongside other model parameters. This approach is used in BERT and GPT-2. The advantage is that the model can learn position representations optimized for the specific task and data distribution, potentially capturing patterns that sinusoidal encodings cannot express. The disadvantage is memory cost: for BERT with $n_{\max} = 512$ and $d_{\text{model}} = 768$, the positional embeddings require $512 \times 768 \times 4 = 1,572,864$ bytes (1.5 MB) in FP32. More critically, learned positional embeddings do not extrapolate well to longer sequences—if the model is trained on sequences up to length 512, it has never seen positional embeddings for positions 513 and beyond, and these positions must be either extrapolated (often poorly) or the model must be fine-tuned on longer sequences.

Relative positional encoding, used in T5 and Transformer-XL, encodes the relative distance between positions rather than their absolute positions. Instead of adding positional information to the input embeddings, relative position information is incorporated directly into the attention computation. For positions i and j , a learned bias b_{i-j} is added to the attention score, where the bias depends only on the relative distance $i-j$. This requires learning biases for relative distances up to some maximum, typically

± 128 or ± 256 . The memory cost is $O(d_{\text{rel}})$ where d_{rel} is the maximum relative distance, much smaller than the $O(n_{\text{max}} \times d_{\text{model}})$ cost of learned absolute positional embeddings. Relative positional encoding extrapolates well to longer sequences because the model learns to interpret relative distances, which remain meaningful regardless of absolute sequence length. T5 uses a simplified form where relative position biases are shared across attention heads and bucketed into logarithmically-spaced bins, further reducing memory requirements.

Rotary Positional Encoding (RoPE), introduced in RoFormer and used in LLaMA and GPT-NeoX, applies rotation matrices to the query and key vectors based on their positions. For position m , the query and key vectors are rotated by angle $m\theta$ where θ depends on the dimension. Mathematically, for each pair of dimensions $(2i, 2i + 1)$, the rotation is:

$$\begin{bmatrix} q_{2i}^{(m)} \\ q_{2i+1}^{(m)} \end{bmatrix} = \begin{bmatrix} \cos(m\theta_i) & -\sin(m\theta_i) \\ \sin(m\theta_i) & \cos(m\theta_i) \end{bmatrix} \begin{bmatrix} q_{2i} \\ q_{2i+1} \end{bmatrix} \quad (8.31)$$

where $\theta_i = 10000^{-2i/d}$. The key insight is that the dot product between rotated queries and keys naturally encodes relative position: $\mathbf{q}^{(m)} \cdot \mathbf{k}^{(n)} = \mathbf{q} \cdot \mathbf{R}_{\theta}^{n-m} \mathbf{k}$, where $\mathbf{R}_{\theta}^{n-m}$ is a rotation by angle $(n-m)\theta$. RoPE requires no additional parameters—the rotation is computed on-the-fly during attention. It extrapolates excellently to longer sequences because the rotation angles scale linearly with position, and the model learns to interpret relative rotations. RoPE has become the standard for large language models due to its parameter efficiency and strong extrapolation properties.

ALiBi (Attention with Linear Biases), used in BLOOM, adds a simple linear bias to attention scores based on position distance. For query position i attending to key position j , a bias $-m \cdot |i - j|$ is added to the attention score, where m is a head-specific slope. Different heads use different slopes, typically $m = 2^{-8/h}, 2^{-16/h}, \dots, 2^{-8}$ for h heads. This penalizes attention to distant positions, with the penalty strength varying across heads. ALiBi requires no parameters and no additional computation beyond adding the bias. It extrapolates remarkably well: BLOOM was trained on sequences of length 2048 but can generate coherent text at lengths exceeding 8000 tokens. The linear bias naturally extends to any sequence length, and the model learns to work within this inductive bias.

The following table summarizes the trade-offs between positional encoding methods:

Type	Memory	Extrapolation	Used In
Sinusoidal	None	Good	Original Transformer
Learned	$n_{\text{max}} \times d$	Poor	BERT, GPT-2
Relative	$O(d_{\text{rel}})$	Good	T5, Transformer-XL
RoPE	None	Excellent	LLaMA, GPT-NeoX
ALiBi	None	Excellent	BLOOM

The trend in recent large language models has been toward parameter-free methods with strong extrapolation: RoPE and ALiBi dominate current architectures. These methods avoid the memory cost of learned positional embeddings while providing better length generalization than sinusoidal encodings. For practitioners, the choice depends on the application: if sequences will always be shorter than the training maximum, learned embeddings may provide slightly better performance. If length generalization is important, RoPE or ALiBi are superior choices.

8.4 Computational Complexity

8.4.1 Memory Complexity Analysis

The memory requirements of self-attention are dominated by the attention matrices, which scale quadratically with sequence length. For a batch of B sequences, each of length n , with h attention heads, the attention matrices $\mathbf{A}^{(i)} \in \mathbb{R}^{n \times n}$ for $i = 1, \dots, h$ require $O(Bhn^2)$ memory. In FP32, this amounts to $Bhn^2 \times 4$ bytes. For BERT-base with $B = 32$, $h = 12$, and $n = 512$, the attention matrices consume $32 \times 12 \times 512^2 \times 4 = 402,653,184$ bytes, or approximately 384 MB. This quadratic scaling means that doubling the sequence length quadruples the memory requirement: for $n = 1024$, the attention matrices would require 1.5 GB, and for $n = 2048$, they would require 6 GB.

In contrast, the QKV projection matrices and their outputs scale linearly with sequence length. The query, key, and value matrices each have shape $B \times n \times d_k$ for each head, requiring $3Bhnd_k \times 4$ bytes total across all heads. For BERT-base with $d_k = 64$, this amounts to $3 \times 32 \times 12 \times 512 \times 64 \times 4 = 150,994,944$ bytes, or approximately 144 MB. The linear scaling means that doubling the sequence length only doubles this memory requirement.

The crossover point where attention matrices dominate total memory usage depends on the model dimensions. Attention matrices require Bhn^2 elements, while QKV matrices require $3Bhnd_k$ elements. Attention dominates when $Bhn^2 > 3Bhnd_k$, which simplifies to $n > 3d_k$. For BERT-base with $d_k = 64$, attention dominates when $n > 192$ —essentially always, since typical sequence lengths are 512. For models with larger d_k , the crossover occurs at longer sequences. However, since $d_k = d_{\text{model}}/h$ and typical architectures use $h = 12$ to $h = 96$, the value of d_k is usually in the range 64 to 128, meaning attention matrices dominate for sequences longer than a few hundred tokens.

8.4.2 Time Complexity Breakdown

The time complexity of self-attention can be decomposed into several operations, each with different scaling properties. The QKV projections involve three matrix multiplications \mathbf{XW}^Q , \mathbf{XW}^K , and \mathbf{XW}^V , where $\mathbf{X} \in \mathbb{R}^{Bn \times d_{\text{model}}}$ and each weight matrix has shape $d_{\text{model}} \times d_k$. For h heads, the total complexity is $O(3Bhnd_{\text{model}}d_k) = O(Bhnd_{\text{model}}^2)$ since $hd_k = d_{\text{model}}$. This is linear in sequence length n but quadratic in model dimension d_{model} .

Computing the attention scores \mathbf{QK}^\top requires a batch matrix multiplication with dimensions $(Bh \times n \times d_k) \times (Bh \times d_k \times n)$, resulting in complexity $O(Bhn^2d_k)$. This is quadratic in sequence length and linear in head dimension. The softmax operation over the attention scores has complexity $O(Bhn^2)$, dominated by the exponential and normalization computations. Finally, applying the attention weights to the values \mathbf{AV} has complexity $O(Bhn^2d_v)$, again quadratic in sequence length. The output projection $[\text{head}_1; \dots; \text{head}_h]\mathbf{W}^O$ has complexity $O(Bnd_{\text{model}}^2)$, linear in sequence length.

Summing these components, the total complexity is:

$$O(Bnd_{\text{model}}^2) + O(Bhn^2d_k) + O(Bhn^2) + O(Bhn^2d_v) + O(Bnd_{\text{model}}^2) = O(Bnd_{\text{model}}^2 + Bhn^2d_k) \quad (8.32)$$

Since $hd_k = d_{\text{model}}$, this simplifies to $O(Bnd_{\text{model}}^2 + Bn^2d_{\text{model}})$. The relative importance of these terms depends on the ratio n/d_{model} . When $n < d_{\text{model}}$, the $O(Bnd_{\text{model}}^2)$ term from the linear projections dominates, and the feed-forward network (which also has $O(Bnd_{\text{model}}^2)$ complexity) is the computational bottleneck. When $n > d_{\text{model}}$, the $O(Bn^2d_{\text{model}})$ term from attention dominates, and attention becomes the bottleneck.

For BERT-base with $d_{\text{model}} = 768$, attention dominates when $n > 768$. Since BERT uses maximum sequence length 512, the model is in the regime where linear projections and feed-forward networks dominate. For GPT-3 with $d_{\text{model}} = 12,288$, attention would only dominate for sequences longer than 12,288 tokens—far beyond the typical context length of 2048 tokens. This explains why efficient attention mechanisms (Chapter 16) focus on reducing the $O(n^2)$ term: for very long sequences, this term becomes prohibitive, but for typical sequence lengths in large models, the linear terms are actually more expensive.

8.4.3 Scaling Experiments

To illustrate the scaling behavior empirically, consider BERT-base with batch size 1 on an NVIDIA A100 GPU. The following measurements show forward pass time for different sequence lengths:

Sequence Length	Time (ms)	Bottleneck
128	2.1	FFN dominates
256	3.8	FFN dominates
512	8.5	Balanced
1024	28.3	Attention dominates
2048	98.7	Attention dominates
4096	367	Attention dominates

For short sequences (128, 256), the time scales approximately linearly, indicating that the $O(nd_{\text{model}}^2)$ terms dominate. At sequence length 512, the scaling begins to show quadratic behavior. For long sequences (1024, 2048, 4096), the time scales quadratically: doubling from 1024 to 2048 increases time by $98.7/28.3 \approx 3.5\times$, and doubling again to 4096 increases time by $367/98.7 \approx 3.7\times$. The slight deviation from exactly $4\times$ is due to the linear terms and memory bandwidth effects, but the quadratic scaling is clearly visible.

These measurements demonstrate why long-context transformers require specialized attention mechanisms. Extending BERT-base to sequence length 8192 would require approximately $367 \times 4 \approx 1,468$ ms per forward pass, or 1.5 seconds—prohibitively slow for interactive applications. The memory requirement would be $32 \times 12 \times 8192^2 \times 4/(1024^3) \approx 96$ GB for attention matrices alone with batch size 32, exceeding the capacity of even the largest single GPUs. This fundamental scaling limitation motivates the development of sparse attention, linear attention, and other efficient variants discussed in Chapter 16.

8.5 Causal (Masked) Self-Attention

Autoregressive language models like GPT generate text sequentially, predicting each token based only on previous tokens. During training, the entire sequence is provided as input, but the model must not be allowed to “see” future tokens when predicting each position—this would constitute cheating, as the model would have access to information unavailable during generation. Causal masking enforces this constraint by preventing each position from attending to subsequent positions in the sequence.

Definition 8.4 (Causal Mask). Create mask matrix $\mathbf{M} \in \mathbb{R}^{n \times n}$:

$$M_{ij} = \begin{cases} 0 & \text{if } j \leq i \\ -\infty & \text{if } j > i \end{cases} \quad (8.33)$$

Apply before softmax:

$$\mathbf{A} = \text{softmax} \left(\frac{\mathbf{QK}^\top + \mathbf{M}}{\sqrt{d_k}} \right) \quad (8.34)$$

After softmax, $\exp(-\infty) = 0$, so position i cannot attend to positions $j > i$.

Example 8.4 (Causal Mask for Length 4).

$$\mathbf{M} = \begin{bmatrix} 0 & -\infty & -\infty & -\infty \\ 0 & 0 & -\infty & -\infty \\ 0 & 0 & 0 & -\infty \\ 0 & 0 & 0 & 0 \end{bmatrix} \quad (8.35)$$

Position 0 attends only to itself. Position 1 attends to positions 0, 1. Position 3 attends to all positions 0, 1, 2, 3.

This ensures autoregressive property for language modeling.

8.5.1 Efficient Causal Mask Implementation

The naive implementation of causal masking stores a full $n \times n$ mask matrix in memory. For sequence length 2048, this requires $2048^2 \times 4 = 16,777,216$ bytes (16 MB) per sequence in FP32. With batch size 32, this amounts to 512 MB just for the mask—a significant memory overhead. Moreover, the mask must be added to the attention scores before softmax, requiring a memory read of the mask matrix.

Efficient implementations compute the mask on-the-fly during the attention computation rather than storing it explicitly. Modern deep learning frameworks support this through boolean masking or by directly computing the upper triangular structure. For example, in PyTorch, the operation `torch.triu(scores, diagonal=1).fill_(-float('inf'))` modifies the attention scores in-place without allocating a separate mask matrix. This reduces memory usage to zero for the mask itself, though the attention scores matrix must still be stored.

Flash Attention takes this optimization further by fusing the masking operation with the attention computation and tiling the computation to fit in SRAM. Instead of computing the full attention matrix, materializing it in global memory, applying the mask, and then computing softmax, Flash Attention computes attention in tiles that fit in on-chip memory. For each tile, the mask is computed on-the-fly, attention is computed, and the result is written back to global memory. This approach reduces memory usage from $O(n^2)$ to $O(n)$ and provides 2-4 \times speedup for long sequences by minimizing global memory traffic.

The impact of causal masking differs between training and inference. During training, the entire sequence is processed in parallel, with masking ensuring that each position only attends to previous positions. The forward pass computes outputs for all positions simultaneously, and the backward pass computes gradients for all positions simultaneously. The masking is explicit in the attention computation. During inference, text generation is inherently sequential: we generate one token at a time, appending it to the context and generating the next token. In this setting, the masking is implicit—when generating position t , we only have tokens $0, \dots, t-1$ available, so there are no future tokens to mask. However, naive inference would recompute attention over the entire sequence for each new token, resulting in $O(n^2)$ complexity for generating n tokens. Key-value caching addresses this by storing the key and value vectors for all previous tokens, allowing each new token to attend to the cached keys and values without recomputation. This reduces inference complexity to $O(n)$ for generating n tokens, at the cost of $O(nd_{\text{model}})$ memory for the cache.

8.6 Attention Patterns and Interpretability

Analysis of trained transformer models reveals that different attention heads learn to capture different types of linguistic relationships. Some heads focus on syntactic structure, attending strongly between words that have grammatical dependencies such as subject-verb agreement or determiner-noun relationships. For example, in the sentence "The cat that chased the mouse was hungry," a syntactic head might show strong attention from "was" to "cat" (the subject), skipping over the relative clause. Other heads capture semantic relationships, attending between words with similar meanings or words that are topically related. In a sentence about cooking, a semantic head might show attention between "recipe," "ingredients," and "oven," even if these words are not syntactically related.

Positional heads exhibit attention patterns based primarily on token distance rather than content. Some heads attend primarily to adjacent tokens, capturing local context. Others attend to tokens at specific relative positions, such as attending to the previous token or to tokens at fixed offsets. These positional patterns can be useful for tasks like copying or for capturing regular linguistic structures. Rare word heads show distinctive behavior where attention is concentrated on infrequent tokens, potentially allowing the model to give special processing to unusual or important words that might otherwise be overwhelmed by common function words.

Attention visualization provides insight into model behavior by displaying the attention weights as heatmaps or graphs. For a given input sentence, we can visualize the attention distribution for each head in each layer, showing which tokens each position attends to. These visualizations often reveal interpretable patterns: early layers tend to focus on local, syntactic relationships, while later layers capture more abstract, semantic relationships. However, interpretation must be approached with caution—attention weights show where the model looks, but not necessarily what information is extracted or how it is used. High attention weight does not necessarily imply high importance for the final prediction.

Research on attention head importance has shown that many heads can be pruned without significant performance degradation. In BERT-base with 144 attention heads (12 heads per layer \times 12 layers), pruning 40-50% of heads typically reduces downstream task accuracy by less than 1%. This

suggests substantial redundancy in the multi-head attention mechanism. Some heads are consistently important across tasks—often those capturing syntactic relationships or attending to special tokens like [CLS] or [SEP]. Other heads appear to be less critical, and their removal has minimal impact. This redundancy may serve an important role during training by providing multiple gradient pathways and helping optimization, even if the final model does not require all heads for inference.

8.7 Hardware-Specific Optimizations

8.7.1 Flash Attention

Flash Attention represents a fundamental rethinking of how attention is computed on modern GPUs. The standard attention implementation computes the full attention matrix $\mathbf{A} = \text{softmax}(\mathbf{Q}\mathbf{K}^\top/\sqrt{d_k})$ and materializes it in GPU global memory before multiplying by \mathbf{V} . For long sequences, this attention matrix is large—for sequence length 4096, a single attention head requires $4096^2 \times 4 = 67$ MB in FP32. Reading and writing this matrix to global memory becomes the performance bottleneck, as global memory bandwidth (approximately 1.5 TB/s on an A100) is much lower than compute throughput (312 TFLOPS for FP16).

Flash Attention addresses this by tiling the attention computation to fit in SRAM, the fast on-chip memory available on each streaming multiprocessor. SRAM has much higher bandwidth (approximately 19 TB/s on A100) but limited capacity (192 KB per SM, totaling about 40 MB across all SMs). The key insight is that attention can be computed in blocks: we partition the query, key, and value matrices into tiles, load each tile into SRAM, compute attention for that tile, and accumulate the results. The attention matrix is never fully materialized in global memory—only the tiles currently being processed reside in SRAM.

The tiling strategy works as follows. Partition the queries into blocks $\mathbf{Q}_1, \dots, \mathbf{Q}_T$ and the keys and values into blocks $\mathbf{K}_1, \mathbf{V}_1, \dots, \mathbf{K}_T, \mathbf{V}_T$. For each query block \mathbf{Q}_i , iterate over all key-value blocks $(\mathbf{K}_j, \mathbf{V}_j)$, computing the attention contribution $\text{softmax}(\mathbf{Q}_i\mathbf{K}_j^\top/\sqrt{d_k})\mathbf{V}_j$ in SRAM and accumulating the results. The softmax normalization requires special handling since we compute it in blocks—we maintain running statistics (maximum and sum of exponentials) and update them as we process each block, then renormalize at the end. This online softmax algorithm ensures numerical stability while avoiding materialization of the full attention matrix.

The memory usage of Flash Attention is $O(n)$ rather than $O(n^2)$, as we only store the query, key, and value matrices (each $O(n \times d)$) and the output, not the attention matrix. The computational cost remains the same—we perform the same number of FLOPs as standard attention—but the memory traffic is dramatically reduced. For sequence length 4096 with $d_k = 64$, standard attention reads/writes approximately $4096^2 \times 4 = 67$ MB for the attention matrix, while Flash Attention reads/writes only the QKV matrices, approximately $3 \times 4096 \times 64 \times 4 = 3$ MB. This $20\times$ reduction in memory traffic translates to $2\text{--}4\times$ speedup in practice, with larger speedups for longer sequences where memory bandwidth is the primary bottleneck.

8.7.2 Fused Kernels

Kernel fusion combines multiple operations into a single GPU kernel, reducing memory traffic by keeping intermediate results in registers or shared memory rather than writing them to global memory. For attention, a common fusion is combining the softmax operation with the attention score computation and the multiplication by values. The standard implementation computes $\mathbf{S} = \mathbf{Q}\mathbf{K}^\top$, writes \mathbf{S} to global memory, launches a separate kernel to compute $\mathbf{A} = \text{softmax}(\mathbf{S})$, writes \mathbf{A} to global memory, and launches another kernel to compute $\mathbf{O} = \mathbf{A}\mathbf{V}$. Each write and read to global memory incurs latency and consumes bandwidth.

A fused attention kernel computes all these operations in a single kernel launch. The kernel loads tiles of \mathbf{Q} , \mathbf{K} , and \mathbf{V} into shared memory, computes attention scores in registers, applies softmax, multiplies by values, and writes the final output—all without intermediate global memory traffic. This fusion reduces memory bandwidth requirements by approximately $2\times$, as we eliminate the reads and

writes of \mathbf{S} and \mathbf{A} . The speedup is typically $1.5\text{--}2\times$ for attention-dominated workloads, with larger benefits for smaller batch sizes where memory bandwidth is the primary bottleneck.

Fused kernels require careful implementation to maximize occupancy and minimize register pressure. The kernel must balance the tile size (larger tiles reduce global memory traffic but increase shared memory and register usage) with occupancy (the number of thread blocks that can run concurrently on each SM). Modern deep learning frameworks like PyTorch and TensorFlow provide fused attention implementations through libraries like cuDNN and custom CUDA kernels, making these optimizations accessible without manual kernel development.

8.7.3 Tensor Core Optimization

Tensor Cores on NVIDIA GPUs provide specialized hardware for matrix multiplication, achieving much higher throughput than standard CUDA cores for reduced-precision formats. To fully utilize Tensor Cores, matrix dimensions should be multiples of 8 for FP16 or multiples of 16 for INT8. For attention, this means padding d_k , n , and batch size to these multiples when necessary. For example, if $d_k = 63$, padding to 64 wastes 1.6% of computation but enables Tensor Core usage, providing a net speedup of $10\text{--}15\times$.

The WMMA (Warp Matrix Multiply-Accumulate) API provides access to Tensor Cores from CUDA code. A warp (32 threads) cooperatively loads matrix tiles into registers, performs matrix multiplication using Tensor Cores, and stores the result. For attention, the query-key multiplication \mathbf{QK}^\top and the attention-value multiplication \mathbf{AV} are both matrix multiplications that can leverage Tensor Cores. Achieving 70-80% of peak TFLOPS requires careful attention to data layout (row-major vs column-major), tile sizes, and memory access patterns to ensure coalesced loads and stores.

In practice, modern deep learning frameworks handle Tensor Core optimization automatically for standard operations like matrix multiplication. However, custom attention implementations or fused kernels may require explicit use of WMMA or the higher-level cuBLAS library to achieve peak performance. The key takeaway for practitioners is that attention performance depends critically on matrix dimensions being multiples of 8 or 16, and that padding dimensions to meet this requirement is almost always worthwhile.

8.8 Memory-Efficient Attention Variants

The quadratic memory and time complexity of standard attention motivates the development of approximate attention mechanisms that reduce complexity while maintaining most of the modeling power. These variants make different trade-offs between accuracy, efficiency, and implementation complexity.

8.8.1 Sparse Attention

Sparse attention restricts each query to attend to only a subset of keys, reducing the attention matrix from dense $O(n^2)$ to sparse $O(ns)$ where $s \ll n$ is the sparsity pattern size. The challenge is choosing which positions to attend to. Local attention restricts each position to attend only to a window of nearby positions, typically $\pm w$ positions for window size w . This captures local context efficiently with $O(nw)$ complexity, but loses long-range dependencies. Strided attention attends to every k -th position, allowing long-range connections with $O(n^2/k)$ complexity. Random attention attends to a random subset of positions, providing a probabilistic approximation with $O(ns)$ complexity where s is the number of random positions sampled.

The Sparse Transformer combines local and strided attention in different heads, allowing some heads to capture local patterns while others capture long-range patterns. For sequence length n , using \sqrt{n} local positions and \sqrt{n} strided positions gives $O(n\sqrt{n})$ complexity—a significant improvement over $O(n^2)$ for long sequences. Longformer extends this with a combination of local attention (for all positions), global attention (for special tokens like [CLS]), and dilated attention (strided with increasing stride in deeper layers). These patterns are task-specific: document classification might use global attention on the [CLS] token, while question answering might use global attention on question tokens.

The accuracy trade-off depends on the task and sparsity pattern. For tasks requiring primarily local context (e.g., language modeling with local coherence), sparse attention with window size 512 typically loses less than 1% accuracy compared to full attention. For tasks requiring long-range reasoning (e.g., document-level question answering), the accuracy loss can be 2-5% unless the sparsity pattern is carefully designed to preserve critical long-range connections. The memory savings are substantial: for sequence length 4096 with window size 512, sparse attention uses $4096 \times 512 = 2,097,152$ elements instead of $4096^2 = 16,777,216$ elements, an $8\times$ reduction.

8.8.2 Linear Attention

Linear attention approximates the softmax attention mechanism using kernel methods to achieve $O(nd^2)$ complexity instead of $O(n^2d)$. The key insight is that the attention output can be rewritten using the associative property of matrix multiplication:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}(\mathbf{Q}\mathbf{K}^\top)\mathbf{V} = \frac{\exp(\mathbf{Q})(\exp(\mathbf{K})^\top\mathbf{V})}{\exp(\mathbf{Q})(\exp(\mathbf{K})^\top\mathbf{1})} \quad (8.36)$$

where \exp is applied element-wise. By computing $\exp(\mathbf{K})^\top\mathbf{V}$ first (complexity $O(nd^2)$) and then multiplying by $\exp(\mathbf{Q})$ (complexity $O(nd^2)$), we avoid the $O(n^2)$ attention matrix.

The approximation lies in replacing the softmax kernel $\exp(\mathbf{q} \cdot \mathbf{k})$ with a more efficient kernel $\phi(\mathbf{q}) \cdot \phi(\mathbf{k})$ where ϕ is a feature map. Common choices include $\phi(\mathbf{x}) = \text{elu}(\mathbf{x}) + 1$ (ensuring non-negativity) or random Fourier features. The attention becomes:

$$\text{LinearAttn}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \frac{\phi(\mathbf{Q})(\phi(\mathbf{K})^\top\mathbf{V})}{\phi(\mathbf{Q})(\phi(\mathbf{K})^\top\mathbf{1})} \quad (8.37)$$

This can be computed in $O(nd^2)$ time by first computing $\mathbf{S} = \phi(\mathbf{K})^\top\mathbf{V} \in \mathbb{R}^{d \times d}$ (cost $O(nd^2)$), then $\phi(\mathbf{Q})\mathbf{S}$ (cost $O(nd^2)$). The memory requirement is $O(nd)$ for the QKV matrices and $O(d^2)$ for \mathbf{S} , avoiding the $O(n^2)$ attention matrix.

The accuracy trade-off is more significant than for sparse attention. Linear attention typically loses 1-3% accuracy on language modeling tasks compared to softmax attention, as the kernel approximation does not perfectly capture the softmax distribution. The approximation is particularly poor for distributions with sharp peaks (high attention to a single position), which are common in tasks like copying or attending to specific keywords. However, for tasks where attention is more diffuse, linear attention can be nearly as accurate as softmax attention while being much faster for long sequences. For sequence length 8192, linear attention is approximately $4\times$ faster than standard attention and uses $8\times$ less memory.

8.8.3 Low-Rank Attention

Low-rank attention factorizes the attention matrix as $\mathbf{A} = \mathbf{U}\mathbf{V}^\top$ where $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{n \times r}$ and $r \ll n$ is the rank. This reduces memory from $O(n^2)$ to $O(nr)$ and computation from $O(n^2d)$ to $O(nrd)$. Linformer implements this by projecting keys and values to a lower-dimensional space: $\mathbf{K}_{\text{proj}} = \mathbf{E}\mathbf{K}$ and $\mathbf{V}_{\text{proj}} = \mathbf{F}\mathbf{V}$ where $\mathbf{E}, \mathbf{F} \in \mathbb{R}^{r \times n}$ are projection matrices. The attention is then computed as:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}_{\text{proj}}, \mathbf{V}_{\text{proj}}) = \text{softmax}(\mathbf{Q}\mathbf{K}_{\text{proj}}^\top)\mathbf{V}_{\text{proj}} \quad (8.38)$$

The attention matrix has shape $n \times r$ instead of $n \times n$, reducing memory and computation.

The projection matrices \mathbf{E} and \mathbf{F} can be learned or fixed (e.g., random projections or selecting every n/r -th position). Learned projections provide better accuracy but add parameters and require careful initialization. The rank r controls the trade-off: $r = 256$ typically provides good accuracy for sequence lengths up to 4096, while $r = 64$ may suffice for shorter sequences. The accuracy loss is typically 1-2% for appropriate choice of r , as the low-rank approximation captures the dominant patterns in the attention matrix while discarding less important details.

These memory-efficient attention variants are explored in greater depth in Chapter 16, where we discuss their application to long-context transformers and provide implementation details. The key

insight is that the $O(n^2)$ complexity of standard attention is not fundamental to the transformer architecture—it is a consequence of the specific attention mechanism used. By carefully approximating or restricting attention, we can build transformers that scale to much longer sequences while maintaining most of the modeling power of full attention.

8.9 Exercises

Exercise 8.1. For GPT-2 ($d_{\text{model}} = 1024$, $h = 16$, $n = 1024$): (1) Compute attention matrix memory in MB (float32), (2) Count parameters in one multi-head attention layer, (3) Estimate FLOPs for single forward pass.

Exercise 8.2. Implement multi-head attention in PyTorch. Test with batch size 32, sequence length 20, $d_{\text{model}} = 128$, 4 heads. Verify output shape and parameter count.

Exercise 8.3. Show that sinusoidal positional encoding allows computing $\text{PE}_{\text{pos}+k}$ as linear function of PE_{pos} for any offset k .

Exercise 8.4. Compare attention weights with and without positional encoding. Show numerically how word order affects attention without PE.

Chapter 9

Attention Variants and Mechanisms

Chapter Overview

Beyond standard scaled dot-product attention, numerous variants have been developed for specific use cases and improved efficiency. This chapter explores cross-attention for encoder-decoder models, soft vs hard attention, attention with relative position representations, and practical considerations for implementing attention mechanisms.

Learning Objectives

1. Distinguish between self-attention and cross-attention
2. Understand relative position representations
3. Implement attention with different scoring functions
4. Apply attention masking for various scenarios
5. Understand attention dropout and layer normalization
6. Visualize and interpret attention patterns

9.1 Cross-Attention

Definition 9.1 (Cross-Attention). In encoder-decoder architectures, decoder attends to encoder output via cross-attention:

$$\mathbf{Q} = \mathbf{X}_{\text{dec}} \mathbf{W}^Q \quad (\text{queries from decoder}) \quad (9.1)$$

$$\mathbf{K} = \mathbf{X}_{\text{enc}} \mathbf{W}^K \quad (\text{keys from encoder}) \quad (9.2)$$

$$\mathbf{V} = \mathbf{X}_{\text{enc}} \mathbf{W}^V \quad (\text{values from encoder}) \quad (9.3)$$

$$\text{CrossAttn}(\mathbf{X}_{\text{dec}}, \mathbf{X}_{\text{enc}}) = \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}} \right) \mathbf{V} \quad (9.4)$$

Dimensions:

- Decoder input: $\mathbf{X}_{\text{dec}} \in \mathbb{R}^{m \times d}$ (m decoder positions)
- Encoder output: $\mathbf{X}_{\text{enc}} \in \mathbb{R}^{n \times d}$ (n encoder positions)
- Attention matrix: $\mathbf{A} \in \mathbb{R}^{m \times n}$ (decoder \times encoder)
- Output: $\mathbb{R}^{m \times d_v}$ (same decoder length)

Example 9.1 (Machine Translation Cross-Attention). English source: "The cat sat" (3 tokens encoded to $\mathbf{X}_{\text{enc}} \in \mathbb{R}^{3 \times 512}$)

French target: "Le chat" (2 tokens so far, $\mathbf{X}_{\text{dec}} \in \mathbb{R}^{2 \times 512}$)

Cross-attention computes:

$$\mathbf{A} = \begin{bmatrix} \alpha_{1,1} & \alpha_{1,2} & \alpha_{1,3} \\ \alpha_{2,1} & \alpha_{2,2} & \alpha_{2,3} \end{bmatrix} \in \mathbb{R}^{2 \times 3} \quad (9.5)$$

where $\alpha_{1,j}$ = attention from decoder position 1 ("Le") to encoder position j .

When generating "Le" (the), model should attend strongly to "The" in source.

When generating "chat" (cat), model should attend strongly to "cat" in source.

9.1.1 Transformer Decoder Attention Layers

A transformer decoder block contains **three** attention mechanisms:

1. **Masked self-attention:** Decoder attends to previous decoder positions

$$\mathbf{Q} = \mathbf{K} = \mathbf{V} = \mathbf{X}_{\text{dec}} \quad (\text{with causal mask}) \quad (9.6)$$

2. **Cross-attention:** Decoder attends to encoder output

$$\mathbf{Q} = \mathbf{X}_{\text{dec}}, \quad \mathbf{K} = \mathbf{V} = \mathbf{X}_{\text{enc}} \quad (9.7)$$

3. **Feed-forward:** Position-wise MLP (not attention)

Key Point 9.1. *Encoder-only models (BERT) use only self-attention. Decoder-only models (GPT) use only masked self-attention. Encoder-decoder models (T5, BART) use all three mechanisms.*

9.2 Relative Position Representations

Problem with absolute positions: Model learns positions 0-512 during training. How to handle position 600 at inference?

Solution: Relative position representations—encode distance between positions, not absolute positions.

9.2.1 Shaw et al. Relative Attention

Definition 9.2 (Relative Position Attention). Modify attention scores to include relative position information:

$$e_{ij} = \frac{\mathbf{q}_i^\top \mathbf{k}_j}{\sqrt{d_k}} + \mathbf{q}_i^\top \mathbf{K}_{i-j} \quad (9.8)$$

where $\mathbf{K}_{i-j} \in \mathbb{R}^{d_k}$ encodes relative position $i - j$ (clipped to maximum distance).

Advantages:

- Generalize to longer sequences
- Model learns distance-based patterns
- More parameter efficient

9.2.2 T5 Relative Position Bias

T5 uses even simpler approach—add learned bias based on relative position:

$$\mathbf{A}_{ij} = \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}} + \mathbf{B} \right)_{ij} \quad (9.9)$$

where B_{ij} depends only on $|i - j|$ (bucketed by distance).

9.3 Alternative Attention Scoring Functions

Beyond scaled dot-product, various scoring functions exist:

9.3.1 Additive (Bahdanau)

$$\text{score}(\mathbf{q}, \mathbf{k}) = \mathbf{v}^\top \tanh(\mathbf{W}_1 \mathbf{q} + \mathbf{W}_2 \mathbf{k}) \quad (9.10)$$

9.3.2 Multiplicative (Luong)

$$\text{score}(\mathbf{q}, \mathbf{k}) = \mathbf{q}^\top \mathbf{W} \mathbf{k} \quad (9.11)$$

9.3.3 Scaled Dot-Product (Transformers)

$$\text{score}(\mathbf{q}, \mathbf{k}) = \frac{\mathbf{q}^\top \mathbf{k}}{\sqrt{d_k}} \quad (9.12)$$

9.3.4 General

$$\text{score}(\mathbf{q}, \mathbf{k}) = \mathbf{q}^\top \mathbf{W} \mathbf{k} \quad (9.13)$$

Comparison:

- **Additive:** More parameters, handles different dimensions
- **Dot-product:** Efficient, used in transformers
- **General:** Flexible but more parameters

9.4 Attention Masking

9.4.1 Padding Mask

For variable-length sequences in batch, mask padding tokens:

$$M_{ij} = \begin{cases} 0 & \text{if position } j \text{ is valid} \\ -\infty & \text{if position } j \text{ is padding} \end{cases} \quad (9.14)$$

Example 9.2 (Padding Mask). Batch with sequences of length [5, 7, 4], padded to length 7:

$$\text{Seq 1: } [w_1, w_2, w_3, w_4, w_5, \text{PAD}, \text{PAD}] \quad (9.15)$$

$$\text{Seq 2: } [w_1, w_2, w_3, w_4, w_5, w_6, w_7] \quad (9.16)$$

$$\text{Seq 3: } [w_1, w_2, w_3, w_4, \text{PAD}, \text{PAD}, \text{PAD}] \quad (9.17)$$

Mask for Seq 1:

$$[0, 0, 0, 0, 0, -\infty, -\infty] \quad (9.18)$$

Prevents attending to padding tokens.

9.4.2 Combined Masks

For decoder, combine causal mask and padding mask:

$$\mathbf{M}_{\text{total}} = \mathbf{M}_{\text{causal}} + \mathbf{M}_{\text{padding}} \quad (9.19)$$

Element-wise, use most restrictive: if either mask blocks, result blocks.

9.5 Attention Dropout

Apply dropout to attention weights for regularization:

$$\mathbf{A} = \text{Dropout} \left(\text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}} \right) \right) \quad (9.20)$$

Typical dropout rate: 0.1 (10%)

Effect: Randomly zero out some attention connections, preventing over-reliance on specific positions.

9.6 Layer Normalization with Attention

Two architectures for combining attention with layer norm:

9.6.1 Post-Norm (Original Transformer)

$$\mathbf{h} = \mathbf{X} + \text{MultiHeadAttn}(\mathbf{X}) \quad (9.21)$$

$$\mathbf{Z} = \text{LayerNorm}(\mathbf{h}) \quad (9.22)$$

9.6.2 Pre-Norm (More Common Now)

$$\mathbf{h} = \mathbf{X} + \text{MultiHeadAttn}(\text{LayerNorm}(\mathbf{X})) \quad (9.23)$$

$$\mathbf{Z} = \mathbf{h} \quad (9.24)$$

Pre-norm advantages:

- More stable training
- Easier gradient flow
- Used in GPT-2, GPT-3, modern transformers

9.7 Visualizing Attention

Attention weights $\mathbf{A} \in \mathbb{R}^{n \times n}$ reveal what model attends to:

9.7.1 Attention Heatmaps

For sentence "The cat sat on the mat":

- Row i : attention distribution when processing token i
- Bright cell (i, j) : token i strongly attends to token j

Patterns observed:

- Diagonal: Attending to self
- Vertical lines: Attending to specific important words (e.g., subject, verb)
- Symmetric patterns: Mutual attention between related words
- Head-specific patterns: Different heads learn different relationships

9.7.2 Interpreting Multiple Heads

In 12-head attention, different heads specialize:

- Some heads attend to adjacent words (local syntax)
- Some heads attend to distant words (long-range dependencies)
- Some heads attend to specific parts of speech
- Some heads attend based on semantic similarity

Attention weights are NOT necessarily model explanations! High attention doesn't always mean high importance for prediction. Attention shows where model looks, not why decisions are made.

9.8 Practical Implementation Considerations

9.8.1 Memory-Efficient Attention

For very long sequences, store attention matrix in chunks:

1. Compute \mathbf{QK}^\top for chunk of queries
2. Apply softmax
3. Multiply by \mathbf{V} chunk
4. Accumulate results

Reduces peak memory from $O(n^2)$ to $O(nc)$ where c is chunk size.

9.8.2 Fused Attention Kernels

Modern implementations fuse operations:

$$\mathbf{QK}^\top \rightarrow \text{Scale} \rightarrow \text{Mask} \rightarrow \text{Softmax} \rightarrow \text{Dropout} \rightarrow \text{multiply } \mathbf{V} \quad (9.25)$$

Single fused kernel faster than separate operations (fewer memory transfers).

Example: FlashAttention achieves 2-4x speedup through fused operations and memory hierarchy optimization.

9.9 Efficient Attention Variants

The standard self-attention mechanism has computational complexity $O(n^2d)$ and memory complexity $O(n^2)$, where n is the sequence length and d is the model dimension. This quadratic scaling in sequence length becomes prohibitive for long sequences. For a sequence of length 4096 with 12 attention heads, the attention matrices alone require $12 \times 4096^2 \times 4 = 805$ MB in FP32 format per example. With batch size 32, this amounts to 25.8 GB just for attention weights, exceeding the memory capacity of most GPUs. This fundamental limitation has motivated extensive research into efficient attention variants that reduce the quadratic complexity while maintaining model quality.

The key insight underlying efficient attention is that not all token pairs require equal attention. In practice, attention patterns often exhibit structure—tokens primarily attend to nearby tokens, specific global tokens, or sparse subsets of the sequence. By exploiting this structure, efficient attention mechanisms can dramatically reduce computational and memory requirements while preserving most of the modeling capacity of full attention. The following sections examine the major classes of efficient attention variants, analyzing their complexity trade-offs, implementation considerations, and practical use cases.

9.9.1 Local Attention

Local attention restricts each token to attend only to tokens within a fixed window around its position, rather than attending to all tokens in the sequence. For a window size w , token at position i attends only to positions $[i - w/2, i + w/2]$. This reduces the attention matrix from $n \times n$ to $n \times w$, yielding linear scaling in sequence length.

The computational complexity of local attention is $O(nwd)$, where n is sequence length, w is window size, and d is model dimension. Compared to standard attention's $O(n^2d)$, this represents a reduction factor of n/w . For a sequence of length 4096 with window size 256, local attention is 16 times faster than full attention. The memory complexity similarly reduces from $O(n^2)$ to $O(nw)$, enabling much longer sequences to fit in GPU memory. For the same 4096-token sequence with 12 heads, local attention with window 256 requires only $12 \times 4096 \times 256 \times 4 = 50.3$ MB per example, a 16-fold reduction from the 805 MB required by full attention.

The primary trade-off of local attention is the loss of long-range dependencies. Tokens separated by more than $w/2$ positions cannot directly attend to each other, requiring information to propagate through multiple layers. In practice, this limitation is often acceptable. Many natural language tasks exhibit strong locality—syntactic dependencies are typically short-range, and semantic relationships can be captured through multiple layers of local attention. Empirical studies show that local attention with window size 256-512 typically achieves 98-99% of full attention's accuracy on language modeling tasks, while enabling sequences 10-20 times longer.

The Longformer architecture demonstrates effective use of local attention for document-level understanding. Longformer combines local windowed attention for most tokens with global attention for special tokens like [CLS] and task-specific tokens. This hybrid approach maintains $O(n)$ complexity while allowing critical tokens to aggregate information from the entire sequence. On document classification tasks with 4096-token inputs, Longformer achieves comparable accuracy to BERT while processing sequences 8 times longer. The local attention pattern also enables efficient implementation on GPUs through blocked matrix operations, achieving 2-3x speedup over naive implementations.

9.9.2 Sparse Attention

Sparse attention generalizes local attention by allowing each token to attend to a sparse subset of positions according to a predefined pattern, rather than a contiguous window. The key insight is that attention patterns in trained transformers often exhibit structure—certain positions are consistently important while others receive minimal attention. By designing sparsity patterns that capture this structure, sparse attention can dramatically reduce computation while maintaining model quality.

Several sparsity patterns have proven effective in practice. Strided attention divides the sequence into blocks and allows each token to attend within its block and to every k -th token globally, where k is the stride. This pattern captures both local context and evenly-spaced global context. Fixed

attention combines local attention with attention to a fixed set of global tokens, similar to Longformer. Learned sparse attention uses a separate network to predict which positions each token should attend to, adapting the sparsity pattern to the input. The Sparse Transformer architecture uses a factorized attention pattern where each token attends to positions in a strided pattern in one head and a local pattern in another head, allowing information to flow efficiently across the sequence.

The computational complexity of sparse attention is $O(n\sqrt{n}d)$ for typical sparsity patterns, where each token attends to approximately \sqrt{n} other tokens. This represents a substantial improvement over full attention's $O(n^2d)$, particularly for long sequences. For a sequence of length 4096, sparse attention with $\sqrt{n} = 64$ positions per token is 64 times faster than full attention. The memory complexity is similarly $O(n\sqrt{n})$, enabling sequences that would be impossible with full attention. For 4096 tokens with 12 heads, sparse attention requires approximately $12 \times 4096 \times 64 \times 4 = 12.6$ MB per example, a 64-fold reduction from full attention's 805 MB.

The accuracy trade-off of sparse attention depends critically on the choice of sparsity pattern. Well-designed patterns that align with the task's dependency structure can achieve 97-99% of full attention's accuracy. The Sparse Transformer achieves perplexity within 0.1 of full attention on language modeling while using only \sqrt{n} attention per token. BigBird, which combines local, global, and random attention patterns, matches BERT's accuracy on question answering and document classification while processing sequences up to 8 times longer. However, poorly chosen sparsity patterns can significantly degrade accuracy, particularly on tasks requiring long-range reasoning.

Implementation of sparse attention on GPUs presents challenges because modern GPUs are optimized for dense matrix operations. Sparse matrix multiplication is less efficient than dense multiplication due to irregular memory access patterns and reduced arithmetic intensity. Specialized kernels and libraries like cuSPARSE can partially mitigate this, but sparse attention typically achieves only 50-70% of the theoretical speedup in practice. Recent work on block-sparse attention, which operates on blocks of the attention matrix rather than individual elements, achieves better GPU utilization by maintaining some regularity in memory access patterns. The Triton framework enables efficient implementation of custom sparse attention patterns through automatic optimization of memory access.

9.9.3 Linear Attention

Linear attention achieves $O(nd^2)$ complexity by reformulating the attention computation to avoid explicitly constructing the $n \times n$ attention matrix. The key insight is that attention can be viewed as a kernel operation, and by choosing an appropriate kernel function, the computation can be reordered to compute the output directly without materializing the full attention matrix.

The standard attention computation is:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}(\mathbf{Q}\mathbf{K}^\top)\mathbf{V} \quad (9.26)$$

This requires computing $\mathbf{Q}\mathbf{K}^\top \in \mathbb{R}^{n \times n}$ before applying softmax and multiplying by \mathbf{V} . Linear attention approximates the softmax kernel with a feature map $\phi: \mathbb{R}^{d_k} \rightarrow \mathbb{R}^{d'}$ such that:

$$\text{softmax}(\mathbf{q}^\top \mathbf{k}) \approx \phi(\mathbf{q})^\top \phi(\mathbf{k}) \quad (9.27)$$

With this approximation, attention becomes:

$$\text{LinearAttn}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \phi(\mathbf{Q})(\phi(\mathbf{K})^\top \mathbf{V}) \quad (9.28)$$

The crucial observation is that the parentheses can be reordered. Instead of computing $\phi(\mathbf{Q})\phi(\mathbf{K})^\top$ (which is $n \times n$) and then multiplying by \mathbf{V} , we first compute $\phi(\mathbf{K})^\top \mathbf{V} \in \mathbb{R}^{d' \times d_v}$ and then multiply by $\phi(\mathbf{Q})$. This reordering changes complexity from $O(n^2d)$ to $O(nd'^2)$, where d' is the feature dimension (typically equal to d_k).

The computational savings of linear attention are substantial for long sequences. For sequence length 4096 and model dimension 768, standard attention requires approximately $4096^2 \times 768 = 12.9$ billion operations per head, while linear attention requires $4096 \times 768^2 = 2.4$ billion operations—a 5.4x reduction. The memory complexity is even more favorable: linear attention requires only $O(nd)$ memory for the intermediate $\phi(\mathbf{K})^\top \mathbf{V}$ matrix, compared to $O(n^2)$ for the full attention matrix. For

4096 tokens with 12 heads, linear attention requires approximately $12 \times 768 \times 768 \times 4 = 28.3$ MB, compared to 805 MB for full attention—a 28-fold reduction.

The primary challenge of linear attention is choosing a feature map ϕ that accurately approximates the softmax kernel while remaining computationally efficient. The Performer architecture uses random Fourier features with $\phi(\mathbf{x}) = \exp(\mathbf{x}^2/2)[\cos(\omega_1^\top \mathbf{x}), \sin(\omega_1^\top \mathbf{x}), \dots]$ where ω_i are random projection vectors. This provides an unbiased approximation of the softmax kernel with controllable accuracy based on the number of random features. The Linear Transformer uses a simpler feature map $\phi(\mathbf{x}) = \text{elu}(\mathbf{x}) + 1$, which is faster to compute but provides a looser approximation.

The accuracy trade-off of linear attention is more significant than local or sparse attention. Empirical studies show that linear attention typically achieves 95-98% of full attention’s accuracy on language modeling, with larger degradation on tasks requiring precise attention patterns. The approximation error is particularly noticeable for small attention weights—the softmax function’s sharp peaking is difficult to approximate with simple feature maps. However, for applications where extreme sequence length is critical, such as processing entire books or long-form video, the 2-5% accuracy loss is often acceptable given the dramatic computational savings. Recent work on learned feature maps and adaptive kernel approximations aims to close this accuracy gap while maintaining linear complexity.

9.9.4 Low-Rank Attention

Low-rank attention exploits the observation that attention matrices in trained transformers often have low effective rank—most of the variance is captured by a small number of singular values. By explicitly factorizing the attention computation through a low-dimensional bottleneck, low-rank attention reduces complexity from $O(n^2d)$ to $O(nrd)$, where r is the rank and typically $r \ll n$.

The Linformer architecture implements low-rank attention by projecting the keys and values to a lower-dimensional space before computing attention. Specifically, Linformer adds projection matrices $\mathbf{E}, \mathbf{F} \in \mathbb{R}^{r \times n}$ that reduce the sequence length dimension:

$$\text{LinformerAttn}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax} \left(\frac{\mathbf{Q}(\mathbf{E}\mathbf{K})^\top}{\sqrt{d_k}} \right) (\mathbf{F}\mathbf{V}) \quad (9.29)$$

The key insight is that $\mathbf{E}\mathbf{K} \in \mathbb{R}^{r \times d_k}$ and $\mathbf{F}\mathbf{V} \in \mathbb{R}^{r \times d_v}$ have reduced sequence length r instead of n . The attention matrix is now $n \times r$ instead of $n \times n$, reducing both computation and memory by a factor of n/r .

For sequence length 4096 and rank 256, low-rank attention reduces computation from $4096^2 \times 768 = 12.9$ billion operations to $4096 \times 256 \times 768 = 805$ million operations per head—a 16-fold reduction. The memory savings are equally dramatic: the attention matrix requires $4096 \times 256 \times 4 = 4.2$ MB per head instead of $4096^2 \times 4 = 67.1$ MB, a 16-fold reduction. With 12 heads, total attention memory drops from 805 MB to 50.3 MB per example.

The accuracy of low-rank attention depends on the choice of rank r and the projection matrices \mathbf{E} and \mathbf{F} . Linformer uses learned projection matrices that are shared across all layers, reducing the parameter overhead. Empirical studies show that rank $r = 256$ achieves 96-98% of full attention’s accuracy for sequences up to 4096 tokens, with minimal degradation on most language understanding tasks. The accuracy loss is more pronounced for tasks requiring fine-grained attention patterns, such as coreference resolution or syntactic parsing, where the low-rank approximation may miss subtle dependencies.

An important consideration for low-rank attention is that the projection matrices \mathbf{E} and \mathbf{F} introduce additional parameters and computation. For rank r and sequence length n , the projections add $2rn$ parameters per layer. However, these projections can be implemented efficiently as 1D convolutions or learned position-wise projections, and the parameter cost is typically small compared to the savings in attention computation. The projection operations themselves require $O(rnd)$ computation, which is negligible compared to the $O(n^2d)$ cost of full attention for $r \ll n$.

9.9.5 Comprehensive Complexity Comparison

Understanding the trade-offs between different attention variants requires examining multiple dimensions: computational complexity, memory requirements, accuracy preservation, and practical imple-

mentation efficiency. The following analysis provides concrete comparisons across these dimensions for typical transformer configurations.

Table 9.1: Complexity comparison of attention variants for sequence length n , model dimension d , window size w , and rank r . Accuracy percentages are relative to full attention on language modeling tasks.

Variant	Time	Memory	Accuracy	Max Length	Use Case
Full Attention	$O(n^2d)$	$O(n^2)$	100%	512-1024	Standard tasks
Local Attention	$O(nwd)$	$O(nw)$	98-99%	4096-8192	Document processing
Sparse Attention	$O(n\sqrt{nd})$	$O(n\sqrt{n})$	97-99%	8192-16384	Long documents
Linear Attention	$O(nd^2)$	$O(nd)$	95-98%	16384+	Extreme length
Low-Rank Attention	$O(nrd)$	$O(nr)$	96-98%	4096-8192	Compression

To make these complexity bounds concrete, consider processing sequences of varying lengths with BERT-base configuration ($d = 768$, 12 heads, $d_k = 64$ per head). The following table shows actual memory requirements for attention matrices across different sequence lengths and attention variants.

Table 9.2: Memory requirements (MB) for attention matrices with 12 heads, batch size 1, FP32 precision. Window size $w = 256$, rank $r = 256$ for applicable variants.

Variant	n=512	n=4096	n=8192	n=16384
Full Attention	12.6 MB	805 MB	3.2 GB	12.9 GB
Local Attention ($w = 256$)	6.3 MB	50.3 MB	101 MB	201 MB
Sparse Attention (\sqrt{n})	1.1 MB	12.6 MB	35.7 MB	101 MB
Linear Attention	0.3 MB	2.3 MB	4.7 MB	9.4 MB
Low-Rank ($r = 256$)	6.3 MB	50.3 MB	101 MB	201 MB

The memory savings become dramatic for long sequences. At 16,384 tokens, full attention requires 12.9 GB per example—impossible to fit on most GPUs even with batch size 1. Local attention reduces this to 201 MB, enabling batch size 32 on a 40 GB A100 GPU. Linear attention requires only 9.4 MB, enabling batch sizes of several hundred even for very long sequences.

The computational cost comparison is equally striking. For a sequence of 8192 tokens with $d = 768$ and 12 heads, full attention requires approximately 48.3 billion floating-point operations (FLOPs) per layer. Local attention with window 256 reduces this to 3.0 billion FLOPs (16x speedup), sparse attention to 6.0 billion FLOPs (8x speedup), linear attention to 4.5 billion FLOPs (10.7x speedup), and low-rank attention to 3.0 billion FLOPs (16x speedup). On an NVIDIA A100 GPU with 312 TFLOPS of FP16 throughput, full attention takes approximately 0.15 ms per layer, while efficient variants take 10-20 microseconds—enabling much faster inference and training.

The accuracy trade-offs vary by task and sequence length. For sequences up to 2048 tokens, local attention with window 512 typically matches full attention within 0.5% on language modeling perplexity. Sparse attention with well-designed patterns achieves similar accuracy. Linear attention shows 2-3% degradation, while low-rank attention with rank 256 shows 1-2% degradation. For longer sequences exceeding 4096 tokens, the accuracy gaps widen slightly, but efficient variants remain highly competitive. Importantly, the accuracy loss is often task-dependent—some tasks like document classification are more tolerant of approximate attention than tasks like machine translation or question answering that require precise alignment.

9.9.6 Implementation Considerations

Implementing efficient attention variants requires careful consideration of hardware characteristics, numerical stability, and software frameworks. The theoretical complexity improvements do not always translate directly to wall-clock speedups due to GPU architecture constraints and implementation details.

Modern GPUs achieve peak performance on dense matrix multiplications with dimensions that are multiples of 16 or 32 (for tensor cores). Sparse attention patterns that result in irregular memory

access or non-aligned dimensions can suffer significant performance degradation. For example, a naive implementation of sparse attention with random sparsity patterns may achieve only 30-40% of the theoretical speedup due to poor memory coalescing and reduced arithmetic intensity. Block-sparse patterns that operate on 16x16 or 32x32 blocks achieve much better GPU utilization, typically reaching 60-80% of theoretical speedup.

Memory bandwidth is often the limiting factor for attention computation, particularly for efficient variants. The attention mechanism is memory-bound rather than compute-bound for typical sequence lengths—the GPU spends more time loading data from memory than performing arithmetic operations. This means that reducing the number of operations (FLOPs) does not always proportionally reduce runtime. Efficient implementations must minimize memory transfers through kernel fusion, where multiple operations are combined into a single GPU kernel that keeps intermediate results in fast on-chip memory. FlashAttention demonstrates this principle by fusing the attention computation (\mathbf{QK}^\top , softmax, multiply by \mathbf{V}) into a single kernel that never materializes the full attention matrix in global memory, achieving 2-4x speedup over standard implementations even for full attention.

Numerical stability is a critical concern for efficient attention variants. The softmax operation in attention is numerically sensitive—subtracting the maximum value before exponentiation is essential to prevent overflow. Linear attention approximations must carefully handle the feature map computation to avoid numerical issues. The Performer’s random Fourier features require computing exponentials of potentially large values, necessitating careful scaling and normalization. Low-rank attention must ensure that the projection matrices are well-conditioned to avoid amplifying numerical errors.

Framework support for efficient attention varies significantly. PyTorch and TensorFlow provide optimized implementations of standard attention through `torch.nn.MultiheadAttention` and `tf.keras.layers.MultiHeadAttention`, but efficient variants often require custom implementations. The xFormers library provides optimized implementations of several efficient attention variants, including memory-efficient attention and block-sparse attention. The Triton framework enables writing custom GPU kernels in Python that achieve performance comparable to hand-written CUDA, making it easier to implement and experiment with novel attention patterns. For production deployment, specialized libraries like FasterTransformer and TensorRT provide highly optimized implementations of common attention variants with automatic kernel selection based on input dimensions and hardware capabilities.

9.10 Exercises

Exercise 9.1. Implement cross-attention layer in PyTorch. Test with encoder output (length 10, dim 128) and decoder input (length 7, dim 128). Verify attention matrix shape is 7×10 .

Exercise 9.2. Calculate the memory requirements for attention matrices in a BERT-base model (12 heads, $d_{\text{model}} = 768$) processing sequences of length 512, 2048, and 4096 tokens. Compare full attention, local attention with window size 256, and linear attention. How much memory is saved at each sequence length?

Exercise 9.3. Implement local attention with window size $w = 128$ for a sequence of length 1024. Compare the computational cost (FLOPs) and memory usage to full attention. Measure actual runtime on GPU and explain any discrepancy between theoretical and observed speedup.

Exercise 9.4. Design a sparse attention pattern for document understanding that combines local attention (window 64), strided attention (stride 128), and global attention to the first token. Calculate the number of attention connections per token and total memory requirements for a 4096-token sequence. What percentage of full attention’s connections does this pattern use?

Exercise 9.5. Implement linear attention using the feature map $\phi(\mathbf{x}) = \text{elu}(\mathbf{x}) + 1$. Compare attention patterns to standard softmax attention on a sample sequence. Measure the approximation error and identify cases where linear attention diverges most from full attention.

Exercise 9.6. For a transformer with 24 layers processing 8192-token sequences, calculate the total memory required for attention matrices using: (1) full attention, (2) local attention with window 512, (3) sparse attention with \sqrt{n} connections per token, (4) linear attention, and (5) low-rank attention with rank 256. Assume 12 heads, $d_{\text{model}} = 1024$, batch size 8, and FP16 precision.

Exercise 9.7. Implement relative position bias as in T5. Use buckets: [0, 1, 2, 3, 4, 5-7, 8-15, 16-31, 32+]. Show how attention scores change with relative distance and compare to absolute position encodings.

Exercise 9.8. Analyze the trade-off between window size and accuracy for local attention. Train a small transformer on a language modeling task with window sizes [64, 128, 256, 512, full]. Plot perplexity vs window size and identify the point of diminishing returns. How does this relate to the average dependency length in the dataset?

Exercise 9.9. Create visualization showing: (1) Self-attention patterns for sentence "The quick brown fox jumps", (2) Effect of causal masking, (3) Difference between heads 1 and 12 in multi-head attention. What patterns emerge?

Exercise 9.10. Compare computational cost of: (1) Additive (Bahdanau) attention, (2) Multiplicative attention, (3) Scaled dot-product attention. For $n = 512$, $d_k = 64$, which is most efficient? How does the ranking change for $n = 4096$?

Part IV

Transformer Architecture

Chapter 10

The Transformer Model

Chapter Overview

The Transformer architecture, introduced in "Attention is All You Need" (Vaswani et al., 2017), revolutionized deep learning by replacing recurrence with pure attention mechanisms. This chapter presents the complete transformer architecture, combining all attention mechanisms from previous chapters into a powerful encoder-decoder model.

We develop the transformer from bottom to top: starting with the attention layer, building encoder and decoder blocks, and assembling the full architecture. We provide complete mathematical specifications, dimension tracking, and parameter counts for standard transformer configurations.

Learning Objectives

1. Understand the complete transformer encoder-decoder architecture
2. Implement position-wise feed-forward networks
3. Apply layer normalization and residual connections
4. Compute output dimensions through the entire network
5. Count parameters for transformer models (BERT-base, GPT-2)
6. Understand training objectives for different transformer variants

10.1 Transformer Architecture Overview

10.1.1 High-Level Structure

The transformer architecture represents a fundamental departure from the recurrent and convolutional architectures that dominated sequence modeling before 2017. At its core, the transformer is an encoder-decoder architecture that processes sequences entirely through attention mechanisms, eliminating the sequential dependencies that made RNNs difficult to parallelize. The encoder processes the input sequence and produces contextualized representations where each position has attended to all other positions in the input. The decoder then generates the output sequence autoregressively, attending both to its own previously generated tokens and to the encoder's output through a cross-attention mechanism. This design enables the model to capture long-range dependencies without the vanishing gradient problems that plague recurrent architectures, while simultaneously allowing massive parallelization during training.

The key innovation that makes transformers practical is the elimination of recurrence in favor of pure attention mechanisms. In an RNN, processing a sequence of length n requires n sequential steps, each depending on the previous hidden state. This sequential dependency means that even with unlimited computational resources, the time complexity remains $O(n)$ because operations cannot be parallelized across time steps. The transformer, by contrast, computes attention between all pairs

of positions simultaneously, requiring only $O(1)$ sequential operations regardless of sequence length. For a sequence of length 512, this means the difference between 512 sequential steps (RNN) and a single parallel operation (transformer). On modern GPUs with thousands of cores, this parallelization advantage translates to training speedups of 10-100 \times compared to recurrent architectures.

The transformer achieves this parallelization through multi-head self-attention, which allows each position to attend to all positions in a single operation. For an input sequence $\mathbf{X} \in \mathbb{R}^{n \times d_{\text{model}}}$, the self-attention mechanism computes attention scores between all n^2 pairs of positions simultaneously, producing an output of the same shape $\mathbb{R}^{n \times d_{\text{model}}}$. This operation is entirely parallelizable across both the batch dimension and the sequence dimension, making it ideally suited for GPU acceleration. The multi-head aspect further enhances expressiveness by allowing the model to attend to different representation subspaces simultaneously—one head might capture syntactic relationships while another captures semantic similarity.

However, pure attention mechanisms lack an inherent notion of sequence order. Unlike RNNs where position information is implicit in the sequential processing, transformers must explicitly encode positional information. This is achieved through positional encodings that are added to the input embeddings, providing each position with a unique signature that the attention mechanism can use to distinguish positions. The original transformer uses sinusoidal positional encodings, though learned positional embeddings have also proven effective. This explicit position encoding is crucial: without it, the transformer would be permutation-invariant, treating "the cat sat" identically to "sat cat the."

The transformer architecture also incorporates residual connections and layer normalization at every sub-layer, forming the pattern $\text{LayerNorm}(x + \text{Sublayer}(x))$ throughout the network. These residual connections serve multiple purposes: they provide direct gradient pathways that enable training of very deep networks (the original transformer uses 6 layers, but modern variants scale to 96 layers in GPT-3), they allow the model to learn incremental refinements rather than complete transformations at each layer, and they stabilize training by preventing the exploding or vanishing gradient problems that can occur in deep networks. Layer normalization, applied after each residual connection, normalizes activations across the feature dimension, ensuring stable activation distributions throughout the network regardless of batch size.

The position-wise feed-forward network, applied after each attention layer, provides additional representational capacity through a simple two-layer network with a ReLU or GELU activation. This network is applied independently to each position, meaning it doesn't mix information across positions (unlike attention). The feed-forward network typically expands the representation to a higher dimension (usually $4 \times d_{\text{model}}$) before projecting back down, creating a bottleneck architecture that encourages the model to learn compressed representations. For BERT-base with $d_{\text{model}} = 768$, the feed-forward network expands to $d_{ff} = 3072$ dimensions, and this expansion-projection accounts for approximately two-thirds of the parameters in each transformer layer.

Key Point 10.1. *Transformers achieve $O(1)$ sequential operations compared to $O(n)$ for RNNs, enabling massive parallelization during training. For a sequence of length 512 on a GPU with 10,000 cores, this means the difference between 512 sequential steps and a single parallel operation, yielding training speedups of 10-100 \times in practice. This parallelization advantage is the primary reason transformers have replaced RNNs as the dominant architecture for sequence modeling.*

10.2 Transformer Encoder

10.2.1 Single Encoder Layer

A transformer encoder layer consists of two main sub-layers: multi-head self-attention followed by a position-wise feed-forward network, with residual connections and layer normalization applied around each sub-layer. This architecture enables the encoder to build increasingly sophisticated representations of the input sequence as information flows through multiple layers. The self-attention mechanism allows each position to gather information from all other positions, creating contextualized representations

where the meaning of each token depends on its surrounding context. The feed-forward network then processes each position independently, applying a non-linear transformation that enhances the model's representational capacity.

The residual connections are crucial for enabling gradient flow through deep networks. Without them, gradients would need to flow through multiple attention and feed-forward layers, potentially vanishing or exploding. With residual connections, gradients have a direct path from the output back to the input of each layer, ensuring stable training even for very deep transformers. The layer normalization, applied after adding the residual, normalizes the activations across the feature dimension, maintaining stable activation distributions throughout the network. This combination of residual connections and layer normalization is what enables transformers to scale to dozens or even hundreds of layers.

Definition 10.1 (Transformer Encoder Layer). An encoder layer applies multi-head self-attention followed by feed-forward network, with residual connections and layer normalization. For input $\mathbf{X} \in \mathbb{R}^{B \times n \times d_{\text{model}}}$ where B is batch size, n is sequence length, and d_{model} is model dimension:

Step 1: Multi-Head Self-Attention

$$\mathbf{h}^{(1)} = \text{LayerNorm}(\mathbf{X} + \text{MultiHeadAttn}(\mathbf{X}, \mathbf{X}, \mathbf{X})) \quad (10.1)$$

where the output maintains shape $\mathbb{R}^{B \times n \times d_{\text{model}}}$.

Step 2: Position-wise Feed-Forward

$$\mathbf{h}^{(2)} = \text{LayerNorm}(\mathbf{h}^{(1)} + \text{FFN}(\mathbf{h}^{(1)})) \quad (10.2)$$

where the output again maintains shape $\mathbb{R}^{B \times n \times d_{\text{model}}}$.

The feed-forward network is defined as:

$$\text{FFN}(\mathbf{x}) = \mathbf{W}_2 \cdot \text{ReLU}(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2 \quad (10.3)$$

with $\mathbf{W}_1 \in \mathbb{R}^{d_{\text{model}} \times d_{\text{ff}}}$, $\mathbf{W}_2 \in \mathbb{R}^{d_{\text{ff}} \times d_{\text{model}}}$, and typically $d_{\text{ff}} = 4 \times d_{\text{model}}$.

The dimension tracking through an encoder layer reveals important properties about memory consumption and computational cost. The input $\mathbf{X} \in \mathbb{R}^{B \times n \times d_{\text{model}}}$ is first projected to queries, keys, and values, each with shape $\mathbb{R}^{B \times n \times d_{\text{model}}}$. For multi-head attention with h heads, these are reshaped to $\mathbb{R}^{B \times h \times n \times d_k}$ where $d_k = d_{\text{model}}/h$. The attention scores form a matrix $\mathbb{R}^{B \times h \times n \times n}$, and this quadratic term in sequence length is what dominates memory consumption for long sequences. After attention, the output is projected back to $\mathbb{R}^{B \times n \times d_{\text{model}}}$, added to the residual, and normalized.

The feed-forward network then expands each position's representation from d_{model} to d_{ff} dimensions before projecting back down. For BERT-base with $d_{\text{model}} = 768$ and $d_{\text{ff}} = 3072$, this means each position's representation temporarily expands to $4 \times$ its original size. This expansion creates a bottleneck that forces the model to learn compressed representations, similar to the hidden layer in an autoencoder. The intermediate activations $\mathbb{R}^{B \times n \times d_{\text{ff}}}$ consume significant memory during training—for batch size 32 and sequence length 512, this amounts to $32 \times 512 \times 3072 \times 4 = 201$ MB per layer in FP32, and with 12 layers in BERT-base, the feed-forward activations alone consume 2.4 GB of GPU memory.

Example 10.1 (BERT-base Encoder Layer). BERT-base uses the following configuration, which has become a standard baseline for many transformer models:

- Model dimension: $d_{\text{model}} = 768$
- Attention heads: $h = 12$, so $d_k = d_v = 768/12 = 64$ per head
- Feed-forward dimension: $d_{\text{ff}} = 3072$ (exactly $4 \times d_{\text{model}}$)

- Sequence length: $n = 512$ (maximum)
- Batch size: $B = 32$ (typical for training)

Dimension tracking through the layer:

Input: $\mathbf{X} \in \mathbb{R}^{32 \times 512 \times 768}$ (batch \times sequence \times model dimension)

Multi-Head Attention:

$$\text{Q, K, V projections: } \mathbb{R}^{32 \times 512 \times 768} \rightarrow \mathbb{R}^{32 \times 512 \times 768} \quad (10.4)$$

$$\text{Reshape for heads: } \mathbb{R}^{32 \times 512 \times 768} \rightarrow \mathbb{R}^{32 \times 12 \times 512 \times 64} \quad (10.5)$$

$$\text{Attention scores: } \mathbb{R}^{32 \times 12 \times 512 \times 512} \quad (\text{quadratic in } n!) \quad (10.6)$$

$$\text{Attention output: } \mathbb{R}^{32 \times 12 \times 512 \times 64} \quad (10.7)$$

$$\text{Concatenate heads: } \mathbb{R}^{32 \times 512 \times 768} \quad (10.8)$$

$$\text{Output projection: } \mathbb{R}^{32 \times 512 \times 768} \quad (10.9)$$

The attention scores matrix $\mathbb{R}^{32 \times 12 \times 512 \times 512}$ requires $32 \times 12 \times 512 \times 512 \times 4 = 402$ MB in FP32. This quadratic scaling means that doubling the sequence length to 1024 would require 1.6 GB just for attention scores in a single layer.

Feed-Forward Network:

$$\text{First projection: } \mathbb{R}^{32 \times 512 \times 768} \xrightarrow{\mathbf{W}_1} \mathbb{R}^{32 \times 512 \times 3072} \quad (10.10)$$

$$\text{ReLU activation: } \mathbb{R}^{32 \times 512 \times 3072} \rightarrow \mathbb{R}^{32 \times 512 \times 3072} \quad (10.11)$$

$$\text{Second projection: } \mathbb{R}^{32 \times 512 \times 3072} \xrightarrow{\mathbf{W}_2} \mathbb{R}^{32 \times 512 \times 768} \quad (10.12)$$

The intermediate activations $\mathbb{R}^{32 \times 512 \times 3072}$ require $32 \times 512 \times 3072 \times 4 = 201$ MB in FP32.

Parameter count breakdown:

$$\text{Multi-head attention: } 4 \times 768^2 = 2,359,296 \quad (\text{Q, K, V, O projections}) \quad (10.13)$$

$$\text{Feed-forward network: } 768 \times 3072 + 3072 + 3072 \times 768 + 768 \quad (10.14)$$

$$= 2,359,296 + 3,072 + 2,359,296 + 768 \quad (10.15)$$

$$= 4,722,432 \quad (10.16)$$

$$\text{Layer normalization (2}\times\text{): } 2 \times 2 \times 768 = 3,072 \quad (\text{scale } \gamma \text{ and shift } \beta) \quad (10.17)$$

Total per encoder layer: $2,359,296 + 4,722,432 + 3,072 = 7,084,800$ parameters

This reveals that the feed-forward network contains approximately twice as many parameters as the attention mechanism (4.7M vs 2.4M), despite attention being conceptually more complex. This is because the feed-forward network's expansion to $4 \times d_{\text{model}}$ dimensions creates two large weight matrices, while attention's parameters are distributed across four projections of size $d_{\text{model}} \times d_{\text{model}}$.

Memory requirements during training:

$$\text{Parameters (FP32): } 7,084,800 \times 4 = 28.3 \text{ MB} \quad (10.18)$$

$$\text{Gradients (FP32): } 7,084,800 \times 4 = 28.3 \text{ MB} \quad (10.19)$$

$$\text{Adam optimizer states: } 7,084,800 \times 8 = 56.7 \text{ MB} \quad (10.20)$$

$$\text{Attention scores: } 402 \text{ MB} \quad (10.21)$$

$$\text{FFN intermediate: } 201 \text{ MB} \quad (10.22)$$

$$\text{Total per layer: } \approx 716 \text{ MB} \quad (10.23)$$

For BERT-base with 12 encoder layers, this amounts to approximately 8.6 GB just for the encoder layers, not including embeddings or other activations. This explains why training BERT-base requires GPUs with at least 16 GB of memory.

10.2.2 Complete Encoder Stack

The complete transformer encoder stacks N identical encoder layers, with each layer's output serving as input to the next layer. This stacking enables the model to build increasingly abstract representations: early layers might capture local syntactic patterns, middle layers might identify semantic relationships, and later layers might encode task-specific features. The depth of the network is crucial for performance—BERT-base uses 12 layers, BERT-large uses 24 layers, and GPT-3 uses 96 layers. However, deeper networks require more careful optimization, including learning rate warmup, gradient clipping, and appropriate weight initialization.

Definition 10.2 (Transformer Encoder). Stack N encoder layers, with input embeddings and positional encodings added at the bottom:

$$\mathbf{X}^{(0)} = \text{Embedding}(\text{input}) + \text{PositionalEncoding} \quad (10.24)$$

where $\text{Embedding} \in \mathbb{R}^{V \times d_{\text{model}}}$ maps vocabulary indices to dense vectors, and $\text{PositionalEncoding} \in \mathbb{R}^{n_{\text{max}} \times d_{\text{model}}}$ provides position information.

Then apply N encoder layers sequentially:

$$\mathbf{X}^{(\ell)} = \text{EncoderLayer}^{(\ell)}(\mathbf{X}^{(\ell-1)}) \quad \text{for } \ell = 1, \dots, N \quad (10.25)$$

The final encoder output $\mathbf{X}^{(N)} \in \mathbb{R}^{B \times n \times d_{\text{model}}}$ contains contextualized representations of the input sequence.

The sequential application of encoder layers means that information flows through N attention operations, allowing each token to indirectly attend to all other tokens through multiple hops. In a 12-layer encoder, information can propagate across the entire sequence through 12 levels of attention, enabling the model to capture very long-range dependencies. However, this sequential stacking also means that encoder layers cannot be parallelized—layer ℓ must wait for layer $\ell - 1$ to complete. The parallelization in transformers occurs within each layer (across batch and sequence dimensions), not across layers.

Example 10.2 (BERT-base Complete Encoder). BERT-base represents the standard configuration that has been widely adopted and serves as a baseline for many NLP tasks. The architecture is:

- Layers: $N = 12$
- Model dimension: $d_{\text{model}} = 768$
- Attention heads: $h = 12$
- Feed-forward dimension: $d_{ff} = 3072$
- Vocabulary size: $V = 30,000$ (WordPiece tokenization)
- Maximum sequence length: $n_{\text{max}} = 512$

Complete parameter count breakdown:

$$\text{Token embeddings: } 30,000 \times 768 = 23,040,000 \quad (10.26)$$

$$\text{Position embeddings: } 512 \times 768 = 393,216 \quad (10.27)$$

$$\text{Token type embeddings: } 2 \times 768 = 1,536 \quad (\text{for segment A/B}) \quad (10.28)$$

$$\text{Embedding layer norm: } 2 \times 768 = 1,536 \quad (10.29)$$

$$12 \text{ encoder layers: } 12 \times 7,084,800 = 85,017,600 \quad (10.30)$$

$$\text{Pooler (for classification): } 768 \times 768 + 768 = 590,592 \quad (10.31)$$

$$\text{Total: } 109,044,480 \approx \mathbf{110M \text{ parameters}} \quad (10.32)$$

This matches the reported BERT-base size of 110M parameters. Notice that the embeddings account for approximately 21% of the total parameters (23M out of 110M), while the transformer layers account for 78%. This ratio changes dramatically for larger vocabularies—models with 50,000 token vocabularies would have embeddings consuming 35% of parameters, motivating techniques like vocabulary pruning or shared embeddings.

Memory requirements for training (batch size 32, sequence length 512):

$$\text{Parameters (FP32): } 110,000,000 \times 4 = 440 \text{ MB} \quad (10.33)$$

$$\text{Gradients (FP32): } 110,000,000 \times 4 = 440 \text{ MB} \quad (10.34)$$

$$\text{Adam optimizer states: } 110,000,000 \times 8 = 880 \text{ MB} \quad (10.35)$$

$$\text{Activations (estimated): } \approx 12 \text{ GB} \quad (10.36)$$

$$\text{Total: } \approx 13.8 \text{ GB} \quad (10.37)$$

The activation memory dominates, consuming approximately 87% of total memory. This is why techniques like gradient checkpointing (recomputing activations during backward pass instead of storing them) can reduce memory consumption by 50-70% at the cost of 20-30% slower training.

Training throughput on NVIDIA A100 GPU:

The A100 provides 312 TFLOPS of FP16 compute with Tensor Cores. For BERT-base, a single forward pass with batch size 32 and sequence length 512 requires approximately:

$$\text{FLOPs per layer: } 24nd_{\text{model}}^2 + 4n^2d_{\text{model}} \quad (10.38)$$

$$= 24 \times 512 \times 768^2 + 4 \times 512^2 \times 768 \quad (10.39)$$

$$= 7.26 \text{ GFLOPs} \quad (10.40)$$

$$\text{Total for 12 layers: } 12 \times 7.26 = 87.1 \text{ GFLOPs} \quad (10.41)$$

$$\text{With embeddings and overhead: } \approx 100 \text{ GFLOPs} \quad (10.42)$$

At 312 TFLOPS, this suggests a forward pass should take $100/312,000 = 0.32$ milliseconds. In practice, memory bandwidth limitations and kernel launch overhead mean actual forward pass time is approximately 5-10 milliseconds, achieving 10-20% of peak FLOPS. With backward pass taking approximately $2\times$ forward pass time, a complete training step takes 15-30 milliseconds, yielding throughput of 30-60 training steps per second, or approximately 500,000-1,000,000 tokens per second.

10.3 Position-wise Feed-Forward Networks

The position-wise feed-forward network represents the second major component of each transformer layer, complementing the attention mechanism with additional non-linear transformations. While attention allows positions to exchange information and build contextualized representations, the feed-forward network processes each position independently, applying the same learned transformation to every position in the sequence. This independence is what makes it "position-wise"—the network applied to position i is identical to the network applied to position j , with no parameter sharing or information

flow between positions.

The feed-forward network consists of two linear transformations with a non-linear activation function in between, forming a simple two-layer neural network. The first layer expands the representation from d_{model} dimensions to a larger dimension d_{ff} (typically $4 \times d_{\text{model}}$), applies an activation function, and then the second layer projects back down to d_{model} dimensions. This expansion-and-contraction creates a bottleneck architecture similar to an autoencoder, forcing the model to learn compressed representations that capture the most important features. The expansion factor of $4\times$ is a design choice from the original transformer paper that has been widely adopted, though some recent models experiment with different ratios.

Definition 10.3 (Position-wise FFN). For input $\mathbf{X} \in \mathbb{R}^{B \times n \times d_{\text{model}}}$, apply the same two-layer network independently to each position:

$$\text{FFN}(\mathbf{x}) = \max(0, \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2 \quad (10.43)$$

where $\mathbf{W}_1 \in \mathbb{R}^{d_{\text{model}} \times d_{ff}}$, $\mathbf{b}_1 \in \mathbb{R}^{d_{ff}}$, $\mathbf{W}_2 \in \mathbb{R}^{d_{ff} \times d_{\text{model}}}$, and $\mathbf{b}_2 \in \mathbb{R}^{d_{\text{model}}}$.

For a sequence $\mathbf{X} \in \mathbb{R}^{B \times n \times d_{\text{model}}}$, apply to each position independently:

$$\text{FFN}(\mathbf{X})_{i,:} = \text{FFN}(\mathbf{X}_{i,:}) \quad \text{for } i = 1, \dots, n \quad (10.44)$$

The output maintains the same shape as the input: $\mathbb{R}^{B \times n \times d_{\text{model}}}$.

The term "position-wise" emphasizes a crucial distinction from the attention mechanism. In attention, every position attends to every other position, creating $O(n^2)$ interactions. In the feed-forward network, each position is processed completely independently, creating only $O(n)$ operations. This means the feed-forward network is embarrassingly parallel—all n positions can be processed simultaneously with no dependencies. In practice, this is implemented as a single matrix multiplication: the input $\mathbf{X} \in \mathbb{R}^{B \times n \times d_{\text{model}}}$ is reshaped to $\mathbb{R}^{Bn \times d_{\text{model}}}$, multiplied by \mathbf{W}_1 , activated, multiplied by \mathbf{W}_2 , and reshaped back to $\mathbb{R}^{B \times n \times d_{\text{model}}}$.

The choice of activation function significantly impacts model performance and training dynamics. The original transformer used ReLU activation, which is simple and computationally efficient but can suffer from "dying ReLU" problems where neurons become permanently inactive. BERT and GPT introduced the GELU (Gaussian Error Linear Unit) activation, which provides a smoother, probabilistic alternative to ReLU. GELU is defined as $\text{GELU}(x) = x \cdot \Phi(x)$ where $\Phi(x)$ is the cumulative distribution function of the standard normal distribution. In practice, GELU is approximated as $\text{GELU}(x) \approx 0.5x(1 + \tanh[\sqrt{2/\pi}(x + 0.044715x^3)])$. Empirically, GELU tends to provide slightly better performance than ReLU for transformer models, though the difference is often small.

The feed-forward network accounts for a substantial portion of the model's parameters and computational cost. For BERT-base with $d_{\text{model}} = 768$ and $d_{ff} = 3072$, each feed-forward network contains $768 \times 3072 + 3072 \times 768 = 4.7\text{M}$ parameters, compared to $4 \times 768^2 = 2.4\text{M}$ parameters in the attention mechanism. This means approximately two-thirds of each layer's parameters are in the feed-forward network. Similarly, for short sequences where $n < d_{\text{model}}$, the feed-forward network dominates computational cost. For BERT-base with sequence length 512, the feed-forward network requires $2 \times 512 \times 768 \times 3072 = 2.4$ GFLOPs per layer, while attention requires $8 \times 512 \times 768^2 + 4 \times 512^2 \times 768 = 3.2$ GFLOPs. The crossover point occurs around $n = 2d_{\text{model}}$ —for longer sequences, attention dominates; for shorter sequences, the feed-forward network dominates.

Example 10.3 (Feed-Forward Network Dimensions and Memory). For BERT-base with $d_{\text{model}} = 768$, $d_{ff} = 3072$, batch size $B = 32$, and sequence length $n = 512$:

Dimension tracking:

$$\text{Input: } \mathbf{X} \in \mathbb{R}^{32 \times 512 \times 768} \quad (10.45)$$

$$\text{First projection: } \mathbf{X}\mathbf{W}_1 + \mathbf{b}_1 \in \mathbb{R}^{32 \times 512 \times 3072} \quad (10.46)$$

$$\text{After ReLU/GELU: } \mathbb{R}^{32 \times 512 \times 3072} \quad (10.47)$$

$$\text{Second projection: } \mathbf{X}\mathbf{W}_2 + \mathbf{b}_2 \in \mathbb{R}^{32 \times 512 \times 768} \quad (10.48)$$

$$\text{Output: } \mathbb{R}^{32 \times 512 \times 768} \quad (10.49)$$

Memory requirements:

$$\text{Input activations: } 32 \times 512 \times 768 \times 4 = 50.3 \text{ MB (FP32)} \quad (10.50)$$

$$\text{Intermediate activations: } 32 \times 512 \times 3072 \times 4 = 201.3 \text{ MB (FP32)} \quad (10.51)$$

$$\text{Output activations: } 32 \times 512 \times 768 \times 4 = 50.3 \text{ MB (FP32)} \quad (10.52)$$

$$\text{Parameters } (\mathbf{W}_1, \mathbf{W}_2): (768 \times 3072 + 3072 \times 768) \times 4 = 18.9 \text{ MB (FP32)} \quad (10.53)$$

The intermediate activations at dimension $d_{ff} = 3072$ consume $4\times$ the memory of the input/output activations at dimension $d_{\text{model}} = 768$. For a 12-layer BERT model, the feed-forward intermediate activations across all layers consume $12 \times 201.3 = 2.4$ GB of memory during training. This is why gradient checkpointing, which recomputes these activations during the backward pass instead of storing them, can significantly reduce memory consumption.

Computational cost:

$$\text{First projection: } Bn \times d_{\text{model}} \times d_{ff} = 32 \times 512 \times 768 \times 3072 = 38.7 \text{ GFLOPs} \quad (10.54)$$

$$\text{Second projection: } Bn \times d_{ff} \times d_{\text{model}} = 32 \times 512 \times 3072 \times 768 = 38.7 \text{ GFLOPs} \quad (10.55)$$

$$\text{Total: } 77.4 \text{ GFLOPs per layer} \quad (10.56)$$

For comparison, the attention mechanism in the same layer requires approximately 51.5 GFLOPs (including Q, K, V projections, attention computation, and output projection). This means the feed-forward network accounts for 60% of the computational cost per layer for this configuration.

Alternative activation functions: While ReLU and GELU are most common, other activation functions have been explored for transformers. The Swish activation $\text{Swish}(x) = x \cdot \sigma(\beta x)$ where σ is the sigmoid function, provides similar properties to GELU. The GLU (Gated Linear Unit) family, including $\text{GLU}(x) = (x\mathbf{W}_1) \odot \sigma(x\mathbf{W}_2)$, uses gating mechanisms similar to LSTMs. Recent work has also explored learned activation functions that adapt during training. However, GELU remains the most widely adopted choice for modern transformers due to its balance of performance and computational efficiency.

10.4 Transformer Decoder

10.4.1 Single Decoder Layer

The transformer decoder extends the encoder architecture with an additional cross-attention mechanism that allows the decoder to attend to the encoder's output. While the encoder uses only self-attention to build contextualized representations of the input, the decoder must perform three distinct operations: masked self-attention on the target sequence, cross-attention to the source sequence, and position-wise feed-forward transformation. This three-sublayer structure enables the decoder to generate output sequences that are conditioned on both the previously generated tokens and the encoded input sequence.

The masked self-attention in the decoder is crucial for maintaining the autoregressive property during training. Unlike the encoder's bidirectional self-attention where each position can attend to all positions, the decoder's self-attention must be causal—position i can only attend to positions $j \leq i$. This masking ensures that the model cannot "cheat" by looking at future tokens during training. Without this mask, the model could simply copy the target sequence during training without learning

to generate it. The mask is implemented by setting attention scores for future positions to $-\infty$ before the softmax, ensuring they receive zero attention weight.

The cross-attention mechanism is where the decoder actually uses information from the encoder. In cross-attention, the queries come from the decoder's hidden states (representing "what information do I need?"), while the keys and values come from the encoder's output (representing "what information is available from the source?"). This asymmetry allows the decoder to selectively focus on relevant parts of the source sequence when generating each target token. For machine translation, this might mean attending to the source word being translated; for summarization, it might mean attending to the most salient sentences in the document.

Definition 10.4 (Transformer Decoder Layer). A decoder layer has three sub-layers, each with residual connections and layer normalization. For input $\mathbf{Y} \in \mathbb{R}^{B \times m \times d_{\text{model}}}$ (target sequence) and encoder output $\mathbf{X}_{\text{enc}} \in \mathbb{R}^{B \times n \times d_{\text{model}}}$ (source sequence):

Step 1: Masked Self-Attention

$$\mathbf{h}^{(1)} = \text{LayerNorm}(\mathbf{Y} + \text{MaskedMultiHeadAttn}(\mathbf{Y}, \mathbf{Y}, \mathbf{Y})) \quad (10.57)$$

where the attention mask prevents position i from attending to positions $j > i$.

Step 2: Cross-Attention to Encoder

$$\mathbf{h}^{(2)} = \text{LayerNorm}(\mathbf{h}^{(1)} + \text{MultiHeadAttn}(\mathbf{h}^{(1)}, \mathbf{X}_{\text{enc}}, \mathbf{X}_{\text{enc}})) \quad (10.58)$$

where queries come from $\mathbf{h}^{(1)}$ and keys/values come from \mathbf{X}_{enc} .

Step 3: Feed-Forward

$$\mathbf{h}^{(3)} = \text{LayerNorm}(\mathbf{h}^{(2)} + \text{FFN}(\mathbf{h}^{(2)})) \quad (10.59)$$

The output $\mathbf{h}^{(3)} \in \mathbb{R}^{B \times m \times d_{\text{model}}}$ maintains the target sequence length m .

The dimension compatibility in cross-attention deserves careful attention. The decoder hidden states $\mathbf{h}^{(1)} \in \mathbb{R}^{B \times m \times d_{\text{model}}}$ are projected to queries $\mathbf{Q} \in \mathbb{R}^{B \times m \times d_{\text{model}}}$, while the encoder output $\mathbf{X}_{\text{enc}} \in \mathbb{R}^{B \times n \times d_{\text{model}}}$ is projected to keys $\mathbf{K} \in \mathbb{R}^{B \times n \times d_{\text{model}}}$ and values $\mathbf{V} \in \mathbb{R}^{B \times n \times d_{\text{model}}}$. The attention scores are computed as $\mathbf{QK}^T \in \mathbb{R}^{B \times m \times n}$, creating a rectangular attention matrix where each of the m target positions attends to all n source positions. This is different from self-attention where the attention matrix is square ($n \times n$ for encoder, $m \times m$ for decoder self-attention).

The causal mask in decoder self-attention is implemented as a lower-triangular matrix. For a sequence of length $m = 5$, the mask looks like:

$$\text{Mask} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 & 1 \end{bmatrix} \quad (10.60)$$

where 1 indicates positions that can be attended to and 0 indicates positions that must be masked. In practice, the zeros are replaced with $-\infty$ before the softmax operation, ensuring masked positions receive zero attention weight. This mask is applied to the attention scores before softmax: $\text{softmax}(\mathbf{QK}^T / \sqrt{d_k} + \text{Mask})$.

Example 10.4 (Decoder Layer Dimension Tracking). For a translation task with source sequence length $n = 20$ (e.g., "The cat sat on the mat") and target sequence length $m = 15$ (e.g., "Le chat était assis"), using BERT-base dimensions ($d_{\text{model}} = 768$, $h = 12$, $d_{ff} = 3072$), batch size $B = 32$:

Inputs:

$$\text{Decoder input: } \mathbf{Y} \in \mathbb{R}^{32 \times 15 \times 768} \quad (10.61)$$

$$\text{Encoder output: } \mathbf{X}_{\text{enc}} \in \mathbb{R}^{32 \times 20 \times 768} \quad (10.62)$$

Masked Self-Attention:

$$\text{Q, K, V from } \mathbf{Y} : \mathbb{R}^{32 \times 15 \times 768} \quad (10.63)$$

$$\text{Attention scores: } \mathbb{R}^{32 \times 12 \times 15 \times 15} \quad (\text{square, causal masked}) \quad (10.64)$$

$$\text{Output: } \mathbb{R}^{32 \times 15 \times 768} \quad (10.65)$$

The attention scores matrix $\mathbb{R}^{32 \times 12 \times 15 \times 15}$ requires $32 \times 12 \times 15 \times 15 \times 4 = 3.5$ MB in FP32. This is much smaller than encoder self-attention because the target sequence is shorter than the source sequence in this example.

Cross-Attention:

$$\text{Q from } \mathbf{h}^{(1)} : \mathbb{R}^{32 \times 15 \times 768} \quad (10.66)$$

$$\text{K, V from } \mathbf{X}_{\text{enc}} : \mathbb{R}^{32 \times 20 \times 768} \quad (10.67)$$

$$\text{Attention scores: } \mathbb{R}^{32 \times 12 \times 15 \times 20} \quad (\text{rectangular!}) \quad (10.68)$$

$$\text{Output: } \mathbb{R}^{32 \times 15 \times 768} \quad (10.69)$$

The cross-attention scores $\mathbb{R}^{32 \times 12 \times 15 \times 20}$ require $32 \times 12 \times 15 \times 20 \times 4 = 4.6$ MB in FP32. Notice this is rectangular: 15 target positions attending to 20 source positions.

Feed-Forward Network:

$$\text{Input: } \mathbb{R}^{32 \times 15 \times 768} \quad (10.70)$$

$$\text{Intermediate: } \mathbb{R}^{32 \times 15 \times 3072} \quad (10.71)$$

$$\text{Output: } \mathbb{R}^{32 \times 15 \times 768} \quad (10.72)$$

The intermediate activations require $32 \times 15 \times 3072 \times 4 = 59.0$ MB in FP32.

10.4.2 Complete Decoder Stack

The complete decoder stacks N decoder layers, with each layer attending to both the previous decoder layer's output and the encoder's final output. This stacking enables the decoder to build increasingly sophisticated representations of the target sequence, conditioned on the source sequence. The encoder output \mathbf{X}_{enc} is reused by every decoder layer—it's computed once by the encoder and then fed into all N decoder layers. This means the encoder output must be stored in memory throughout the decoder's computation, contributing to memory requirements.

Definition 10.5 (Transformer Decoder). Stack N decoder layers, with target embeddings and positional encodings at the bottom:

$$\mathbf{Y}^{(0)} = \text{Embedding}(\text{target}) + \text{PositionalEncoding} \quad (10.73)$$

Then apply N decoder layers sequentially, each attending to the encoder output:

$$\mathbf{Y}^{(\ell)} = \text{DecoderLayer}^{(\ell)}(\mathbf{Y}^{(\ell-1)}, \mathbf{X}_{\text{enc}}) \quad \text{for } \ell = 1, \dots, N \quad (10.74)$$

The final decoder output $\mathbf{Y}^{(N)} \in \mathbb{R}^{B \times m \times d_{\text{model}}}$ is projected to vocabulary logits:

$$\text{logits} = \mathbf{Y}^{(N)} \mathbf{W}_{\text{out}} + \mathbf{b}_{\text{out}} \in \mathbb{R}^{B \times m \times V} \quad (10.75)$$

where $\mathbf{W}_{\text{out}} \in \mathbb{R}^{d_{\text{model}} \times V}$ and V is the vocabulary size.

During training, the entire target sequence is processed in parallel using teacher forcing—the model receives the ground-truth previous tokens rather than its own predictions. The causal mask ensures that position i cannot attend to future positions, maintaining the autoregressive property even though all positions are computed simultaneously. This parallel training is a major advantage over RNN decoders, which must process the target sequence sequentially even during training.

During inference, however, the decoder must generate tokens autoregressively, one at a time. At step t , the decoder has generated tokens y_1, \dots, y_{t-1} and must predict y_t . This requires running the decoder with input sequence length $t-1$, computing attention over all previously generated tokens. For a target sequence of length m , this requires m forward passes through the decoder, making inference much slower than training. This is why techniques like KV caching (storing computed key and value projections) are crucial for efficient inference.

Example 10.5 (Decoder Layer Parameter Count). For BERT-base dimensions ($d_{\text{model}} = 768$, $h = 12$, $d_{\text{ff}} = 3072$), a decoder layer contains:

Masked self-attention:

$$\text{Q, K, V, O projections: } 4 \times 768^2 = 2,359,296 \quad (10.76)$$

Cross-attention:

$$\text{Q, K, V, O projections: } 4 \times 768^2 = 2,359,296 \quad (10.77)$$

Feed-forward network:

$$\mathbf{W}_1, \mathbf{b}_1, \mathbf{W}_2, \mathbf{b}_2 : 768 \times 3072 + 3072 + 3072 \times 768 + 768 = 4,722,432 \quad (10.78)$$

Layer normalization (3 instances):

$$\text{Scale and shift parameters: } 3 \times 2 \times 768 = 4,608 \quad (10.79)$$

Total per decoder layer: $2,359,296 + 2,359,296 + 4,722,432 + 4,608 = 9,445,632$ parameters

This is approximately 33% more parameters than an encoder layer (9.4M vs 7.1M) due to the additional cross-attention mechanism. For a 6-layer decoder, this amounts to $6 \times 9,445,632 = 56.7\text{M}$ parameters, compared to $6 \times 7,084,800 = 42.5\text{M}$ for a 6-layer encoder.

Example 10.6 (Autoregressive Generation Memory). During autoregressive generation, the decoder must recompute attention over all previously generated tokens at each step. For a target sequence of length $m = 100$, generating the final token requires:

Without KV caching:

- Process sequence of length 100
- Compute Q, K, V for all 100 positions
- Compute attention scores $\mathbb{R}^{100 \times 100}$
- Total: 100 forward passes through decoder, each processing increasing sequence lengths

With KV caching:

- Store K, V from previous steps: $\mathbb{R}^{99 \times 768}$ per layer

- At step 100, compute only Q for new position: $\mathbb{R}^{1 \times 768}$
- Concatenate with cached K, V: $\mathbb{R}^{100 \times 768}$
- Compute attention scores $\mathbb{R}^{1 \times 100}$ (only for new position)
- Total: 100 forward passes, but each processes only 1 new position

For BERT-base dimensions with 12 decoder layers, the KV cache requires:

$$\text{Per layer: } 2 \times 100 \times 768 \times 4 = 614 \text{ KB (FP32)} \quad (10.80)$$

$$\text{All 12 layers: } 12 \times 614 = 7.4 \text{ MB} \quad (10.81)$$

This modest memory cost (7.4 MB for 100 tokens) enables approximately $50\times$ speedup in generation, reducing generation time from several seconds to tens of milliseconds for typical sequences.

10.5 Computational Complexity and Hardware Analysis

10.5.1 FLOPs Analysis

Understanding the computational complexity of transformers is essential for predicting training time, estimating hardware requirements, and identifying optimization opportunities. The transformer's computational cost is dominated by matrix multiplications in the attention mechanism and feed-forward network, with the relative importance depending on sequence length. For short sequences, the feed-forward network dominates; for long sequences, attention dominates due to its quadratic scaling.

For a single transformer encoder layer processing a sequence of length n with model dimension d_{model} , feed-forward dimension d_{ff} , and h attention heads (where $d_k = d_{\text{model}}/h$), the computational cost breaks down as follows. The multi-head attention requires four matrix multiplications for Q, K, V, and output projections, each costing $2nd_{\text{model}}^2$ FLOPs (the factor of 2 accounts for both multiplication and addition in matrix multiplication). The attention score computation \mathbf{QK}^T requires $2n^2d_{\text{model}}$ FLOPs, and the attention-weighted sum $\text{Attn}\mathbf{V}$ requires another $2n^2d_{\text{model}}$ FLOPs. The feed-forward network requires $2nd_{\text{model}}d_{\text{ff}}$ FLOPs for the first projection and another $2nd_{\text{model}}d_{\text{ff}}$ FLOPs for the second projection.

Definition 10.6 (Transformer Layer Computational Cost). For a single encoder layer with input $\mathbf{X} \in \mathbb{R}^{B \times n \times d_{\text{model}}}$:

Multi-head attention:

$$\text{Q, K, V projections: } 3 \times 2Bnd_{\text{model}}^2 = 6Bnd_{\text{model}}^2 \quad (10.82)$$

$$\text{Attention scores } (\mathbf{QK}^T): 2Bn^2d_{\text{model}} \quad (10.83)$$

$$\text{Attention-weighted sum: } 2Bn^2d_{\text{model}} \quad (10.84)$$

$$\text{Output projection: } 2Bnd_{\text{model}}^2 \quad (10.85)$$

$$\text{Total attention: } 8Bnd_{\text{model}}^2 + 4Bn^2d_{\text{model}} \quad (10.86)$$

Feed-forward network (assuming $d_{\text{ff}} = 4d_{\text{model}}$):

$$\text{First projection: } 2Bnd_{\text{model}}d_{\text{ff}} = 8Bnd_{\text{model}}^2 \quad (10.87)$$

$$\text{Second projection: } 2Bnd_{\text{ff}}d_{\text{model}} = 8Bnd_{\text{model}}^2 \quad (10.88)$$

$$\text{Total FFN: } 16Bnd_{\text{model}}^2 \quad (10.89)$$

Total per encoder layer:

$$\text{FLOPs} = 24Bnd_{\text{model}}^2 + 4Bn^2d_{\text{model}} \quad (10.90)$$

The crossover point where attention and FFN have equal cost occurs when $8nd_{\text{model}}^2 = 4n^2d_{\text{model}}$, which simplifies to $n = 2d_{\text{model}}$. For $n < 2d_{\text{model}}$, the feed-forward network dominates; for $n > 2d_{\text{model}}$, attention dominates.

This crossover point has important implications for model design and optimization. For BERT-base with $d_{\text{model}} = 768$, the crossover occurs at $n = 1536$ tokens. Since BERT uses maximum sequence length 512, the feed-forward network accounts for approximately 67% of computation per layer. For longer-sequence models like Longformer or BigBird that process 4096 tokens, attention accounts for approximately 80% of computation. This explains why efficient attention mechanisms focus on reducing the $O(n^2)$ term—it's the bottleneck for long sequences.

Example 10.7 (BERT-base Training FLOPs). For BERT-base with 12 encoder layers, $d_{\text{model}} = 768$, batch size $B = 32$, sequence length $n = 512$:

Per encoder layer:

$$\text{Attention: } 8 \times 32 \times 512 \times 768^2 + 4 \times 32 \times 512^2 \times 768 \quad (10.91)$$

$$= 77.3 \text{ GFLOPs} + 16.1 \text{ GFLOPs} = 93.4 \text{ GFLOPs} \quad (10.92)$$

$$\text{Feed-forward: } 16 \times 32 \times 512 \times 768^2 = 154.6 \text{ GFLOPs} \quad (10.93)$$

$$\text{Total per layer: } 248.0 \text{ GFLOPs} \quad (10.94)$$

All 12 layers: $12 \times 248.0 = 2,976 \text{ GFLOPs} \approx 3.0 \text{ TFLOPs}$ per training step

Backward pass: Approximately $2 \times$ forward pass = 6.0 TFLOPs

Total per training step: $3.0 + 6.0 = 9.0 \text{ TFLOPs}$

Training time on NVIDIA A100 (312 TFLOPS FP16):

Theoretical minimum: $9.0/312 = 28.8$ milliseconds per step

In practice, memory bandwidth limitations, kernel launch overhead, and non-matrix operations reduce efficiency to approximately 40-50% of peak FLOPs, yielding actual training time of 60-75 milliseconds per step. This corresponds to throughput of 13-17 training steps per second, or approximately 210,000-270,000 tokens per second.

For the full BERT-base training (1 million steps), this amounts to:

$$\text{Total FLOPs: } 9.0 \times 10^{12} \times 10^6 = 9.0 \times 10^{18} \text{ FLOPs} \quad (10.95)$$

$$\text{Training time on A100: } \frac{9.0 \times 10^{18}}{312 \times 10^{12} \times 0.45} \approx 64,000 \text{ seconds} \approx 18 \text{ hours} \quad (10.96)$$

This assumes 45% efficiency and continuous training. In practice, BERT-base training takes approximately 3-4 days on 16 V100 GPUs (equivalent to 1-2 days on 16 A100 GPUs), accounting for data loading, checkpointing, and other overhead.

10.5.2 Memory Bandwidth Considerations

While FLOPs provide a theoretical upper bound on training speed, memory bandwidth often becomes the practical bottleneck. Modern GPUs have enormous computational capacity but limited memory bandwidth. The NVIDIA A100 provides 312 TFLOPS of FP16 compute but only 1.6 TB/s of memory bandwidth. For operations to be compute-bound (limited by FLOPs rather than memory), they must have high arithmetic intensity—the ratio of FLOPs to bytes transferred.

Matrix multiplication has arithmetic intensity $O(n)$ for $n \times n$ matrices, making it compute-bound for large matrices. However, element-wise operations like activation functions, layer normalization, and residual additions have arithmetic intensity $O(1)$, making them memory-bound. For transformers,

the large matrix multiplications in attention and feed-forward networks are typically compute-bound, while the numerous element-wise operations between them are memory-bound. This is why kernel fusion—combining multiple operations into a single kernel to reduce memory transfers—is crucial for transformer efficiency.

Example 10.8 (Memory Bandwidth Analysis). For BERT-base layer with $d_{\text{model}} = 768$, batch size 32, sequence length 512:

Feed-forward first projection:

$$\text{Input: } 32 \times 512 \times 768 \times 2 = 25.2 \text{ MB (FP16)} \quad (10.97)$$

$$\text{Weight: } 768 \times 3072 \times 2 = 4.7 \text{ MB (FP16)} \quad (10.98)$$

$$\text{Output: } 32 \times 512 \times 3072 \times 2 = 100.7 \text{ MB (FP16)} \quad (10.99)$$

$$\text{Total memory: } 130.6 \text{ MB} \quad (10.100)$$

$$\text{FLOPs: } 2 \times 32 \times 512 \times 768 \times 3072 = 77.3 \text{ GFLOPs} \quad (10.101)$$

Arithmetic intensity: $77.3 \times 10^9 / (130.6 \times 10^6) = 592 \text{ FLOPs/byte}$

Time on A100:

$$\text{Compute-bound: } 77.3 / 312,000 = 0.25 \text{ ms} \quad (10.102)$$

$$\text{Memory-bound: } 130.6 / 1,600,000 = 0.08 \text{ ms} \quad (10.103)$$

Since compute time exceeds memory time, this operation is compute-bound. The GPU's computational capacity is the bottleneck, not memory bandwidth.

Layer normalization:

$$\text{Input/output: } 32 \times 512 \times 768 \times 2 \times 2 = 50.3 \text{ MB (read + write)} \quad (10.104)$$

$$\text{FLOPs: } \approx 32 \times 512 \times 768 \times 10 = 1.3 \text{ GFLOPs (approximate)} \quad (10.105)$$

Arithmetic intensity: $1.3 \times 10^9 / (50.3 \times 10^6) = 26 \text{ FLOPs/byte}$

Time on A100:

$$\text{Compute-bound: } 1.3 / 312,000 = 0.004 \text{ ms} \quad (10.106)$$

$$\text{Memory-bound: } 50.3 / 1,600,000 = 0.031 \text{ ms} \quad (10.107)$$

Since memory time exceeds compute time, layer normalization is memory-bound. The GPU's memory bandwidth is the bottleneck, not computational capacity. This is why fusing layer normalization with adjacent operations can significantly improve performance.

10.5.3 Scaling to Large Models

As transformer models scale from millions to billions of parameters, the computational and memory requirements grow dramatically. GPT-3 with 175 billion parameters requires approximately 700 GB of memory just to store the parameters in FP32 (or 350 GB in FP16), far exceeding the capacity of any single GPU. This necessitates model parallelism, where the model is split across multiple GPUs. The three main parallelism strategies are data parallelism (different GPUs process different batches), model parallelism (different GPUs hold different parts of the model), and pipeline parallelism (different GPUs process different layers).

Example 10.9 (GPT-3 Scale Analysis). GPT-3 (175B parameters) uses 96 layers, $d_{\text{model}} = 12,288$, $h = 96$ heads, $d_{ff} = 49,152$:

Parameter count per layer:

$$\text{Self-attention: } 4 \times 12,288^2 = 604,045,824 \approx 604\text{M} \quad (10.108)$$

$$\text{Feed-forward: } 2 \times 12,288 \times 49,152 = 1,208,091,648 \approx 1,208\text{M} \quad (10.109)$$

$$\text{Total per layer: } \approx 1,812\text{M parameters} \quad (10.110)$$

Total model: $96 \times 1,812 = 173,952\text{M} \approx 174\text{B}$ parameters (plus embeddings $\approx 1\text{B}$)

Memory requirements:

$$\text{Parameters (FP16): } 175 \times 10^9 \times 2 = 350 \text{ GB} \quad (10.111)$$

$$\text{Gradients (FP16): } 350 \text{ GB} \quad (10.112)$$

$$\text{Adam states (FP32): } 175 \times 10^9 \times 8 = 1,400 \text{ GB} \quad (10.113)$$

$$\text{Activations (batch 1, seq 2048): } \approx 60 \text{ GB} \quad (10.114)$$

$$\text{Total: } \approx 2,160 \text{ GB} \quad (10.115)$$

This requires at minimum 28 A100 GPUs (80 GB each) just to store parameters and optimizer states, not including activations. In practice, GPT-3 training used hundreds of GPUs with sophisticated parallelism strategies.

Training cost estimate:

For 300 billion tokens (GPT-3 training corpus):

$$\text{FLOPs per token: } \approx 6 \times 175 \times 10^9 = 1.05 \times 10^{12} \text{ FLOPs} \quad (10.116)$$

$$\text{Total FLOPs: } 300 \times 10^9 \times 1.05 \times 10^{12} = 3.15 \times 10^{23} \text{ FLOPs} \quad (10.117)$$

At 40% efficiency on A100 (312 TFLOPS FP16):

$$\text{GPU-hours: } \frac{3.15 \times 10^{23}}{312 \times 10^{12} \times 0.4 \times 3600} \approx 700,000 \text{ GPU-hours} \quad (10.118)$$

With 1024 A100 GPUs, this amounts to approximately 28 days of training. At cloud pricing of approximately \$2.50/hour per A100, the compute cost is approximately \$1.75 million, not including data storage, networking, or engineering costs.

10.6 Complete Transformer Architecture

10.6.1 Full Encoder-Decoder Model

[H] Transformer Forward Pass Source sequence $\mathbf{x} = [x_1, \dots, x_n]$, target sequence $\mathbf{y} = [y_1, \dots, y_m]$
 Predicted probabilities for each target position

Encoder $\mathbf{X}_{\text{emb}} = \text{Embedding}(\mathbf{x})$

$\mathbf{X}^{(0)} = \mathbf{X}_{\text{emb}} + \text{PositionalEncoding}(\text{positions})$

$\ell = 1 \dots N_{\text{enc}} \quad \mathbf{X}^{(\ell)} = \text{EncoderLayer}^{(\ell)}(\mathbf{X}^{(\ell-1)}) \quad \mathbf{X}_{\text{enc}} = \mathbf{X}^{(N_{\text{enc}})}$

Decoder $\mathbf{Y}_{\text{emb}} = \text{Embedding}(\mathbf{y})$

$\mathbf{Y}^{(0)} = \mathbf{Y}_{\text{emb}} + \text{PositionalEncoding}(\text{positions})$

$\ell = 1 \dots N_{\text{dec}} \quad \mathbf{Y}^{(\ell)} = \text{DecoderLayer}^{(\ell)}(\mathbf{Y}^{(\ell-1)}, \mathbf{X}_{\text{enc}}) \quad \mathbf{Y}_{\text{dec}} = \mathbf{Y}^{(N_{\text{dec}})}$

Output Projection logits $= \mathbf{Y}_{\text{dec}} \mathbf{W}_{\text{out}} + \mathbf{b}_{\text{out}}$ where $\mathbf{W}_{\text{out}} \in \mathbb{R}^{d_{\text{model}} \times V}$
 probs $= \text{softmax}(\text{logits})$

probs

10.6.2 Original Transformer Configuration

”Attention is All You Need” base model:

- Encoder layers: $N_{\text{enc}} = 6$

- Decoder layers: $N_{\text{dec}} = 6$
- Model dimension: $d_{\text{model}} = 512$
- Attention heads: $h = 8$
- Feed-forward dimension: $d_{\text{ff}} = 2048$
- Dropout rate: $p = 0.1$

Parameter count:

$$\text{Encoder (6 layers): } 6 \times (\text{attn} + \text{FFN}) \approx 25M \quad (10.119)$$

$$\text{Decoder (6 layers): } 6 \times (2 \times \text{attn} + \text{FFN}) \approx 31M \quad (10.120)$$

$$\text{Embeddings: varies by vocabulary} \quad (10.121)$$

$$\text{Total (excluding embeddings): } \approx \mathbf{56M \text{ parameters}} \quad (10.122)$$

10.7 Residual Connections and Layer Normalization

10.7.1 Residual Connections

Residual connections, also known as skip connections, are fundamental to enabling the training of deep transformer networks. Without residual connections, gradients would need to flow through dozens of attention and feed-forward layers during backpropagation, leading to vanishing or exploding gradients that make optimization extremely difficult. The residual connection provides a direct path from each layer's output back to its input, allowing gradients to flow unimpeded through the network. This gradient highway ensures that even the earliest layers receive meaningful gradient signals, enabling effective training of networks with 96 layers (GPT-3) or more.

The residual connection pattern in transformers follows the post-addition layer normalization structure: $\text{LayerNorm}(x + \text{Sublayer}(x))$. This means the sublayer's output is added to its input before normalization. The addition operation has a gradient of 1 with respect to both operands, so during backpropagation, gradients flow both through the sublayer (learning to refine representations) and directly through the residual connection (providing a gradient highway). This dual path enables the network to learn both identity mappings (when the sublayer output is near zero) and complex transformations (when the sublayer output is large).

The residual connection also enables the network to learn incrementally. Early in training, the sublayer outputs are typically small due to weight initialization, so the network effectively starts as a near-identity function. As training progresses, the sublayers learn to make increasingly sophisticated transformations, building on the representations from previous layers. This incremental learning is much more stable than trying to learn the complete transformation from scratch. For a 12-layer BERT model, each layer can focus on learning a small refinement rather than a complete transformation, making optimization tractable.

10.7.2 Layer Normalization

Layer normalization stabilizes training by normalizing activations across the feature dimension, ensuring that each layer receives inputs with consistent statistics regardless of how previous layers' parameters change during training. Unlike batch normalization, which normalizes across the batch dimension and is commonly used in convolutional networks, layer normalization normalizes across features for each example independently. This independence from batch size is crucial for transformers, which often use small batch sizes during inference or fine-tuning, and for handling variable-length sequences where batch normalization's statistics would be unreliable.

Definition 10.7 (Layer Normalization). For input $\mathbf{x} \in \mathbb{R}^d$, layer normalization computes mean and variance across the feature dimension, then normalizes and applies learned affine transforma-

tion:

$$\mu = \frac{1}{d} \sum_{i=1}^d x_i \quad (10.123)$$

$$\sigma^2 = \frac{1}{d} \sum_{i=1}^d (x_i - \mu)^2 \quad (10.124)$$

$$\hat{x}_i = \frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}} \quad (10.125)$$

$$y_i = \gamma \hat{x}_i + \beta \quad (10.126)$$

where $\gamma, \beta \in \mathbb{R}^d$ are learnable scale and shift parameters, and $\epsilon \approx 10^{-5}$ prevents division by zero.

For a batch of sequences $\mathbf{X} \in \mathbb{R}^{B \times n \times d}$, layer normalization is applied independently to each of the $B \times n$ vectors, normalizing across the d features.

The learned parameters γ and β allow the network to undo the normalization if beneficial. If $\gamma_i = \sqrt{\sigma^2 + \epsilon}$ and $\beta_i = \mu$, the normalization is completely undone. In practice, the network learns appropriate values that balance normalization's stabilizing effect with the flexibility to learn arbitrary distributions.

Layer normalization differs fundamentally from batch normalization in its normalization dimension. Batch normalization computes statistics across the batch dimension (normalizing each feature across all examples in the batch), making it dependent on batch size and batch composition. Layer normalization computes statistics across the feature dimension (normalizing all features for each example independently), making it independent of batch size. For transformers processing variable-length sequences with potentially small batch sizes, this independence is essential. A batch size of 1 works perfectly with layer normalization but would be problematic for batch normalization.

10.7.3 Pre-Norm vs Post-Norm

The placement of layer normalization relative to the residual connection significantly impacts training dynamics. The original transformer paper used post-norm: $\text{LayerNorm}(x + \text{Sublayer}(x))$, where normalization is applied after adding the residual. More recent models like GPT-2 and GPT-3 use pre-norm: $x + \text{LayerNorm}(\text{Sublayer}(x))$, where normalization is applied before the sublayer, and the residual connection bypasses normalization entirely.

Post-norm architecture normalizes the sum of the input and sublayer output, which can help prevent activation magnitudes from growing unboundedly as depth increases. However, post-norm requires careful learning rate warmup and can be unstable for very deep networks. The gradients must flow through the layer normalization operation, which can introduce additional numerical instabilities. BERT uses post-norm with 12-24 layers successfully, but scaling to 96+ layers becomes challenging.

Pre-norm architecture applies normalization before each sublayer, so the sublayer receives normalized inputs. The residual connection then adds the sublayer output directly to the (unnormalized) input, bypassing the normalization. This provides a cleaner gradient path through the residual connection and tends to be more stable for very deep networks. GPT-2 and GPT-3 use pre-norm, enabling training of 48-96 layer models without learning rate warmup. The trade-off is that pre-norm may achieve slightly lower final performance than post-norm for shallow networks, but this difference diminishes for deeper networks where pre-norm's stability advantages dominate.

Example 10.10 (Layer Normalization Computation). For a single position's representation $\mathbf{x} \in \mathbb{R}^{768}$ from BERT-base:

Input: $\mathbf{x} = [0.5, -0.3, 1.2, \dots]$ (768 values)

Compute statistics:

$$\mu = \frac{1}{768} \sum_{i=1}^{768} x_i = 0.15 \quad (\text{example value}) \quad (10.127)$$

$$\sigma^2 = \frac{1}{768} \sum_{i=1}^{768} (x_i - 0.15)^2 = 0.42 \quad (\text{example value}) \quad (10.128)$$

$$\sigma = \sqrt{0.42 + 10^{-5}} = 0.648 \quad (10.129)$$

Normalize:

$$\hat{x}_1 = \frac{0.5 - 0.15}{0.648} = 0.540 \quad (10.130)$$

$$\hat{x}_2 = \frac{-0.3 - 0.15}{0.648} = -0.694 \quad (10.131)$$

$$\hat{x}_3 = \frac{1.2 - 0.15}{0.648} = 1.620 \quad (10.132)$$

$$\vdots \quad (10.133)$$

The normalized values $\hat{\mathbf{x}}$ have mean 0 and variance 1 across the 768 dimensions.

Apply learned affine transformation:

$$y_1 = \gamma_1 \times 0.540 + \beta_1 \quad (10.134)$$

$$y_2 = \gamma_2 \times (-0.694) + \beta_2 \quad (10.135)$$

$$y_3 = \gamma_3 \times 1.620 + \beta_3 \quad (10.136)$$

$$\vdots \quad (10.137)$$

where $\gamma, \beta \in \mathbb{R}^{768}$ are learned during training. Initially, γ is typically initialized to 1 and β to 0, making layer normalization initially act as pure normalization.

Memory and computation:

- Parameters: $2 \times 768 = 1,536$ (scale and shift)
- FLOPs per position: $\approx 10 \times 768 = 7,680$ (mean, variance, normalize, scale, shift)
- For batch 32, sequence 512: $32 \times 512 \times 7,680 = 126$ MFLOPs

Layer normalization is computationally cheap compared to attention or feed-forward networks, but it's memory-bound rather than compute-bound, so kernel fusion with adjacent operations is important for efficiency.

10.8 Training Objectives

10.8.1 Sequence-to-Sequence Training

For machine translation, minimize cross-entropy loss:

$$\mathcal{L} = - \sum_{t=1}^m \log P(y_t | y_{<t}, \mathbf{x}; \theta) \quad (10.138)$$

Teacher forcing: During training, use ground-truth previous tokens $y_{<t}$, not model predictions.

10.8.2 Autoregressive Generation

At inference, generate one token at a time: [H] Autoregressive Decoding Source sequence \mathbf{x} , max length T Generated sequence \mathbf{y}

Encode source: $\mathbf{X}_{\text{enc}} = \text{Encoder}(\mathbf{x})$

Initialize: $\mathbf{y} = [\text{BOS}]$ (begin-of-sequence token)

$t = 1$ T $\text{probs}_t = \text{Decoder}(\mathbf{y}, \mathbf{X}_{\text{enc}})$

$y_t = \arg \max(\text{probs}_t)$ (or sample from distribution)

Append y_t to \mathbf{y}

$y_t = \text{EOS}$ **break** (end-of-sequence token) \mathbf{y}

10.9 Transformer Variants: Architectural Patterns

While the original transformer uses both an encoder and decoder for sequence-to-sequence tasks, subsequent research has shown that encoder-only and decoder-only architectures can be highly effective for specific task families. These three architectural patterns—encoder-only, decoder-only, and encoder-decoder—represent different trade-offs between bidirectional context, autoregressive generation, and architectural complexity. Understanding these trade-offs is essential for choosing the right architecture for a given application.

10.9.1 Encoder-Only Architecture (BERT)

Encoder-only models use bidirectional self-attention throughout, allowing each position to attend to all other positions in both directions. This bidirectional context is ideal for understanding tasks where the model needs to build rich representations of the input but doesn't need to generate output sequences. BERT (Bidirectional Encoder Representations from Transformers) exemplifies this architecture, using 12 or 24 encoder layers with no decoder. The model is pre-trained using masked language modeling, where random tokens are masked and the model must predict them using bidirectional context.

The key advantage of encoder-only models is computational efficiency for understanding tasks. Since all positions can attend to all other positions, the entire sequence is processed in a single forward pass with full parallelization. For a classification task with sequence length 512, BERT requires one forward pass through 12 layers, computing attention over all 512^2 position pairs simultaneously. This is dramatically faster than autoregressive generation, which would require 512 sequential forward passes.

Encoder-only models excel at tasks requiring deep understanding of input text: sentiment classification, named entity recognition, question answering (when the answer is a span in the input), and semantic similarity. They are less suitable for generation tasks, though they can be adapted for generation through techniques like iterative refinement or by using the encoder representations to condition a separate decoder. The bidirectional attention means the model cannot be used for standard autoregressive generation without modification.

10.9.2 Decoder-Only Architecture (GPT)

Decoder-only models use causal (masked) self-attention, where each position can only attend to previous positions. This maintains the autoregressive property essential for generation: the model predicts each token based only on previous tokens, never "cheating" by looking ahead. GPT (Generative Pre-trained Transformer) exemplifies this architecture, using 12-96 decoder layers with no encoder. The model is pre-trained using standard language modeling, predicting the next token given all previous tokens.

The key advantage of decoder-only models is their simplicity and flexibility. With no encoder-decoder cross-attention, the architecture is simpler and has fewer parameters than an equivalent encoder-decoder model. More importantly, decoder-only models can handle both understanding and generation tasks through careful prompting. For classification, the model generates the class label; for question answering, it generates the answer; for translation, it generates the target language text. This unified interface enables few-shot learning, where the model learns new tasks from just a few examples in the prompt.

Decoder-only models excel at generation tasks: text completion, dialogue, creative writing, code generation, and few-shot learning. They can also handle understanding tasks by framing them as generation problems, though this may be less parameter-efficient than encoder-only models. The causal attention means generation is inherently sequential—generating 512 tokens requires 512 forward passes—but techniques like KV caching make this practical. Modern large language models like GPT-3, GPT-4, and LLaMA all use decoder-only architectures due to their flexibility and scaling properties.

10.9.3 Encoder-Decoder Architecture (T5, BART)

Encoder-decoder models combine both architectural patterns: a bidirectional encoder processes the input, and a causal decoder generates the output while attending to the encoder through cross-attention. This is the original transformer architecture from “Attention is All You Need,” and it remains optimal for sequence-to-sequence tasks where the input and output are distinct sequences. T5 (Text-to-Text Transfer Transformer) and BART (Bidirectional and Autoregressive Transformer) exemplify modern encoder-decoder models.

The key advantage of encoder-decoder models is their explicit separation of understanding and generation. The encoder can use bidirectional attention to build rich representations of the input without worrying about causality, while the decoder can focus on generation while attending to relevant parts of the input through cross-attention. This separation is particularly valuable for tasks like translation, where the source and target languages have different structures, or summarization, where the output is much shorter than the input.

Encoder-decoder models excel at sequence-to-sequence tasks: machine translation, summarization, question answering (when generating free-form answers), and text simplification. They require approximately twice the parameters of encoder-only or decoder-only models of similar capacity (due to having both encoder and decoder stacks), but this investment pays off for tasks requiring both deep understanding and flexible generation. The cross-attention mechanism provides interpretability, showing which source positions the model attends to when generating each target token.

Example 10.11 (Architecture Comparison: BERT vs GPT-2 vs T5-base). **BERT-base (Encoder-only):**

- Layers: 12 encoder layers
- Dimensions: $d_{\text{model}} = 768$, $h = 12$, $d_{\text{ff}} = 3072$
- Parameters: 110M (embeddings + 12 encoder layers)
- Attention: Bidirectional self-attention in all layers
- Pre-training: Masked language modeling (predict masked tokens)
- Inference: Single forward pass for sequence length n
- Best for: Classification, NER, extractive QA

GPT-2 (Decoder-only):

- Layers: 12 decoder layers (but no cross-attention, so effectively simplified decoders)
- Dimensions: $d_{\text{model}} = 768$, $h = 12$, $d_{\text{ff}} = 3072$
- Parameters: 117M (embeddings + 12 decoder layers without cross-attention)
- Attention: Causal self-attention in all layers
- Pre-training: Autoregressive language modeling (predict next token)

- Inference: n forward passes for sequence length n (autoregressive)
- Best for: Text generation, few-shot learning, dialogue

T5-base (Encoder-Decoder):

- Layers: 12 encoder layers + 12 decoder layers
- Dimensions: $d_{\text{model}} = 768$, $h = 12$, $d_{ff} = 3072$
- Parameters: 220M (embeddings + 12 encoders + 12 decoders with cross-attention)
- Attention: Bidirectional in encoder, causal + cross-attention in decoder
- Pre-training: Span corruption (predict masked spans)
- Inference: One encoder pass + m decoder passes for output length m
- Best for: Translation, summarization, generative QA

Parameter breakdown comparison:

BERT encoder layer:	7.1M parameters	(10.139)
GPT-2 decoder layer (no cross-attn):	7.1M parameters	(10.140)
T5 encoder layer:	7.1M parameters	(10.141)
T5 decoder layer (with cross-attn):	9.4M parameters	(10.142)

Notice that GPT-2's "decoder" layers are actually simpler than true decoder layers because they lack cross-attention. This makes GPT-2 and BERT have similar parameter counts despite different attention patterns.

Computational cost for sequence length 512:

BERT (understanding):	1 forward pass \times 12 layers = 12 layer passes	(10.143)
GPT-2 (generation):	512 forward passes \times 12 layers = 6,144 layer passes	(10.144)
T5 (translation):	1 encoder pass \times 12 + 512 decoder passes \times 12	(10.145)
	= 12 + 6,144 = 6,156 layer passes	(10.146)

This illustrates why generation is much slower than understanding: autoregressive decoding requires hundreds of sequential forward passes, while understanding requires just one parallel forward pass.

10.9.4 Choosing the Right Architecture

The choice between encoder-only, decoder-only, and encoder-decoder architectures depends on the task requirements and deployment constraints. For pure understanding tasks (classification, entity recognition, span-based QA), encoder-only models like BERT provide the best parameter efficiency and inference speed. For pure generation tasks (text completion, creative writing, code generation), decoder-only models like GPT provide simplicity and flexibility. For sequence-to-sequence tasks with distinct input and output (translation, summarization, generative QA), encoder-decoder models like T5 provide the best performance despite higher parameter counts.

Recent trends favor decoder-only architectures for their versatility. Large language models like GPT-3, GPT-4, PaLM, and LLaMA all use decoder-only architectures, handling both understanding and generation through prompting. This unified architecture simplifies deployment (one model for all tasks) and enables few-shot learning (learning new tasks from examples in the prompt). However, for specific applications where understanding or sequence-to-sequence performance is critical, encoder-only or encoder-decoder models may still provide better parameter efficiency and performance.

The computational trade-offs also matter for deployment. Encoder-only models are fastest for

understanding tasks, requiring one forward pass regardless of sequence length. Decoder-only models are slower for generation due to autoregressive decoding, but KV caching makes this practical. Encoder-decoder models combine both costs: one encoder pass plus autoregressive decoder passes. For latency-sensitive applications, these computational differences can be decisive.

10.10 Exercises

Exercise 10.1. For transformer with $N = 6$, $d_{\text{model}} = 512$, $h = 8$, $d_{ff} = 2048$, $V = 32000$:

1. Calculate total parameters in encoder
2. Calculate total parameters in decoder
3. What percentage are in embeddings vs transformer layers?
4. How does this change if vocabulary increases to 50,000?

Exercise 10.2. Implement single transformer encoder layer in PyTorch. Test with batch size 16, sequence length 64, $d_{\text{model}} = 256$. Verify output shape and gradient flow through residual connections.

Exercise 10.3. Compare memory and computation for:

1. Encoder processing sequence length 1024
2. Decoder generating 1024 tokens autoregressively

Why is decoding slower? How many forward passes required?

Exercise 10.4. Show that layer normalization is invariant to input scale: if $\mathbf{x}' = c\mathbf{x}$ for constant $c > 0$, then $\text{LayerNorm}(\mathbf{x}') = \text{LayerNorm}(\mathbf{x})$ (ignoring learnable γ, β).

Chapter 11

Training Transformers

Chapter Overview

Training transformers requires specialized techniques beyond standard optimization. This chapter provides comprehensive coverage of transformer training procedures, from loss functions and backpropagation through the architecture to optimization algorithms, learning rate schedules, and hardware-efficient training strategies. We examine why transformers need warmup, how mixed precision training reduces memory consumption, when to use gradient accumulation and checkpointing, and how distributed training enables models that exceed single-GPU capacity. Throughout, we provide detailed hardware analysis, memory calculations, and practical guidance drawn from training state-of-the-art models like BERT, GPT-2, and GPT-3.

Learning Objectives

1. Understand training objectives and loss functions for different transformer architectures
2. Analyze gradient flow and backpropagation through transformer layers
3. Implement optimization algorithms (Adam, AdamW, LAMB) with appropriate hyperparameters
4. Apply learning rate schedules with warmup and decay
5. Use mixed precision training to reduce memory and accelerate training
6. Apply gradient accumulation and checkpointing for memory-constrained scenarios
7. Understand distributed training strategies for large-scale models
8. Select appropriate batch sizes and sequence lengths based on hardware constraints
9. Apply regularization techniques to prevent overfitting
10. Estimate training time and costs for transformer models

11.1 Training Objectives and Loss Functions

The training objective fundamentally shapes how a transformer learns and what capabilities it develops. Different transformer architectures employ distinct training objectives tailored to their intended use cases, from masked language modeling in BERT to causal language modeling in GPT to sequence-to-sequence learning in T5. Understanding these objectives in depth—including their mathematical formulations, computational requirements, and practical implications—is essential for training transformers effectively.

11.1.1 Masked Language Modeling

Masked language modeling, introduced by BERT, trains the model to predict randomly masked tokens based on bidirectional context. This objective enables the model to learn rich representations that capture relationships in both directions, making it particularly effective for tasks requiring understanding of complete sentences or documents.

The masking strategy is more sophisticated than simply replacing tokens with a special [MASK] symbol. BERT’s approach selects 15% of tokens for prediction, but handles them in three different ways: 80% are replaced with [MASK], 10% are replaced with random tokens from the vocabulary, and 10% are left unchanged. This strategy prevents the model from simply memorizing that [MASK] tokens need prediction and forces it to maintain representations for all tokens, since any token might need to be predicted. The random token replacement encourages the model to use context to correct errors, while leaving some tokens unchanged helps the model learn that not all tokens are corrupted.

The loss function for masked language modeling is cross-entropy computed only over the masked positions. For a sequence $\mathbf{x} = (x_1, \dots, x_n)$ with masked positions $M \subseteq \{1, \dots, n\}$, the loss is:

$$L_{\text{MLM}} = -\frac{1}{|M|} \sum_{i \in M} \log P(x_i | \mathbf{x}_{\setminus M}) \quad (11.1)$$

where $\mathbf{x}_{\setminus M}$ denotes the sequence with masked positions corrupted according to the strategy above. The model outputs logits $\mathbf{z}_i \in \mathbb{R}^V$ for each position i , where V is the vocabulary size, and the probability distribution is obtained via softmax: $P(x_i | \mathbf{x}_{\setminus M}) = \text{softmax}(\mathbf{z}_i)_{x_i}$.

The computational and memory implications of this loss are significant. For vocabulary size $V = 30,000$, sequence length $n = 512$, and batch size $B = 32$, the output logits tensor has shape $\mathbb{R}^{32 \times 512 \times 30000}$, requiring $32 \times 512 \times 30,000 \times 4 = 1,966,080,000$ bytes, or approximately 1.97 GB of memory just for the logits in FP32. This massive memory footprint explains why the output projection and softmax computation often become bottlenecks during training. The memory requirement can be reduced by computing the loss in chunks (processing subsets of positions at a time) or by using mixed precision training where logits are computed in FP16, though care must be taken to maintain numerical stability in the softmax operation.

In practice, BERT-base masks approximately 77 tokens per sequence (15% of 512), so the loss is computed over $32 \times 77 = 2,464$ predictions per batch. The cross-entropy computation requires exponentiating 30,000 logits for each prediction to compute the softmax denominator, then taking the logarithm of the target class probability. Modern implementations optimize this by fusing the softmax and cross-entropy operations and by using numerically stable implementations that subtract the maximum logit before exponentiation to prevent overflow.

11.1.2 Causal Language Modeling

Causal language modeling, used in GPT and other decoder-only models, trains the model to predict the next token given all previous tokens. Unlike masked language modeling, which uses bidirectional context, causal language modeling uses only left-to-right context, enforced through causal attention masks that prevent positions from attending to future positions.

The training objective is to maximize the likelihood of each token given its preceding context. For a sequence $\mathbf{x} = (x_1, \dots, x_n)$, the loss is:

$$L_{\text{CLM}} = -\frac{1}{n} \sum_{i=1}^n \log P(x_i | x_1, \dots, x_{i-1}) \quad (11.2)$$

This formulation means that every position in the sequence contributes to the loss, unlike masked language modeling where only 15% of positions contribute. For a batch of 32 sequences of length 512, we compute loss over $32 \times 512 = 16,384$ predictions, compared to only 2,464 for BERT’s masked language modeling. This makes causal language modeling more sample-efficient in terms of predictions per sequence, though the unidirectional context may be less informative than bidirectional context for some tasks.

A crucial distinction exists between training and inference for causal language models. During training, we use teacher forcing: the model receives the ground-truth previous tokens as input, even if it would have predicted different tokens. This enables parallel computation of the loss across all positions in a sequence, since we can compute $P(x_i|x_1, \dots, x_{i-1})$ for all i simultaneously using causal masking. During inference, however, generation is autoregressive: the model generates one token at a time, using its own predictions as input for subsequent positions. This sequential generation process is much slower than parallel training, which motivates optimizations like KV caching (discussed in Chapter 12).

The memory requirements for causal language modeling are similar to masked language modeling: the output logits tensor for batch size 32, sequence length 512, and vocabulary size 50,257 (GPT-2’s vocabulary) requires $32 \times 512 \times 50,257 \times 4 = 3,296,019,456$ bytes, or approximately 3.3 GB in FP32. However, since we compute loss over all positions rather than just 15%, the gradient computation is more expensive. The backward pass through the output projection receives gradients from all 16,384 predictions rather than just 2,464, increasing the gradient computation cost proportionally.

11.1.3 Sequence-to-Sequence Training

Sequence-to-sequence models like T5 and BART use encoder-decoder architectures where the encoder processes the input sequence bidirectionally and the decoder generates the output sequence autoregressively. The training objective combines aspects of both masked and causal language modeling: the encoder can use bidirectional attention over the input, while the decoder uses causal attention over the output sequence and cross-attention to the encoder’s representations.

The loss function for sequence-to-sequence training is computed over the target sequence. For input sequence $\mathbf{x} = (x_1, \dots, x_n)$ and target sequence $\mathbf{y} = (y_1, \dots, y_m)$:

$$L_{\text{seq2seq}} = -\frac{1}{m} \sum_{j=1}^m \log P(y_j|y_1, \dots, y_{j-1}, \mathbf{x}) \quad (11.3)$$

Like causal language modeling, sequence-to-sequence training uses teacher forcing during training: the decoder receives the ground-truth previous target tokens as input, enabling parallel computation of the loss. This differs from inference, where the decoder must generate tokens sequentially using its own predictions.

The memory requirements for sequence-to-sequence models are higher than encoder-only or decoder-only models because both encoder and decoder activations must be stored. For T5-base with input length 512, target length 512, and batch size 32, we must store encoder activations ($32 \times 512 \times 768$ per layer), decoder activations ($32 \times 512 \times 768$ per layer), and cross-attention activations ($32 \times 12 \times 512 \times 512$ for attention matrices between decoder and encoder). The total activation memory is roughly 1.5-2× that of an encoder-only model of the same size.

Different sequence-to-sequence models use different input corruption strategies. T5 uses span corruption, where contiguous spans of tokens are replaced with sentinel tokens and the model must predict the original spans. BART uses a variety of corruption strategies including token masking, token deletion, sentence permutation, and document rotation. These diverse corruption strategies help the model learn robust representations that generalize across different types of noise and transformations.

11.2 Backpropagation Through Transformers

Understanding how gradients flow through the transformer architecture is essential for diagnosing training issues, designing better architectures, and implementing custom training procedures. The transformer’s combination of attention mechanisms, residual connections, layer normalization, and feed-forward networks creates a complex gradient flow pattern that differs fundamentally from simpler architectures like MLPs or CNNs.

11.2.1 Gradient Flow Analysis

Backpropagation through a transformer begins at the output and flows backward through each component. For a language modeling task, the loss L is computed from the output logits, and we must compute gradients with respect to all parameters in the model. The gradient flow follows the reverse path of the forward computation, with each operation contributing its Jacobian to the chain rule.

The output projection layer maps the final transformer layer's output to vocabulary logits. For output $\mathbf{h}_n \in \mathbb{R}^{d_{\text{model}}}$ at position n and output weight matrix $\mathbf{W}^{\text{out}} \in \mathbb{R}^{d_{\text{model}} \times V}$, the logits are $\mathbf{z}_n = \mathbf{W}^{\text{out}^\top} \mathbf{h}_n$. The gradient of the loss with respect to the output weights is:

$$\frac{\partial L}{\partial \mathbf{W}^{\text{out}}} = \sum_{i=1}^n \mathbf{h}_i \frac{\partial L}{\partial \mathbf{z}_i}^\top \quad (11.4)$$

where $\frac{\partial L}{\partial \mathbf{z}_i} \in \mathbb{R}^V$ is the gradient from the softmax and cross-entropy loss. This gradient matrix has the same shape as \mathbf{W}^{out} : $\mathbb{R}^{d_{\text{model}} \times V}$. For BERT-base with $d_{\text{model}} = 768$ and $V = 30,000$, this gradient requires $768 \times 30,000 \times 4 = 92,160,000$ bytes (92 MB) in FP32.

The gradient with respect to the output representations is:

$$\frac{\partial L}{\partial \mathbf{h}_i} = \mathbf{W}^{\text{out}} \frac{\partial L}{\partial \mathbf{z}_i} \quad (11.5)$$

This gradient then flows backward through each transformer layer. Within a layer, the gradient must flow through the feed-forward network, the second residual connection and layer normalization, the attention mechanism, and the first residual connection and layer normalization.

11.2.2 Gradients Through Residual Connections

Residual connections are crucial for training deep transformers because they provide "gradient highways" that allow gradients to flow directly through many layers without vanishing. Consider a residual block with function F :

$$\mathbf{y} = \mathbf{x} + F(\mathbf{x}) \quad (11.6)$$

The gradient with respect to the input is:

$$\frac{\partial L}{\partial \mathbf{x}} = \frac{\partial L}{\partial \mathbf{y}} + \frac{\partial L}{\partial \mathbf{y}} \frac{\partial F(\mathbf{x})}{\partial \mathbf{x}} \quad (11.7)$$

The first term $\frac{\partial L}{\partial \mathbf{y}}$ is the direct gradient path that bypasses the function F entirely. This ensures that even if $\frac{\partial F(\mathbf{x})}{\partial \mathbf{x}}$ becomes very small (vanishing gradients) or very large (exploding gradients), the gradient $\frac{\partial L}{\partial \mathbf{x}}$ still receives the direct contribution $\frac{\partial L}{\partial \mathbf{y}}$. This is why transformers can be trained with many layers (BERT-large has 24 layers, GPT-3 has 96 layers) without suffering from vanishing gradients that plagued early deep networks.

For a transformer with L layers, the gradient from the output to the input has 2^L paths through the network: at each layer, the gradient can either flow through the residual connection (direct path) or through the attention/FFN (indirect path). This exponential number of paths creates a rich gradient flow that helps training, though in practice most gradient flows through the shorter paths that use more residual connections.

11.2.3 Gradients Through Layer Normalization

Layer normalization normalizes activations across the feature dimension, computing mean and variance for each position independently. For input $\mathbf{x} \in \mathbb{R}^d$, layer normalization computes:

$$\mathbf{y} = \frac{\mathbf{x} - \mu}{\sqrt{\sigma^2 + \epsilon}} \odot \gamma + \beta \quad (11.8)$$

where $\mu = \frac{1}{d} \sum_{i=1}^d x_i$, $\sigma^2 = \frac{1}{d} \sum_{i=1}^d (x_i - \mu)^2$, and $\gamma, \beta \in \mathbb{R}^d$ are learned scale and shift parameters.

The gradient computation for layer normalization is complex because the normalization couples all dimensions: changing one input element affects the mean and variance, which affects all output elements. The gradient with respect to the input is:

$$\frac{\partial L}{\partial \mathbf{x}} = \frac{\gamma}{\sqrt{\sigma^2 + \epsilon}} \left(\frac{\partial L}{\partial \mathbf{y}} - \frac{1}{d} \sum_{j=1}^d \frac{\partial L}{\partial y_j} - \frac{\mathbf{x} - \mu}{\sigma^2 + \epsilon} \frac{1}{d} \sum_{j=1}^d \frac{\partial L}{\partial y_j} (x_j - \mu) \right) \quad (11.9)$$

This gradient has three terms: the direct gradient scaled by the normalization factor, a mean-centering term, and a variance-correction term. The complexity of this gradient is why layer normalization is sometimes replaced with simpler alternatives like RMSNorm in some recent models, though layer normalization generally provides better training stability.

The learned parameters γ and β have simple gradients:

$$\frac{\partial L}{\partial \gamma} = \frac{\partial L}{\partial \mathbf{y}} \odot \frac{\mathbf{x} - \mu}{\sqrt{\sigma^2 + \epsilon}}, \quad \frac{\partial L}{\partial \beta} = \frac{\partial L}{\partial \mathbf{y}} \quad (11.10)$$

Layer normalization helps gradient flow by preventing activations from becoming too large or too small, which would cause gradients to vanish or explode. By maintaining normalized activations throughout the network, layer normalization ensures that gradients remain in a reasonable range, facilitating stable training.

11.2.4 Gradients Through Attention

The attention mechanism involves several matrix multiplications and a softmax operation, each contributing to the gradient computation. For self-attention with queries \mathbf{Q} , keys \mathbf{K} , and values \mathbf{V} :

$$\mathbf{O} = \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}} \right) \mathbf{V} \quad (11.11)$$

Working backward, the gradient with respect to the values is:

$$\frac{\partial L}{\partial \mathbf{V}} = \mathbf{A}^\top \frac{\partial L}{\partial \mathbf{O}} \quad (11.12)$$

where $\mathbf{A} = \text{softmax}(\mathbf{Q}\mathbf{K}^\top / \sqrt{d_k})$ is the attention matrix. This is a matrix multiplication of shape $(n \times n)^\top \times (n \times d_v) = (n \times d_v)$, matching the shape of \mathbf{V} .

The gradient with respect to the attention matrix is:

$$\frac{\partial L}{\partial \mathbf{A}} = \frac{\partial L}{\partial \mathbf{O}} \mathbf{V}^\top \quad (11.13)$$

This has shape $(n \times d_v) \times (d_v \times n) = (n \times n)$, matching the attention matrix shape.

The gradient must then flow through the softmax operation. For softmax output $\mathbf{a} = \text{softmax}(\mathbf{s})$, the Jacobian is:

$$\frac{\partial a_i}{\partial s_j} = a_i(\delta_{ij} - a_j) \quad (11.14)$$

where δ_{ij} is the Kronecker delta. This means the gradient with respect to the pre-softmax scores is:

$$\frac{\partial L}{\partial s_i} = \sum_j \frac{\partial L}{\partial a_j} a_j (\delta_{ij} - a_i) = a_i \left(\frac{\partial L}{\partial a_i} - \sum_j \frac{\partial L}{\partial a_j} a_j \right) \quad (11.15)$$

This computation must be performed for each row of the attention matrix independently, since softmax is applied row-wise.

Finally, gradients flow to the query and key projections. The gradient with respect to queries is:

$$\frac{\partial L}{\partial \mathbf{Q}} = \frac{1}{\sqrt{d_k}} \frac{\partial L}{\partial \mathbf{S}} \mathbf{K} \quad (11.16)$$

where $\mathbf{S} = \mathbf{QK}^\top / \sqrt{d_k}$ are the pre-softmax scores. The gradient with respect to keys is:

$$\frac{\partial L}{\partial \mathbf{K}} = \frac{1}{\sqrt{d_k}} \frac{\partial L}{\partial \mathbf{S}}^\top \mathbf{Q} \quad (11.17)$$

These gradients then flow through the projection matrices \mathbf{W}^Q , \mathbf{W}^K , and \mathbf{W}^V . For the query projection $\mathbf{Q} = \mathbf{XW}^Q$:

$$\frac{\partial L}{\partial \mathbf{W}^Q} = \mathbf{X}^\top \frac{\partial L}{\partial \mathbf{Q}} \quad (11.18)$$

This gradient has shape $(d_{\text{model}} \times n) \times (n \times d_k) = (d_{\text{model}} \times d_k)$, matching \mathbf{W}^Q . For BERT-base with $d_{\text{model}} = 768$ and $d_k = 64$, this requires $768 \times 64 \times 4 = 196,608$ bytes (197 KB) per head, or $12 \times 197 = 2.4$ MB for all 12 heads.

11.2.5 Gradients Through Feed-Forward Networks

The feed-forward network consists of two linear transformations with a non-linear activation (typically GELU) in between:

$$\text{FFN}(\mathbf{x}) = \mathbf{W}_2 \text{GELU}(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2 \quad (11.19)$$

The gradient with respect to the second layer weights is:

$$\frac{\partial L}{\partial \mathbf{W}_2} = \mathbf{h}^\top \frac{\partial L}{\partial \mathbf{y}} \quad (11.20)$$

where $\mathbf{h} = \text{GELU}(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1)$ is the intermediate activation. For BERT-base with $d_{\text{ff}} = 3072$ and $d_{\text{model}} = 768$, this gradient has shape (3072×768) and requires $3072 \times 768 \times 4 = 9,437,184$ bytes (9.4 MB) in FP32.

The gradient flows through the GELU activation. GELU is defined as:

$$\text{GELU}(x) = x\Phi(x) \quad (11.21)$$

where $\Phi(x)$ is the cumulative distribution function of the standard normal distribution. The derivative is:

$$\text{GELU}'(x) = \Phi(x) + x\phi(x) \quad (11.22)$$

where $\phi(x)$ is the probability density function. The gradient with respect to the pre-activation is:

$$\frac{\partial L}{\partial (\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1)} = \frac{\partial L}{\partial \mathbf{h}} \odot \text{GELU}'(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) \quad (11.23)$$

Finally, the gradient with respect to the first layer weights is:

$$\frac{\partial L}{\partial \mathbf{W}_1} = \mathbf{x}^\top \frac{\partial L}{\partial (\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1)} \quad (11.24)$$

This has shape $(d_{\text{model}} \times d_{\text{ff}}) = (768 \times 3072)$, also requiring 9.4 MB in FP32.

11.2.6 Computational Cost of Backpropagation

The backward pass through a transformer requires approximately twice the FLOPs of the forward pass. This factor of two arises because each matrix multiplication $\mathbf{Y} = \mathbf{XW}$ in the forward pass requires two matrix multiplications in the backward pass: $\frac{\partial L}{\partial \mathbf{W}} = \mathbf{X}^\top \frac{\partial L}{\partial \mathbf{Y}}$ and $\frac{\partial L}{\partial \mathbf{X}} = \frac{\partial L}{\partial \mathbf{Y}} \mathbf{W}^\top$. Each of these backward matrix multiplications has similar computational cost to the forward multiplication.

For BERT-base with 96.6 GFLOPs per forward pass, the backward pass requires approximately $2 \times 96.6 = 193.2$ GFLOPs. A complete training step (forward pass + backward pass) thus requires approximately $96.6 + 193.2 = 289.8$ GFLOPs, or roughly three times the forward pass cost. This $3 \times$ factor is a useful rule of thumb for estimating training costs from inference costs.

The memory requirements for backpropagation are substantial because all intermediate activations from the forward pass must be stored to compute gradients. For BERT-base with batch size 32 and sequence length 512, the activations require approximately 12 GB as analyzed in Chapter 12. This activation memory often dominates the total memory consumption during training, which motivates techniques like gradient checkpointing that trade computation for memory by recomputing activations during the backward pass.

11.3 Optimization Algorithms

The choice of optimization algorithm significantly impacts transformer training dynamics, convergence speed, and final model quality. While stochastic gradient descent (SGD) with momentum works well for many deep learning tasks, transformers benefit particularly from adaptive learning rate methods that adjust the learning rate for each parameter based on gradient statistics. The Adam family of optimizers has become the de facto standard for transformer training, with variants like AdamW and LAMB addressing specific challenges in large-scale training.

11.3.1 Adam Optimizer

Adam (Adaptive Moment Estimation) maintains exponential moving averages of both the gradient (first moment) and the squared gradient (second moment) for each parameter. These statistics enable adaptive per-parameter learning rates that automatically adjust based on the gradient history, helping with the varying scales of gradients across different layers and components of the transformer.

The Adam algorithm maintains two state vectors for each parameter \mathbf{w} : the first moment \mathbf{m} (exponential moving average of gradients) and the second moment \mathbf{v} (exponential moving average of squared gradients). At each training step t with gradient \mathbf{g}_t :

$$\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t \quad (11.25)$$

$$\mathbf{v}_t = \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \mathbf{g}_t^2 \quad (11.26)$$

where β_1 and β_2 are decay rates (typically $\beta_1 = 0.9$ and $\beta_2 = 0.999$). The squared gradient \mathbf{g}_t^2 is computed element-wise.

Because \mathbf{m} and \mathbf{v} are initialized to zero, they are biased toward zero, especially in early training steps. Adam corrects this bias by computing bias-corrected estimates:

$$\hat{\mathbf{m}}_t = \frac{\mathbf{m}_t}{1 - \beta_1^t} \quad (11.27)$$

$$\hat{\mathbf{v}}_t = \frac{\mathbf{v}_t}{1 - \beta_2^t} \quad (11.28)$$

The parameter update is then:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \frac{\hat{\mathbf{m}}_t}{\sqrt{\hat{\mathbf{v}}_t} + \epsilon} \quad (11.29)$$

where η is the learning rate and ϵ is a small constant (typically 10^{-8}) for numerical stability.

The adaptive learning rate $\frac{\eta}{\sqrt{\hat{\mathbf{v}}_t} + \epsilon}$ is larger for parameters with small historical gradients and smaller for parameters with large historical gradients. This adaptation is particularly beneficial for transformers because different components have vastly different gradient scales. Embedding layers, which are updated sparsely (only for tokens present in the batch), benefit from larger effective learning rates, while frequently updated parameters in the attention and FFN layers benefit from smaller effective learning rates that prevent overshooting.

The memory requirements for Adam are substantial: for each parameter, we must store the parameter itself, the gradient, the first moment, and the second moment. For a model with P parameters in FP32, Adam requires:

- Parameters: $P \times 4$ bytes
- Gradients: $P \times 4$ bytes
- First moments: $P \times 4$ bytes
- Second moments: $P \times 4$ bytes
- Total: $16P$ bytes

For BERT-base with 110 million parameters, Adam requires $110,000,000 \times 16 = 1,760,000,000$ bytes, or 1.76 GB, just for the optimizer state. This is four times the memory required for the parameters alone, and this overhead grows linearly with model size. For GPT-3 with 175 billion parameters, Adam would require $175,000,000,000 \times 16 = 2,800$ GB just for parameters and optimizer states, necessitating distributed training strategies that shard the optimizer state across multiple GPUs.

11.3.2 AdamW: Decoupled Weight Decay

AdamW modifies Adam by decoupling weight decay from the gradient-based update. In standard Adam with L2 regularization, the weight decay is incorporated into the gradient: $\mathbf{g}_t = \nabla L(\mathbf{w}_t) + \lambda \mathbf{w}_t$, where λ is the regularization coefficient. This means the weight decay is affected by the adaptive learning rate, which can lead to unexpected behavior.

AdamW instead applies weight decay directly to the parameters after the adaptive update:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \frac{\hat{\mathbf{m}}_t}{\sqrt{\hat{\mathbf{v}}_t + \epsilon}} - \eta \lambda \mathbf{w}_t \quad (11.30)$$

This decoupling means that weight decay acts as a true regularizer, shrinking parameters toward zero at a rate proportional to the learning rate, independent of the gradient statistics. In practice, this leads to better generalization, particularly for transformers where different parameters have very different gradient scales.

The typical weight decay coefficient for transformer training is $\lambda = 0.01$. However, weight decay is usually not applied to all parameters. Biases and layer normalization parameters (the scale γ and shift β parameters) are typically excluded from weight decay, as regularizing these parameters can hurt performance. The exclusion is implemented by maintaining separate parameter groups in the optimizer, with different weight decay settings for each group.

AdamW has become the standard optimizer for training transformers, used in BERT, GPT-2, GPT-3, T5, and most other modern models. The improved generalization from decoupled weight decay often allows training with higher learning rates, which can accelerate convergence. The memory requirements are identical to Adam: $16P$ bytes for a model with P parameters in FP32.

11.3.3 LAMB: Large Batch Training

LAMB (Layer-wise Adaptive Moments optimizer for Batch training) extends Adam to enable training with very large batch sizes, up to 64,000 or more. Large batch training is desirable because it improves hardware utilization and reduces training time by processing more examples in parallel, but naive scaling of the batch size often hurts convergence and final model quality.

The key insight of LAMB is to compute layer-wise learning rates that adapt based on the ratio of parameter norm to gradient norm within each layer. For layer l with parameters $\mathbf{w}^{(l)}$ and Adam update $\mathbf{u}^{(l)} = \frac{\hat{\mathbf{m}}^{(l)}}{\sqrt{\hat{\mathbf{v}}^{(l)} + \epsilon}} + \lambda \mathbf{w}^{(l)}$, LAMB computes:

$$\phi^{(l)} = \frac{\|\mathbf{w}^{(l)}\|_2}{\|\mathbf{u}^{(l)}\|_2} \quad (11.31)$$

The parameter update is then:

$$\mathbf{w}_{t+1}^{(l)} = \mathbf{w}_t^{(l)} - \eta \phi^{(l)} \mathbf{u}^{(l)} \quad (11.32)$$

This layer-wise adaptation ensures that the update magnitude is proportional to the parameter magnitude within each layer, preventing some layers from being updated too aggressively while others are updated too conservatively. This is particularly important for large batch training because large batches produce more accurate gradient estimates, which can lead to overly aggressive updates without proper scaling.

LAMB enabled training BERT-large to the same accuracy as the original paper in just 76 minutes using a batch size of 65,536 on 1,024 TPU v3 chips, compared to several days with standard batch sizes. The ability to use such large batches dramatically reduces training time for large-scale models, though it requires access to substantial computational resources to realize the benefits.

The memory requirements for LAMB are similar to Adam and AdamW: $16P$ bytes for a model with P parameters in FP32. The additional computation for layer-wise norm calculations is negligible compared to the forward and backward passes.

11.3.4 Optimizer Memory Comparison

Different optimizers have different memory footprints, which can be a critical consideration for large models:

- **SGD (no momentum):** $8P$ bytes (parameters + gradients in FP32)
- **SGD with momentum:** $12P$ bytes (parameters + gradients + momentum in FP32)
- **Adam/AdamW/LAMB:** $16P$ bytes (parameters + gradients + first moment + second moment in FP32)

For BERT-base with 110 million parameters:

- SGD: $110\text{M} \times 8 = 880$ MB
- SGD with momentum: $110\text{M} \times 12 = 1,320$ MB
- Adam/AdamW/LAMB: $110\text{M} \times 16 = 1,760$ MB

The additional memory overhead of Adam-family optimizers (880 MB compared to SGD) is usually worthwhile because the adaptive learning rates lead to faster convergence and better final performance. However, for very large models where memory is at a premium, techniques like ZeRO (Zero Redundancy Optimizer) can shard the optimizer state across multiple GPUs to reduce per-GPU memory requirements.

11.4 Learning Rate Schedules

Learning rate schedules are critical for transformer training, perhaps more so than for other architectures. Transformers are sensitive to the learning rate, and using a constant learning rate throughout training typically leads to poor results. The standard approach combines a warmup phase, where the learning rate increases from zero to a maximum value, with a decay phase, where the learning rate gradually decreases. This schedule helps stabilize early training and enables continued improvement in later training.

11.4.1 The Necessity of Warmup

Learning rate warmup is essential for stable transformer training. Without warmup, using the full learning rate from the beginning often causes training to diverge or get stuck in poor local minima. The instability arises from the interaction between large initial gradients and Adam’s adaptive learning rates.

In the first few training steps, Adam’s second moment estimates \mathbf{v} are very small because they are initialized to zero and have not yet accumulated gradient statistics. This means the effective learning rate $\frac{\eta}{\sqrt{\mathbf{v} + \epsilon}}$ is very large, potentially much larger than the nominal learning rate η . When combined with large gradients that are common early in training (when the model’s predictions are random and the loss is high), these large effective learning rates can cause parameter updates that are far too aggressive, leading to numerical instability or divergence.

Warmup solves this problem by starting with a very small learning rate and gradually increasing it over the first W steps (typically 10% of total training steps). During warmup, the learning rate at step t is:

$$\eta_t = \eta_{\max} \cdot \frac{t}{W} \quad (11.33)$$

This linear increase gives Adam’s moment estimates time to accumulate meaningful statistics while preventing overly aggressive updates. By the time the learning rate reaches its maximum value η_{\max} , the optimizer has stabilized and can handle the full learning rate safely.

The warmup period also serves another purpose: it allows the model to learn basic patterns before attempting more complex optimization. In the first few steps, the model learns simple statistics like token frequencies and basic co-occurrence patterns. These foundational patterns provide a stable base for learning more complex relationships later in training.

11.4.2 Warmup Plus Linear Decay

The warmup plus linear decay schedule, used in BERT and many other models, combines linear warmup with linear decay to zero. For total training steps T and warmup steps W :

$$\eta_t = \begin{cases} \eta_{\max} \cdot \frac{t}{W} & \text{if } t \leq W \quad (\text{warmup}) \\ \eta_{\max} \cdot \frac{T-t}{T-W} & \text{if } t > W \quad (\text{decay}) \end{cases} \quad (11.34)$$

The decay phase gradually reduces the learning rate to zero over the remaining training steps. This decay is beneficial because it allows the model to make large updates early in training when far from a good solution, then make progressively smaller updates as it approaches a good solution. The smaller learning rate in late training helps the model settle into a sharper minimum, which often generalizes better.

For BERT-base, the typical configuration is $\eta_{\max} = 1 \times 10^{-4}$, $W = 10,000$ steps, and $T = 1,000,000$ steps. This means the learning rate increases linearly from 0 to 10^{-4} over the first 10,000 steps (1% of training), then decreases linearly from 10^{-4} to 0 over the remaining 990,000 steps. The warmup period is relatively short, but it is crucial for stable training.

Different models use different maximum learning rates based on their size and architecture. GPT-2 uses $\eta_{\max} = 2.5 \times 10^{-4}$, slightly higher than BERT. GPT-3 uses $\eta_{\max} = 6 \times 10^{-5}$, lower than smaller models, reflecting the general trend that larger models require smaller learning rates for stable training. The warmup period for GPT-3 is 375 million tokens, which corresponds to a different number of steps depending on the batch size and sequence length.

11.4.3 Inverse Square Root Decay

The original "Attention is All You Need" paper used a different schedule that combines warmup with inverse square root decay:

$$\eta_t = d_{\text{model}}^{-0.5} \cdot \min(t^{-0.5}, t \cdot W^{-1.5}) \quad (11.35)$$

This schedule has two phases. During warmup ($t \leq W$), the learning rate increases linearly:

$$\eta_t = d_{\text{model}}^{-0.5} \cdot t \cdot W^{-1.5} = d_{\text{model}}^{-0.5} \cdot W^{-0.5} \cdot \frac{t}{W} \quad (11.36)$$

After warmup ($t > W$), the learning rate decays as the inverse square root of the step number:

$$\eta_t = d_{\text{model}}^{-0.5} \cdot t^{-0.5} \quad (11.37)$$

The inverse square root decay is slower than linear decay, maintaining a higher learning rate for longer. This can be beneficial for very long training runs where continued exploration is desirable. The original Transformer used $W = 4,000$ warmup steps and $d_{\text{model}} = 512$, giving a peak learning rate of $512^{-0.5} \cdot 4000^{-0.5} \approx 0.00070$.

The inverse square root schedule is less commonly used than linear decay in modern transformers, but it remains popular for some applications, particularly in machine translation where the original Transformer architecture is still widely used.

11.4.4 Cosine Annealing

Cosine annealing provides a smooth decay curve that starts slowly, accelerates in the middle, and slows again near the end. After warmup, the learning rate follows a cosine curve:

$$\eta_t = \eta_{\min} + \frac{1}{2}(\eta_{\max} - \eta_{\min}) \left(1 + \cos \left(\pi \frac{t - W}{T - W} \right) \right) \quad (11.38)$$

where η_{\min} is the minimum learning rate (often 0 or $0.1\eta_{\max}$). At the start of decay ($t = W$), the cosine term is $\cos(0) = 1$, giving $\eta_W = \eta_{\max}$. At the end of training ($t = T$), the cosine term is $\cos(\pi) = -1$, giving $\eta_T = \eta_{\min}$.

The smooth decay of cosine annealing can provide better final performance than linear decay, particularly for tasks where the model benefits from extended fine-tuning at low learning rates. The slower initial decay allows the model to continue exploring, while the accelerated decay in the middle helps the model converge, and the slow final decay allows careful refinement.

Cosine annealing is popular in computer vision (where it was originally developed) and has been adopted for some transformer training, particularly in vision transformers and multimodal models. However, linear decay remains more common for language models.

11.5 Mixed Precision Training

Mixed precision training is one of the most impactful optimizations for transformer training, reducing memory consumption and accelerating computation by leveraging lower-precision arithmetic. The technique uses 16-bit floating point (FP16 or BF16) for most operations while maintaining 32-bit floating point (FP32) master weights for numerical stability. This combination achieves substantial speedups on modern hardware while preserving training dynamics and final model quality.

11.5.1 FP16 Training Algorithm

Mixed precision training with FP16 maintains two copies of the model parameters: an FP16 copy used for forward and backward passes, and an FP32 master copy used for parameter updates. The algorithm proceeds as follows:

1. **Forward pass:** Convert FP32 master weights to FP16, perform all forward computations in FP16, producing FP16 activations
2. **Loss computation:** Compute loss in FP16, then scale the loss by a large factor S (typically 1024 or dynamically adjusted)
3. **Backward pass:** Compute gradients in FP16 using the scaled loss, producing FP16 gradients that are also scaled by S
4. **Gradient unscaling:** Divide FP16 gradients by S to recover the true gradient scale
5. **Gradient conversion:** Convert unscaled FP16 gradients to FP32
6. **Parameter update:** Update FP32 master weights using FP32 gradients and the optimizer
7. **Repeat:** Copy updated FP32 weights to FP16 for the next iteration

The loss scaling step is crucial for preventing gradient underflow. FP16 has a much smaller representable range than FP32: the smallest positive normal number in FP16 is approximately 6×10^{-5} , compared to 1.2×10^{-38} in FP32. Gradients in deep networks are often very small, particularly in later layers or after many training steps. Without scaling, these small gradients would underflow to zero in FP16, preventing the corresponding parameters from being updated.

By scaling the loss by a factor S before backpropagation, all gradients are also scaled by S (due to the chain rule). This shifts the gradient values into the representable range of FP16. After the backward pass, we divide by S to recover the true gradient values. The scaling and unscaling operations

are mathematically equivalent to computing gradients in FP32, but they allow the actual gradient computation to occur in FP16, leveraging faster FP16 hardware.

The scaling factor S can be fixed (typically 1024 or 2048) or dynamic. Dynamic loss scaling starts with a large scaling factor and reduces it if gradient overflow is detected (indicated by NaN or Inf values in the gradients). If training proceeds without overflow for a certain number of steps, the scaling factor is increased. This adaptive approach maximizes the use of FP16's range while preventing overflow.

11.5.2 Memory Savings

Mixed precision training reduces memory consumption primarily through smaller activations. The memory breakdown for mixed precision training is:

- **FP16 parameters (forward/backward):** $2P$ bytes
- **FP32 master parameters:** $4P$ bytes
- **FP32 gradients:** $4P$ bytes
- **FP32 optimizer states (Adam):** $8P$ bytes (first and second moments)
- **FP16 activations:** $A/2$ bytes (where A is FP32 activation memory)

The total is $18P + A/2$ bytes, compared to $16P + A$ bytes for FP32 training. Surprisingly, mixed precision uses slightly more memory for parameters and optimizer states ($18P$ vs $16P$) because we maintain both FP16 and FP32 copies of the parameters. However, the activation memory is halved ($A/2$ vs A), and since activations typically dominate memory consumption, mixed precision usually provides substantial overall savings.

For BERT-base with 110 million parameters, batch size 32, and sequence length 512:

FP32 training:

$$\text{Parameters + gradients + optimizer: } 110\text{M} \times 16 = 1,760 \text{ MB} \quad (11.39)$$

$$\text{Activations: } \approx 12,000 \text{ MB} \quad (11.40)$$

$$\text{Total: } 13,760 \text{ MB} \approx 13.8 \text{ GB} \quad (11.41)$$

Mixed precision training:

$$\text{FP16 parameters: } 110\text{M} \times 2 = 220 \text{ MB} \quad (11.42)$$

$$\text{FP32 master + gradients + optimizer: } 110\text{M} \times 16 = 1,760 \text{ MB} \quad (11.43)$$

$$\text{FP16 activations: } \approx 6,000 \text{ MB} \quad (11.44)$$

$$\text{Total: } 7,980 \text{ MB} \approx 8.0 \text{ GB} \quad (11.45)$$

Mixed precision saves $13.8 - 8.0 = 5.8$ GB, a 42% reduction. This memory saving enables larger batch sizes or longer sequences on the same hardware, directly improving training efficiency.

11.5.3 Hardware Acceleration

Modern GPUs provide dedicated hardware for accelerated FP16 computation. NVIDIA's Tensor Cores, available on Volta (V100), Turing (RTX 20xx), Ampere (A100, RTX 30xx), and newer architectures, can perform FP16 matrix multiplications at twice the throughput of FP32 operations.

For the NVIDIA A100 GPU:

- **FP32 performance:** 156 TFLOPS (teraflops)
- **FP16 performance (Tensor Cores):** 312 TFLOPS
- **Theoretical speedup:** $2\times$

In practice, the speedup is typically $1.5\text{--}1.8\times$ rather than the full $2\times$ because:

- Not all operations benefit from FP16 (e.g., layer normalization, softmax, and other element-wise operations may still run in FP32 for numerical stability)
- Memory bandwidth limitations can bottleneck performance, particularly for small batch sizes
- Overhead from data type conversions and loss scaling
- Non-matrix operations (activations, normalizations) don't use Tensor Cores

For BERT-base training on an A100 GPU, mixed precision typically provides a $1.6\times$ speedup, reducing training time from approximately 4 days to 2.5 days on the same hardware. This speedup, combined with the memory savings that enable larger batch sizes, makes mixed precision training essential for efficient transformer training.

11.5.4 BF16: An Alternative to FP16

BF16 (bfloat16) is an alternative 16-bit format that maintains the same exponent range as FP32 (8 bits) while reducing the mantissa precision (7 bits, compared to 10 bits in FP16). This design choice provides better numerical stability than FP16 at the cost of slightly lower precision.

The key advantage of BF16 is that it can represent the same range of values as FP32, from approximately 10^{-38} to 10^{38} . This eliminates the need for loss scaling because gradients are unlikely to underflow in BF16's range. The training algorithm simplifies to:

1. Forward pass in BF16
2. Loss computation in BF16 (no scaling needed)
3. Backward pass in BF16
4. Convert BF16 gradients to FP32
5. Update FP32 master weights

BF16 is supported on Google's TPUs (v2, v3, v4), NVIDIA A100 GPUs, and newer hardware. For transformers, BF16 often provides similar or slightly better results than FP16 with less tuning required, since the loss scaling factor doesn't need to be adjusted. However, FP16 remains more widely supported across different hardware platforms.

The memory savings and computational speedups for BF16 are similar to FP16: activations are halved, and Tensor Cores provide approximately $2\times$ theoretical speedup ($1.5\text{--}1.8\times$ in practice). The choice between FP16 and BF16 often depends on hardware availability and whether loss scaling tuning is problematic for a particular training setup.

11.6 Gradient Accumulation

Gradient accumulation is a technique for achieving large effective batch sizes when GPU memory limits the actual batch size that can be processed in a single forward-backward pass. The technique accumulates gradients over multiple mini-batches before updating parameters, mathematically equivalent to training with a larger batch but with lower memory requirements.

11.6.1 Algorithm and Implementation

The gradient accumulation algorithm processes K mini-batches of size B_{mini} , accumulating their gradients, then performs a single parameter update. The effective batch size is $B_{\text{eff}} = K \times B_{\text{mini}}$.

[H] Gradient Accumulation **Input:** Mini-batch size B_{mini} , accumulation steps K , dataset

[1] Initialize model parameters \mathbf{w} Initialize optimizer each epoch each batch of K mini-batches
 Zero accumulated gradients: $\mathbf{g}_{\text{accum}} = \mathbf{0}$ $k = 1$ to K Load mini-batch \mathcal{B}_k of size B_{mini} Forward pass: compute loss L_k on \mathcal{B}_k Scale loss: $L_k \leftarrow L_k/K$ Backward pass: compute gradients \mathbf{g}_k Accumulate: $\mathbf{g}_{\text{accum}} \leftarrow \mathbf{g}_{\text{accum}} + \mathbf{g}_k$ Update parameters using $\mathbf{g}_{\text{accum}}$ Zero gradients for next accumulation

The loss scaling by $1/K$ ensures that the accumulated gradient has the correct magnitude. Without this scaling, the accumulated gradient would be K times larger than the gradient from a single batch of size B_{eff} , leading to overly aggressive parameter updates.

In PyTorch, gradient accumulation is implemented by simply not calling `optimizer.zero_grad()` after each mini-batch. Gradients accumulate automatically because PyTorch adds new gradients to existing gradients by default:

```
optimizer.zero_grad()
for k in range(accumulation_steps):
    batch = next(dataloader)
    loss = model(batch) / accumulation_steps
    loss.backward() # Accumulates gradients

optimizer.step() # Update parameters
```

11.6.2 Trade-offs and Considerations

Gradient accumulation is mathematically equivalent to training with a larger batch size, but it has different computational characteristics. The key trade-offs are:

Memory: Gradient accumulation requires only the memory for a single mini-batch of size B_{mini} , not the full effective batch size B_{eff} . This is the primary benefit—it enables training with large effective batch sizes on memory-constrained hardware.

Computation time: Gradient accumulation is slower than true large-batch training because the mini-batches are processed sequentially rather than in parallel. For K accumulation steps, we perform K forward passes and K backward passes before a single parameter update. If we could fit the full batch in memory, we would perform 1 forward pass and 1 backward pass, processing K times more data in parallel.

The time overhead is typically 10-20% compared to true large-batch training, arising from:

- Reduced parallelism: processing mini-batches sequentially rather than in parallel
- Increased overhead: K forward-backward passes have more overhead than 1 pass
- Memory bandwidth: loading model parameters K times rather than once

Batch normalization incompatibility: Gradient accumulation is incompatible with batch normalization because batch normalization computes statistics over the mini-batch, not the effective batch. Each mini-batch has different statistics, leading to incorrect normalization. Fortunately, transformers use layer normalization rather than batch normalization, so this is not a concern for transformer training.

11.6.3 Practical Example

Consider training BERT-base where we want an effective batch size of 512, but GPU memory only allows batch size 32. We use gradient accumulation with $K = 512/32 = 16$ steps.

Memory requirements:

- Without accumulation (batch 512): ≈ 220 GB (exceeds any single GPU)
- With accumulation (batch 32): ≈ 13.8 GB (fits on V100 16GB)

Training time comparison:

- True batch 512 (if it fit): 1 forward + 1 backward = 2 passes
- Gradient accumulation: 16 forward + 16 backward = 32 passes

The gradient accumulation approach requires $16\times$ more passes, but each pass is faster because it processes less data. The total time is approximately 15% longer than true batch 512 would be, but it's feasible on available hardware.

When to use gradient accumulation:

- When the desired batch size exceeds GPU memory capacity
- When trying to match published training recipes that use large batches
- When larger batches improve convergence (common for transformers)
- When training time is less critical than achieving good final performance

When not to use gradient accumulation:

- When the mini-batch size is already optimal for convergence
- When training time is critical and larger batches don't improve convergence
- When the overhead (15-20%) is unacceptable

For BERT-base, gradient accumulation is commonly used to achieve effective batch sizes of 256-512, which provide better convergence than smaller batches. The time overhead is acceptable given the improved final performance.

11.7 Gradient Checkpointing

Gradient checkpointing, also called activation checkpointing, is a memory-computation trade-off technique that dramatically reduces activation memory at the cost of increased training time. Instead of storing all intermediate activations during the forward pass for use in backpropagation, gradient checkpointing stores only a subset of activations (typically at layer boundaries) and recomputes the remaining activations during the backward pass as needed.

11.7.1 The Memory-Computation Trade-off

Standard backpropagation requires storing all intermediate activations from the forward pass because computing gradients requires both the gradients flowing backward and the activations from the forward pass. For a transformer with L layers, batch size B , and sequence length n , the activation memory scales as $O(LBnd_{\text{model}})$ for linear terms and $O(LBhn^2)$ for attention matrices. As analyzed in Chapter 12, this activation memory often dominates total memory consumption, particularly for large batch sizes or long sequences.

Gradient checkpointing reduces activation memory by storing only activations at layer boundaries (the input to each transformer layer) and discarding all intermediate activations within layers. During the backward pass, when gradients need to flow through a layer, the forward computation for that layer is re-executed to reconstruct the intermediate activations needed for gradient computation. This recomputation happens on-the-fly during backpropagation, so the intermediate activations are used immediately and then discarded.

The memory savings are substantial. Without checkpointing, we store activations for every operation: QKV projections, attention scores, attention outputs, FFN intermediate activations, layer norm outputs, and residual connections. With checkpointing, we store only the layer inputs. For a typical transformer layer, this reduces activation memory by approximately 80%, storing only 1-2 tensors per layer instead of 8-10 tensors.

The computational cost is the price for these memory savings. Each layer's forward computation must be executed twice: once during the forward pass (with activations discarded) and once during the backward pass (to reconstruct activations for gradient computation). This doubles the forward computation cost, but the backward pass cost remains the same. Since the backward pass already costs approximately $2\times$ the forward pass, the total cost increases from $3\times$ to $4\times$ the forward pass, a 33% increase in training time. In practice, the overhead is typically 20-30% due to optimizations and the fact that some operations (like attention softmax) are relatively cheap to recompute.

11.7.2 Implementation Strategies

The most common checkpointing strategy is to checkpoint at transformer layer boundaries. For a model with L layers, we store $L + 1$ activation tensors (the input to each layer plus the final output), rather than storing all intermediate activations within layers.

In PyTorch, gradient checkpointing is implemented using `torch.utils.checkpoint.checkpoint`, which wraps a function and handles the recomputation automatically:

```
from torch.utils.checkpoint import checkpoint

class TransformerLayer(nn.Module):
    def forward(self, x):
        # Use checkpointing for this layer
        return checkpoint(self._forward, x)

    def _forward(self, x):
        # Actual layer computation
        # Attention
        attn_out = self.attention(x)
        x = x + self.dropout(attn_out)
        x = self.layer_norm1(x)

        # Feed-forward
        ffn_out = self.ffn(x)
        x = x + self.dropout(ffn_out)
        x = self.layer_norm2(x)

        return x
```

During the forward pass, PyTorch executes `_forward` but doesn't store intermediate activations. During the backward pass, when gradients reach this layer, PyTorch re-executes `_forward` with the saved input `x`, reconstructing the intermediate activations needed for gradient computation.

An alternative strategy is selective checkpointing, where only some layers are checkpointed. This provides a middle ground between memory and computation. For example, checkpointing every other layer reduces activation memory by approximately 50% while increasing training time by only 10-15%. This can be optimal when memory is tight but not critically constrained.

11.7.3 Practical Impact

The impact of gradient checkpointing is best illustrated with concrete examples. For GPT-2 (small) with 12 layers, $d_{\text{model}} = 768$, sequence length 1024, and batch size 32:

Without checkpointing:

$$\text{Activation memory per layer: } \approx 85 \text{ MB} \quad (11.46)$$

$$\text{Total activation memory (12 layers): } \approx 1,020 \text{ MB} \approx 1 \text{ GB per sequence} \quad (11.47)$$

$$\text{Batch size 32: } 32 \text{ GB} \quad (11.48)$$

This exceeds the memory of most GPUs when combined with parameters and optimizer states.

With checkpointing:

$$\text{Stored activations (layer inputs only): } 13 \times 32 \times 1024 \times 768 \times 4 \approx 1,308 \text{ MB} \quad (11.49)$$

$$\text{Reduction: } 32,000 \text{ MB} \rightarrow 1,308 \text{ MB} \quad (96\% \text{reduction!}) \quad (11.50)$$

This dramatic reduction enables training with much larger batch sizes or longer sequences on the same hardware. For GPT-2 on an NVIDIA V100 (16 GB), checkpointing enables increasing the batch size from approximately 4 to 20, a $5\times$ improvement.

Training time impact:

- Without checkpointing: 100% (baseline)
- With checkpointing: 125% (25% slower)

The 25% time increase is usually acceptable given the $5\times$ increase in batch size, which often improves convergence and reduces the total number of steps needed for training.

11.7.4 When to Use Gradient Checkpointing

Gradient checkpointing is most beneficial in specific scenarios:

Use checkpointing when:

- Training with long sequences (e.g., $n > 1024$) where activation memory dominates
- GPU memory is the limiting factor preventing larger batch sizes
- The model is very deep (many layers) and activation memory scales linearly with depth
- Training time is less critical than maximizing batch size or sequence length
- Combined with mixed precision, checkpointing enables training that would otherwise be impossible

Avoid checkpointing when:

- Memory is not constrained and the 20-30% time overhead is unacceptable
- Training with short sequences and small batch sizes where activation memory is already manageable
- Optimizing for minimum training time rather than maximum throughput
- The model is shallow enough that activation memory is not the bottleneck

For most transformer training, particularly for models with more than 12 layers or sequences longer than 512 tokens, gradient checkpointing is beneficial. The memory savings enable configurations that would otherwise be impossible, and the time overhead is modest compared to the benefits.

11.8 Distributed Training Strategies

As transformer models grow beyond the capacity of single GPUs, distributed training becomes essential. Different distributed training strategies partition the model, data, or optimizer state across multiple GPUs, each with distinct trade-offs in terms of memory reduction, communication overhead, and implementation complexity. Understanding these strategies is crucial for training large-scale models efficiently.

11.8.1 Data Parallelism

Data parallelism is the simplest and most widely used distributed training strategy. The model is replicated on each GPU, and each GPU processes a different subset of the training batch. After computing gradients locally, the GPUs synchronize their gradients using an AllReduce operation, then each GPU updates its local copy of the model with the averaged gradients.

The algorithm proceeds as follows:

1. Each GPU has a complete copy of the model
2. The global batch is split across GPUs: GPU i processes mini-batch \mathcal{B}_i
3. Each GPU performs forward and backward passes independently, computing local gradients \mathbf{g}_i

4. AllReduce operation computes the average gradient: $\bar{\mathbf{g}} = \frac{1}{N} \sum_{i=1}^N \mathbf{g}_i$ where N is the number of GPUs
5. Each GPU updates its model using $\bar{\mathbf{g}}$
6. All GPUs now have identical models (up to floating-point precision)

Data parallelism scales efficiently to 8-16 GPUs on a single node (connected via NVLink or PCIe) because the communication overhead is relatively small compared to computation. For BERT-base with 110M parameters, the AllReduce operation must communicate $110\text{M} \times 4 = 440$ MB of gradients. On NVLink (300 GB/s bandwidth), this takes approximately $440 \text{ MB} / 300 \text{ GB/s} \approx 1.5$ ms, which is small compared to the forward-backward computation time of 10-20 ms per batch.

However, data parallelism does not reduce memory requirements per GPU—each GPU still stores the complete model, optimizer states, and activations for its mini-batch. This limits the size of models that can be trained with data parallelism alone. For GPT-3 with 175B parameters requiring 700 GB in FP32, data parallelism is insufficient because no single GPU has enough memory for the complete model.

11.8.2 Model Parallelism

Model parallelism splits the model across multiple GPUs, with different layers residing on different devices. For a model with L layers split across N GPUs, each GPU stores approximately L/N layers. This reduces per-GPU memory proportionally to the number of GPUs.

The forward pass proceeds sequentially: GPU 1 processes the input through its layers, sends activations to GPU 2, which processes through its layers, and so on. The backward pass proceeds in reverse: GPU N computes gradients for its layers, sends gradients to GPU $N - 1$, which computes gradients for its layers, and so on.

The primary challenge with model parallelism is the pipeline bubble problem. While GPU 1 is processing the next batch, GPUs 2 through N are idle, waiting for activations from GPU 1. Similarly, during the backward pass, GPU N finishes first and sits idle while earlier GPUs complete their backward passes. This sequential execution leads to poor GPU utilization, with each GPU active only $1/N$ of the time in the worst case.

Model parallelism is necessary when a single layer or the complete model exceeds single-GPU memory, but it should be combined with other strategies to improve utilization. For GPT-3, model parallelism alone would require hundreds of GPUs and would have terrible utilization due to pipeline bubbles.

11.8.3 Pipeline Parallelism

Pipeline parallelism improves upon model parallelism by splitting each batch into micro-batches and pipelining their execution across GPUs. Instead of processing one batch completely before starting the next, pipeline parallelism processes multiple micro-batches concurrently, with different micro-batches at different stages of the pipeline.

For example, with 4 GPUs and 4 micro-batches:

- Time 1: GPU 1 processes micro-batch 1 (forward)
- Time 2: GPU 1 processes micro-batch 2 (forward), GPU 2 processes micro-batch 1 (forward)
- Time 3: GPU 1 processes micro-batch 3 (forward), GPU 2 processes micro-batch 2 (forward), GPU 3 processes micro-batch 1 (forward)
- Time 4: All GPUs are active, processing different micro-batches

This pipelining significantly reduces idle time. The pipeline bubble (time when some GPUs are idle) is proportional to the number of GPUs divided by the number of micro-batches. With N GPUs and M micro-batches, the bubble fraction is approximately N/M . Using $M = 4N$ micro-batches reduces the bubble to 25%, achieving 75% utilization.

Pipeline parallelism implementations like GPipe and PipeDream differ in how they handle gradient computation and weight updates. GPipe uses synchronous updates, accumulating gradients from all micro-batches before updating weights. PipeDream uses asynchronous updates, updating weights after each micro-batch, which can improve throughput but requires careful handling of weight versions.

11.8.4 Tensor Parallelism

Tensor parallelism, pioneered by Megatron-LM, splits individual layers across multiple GPUs rather than splitting the model layer-wise. For attention and feed-forward layers, the computation can be partitioned across GPUs with minimal communication.

For the attention mechanism, the heads can be split across GPUs. With h heads and N GPUs, each GPU computes h/N heads independently. The only communication required is an AllReduce after computing the attention output, to sum the contributions from all heads.

For the feed-forward network, the first linear layer $\mathbf{W}_1 \in \mathbb{R}^{d_{\text{model}} \times d_{ff}}$ can be column-partitioned across GPUs. Each GPU computes a subset of the d_{ff} intermediate activations. The GELU activation is applied independently on each GPU. The second linear layer $\mathbf{W}_2 \in \mathbb{R}^{d_{ff} \times d_{\text{model}}}$ is row-partitioned, and an AllReduce sums the outputs from all GPUs.

Tensor parallelism achieves $N\times$ memory reduction with only two AllReduce operations per layer (one for attention, one for FFN). The communication volume is $O(Bnd_{\text{model}})$ per layer, which is much smaller than the $O(P)$ communication required for data parallelism (where P is the number of parameters).

Tensor parallelism is particularly effective for very large layers. For GPT-3 with $d_{\text{model}} = 12,288$ and $d_{ff} = 49,152$, a single FFN layer has $2 \times 12,288 \times 49,152 \approx 1.2\text{B}$ parameters, requiring 4.8 GB in FP32. Splitting across 8 GPUs reduces this to 600 MB per GPU, making the layer tractable.

11.8.5 ZeRO: Zero Redundancy Optimizer

ZeRO (Zero Redundancy Optimizer) is a family of optimizations that reduce memory by sharding optimizer states, gradients, and parameters across GPUs while maintaining the computational efficiency of data parallelism. ZeRO has three stages, each providing progressively more memory reduction:

ZeRO Stage 1: Optimizer State Partitioning

Each GPU stores only $1/N$ of the optimizer states (first and second moments for Adam). During the optimizer step, each GPU updates only its partition of the parameters. This reduces optimizer memory by $N\times$ with minimal communication overhead.

For BERT-base with 110M parameters and 8 GPUs:

- Without ZeRO: Each GPU stores 880 MB of optimizer states
- With ZeRO-1: Each GPU stores $880/8 = 110$ MB of optimizer states
- Memory saved: 770 MB per GPU

ZeRO Stage 2: Gradient Partitioning

In addition to optimizer states, gradients are also partitioned. Each GPU computes gradients for all parameters during backpropagation but only retains the gradients for its partition, discarding the rest. This reduces gradient memory by $N\times$.

For BERT-base with 8 GPUs:

- Without ZeRO: Each GPU stores 440 MB of gradients
- With ZeRO-2: Each GPU stores $440/8 = 55$ MB of gradients
- Total memory saved: $770 + 385 = 1,155$ MB per GPU

ZeRO Stage 3: Parameter Partitioning

The most aggressive stage partitions the parameters themselves. Each GPU stores only $1/N$ of the parameters. During the forward pass, each GPU gathers the parameters it needs from other GPUs,

computes its portion of the forward pass, then discards the gathered parameters. The backward pass proceeds similarly.

For BERT-base with 8 GPUs:

- Without ZeRO: Each GPU stores 440 MB of parameters
- With ZeRO-3: Each GPU stores $440/8 = 55$ MB of parameters
- Total memory saved: $770 + 385 + 385 = 1,540$ MB per GPU

ZeRO-3 enables training models that wouldn't fit on any single GPU by distributing all memory across the cluster. For GPT-3 with 175B parameters requiring 700 GB in FP32, ZeRO-3 across 64 A100 GPUs (80 GB each) reduces per-GPU memory to $700/64 \approx 11$ GB, making training feasible.

The communication overhead of ZeRO increases with each stage. ZeRO-1 has minimal overhead (only during optimizer step). ZeRO-2 adds gradient communication (similar to data parallelism). ZeRO-3 adds parameter communication during forward and backward passes, which can be significant but is often acceptable given the memory savings.

11.8.6 Comparison of Strategies

Strategy	Memory Reduction	Communication	Use Case
Data Parallel	None	Gradients	Small models, many GPUs
Model Parallel	$N \times$	Activations	Large models, sequential
Pipeline Parallel	$N \times$	Activations	Very large models
Tensor Parallel	$N \times$	Activations (small)	Huge layers
ZeRO Stage 1	$4 \times$	Minimal	Optimizer memory bound
ZeRO Stage 2	$8 \times$	Gradients	Gradient memory bound
ZeRO Stage 3	$N \times$	All	Extreme scale

Table 11.1: Comparison of distributed training strategies for N GPUs

In practice, large-scale training often combines multiple strategies. GPT-3 training used a combination of data parallelism, model parallelism, and pipeline parallelism across thousands of GPUs. Modern frameworks like DeepSpeed and Megatron-LM provide implementations of these strategies that can be combined flexibly based on model size and available hardware.

11.9 Batch Size and Sequence Length Selection

Selecting appropriate batch sizes and sequence lengths is crucial for efficient transformer training. These choices directly impact memory consumption, training throughput, convergence behavior, and final model quality. The optimal configuration depends on the interplay between hardware constraints, model architecture, and training objectives.

11.9.1 Batch Size Considerations

Batch size affects both computational efficiency and optimization dynamics. Larger batches improve GPU utilization by amortizing the cost of loading model parameters and by providing more parallelism for matrix operations. Modern GPUs achieve peak performance with large matrix multiplications, and larger batches create larger matrices that better utilize the hardware.

For BERT-base on an NVIDIA A100, throughput (tokens processed per second) increases significantly with batch size:

- Batch size 8: $\approx 15,000$ tokens/sec (30% GPU utilization)
- Batch size 32: $\approx 50,000$ tokens/sec (80% GPU utilization)
- Batch size 64: $\approx 70,000$ tokens/sec (90% GPU utilization)

- Batch size 128: $\approx 75,000$ tokens/sec (95% GPU utilization)

Beyond batch size 64, the throughput gains diminish because the GPU is already well-utilized. The optimal batch size for throughput is typically where GPU utilization reaches 85-95%, which depends on the model size and sequence length.

However, larger batches are not always better for optimization. Very large batches can hurt generalization, a phenomenon known as the "generalization gap." The intuition is that large batches provide very accurate gradient estimates, which can lead the optimizer to sharp minima that don't generalize well. Smaller batches provide noisier gradients that help the optimizer find flatter minima with better generalization.

The relationship between batch size and generalization is complex and depends on the learning rate schedule and total training budget. Research has shown that the generalization gap can be mitigated by:

- Scaling the learning rate proportionally with batch size (linear scaling rule)
- Extending the warmup period for larger batches
- Training for more steps to compensate for fewer parameter updates

For transformer training, batch sizes of 256-2048 are typical. BERT-base uses an effective batch size of 256 (32 per GPU \times 8 GPUs). GPT-2 uses batch sizes of 512-1024. GPT-3 uses batch sizes up to 3.2 million tokens (approximately 1600 sequences of length 2048), enabled by LAMB optimizer and massive parallelism.

11.9.2 Memory Scaling with Batch Size

Memory consumption scales linearly with batch size for most components. For BERT-base with sequence length 512:

$$\text{Batch size 8: } \approx 3.5 \text{ GB} \quad (11.51)$$

$$\text{Batch size 16: } \approx 6.8 \text{ GB} \quad (11.52)$$

$$\text{Batch size 32: } \approx 13.8 \text{ GB} \quad (11.53)$$

$$\text{Batch size 64: } \approx 27.6 \text{ GB} \quad (11.54)$$

The linear scaling means that doubling the batch size doubles the memory requirement. This quickly exceeds single-GPU capacity, necessitating either gradient accumulation (to simulate large batches with small physical batches) or distributed training (to split the batch across multiple GPUs).

The memory breakdown for batch size 32 is approximately:

- Parameters + optimizer: 1.76 GB (independent of batch size)
- Activations: 12 GB (scales linearly with batch size)

Since activations dominate, techniques that reduce activation memory (mixed precision, gradient checkpointing) have a large impact on the maximum feasible batch size.

11.9.3 Sequence Length Considerations

Sequence length has a more complex impact on memory and computation than batch size. The attention mechanism's quadratic scaling means that memory and computation grow as $O(n^2)$ for sequence length n , while other components grow linearly as $O(n)$.

For BERT-base with batch size 32, memory consumption varies dramatically with sequence length:

$$\text{Sequence length 128: } \approx 3.5 \text{ GB} \quad (11.55)$$

$$\text{Sequence length 256: } \approx 6.2 \text{ GB} \quad (11.56)$$

$$\text{Sequence length 512: } \approx 13.8 \text{ GB} \quad (11.57)$$

$$\text{Sequence length 1024: } \approx 42 \text{ GB} \quad (11.58)$$

Doubling the sequence length from 512 to 1024 roughly triples the memory (not quadruples, because some components scale linearly). The attention matrices grow quadratically: for 12 heads, the attention memory is $32 \times 12 \times n^2 \times 4$ bytes. At $n = 512$, this is 403 MB; at $n = 1024$, this is 1.6 GB; at $n = 2048$, this is 6.4 GB.

The quadratic scaling limits practical sequence lengths. BERT uses $n = 512$, GPT-2 uses $n = 1024$, GPT-3 uses $n = 2048$. Longer sequences require either:

- Efficient attention mechanisms (sparse attention, linear attention) that reduce the $O(n^2)$ complexity
- Gradient checkpointing to reduce activation memory
- Smaller batch sizes to fit within memory constraints
- More powerful GPUs with larger memory

The choice of sequence length depends on the task. For tasks requiring long-range dependencies (document classification, long-form generation), longer sequences are beneficial despite the computational cost. For tasks with local dependencies (named entity recognition, part-of-speech tagging), shorter sequences may suffice.

11.9.4 Dynamic Batching

Dynamic batching groups sequences of similar length together to minimize padding waste. In a typical batch, sequences have varying lengths, and all sequences are padded to the length of the longest sequence in the batch. This padding wastes computation and memory on padding tokens that don't contribute to learning.

For example, if a batch contains sequences of lengths [128, 256, 512, 512], all sequences are padded to 512, wasting:

$$(512 - 128) + (512 - 256) + 0 + 0 = 640 \text{ tokens} \quad (11.59)$$

Out of $4 \times 512 = 2048$ total tokens, 640 (31%) are padding.

Dynamic batching sorts sequences by length and groups similar lengths together. This reduces padding significantly. If we instead batch [128, 128, 128, 128] and [512, 512, 512, 512] separately, there's no padding waste within each batch.

The throughput improvement from dynamic batching can be substantial:

- Without dynamic batching: 50,000 tokens/sec (including padding)
- With dynamic batching: 70,000 tokens/sec (40% improvement)

The improvement depends on the length distribution in the dataset. For datasets with highly variable lengths, dynamic batching can provide 2-3 \times throughput improvements. For datasets with uniform lengths, the benefit is minimal.

Dynamic batching is implemented by sorting the dataset by sequence length before creating batches, or by using a bucketing strategy that assigns sequences to length buckets and samples batches from within buckets. Most modern training frameworks (Hugging Face Transformers, fairseq) support dynamic batching.

11.9.5 Practical Guidelines

Based on the analysis above, practical guidelines for batch size and sequence length selection are:

For batch size:

- Start with the largest batch size that fits in GPU memory
- If memory-constrained, use gradient accumulation to achieve larger effective batch sizes
- For BERT-base on V100 (16 GB): batch size 16-32 with sequence length 512
- For BERT-base on A100 (40 GB): batch size 32-64 with sequence length 512
- Scale learning rate proportionally when increasing batch size
- Extend warmup period for very large batches (≥ 1024)

For sequence length:

- Use the longest sequence length that fits in memory and is relevant for the task
- For memory-constrained scenarios, reduce batch size rather than sequence length if long context is important
- Use gradient checkpointing to enable longer sequences
- Consider efficient attention mechanisms for sequences longer than 2048
- Use dynamic batching to reduce padding waste

Memory-constrained optimization:

1. Enable mixed precision training (FP16/BF16): 40-50% memory reduction
2. Enable gradient checkpointing: 80% activation memory reduction
3. Use gradient accumulation: simulate large batches with small physical batches
4. Reduce sequence length if task permits
5. Use dynamic batching to reduce padding waste

These techniques can be combined. For example, BERT-base with mixed precision + gradient checkpointing can train with batch size 128 and sequence length 512 on a V100 (16 GB), compared to batch size 16 without these optimizations.

11.10 Regularization Techniques

Regularization prevents overfitting by constraining the model's capacity or adding noise during training. Transformers, with their large parameter counts, are particularly susceptible to overfitting on small datasets. Effective regularization enables transformers to generalize well from training data to unseen examples.

11.10.1 Dropout

Dropout randomly sets activations to zero during training with probability p , forcing the model to learn robust features that don't rely on any single activation. During inference, dropout is disabled, and activations are scaled by $(1 - p)$ to maintain the expected magnitude.

In transformers, dropout is applied at multiple locations:

Attention dropout: Applied to the attention weights after softmax, before multiplying by values:

$$\mathbf{O} = \text{Dropout}(\text{softmax}(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}))\mathbf{V} \quad (11.60)$$

This prevents the model from relying too heavily on specific attention patterns, encouraging it to learn diverse attention strategies.

Residual dropout: Applied to the output of each sub-layer before adding to the residual connection:

$$\mathbf{y} = \mathbf{x} + \text{Dropout}(\text{Sublayer}(\mathbf{x})) \quad (11.61)$$

This regularizes the transformations learned by attention and feed-forward layers.

Embedding dropout: Applied to the sum of token embeddings and positional encodings:

$$\mathbf{x} = \text{Dropout}(\text{TokenEmbed}(x) + \text{PositionalEncoding}(x)) \quad (11.62)$$

This prevents overfitting to specific token representations.

Typical dropout rates for transformers are relatively low compared to other architectures. BERT uses $p = 0.1$ (10% dropout) for all dropout locations. GPT-2 also uses $p = 0.1$. Larger models sometimes use even lower dropout rates ($p = 0.05$ or less) because their increased capacity provides implicit regularization.

The dropout rate should be tuned based on the dataset size and model capacity. For small datasets (thousands of examples), higher dropout rates ($p = 0.2$ or $p = 0.3$) may be beneficial. For large datasets (millions of examples), lower dropout rates ($p = 0.1$ or less) are typically sufficient.

11.10.2 Weight Decay

Weight decay adds an L2 penalty to the loss function, encouraging parameters to remain small. In the context of AdamW (the standard optimizer for transformers), weight decay is applied directly to parameters rather than through the gradient:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \frac{\hat{\mathbf{m}}_t}{\sqrt{\hat{\mathbf{v}}_t + \epsilon}} - \eta \lambda \mathbf{w}_t \quad (11.63)$$

The weight decay coefficient λ controls the strength of regularization. Typical values for transformer training are $\lambda = 0.01$ or $\lambda = 0.001$. BERT uses $\lambda = 0.01$, which provides moderate regularization without overly constraining the model.

Weight decay is not applied uniformly to all parameters. Biases and layer normalization parameters (scale γ and shift β) are typically excluded from weight decay. The reasoning is that these parameters control the scale and offset of activations rather than the complexity of learned features, and regularizing them can hurt performance. In practice, this exclusion is implemented by creating separate parameter groups in the optimizer with different weight decay settings.

The interaction between weight decay and learning rate is important. Because weight decay is applied with coefficient $\eta\lambda$, the effective regularization strength increases with the learning rate. During warmup, when the learning rate is small, weight decay has minimal effect. As the learning rate increases, weight decay becomes stronger. During decay, as the learning rate decreases, weight decay weakens. This dynamic regularization schedule often works well in practice.

11.10.3 Label Smoothing

Label smoothing replaces hard one-hot targets with soft targets that assign small probabilities to incorrect classes. For a classification problem with vocabulary size V and true class y , the smoothed target distribution is:

$$q(k) = \begin{cases} 1 - \epsilon + \frac{\epsilon}{V} & \text{if } k = y \\ \frac{\epsilon}{V} & \text{if } k \neq y \end{cases} \quad (11.64)$$

where ϵ is the smoothing parameter, typically $\epsilon = 0.1$.

Label smoothing prevents the model from becoming overconfident in its predictions. Without smoothing, the model is trained to assign probability 1 to the correct class and probability 0 to all other classes. This can lead to overconfident predictions that don't reflect the model's true uncertainty. With smoothing, the model is trained to assign high probability to the correct class but also small probabilities to other classes, leading to better-calibrated predictions.

For language modeling with vocabulary size 30,000 and $\epsilon = 0.1$:

$$\text{Correct class: } q(y) = 1 - 0.1 + \frac{0.1}{30000} = 0.900003 \quad (11.65)$$

$$\text{Incorrect classes: } q(k) = \frac{0.1}{30000} = 0.0000033 \quad (11.66)$$

The smoothed target assigns 90% probability to the correct class and distributes the remaining 10% uniformly across all classes.

Label smoothing is particularly beneficial for tasks with ambiguous labels or where multiple outputs could be considered correct. In machine translation, for example, multiple translations may be valid, and label smoothing encourages the model to consider alternatives rather than committing entirely to the reference translation.

The cross-entropy loss with label smoothing is:

$$L = - \sum_{k=1}^V q(k) \log p(k) = -(1 - \epsilon) \log p(y) - \frac{\epsilon}{V} \sum_{k=1}^V \log p(k) \quad (11.67)$$

The second term is the negative entropy of the predicted distribution, which encourages the model to maintain some uncertainty rather than collapsing to a single prediction.

11.10.4 Gradient Clipping

Gradient clipping prevents exploding gradients by limiting the norm of the gradient vector. If the gradient norm exceeds a threshold θ , the gradient is scaled down:

$$\mathbf{g} \leftarrow \begin{cases} \mathbf{g} & \text{if } \|\mathbf{g}\|_2 \leq \theta \\ \frac{\theta \mathbf{g}}{\|\mathbf{g}\|_2} & \text{if } \|\mathbf{g}\|_2 > \theta \end{cases} \quad (11.68)$$

The typical threshold for transformer training is $\theta = 1.0$. This value is chosen empirically and works well across different model sizes and tasks.

Gradient clipping is essential for training stability, particularly in the early stages of training when gradients can be very large. Without clipping, occasional large gradients can cause the parameters to jump to regions of the loss landscape with poor gradients, derailing training. With clipping, these large gradients are tamed, allowing training to proceed smoothly.

The clipping threshold should be tuned based on the typical gradient norms observed during training. If gradients are frequently clipped, the threshold may be too low, preventing the model from making necessary large updates. If gradients are rarely clipped, the threshold may be too high, providing insufficient protection against exploding gradients. Monitoring the fraction of steps where clipping occurs (typically 1-5%) helps tune the threshold.

Gradient clipping interacts with the learning rate: with a lower learning rate, gradients have less impact, so clipping is less necessary. With a higher learning rate, clipping becomes more important. The combination of learning rate warmup and gradient clipping provides robust training stability.

11.11 Training Time and Cost Estimates

Understanding the time and financial costs of training transformers is essential for planning research projects and production deployments. Training costs scale dramatically with model size, and accurate estimates help make informed decisions about model architecture, hardware selection, and training strategies.

11.11.1 BERT-base Training

BERT-base, with 110 million parameters, represents a moderately-sized transformer that can be trained on a small cluster of GPUs. The original BERT paper reported training on 16 Cloud TPU chips (equivalent to 16 TPU v2 cores), but the model can also be trained efficiently on NVIDIA GPUs.

Training configuration:

- Hardware: 16× NVIDIA V100 GPUs (16 GB each)
- Batch size: 256 per GPU, 4096 total effective batch size
- Sequence length: 512 tokens
- Training data: 3.3 billion words (approximately 16 GB of text)
- Training steps: 1 million steps
- Optimizer: AdamW with learning rate 1×10^{-4} , warmup 10,000 steps

Computational analysis:

Each training step processes $4096 \times 512 = 2,097,152$ tokens. With 1 million steps, the total training processes approximately 2.1 trillion tokens. Given that the dataset contains 3.3 billion words (approximately 4.4 billion tokens with subword tokenization), the model sees each token approximately 475 times during training.

Each training step requires approximately 290 GFLOPs (96.6 GFLOPs forward + 193.2 GFLOPs backward). With 1 million steps:

$$\text{Total compute} = 1,000,000 \times 290 \times 10^9 = 2.9 \times 10^{17} \text{ FLOPs} \quad (11.69)$$

On V100 GPUs with 125 TFLOPS (FP16 with Tensor Cores) and assuming 70% utilization:

$$\text{Time per step} = \frac{290 \times 10^9}{16 \times 125 \times 10^{12} \times 0.7} \approx 0.21 \text{ seconds} \quad (11.70)$$

Total training time:

$$1,000,000 \times 0.21 \text{ s} = 210,000 \text{ s} \approx 58 \text{ hours} \approx 2.4 \text{ days} \quad (11.71)$$

In practice, training takes approximately 3-4 days accounting for data loading, checkpointing, validation, and other overhead. The original BERT paper reported approximately 4 days of training on TPUs.

Cost estimate:

On cloud platforms (AWS, Google Cloud, Azure), V100 GPU instances cost approximately \$3-4 per GPU-hour. For 16 GPUs over 4 days:

$$\text{Cost} = 16 \times 96 \text{ hours} \times \$3.50 = \$5,376 \quad (11.72)$$

Including storage, data transfer, and other costs, the total cost is approximately \$6,000-7,000. This makes BERT-base training accessible to academic research groups and small companies.

11.11.2 GPT-2 Training

GPT-2 comes in several sizes, with the largest (GPT-2 XL) having 1.5 billion parameters. This model requires more substantial computational resources than BERT-base but remains trainable on a modest cluster.

Training configuration (GPT-2 XL):

- Parameters: 1.5 billion
- Hardware: 32× NVIDIA V100 GPUs
- Training data: 40 GB of text (WebText dataset)
- Sequence length: 1024 tokens
- Batch size: 512 total effective batch size
- Training time: Approximately 1 week

Computational analysis:

GPT-2 XL has 48 layers with $d_{\text{model}} = 1600$ and $d_{\text{ff}} = 6400$. The FLOPs per token are approximately:

$$\text{FLOPs per token} \approx 48 \times (24 \times 1024 \times 1600^2 + 4 \times 1024^2 \times 1600) \approx 6 \times 10^{12} \quad (11.73)$$

With 40 GB of text (approximately 10 billion tokens) and multiple epochs:

$$\text{Total compute} \approx 10^{10} \times 6 \times 10^{12} \times 3 = 1.8 \times 10^{23} \text{ FLOPs} \quad (11.74)$$

This is approximately 600× more compute than BERT-base, reflecting the larger model size and dataset.

Cost estimate:

With 32 V100 GPUs for 7 days:

$$\text{Cost} = 32 \times 168 \text{ hours} \times \$3.50 = \$18,816 \quad (11.75)$$

Including overhead, the total cost is approximately \$20,000-25,000. OpenAI reported spending approximately \$50,000 on compute for GPT-2, which includes experimentation, hyperparameter tuning, and multiple training runs.

11.11.3 GPT-3 Training

GPT-3, with 175 billion parameters, represents the extreme end of transformer training, requiring massive computational resources and sophisticated distributed training strategies.

Training configuration:

- Parameters: 175 billion
- Architecture: 96 layers, $d_{\text{model}} = 12,288$, $d_{\text{ff}} = 49,152$
- Training data: 300 billion tokens (approximately 570 GB of text)
- Sequence length: 2048 tokens
- Hardware: Estimated 10,000+ NVIDIA V100 GPUs (or equivalent)
- Training time: Approximately 1 month

Computational analysis:

The FLOPs per token for GPT-3 are approximately:

$$\text{FLOPs per token} \approx 96 \times (24 \times 2048 \times 12288^2 + 4 \times 2048^2 \times 12288) \approx 7 \times 10^{14} \quad (11.76)$$

With 300 billion tokens:

$$\text{Total compute} \approx 3 \times 10^{11} \times 7 \times 10^{14} \times 3 = 6.3 \times 10^{26} \text{ FLOPs} \quad (11.77)$$

This is approximately 2 million times more compute than BERT-base, illustrating the exponential scaling of training costs with model size.

Cost estimate:

The exact hardware configuration for GPT-3 training has not been publicly disclosed, but estimates suggest:

- Compute cost: \$4-12 million (depending on hardware and efficiency)
- Energy consumption: Approximately 1,287 MWh
- Carbon footprint: Approximately 552 metric tons CO₂ equivalent (depending on energy source)

These estimates are based on the reported compute of 3.14×10^{23} FLOPs (petaflop-days) and typical cloud GPU pricing. The actual cost to OpenAI may be lower due to optimized infrastructure and bulk pricing, but the order of magnitude illustrates the massive investment required for training such large models.

11.11.4 Scaling Laws

Research on scaling laws for language models has revealed predictable relationships between model size, dataset size, compute budget, and performance. These laws enable estimation of training costs for models of different sizes.

Key scaling relationships:

Compute scaling: Doubling the model size (number of parameters) requires approximately $4\times$ the compute for the same amount of training data. This quadratic scaling arises because:

- FLOPs scale linearly with parameters: $\text{FLOPs} \propto P$
- Optimal training data scales linearly with parameters: $\text{Data} \propto P$
- Total compute is $\text{FLOPs} \times \text{Data}$: $\text{Compute} \propto P^2$

Data scaling: Doubling the dataset size requires approximately $2\times$ the compute (assuming model size is fixed). This linear scaling is straightforward: processing twice as much data requires twice as many training steps.

Optimal allocation: For a fixed compute budget C , the optimal allocation between model size P and dataset size D follows:

$$P \propto C^{0.73}, \quad D \propto C^{0.27} \quad (11.78)$$

This means that as compute increases, most of the additional compute should go toward larger models rather than more data. For example, increasing compute by $10\times$ should increase model size by approximately $5.4\times$ and dataset size by approximately $1.9\times$.

These scaling laws have important implications for training strategy. For a given compute budget, training a larger model on less data often yields better performance than training a smaller model on more data. This insight has driven the trend toward ever-larger models like GPT-3, GPT-4, and beyond.

11.11.5 Cost-Performance Trade-offs

The relationship between training cost and model performance is not linear. Initial improvements are relatively cheap, but achieving state-of-the-art performance requires exponentially increasing compute.

Example progression:

- BERT-base (110M params): \$7,000, strong performance on many tasks
- BERT-large (340M params): \$25,000, 2-3% improvement over BERT-base
- GPT-2 XL (1.5B params): \$50,000, significant improvement in generation quality
- GPT-3 (175B params): \$4-12 million, state-of-the-art few-shot learning

The cost increases by 3-4 orders of magnitude from BERT-base to GPT-3, while performance improvements, though substantial, are more modest. This diminishing return on investment means that the choice of model size should be driven by the specific application requirements and available budget.

For many applications, smaller models like BERT-base or GPT-2 provide excellent performance at a fraction of the cost of the largest models. Fine-tuning these models on task-specific data often yields better results than using much larger models without fine-tuning. The trend toward efficient training methods (distillation, pruning, quantization) aims to achieve strong performance with lower training costs.

11.12 Practical Training Recipe

This section provides a comprehensive, step-by-step guide for training a transformer model, synthesizing the techniques and considerations discussed throughout the chapter. This recipe is based on best practices from training BERT, GPT-2, and other successful models, adapted for practical use.

11.12.1 Data Preparation

Effective training begins with proper data preparation. The quality and format of training data significantly impact model performance and training efficiency.

Tokenization: Use subword tokenization (BPE, WordPiece, or SentencePiece) to balance vocabulary size and representation quality. For English, a vocabulary size of 30,000-50,000 works well. Train the tokenizer on a representative sample of your data (at least 1 million sentences) to ensure good coverage. The tokenizer should handle rare words, numbers, and special characters appropriately.

Sequence packing: Combine multiple short documents into single sequences to minimize padding waste. For example, if training with sequence length 512, pack documents separated by special tokens until reaching 512 tokens. This improves efficiency by ensuring most tokens in each sequence are meaningful rather than padding.

Data augmentation: For tasks where data is limited, consider augmentation strategies like back-translation, synonym replacement, or random insertion/deletion of tokens. However, for large-scale pretraining, augmentation is typically unnecessary and may hurt performance by introducing noise.

Data filtering: Remove low-quality examples (duplicates, non-linguistic content, extremely short or long sequences) to improve training efficiency. For web-scraped data, filter by language, remove boilerplate content, and deduplicate at the document level.

11.12.2 Model Initialization

Proper initialization is crucial for training stability and convergence speed. Poor initialization can lead to vanishing or exploding gradients, slow convergence, or failure to train at all.

Weight initialization: Use Xavier (Glorot) initialization for linear layers:

$$W_{ij} \sim \mathcal{N}\left(0, \frac{2}{d_{\text{in}} + d_{\text{out}}}\right) \quad (11.79)$$

where d_{in} and d_{out} are the input and output dimensions. This initialization maintains variance across layers, preventing vanishing or exploding activations.

Bias initialization: Initialize all biases to zero. This is standard practice and works well for transformers.

Embedding initialization: Initialize token embeddings with Xavier initialization. Position embeddings can be initialized randomly or with sinusoidal patterns (as in the original Transformer). Random initialization is more common in modern models and allows the model to learn task-specific positional patterns.

Layer normalization initialization: Initialize scale parameters γ to 1 and shift parameters β to 0. This makes layer normalization initially act as the identity function, allowing gradients to flow freely in early training.

Output layer initialization: For the output projection to vocabulary logits, use Xavier initialization with a smaller scale (multiply by 0.5 or 0.1) to prevent overly confident initial predictions that can destabilize training.

11.12.3 Hyperparameter Selection

Choosing appropriate hyperparameters is critical for successful training. These recommendations are based on extensive empirical experience with transformer training.

Learning rate: Start with $\eta_{\text{max}} = 1 \times 10^{-4}$ to 3×10^{-4} for Adam/AdamW. Larger models typically require smaller learning rates. For BERT-base, use 1×10^{-4} . For GPT-2, use 2.5×10^{-4} . For models larger than 1B parameters, use 6×10^{-5} or smaller.

Warmup: Use 10% of total training steps for warmup, or at least 1,000 steps. For very large models or large batch sizes, extend warmup to 20% of steps. The warmup period should be long enough for Adam's moment estimates to stabilize.

Batch size: Use the largest batch size that fits in memory, typically 256-2048 for transformers. If memory-constrained, use gradient accumulation to achieve larger effective batch sizes. Scale the learning rate proportionally when increasing batch size beyond 256.

Weight decay: Use $\lambda = 0.01$ for AdamW. Exclude biases and layer normalization parameters from weight decay. This moderate regularization prevents overfitting without overly constraining the model.

Dropout: Use $p = 0.1$ for all dropout locations (attention, residual, embedding). For small datasets, increase to $p = 0.2$ or $p = 0.3$. For very large models (≥ 10 B parameters), consider reducing to $p = 0.05$.

Gradient clipping: Use threshold $\theta = 1.0$. Monitor the fraction of steps where clipping occurs (should be 1-5%). If clipping occurs more frequently, consider reducing the learning rate.

Adam hyperparameters: Use $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-8}$. These values work well across different model sizes and tasks.

11.12.4 Training Loop

The training loop orchestrates data loading, forward-backward passes, optimization, and monitoring. A well-structured training loop is essential for efficient and stable training.

Mixed precision: Enable FP16 or BF16 mixed precision training to reduce memory and accelerate computation. Use automatic mixed precision (AMP) libraries like PyTorch's `torch.cuda.amp` or NVIDIA Apex to handle the complexity of loss scaling and data type conversions.

Gradient accumulation: If using gradient accumulation, ensure the loss is scaled by $1/K$ where K is the number of accumulation steps. This ensures the accumulated gradient has the correct magnitude.

Gradient checkpointing: Enable gradient checkpointing if memory-constrained, particularly for long sequences or deep models. The 20-30% time overhead is usually acceptable given the memory savings.

Learning rate schedule: Implement the chosen schedule (warmup + linear decay, warmup + cosine, etc.) and update the learning rate at each step. Most optimization libraries provide schedulers that handle this automatically.

Validation: Evaluate on a validation set every N steps (typically every 1,000-10,000 steps depending on dataset size). Compute validation loss and task-specific metrics. Use validation performance to detect overfitting and select the best checkpoint.

Checkpointing: Save model checkpoints regularly (every 10,000-50,000 steps) to enable recovery from failures and to preserve the best model. Save optimizer state along with model parameters to enable seamless resumption of training.

11.12.5 Monitoring and Debugging

Effective monitoring helps detect issues early and guides hyperparameter tuning. Track these metrics throughout training:

Training loss: Should decrease steadily. If loss plateaus early, the learning rate may be too low. If loss spikes or diverges, the learning rate may be too high or gradient clipping may be insufficient.

Validation loss: Should track training loss initially, then diverge as the model begins to overfit. If validation loss increases while training loss decreases, increase regularization (dropout, weight decay) or reduce model capacity.

Perplexity: For language modeling, perplexity = $\exp(\text{loss})$ provides an interpretable metric. Lower perplexity indicates better predictions. BERT-base achieves perplexity around 3-4 on masked language modeling.

Learning rate: Monitor the current learning rate to verify the schedule is working correctly. The learning rate should increase during warmup, then decrease during decay.

Gradient norm: Track the norm of the gradient vector. Typical values are 0.1-10. Very small gradients (<0.01) may indicate vanishing gradients or a learning rate that's too low. Very large gradients (>100) may indicate exploding gradients or a learning rate that's too high.

Parameter norm: Track the norm of the parameter vector. This should increase gradually during training as the model learns. Sudden jumps may indicate instability.

GPU memory usage: Monitor memory consumption to ensure you're using available memory efficiently. If memory usage is much lower than GPU capacity, consider increasing batch size.

Throughput: Track tokens processed per second. This helps identify performance bottlenecks and measure the impact of optimizations.

Common issues and solutions:

Loss doesn't decrease: Check learning rate (may be too low), verify data is loading correctly, check initialization (may be poor), ensure gradients are flowing (check gradient norms).

Loss spikes or diverges: Reduce learning rate, increase warmup period, enable or strengthen gradient clipping, check for data quality issues (corrupted examples, extreme outliers).

Training is slow: Enable mixed precision if not already enabled, increase batch size if memory allows, use gradient accumulation to increase effective batch size, profile to identify bottlenecks (data loading, computation, communication).

Out of memory: Reduce batch size, enable gradient checkpointing, reduce sequence length, enable mixed precision, use gradient accumulation to maintain effective batch size.

Poor generalization: Increase regularization (dropout, weight decay), use label smoothing, train on more data, reduce model capacity, use data augmentation.

This comprehensive training recipe provides a solid foundation for training transformers. While specific details may need adjustment based on the task, dataset, and available hardware, these guidelines capture best practices that have proven effective across a wide range of transformer training scenarios.

11.13 Exercises

Exercise 11.1. Implement the complete mixed precision training algorithm for a small transformer. Compare memory consumption and training time with FP32 training. Experiment with different loss scaling factors and observe their impact on training stability.

Exercise 11.2. For BERT-base with batch size 32 and sequence length 512, calculate the exact memory requirements for: (a) parameters and optimizer states (AdamW), (b) activations for each layer type, (c) total memory with and without gradient checkpointing. Verify your calculations by profiling actual memory usage during training.

Exercise 11.3. Implement gradient accumulation to achieve an effective batch size of 512 with physical batch size 32. Measure the training time overhead compared to true batch size 512 (if it fits in memory). Verify that the training dynamics are identical by comparing loss curves.

Exercise 11.4. Train a small transformer (6 layers, $d_{\text{model}} = 256$) with different learning rate schedules: (a) warmup + linear decay, (b) warmup + inverse square root decay, (c) warmup + cosine annealing. Compare convergence speed and final performance. Plot the learning rate curves and loss curves.

Exercise 11.5. Implement data parallelism for training on 4 GPUs. Measure the speedup compared to single-GPU training. Calculate the communication overhead by comparing the time spent in AllReduce operations versus computation. Experiment with different batch sizes and observe how they affect the computation-to-communication ratio.

Exercise 11.6. Analyze the impact of different regularization techniques on a small transformer trained on a limited dataset (10,000 examples). Compare: (a) no regularization, (b) dropout only, (c) weight decay only, (d) dropout + weight decay, (e) dropout + weight decay + label smoothing. Measure training loss, validation loss, and generalization gap.

Exercise 11.7. Estimate the training time and cost for a GPT-2 medium model (345M parameters) on your available hardware. Calculate: (a) FLOPs per training step, (b) expected throughput (tokens/sec), (c) total training time for 10B tokens, (d) estimated cost on cloud platforms. Compare your estimates with actual training runs.

Exercise 11.8. Implement dynamic batching to minimize padding waste. Compare throughput (tokens/sec) with and without dynamic batching on a dataset with variable-length sequences. Measure the padding fraction in each case and calculate the theoretical maximum speedup from eliminating padding.

Chapter 12

Computational Analysis of Transformers

Chapter Overview

Understanding computational requirements is crucial for deploying transformers. This chapter analyzes time and space complexity, memory footprints, and inference costs. We derive exact FLOP counts, memory requirements, and scaling laws for transformers of different sizes.

Learning Objectives

1. Calculate FLOPs for transformer forward and backward passes
2. Analyze memory requirements for training and inference
3. Understand scaling laws for model size, data, and compute
4. Optimize inference through batching and caching
5. Estimate training time and costs for large models

12.1 Computational Complexity

Understanding the computational complexity of transformers is essential for making informed decisions about model architecture, hardware requirements, and deployment strategies. The transformer’s computational profile differs fundamentally from recurrent architectures, trading sequential dependencies for quadratic memory scaling—a trade-off that profoundly impacts both training and inference.

12.1.1 Self-Attention Complexity

Self-attention is the defining operation of transformers, and its computational characteristics determine much of the model’s behavior. For a sequence of length n with model dimension d_{model} , we analyze each component of the attention mechanism in detail.

QKV Projections: The first step projects the input $\mathbf{X} \in \mathbb{R}^{n \times d_{\text{model}}}$ into query, key, and value spaces. Each projection is a matrix multiplication:

$$\mathbf{Q} = \mathbf{X}\mathbf{W}^Q, \quad \mathbf{K} = \mathbf{X}\mathbf{W}^K, \quad \mathbf{V} = \mathbf{X}\mathbf{W}^V \quad (12.1)$$

where $\mathbf{W}^Q, \mathbf{W}^K, \mathbf{W}^V \in \mathbb{R}^{d_{\text{model}} \times d_k}$ (typically $d_k = d_{\text{model}}/h$ for h heads).

Each matrix multiplication $\mathbf{X}\mathbf{W}$ requires $2nd_{\text{model}}d_k$ floating-point operations (FLOPs): for each of $n \times d_k$ output elements, we perform d_{model} multiply-add operations. With three projections:

$$\text{QKV FLOPs} = 3 \times 2nd_{\text{model}}d_k = 6nd_{\text{model}}d_k \quad (12.2)$$

For the common case where $d_k = d_{\text{model}}$ (single-head or considering all heads together):

$$\text{QKV FLOPs} = 6nd_{\text{model}}^2 \quad (12.3)$$

Why this matters for hardware: These are dense matrix multiplications, which achieve high utilization on modern GPUs. NVIDIA A100 GPUs can perform up to 312 TFLOPS (FP16 with Tensor Cores), meaning these projections are typically compute-bound rather than memory-bound. However, for small batch sizes or short sequences, the operations may become memory-bandwidth limited, achieving only 10-20% of peak FLOPS.

Attention Score Computation: Computing $\mathbf{S} = \mathbf{Q}\mathbf{K}^\top$ involves multiplying $\mathbf{Q} \in \mathbb{R}^{n \times d_k}$ by $\mathbf{K}^\top \in \mathbb{R}^{d_k \times n}$, producing $\mathbf{S} \in \mathbb{R}^{n \times n}$:

$$\text{Score FLOPs} = 2n^2d_k \quad (12.4)$$

The attention matrix \mathbf{S} has n^2 elements, and computing each requires d_k multiply-add operations. This quadratic scaling in sequence length is the fundamental bottleneck for long-context transformers.

Dimension tracking example: For BERT-base with $n = 512$, $d_k = 64$ (per head), and $h = 12$ heads:

$$\mathbf{Q}^{(i)} \in \mathbb{R}^{512 \times 64} \quad (\text{one head}) \quad (12.5)$$

$$\mathbf{K}^{(i)\top} \in \mathbb{R}^{64 \times 512} \quad (12.6)$$

$$\mathbf{S}^{(i)} = \mathbf{Q}^{(i)}\mathbf{K}^{(i)\top} \in \mathbb{R}^{512 \times 512} \quad (262,144 \text{ elements!}) \quad (12.7)$$

Across 12 heads, we compute 12 separate 512×512 attention matrices, requiring:

$$12 \times 2 \times 512^2 \times 64 = 402,653,184 \text{ FLOPs} \approx 403 \text{ MFLOPs} \quad (12.8)$$

Hardware implications: The attention matrix requires n^2 memory per head. For $n = 512$ and 12 heads with FP32:

$$12 \times 512^2 \times 4 \text{ bytes} = 12,582,912 \text{ bytes} \approx 12 \text{ MB} \quad (12.9)$$

This seems modest, but for $n = 2048$ (GPT-2): $12 \times 2048^2 \times 4 = 201 \text{ MB}$ per sequence. With batch size 32: 6.4 GB just for attention matrices! This is why long-context models require substantial GPU memory.

Softmax and Scaling: Applying softmax to each row of \mathbf{S} requires $O(n^2)$ operations (exponentials and normalization), which is negligible compared to the matrix multiplications but can become significant for very long sequences due to memory access patterns.

Attention Output: Computing $\mathbf{O} = \mathbf{A}\mathbf{V}$ multiplies the attention weights $\mathbf{A} \in \mathbb{R}^{n \times n}$ by values $\mathbf{V} \in \mathbb{R}^{n \times d_v}$:

$$\text{Output FLOPs} = 2n^2d_v \quad (12.10)$$

Again, this scales quadratically with sequence length. For $d_v = d_k$:

$$\text{Attention output FLOPs} = 2n^2d_k \quad (12.11)$$

Output Projection: Finally, concatenated head outputs are projected back to model dimension:

$$\text{Output projection FLOPs} = 2n(hd_k)d_{\text{model}} = 2nd_{\text{model}}^2 \quad (12.12)$$

(assuming $hd_k = d_{\text{model}}$).

Total Self-Attention FLOPs:

$$\text{Total} = 6nd_{\text{model}}^2 + 2n^2d_kh + 2n^2d_vh + 2nd_{\text{model}}^2 = 8nd_{\text{model}}^2 + 4n^2d_{\text{model}} \quad (12.13)$$

For typical configurations where $d_k = d_v = d_{\text{model}}/h$:

$$\boxed{\text{Self-Attention FLOPs} = 8nd_{\text{model}}^2 + 4n^2d_{\text{model}}} \quad (12.14)$$

Complexity regime analysis: The relative importance of the two terms depends on the ratio n/d_{model} :

- **Short sequences** ($n \ll d_{\text{model}}$): The $8nd_{\text{model}}^2$ term dominates. For BERT-base with $n = 128$, $d = 768$: $8 \times 128 \times 768^2 \approx 603\text{M}$ vs $4 \times 128^2 \times 768 \approx 50\text{M}$. The projections dominate.
- **Long sequences** ($n \gg d_{\text{model}}$): The $4n^2d_{\text{model}}$ term dominates. For $n = 8192$, $d = 768$: $8 \times 8192 \times 768^2 \approx 38.7\text{G}$ vs $4 \times 8192^2 \times 768 \approx 206\text{G}$. The attention computation dominates.
- **Crossover point**: When $8nd_{\text{model}}^2 \approx 4n^2d_{\text{model}}$, solving gives $n \approx 2d_{\text{model}}$. For $d = 768$, this occurs around $n \approx 1536$.

This analysis explains why efficient attention mechanisms (Chapter 16) focus on reducing the $O(n^2)$ term for long-context applications.

12.1.2 Feed-Forward Network Complexity

The position-wise feed-forward network (FFN) in each transformer layer typically expands the representation to a higher dimension before projecting back. This two-layer network with GELU or ReLU activation is applied independently to each position in the sequence.

Architecture: For input $\mathbf{X} \in \mathbb{R}^{n \times d_{\text{model}}}$:

$$\mathbf{H} = \text{GELU}(\mathbf{X}\mathbf{W}_1 + \mathbf{b}_1) \quad \mathbf{W}_1 \in \mathbb{R}^{d_{\text{model}} \times d_{\text{ff}}}, \quad \mathbf{H} \in \mathbb{R}^{n \times d_{\text{ff}}} \quad (12.15)$$

$$\mathbf{Y} = \mathbf{H}\mathbf{W}_2 + \mathbf{b}_2 \quad \mathbf{W}_2 \in \mathbb{R}^{d_{\text{ff}} \times d_{\text{model}}}, \quad \mathbf{Y} \in \mathbb{R}^{n \times d_{\text{model}}} \quad (12.16)$$

The intermediate dimension d_{ff} is typically $4d_{\text{model}}$ in standard transformers (BERT, GPT), though some models use different ratios. This expansion allows the network to learn complex non-linear transformations.

First Projection FLOPs: Computing $\mathbf{X}\mathbf{W}_1$ requires:

$$\text{First projection} = 2n \cdot d_{\text{model}} \cdot d_{\text{ff}} \quad (12.17)$$

For $d_{\text{ff}} = 4d_{\text{model}}$:

$$\text{First projection} = 2n \cdot d_{\text{model}} \cdot 4d_{\text{model}} = 8nd_{\text{model}}^2 \quad (12.18)$$

Second Projection FLOPs: Computing $\mathbf{H}\mathbf{W}_2$ requires:

$$\text{Second projection} = 2n \cdot d_{\text{ff}} \cdot d_{\text{model}} = 8nd_{\text{model}}^2 \quad (12.19)$$

Total FFN FLOPs:

$$\boxed{\text{FFN FLOPs} = 16nd_{\text{model}}^2 \quad (\text{for } d_{\text{ff}} = 4d_{\text{model}})} \quad (12.20)$$

Activation function: GELU requires additional operations (exponentials, multiplications) but these are $O(nd_{\text{ff}})$, negligible compared to the matrix multiplications.

Why FFN dominates computation: Comparing FFN to self-attention:

$$\text{FFN: } 16nd_{\text{model}}^2 \quad (12.21)$$

$$\text{Attention: } 8nd_{\text{model}}^2 + 4n^2d_{\text{model}} \quad (12.22)$$

For typical sequence lengths where $n < 2d_{\text{model}}$, the FFN requires roughly twice the FLOPs of attention! This is why some efficient transformer variants (e.g., mixture-of-experts) focus on making the FFN more efficient.

Memory and bandwidth considerations: The FFN intermediate activations $\mathbf{H} \in \mathbb{R}^{n \times d_{\text{ff}}}$ must be stored for backpropagation. For BERT-base with $n = 512$, $d_{\text{ff}} = 3072$:

$$512 \times 3072 \times 4 \text{ bytes} = 6,291,456 \text{ bytes} \approx 6 \text{ MB per layer} \quad (12.23)$$

With 12 layers and batch size 32: $6 \times 12 \times 32 = 2.3 \text{ GB}$ just for FFN intermediate activations. This is a significant portion of training memory.

Hardware utilization: FFN matrix multiplications are highly regular and achieve excellent GPU utilization (often 70-90% of peak FLOPS on modern GPUs). The operations are:

- **Compute-bound** for reasonable batch sizes and sequence lengths
- **Well-suited for Tensor Cores** on NVIDIA GPUs (FP16/BF16 operations)
- **Easily parallelizable** across the sequence dimension

On an NVIDIA A100 GPU (312 TFLOPS FP16), computing the FFN for BERT-base with batch size 32 and $n = 512$:

$$\text{FLOPs} = 32 \times 16 \times 512 \times 768^2 \approx 154 \text{ GFLOPs} \quad (12.24)$$

$$\text{Time} \approx \frac{154 \text{ GFLOPs}}{312 \times 0.8 \text{ TFLOPs}} \approx 0.62 \text{ ms} \quad (12.25)$$

(assuming 80% utilization).

12.1.3 Per-Layer Total Complexity

Combining self-attention and FFN, a complete transformer layer requires:

$$\boxed{\text{Transformer layer} = (8nd_{\text{model}}^2 + 4n^2d_{\text{model}}) + 16nd_{\text{model}}^2 = 24nd_{\text{model}}^2 + 4n^2d_{\text{model}} \text{ FLOPs}} \quad (12.26)$$

Additional operations: Layer normalization, residual connections, and dropout add $O(nd_{\text{model}})$ operations, which are negligible compared to the matrix multiplications.

Breakdown by component:

- **FFN:** $16nd_{\text{model}}^2$ (typically 60-70% of layer FLOPs for short sequences)
- **Attention projections:** $8nd_{\text{model}}^2$ (typically 25-35%)
- **Attention computation:** $4n^2d_{\text{model}}$ (grows with sequence length)

This breakdown is crucial for optimization: for short sequences, optimizing FFN yields the largest gains; for long sequences, efficient attention mechanisms become critical.

Example 12.1 (BERT-base Single Layer: Detailed Analysis). BERT-base parameters: $n = 512$, $d_{\text{model}} = 768$, $h = 12$, $d_k = d_v = 64$, $d_{ff} = 3072$

Self-Attention Computation:

QKV Projections (all heads):

$$3 \times 2 \times 512 \times 768 \times 768 = 1,811,939,328 \approx 1.81 \text{ GFLOPs} \quad (12.27)$$

Attention scores (\mathbf{QK}^\top for 12 heads):

$$12 \times 2 \times 512^2 \times 64 = 402,653,184 \approx 403 \text{ MFLOPs} \quad (12.28)$$

Attention output (\mathbf{AV} for 12 heads):

$$12 \times 2 \times 512^2 \times 64 = 402,653,184 \approx 403 \text{ MFLOPs} \quad (12.29)$$

Output projection:

$$2 \times 512 \times 768 \times 768 = 603,979,776 \approx 604 \text{ MFLOPs} \quad (12.30)$$

Total self-attention:

$$1.81 + 0.40 + 0.40 + 0.60 = 3.21 \text{ GFLOPs} \quad (12.31)$$

Feed-Forward Network:

First projection (\mathbf{XW}_1 , where $\mathbf{W}_1 \in \mathbb{R}^{768 \times 3072}$):

$$2 \times 512 \times 768 \times 3072 = 2,415,919,104 \approx 2.42 \text{ GFLOPs} \quad (12.32)$$

Second projection (\mathbf{HW}_2 , where $\mathbf{W}_2 \in \mathbb{R}^{3072 \times 768}$):

$$2 \times 512 \times 3072 \times 768 = 2,415,919,104 \approx 2.42 \text{ GFLOPs} \quad (12.33)$$

Total FFN:

$$2.42 + 2.42 = 4.84 \text{ GFLOPs} \quad (12.34)$$

Complete Layer:

$$3.21 + 4.84 = 8.05 \text{ GFLOPs per layer} \quad (12.35)$$

Full 12-Layer BERT-base:

$$12 \times 8.05 = 96.6 \text{ GFLOPs per forward pass} \quad (12.36)$$

Training (forward + backward): Backward pass requires approximately $2 \times$ the FLOPs of forward pass:

$$\text{Training step} \approx 3 \times 96.6 = 290 \text{ GFLOPs} \quad (12.37)$$

Hardware timing on NVIDIA A100 (312 TFLOPS FP16):

Assuming 70% utilization (realistic for mixed operations):

$$\text{Forward pass time} \approx \frac{96.6 \text{ GFLOPS}}{312 \times 0.7 \text{ TFLOPS}} \approx 0.44 \text{ ms} \quad (12.38)$$

For batch size 32:

$$\text{Batch forward time} \approx 32 \times 0.44 = 14 \text{ ms} \quad (12.39)$$

Training throughput: With batch size 32, sequence length 512:

$$\text{Tokens per second} = \frac{32 \times 512}{14 \text{ ms} \times 3} \approx 390,000 \text{ tokens/sec} \quad (12.40)$$

This analysis shows why BERT-base training is feasible on single GPUs, while larger models require distributed training.

Memory bandwidth considerations: The A100 has 1.6 TB/s memory bandwidth. Loading model parameters ($110\text{M} \times 4 \text{ bytes} = 440 \text{ MB}$) takes:

$$\text{Parameter load time} = \frac{440 \text{ MB}}{1600 \text{ GB/s}} \approx 0.28 \text{ ms} \quad (12.41)$$

This is comparable to compute time, indicating that for small batch sizes, the model can become memory-bandwidth bound rather than compute-bound. Larger batch sizes amortize parameter loading across more computation, improving utilization.

12.1.4 Complexity Analysis

Theorem 12.1 (Transformer Complexity). *For L layers, sequence length n , dimension d :*

Time complexity: $O(Ln^2d + Lnd^2)$

Space complexity: $O(Ln^2 + Lnd)$

Comparison with RNN:

- RNN: $O(Lnd^2)$ time, $O(Ld^2)$ space
- Transformer: Quadratic in n but parallel; RNN sequential

Bottleneck regimes:

- Short sequences ($n < d$): FFN dominates, $O(Lnd^2)$
- Long sequences ($n > d$): Attention dominates, $O(Ln^2d)$

12.2 Memory Requirements

Memory is often the limiting factor in training and deploying large transformer models. Understanding memory requirements at a granular level enables informed decisions about model architecture, batch sizes, and hardware selection. We analyze memory consumption across four categories: model parameters, gradients, optimizer states, and activations.

12.2.1 Model Parameters

Model parameters must be stored in GPU memory during both training and inference. The memory footprint depends on the numerical precision used.

Precision options:

- **FP32 (float32):** 4 bytes per parameter, standard precision
- **FP16 (float16):** 2 bytes per parameter, half precision
- **BF16 (bfloat16):** 2 bytes per parameter, better range than FP16
- **INT8:** 1 byte per parameter, quantized inference

For BERT-base with 110 million parameters:

$$\text{FP32: } 110,000,000 \times 4 = 440,000,000 \text{ bytes} = 440 \text{ MB} \quad (12.42)$$

$$\text{FP16/BF16: } 110,000,000 \times 2 = 220,000,000 \text{ bytes} = 220 \text{ MB} \quad (12.43)$$

$$\text{INT8: } 110,000,000 \times 1 = 110,000,000 \text{ bytes} = 110 \text{ MB} \quad (12.44)$$

Parameter breakdown for BERT-base:

$$\text{Token embeddings: } V \times d_{\text{model}} = 30,000 \times 768 = 23,040,000 \text{ params} \quad (12.45)$$

$$\text{Position embeddings: } n_{\text{max}} \times d_{\text{model}} = 512 \times 768 = 393,216 \text{ params} \quad (12.46)$$

$$\text{Segment embeddings: } 2 \times d_{\text{model}} = 2 \times 768 = 1,536 \text{ params} \quad (12.47)$$

Per transformer layer:

$$\text{Self-attention: } 4 \times d_{\text{model}}^2 = 4 \times 768^2 = 2,359,296 \text{ params} \quad (12.48)$$

$$\text{FFN: } 2 \times d_{\text{model}} \times d_{\text{ff}} = 2 \times 768 \times 3072 = 4,718,592 \text{ params} \quad (12.49)$$

$$\text{Layer norms: } 4 \times d_{\text{model}} = 4 \times 768 = 3,072 \text{ params} \quad (12.50)$$

$$\text{Total per layer: } 7,080,960 \text{ params} \quad (12.51)$$

12 layers:

$$12 \times 7,080,960 = 84,971,520 \text{ params} \quad (12.52)$$

Total BERT-base:

$$23,040,000 + 393,216 + 1,536 + 84,971,520 = 108,406,272 \approx 110\text{M params} \quad (12.53)$$

In FP32: $110\text{M} \times 4 = 440 \text{ MB}$

Larger models scale dramatically:

$$\text{GPT-2 (1.5B): } 1,500,000,000 \times 4 = 6,000 \text{ MB} = 6 \text{ GB (FP32)} \quad (12.54)$$

$$\text{GPT-3 (175B): } 175,000,000,000 \times 4 = 700,000 \text{ MB} = 700 \text{ GB (FP32)} \quad (12.55)$$

GPT-3 in FP32 requires 700 GB just for parameters—far exceeding single GPU memory (A100 has 80 GB). This necessitates:

- **Model parallelism:** Split model across multiple GPUs
- **Mixed precision:** Use FP16/BF16 (350 GB for GPT-3)
- **Quantization:** INT8 inference (175 GB for GPT-3)

12.2.2 Activation Memory

During training, intermediate activations must be stored for backpropagation. Activation memory scales with batch size and sequence length, often dominating memory consumption.

Activations per transformer layer:

- **Input to layer:** $B \times n \times d_{\text{model}}$
- **Query, Key, Value:** $3 \times B \times n \times d_{\text{model}}$
- **Attention scores:** $B \times h \times n \times n$ (quadratic in sequence length!)
- **Attention output:** $B \times n \times d_{\text{model}}$
- **FFN intermediate:** $B \times n \times d_{\text{ff}}$
- **Layer norm activations:** $2 \times B \times n \times d_{\text{model}}$

Total activation memory per layer (approximate):

$$\text{Memory} \approx B \times n \times (8d_{\text{model}} + d_{\text{ff}}) + B \times h \times n^2 \quad (12.56)$$

For BERT-base ($B = 32$, $n = 512$, $d_{\text{model}} = 768$, $h = 12$, $d_{\text{ff}} = 3072$):

$$\text{Linear terms: } 32 \times 512 \times (8 \times 768 + 3072) \times 4 \text{ bytes} \quad (12.57)$$

$$= 32 \times 512 \times 9,216 \times 4 = 603,979,776 \text{ bytes} \approx 604 \text{ MB} \quad (12.58)$$

$$\text{Attention matrices: } 32 \times 12 \times 512^2 \times 4 = 402,653,184 \text{ bytes} \approx 403 \text{ MB} \quad (12.59)$$

$$\text{Total per layer: } \approx 1,007 \text{ MB} \approx 1 \text{ GB} \quad (12.60)$$

For 12 layers: $12 \times 1 \text{ GB} = 12 \text{ GB}$ just for activations!

Impact of sequence length: The attention matrix term $B \times h \times n^2$ grows quadratically. For $n = 2048$ (4× longer):

$$32 \times 12 \times 2048^2 \times 4 = 6,442,450,944 \text{ bytes} \approx 6.4 \text{ GB per layer} \quad (12.61)$$

For 12 layers: 77 GB just for attention matrices—nearly filling an A100 GPU!

This quadratic scaling is why:

- Long-context models require gradient checkpointing (recompute activations during backward pass)
- Efficient attention mechanisms (sparse, linear) are crucial for long sequences
- Batch sizes must be reduced for longer sequences

Gradient checkpointing trade-off: Recomputing activations during backward pass:

- **Memory savings:** Reduce activation memory by $\sim 80\%$
- **Compute cost:** Increase training time by $\sim 20\text{-}30\%$
- **When to use:** When memory-constrained, especially for long sequences

Example 12.2 (GPT-2 Activation Memory: Complete Analysis). GPT-2 (small): $L = 12$, $d_{\text{model}} = 768$, $h = 12$, $d_k = 64$, $d_{\text{ff}} = 3072$, $n = 1024$

Per-layer activation breakdown (batch size $B = 1$):

QKV projections:

$$3 \times 1024 \times 768 \times 4 = 9,437,184 \text{ bytes} \approx 9.4 \text{ MB} \quad (12.62)$$

Attention matrices (12 heads):

$$12 \times 1024^2 \times 4 = 50,331,648 \text{ bytes} \approx 50.3 \text{ MB} \quad (12.63)$$

This is the dominant term! For $n = 2048$: $12 \times 2048^2 \times 4 = 201 \text{ MB}$ ($4\times$ larger).

Attention output:

$$1024 \times 768 \times 4 = 3,145,728 \text{ bytes} \approx 3.1 \text{ MB} \quad (12.64)$$

FFN intermediate:

$$1024 \times 3072 \times 4 = 12,582,912 \text{ bytes} \approx 12.6 \text{ MB} \quad (12.65)$$

Layer norm and residuals:

$$3 \times 1024 \times 768 \times 4 = 9,437,184 \text{ bytes} \approx 9.4 \text{ MB} \quad (12.66)$$

Total per layer:

$$9.4 + 50.3 + 3.1 + 12.6 + 9.4 = 84.8 \text{ MB} \quad (12.67)$$

12 layers: $12 \times 84.8 = 1,018 \text{ MB} \approx 1 \text{ GB}$ for single sequence

Batch size scaling:

$$B = 8 : 8 \text{ GB} \quad (12.68)$$

$$B = 16 : 16 \text{ GB} \quad (12.69)$$

$$B = 32 : 32 \text{ GB} \quad (12.70)$$

$$B = 64 : 64 \text{ GB} \quad (12.71)$$

Hardware implications:

- **NVIDIA V100 (16 GB):** Maximum batch size $\approx 12 - 14$ (accounting for parameters and optimizer states)
- **NVIDIA A100 (40 GB):** Maximum batch size $\approx 30 - 35$
- **NVIDIA A100 (80 GB):** Maximum batch size $\approx 70 - 75$

Gradient checkpointing impact: With checkpointing, only store activations at layer boundaries, recompute within layers during backward pass:

$$\text{Memory reduction} \approx 80\% \Rightarrow 1 \text{ GB} \rightarrow 200 \text{ MB per sequence} \quad (12.72)$$

This allows batch size 64 on V100 (16 GB), but increases training time by $\sim 25\%$.

Mixed precision training: Using FP16 for activations (FP32 for parameters):

$$\text{Activation memory} \rightarrow 1 \text{ GB}/2 = 500 \text{ MB per sequence} \quad (12.73)$$

Combined with gradient checkpointing: $500 \times 0.2 = 100 \text{ MB}$ per sequence, enabling very large batch sizes.

12.2.3 Training Memory Budget

Training requires memory for parameters, gradients, optimizer states, and activations. Understanding this breakdown is essential for selecting hardware and configuring training.

Total training memory:

$$\text{Memory}_{\text{total}} = \text{Parameters} + \text{Gradients} + \text{Optimizer States} + \text{Activations} \quad (12.74)$$

For AdamW optimizer (most common for transformers):

- **Model parameters:** P parameters $\times 4$ bytes (FP32) $= 4P$ bytes

- **Gradients:** P parameters \times 4 bytes = $4P$ bytes
- **First moment (momentum):** P parameters \times 4 bytes = $4P$ bytes
- **Second moment (variance):** P parameters \times 4 bytes = $4P$ bytes
- **Activations:** A bytes (depends on batch size, sequence length, model depth)

Total: $16P + A$ bytes

Mixed precision training (FP16/BF16 with FP32 master weights):

- **FP16 parameters (forward/backward):** $2P$ bytes
- **FP32 master parameters:** $4P$ bytes
- **FP32 gradients:** $4P$ bytes
- **FP32 optimizer states:** $8P$ bytes
- **FP16 activations:** $A/2$ bytes

Total: $18P + A/2$ bytes

Surprisingly, mixed precision uses slightly MORE memory for parameters/optimizer (18P vs 16P) but saves significantly on activations ($A/2$ vs A). Since activations often dominate, mixed precision typically reduces total memory.

Example 12.3 (BERT-base Training Memory Budget). BERT-base: 110M parameters, batch size 32, sequence length 512

FP32 training:

$$\text{Parameters: } 110\text{M} \times 4 = 440 \text{ MB} \quad (12.75)$$

$$\text{Gradients: } 110\text{M} \times 4 = 440 \text{ MB} \quad (12.76)$$

$$\text{Adam states (2}\times\text{): } 2 \times 110\text{M} \times 4 = 880 \text{ MB} \quad (12.77)$$

$$\text{Activations: } 32 \times 12 \times 1 \text{ GB} = 12 \text{ GB} \quad (12.78)$$

$$\text{Total: } 440 + 440 + 880 + 12,000 = 13,760 \text{ MB} \approx 13.8 \text{ GB} \quad (12.79)$$

Fits on: NVIDIA V100 (16 GB), A100 (40/80 GB), RTX 3090 (24 GB)

Mixed precision training:

$$\text{FP16 parameters: } 110\text{M} \times 2 = 220 \text{ MB} \quad (12.80)$$

$$\text{FP32 master + gradients + Adam: } 110\text{M} \times 16 = 1,760 \text{ MB} \quad (12.81)$$

$$\text{FP16 activations: } 12,000/2 = 6,000 \text{ MB} \quad (12.82)$$

$$\text{Total: } 220 + 1,760 + 6,000 = 7,980 \text{ MB} \approx 8 \text{ GB} \quad (12.83)$$

Mixed precision saves: $13.8 - 8 = 5.8 \text{ GB}$ (42% reduction)

With gradient checkpointing: Activations reduced by 80%:

$$220 + 1,760 + 1,200 = 3,180 \text{ MB} \approx 3.2 \text{ GB} \quad (12.84)$$

This enables batch size 128 on V100 (16 GB)!

Example 12.4 (GPT-3 Training Memory Requirements). GPT-3: 175B parameters, sequence length 2048

Parameters and optimizer (FP32):

$$175\text{B} \times 16 = 2,800 \text{ GB} \quad (12.85)$$

Activations (batch size 1, single sequence):

$$\text{Per layer: } \approx 2048 \times (8 \times 12,288 + 4 \times 12,288) + 96 \times 2048^2 \quad (12.86)$$

$$\approx 2048 \times 147,456 + 402,653,184 \quad (12.87)$$

$$\approx 704 \text{ MB per layer} \quad (12.88)$$

96 layers: $96 \times 704 \text{ MB} \approx 68 \text{ GB}$ per sequence

Total for batch size 1: $2,800 + 68 = 2,868 \text{ GB}$

Hardware requirements:

- **Single A100 (80 GB):** Impossible—need 36 GPUs just for parameters!
- **Model parallelism:** Split across 8 GPUs: $2,868/8 = 359 \text{ GB}$ per GPU—still too large!
- **Mixed precision + model parallelism:** $\approx 1,500 \text{ GB total}/8 = 188 \text{ GB}$ per GPU—still too large!
- **Mixed precision + model parallelism + gradient checkpointing:** $\approx 800 \text{ GB}/8 = 100 \text{ GB}$ per GPU—still exceeds A100!

Actual GPT-3 training: Used ZeRO optimizer (shards optimizer states across GPUs) + model parallelism + pipeline parallelism across thousands of GPUs.

This example illustrates why training models beyond $\sim 10\text{B}$ parameters requires sophisticated distributed training strategies.

12.2.4 Hardware Selection Guide**GPU memory requirements by model size (mixed precision + gradient checkpointing):**

Model Size	Parameters	Min GPU Memory	Recommended GPU
Small	100M	8 GB	RTX 3070, V100
Base	300M	12 GB	RTX 3080, V100
Large	1B	24 GB	RTX 3090, A5000
XL	3B	40 GB	A100 (40 GB)
XXL	10B	80 GB	A100 (80 GB)
175B (GPT-3)	175B	$8 \times \text{A100 (80 GB)}$	Multi-node cluster

Table 12.1: GPU memory requirements for training transformer models (batch size 8-16, sequence length 512-1024)

Inference memory requirements (FP16):

- **Parameters only:** $2P$ bytes
- **KV cache (autoregressive):** $2 \times L \times h \times n_{\max} \times d_k \times B$ bytes
- **Activations (single forward pass):** Minimal compared to training

For GPT-2 (117M params) inference:

$$\text{Parameters: } 117\text{M} \times 2 = 234 \text{ MB} \quad (12.89)$$

$$\text{KV cache (batch 1, } n = 1024\text{): } 2 \times 12 \times 12 \times 1024 \times 64 \times 2 = 38 \text{ MB} \quad (12.90)$$

$$\text{Total: } \approx 300 \text{ MB} \quad (12.91)$$

GPT-2 inference easily fits on consumer GPUs or even CPUs!

12.3 Inference Optimization

Inference optimization is critical for deploying transformers in production. Unlike training, which prioritizes throughput (tokens/second across large batches), inference prioritizes latency (time to generate a single response) while maintaining reasonable throughput. We analyze key optimization techniques and their trade-offs.

12.3.1 KV Caching for Autoregressive Decoding

Autoregressive generation (used in GPT, decoder-only models) generates tokens sequentially, where each new token attends to all previous tokens. Naive implementation recomputes attention for all previous positions at each step—highly inefficient.

Problem analysis: Generating sequence of length T tokens:

- **Step 1:** Compute attention for position 1 (attends to position 1)
- **Step 2:** Compute attention for position 2 (attends to positions 1-2)
- **Step 3:** Compute attention for position 3 (attends to positions 1-3)
- **Step T :** Compute attention for position T (attends to positions 1- T)

Total attention computations: $\sum_{t=1}^T t = \frac{T(T+1)}{2} \approx \frac{T^2}{2}$

For $T = 1000$ tokens: $\approx 500,000$ attention computations!

KV Caching solution: Key and value projections depend only on input tokens, not on the query position. Cache \mathbf{K} and \mathbf{V} from previous steps:

[H] Autoregressive Generation with KV Caching Prompt tokens $\mathbf{x}_1, \dots, \mathbf{x}_p$, max length T Generated sequence $\mathbf{x}_1, \dots, \mathbf{x}_T$

Initialize cache $\text{cache}_K = []$, $\text{cache}_V = []$

Process prompt $t = 1 \dots p$ $\mathbf{k}_t = \mathbf{W}^K \mathbf{x}_t$, $\mathbf{v}_t = \mathbf{W}^V \mathbf{x}_t$

Append \mathbf{k}_t to cache_K , \mathbf{v}_t to cache_V

$\mathbf{q}_t = \mathbf{W}^Q \mathbf{x}_t$

Compute attention using \mathbf{q}_t and all cached keys/values

Generate \mathbf{h}_t

Generate new tokens $t = p + 1 \dots T$ Sample \mathbf{x}_t from \mathbf{h}_{t-1}

$\mathbf{k}_t = \mathbf{W}^K \mathbf{x}_t$, $\mathbf{v}_t = \mathbf{W}^V \mathbf{x}_t$

Append \mathbf{k}_t to cache_K , \mathbf{v}_t to cache_V

$\mathbf{q}_t = \mathbf{W}^Q \mathbf{x}_t$

Compute attention: $\text{Attention}(\mathbf{q}_t, \text{cache}_K, \text{cache}_V)$

Generate \mathbf{h}_t

Computational savings: With caching, each step computes attention once (not recomputing previous positions):

$$\text{Total computations} = T \quad (\text{vs. } \frac{T^2}{2} \text{ without caching}) \quad (12.92)$$

Speedup: For $T = 1000$: $\frac{500,000}{1,000} = 500\times$ faster!

Memory cost: Store keys and values for all positions and layers:

$$\text{KV cache size} = 2 \times L \times h \times T \times d_k \times \text{sizeof(float)} \quad (12.93)$$

For GPT-2 ($L = 12$, $h = 12$, $d_k = 64$, FP16):

$$2 \times 12 \times 12 \times T \times 64 \times 2 = 36,864 \times T \text{ bytes} \quad (12.94)$$

Memory scaling with sequence length:

$$T = 512 : 36,864 \times 512 = 18,874,368 \text{ bytes} \approx 19 \text{ MB} \quad (12.95)$$

$$T = 1024 : 36,864 \times 1024 = 37,748,736 \text{ bytes} \approx 38 \text{ MB} \quad (12.96)$$

$$T = 2048 : 36,864 \times 2048 = 75,497,472 \text{ bytes} \approx 75 \text{ MB} \quad (12.97)$$

$$T = 4096 : 36,864 \times 4096 = 150,994,944 \text{ bytes} \approx 151 \text{ MB} \quad (12.98)$$

For GPT-3 ($L = 96$, $h = 96$, $d_k = 128$, $T = 2048$, FP16):

$$2 \times 96 \times 96 \times 2048 \times 128 \times 2 = 9,663,676,416 \text{ bytes} \approx 9.7 \text{ GB per sequence} \quad (12.99)$$

Batch inference with KV cache: For batch size B :

$$\text{Total KV cache} = B \times 2 \times L \times h \times T \times d_k \times \text{sizeof(float)} \quad (12.100)$$

For GPT-3 with $B = 8$: $8 \times 9.7 \text{ GB} = 77.6 \text{ GB}$ —nearly filling an A100 (80 GB)!

This is why large-scale inference services:

- Use smaller batch sizes for long contexts
- Implement dynamic batching (group requests of similar lengths)
- Use quantization (INT8) to reduce cache size by $2\text{-}4\times$

12.3.2 Batched Inference

Processing multiple sequences simultaneously increases GPU utilization and throughput.

Single sequence inference: For GPT-2 generating 100 tokens:

- **Compute:** $\approx 100 \times 8 \text{ GFLOPs} = 800 \text{ GFLOPs}$
- **Time on A100:** $\frac{800 \text{ GFLOPs}}{312 \text{ TFLOPs} \times 0.3} \approx 8.5 \text{ ms}$
- **GPU utilization:** $\approx 30\%$ (memory-bound, not compute-bound)

Batched inference (batch size 32):

- **Compute:** $32 \times 800 \text{ GFLOPs} = 25,600 \text{ GFLOPs}$
- **Time on A100:** $\frac{25,600 \text{ GFLOPs}}{312 \text{ TFLOPs} \times 0.7} \approx 117 \text{ ms}$
- **GPU utilization:** $\approx 70\%$ (much better!)
- **Throughput:** $\frac{32 \times 100}{117 \text{ ms}} \approx 27,350 \text{ tokens/sec}$

Latency vs. throughput trade-off:

- **Batch size 1:** Latency = 8.5 ms, throughput = 11,765 tokens/sec
- **Batch size 32:** Latency = 117 ms ($13.8\times$ worse), throughput = 27,350 tokens/sec ($2.3\times$ better)

Padding challenge: Sequences in a batch must have the same length. Shorter sequences are padded, wasting computation:

- Sequence lengths: [512, 256, 128, 64]
- Padded to: [512, 512, 512, 512]
- Wasted computation: $(512 - 256) + (512 - 128) + (512 - 64) = 1024 \text{ positions}$ (50%!)

Solutions:

- **Dynamic batching:** Group sequences of similar lengths
- **Bucket batching:** Pre-defined length buckets (128, 256, 512, 1024)
- **Packed sequences:** Concatenate sequences without padding (requires careful attention masking)

12.3.3 Quantization for Inference

Quantization reduces memory and increases throughput by using lower-precision arithmetic.

Precision options:

- **FP32:** 4 bytes, full precision
- **FP16/BF16:** 2 bytes, half precision (1.5-2× speedup)
- **INT8:** 1 byte, 8-bit integer (2-4× speedup, 4× memory reduction)
- **INT4:** 0.5 bytes, 4-bit integer (4-8× speedup, 8× memory reduction)

INT8 quantization: Map FP32 weights $w \in [-w_{\max}, w_{\max}]$ to INT8 $w_q \in [-128, 127]$:

$$w_q = \text{round} \left(\frac{w}{w_{\max}} \times 127 \right) \quad (12.101)$$

Dequantize during computation:

$$w \approx \frac{w_q \times w_{\max}}{127} \quad (12.102)$$

Quantization impact on GPT-2:

$$\text{FP32: } 117\text{M} \times 4 = 468 \text{ MB} \quad (12.103)$$

$$\text{FP16: } 117\text{M} \times 2 = 234 \text{ MB} \quad (2 \times \text{ reduction}) \quad (12.104)$$

$$\text{INT8: } 117\text{M} \times 1 = 117 \text{ MB} \quad (4 \times \text{ reduction}) \quad (12.105)$$

$$\text{INT4: } 117\text{M} \times 0.5 = 58.5 \text{ MB} \quad (8 \times \text{ reduction}) \quad (12.106)$$

Accuracy trade-offs:

- **FP16/BF16:** Negligible accuracy loss (0.1% perplexity increase)
- **INT8:** Small accuracy loss (0.5-2% perplexity increase) with calibration
- **INT4:** Moderate accuracy loss (2-5% perplexity increase), requires careful quantization

Hardware support:

- **NVIDIA Tensor Cores:** Accelerate FP16/BF16 (up to 2× speedup)
- **NVIDIA INT8 Tensor Cores:** Accelerate INT8 (up to 4× speedup)
- **CPU AVX-512 VNNI:** Accelerate INT8 on CPUs

12.3.4 Model Distillation

Train smaller "student" model to mimic larger "teacher" model:

- **DistilBERT:** 66M params (vs. BERT-base 110M), 97% performance, 2× faster
- **TinyBERT:** 14M params, 96% performance, 7× faster

Distillation enables deployment on resource-constrained devices (mobile, edge).

12.3.5 Inference Optimization Summary

Combined optimizations: KV caching + FP16 + batching + INT8 can achieve 1000× speedup with minimal accuracy loss!

Technique	Speedup	Memory Reduction	Accuracy Impact
KV Caching	100-500×	-50% (cache overhead)	None
Batching (32×	2-3×	throughput	None
FP16/BF16	1.5-2×	2×	Negligible
INT8 Quantization	2-4×	4×	Small (0.5-2%)
INT4 Quantization	4-8×	8×	Moderate (2-5%)
Distillation	2-7×	2-8×	Small (3-4%)

Table 12.2: Inference optimization techniques and their trade-offs

12.4 Scaling Laws

12.4.1 Kaplan et al. Scaling Laws

Performance scales as power law with model size N , dataset size D , and compute C :

$$L(N, D, C) \approx \left(\frac{N_c}{N}\right)^{\alpha_N} + \left(\frac{D_c}{D}\right)^{\alpha_D} + \left(\frac{C_c}{C}\right)^{\alpha_C} \quad (12.107)$$

Key findings:

- Larger models are more sample-efficient
- Compute-optimal: Balance model size and data
- Doubling compute \rightarrow consistent loss reduction

12.4.2 Chinchilla Scaling Laws

For fixed compute budget, optimal allocation:

$$N_{\text{optimal}} \propto C^{0.5}, \quad D_{\text{optimal}} \propto C^{0.5} \quad (12.108)$$

Implication: Many large models (GPT-3) are over-parameterized and under-trained! Chinchilla (70B params, more data) outperforms Gopher (280B params, less data).

12.5 Exercises

Exercise 12.1. Calculate FLOPs for GPT-3 (175B parameters, $L = 96$, $d = 12288$, $h = 96$, $n = 2048$) for: (1) Single forward pass, (2) Generating 100 tokens autoregressively, (3) Training on 1 trillion tokens.

Exercise 12.2. Estimate memory for training 1.3B parameter model with batch size 64, sequence length 2048. What GPU memory required? How to fit on A100 (80GB)?

Exercise 12.3. Implement KV caching for GPT-2. Measure speedup for generating 256 tokens. Plot generation time vs sequence length with/without caching.

Exercise 12.4. For fixed compute budget $C = 10^{24}$ FLOPs: Use Chinchilla scaling to find optimal model size and data size. Compare with GPT-3 allocation.

Part V

Modern Transformer Variants

Chapter 13

BERT: Bidirectional Encoder Representations

Chapter Overview

BERT (Bidirectional Encoder Representations from Transformers) revolutionized NLP by introducing effective bidirectional pre-training. This chapter covers BERT's architecture, pre-training objectives (masked language modeling and next sentence prediction), fine-tuning strategies, and variants (RoBERTa, ALBERT, DistilBERT).

Learning Objectives

1. Understand BERT's encoder-only architecture
2. Implement masked language modeling (MLM)
3. Apply BERT to downstream tasks via fine-tuning
4. Compare BERT variants and their improvements
5. Analyze BERT's learned representations
6. Understand limitations and failure modes

13.1 BERT Architecture

13.1.1 Model Specification

Definition 13.1 (BERT Model). BERT is a stack of transformer encoder layers with:

- **Input:** Token + Segment + Position embeddings
- **Processing:** L transformer encoder layers
- **Output:** Contextualized representations for all tokens

BERT represents a fundamental shift in how we approach natural language understanding by using bidirectional context throughout the entire model. Unlike autoregressive language models such as GPT that can only attend to previous tokens, BERT's encoder-only architecture allows each token to attend to all other tokens in the sequence simultaneously. This bidirectional attention enables BERT to build richer contextual representations that capture both left and right context, making it particularly effective for understanding tasks like question answering, named entity recognition, and text classification.

The architecture comes in two primary configurations that trade off between computational efficiency and model capacity. BERT-base uses 12 transformer encoder layers with hidden dimension $d_{\text{model}} = 768$, employing 12 attention heads where each head operates on dimension $d_k = d_v = 64$. The feed-forward network in each layer expands to dimension $d_{ff} = 3072$, following the standard $4\times$ expansion ratio. This configuration contains approximately 110 million parameters and was designed to be trainable on a modest cluster of TPUs while still achieving strong performance across a wide range of tasks.

BERT-large doubles the model depth to 24 layers and increases the hidden dimension to $d_{\text{model}} = 1024$ with 16 attention heads of dimension $d_k = d_v = 64$ each. The feed-forward dimension scales proportionally to $d_{ff} = 4096$, maintaining the $4\times$ expansion ratio. With approximately 340 million parameters, BERT-large achieves substantially better performance on challenging tasks but requires significantly more computational resources for both training and inference. The parameter count scales roughly quadratically with hidden dimension due to the d_{model}^2 terms in attention projections and feed-forward networks.

BERT-base specifications:

- Layers: $L = 12$
- Hidden size: $d_{\text{model}} = 768$
- Attention heads: $h = 12$, head dimension: $d_k = d_v = 64$
- Feed-forward size: $d_{ff} = 3072$
- Maximum sequence length: $n_{\text{max}} = 512$
- Vocabulary size: $V \approx 30,000$ (WordPiece)
- Total parameters: $\approx 110\text{M}$

BERT-large specifications:

- Layers: $L = 24$
- Hidden size: $d_{\text{model}} = 1024$
- Attention heads: $h = 16$, head dimension: $d_k = d_v = 64$
- Feed-forward size: $d_{ff} = 4096$
- Maximum sequence length: $n_{\text{max}} = 512$
- Vocabulary size: $V \approx 30,000$ (WordPiece)
- Total parameters: $\approx 340\text{M}$

13.1.2 Parameter Breakdown and Memory Requirements

Understanding the parameter distribution across BERT’s components is essential for memory planning and optimization. The embedding layer typically consumes a substantial fraction of total parameters, particularly for models with large vocabularies. For BERT-base, the token embeddings alone account for $V \times d_{\text{model}} = 30,000 \times 768 = 23,040,000$ parameters, representing approximately 21% of the model. Position embeddings add another $n_{\text{max}} \times d_{\text{model}} = 512 \times 768 = 393,216$ parameters, while segment embeddings contribute a negligible $2 \times 768 = 1,536$ parameters for distinguishing sentence pairs.

Each transformer encoder layer contains parameters in three main components. The multi-head self-attention mechanism requires four projection matrices: query, key, value, and output projections, each of dimension $d_{\text{model}} \times d_{\text{model}}$. This yields $4 \times 768^2 = 2,359,296$ parameters per layer for BERT-base. The feed-forward network dominates the per-layer parameter count with two projections: the expansion from d_{model} to d_{ff} and the contraction back, totaling $2 \times 768 \times 3072 = 4,718,592$ parameters. Layer normalization adds $4 \times d_{\text{model}} = 3,072$ parameters (two layer norms per layer, each with scale and shift

parameters). Summing these components gives approximately 7.1 million parameters per layer, or 85 million parameters across all 12 layers.

The memory footprint depends critically on numerical precision. In standard FP32 (32-bit floating point), each parameter requires 4 bytes, so BERT-base occupies $110,000,000 \times 4 = 440$ MB. Modern training typically uses mixed precision with FP16 or BF16 (16-bit formats) for activations and gradients while maintaining FP32 master weights for numerical stability. This reduces the working memory for forward and backward passes to $110,000,000 \times 2 = 220$ MB for the model parameters, though the optimizer still maintains FP32 copies. For inference, pure FP16 weights require only 220 MB, enabling BERT-base to run comfortably on consumer GPUs with 8-16 GB of memory.

BERT-large’s parameter distribution follows the same structure but scales significantly. Token embeddings remain at 23 million parameters (vocabulary size unchanged), but each layer now contains $4 \times 1024^2 = 4,194,304$ attention parameters and $2 \times 1024 \times 4096 = 8,388,608$ feed-forward parameters, totaling approximately 12.6 million parameters per layer. With 24 layers, the transformer stack contributes 302 million parameters. The total of 340 million parameters requires 1.36 GB in FP32 or 680 MB in FP16. This larger footprint means BERT-large training typically requires GPUs with at least 16 GB of memory (such as NVIDIA V100 or A100), and inference benefits from GPUs with 12+ GB to accommodate reasonable batch sizes.

13.1.3 Input Representation

$$\text{Input} = \text{TokenEmb} + \text{SegmentEmb} + \text{PositionEmb} \quad (13.1)$$

Token Embeddings: WordPiece tokenization, vocabulary $\approx 30,000$

Segment Embeddings: Distinguish sentence A vs B (for sentence-pair tasks)

$$\text{SegEmb}(i) = \begin{cases} \mathbf{e}_A & \text{if token } i \text{ in sentence A} \\ \mathbf{e}_B & \text{if token } i \text{ in sentence B} \end{cases} \quad (13.2)$$

Position Embeddings: Learned absolute positions (not sinusoidal)

Special tokens:

- [CLS]: Start of sequence, used for classification
- [SEP]: Separate sentences
- [MASK]: Masked token for MLM
- [PAD]: Padding

Example 13.1 (BERT Input). Sentence pair: "The cat sat" and "It was tired"

Tokenized:

$$[\text{CLS}], \text{The}, \text{cat}, \text{sat}, [\text{SEP}], \text{It}, \text{was}, \text{tired}, [\text{SEP}] \quad (13.3)$$

Segment IDs:

$$[0, 0, 0, 0, 0, 1, 1, 1, 1] \quad (13.4)$$

Position IDs:

$$[0, 1, 2, 3, 4, 5, 6, 7, 8] \quad (13.5)$$

13.2 Pre-Training Objectives

13.2.1 Masked Language Modeling (MLM)

Definition 13.2 (Masked Language Modeling). Randomly mask 15% of tokens and predict them:

1. Select 15% of tokens
2. Of selected tokens:
 - 80%: Replace with [MASK]
 - 10%: Replace with random token
 - 10%: Keep original
3. Predict original tokens

Masked Language Modeling represents BERT’s core pre-training objective and the key innovation that enables bidirectional pre-training. Unlike traditional left-to-right language modeling where the model can only condition on previous tokens, MLM randomly masks a subset of input tokens and trains the model to predict the original tokens based on bidirectional context. This approach allows BERT to learn deep bidirectional representations that capture both left and right context simultaneously, which proves crucial for understanding tasks.

The masking strategy employs a carefully designed 80-10-10 split that prevents the model from simply memorizing the training data or developing trivial solutions. When a token is selected for masking (15% of all tokens), it is replaced with the special [MASK] token 80% of the time, creating the primary training signal. However, if the model only ever saw [MASK] tokens during training, it would learn to rely exclusively on this special token and fail to generalize during fine-tuning, where [MASK] tokens never appear. To address this mismatch, 10% of selected tokens are replaced with random tokens from the vocabulary, forcing the model to maintain robust representations even when the input contains noise. The remaining 10% of selected tokens are kept unchanged, encouraging the model to preserve meaningful representations for all tokens rather than only attending to masked positions.

Objective:

$$\mathcal{L}_{\text{MLM}} = - \sum_{i \in \mathcal{M}} \log P(x_i | \mathbf{x}_{\setminus \mathcal{M}}) \quad (13.6)$$

where \mathcal{M} is set of masked positions and $\mathbf{x}_{\setminus \mathcal{M}}$ are unmasked tokens.

The computational cost of MLM is substantial but manageable. For each training example, only 15% of tokens contribute to the loss, meaning the model must process approximately 6.7 times as many tokens as a standard language model to see the same number of training signals. However, this cost is offset by the ability to process entire sequences in parallel rather than autoregressively. For a sequence of length 512 with 15% masking, approximately 77 tokens are masked per example. The prediction head for each masked token performs a matrix multiplication with the vocabulary matrix $\mathbf{W}_{\text{vocab}} \in \mathbb{R}^{d_{\text{model}} \times V}$ followed by softmax over $V \approx 30,000$ tokens, requiring $2 \times 77 \times 768 \times 30,000 \approx 3.6$ billion FLOPs per example—a small fraction of the 96.6 billion FLOPs required for the BERT-base forward pass itself.

Example 13.2 (MLM Example). Original: "The cat sat on the mat"

Step 1: Select 15%: positions 2, 5

Step 2: Apply masking strategy:

- Position 2 ("cat"): Replace with [MASK] (80% case)
- Position 5 ("the"): Keep original (10% case)

Input: "The [MASK] sat on the mat"

Targets: Predict "cat" at position 2, "the" at position 5

Output layer:

$$\text{logits}_2 = \mathbf{h}_2 \mathbf{W}_{\text{vocab}} \quad \text{where } \mathbf{h}_2 \in \mathbb{R}^{768} \quad (13.7)$$

$$P(\text{token}|\text{position } 2) = \text{softmax}(\text{logits}_2) \quad (13.8)$$

Why this masking strategy?

- 80% [MASK]: Standard masking
- 10% random: Prevents over-reliance on [MASK] token
- 10% original: Encourages model to maintain representations

13.2.2 Next Sentence Prediction (NSP)

Definition 13.3 (Next Sentence Prediction). Binary classification: Does sentence B follow sentence A?

$$P(\text{IsNext} | [\text{CLS}]) = \sigma(\mathbf{W}_{\text{NSP}} \mathbf{h}_{[\text{CLS}]} + \mathbf{b}_{\text{NSP}}) \quad (13.9)$$

Training data:

- 50%: B actually follows A (label: IsNext)
- 50%: B is random sentence (label: NotNext)

NSP Loss:

$$\mathcal{L}_{\text{NSP}} = -\log P(y_{\text{NSP}} | [\text{CLS}]) \quad (13.10)$$

Total pre-training loss:

$$\mathcal{L} = \mathcal{L}_{\text{MLM}} + \mathcal{L}_{\text{NSP}} \quad (13.11)$$

Key Point 13.1. *Later work (RoBERTa) showed NSP provides minimal benefit. Modern models often use only MLM or variants like span corruption.*

Next Sentence Prediction was introduced as a secondary pre-training objective to help BERT learn relationships between sentence pairs, which is crucial for tasks like question answering and natural language inference. The task takes two sentences A and B, where 50% of the time B is the actual next sentence that follows A in the corpus, and 50% of the time B is a random sentence from elsewhere. The model must predict whether B genuinely follows A by using the representation of the special [CLS] token, which is prepended to every input sequence and serves as an aggregate sequence representation.

The motivation for NSP was that many downstream tasks require understanding relationships between sentence pairs, and pre-training on this objective should provide useful inductive biases. However, subsequent research, particularly the RoBERTa paper, demonstrated that NSP provides minimal or even negative benefit to downstream task performance. The likely explanation is that NSP is too easy a task—the model can often distinguish random sentences from consecutive sentences based on topic coherence alone, without learning the deeper inter-sentence relationships that would transfer to downstream tasks. RoBERTa removed NSP entirely and instead trained with longer sequences and dynamic masking, achieving substantially better performance. Modern BERT-style models typically omit NSP or replace it with more challenging objectives like sentence order prediction (used in ALBERT) or span corruption (used in T5).

13.3 Training Details and Computational Cost

13.3.1 Hardware and Training Configuration

BERT’s original training represents a landmark in demonstrating that large-scale pre-training on commodity hardware clusters could produce models with broad applicability. BERT-base was trained on 16 Cloud TPU v3 chips, where each chip contains 2 cores for a total of 32 TPU cores. However, the paper reports using 4 Cloud TPU v3 Pods with 16 chips each, giving 64 TPU cores total. Each TPU v3 core provides approximately 123 TFLOPS of bfloat16 performance, yielding a combined peak performance of roughly 7.9 PFLOPS for the full training setup. In practice, achieving even 40-50% utilization on such distributed training is considered excellent, suggesting an effective compute rate of 3-4 PFLOPS during training.

The training configuration uses a batch size of 256 sequences, each of maximum length 512 tokens, for a total of 131,072 tokens per batch. This large batch size is essential for stable training with the Adam optimizer and enables efficient utilization of the TPU hardware, which achieves best performance with large matrix operations. The learning rate schedule employs a warmup phase over the first 10,000 steps where the learning rate increases linearly from 0 to the peak learning rate of 1×10^{-4} , followed by linear decay to 0 over the remaining training steps. This warmup prevents the large initial gradients from destabilizing training, while the decay helps the model converge to a better optimum.

BERT-base training runs for 1 million steps with this configuration, requiring approximately 4 days of continuous training on the 64 TPU cores. Each training step processes 256 sequences of 512 tokens, so the total training data comprises $1,000,000 \times 256 \times 512 = 131$ billion tokens. The training corpus consists of BooksCorpus (800 million words) and English Wikipedia (2.5 billion words), totaling approximately 3.3 billion words or roughly 4.4 billion tokens after WordPiece tokenization. This means the model sees each token approximately 30 times during training, providing sufficient repetition for the model to learn robust representations while maintaining diversity through the random masking strategy.

13.3.2 Computational Cost Analysis

The computational cost of BERT pre-training can be estimated from the FLOPs per training step and the total number of steps. As analyzed in Chapter 12, BERT-base requires approximately 96.6 billion FLOPs for a forward pass on a single sequence of length 512. The backward pass requires roughly twice the FLOPs of the forward pass, giving approximately 290 billion FLOPs per training step per sequence. With batch size 256, each training step requires $290 \times 256 = 74.2$ trillion FLOPs. Over 1 million training steps, the total computation is approximately $74.2 \times 10^{15} \times 10^6 = 7.42 \times 10^{22}$ FLOPs, or 74.2 zettaFLOPs.

At an effective compute rate of 3.5 PFLOPS (assuming 45% utilization of the 7.9 PFLOPS peak), each training step takes approximately $74.2 \times 10^{12} / (3.5 \times 10^{15}) = 21.2$ milliseconds. Over 1 million steps, this yields $21.2 \times 10^6 = 21.2$ million milliseconds, or approximately 5.9 hours of pure compute time. However, the reported 4-day training time includes data loading, checkpointing, and distributed communication overhead, which can easily account for a factor of $16\times$ slowdown in practice. This discrepancy highlights the importance of efficient data pipelines and communication strategies in distributed training.

The estimated cost of BERT-base pre-training in 2018 was approximately \$7,000 using Google Cloud TPU pricing. This relatively modest cost (compared to later models like GPT-3, which cost millions of dollars to train) made BERT accessible to academic research groups and smaller companies, democratizing access to large-scale pre-trained models. The cost breaks down to roughly \$1.75 per hour for the TPU Pod configuration, and 4 days of training equals 96 hours, yielding $96 \times \$1.75 \approx \168 . However, the \$7,000 figure likely includes multiple training runs, hyperparameter tuning, and ablation studies rather than a single successful training run.

BERT-large requires substantially more computation due to its larger size. With 340 million parameters compared to BERT-base’s 110 million, and 24 layers instead of 12, BERT-large requires approximately 3-4 \times the FLOPs per training step. The original paper reports training BERT-large for 1 million steps on 64 TPU cores, taking approximately 4 days as well, though this likely involved more

aggressive optimization or different batch sizes. The estimated cost for BERT-large pre-training was around \$20,000-\$25,000, reflecting the increased computational requirements.

13.4 Fine-Tuning BERT

13.4.1 Fine-Tuning Procedure and Memory Requirements

Fine-tuning BERT for downstream tasks represents one of the model's key advantages: the pre-trained representations can be adapted to specific tasks with relatively small amounts of labeled data and modest computational resources. The fine-tuning process adds a task-specific head on top of the pre-trained BERT model and trains the entire model end-to-end on the target task. This approach typically requires only 2-4 epochs of training on task-specific data, taking minutes to hours on a single GPU depending on dataset size.

The memory requirements for fine-tuning are substantially lower than pre-training because fine-tuning typically uses smaller batch sizes and shorter sequences. For BERT-base fine-tuning on a classification task with batch size 16 and sequence length 128, the memory breakdown is approximately: 440 MB for model parameters (FP32), 440 MB for gradients, 880 MB for Adam optimizer states, and roughly 2 GB for activations. This totals approximately 3.8 GB, comfortably fitting on consumer GPUs with 8 GB of memory like the RTX 2080 or RTX 3070. Using mixed precision training reduces this to approximately 2.5 GB, enabling batch sizes of 32-48 on the same hardware.

The typical hyperparameter configuration for BERT fine-tuning uses a much smaller learning rate than pre-training to avoid catastrophically forgetting the pre-trained representations. Learning rates in the range of 2×10^{-5} to 5×10^{-5} work well for most tasks, with a linear warmup over 10% of training steps followed by linear decay. The batch size typically ranges from 16 to 32 for most tasks, though larger batch sizes (64-128) can improve performance on tasks with abundant training data. Fine-tuning for 2-4 epochs is usually sufficient, as longer training often leads to overfitting on the task-specific data.

13.4.2 Classification Tasks

For sequence classification (sentiment, topic, etc.):

1. Add classification head on [CLS] token

$$\text{logits} = \mathbf{W}_{\text{cls}} \mathbf{h}_{[\text{CLS}]} + \mathbf{b}_{\text{cls}} \quad (13.12)$$

2. Fine-tune entire model end-to-end

Example 13.3 (Sentiment Classification). Task: Binary sentiment (positive/negative)

Input: "This movie was amazing!" \rightarrow [CLS] This movie was amazing ! [SEP]

BERT encoding: $\mathbf{h}_{[\text{CLS}]} \in \mathbb{R}^{768}$

Classification head:

$$\text{logits} = \mathbf{W} \mathbf{h}_{[\text{CLS}]} + \mathbf{b} \quad \text{where } \mathbf{W} \in \mathbb{R}^{2 \times 768} \quad (13.13)$$

Prediction:

$$P(\text{positive}) = \text{softmax}(\text{logits})_1 \quad (13.14)$$

Fine-tuning: Train on labeled sentiment data for 2-4 epochs with small learning rate (2×10^{-5}).

13.4.3 Token-Level Tasks

For named entity recognition (NER), POS tagging:

1. Add classification head on each token

$$\text{logits}_i = \mathbf{W}_{\text{token}} \mathbf{h}_i + \mathbf{b}_{\text{token}} \quad (13.15)$$

2. Predict label for each token independently

13.4.4 Question Answering (SQuAD)

For span-based QA:

1. Input: [CLS] Question [SEP] Context [SEP]
2. Predict start and end positions in context

$$P_{\text{start}}(i) = \text{softmax}(\mathbf{h}_i^\top \mathbf{s}) \quad (13.16)$$

$$P_{\text{end}}(i) = \text{softmax}(\mathbf{h}_i^\top \mathbf{e}) \quad (13.17)$$

where $\mathbf{s}, \mathbf{e} \in \mathbb{R}^{768}$ are learned vectors.

13.5 BERT Variants

13.5.1 RoBERTa (Robustly Optimized BERT)

RoBERTa, introduced by Facebook AI Research in 2019, demonstrated that BERT was significantly undertrained and that careful attention to training procedures could yield substantial improvements without architectural changes. The key insight was that many of BERT’s design choices were not carefully ablated, and several modifications to the training procedure could dramatically improve performance. RoBERTa achieves state-of-the-art results on GLUE, SQuAD, and RACE benchmarks by training the same architecture as BERT-base and BERT-large with improved training procedures.

The first major change removes the Next Sentence Prediction objective entirely, training only with masked language modeling. Ablation studies showed that NSP actually hurt performance on downstream tasks, likely because the task is too simple and doesn’t provide useful training signal. Instead, RoBERTa uses full-length sequences of up to 512 tokens sampled contiguously from documents, allowing the model to learn longer-range dependencies without the artificial sentence-pair structure. This change alone improves downstream task performance by 1-2% across most benchmarks.

Dynamic masking represents another crucial improvement over BERT’s static masking. BERT generates masked training examples once during data preprocessing and uses the same masked version throughout all epochs of training. This means the model sees the exact same masked examples multiple times, potentially leading to memorization. RoBERTa instead generates new masking patterns every time a sequence is fed to the model, ensuring that the model never sees the exact same masked example twice across the entire training run. This dynamic masking provides more diverse training signal and improves generalization, particularly for longer training runs.

The training scale increases dramatically compared to BERT. RoBERTa uses batch sizes of 8,192 sequences (32× larger than BERT’s 256), enabled by gradient accumulation across multiple steps. Larger batch sizes improve training stability and allow for higher learning rates, accelerating convergence. The training data expands from BERT’s 16 GB of text (BooksCorpus + Wikipedia) to 160 GB, incorporating CC-News (76 GB), OpenWebText (38 GB), and Stories (31 GB) in addition to the original sources. This 10× increase in training data provides substantially more diverse examples for the model to learn from.

Perhaps most importantly, RoBERTa trains for much longer than BERT. While BERT-base trains for 1 million steps, RoBERTa trains for 500,000 steps with the larger batch size, corresponding to processing 4× more tokens overall. Some RoBERTa variants train for even longer, up to 1 million steps with the large batch size. This extended training allows the model to better learn the training objective and develop more robust representations. The computational cost of RoBERTa training is approximately 10-15× higher than BERT due to the combination of more data, larger batches, and longer training, estimated at \$50,000-\$100,000 for the full training run.

The results demonstrate the value of these improvements. RoBERTa-base matches or exceeds BERT-large performance on most tasks despite having the same architecture as BERT-base, and RoBERTa-large achieves new state-of-the-art results across GLUE, SQuAD 2.0, and RACE. The improvements are particularly pronounced on tasks requiring deep language understanding, such as natural language inference and reading comprehension, where RoBERTa-large improves over BERT-large by 2-4% absolute.

13.5.2 ALBERT (A Lite BERT)

ALBERT addresses BERT’s memory consumption and training time through parameter sharing and factorized embeddings, achieving comparable performance with far fewer parameters. The key insight is that many of BERT’s parameters are redundant, and careful parameter sharing can maintain model capacity while dramatically reducing memory requirements. ALBERT-xxlarge achieves similar performance to BERT-large with only 235 million parameters compared to BERT-large’s 340 million, and the parameter sharing enables training on hardware that couldn’t accommodate the full BERT-large model.

Factorized embedding parameterization represents the first major innovation. In BERT, the token embedding dimension equals the hidden dimension ($V \times d_{\text{model}}$), meaning the vocabulary matrix for BERT-base contains $30,000 \times 768 = 23$ million parameters. ALBERT observes that token embeddings are meant to learn context-independent representations, while hidden layers learn context-dependent representations, so these don’t need to share the same dimension. ALBERT instead projects the vocabulary to a smaller embedding dimension E (typically 128), then projects from E to d_{model} . This factorization reduces embedding parameters from $V \times d_{\text{model}}$ to $V \times E + E \times d_{\text{model}}$. For ALBERT-base with $E = 128$: $30,000 \times 128 + 128 \times 768 = 3,938,304$ parameters, an 83% reduction from BERT’s 23 million embedding parameters.

Cross-layer parameter sharing takes the reduction further by using the same parameters for all transformer layers. Instead of having unique parameters for each of the 12 or 24 layers, ALBERT uses a single set of layer parameters that is applied repeatedly. This reduces the transformer stack parameters by a factor of L (the number of layers). For ALBERT-base, the 12-layer transformer stack requires only 7.1 million parameters (one layer’s worth) instead of BERT-base’s 85 million parameters. The memory savings are dramatic: ALBERT-base contains only 12 million parameters total compared to BERT-base’s 110 million, an 89% reduction.

The performance impact of parameter sharing is surprisingly modest. Cross-layer parameter sharing causes a small performance degradation (typically 1-2% on downstream tasks) compared to models with unique layer parameters, but this is far less than the 89% parameter reduction would suggest. The likely explanation is that the transformer layers learn similar functions across depths, so sharing parameters doesn’t severely limit model capacity. The factorized embeddings actually improve performance slightly by preventing the model from overfitting to the vocabulary and encouraging more robust token representations.

ALBERT also replaces Next Sentence Prediction with Sentence Order Prediction (SOP), a more challenging task. Instead of distinguishing consecutive sentences from random sentences (which can often be solved by topic matching), SOP requires the model to distinguish the correct sentence order from the reversed order. Given sentences A and B that appear consecutively in the corpus, 50% of examples present them as A-B (positive) and 50% as B-A (negative). This task requires understanding inter-sentence coherence and cannot be solved by topic matching alone, providing a more useful training signal than NSP.

The memory and speed implications are substantial. ALBERT-base requires only 48 MB in FP32 (12 million parameters \times 4 bytes) compared to BERT-base’s 440 MB, enabling training with much larger batch sizes on the same hardware. However, parameter sharing doesn’t reduce computation—each layer still performs the same number of FLOPs, just with shared parameters. This means ALBERT-base has similar training time per step as BERT-base despite the parameter reduction. The primary benefit is memory efficiency, not computational efficiency. ALBERT-xxlarge, with $d_{\text{model}} = 4096$ and 12 layers, contains 235 million parameters but achieves performance comparable to BERT-large (340 million parameters) on most benchmarks.

13.5.3 DistilBERT: Knowledge Distillation for Compression

DistilBERT demonstrates that knowledge distillation can compress BERT to a fraction of its size while retaining most of its performance, making deployment feasible on resource-constrained devices. The core idea is to train a smaller “student” model to mimic the behavior of the larger “teacher” BERT model, transferring the knowledge encoded in the teacher’s parameters to the more compact student architecture. DistilBERT achieves 97% of BERT-base’s performance on GLUE while being 40% smaller and 60% faster, making it practical for production deployment where latency and memory are critical.

The DistilBERT architecture uses 6 transformer layers instead of BERT-base’s 12, with the same hidden dimension of 768 and 12 attention heads. This halving of depth reduces parameters from 110 million to approximately 66 million, a 40% reduction. The parameter savings come entirely from the transformer layers (6 layers \times 7.1M parameters = 42.6M instead of 85M), while the embedding layer remains unchanged at 23 million parameters. The reduced depth means each forward pass requires only 6 layers of computation instead of 12, directly translating to a $2\times$ speedup in the ideal case. In practice, the speedup is approximately $1.6\text{--}1.7\times$ due to overhead from embedding lookups and the final prediction layer, which don’t scale with depth.

The distillation loss combines three components to transfer knowledge from teacher to student. The first component is the standard cross-entropy loss between the student’s predictions and the true labels, ensuring the student learns the correct task. The second component is the distillation loss, which minimizes the KL divergence between the student’s output distribution and the teacher’s output distribution. The teacher’s “soft” predictions (probability distributions over the vocabulary) contain more information than the hard labels alone—for example, the teacher might assign 0.7 probability to the correct token, 0.2 to a synonym, and 0.1 to other tokens, revealing semantic relationships. The student learns from this richer signal by matching the teacher’s distribution.

Distillation loss:

$$\mathcal{L} = \alpha \mathcal{L}_{\text{CE}}(\text{student}, \text{labels}) + (1 - \alpha) \mathcal{L}_{\text{KD}}(\text{student}, \text{teacher}) \quad (13.18)$$

where:

$$\mathcal{L}_{\text{KD}} = \text{KL}(\text{softmax}(z_s/T) \parallel \text{softmax}(z_t/T)) \quad (13.19)$$

The temperature parameter T (typically 2-4) softens the probability distributions, amplifying the differences between high-probability and low-probability tokens. Higher temperatures make the distributions more uniform, providing more information about the teacher’s uncertainty. The coefficient α (typically 0.5) balances the two objectives. The third component (not shown in the simplified equation) is a cosine embedding loss that encourages the student’s hidden states to align with the teacher’s hidden states, providing additional supervision beyond the output distributions.

The training procedure initializes the student by taking every other layer from the pre-trained BERT teacher, providing a warm start that accelerates convergence. The student is then trained on the same masked language modeling task as BERT, but with the teacher’s soft targets providing additional supervision. Training DistilBERT requires approximately 90 hours on 8 NVIDIA V100 GPUs, substantially less than the 4 days on 64 TPU cores required for BERT-base pre-training. The reduced training time reflects both the smaller model size and the benefit of learning from the teacher’s predictions rather than learning from scratch.

The performance-efficiency trade-off is compelling for production deployment. DistilBERT retains 97% of BERT-base’s performance on GLUE tasks, with the largest drops on tasks requiring deep reasoning (2-3% degradation on MNLI and QQP) and smaller drops on simpler tasks (0-1% on SST-2 and MRPC). The 40% parameter reduction translates directly to memory savings: DistilBERT requires 264 MB in FP32 compared to BERT-base’s 440 MB, enabling deployment on mobile devices and edge hardware. The 60% speedup ($1.6\text{--}1.7\times$ faster) reduces inference latency from approximately 14 ms to 8-9 ms per sequence on an NVIDIA V100, crucial for real-time applications.

13.5.4 Memory and Speed Comparisons

Comparing BERT variants across memory footprint and inference speed reveals clear trade-offs between model capacity and deployment efficiency. The following analysis uses BERT-base as the baseline and

measures relative performance on GLUE benchmark tasks, memory consumption in FP32, and inference throughput on an NVIDIA V100 GPU with batch size 1 and sequence length 128.

Parameter counts and memory:

- **BERT-base:** 110M parameters, 440 MB (FP32), 220 MB (FP16)
- **BERT-large:** 340M parameters, 1.36 GB (FP32), 680 MB (FP16)
- **RoBERTa-base:** 125M parameters, 500 MB (FP32), 250 MB (FP16) — slightly larger due to different vocabulary
- **RoBERTa-large:** 355M parameters, 1.42 GB (FP32), 710 MB (FP16)
- **ALBERT-base:** 12M parameters, 48 MB (FP32), 24 MB (FP16) — 89% reduction from BERT-base
- **ALBERT-xxlarge:** 235M parameters, 940 MB (FP32), 470 MB (FP16)
- **DistilBERT:** 66M parameters, 264 MB (FP32), 132 MB (FP16) — 40% reduction from BERT-base

Inference speed (sequences/second on V100, batch size 1, sequence length 128):

- **BERT-base:** ~140 sequences/sec (baseline)
- **BERT-large:** ~50 sequences/sec (2.8× slower due to larger dimensions and more layers)
- **RoBERTa-base:** ~135 sequences/sec (similar to BERT-base)
- **ALBERT-base:** ~140 sequences/sec (same speed as BERT-base despite fewer parameters—computation unchanged)
- **DistilBERT:** ~230 sequences/sec (1.6× faster due to half the layers)

Performance on GLUE (average score):

- **BERT-base:** 78.5 (baseline)
- **BERT-large:** 82.1 (+3.6 points)
- **RoBERTa-base:** 80.4 (+1.9 points — matches BERT-large with base architecture)
- **RoBERTa-large:** 84.3 (+5.8 points — new state-of-the-art)
- **ALBERT-base:** 77.2 (-1.3 points — slight degradation from parameter sharing)
- **ALBERT-xxlarge:** 82.3 (+3.8 points — matches BERT-large with fewer parameters)
- **DistilBERT:** 76.1 (-2.4 points — retains 97% of performance)

These comparisons reveal distinct use cases for each variant. RoBERTa offers the best performance when computational resources are available, making it ideal for offline processing and high-accuracy applications. ALBERT provides excellent memory efficiency for training and deployment on memory-constrained hardware, though it doesn't reduce inference time. DistilBERT offers the best balance of speed and performance for latency-sensitive applications like real-time search and interactive systems. BERT-base remains a strong baseline that balances performance, speed, and memory for most applications.

13.6 Hardware Requirements and Deployment

13.6.1 GPU Memory Requirements

Understanding GPU memory requirements is essential for selecting appropriate hardware for BERT training and inference. The memory footprint depends on whether the model is being trained or used for inference, the batch size, sequence length, and numerical precision. For training, memory must accommodate model parameters, gradients, optimizer states, and activations, while inference requires only parameters and a single forward pass of activations.

BERT-base training in FP32 with batch size 32 and sequence length 512 requires approximately 13.8 GB of memory, as detailed in Chapter 12. This breaks down to 440 MB for parameters, 440 MB for gradients, 880 MB for Adam optimizer states (first and second moments), and approximately 12 GB for activations across the 12 layers. This memory requirement fits comfortably on NVIDIA V100 GPUs with 16 GB of memory, RTX 3090 GPUs with 24 GB, or A100 GPUs with 40-80 GB. However, the batch size must be reduced for GPUs with less memory—an RTX 3080 with 10 GB can train BERT-base with batch size 16-20, while an RTX 3070 with 8 GB is limited to batch size 8-12.

Mixed precision training with FP16 or BF16 dramatically reduces memory consumption by storing activations in 16-bit format while maintaining FP32 master weights for numerical stability. For BERT-base, mixed precision reduces total memory from 13.8 GB to approximately 8 GB, enabling batch size 32 on GPUs with 12 GB of memory or batch size 64 on GPUs with 24 GB. The memory savings come primarily from activations, which are halved from 12 GB to 6 GB, while the parameter and optimizer memory increases slightly from 1.76 GB to 1.98 GB due to maintaining both FP16 and FP32 copies of parameters. Gradient checkpointing can further reduce memory by recomputing activations during the backward pass rather than storing them, reducing activation memory by approximately 80% at the cost of 20-30% longer training time.

BERT-large training requires substantially more memory due to its larger hidden dimension and greater depth. In FP32 with batch size 32 and sequence length 512, BERT-large requires approximately 32 GB of memory: 1.36 GB for parameters, 1.36 GB for gradients, 2.72 GB for optimizer states, and approximately 26 GB for activations. This necessitates GPUs with at least 32 GB of memory, such as the NVIDIA V100 (32 GB variant) or A100 (40-80 GB). Mixed precision reduces this to approximately 18 GB, enabling training on A100 40 GB GPUs with batch size 16-24. Consumer GPUs like the RTX 3090 (24 GB) can train BERT-large with mixed precision and gradient checkpointing at batch size 8-12, though training time increases significantly.

Inference memory requirements are far more modest because they don't include gradients, optimizer states, or stored activations for backpropagation. BERT-base inference in FP32 requires only 440 MB for parameters plus approximately 200-300 MB for a single forward pass of activations with batch size 1, totaling under 1 GB. In FP16, this drops to approximately 400 MB total, enabling BERT-base inference on virtually any GPU, including mobile GPUs and edge devices. BERT-large inference requires approximately 1.5 GB in FP32 or 800 MB in FP16, still easily fitting on consumer GPUs. The primary consideration for inference is batch size: larger batch sizes improve throughput but increase activation memory linearly. A V100 with 16 GB can run BERT-base inference with batch size 128-256 in FP16, achieving throughput of 15,000-20,000 sequences per second.

13.6.2 Batch Size Limits by GPU Type

The maximum batch size for BERT training varies significantly across GPU types, directly impacting training throughput and efficiency. Larger batch sizes improve GPU utilization by amortizing memory bandwidth costs across more computation, but they're limited by available memory. The following analysis assumes mixed precision training with sequence length 512 for BERT-base and BERT-large.

BERT-base maximum batch sizes (mixed precision, sequence length 512):

- **RTX 3060 (12 GB):** Batch size 24-28 without gradient checkpointing, 48-56 with checkpointing
- **RTX 3070 (8 GB):** Batch size 14-18 without gradient checkpointing, 28-36 with checkpointing
- **RTX 3080 (10 GB):** Batch size 18-22 without gradient checkpointing, 36-44 with checkpointing

- **RTX 3090 (24 GB):** Batch size 48-56 without gradient checkpointing, 96-112 with checkpointing
- **V100 (16 GB):** Batch size 28-32 without gradient checkpointing, 56-64 with checkpointing
- **V100 (32 GB):** Batch size 64-72 without gradient checkpointing, 128-144 with checkpointing
- **A100 (40 GB):** Batch size 80-96 without gradient checkpointing, 160-192 with checkpointing
- **A100 (80 GB):** Batch size 160-192 without gradient checkpointing, 320-384 with checkpointing

BERT-large maximum batch sizes (mixed precision, sequence length 512):

- **RTX 3090 (24 GB):** Batch size 12-16 without gradient checkpointing, 24-32 with checkpointing
- **V100 (16 GB):** Batch size 6-8 without gradient checkpointing, 12-16 with checkpointing
- **V100 (32 GB):** Batch size 16-20 without gradient checkpointing, 32-40 with checkpointing
- **A100 (40 GB):** Batch size 20-24 without gradient checkpointing, 40-48 with checkpointing
- **A100 (80 GB):** Batch size 48-56 without gradient checkpointing, 96-112 with checkpointing

These batch size limits have direct implications for training efficiency. Smaller batch sizes reduce GPU utilization because the model spends more time on memory transfers relative to computation. For BERT-base on an RTX 3070 with batch size 16, GPU utilization typically reaches only 50-60% of peak FLOPS, while an A100 with batch size 96 can achieve 70-80% utilization. Gradient accumulation can simulate larger batch sizes by accumulating gradients over multiple forward-backward passes before updating parameters, enabling effective batch sizes of 128-256 even on GPUs limited to batch size 16-32 per step. However, gradient accumulation increases training time proportionally to the accumulation steps.

13.6.3 Inference Speed Analysis

Inference speed determines the feasibility of deploying BERT in production systems where latency and throughput are critical. We measure inference speed in two ways: latency (time per sequence for batch size 1, important for interactive applications) and throughput (sequences per second for large batches, important for offline processing). The following measurements use sequence length 128, which is typical for many classification and NER tasks, and FP16 precision on NVIDIA GPUs.

BERT-base inference latency (batch size 1, sequence length 128):

- **V100:** 7.2 ms per sequence (139 sequences/sec)
- **A100:** 3.8 ms per sequence (263 sequences/sec) — 1.9× faster than V100
- **RTX 3090:** 5.1 ms per sequence (196 sequences/sec)
- **CPU (Intel Xeon Gold 6248):** 45-60 ms per sequence (17-22 sequences/sec) — 10-15× slower than GPU

BERT-base inference throughput (batch size 128, sequence length 128):

- **V100:** 18,000-20,000 sequences/sec
- **A100:** 35,000-40,000 sequences/sec — 2× faster than V100
- **RTX 3090:** 25,000-28,000 sequences/sec

BERT-large inference latency (batch size 1, sequence length 128):

- **V100:** 20.5 ms per sequence (49 sequences/sec)
- **A100:** 10.8 ms per sequence (93 sequences/sec) — 1.9× faster than V100

- **RTX 3090:** 14.2 ms per sequence (70 sequences/sec)

BERT-large inference throughput (batch size 64, sequence length 128):

- **V100:** 6,500-7,500 sequences/sec
- **A100:** 12,000-14,000 sequences/sec — $1.9\times$ faster than V100
- **RTX 3090:** 9,000-10,500 sequences/sec

The A100’s superior performance comes from its higher memory bandwidth (1.6 TB/s vs V100’s 900 GB/s) and more powerful Tensor Cores (312 TFLOPS FP16 vs V100’s 125 TFLOPS). For BERT inference, which is often memory-bandwidth bound due to loading model parameters, the A100’s bandwidth advantage is particularly valuable. The approximately $2\times$ speedup of A100 over V100 holds across different batch sizes and model sizes, making the A100 the preferred choice for production BERT deployment when latency is critical.

Sequence length significantly impacts inference speed due to the quadratic scaling of attention computation. For BERT-base on a V100, increasing sequence length from 128 to 512 ($4\times$ longer) increases latency from 7.2 ms to approximately 18 ms ($2.5\times$ slower), less than the $4\times$ that pure quadratic scaling would suggest because the feed-forward network and embedding layers don’t scale quadratically. For very long sequences approaching the 512 token maximum, attention computation dominates and the scaling approaches quadratic. This explains why efficient attention mechanisms (Chapter 16) focus on reducing the $O(n^2)$ attention complexity for long-context applications.

13.7 Analysis and Interpretability

13.7.1 What BERT Learns

Lower layers: Syntactic information (POS tags, parse trees)

Middle layers: Semantic information (word sense, entity types)

Upper layers: Task-specific information

Attention patterns:

- Some heads attend to next token (language modeling pattern)
- Some heads attend to syntactic relations (e.g., verbs to subjects)
- Some heads attend broadly (averaging)

13.7.2 Probing Tasks

Test what linguistic information is encoded:

- Surface: Sentence length, word order
- Syntactic: POS tags, dependency labels, constituency trees
- Semantic: Named entities, semantic roles, coreference

Method: Train linear classifier on frozen BERT representations

Result: BERT captures surprisingly rich linguistic structure!

13.8 Exercises

Exercise 13.1. Implement masked language modeling. For sentence “The quick brown fox jumps”, mask 15% of tokens and compute MLM loss. Show prediction probabilities for masked positions.

Exercise 13.2. Fine-tune BERT-base on binary classification with 10,000 examples. Compare learning curves for: (1) Training only classification head, (2) Fine-tuning all layers. Which converges faster? Which achieves better performance?

Exercise 13.3. Compare parameter counts for BERT-base, RoBERTa-base, ALBERT-base, DistilBERT. For each, calculate: (1) Total parameters, (2) Memory footprint (FP32), (3) Inference FLOPs for sequence length 128.

Exercise 13.4. Visualize attention patterns for multi-head attention in BERT. For sentence "The cat that chased the mouse ran away", identify heads that capture: (1) Adjacent words, (2) Subject-verb relations, (3) Long-range dependencies.

Chapter 14

GPT: Generative Pre-Training

Chapter Overview

GPT (Generative Pre-trained Transformer) pioneered decoder-only transformer architectures for autoregressive language modeling. This chapter traces the evolution from GPT-1 through GPT-4, covering architecture, pre-training, scaling, few-shot learning, and emergent abilities.

Learning Objectives

1. Understand GPT's decoder-only architecture
2. Implement autoregressive language modeling
3. Apply in-context learning and few-shot prompting
4. Analyze scaling laws and emergent abilities
5. Compare GPT variants (GPT-1, GPT-2, GPT-3, GPT-4)
6. Understand instruction tuning and RLHF

14.1 GPT Architecture

14.1.1 Decoder-Only Transformers

The GPT architecture represents a fundamental departure from the encoder-decoder paradigm that dominated sequence-to-sequence models. Rather than using separate encoder and decoder stacks, GPT employs only transformer decoder blocks, creating a purely autoregressive language model. This architectural choice has profound implications for both the model's capabilities and its computational characteristics.

The core innovation lies in the attention mechanism's masking pattern. GPT uses causal masking, which prevents each position from attending to future positions in the sequence. Mathematically, when computing attention scores $\mathbf{S} = \mathbf{QK}^\top$, a mask is applied such that $S_{ij} = -\infty$ for all $j > i$. After the softmax operation, these masked positions have zero attention weight, ensuring that the representation at position i depends only on tokens at positions 1 through i . This causal constraint is essential for autoregressive generation, where the model must predict the next token without access to future context.

Unlike the original transformer architecture which included cross-attention layers to attend from decoder to encoder, GPT eliminates cross-attention entirely. Each decoder block contains only a masked self-attention layer followed by a position-wise feed-forward network. This simplification reduces architectural complexity while maintaining the transformer's parallel processing advantages. The self-attention layer allows each position to gather information from all previous positions simultaneously, avoiding the sequential bottleneck of recurrent networks.

GPT-2 and later versions introduced an important architectural refinement: pre-normalization. Rather than applying layer normalization after each sub-layer (post-norm), pre-norm applies normalization before the attention and feed-forward operations. This seemingly minor change significantly improves training stability for deep networks. In the pre-norm configuration, the residual path carries the original signal without normalization, providing a clean gradient path during backpropagation. This enables training of much deeper models without the gradient instability that plagued earlier architectures.

Definition 14.1 (GPT Architecture). GPT uses transformer decoder blocks with:

- **Masked self-attention:** Causal masking (no future tokens)
- **No cross-attention:** Decoder-only (vs encoder-decoder)
- **Position-wise FFN:** Same as standard transformer
- **Pre-norm:** Layer norm before sub-layers (GPT-2+)

The distinction between GPT and BERT architectures illuminates different modeling philosophies. BERT employs bidirectional attention, allowing each position to attend to the entire sequence including future tokens. This bidirectionality enables rich contextual representations ideal for understanding tasks like classification and question answering. However, bidirectional attention is incompatible with autoregressive generation—the model cannot predict the next token if it has already seen it. GPT’s unidirectional causal attention sacrifices bidirectional context but gains the ability to generate coherent text autoregressively. This trade-off reflects the fundamental tension between understanding (BERT) and generation (GPT) in language modeling.

14.1.2 GPT Model Sizes

The evolution of GPT models demonstrates the remarkable scaling properties of transformer architectures. Each generation increased model capacity by orders of magnitude, revealing new capabilities that emerged only at larger scales. Understanding the progression from GPT-1 through GPT-3 provides insight into the relationship between model size and performance.

GPT-1, introduced in 2018, established the decoder-only pre-training paradigm with 117 million parameters. The architecture used 12 transformer layers with hidden dimension $d = 768$ and 12 attention heads, processing sequences up to 512 tokens. While modest by today’s standards, GPT-1 demonstrated that unsupervised pre-training on large text corpora followed by task-specific fine-tuning could achieve strong performance across diverse NLP tasks. The model was trained on BookCorpus, a dataset of approximately 7,000 unpublished books containing 800 million words. This training data, while substantial for 2018, would be considered quite limited compared to later models.

GPT-2, released in 2019, expanded the scaling experiment by training four model sizes ranging from 117 million to 1.5 billion parameters. The smallest GPT-2 matched GPT-1’s architecture, while GPT-2 XL scaled to 48 layers with hidden dimension $d = 1600$ and 25 attention heads. The context window doubled to 1024 tokens, enabling the model to maintain coherence over longer passages. More significantly, GPT-2 was trained on WebText, a dataset of 40 GB containing 8 million web pages. This diverse training data, scraped from outbound links on Reddit with at least 3 karma, provided much broader coverage of topics and writing styles than BookCorpus. GPT-2’s key finding was that larger models trained on more diverse data could perform many tasks zero-shot, without any task-specific fine-tuning—a surprising emergent capability.

GPT-3, unveiled in 2020, represented a massive leap to 175 billion parameters. The architecture scaled to 96 layers with hidden dimension $d = 12288$ and 96 attention heads, processing sequences of 2048 tokens. The parameter count increased by more than $100\times$ compared to GPT-2 XL, requiring fundamentally different training infrastructure. GPT-3 was trained on approximately 300 billion tokens drawn from Common Crawl (filtered), WebText2, Books1, Books2, and Wikipedia, totaling roughly 570

GB of text. The training used a single pass through this massive dataset rather than multiple epochs, reflecting the compute-optimal insight that data diversity matters more than repeated exposure to the same examples. GPT-3's most striking capability was few-shot learning: the model could perform new tasks by conditioning on a few examples in the prompt, without any parameter updates. This in-context learning ability scaled dramatically with model size, with GPT-3 175B far outperforming smaller variants.

GPT-4, released in 2023, marked another architectural evolution, though OpenAI disclosed fewer details. Estimates suggest the model uses a mixture-of-experts architecture with 1 to 1.7 trillion total parameters, though only a fraction are active for any given input. The context window expanded dramatically to 8,192 tokens in the standard version and 32,768 tokens in the extended version, enabling the model to process entire documents or codebases. GPT-4 demonstrated significant improvements in reasoning, factual accuracy, and instruction following, suggesting that architectural innovations beyond pure parameter scaling contributed to its capabilities.

GPT-1 (2018):

- Layers: $L = 12$, Hidden: $d = 768$, Heads: $h = 12$
- Parameters: 117M
- Context: 512 tokens

GPT-2 (2019):

- Small: 117M, Medium: 345M, Large: 762M, XL: 1.5B
- GPT-2 XL: $L = 48$, $d = 1600$, $h = 25$
- Context: 1024 tokens

GPT-3 (2020):

- Small: 125M to XL: 175B
- GPT-3 175B: $L = 96$, $d = 12288$, $h = 96$
- Context: 2048 tokens
- Parameters: 175 billion!

GPT-4 (2023):

- Architecture details not fully disclosed
- Estimated: 1-1.7 trillion parameters (mixture of experts)
- Context: 8K (standard), 32K (extended)

Example 14.1 (GPT-2 Small Layer). Configuration: $L = 12$, $d = 768$, $h = 12$, $d_{ff} = 3072$

Understanding the parameter breakdown of GPT-2 Small reveals how transformer capacity is distributed across different components. Each of the 12 decoder layers contains approximately 7 million parameters, with the feed-forward network consuming roughly two-thirds of this total. This distribution reflects the architectural choice to use an expansion factor of 4 in the FFN, where the hidden dimension $d_{ff} = 4 \times d_{\text{model}} = 3072$.

Single decoder layer:

1. Layer norm
2. Masked multi-head attention (12 heads)
3. Residual connection
4. Layer norm

5. Feed-forward ($768 \rightarrow 3072 \rightarrow 768$)

6. Residual connection

The masked multi-head attention mechanism requires four weight matrices: \mathbf{W}^Q , \mathbf{W}^K , \mathbf{W}^V for projecting to query, key, and value spaces, and \mathbf{W}^O for projecting the concatenated head outputs back to model dimension. Each of these matrices has dimensions 768×768 , contributing $4 \times 768^2 = 2,359,296$ parameters. The feed-forward network contains two linear transformations: the first expands from 768 to 3072 dimensions ($768 \times 3072 = 2,359,296$ parameters), and the second projects back from 3072 to 768 dimensions (another $768 \times 3072 = 2,359,296$ parameters), totaling 4,718,592 parameters. Layer normalization adds minimal parameters—just scale and bias terms for each dimension, contributing $2 \times 2 \times 768 = 3,072$ parameters across the two layer norms per block.

Parameters per layer:

$$\text{Attention: } 4 \times 768^2 = 2,359,296 \quad (14.1)$$

$$\text{FFN: } 2 \times 768 \times 3072 = 4,718,592 \quad (14.2)$$

$$\text{Layer norms: } 2 \times 2 \times 768 = 3,072 \quad (14.3)$$

$$\text{Total: } 7,080,960 \approx 7M \quad (14.4)$$

Multiplying by 12 layers yields approximately 85 million parameters in the transformer blocks. The remaining 32 million parameters reside in the token embeddings, which map the vocabulary (typically 50,257 tokens for GPT-2) to the 768-dimensional model space. This embedding matrix alone contains $50,257 \times 768 = 38,597,376$ parameters, though the actual vocabulary size may vary slightly. Position embeddings add another $1024 \times 768 = 786,432$ parameters for the maximum sequence length of 1024 tokens. The final layer norm and output projection (which often shares weights with the token embedding) complete the 117 million parameter total.

12 layers: $\approx 85M$, plus embeddings $\approx 32M = \mathbf{117M \text{ total}}$

14.2 Pre-Training: Autoregressive Language Modeling

14.2.1 Training Objective

Autoregressive language modeling forms the foundation of GPT’s pre-training approach. Unlike masked language modeling used in BERT, which predicts randomly masked tokens using bidirectional context, autoregressive modeling predicts each token based solely on preceding tokens. This objective aligns naturally with text generation tasks and enables the model to learn the statistical structure of language through next-token prediction.

The training objective maximizes the likelihood of each token given all previous tokens in the sequence. For a sequence $\mathbf{x} = [x_1, x_2, \dots, x_n]$, the model learns to maximize the joint probability $P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i | x_1, \dots, x_{i-1})$. Taking the logarithm converts this product into a sum, yielding the standard language modeling loss. This formulation has an elegant interpretation: the model learns to compress the training data by assigning high probability to observed sequences, with the negative log-likelihood measuring the number of bits required to encode the data under the model’s distribution.

Definition 14.2 (Autoregressive Language Modeling). Maximize likelihood of next token given previous context:

$$\mathcal{L} = \sum_{i=1}^n \log P(x_i | x_1, \dots, x_{i-1}; \theta) \quad (14.5)$$

The implementation leverages the transformer’s parallel processing capabilities through teacher forcing. Rather than generating tokens sequentially during training, the entire sequence is processed in

a single forward pass. The input sequence $[x_1, x_2, \dots, x_n]$ is fed to the model, which produces hidden representations for all positions simultaneously. The causal attention mask ensures that position i cannot attend to positions $j > i$, maintaining the autoregressive property despite parallel computation. The model's output at position i is trained to predict token x_{i+1} , creating $n - 1$ training signals from a single sequence of length n . This parallel training is dramatically more efficient than sequential generation, enabling large-scale pre-training on massive text corpora.

The cross-entropy loss is computed at each position by comparing the model's predicted distribution over the vocabulary with the true next token. For position i with hidden state \mathbf{h}_i , the model computes logits $\mathbf{z}_i = \mathbf{h}_i \mathbf{W}_{\text{out}}$ where $\mathbf{W}_{\text{out}} \in \mathbb{R}^{d_{\text{model}} \times V}$ projects to vocabulary size V . Applying softmax yields a probability distribution $P(x_{i+1}|x_1, \dots, x_i) = \text{softmax}(\mathbf{z}_i)$. The loss for this position is $-\log P(x_{i+1}|x_1, \dots, x_i)$, and the total loss sums over all positions. This formulation naturally handles variable-length sequences and provides dense training signal from every token in the corpus.

Implementation:

1. Input: $[x_1, x_2, \dots, x_n]$
2. Target: $[x_2, x_3, \dots, x_{n+1}]$ (shifted by 1)
3. Causal mask: Position i cannot attend to $j > i$
4. Cross-entropy loss at each position

Example 14.2 (GPT Training Example). Sentence: "The cat sat on the mat"

Tokenized: $[T_1, T_2, T_3, T_4, T_5, T_6] = [\text{The}, \text{cat}, \text{sat}, \text{on}, \text{the}, \text{mat}]$

This simple example illustrates how GPT processes a sequence during training. The model receives the tokenized sequence as input and must predict each subsequent token based on the preceding context. At position 1, having seen only "The", the model predicts "cat". At position 2, with context "The cat", it predicts "sat". This continues through the sequence, with each position providing a training signal. The beauty of teacher forcing is that all these predictions occur in parallel during a single forward pass, despite the autoregressive dependency structure.

Training:

$$P(T_2|T_1) = \text{softmax}(\mathbf{h}_1 \mathbf{W}_{\text{out}}) \quad \text{predict "cat"} \quad (14.6)$$

$$P(T_3|T_1, T_2) = \text{softmax}(\mathbf{h}_2 \mathbf{W}_{\text{out}}) \quad \text{predict "sat"} \quad (14.7)$$

$$\vdots \quad (14.8)$$

$$P(T_6|T_1, \dots, T_5) = \text{softmax}(\mathbf{h}_5 \mathbf{W}_{\text{out}}) \quad \text{predict "mat"} \quad (14.9)$$

The loss function sums the negative log-probabilities of the correct tokens at each position. If the model assigns high probability to the correct next token, the loss is low; if it assigns low probability, the loss is high. During backpropagation, gradients flow through all positions simultaneously, updating the model parameters to increase the probability of observed sequences. This dense training signal from every token in the corpus enables efficient learning of language statistics.

Loss:

$$\mathcal{L} = - \sum_{i=1}^5 \log P(T_{i+1}|T_1, \dots, T_i) \quad (14.10)$$

All positions trained simultaneously in parallel (teacher forcing)!

14.2.2 Pre-Training Data

The scale and diversity of pre-training data have proven critical to GPT's capabilities. Each generation of GPT models trained on progressively larger and more diverse text corpora, revealing that data quality and quantity both matter significantly for downstream performance.

GPT-1 was trained on BooksCorpus, a collection of approximately 7,000 unpublished books from various genres including adventure, fantasy, and romance. This dataset contained roughly 800 million words, providing coherent long-form text that helped the model learn narrative structure and long-range dependencies. The choice of books as training data reflected the hypothesis that long-form text with coherent structure would be more valuable than shorter, disconnected documents. However, the relatively narrow domain coverage limited the model’s exposure to diverse topics and writing styles.

GPT-2 marked a significant shift in data philosophy with the creation of WebText, a dataset of 40 GB containing text from 8 million web pages. The data was collected by scraping outbound links from Reddit posts with at least 3 karma, using social curation as a quality filter. This approach yielded much more diverse content spanning news articles, tutorials, discussions, and creative writing across virtually all topics. The 40 GB corpus represented approximately 10 billion tokens, more than an order of magnitude larger than BooksCorpus. This scale and diversity enabled GPT-2 to demonstrate surprising zero-shot capabilities on tasks it had never been explicitly trained to perform.

GPT-3 scaled data collection to unprecedented levels, training on approximately 300 billion tokens drawn from multiple sources. The training mixture included Common Crawl (filtered to remove low-quality content), WebText2 (an expanded version of GPT-2’s dataset), Books1, Books2, and Wikipedia. The total dataset size reached roughly 570 GB of text. Critically, GPT-3 was trained for a single epoch over this massive dataset rather than multiple passes over smaller data. This decision reflected emerging understanding of scaling laws: given fixed compute budget, it is often better to train on more diverse data once than to repeatedly train on the same limited data. The single-epoch approach also reduced the risk of memorizing specific training examples, though concerns about data contamination and memorization remained.

The composition of GPT-3’s training data was carefully weighted, with higher-quality sources sampled more frequently. Common Crawl, despite being the largest source, was downweighted due to quality concerns, while Wikipedia and books received higher sampling rates. This weighting scheme balanced scale with quality, ensuring the model learned from both broad web text and curated high-quality sources. The exact mixing ratios and filtering procedures significantly impacted model performance, though these details were not fully disclosed.

GPT-1: BooksCorpus (7,000 books, \approx 800M words)

GPT-2: WebText (40GB, 8M web pages)

GPT-3: Common Crawl (filtered), WebText2, Books1, Books2, Wikipedia

- Total: \approx 570GB text
- Tokens: \approx 300 billion
- Training: Single pass (not multiple epochs)

14.2.3 Training Infrastructure and Costs

The computational requirements for training GPT models reveal the massive scale of resources needed for state-of-the-art language models. Understanding these requirements is essential for practitioners considering whether to train models from scratch or use pre-trained models.

GPT-2’s training represented a significant but manageable computational investment. The 1.5 billion parameter XL model was trained on 32 TPU v3 chips for approximately one week. TPU v3 chips provide roughly 420 TFLOPS of bfloat16 performance, giving a total cluster capacity of about 13.4 PFLOPS. The training cost was estimated at approximately \$50,000, making it accessible to well-funded research labs and companies but beyond the reach of individual researchers. The batch size was carefully tuned to maximize GPU utilization while maintaining training stability, typically using batch sizes of 512 sequences with 1024 tokens each. The learning rate followed a cosine decay schedule with warmup, starting from a small value to prevent early training instability.

The relatively modest infrastructure requirements for GPT-2 meant that the model could be trained on a single multi-GPU machine or small cluster. This accessibility contributed to widespread experimentation with the architecture and training approach. Researchers could reproduce the training process,

fine-tune on domain-specific data, and explore architectural variations without requiring massive computational resources. The one-week training time also enabled rapid iteration on hyperparameters and training procedures.

GPT-3’s training requirements increased by multiple orders of magnitude, necessitating infrastructure available only to the largest technology companies. The 175 billion parameter model required an estimated 10,000 or more V100 GPUs for approximately one month of training. Each V100 provides 125 TFLOPS of FP16 performance, yielding a total cluster capacity exceeding 1 exaFLOP. The training cost was estimated between \$4 million and \$12 million depending on cloud pricing and hardware utilization assumptions. The energy consumption reached approximately 1,287 MWh, equivalent to the annual electricity usage of over 100 average US households.

This massive scale was necessary for several reasons. First, the 175 billion parameters required substantial memory even with model parallelism—the model weights alone occupy 700 GB in FP32 or 350 GB in FP16. Second, processing 300 billion tokens through such a large model requires enormous computational throughput. Third, maintaining reasonable training time (one month rather than one year) demanded massive parallelism across thousands of GPUs. The training employed sophisticated distributed strategies including data parallelism, model parallelism, and pipeline parallelism to efficiently utilize the hardware.

The infrastructure challenges extended beyond raw compute. Network bandwidth between GPUs became critical, as model parallelism requires frequent communication of activations and gradients. High-bandwidth interconnects like NVIDIA NVLink and InfiniBand were essential for maintaining efficiency. Memory bandwidth also constrained performance, as loading model parameters and activations from GPU memory often became the bottleneck rather than arithmetic operations. These hardware considerations influenced architectural choices, such as the FFN expansion factor and attention head dimensions.

The enormous cost of training GPT-3 has profound implications for the AI research ecosystem. Only a handful of organizations can afford to train models at this scale, concentrating power and capability among well-resourced entities. This has motivated research into more efficient training methods, better scaling laws to predict optimal model sizes, and techniques for adapting pre-trained models to new domains without full retraining. The Chinchilla findings, discussed later, suggest that GPT-3 was actually over-parameterized for its training compute budget—a smaller model trained on more data would have achieved better performance for the same cost.

GPT-2 Training:

- Hardware: 32 TPU v3 chips (≈ 13.4 PFLOPS)
- Training time: ≈ 1 week
- Cost: $\approx \$50,000$
- Batch size: 512 sequences \times 1024 tokens
- Learning rate: Cosine decay with warmup

GPT-3 Training:

- Hardware: 10,000+ V100 GPUs (estimated, > 1 exaFLOP)
- Training time: ≈ 1 month
- Cost: \$4-12 million (estimated)
- Energy consumption: 1,287 MWh
- Requires model parallelism, pipeline parallelism, and data parallelism
- High-bandwidth interconnects (NVLink, InfiniBand) essential

14.3 In-Context Learning and Few-Shot Prompting

14.3.1 Autoregressive Generation with KV Caching

Before exploring in-context learning, we must understand how GPT generates text autoregressively. The generation process differs fundamentally from training, as tokens are produced sequentially rather than in parallel. Naive implementation of autoregressive generation is extremely inefficient, but key-value caching provides dramatic speedups that make interactive generation practical.

During generation, the model produces one token at a time. Starting with a prompt, the model computes attention over all prompt tokens to generate the first new token. Then it appends this token to the sequence and computes attention over all tokens (prompt plus generated) to produce the second token. This continues until reaching a stopping condition like a maximum length or end-of-sequence token. The critical inefficiency is that each generation step recomputes attention for all previous tokens, even though their key and value representations never change.

Consider generating a sequence of length T tokens. The first step processes n_0 prompt tokens, computing keys and values for all positions. The second step processes $n_0 + 1$ tokens, recomputing the same keys and values for the prompt plus computing them for the new token. By step T , we have computed keys and values for the prompt tokens T times, despite them being identical each time. The total computation grows quadratically: $\sum_{t=1}^T (n_0 + t) = Tn_0 + T(T+1)/2 \approx Tn_0 + T^2/2$ forward passes through the attention mechanism.

Key-value caching eliminates this redundancy by storing the computed keys and values for all previous tokens. When generating token t , we only compute keys and values for the new token at position t , then concatenate with the cached keys and values from positions 1 through $t - 1$. The attention computation at position t uses the full key and value matrices, but we avoid recomputing the cached portions. This reduces the computation from quadratic to linear in the generation length.

The memory requirements for KV caching scale with the sequence length, number of layers, and model dimension. For each layer, we must store key and value matrices of shape $[n_{\text{current}}, d_{\text{model}}]$ where n_{current} is the current sequence length. With L layers and hidden dimension d , the cache requires $2 \times L \times n_{\text{current}} \times d$ values. For GPT-2 with 12 layers, dimension 768, and sequence length 1024, the cache occupies $2 \times 12 \times 1024 \times 768 = 18,874,368$ values, or approximately 75 MB in FP32 per sequence. This is modest compared to model parameters (440 MB for GPT-2), but grows linearly with batch size and sequence length.

The generation speed improvement from KV caching is dramatic. Without caching, generating T tokens requires $O(T^2)$ operations. With caching, it requires $O(T)$ operations. For GPT-2 generating 100 tokens, this represents a $50\times$ speedup in theory. In practice, the speedup is somewhat less due to memory bandwidth limitations and the overhead of managing the cache, but $10\text{-}20\times$ speedups are typical. This transforms generation from painfully slow (1-2 tokens per second) to interactive (20-50 tokens per second) on modern GPUs.

Batch generation introduces additional trade-offs. Processing multiple sequences in parallel amortizes the cost of loading model parameters and improves GPU utilization. However, the KV cache memory scales linearly with batch size. For GPT-2 with batch size 32 and sequence length 1024, the cache requires $32 \times 75 \text{ MB} = 2.4 \text{ GB}$. Combined with model parameters and activations, this can exhaust GPU memory. Practitioners must balance batch size against sequence length and model size to fit within memory constraints. Dynamic batching, where sequences of different lengths are grouped together, can improve efficiency by allowing longer sequences when the batch is small and more sequences when they are short.

Generation algorithm with KV caching:

1. Process prompt tokens $[x_1, \dots, x_{n_0}]$ in parallel, computing and caching keys/values for all layers
2. For generation step $t = 1, 2, \dots, T$:
 - (a) Compute keys/values only for new token at position $n_0 + t$
 - (b) Concatenate with cached keys/values from positions 1 to $n_0 + t - 1$
 - (c) Compute attention using full key/value matrices

- (d) Generate next token from output distribution
 - (e) Append new keys/values to cache
3. Return generated sequence $[x_{n_0+1}, \dots, x_{n_0+T}]$

Memory requirements for KV cache:

$$\text{Cache memory} = 2 \times L \times n_{\max} \times d_{\text{model}} \times B \times \text{bytes per value} \quad (14.11)$$

For GPT-2 (12 layers, 768 dim, 1024 tokens, batch 1, FP32):

$$2 \times 12 \times 1024 \times 768 \times 1 \times 4 = 75,497,472 \text{ bytes} \approx 75 \text{ MB} \quad (14.12)$$

For GPT-3 (96 layers, 12288 dim, 2048 tokens, batch 1, FP16):

$$2 \times 96 \times 2048 \times 12288 \times 1 \times 2 = 9,663,676,416 \text{ bytes} \approx 9.7 \text{ GB} \quad (14.13)$$

Generation speed comparison:

- **Without caching:** ~ 1 -2 tokens/sec (recomputes all previous tokens)
- **With caching:** ~ 20 -50 tokens/sec for GPT-2 on V100
- **With caching:** ~ 10 -15 tokens/sec for GPT-3 on A100 (batch 1)
- **Batch generation:** Higher throughput (tokens/sec) but same latency per sequence

14.3.2 Zero-Shot, One-Shot, Few-Shot

14.3.3 Zero-Shot, One-Shot, Few-Shot

GPT-3's most remarkable capability is in-context learning: the ability to perform new tasks by conditioning on examples provided in the prompt, without any parameter updates or gradient descent. This emergent behavior was not explicitly trained for, yet it scales dramatically with model size, suggesting that large language models develop meta-learning capabilities through pre-training alone.

Zero-shot learning provides only a task description without examples. The model must infer the desired behavior from the natural language instruction alone. For translation, a zero-shot prompt might simply state "Translate English to French:" followed by the source text. The model must recognize the task from the instruction and generate an appropriate translation. Zero-shot performance varies widely across tasks—GPT-3 performs well on common tasks like translation and summarization but struggles with specialized or ambiguous tasks where the instruction alone provides insufficient specification.

Zero-shot: Task description only

```
Translate English to French:
sea otter =>
```

One-shot learning adds a single example demonstrating the desired input-output mapping. This single example often dramatically improves performance by clarifying the task format, output style, and level of detail expected. For translation, showing one English-French pair helps the model understand not just that translation is required, but also the desired formality level, whether to include punctuation, and how to handle proper nouns. The improvement from zero-shot to one-shot is often larger than from one-shot to few-shot, suggesting that the first example resolves most of the task ambiguity.

One-shot: One example

```
Translate English to French:
sea otter => loutre de mer
cheese =>
```

Few-shot learning provides multiple examples, typically between 10 and 100 depending on the task complexity and context window size. Additional examples help the model learn task-specific patterns, edge cases, and output formatting. For classification tasks, few-shot examples should cover all classes to avoid bias toward classes seen more frequently. For generation tasks, examples demonstrate the desired output length, style, and structure. The performance improvement from few-shot learning scales with both the number of examples and the model size—larger models extract more information from the same examples.

Few-shot: Multiple examples (typical: 10-100)

```
Translate English to French:
sea otter => loutre de mer
peppermint => menthe poivrée
plush giraffe => girafe en peluche
cheese =>
```

The mechanism underlying in-context learning remains partially mysterious. The model is not performing gradient descent or updating parameters—it processes the prompt in a single forward pass. Instead, the model appears to perform a form of implicit Bayesian inference, using the examples to narrow down the space of possible tasks and then applying the inferred task to the query. The attention mechanism plays a crucial role, allowing later tokens to attend to earlier examples and extract relevant patterns. Larger models have more capacity to represent complex task distributions and perform more sophisticated inference, explaining why few-shot learning improves dramatically with scale.

Key Point 14.1. *GPT-3's key discovery: Large language models can perform tasks through in-context learning without parameter updates! Performance improves with model scale and number of examples.*

The practical implications are profound. In-context learning enables rapid adaptation to new tasks without fine-tuning, which requires labeled data, computational resources, and time. Users can deploy GPT-3 on novel tasks by simply crafting appropriate prompts with examples. This has spawned the field of prompt engineering, where practitioners carefully design prompts to elicit desired behaviors. However, in-context learning has limitations—it cannot match fine-tuned performance on tasks with abundant training data, and it is sensitive to example selection and ordering. The examples must fit within the context window, limiting the amount of task-specific information that can be provided.

14.3.4 Emergent Abilities

As language models scale to billions and hundreds of billions of parameters, they exhibit emergent abilities—capabilities that appear suddenly at certain scale thresholds rather than improving gradually. These emergent behaviors were not explicitly programmed or trained for, yet they arise naturally from the combination of scale, architecture, and training data. Understanding emergence is crucial for predicting what capabilities future models might develop and for identifying the minimum scale required for specific applications.

Few-shot learning itself is an emergent ability. Models with fewer than 1 billion parameters show minimal few-shot learning capability—providing examples in the prompt barely improves performance over zero-shot. Between 1 billion and 10 billion parameters, few-shot learning begins to emerge, with clear improvements from adding examples. By 100 billion parameters, few-shot learning becomes highly effective, with GPT-3 175B demonstrating strong performance on many tasks with just 10-20 examples. This non-linear scaling suggests a phase transition in the model's internal representations, where sufficient capacity enables a qualitatively different form of processing.

Chain-of-thought reasoning represents another striking emergent ability. When prompted to show its reasoning step-by-step before providing an answer, models around 100 billion parameters begin to solve complex multi-step problems that smaller models cannot. For arithmetic word problems, asking

the model to "think step by step" dramatically improves accuracy. The model generates intermediate reasoning steps, then uses those steps to arrive at the final answer. This capability appears suddenly—models below a certain scale show no benefit from chain-of-thought prompting, while larger models show substantial improvements. The emergence suggests that large models develop internal mechanisms for decomposing complex problems into simpler sub-problems.

Complex instruction following emerges only in the largest models. GPT-3 175B can follow multi-part instructions, maintain consistency across long generations, and adapt its behavior based on nuanced prompt details. Smaller models often ignore parts of complex instructions or fail to maintain consistency. This capability is essential for practical applications where users need fine-grained control over model behavior. The emergence of instruction following motivated the development of instruction-tuned models like InstructGPT, which further enhance this capability through supervised fine-tuning and reinforcement learning.

The scaling curve for most capabilities follows a smooth power law—performance improves predictably as model size increases. However, emergent abilities show sharp phase transitions where performance jumps discontinuously at certain scales. This creates challenges for predicting model capabilities: extrapolating from smaller models may underestimate the capabilities of larger models. It also raises questions about what other abilities might emerge at even larger scales. Some researchers hypothesize that abilities like true reasoning, planning, and causal understanding might emerge at scales beyond current models, while others argue that architectural changes or different training objectives are necessary.

The mechanism underlying emergence remains debated. One hypothesis is that emergent abilities require a minimum representational capacity—below this threshold, the model cannot represent the necessary abstractions, while above it, the ability appears. Another hypothesis focuses on the training dynamics: certain capabilities require seeing specific patterns in the training data a minimum number of times, which only occurs when training on massive datasets. A third perspective suggests that emergence is partially an artifact of evaluation metrics—capabilities may improve gradually, but threshold-based metrics (like exact match accuracy) show discontinuous jumps.

Abilities that appear suddenly at certain scales:

- **Few-shot learning:** Emerges around 1B-10B parameters
- **Chain-of-thought reasoning:** Emerges around 100B parameters
- **Complex instruction following:** Largest models

Scaling curve: Performance on many tasks follows smooth power law, but some tasks show sharp phase transitions.

14.4 Scaling Laws

14.4.1 Parameter Scaling

The relationship between model size and performance follows remarkably predictable patterns, enabling researchers to forecast the capabilities of larger models before building them. These scaling laws have become central to modern AI research, guiding decisions about how to allocate computational resources between model size, training data, and training time.

The fundamental scaling law relates model performance, measured by loss on held-out data, to the number of parameters. Empirically, the loss follows a power law:

$$L(N) \approx \left(\frac{N_c}{N} \right)^\alpha \quad (14.14)$$

where N is the number of parameters, N_c is a constant, and $\alpha \approx 0.076$. This relationship holds over multiple orders of magnitude, from millions to hundreds of billions of parameters. The power law implies that every 10× increase in parameters yields a consistent reduction in loss, with no sign of saturation up to the largest models tested.

Performance (measured by loss) scales as:

$$L(N) \approx \left(\frac{N_c}{N}\right)^\alpha \quad (14.15)$$

where N is number of parameters, N_c is constant, $\alpha \approx 0.076$.

The practical implications are profound. The power law allows researchers to predict the performance of a 1 trillion parameter model by extrapolating from experiments with 1 billion and 10 billion parameter models. This predictability has motivated continued scaling efforts, as the returns to scale remain consistent even at enormous sizes. However, the exponent $\alpha \approx 0.076$ means that improvements slow as models grow—achieving the same loss reduction requires exponentially more parameters. Reducing loss by half requires increasing parameters by a factor of $(2)^{1/0.076} \approx 150$, making continued progress increasingly expensive.

The scaling law applies specifically to the pre-training loss, which measures how well the model predicts the next token. Downstream task performance does not always scale as smoothly—some tasks show rapid improvement with scale while others plateau. This discrepancy arises because pre-training loss captures general language understanding, while specific tasks may require capabilities that emerge only at certain scales or that are not well-measured by next-token prediction. Nevertheless, pre-training loss remains the most reliable predictor of overall model capability.

Importantly, the scaling law holds only when other factors are not bottlenecks. If the training data is too small, the model will overfit and the scaling law breaks down. If the training time is too short, the model will not converge and performance will be suboptimal. The scaling laws assume that data and compute are scaled appropriately with model size, a condition that is not always met in practice.

Implications:

- Every 10× increase in parameters → consistent loss reduction
- No sign of saturation up to 175B parameters
- Motivates continued scaling

14.4.2 Compute-Optimal Training

While the parameter scaling law shows that larger models achieve better performance, it does not address the question of how to optimally allocate a fixed compute budget. Should we train a very large model on limited data, or a smaller model on more data? The Chinchilla paper provided a surprising answer that has reshaped thinking about model scaling.

The Chinchilla findings, based on training over 400 language models ranging from 70 million to 16 billion parameters, revealed that for a given compute budget C , the optimal allocation scales both model size and training data:

$$N_{\text{optimal}} \propto C^{0.5}, \quad D_{\text{optimal}} \propto C^{0.5} \quad (14.16)$$

This square-root scaling means that if you increase compute by 100×, you should increase both model size and training data by 10×. Critically, this implies that model size and data should scale equally—doubling compute should double both parameters and training tokens.

Chinchilla findings: For compute budget C , optimal allocation is:

$$N_{\text{optimal}} \propto C^{0.5}, \quad D_{\text{optimal}} \propto C^{0.5} \quad (14.17)$$

Applying this formula to GPT-3 reveals a striking conclusion: the model was significantly over-parameterized for its training compute. GPT-3 used 175 billion parameters trained on 300 billion tokens. According to Chinchilla scaling laws, the same compute budget would be better spent on an 80 billion parameter model trained on 1.4 trillion tokens. This smaller, better-trained model would achieve lower loss and better downstream performance than GPT-3, despite having less than half the parameters.

This finding explains why many large models are over-parameterized and under-trained. The focus on parameter count as a headline metric incentivized building the largest possible models, even when

training data was insufficient. The Chinchilla results suggest that future models should prioritize data quality and quantity alongside parameter scaling. This has motivated efforts to curate larger, higher-quality training datasets and to train models for more steps on existing data.

The compute-optimal scaling also has implications for inference costs. Larger models are more expensive to serve, requiring more memory and compute per token generated. If a smaller, better-trained model achieves the same performance, it will be cheaper to deploy. This economic consideration is increasingly important as language models move from research to production applications serving millions of users.

However, the Chinchilla findings come with caveats. The optimal allocation depends on the relative costs of training versus inference. If inference costs dominate (as in production systems serving many users), a larger model trained on less data may be preferable because it achieves better performance per inference FLOP. The optimal allocation also depends on the availability of high-quality training data—if data is limited or expensive to collect, training a larger model on available data may be the only option.

GPT-3 analysis:

- 175B parameters trained on 300B tokens
- Chinchilla suggests: 80B parameters on 1.4T tokens would be better
- Many large models are over-parameterized, under-trained

The future direction suggested by these findings is clear: smaller models trained on more data. This approach reduces training costs (fewer parameters to update), reduces inference costs (smaller models to serve), and improves performance (better training efficiency). The challenge lies in collecting and curating the massive datasets required—1.4 trillion tokens is nearly $5\times$ the data used for GPT-3, requiring extensive web scraping, filtering, and deduplication. Nevertheless, the Chinchilla findings have fundamentally shifted the scaling paradigm from “bigger is better” to “balanced scaling is optimal.”

14.4.3 Hardware Requirements for Inference

While training requirements determine whether a model can be built, inference requirements determine whether it can be deployed. Understanding the hardware needed to serve GPT models is essential for practitioners considering which models to use in production and for researchers designing new architectures.

GPT-2 with 1.5 billion parameters represents the upper end of models that can be served efficiently on consumer hardware. In FP16 precision, the model parameters occupy $1.5 \times 10^9 \times 2 = 3$ GB of memory. Adding the KV cache for a sequence of 1024 tokens requires approximately 75 MB per sequence, and activations for a single forward pass add another 100-200 MB. A single NVIDIA V100 GPU with 16 GB of memory can comfortably serve GPT-2 with batch sizes of 4-8 sequences, achieving generation speeds of approximately 50 tokens per second per sequence. This makes GPT-2 practical for real-time applications like chatbots, code completion, and interactive writing assistants.

The generation speed of 50 tokens per second on a V100 reflects several factors. The V100 provides 125 TFLOPS of FP16 performance, but actual utilization is typically 30-50% for autoregressive generation due to the sequential nature of the computation and memory bandwidth limitations. Each token generation requires a forward pass through all 48 layers, computing attention over the growing sequence length. With KV caching, the computation per token is roughly constant, but memory bandwidth for loading the cache and model parameters becomes the bottleneck. Batch processing multiple sequences in parallel improves throughput by amortizing parameter loading, but latency per sequence remains constant.

GPT-2 (1.5B) Inference:

- **Memory (FP16):** 3 GB parameters + 75 MB KV cache per sequence + 200 MB activations
- **Hardware:** Single V100 (16 GB) or RTX 3090 (24 GB)
- **Batch size:** 4-8 sequences on V100

- **Generation speed:** ~ 50 tokens/sec per sequence
- **Latency:** ~ 20 ms per token
- **Practical for:** Real-time applications, edge deployment

GPT-3 with 175 billion parameters presents dramatically different challenges. In FP16 precision, the parameters alone require $175 \times 10^9 \times 2 = 350$ GB of memory. No single GPU can hold the entire model—even the largest NVIDIA A100 with 80 GB falls far short. Model parallelism is essential, splitting the model across multiple GPUs. A minimum of $8 \times$ A100 (80 GB) GPUs is required just to hold the parameters, with each GPU storing approximately 44 GB of model weights. The KV cache for GPT-3 with 2048 tokens requires approximately 9.7 GB per sequence, further constraining batch sizes. With 8 GPUs, the total available memory is 640 GB, leaving roughly 290 GB for KV cache and activations after storing parameters—enough for batch sizes of 20-30 sequences.

The generation speed for GPT-3 is significantly slower than GPT-2, despite using more powerful hardware. With batch size 1 on $8 \times$ A100 GPUs, GPT-3 generates approximately 10 tokens per second. The slowdown reflects several factors. First, the model is $100\times$ larger, requiring $100\times$ more computation per token. Second, model parallelism introduces communication overhead—activations must be transferred between GPUs at each layer, consuming bandwidth and adding latency. Third, the larger KV cache requires more memory bandwidth to load at each generation step. Increasing batch size improves throughput (total tokens per second across all sequences) but does not reduce latency per sequence.

GPT-3 (175B) Inference:

- **Memory (FP16):** 350 GB parameters + 9.7 GB KV cache per sequence
- **Hardware:** Minimum $8 \times$ A100 (80 GB), often $16 \times$ for production
- **Model parallelism:** Required—split across GPUs
- **Batch size:** 1-4 sequences per 8-GPU node (memory constrained)
- **Generation speed:** ~ 10 tokens/sec per sequence (batch 1)
- **Latency:** ~ 100 ms per token
- **Cost:** \$0.02-0.06 per 1000 tokens (cloud pricing)

The high cost of GPT-3 inference has motivated extensive optimization efforts. Quantization to INT8 or INT4 reduces memory requirements by $2\text{--}4\times$, enabling larger batch sizes or smaller hardware configurations. However, quantization requires careful calibration to avoid accuracy degradation, and not all operations benefit equally—attention computations are particularly sensitive to reduced precision. Distillation, where a smaller model is trained to mimic GPT-3’s outputs, can achieve 90-95% of the performance with $10\times$ fewer parameters, dramatically reducing inference costs. Sparse models, where only a subset of parameters are active for each input, offer another path to efficiency.

The economics of serving GPT-3 at scale are daunting. A single $8 \times$ A100 node costs approximately \$30,000-50,000 to purchase or \$20-30 per hour to rent from cloud providers. At 10 tokens per second, a single node can serve roughly 36,000 tokens per hour, or 864,000 tokens per day. For applications serving millions of users, dozens or hundreds of nodes are required, with costs reaching millions of dollars per month. This has created a market for inference-optimized models and specialized hardware, as well as prompting research into more efficient architectures that maintain capability while reducing computational requirements.

Why GPT-3 inference is expensive:

- **Memory:** 350 GB parameters require multiple high-end GPUs
- **Compute:** 175B parameters means $100\times$ more FLOPs than GPT-2
- **Communication:** Model parallelism requires high-bandwidth interconnects
- **Latency:** Sequential generation cannot be parallelized across tokens
- **Utilization:** Autoregressive generation achieves 20-40% of peak FLOPS

14.5 Instruction Tuning and RLHF

14.5.1 Instruction Tuning

Fine-tune on (instruction, output) pairs:

Instruction: Summarize the following in one sentence:

[long text]

Output: [one-sentence summary]

InstructGPT / ChatGPT approach:

1. Pre-train with language modeling
2. Supervised fine-tuning on high-quality instructions
3. Train reward model from human preferences
4. Optimize policy with reinforcement learning

14.5.2 RLHF (Reinforcement Learning from Human Feedback)

[H] RLHF Training

Step 1: Supervised Fine-Tuning

- Collect demonstrations: (prompt, high-quality response)
- Fine-tune GPT on demonstrations

Step 2: Reward Model Training

- Generate multiple responses per prompt
- Humans rank responses
- Train reward model $r(x, y)$ to predict rankings

Step 3: RL Fine-Tuning

- Optimize policy π_θ using PPO
- Objective: $\mathbb{E}_{x, y \sim \pi_\theta} [r(x, y)] - \beta \text{KL}(\pi_\theta \| \pi_{\text{ref}})$
- KL penalty prevents divergence from original model

Result: Models better aligned with human preferences, more helpful, honest, and harmless.

14.6 GPT Capabilities and Limitations

14.6.1 Capabilities

Strong:

- Text generation (creative writing, code, dialogue)
- Translation and summarization
- Question answering
- Few-shot learning
- Chain-of-thought reasoning
- Instruction following

14.6.2 Limitations

Weak:

- Factual accuracy (hallucinations)
- Mathematical reasoning (without tools)
- Long-term coherence in very long texts
- True understanding vs pattern matching
- Consistent personality/beliefs

Hallucinations: Model generates plausible but false information with high confidence.

Mitigation strategies:

- Retrieval-augmented generation (RAG)
- Tool use (calculators, search)
- Verification and fact-checking
- Constitutional AI principles

14.7 Exercises

Exercise 14.1. Implement autoregressive language modeling loss. For sequence "The quick brown fox", compute loss with teacher forcing. Compare with exposed schedule where model sees its own predictions.

Exercise 14.2. Estimate training cost for GPT-3 (175B params, 300B tokens):

1. FLOPs per forward pass
2. FLOPs for entire training (forward + backward $\approx 3 \times$ forward)
3. Time on 1024 A100 GPUs (312 TFLOPS each)
4. Cost at \$2/GPU-hour

Exercise 14.3. Implement few-shot prompting. Test GPT-2 on classification task with 0, 1, 5, 10 examples. Plot accuracy vs number of shots. Does performance improve?

Exercise 14.4. Analyze scaling: Train models with [10M, 50M, 100M, 500M] parameters on same data. Plot loss vs parameters on log-log scale. Does it follow power law? Estimate exponent.

Chapter 15

T5 and BART: Encoder-Decoder Architectures

Chapter Overview

T5 (Text-to-Text Transfer Transformer) and BART (Bidirectional and Auto-Regressive Transformers) represent encoder-decoder architectures that combine the strengths of BERT and GPT. This chapter covers their architectures, pre-training objectives, unified text-to-text framework, and applications to sequence-to-sequence tasks.

Learning Objectives

1. Understand encoder-decoder transformer architectures
2. Implement span corruption and denoising objectives
3. Apply text-to-text framework to diverse tasks
4. Compare T5, BART, and other seq2seq transformers
5. Fine-tune for summarization, translation, and question answering
6. Understand prefix LM and mixture of denoisers

15.1 T5: Text-to-Text Transfer Transformer

15.1.1 Unified Text-to-Text Framework

T5 introduces a conceptually elegant framework that reformulates every NLP task as text-to-text transformation. Rather than designing task-specific architectures with classification heads, span prediction layers, or other specialized output structures, T5 treats all tasks uniformly: the model receives text as input and produces text as output. This unification enables a single model architecture and training objective to handle diverse tasks ranging from translation and summarization to classification and question answering.

The text-to-text framework operates by prepending task-specific prefixes to the input text. For translation, the input becomes "translate English to German: That is good", and the model generates "Das ist gut". For summarization, the input is "summarize: [article text]", and the model produces a concise summary. Even classification tasks, which traditionally output discrete labels, are reformulated as text generation: "sst2 sentence: This movie is great" produces the text "positive" rather than a class index. Question answering similarly becomes "question: What is the capital of France? context: Paris is the capital and largest city of France..." with the model generating "Paris" as output.

This unification provides several compelling advantages. First, a single model can handle all tasks without architectural modifications, simplifying deployment and maintenance. Second, the same pre-training objective and fine-tuning procedure apply across tasks, eliminating the need for task-specific

training strategies. Third, the framework enables natural transfer learning across tasks—knowledge learned from translation can potentially benefit summarization, and vice versa. Fourth, evaluation becomes consistent across tasks, as all outputs are text sequences that can be compared using standard metrics. The text-to-text framework represents a philosophical shift toward treating language understanding and generation as a unified capability rather than separate skills requiring different architectures.

Definition 15.1 (Text-to-Text Format). All tasks formulated as: text input \rightarrow text output

- Translation: "translate English to German: That is good" \rightarrow "Das ist gut"
- Summarization: "summarize: [article]" \rightarrow "[summary]"
- Classification: "sst2 sentence: This movie is great" \rightarrow "positive"
- QA: "question: ... context: ..." \rightarrow "[answer]"

15.1.2 T5 Architecture

T5 employs a standard encoder-decoder transformer architecture with several important modifications that distinguish it from the original transformer design. The architecture combines the bidirectional encoding capabilities of BERT with the autoregressive generation capabilities of GPT, creating a model that excels at both understanding input context and generating coherent output sequences.

The encoder processes the input text using fully-visible self-attention, identical to BERT's architecture. Each token in the encoder can attend to all other tokens in the input sequence without any causal masking, enabling the model to build rich bidirectional representations that capture both left and right context. This bidirectional attention is crucial for understanding tasks where the meaning of each token depends on the entire input context. The encoder consists of a stack of transformer layers, each containing multi-head self-attention followed by a position-wise feed-forward network, with residual connections and layer normalization applied in the pre-norm configuration for improved training stability.

The decoder generates the output text autoregressively using causal self-attention, similar to GPT's architecture. Each position in the decoder can only attend to previous positions in the output sequence, ensuring that the model cannot "cheat" by looking at future tokens during generation. Critically, the decoder also includes cross-attention layers that attend to the encoder's output representations. This cross-attention mechanism allows the decoder to focus on relevant parts of the input sequence while generating each output token, enabling the model to perform sequence-to-sequence transformations like translation and summarization where the output depends heavily on specific input content.

T5's most distinctive architectural innovation is its use of relative positional encodings rather than the absolute sinusoidal or learned positional embeddings used in BERT and GPT. Instead of adding position-specific embeddings to the input, T5 computes position-dependent biases that are added to the attention scores. These biases depend only on the relative distance between query and key positions, not their absolute positions in the sequence. The relative position biases are learned during training and shared across all layers, reducing the number of parameters while providing the model with flexible position information. The biases use a bucketing scheme where nearby positions have unique biases but distant positions share biases, reflecting the intuition that precise relative position matters more for nearby tokens than distant ones.

Definition 15.2 (T5 Architecture). T5 uses encoder-decoder transformer with:

- **Encoder:** Fully-visible self-attention (like BERT), no causal masking
- **Decoder:** Causal self-attention (like GPT) plus cross-attention to encoder

- **Positional encoding:** Relative position bias, shared across layers, learned bucket-based distances
- **Normalization:** Pre-norm (layer norm before sub-layers)

Example 15.1 (T5-Base Architecture). Understanding T5-base’s parameter distribution reveals how encoder-decoder architectures allocate capacity between understanding and generation. The model uses 12 encoder layers and 12 decoder layers, each with hidden dimension $d = 768$, 12 attention heads, and feed-forward dimension $d_{ff} = 3072$. The vocabulary contains 32,000 tokens using SentencePiece tokenization, which provides better multilingual coverage and handles rare words more gracefully than WordPiece.

The parameter breakdown shows that the decoder contains more parameters than the encoder despite having the same number of layers and hidden dimensions. This asymmetry arises from the cross-attention mechanism in the decoder, which requires additional weight matrices to project encoder outputs into key and value spaces. Each encoder layer contains approximately 7.1 million parameters: 2.36 million in the self-attention mechanism (four projection matrices of dimension 768×768) and 4.72 million in the feed-forward network (two projections: $768 \rightarrow 3072$ and $3072 \rightarrow 768$). Multiplying by 12 layers yields 85.2 million parameters in the encoder stack.

Each decoder layer contains approximately 9.4 million parameters due to the additional cross-attention mechanism. The causal self-attention contributes 2.36 million parameters, identical to the encoder’s self-attention. The cross-attention layer adds another 2.36 million parameters for its query, key, value, and output projections. The feed-forward network contributes 4.72 million parameters, same as the encoder. Multiplying by 12 decoder layers yields 112.8 million parameters in the decoder stack. The token embeddings add 24.6 million parameters ($32,000 \times 768$), bringing the total to approximately 220 million parameters.

The memory requirements for T5-base depend on the numerical precision used. In FP32, the 220 million parameters occupy $220,000,000 \times 4 = 880$ MB. Mixed precision training with FP16 activations and FP32 master weights reduces the working memory to approximately 440 MB for the model parameters during forward and backward passes, though the optimizer maintains FP32 copies. For inference, pure FP16 weights require only 440 MB, enabling T5-base to run comfortably on GPUs with 8-16 GB of memory. The encoder-decoder architecture requires more memory than encoder-only (BERT) or decoder-only (GPT) models of similar capacity, but the additional cross-attention capability justifies this cost for sequence-to-sequence tasks.

Configuration:

- Encoder layers: $L_{\text{enc}} = 12$, Decoder layers: $L_{\text{dec}} = 12$
- Hidden size: $d = 768$, Attention heads: $h = 12$, FFN dimension: $d_{ff} = 3072$
- Vocabulary: $V = 32,000$ (SentencePiece)
- Parameters: $\approx 220\text{M}$

Parameter breakdown:

$$\text{Embeddings: } 32,000 \times 768 = 24.6\text{M} \quad (15.1)$$

$$\text{Encoder (12 layers): } 12 \times 7.1\text{M} = 85.2\text{M} \quad (15.2)$$

$$\text{Decoder (12 layers): } 12 \times 9.4\text{M} = 112.8\text{M} \quad (15.3)$$

$$\text{Total: } \approx 220\text{M} \quad (15.4)$$

Memory requirements:

- FP32: 880 MB (model parameters only)

- FP16: 440 MB (inference)
- Training (mixed precision, batch size 128, sequence length 512): ≈ 12 GB

Decoder has more parameters due to cross-attention layer.

15.1.3 Pre-Training Objective: Span Corruption

T5 introduces span corruption as its primary pre-training objective, a more sophisticated variant of masked language modeling that better aligns with sequence-to-sequence tasks. Rather than masking individual tokens independently as in BERT, span corruption masks contiguous sequences of tokens and trains the model to predict the entire masked span. This objective encourages the model to learn longer-range dependencies and develop stronger generation capabilities, as the decoder must produce multi-token sequences rather than single tokens.

The span corruption procedure begins by sampling span lengths from a Poisson distribution with parameter $\lambda = 3$, yielding an average span length of 3 tokens. The algorithm then selects spans to mask such that approximately 15% of tokens in the sequence are corrupted, matching BERT's masking rate for fair comparison. Each masked span is replaced with a unique sentinel token (denoted $\langle X \rangle$, $\langle Y \rangle$, $\langle Z \rangle$, etc.), which serves as a placeholder indicating that tokens have been removed at this position. The model must predict the original content of each masked span in the correct order, identified by the sentinel tokens.

The training format differs significantly from BERT's masked language modeling. The encoder receives the corrupted input sequence with sentinel tokens replacing the masked spans. The decoder must generate a sequence containing the sentinel tokens followed by the original content of each span. For example, if the original text is "Thank you for inviting me to your party last week" and spans at positions 3-4 and 8-9 are masked, the encoder input becomes "Thank you $\langle X \rangle$ inviting me to your $\langle Y \rangle$ week". The decoder target is " $\langle X \rangle$ for $\langle Y \rangle$ party last $\langle Z \rangle$ ", where $\langle Z \rangle$ marks the end of the sequence. This format trains the decoder to produce structured output with clear delimiters, a skill that transfers well to downstream generation tasks.

The computational efficiency of span corruption is notable. By masking spans rather than individual tokens, the number of prediction targets decreases while maintaining the same fraction of corrupted tokens. If 15% of tokens are masked in spans of average length 3, only 5% of positions contain sentinel tokens that trigger predictions. This reduces the decoder's generation length compared to predicting every masked token individually, accelerating training. However, the decoder must still generate all the masked tokens, so the total number of tokens predicted remains approximately 15% of the input length. The efficiency gain comes from the reduced number of sentinel tokens that must be processed by the encoder.

The span corruption objective provides several advantages over BERT's masked language modeling for encoder-decoder models. First, it trains the decoder to generate multi-token sequences, developing the autoregressive generation capabilities needed for downstream tasks like summarization and translation. Second, it encourages the model to learn longer-range dependencies, as predicting a span requires understanding the broader context rather than just neighboring tokens. Third, it creates a more challenging task that prevents the model from relying on simple local patterns, forcing it to develop deeper semantic understanding. Fourth, the sentinel token mechanism provides a natural way to structure the decoder's output, which transfers to tasks requiring structured generation.

Definition 15.3 (Span Corruption). Corrupt spans of consecutive tokens, predict them:

1. Sample span lengths from $\text{Poisson}(\lambda = 3)$, average span length 3 tokens
2. Mask 15% of tokens in spans (same total masking rate as BERT)
3. Replace each span with sentinel token $\langle X \rangle$, $\langle Y \rangle$, etc.

4. Encoder processes corrupted input with sentinels
5. Decoder predicts original spans in order, delimited by sentinels

Example 15.2 (Span Corruption Example). **Original:** "Thank you for inviting me to your party last week"

Step 1: Select spans (15% total): positions [3-4] ("for inviting"), [8-9] ("party last")

Step 2: Replace spans with sentinels

Corrupted input (encoder):

Thank you <X> me to your <Y> week

Target output (decoder):

<X> for inviting <Y> party last <Z>

The encoder processes the corrupted sequence, building bidirectional representations that capture the context around each sentinel token. The decoder must generate the sentinel tokens in order, followed by the original content of each span. The <Z> token marks the end of the sequence, training the model to recognize when generation is complete. This structured prediction task requires the model to maintain coherent state across multiple spans, developing the sequential generation capabilities needed for downstream tasks.

Model must predict masked content and sentinel order, requiring understanding of both local context (what words fit in each span) and global structure (the order of spans in the original sequence).

15.1.4 T5 Model Sizes and Scaling

T5 was released in five different sizes to accommodate different computational budgets and performance requirements. This range of model sizes enables practitioners to choose the appropriate trade-off between accuracy and computational cost for their specific use case. The scaling behavior across these sizes provides valuable insights into how encoder-decoder architectures benefit from increased capacity.

T5-Small contains only 60 million parameters with 6 encoder and 6 decoder layers, hidden dimension $d = 512$, and 8 attention heads. This compact model requires approximately 240 MB in FP32 or 120 MB in FP16, making it suitable for deployment on resource-constrained devices or for applications where inference latency is critical. Despite its small size, T5-Small achieves reasonable performance on many tasks, demonstrating that the text-to-text framework and pre-training objective provide strong inductive biases even with limited capacity. Training T5-Small requires approximately 2-3 days on 8 GPUs, making it accessible for academic research and smaller organizations.

T5-Base, with 220 million parameters as detailed previously, represents the standard configuration that balances performance and computational cost. This size is comparable to BERT-base and GPT-2 Small, enabling direct comparisons across architectural paradigms. T5-Base training requires approximately 1 week on 64 TPU cores or equivalent GPU clusters, with an estimated cost of \$10,000-\$15,000 using cloud computing resources. The model achieves strong performance across diverse tasks, often matching or exceeding BERT-large despite having fewer parameters, demonstrating the effectiveness of the encoder-decoder architecture for many applications.

T5-Large scales to 770 million parameters with 24 encoder and 24 decoder layers, hidden dimension $d = 1024$, and 16 attention heads. The parameter count increases by $3.5\times$ compared to T5-Base, requiring approximately 3 GB in FP32 or 1.5 GB in FP16. Training T5-Large demands approximately 2-3 weeks on 256 TPU cores, with estimated costs of \$50,000-\$75,000. The performance improvements over T5-Base are substantial, particularly on challenging tasks requiring deep reasoning or long-range dependencies. However, the inference latency also increases proportionally, making T5-Large more suitable for offline processing or applications where accuracy is paramount.

T5-3B pushes to 3 billion parameters with 24 encoder and 24 decoder layers, hidden dimension $d = 1024$, and 32 attention heads. The increased head count (compared to T5-Large’s 16 heads) allows for more diverse attention patterns without increasing the per-head dimension. The model requires approximately 12 GB in FP32 or 6 GB in FP16, necessitating high-memory GPUs like the A100 (40-80 GB) for training. Training T5-3B takes approximately 1 month on 512 TPU cores, with estimated costs exceeding \$200,000. The performance gains over T5-Large are more modest, suggesting diminishing returns as model size increases, though T5-3B still achieves state-of-the-art results on several benchmarks.

T5-11B represents the largest variant with 11 billion parameters, using 24 encoder and 24 decoder layers, hidden dimension $d = 1024$, and 128 attention heads. The massive increase in attention heads (from 32 to 128) enables extremely fine-grained attention patterns, though each head operates on a smaller dimension ($d_k = 1024/128 = 8$). The model requires approximately 44 GB in FP32 or 22 GB in FP16, necessitating model parallelism across multiple GPUs even for inference. Training T5-11B demands approximately 2-3 months on 1024 TPU cores, with estimated costs exceeding \$1 million. The model achieves the best performance across virtually all tasks, setting new state-of-the-art results on GLUE, SuperGLUE, and SQuAD at the time of release. However, the computational requirements limit its practical deployment to scenarios where maximum accuracy justifies the cost.

T5 Model Sizes:

- **T5-Small:** 60M parameters, 6 enc + 6 dec layers, $d = 512$, 8 heads
 - Memory: 240 MB (FP32), 120 MB (FP16)
 - Training: 2-3 days on 8 GPUs, cost \approx \$2,000
- **T5-Base:** 220M parameters, 12 enc + 12 dec layers, $d = 768$, 12 heads
 - Memory: 880 MB (FP32), 440 MB (FP16)
 - Training: 1 week on 64 TPU cores, cost \approx \$10,000-\$15,000
- **T5-Large:** 770M parameters, 24 enc + 24 dec layers, $d = 1024$, 16 heads
 - Memory: 3 GB (FP32), 1.5 GB (FP16)
 - Training: 2-3 weeks on 256 TPU cores, cost \approx \$50,000-\$75,000
- **T5-3B:** 3 billion parameters, 24 enc + 24 dec layers, $d = 1024$, 32 heads
 - Memory: 12 GB (FP32), 6 GB (FP16)
 - Training: 1 month on 512 TPU cores, cost $>$ \$200,000
- **T5-11B:** 11 billion parameters, 24 enc + 24 dec layers, $d = 1024$, 128 heads
 - Memory: 44 GB (FP32), 22 GB (FP16)
 - Training: 2-3 months on 1024 TPU cores, cost $>$ \$1,000,000

The scaling behavior reveals important insights about encoder-decoder architectures. Performance improves consistently with model size, but the rate of improvement decreases at larger scales. The cost per percentage point of accuracy improvement increases dramatically beyond T5-3B, suggesting that for most practical applications, T5-Base or T5-Large provide the best trade-off between performance and computational cost. The largest models are primarily valuable for research into scaling laws and for applications where even small accuracy improvements justify substantial computational investment.

15.1.5 T5 Training Details

T5’s pre-training represents a massive computational undertaking that required careful optimization of hardware utilization and training procedures. The model was trained on the Colossal Clean Crawled Corpus (C4), a dataset of approximately 750 GB of cleaned English text extracted from Common Crawl. The C4 dataset underwent extensive filtering to remove low-quality content, including deduplication,

language identification to retain only English text, removal of placeholder text and profanity, and filtering of sentences without terminal punctuation. This cleaning process reduced the raw Common Crawl data by approximately 90%, but the resulting corpus provided much higher quality training signal.

The training infrastructure for T5-11B, the largest variant, required 1024 TPU v3 cores running continuously for approximately 2-3 months. Each TPU v3 core provides roughly 123 TFLOPS of bfloat16 performance, yielding a combined peak performance of approximately 126 PFLOPS for the full training cluster. The training used a batch size of 2048 sequences, each of maximum length 512 tokens, for a total of 1,048,576 tokens per batch. This enormous batch size enabled efficient utilization of the TPU hardware and provided stable gradient estimates despite the model’s scale. The learning rate schedule employed a linear warmup over 10,000 steps to a peak learning rate of 10^{-2} , followed by inverse square root decay. The high peak learning rate, much larger than typical for transformer training, was enabled by the large batch size and careful gradient clipping.

The computational cost of T5-11B training is staggering. With 11 billion parameters and processing approximately 1 trillion tokens during training (the C4 dataset seen roughly 1.3 times), the total computation exceeds 10^{24} FLOPs. At an effective compute rate of 50 PFLOPS (assuming 40% utilization of the 126 PFLOPS peak), the training requires approximately $10^{24}/(50 \times 10^{15}) = 20$ million seconds, or roughly 230 days of continuous computation. The reported 2-3 month training time suggests either higher utilization rates or more efficient training procedures than this conservative estimate. The estimated cost exceeds \$1 million using cloud TPU pricing, making T5-11B one of the most expensive models trained at the time of its release in 2019.

T5-Base training is far more accessible, requiring approximately 1 week on 64 TPU v3 cores (128 TPU cores total). The batch size is reduced to 128 sequences of 512 tokens, totaling 65,536 tokens per batch. The training processes approximately 34 billion tokens (the C4 dataset seen once), requiring roughly 10^{21} FLOPs total. At an effective compute rate of 2 PFLOPS, the training takes approximately 5-7 days, matching the reported training time. The estimated cost is \$10,000-\$15,000, making T5-Base training feasible for well-funded academic labs and smaller companies. The more modest computational requirements have enabled widespread experimentation with the T5 architecture and training approach.

The memory requirements during training are substantial due to the encoder-decoder architecture. For T5-11B with batch size 2048 and sequence length 512, the activations alone consume approximately 200-300 GB of memory. The model parameters require 44 GB in FP32, and the optimizer states (Adam maintains first and second moment estimates) require an additional 88 GB. The total memory footprint exceeds 400 GB, necessitating model parallelism across multiple TPU cores. The training employed a combination of data parallelism (different sequences on different cores) and model parallelism (different layers on different cores) to distribute the memory and computation efficiently. The cross-attention mechanism in the decoder requires storing encoder outputs for all sequences in the batch, adding significant memory overhead compared to encoder-only or decoder-only architectures.

T5-11B Training Configuration:

- Hardware: 1024 TPU v3 cores (≈ 126 PFLOPS peak)
- Training time: 2-3 months continuous
- Dataset: C4 (750 GB cleaned text, ≈ 1 trillion tokens)
- Batch size: 2048 sequences \times 512 tokens = 1,048,576 tokens/batch
- Learning rate: 10^{-2} peak with inverse square root decay
- Total computation: $> 10^{24}$ FLOPs
- Estimated cost: $> \$1,000,000$
- Memory: > 400 GB (requires model parallelism)

T5-Base Training Configuration:

- Hardware: 64 TPU v3 chips (128 cores, ≈ 15 PFLOPS peak)

- Training time: 5-7 days
- Dataset: C4 (750 GB, single pass \approx 34 billion tokens)
- Batch size: 128 sequences \times 512 tokens = 65,536 tokens/batch
- Learning rate: 10^{-2} peak with inverse square root decay
- Total computation: $\approx 10^{21}$ FLOPs
- Estimated cost: \$10,000-\$15,000
- Memory: \approx 20-30 GB (fits on single GPU with gradient accumulation)

The training procedures incorporated several optimizations to improve efficiency and stability. Mixed precision training with bfloat16 reduced memory consumption and accelerated computation on TPU hardware. Gradient clipping prevented instability from occasional large gradients. Dropout was applied with rate 0.1 during pre-training to prevent overfitting, though later work (T5.1.1) found that removing dropout during pre-training improved performance. The relative position biases were initialized to small random values and learned during training, converging to patterns that emphasized nearby positions while maintaining some attention to distant positions.

15.2 BART: Denoising Autoencoder

15.2.1 BART Architecture and Design Philosophy

BART (Bidirectional and Auto-Regressive Transformers) represents Facebook AI Research’s approach to combining the strengths of BERT and GPT through a denoising autoencoder framework. While T5 focuses on the text-to-text paradigm with task-specific prefixes, BART emphasizes learning robust representations through diverse corruption strategies during pre-training. The model architecture is conceptually similar to T5—an encoder-decoder transformer—but the pre-training approach and design philosophy differ significantly.

The BART encoder employs fully bidirectional attention identical to BERT, allowing each token to attend to all other tokens in the input sequence. This bidirectional processing enables the encoder to build rich contextual representations that capture dependencies in both directions. The encoder processes corrupted input text, where corruption can take many forms including token masking, deletion, infilling, sentence permutation, or document rotation. The diversity of corruption strategies forces the encoder to learn robust representations that can handle various types of noise and structural perturbations.

The BART decoder uses causal self-attention like GPT, generating output tokens autoregressively from left to right. Each position in the decoder can only attend to previous positions in the output sequence, maintaining the autoregressive property essential for text generation. The decoder also includes cross-attention layers that attend to the encoder’s output representations, enabling it to focus on relevant parts of the corrupted input while reconstructing the original text. This cross-attention mechanism is crucial for tasks like summarization and translation where the output must be grounded in specific input content.

BART-large, the primary configuration, uses 12 encoder layers and 12 decoder layers with hidden dimension $d = 1024$ and 16 attention heads. This configuration is comparable to BERT-large in terms of depth and width, but the encoder-decoder architecture results in more total parameters. The model uses learned absolute positional embeddings rather than T5’s relative position biases or the original transformer’s sinusoidal encodings. The vocabulary contains approximately 50,000 tokens using byte-pair encoding (BPE), providing finer-grained tokenization than T5’s 32,000-token SentencePiece vocabulary.

Definition 15.4 (BART). Bidirectional And Auto-Regressive Transformers:

- Encoder: Bidirectional self-attention (like BERT), processes corrupted input
- Decoder: Autoregressive causal attention (like GPT) plus cross-attention to encoder
- Pre-training: Reconstruct original text from diversely corrupted input
- Position encoding: Learned absolute positional embeddings

15.2.2 BART Parameter Breakdown and Memory Requirements

Understanding BART-large’s parameter distribution reveals how the model allocates capacity across its components. Each encoder layer contains approximately 12.6 million parameters. The self-attention mechanism requires four projection matrices (\mathbf{W}^Q , \mathbf{W}^K , \mathbf{W}^V , \mathbf{W}^O), each of dimension 1024×1024 , contributing $4 \times 1024^2 = 4,194,304$ parameters. The feed-forward network uses expansion factor 4, projecting from 1024 to 4096 dimensions and back, contributing $2 \times 1024 \times 4096 = 8,388,608$ parameters. Layer normalization adds minimal parameters. Multiplying by 12 encoder layers yields approximately 151 million parameters in the encoder stack.

Each decoder layer contains approximately 16.8 million parameters due to the additional cross-attention mechanism. The causal self-attention contributes 4.2 million parameters, identical to the encoder’s self-attention. The cross-attention layer adds another 4.2 million parameters for its query, key, value, and output projections. The feed-forward network contributes 8.4 million parameters, same as the encoder. Multiplying by 12 decoder layers yields approximately 202 million parameters in the decoder stack. The token embeddings add $50,000 \times 1024 = 51,200,000$ parameters, and positional embeddings for sequences up to 1024 tokens add another $1024 \times 1024 = 1,048,576$ parameters. The total reaches approximately 406 million parameters.

The memory requirements for BART-large are substantial. In FP32, the 406 million parameters occupy $406,000,000 \times 4 = 1,624$ MB, or approximately 1.6 GB. Mixed precision training with FP16 activations and FP32 master weights reduces the working memory to approximately 812 MB for the model parameters during forward and backward passes. For inference, pure FP16 weights require only 812 MB, enabling BART-large to run on GPUs with 12-16 GB of memory with reasonable batch sizes. Training with batch size 32 and sequence length 512 requires approximately 20-25 GB of GPU memory, necessitating high-memory GPUs like the V100 (32 GB) or A100 (40-80 GB).

BART-large Configuration:

- Encoder: 12 layers, Decoder: 12 layers
- Hidden: $d = 1024$, Heads: $h = 16$, FFN: $d_{ff} = 4096$
- Vocabulary: $V \approx 50,000$ (BPE)
- Parameters: $\approx 406\text{M}$

Parameter breakdown:

$$\begin{aligned} \text{Embeddings: } & 50,000 \times 1024 + 1024 \times 1024 = 52.2\text{M} & (15.5) \\ \text{Encoder (12 layers): } & 12 \times 12.6\text{M} = 151.2\text{M} & (15.6) \\ \text{Decoder (12 layers): } & 12 \times 16.8\text{M} = 201.6\text{M} & (15.7) \\ \text{Total: } & \approx 406\text{M} & (15.8) \end{aligned}$$

Memory requirements:

- FP32: 1.6 GB (model parameters only)
- FP16: 812 MB (inference)
- Training (mixed precision, batch size 32, sequence length 512): $\approx 20\text{-}25$ GB

15.2.3 Denoising Objectives and Corruption Strategies

BART’s key innovation lies in exploring multiple corruption strategies during pre-training, systematically evaluating which types of noise lead to the most robust and transferable representations. Unlike BERT’s single masking strategy or T5’s span corruption, BART experiments with five different corruption approaches and combinations thereof. This exploration revealed that the choice of corruption strategy significantly impacts downstream task performance, with different strategies providing complementary benefits.

Token masking, borrowed directly from BERT, replaces random tokens with a special [MASK] token. Approximately 15% of tokens are selected and replaced, forcing the model to predict the original tokens based on surrounding context. This strategy is familiar and well-understood, providing a baseline for comparison with other corruption approaches. However, token masking has limitations: the [MASK] token never appears during fine-tuning, creating a train-test mismatch, and the independent masking of tokens doesn’t encourage the model to learn longer-range dependencies or sequential generation capabilities.

Token deletion removes random tokens entirely from the input sequence, forcing the model to determine which positions are missing and what content should fill them. Unlike masking, which provides explicit markers indicating where tokens were removed, deletion requires the model to infer the locations of missing content from the remaining context. This creates a more challenging task that encourages the model to develop robust positional understanding and the ability to detect gaps in the input. For example, deleting "B" and "D" from "A B C D E" yields "A C E", and the model must reconstruct the full sequence "A B C D E" without explicit indicators of where tokens were removed.

Text infilling represents a more sophisticated corruption strategy that combines aspects of span masking and deletion. Spans of text are sampled (with lengths drawn from a Poisson distribution with $\lambda = 3$, similar to T5), but instead of replacing each span with a unique sentinel token, all spans are replaced with a single [MASK] token. This forces the decoder to determine how many tokens to generate for each masked span based on context alone. For example, replacing "B C D E" in "A B C D E F" with a single [MASK] yields "A [MASK] F", and the model must generate "B C D E" without knowing in advance that four tokens are needed. This uncertainty makes text infilling substantially more challenging than T5’s span corruption with explicit sentinel tokens.

Sentence permutation shuffles the order of sentences within a document, requiring the model to reconstruct the original sentence order. This corruption strategy targets document-level structure rather than token-level content, encouraging the model to learn discourse coherence and inter-sentence dependencies. For example, a document with sentences [S1, S2, S3, S4] might be permuted to [S3, S1, S4, S2], and the model must generate the original order [S1, S2, S3, S4]. This task is particularly relevant for summarization and document understanding, where maintaining coherent structure is crucial.

Document rotation selects a random token as the new start of the document and rotates the entire sequence accordingly. The model must identify the true start of the document and generate the original sequence. For example, rotating "A B C D E" at position 3 yields "D E A B C", and the model must recognize that "A" is the true start and generate "A B C D E". This task encourages the model to learn document-level structure and identify natural boundaries, though it proved less effective than other corruption strategies in practice.

The BART paper systematically evaluated these corruption strategies individually and in combination, finding that text infilling combined with sentence permutation provided the best performance across downstream tasks. This combination balances token-level and document-level corruption, encouraging the model to learn both local language patterns and global document structure. The text infilling component develops strong generation capabilities by forcing the model to produce variable-length spans, while sentence permutation develops discourse understanding by requiring the model to reason about inter-sentence relationships.

BART Corruption Strategies:

1. Token Masking: Replace tokens with [MASK] (like BERT)

- 15% of tokens replaced with [MASK]
- Provides explicit markers for missing content

- Baseline strategy for comparison

2. Token Deletion: Remove random tokens entirely

Original: A B C D E

Corrupted: A C E

Target: A B C D E

- Model must infer locations of missing tokens
- More challenging than masking
- Encourages robust positional understanding

3. Text Infilling: Replace spans with single [MASK]

Original: A B C D E F

Corrupted: A [MASK] F

Target: B C D E

- Span lengths sampled from $\text{Poisson}(\lambda = 3)$
- Model must determine span length from context
- More challenging than T5's sentinel-based span corruption

4. Sentence Permutation: Shuffle sentence order

- Targets document-level structure
- Encourages learning of discourse coherence
- Particularly beneficial for summarization

5. Document Rotation: Rotate document, model finds start

- Less effective than other strategies
- Encourages learning of document boundaries

Best combination (BART's final): Text infilling + sentence permutation

- Balances token-level and document-level corruption
- Develops both generation and discourse understanding
- Achieves best performance across diverse downstream tasks

Example 15.3 (BART Pre-training). Original document:

The cat sat on the mat. It was very comfortable.

The dog barked loudly.

After corruption (infilling + permutation):

The dog barked loudly.

The [MASK] comfortable.

Encoder input: Corrupted text

Decoder target: Original complete text

The model must reconstruct the missing span "cat sat on the mat. It was very" and reorder the sentences to match the original document structure. This combined corruption strategy forces the

model to develop both local generation capabilities (filling in missing text) and global discourse understanding (recognizing proper sentence order).

15.2.4 BART Training Details

BART-large was trained on a combination of datasets totaling approximately 160 GB of text, including BooksCorpus, English Wikipedia, CC-News, OpenWebText, and Stories. This diverse training corpus provides broad coverage of topics and writing styles, enabling the model to learn robust representations that transfer well to downstream tasks. The training used 256 NVIDIA V100 GPUs for approximately 2 weeks, with an estimated cost of \$50,000-\$75,000 using cloud computing resources.

The training configuration employed a batch size of 128 sequences with maximum length 1024 tokens, totaling 131,072 tokens per batch. This large batch size enabled stable training with the Adam optimizer and efficient GPU utilization. The learning rate schedule used a polynomial decay from a peak learning rate of 3×10^{-4} with 500 warmup steps. The training processed approximately 50 billion tokens total, seeing the training corpus roughly once. Mixed precision training with FP16 reduced memory consumption and accelerated computation on the V100 GPUs.

The memory requirements during training are substantial due to the encoder-decoder architecture and large batch size. With batch size 128 and sequence length 1024, the activations consume approximately 40-50 GB of memory. The model parameters require 1.6 GB in FP32, and the Adam optimizer states require an additional 3.2 GB. The total memory footprint reaches approximately 50-60 GB, necessitating data parallelism across multiple GPUs. Each GPU processes a subset of the batch, with gradients synchronized across GPUs after each backward pass.

BART-large Training Configuration:

- Hardware: 256 NVIDIA V100 GPUs (32 GB each)
- Training time: ≈ 2 weeks
- Dataset: 160 GB text (BooksCorpus, Wikipedia, CC-News, OpenWebText, Stories)
- Batch size: 128 sequences \times 1024 tokens = 131,072 tokens/batch
- Learning rate: 3×10^{-4} peak with polynomial decay
- Total tokens: ≈ 50 billion
- Estimated cost: \$50,000-\$75,000
- Memory per GPU: ≈ 25 -30 GB (data parallelism across GPUs)

15.3 Encoder-Decoder Efficiency Analysis

15.3.1 Computational Cost of Cross-Attention

Understanding the computational and memory costs of encoder-decoder architectures compared to encoder-only (BERT) or decoder-only (GPT) models is essential for choosing the appropriate architecture for a given task. The key difference lies in the cross-attention mechanism, which enables the decoder to attend to encoder outputs but introduces additional computational and memory overhead.

The cross-attention mechanism in each decoder layer requires computing attention between decoder queries and encoder keys/values. For a decoder sequence of length n_{dec} and encoder sequence of length n_{enc} , the cross-attention computation involves three main steps. First, the decoder hidden states are projected to queries $\mathbf{Q} \in \mathbb{R}^{n_{\text{dec}} \times d}$ using weight matrix $\mathbf{W}^Q \in \mathbb{R}^{d \times d}$, requiring $n_{\text{dec}} \times d^2$ FLOPs. Second, the encoder outputs are projected to keys $\mathbf{K} \in \mathbb{R}^{n_{\text{enc}} \times d}$ and values $\mathbf{V} \in \mathbb{R}^{n_{\text{enc}} \times d}$ using weight matrices $\mathbf{W}^K, \mathbf{W}^V \in \mathbb{R}^{d \times d}$, requiring $2 \times n_{\text{enc}} \times d^2$ FLOPs. Third, the attention scores $\mathbf{S} = \mathbf{QK}^\top$ are computed, requiring $n_{\text{dec}} \times n_{\text{enc}} \times d$ FLOPs, followed by softmax and multiplication with values, requiring another $n_{\text{dec}} \times n_{\text{enc}} \times d$ FLOPs.

The total computational cost of cross-attention per layer is approximately $n_{\text{dec}} \times d^2 + 2 \times n_{\text{enc}} \times d^2 + 2 \times n_{\text{dec}} \times n_{\text{enc}} \times d$ FLOPs. For typical sequence lengths where $n_{\text{dec}} \approx n_{\text{enc}} = n$, this simplifies to $3nd^2 + 2n^2d$ FLOPs. Comparing to self-attention, which requires $4nd^2 + 2n^2d$ FLOPs, cross-attention adds approximately 75% of the cost of self-attention per layer. With L_{dec} decoder layers, the total cross-attention cost is $L_{\text{dec}} \times (3nd^2 + 2n^2d)$ FLOPs.

For T5-Base with 12 decoder layers, $d = 768$, and $n = 512$, the cross-attention computation requires approximately $12 \times (3 \times 512 \times 768^2 + 2 \times 512^2 \times 768) \approx 12 \times (9.1 + 4.0) \times 10^8 = 1.57 \times 10^{11}$ FLOPs per forward pass. This represents approximately 15-20% of the total forward pass computation, a significant but not dominant fraction. The cross-attention cost scales linearly with the number of decoder layers and quadratically with sequence length, making it increasingly expensive for long sequences.

The memory requirements for cross-attention are equally important. The encoder outputs must be stored in memory for all decoder layers to access during cross-attention. For batch size B , encoder sequence length n_{enc} , and hidden dimension d , the encoder outputs require $B \times n_{\text{enc}} \times d$ values. For T5-Base with batch size 32, sequence length 512, and dimension 768, this amounts to $32 \times 512 \times 768 = 12,582,912$ values, or approximately 50 MB in FP32 or 25 MB in FP16. While modest compared to model parameters, this memory scales linearly with batch size and sequence length, becoming significant for large batches or long sequences.

Additionally, the cross-attention mechanism requires storing attention weights $\mathbf{A} \in \mathbb{R}^{n_{\text{dec}} \times n_{\text{enc}}}$ for each head in each layer during training (for backpropagation). With h attention heads and L_{dec} decoder layers, the total attention weight memory is $B \times L_{\text{dec}} \times h \times n_{\text{dec}} \times n_{\text{enc}}$ values. For T5-Base with batch size 32, 12 decoder layers, 12 heads, and sequence length 512, this amounts to $32 \times 12 \times 12 \times 512 \times 512 = 1,207,959,552$ values, or approximately 4.8 GB in FP32 or 2.4 GB in FP16. This memory requirement can become a bottleneck for training with large batch sizes or long sequences.

Cross-attention computational cost per layer:

$$\text{FLOPs}_{\text{cross-attn}} = n_{\text{dec}} \times d^2 + 2 \times n_{\text{enc}} \times d^2 + 2 \times n_{\text{dec}} \times n_{\text{enc}} \times d \quad (15.9)$$

For $n_{\text{dec}} = n_{\text{enc}} = n$:

$$\text{FLOPs}_{\text{cross-attn}} \approx 3nd^2 + 2n^2d \quad (15.10)$$

Memory requirements:

- Encoder outputs: $B \times n_{\text{enc}} \times d$ values (must be stored for all decoder layers)
- Cross-attention weights (training): $B \times L_{\text{dec}} \times h \times n_{\text{dec}} \times n_{\text{enc}}$ values

Example: T5-Base (batch size 32, sequence length 512):

- Cross-attention FLOPs per layer: $\approx 1.3 \times 10^{10}$ FLOPs
- Total cross-attention (12 layers): $\approx 1.6 \times 10^{11}$ FLOPs (15-20% of forward pass)
- Encoder output memory: 50 MB (FP32) or 25 MB (FP16)
- Cross-attention weight memory: 4.8 GB (FP32) or 2.4 GB (FP16)

15.3.2 Comparison: Encoder-Decoder vs Decoder-Only

The choice between encoder-decoder architectures (T5, BART) and decoder-only architectures (GPT) involves fundamental trade-offs in computational efficiency, memory usage, and task suitability. Understanding these trade-offs is essential for practitioners deciding which architecture to use for their specific application.

Decoder-only models like GPT use only causal self-attention, processing sequences autoregressively from left to right. For a sequence of length n , a decoder-only model with L layers requires approximately $L \times (4nd^2 + 2n^2d)$ FLOPs for the forward pass. The memory requirements include model parameters, activations, and KV cache for generation. For GPT-2 with 12 layers, $d = 768$, and $n = 512$, the forward pass requires approximately $12 \times (4 \times 512 \times 768^2 + 2 \times 512^2 \times 768) \approx 1.2 \times 10^{12}$ FLOPs. The KV cache

for generation requires $2 \times L \times n \times d$ values, or approximately 75 MB in FP32 for GPT-2 with sequence length 1024.

Encoder-decoder models like T5 and BART use separate encoder and decoder stacks with cross-attention connecting them. For input sequence length n_{enc} and output sequence length n_{dec} , the encoder requires $L_{\text{enc}} \times (4n_{\text{enc}}d^2 + 2n_{\text{enc}}^2d)$ FLOPs, and the decoder requires $L_{\text{dec}} \times (4n_{\text{dec}}d^2 + 2n_{\text{dec}}^2d + 3n_{\text{dec}}d^2 + 2n_{\text{dec}}n_{\text{enc}}d)$ FLOPs. For T5-Base with $n_{\text{enc}} = n_{\text{dec}} = 512$, the total forward pass requires approximately 2.1×10^{12} FLOPs, roughly $1.75 \times$ more than GPT-2 of similar size. The memory requirements include encoder outputs ($B \times n_{\text{enc}} \times d$) and cross-attention weights, adding 25-50 MB beyond decoder-only models.

The parameter count comparison reveals that encoder-decoder models require more parameters than decoder-only models of similar capacity. T5-Base with 220 million parameters has 12 encoder layers (85M parameters) and 12 decoder layers (113M parameters including cross-attention). GPT-2 with 12 layers and the same hidden dimension contains only 117 million parameters, as it lacks the encoder stack and cross-attention mechanisms. This means encoder-decoder models require approximately $1.9 \times$ more parameters than decoder-only models with the same number of layers and hidden dimension.

However, the computational comparison depends critically on the task. For generation tasks where the input is short and the output is long (e.g., generating a long document from a short prompt), decoder-only models can be more efficient. The encoder-decoder model processes the short input once through the encoder, then generates the long output through the decoder with cross-attention. The decoder-only model must process the entire sequence (input plus generated output) autoregressively, with each new token requiring attention over all previous tokens. For input length n_{in} and output length n_{out} , the decoder-only model requires $\sum_{t=1}^{n_{\text{out}}} (n_{\text{in}} + t) \approx n_{\text{out}} \times n_{\text{in}} + n_{\text{out}}^2/2$ attention operations, while the encoder-decoder model requires n_{in}^2 (encoder) plus n_{out}^2 (decoder self-attention) plus $n_{\text{out}} \times n_{\text{in}}$ (cross-attention). When $n_{\text{out}} \gg n_{\text{in}}$, the encoder-decoder model is more efficient.

For tasks where the input is long and the output is short (e.g., classification or extractive question answering), decoder-only models can be more efficient. The encoder-decoder model must process the long input through the encoder, then generate the short output through the decoder. The decoder-only model processes the input once, then generates the short output. However, encoder-only models like BERT are typically most efficient for these tasks, as they avoid the decoder entirely and use a simple classification head.

The memory efficiency comparison favors decoder-only models for inference, as they avoid storing encoder outputs and cross-attention weights. However, for training with large batch sizes, the difference is less significant, as both architectures require substantial memory for activations and gradients. The KV cache for decoder-only models grows with the total sequence length (input plus output), while encoder-decoder models cache only decoder states, potentially providing memory advantages for long input sequences.

When to use encoder-decoder (T5, BART):

- Sequence-to-sequence tasks: translation, summarization, question answering with generation
- Tasks requiring bidirectional understanding of input: the encoder can attend to the full input context
- Tasks with long input and short output: encoder processes input once, decoder generates short output
- Multi-task learning: text-to-text framework enables unified training across diverse tasks

When to use decoder-only (GPT):

- Pure generation tasks: story generation, dialogue, code generation
- Tasks with short input and long output: decoder-only can be more efficient
- In-context learning: decoder-only models excel at few-shot learning from examples in the prompt
- Simplicity: decoder-only architecture is simpler to implement and deploy

Computational comparison (similar capacity):

- Parameters: Encoder-decoder $\approx 1.9\times$ decoder-only (due to encoder stack and cross-attention)
- FLOPs per forward pass: Encoder-decoder $\approx 1.5\text{-}2\times$ decoder-only (depends on sequence lengths)
- Memory (inference): Decoder-only more efficient (no encoder outputs or cross-attention weights)
- Memory (training): Similar for both architectures with large batch sizes

15.4 Comparing T5 and BART

Aspect	T5	BART
Framework	Text-to-text	Denoising autoencoder
Pre-training	Span corruption	Multiple denoisers
Position encoding	Relative bias	Absolute learned
Vocabulary	32K (SentencePiece)	50K (BPE)
Best for	Unified multi-task	Summarization/generation
Largest size	11B parameters	400M parameters

Performance comparison on GLUE:

- T5-11B: 90.3 (state-of-art at release)
- BART-large: 88.4
- RoBERTa-large: 88.5

Summarization (CNN/DailyMail):

- BART-large: ROUGE-L 44.16 (best)
- T5-base: ROUGE-L 42.05

15.5 Prefix Language Models

15.5.1 Prefix LM Objective

Definition 15.5 (Prefix Language Model). Bidirectional attention on prefix, causal on rest:

- Prefix (input): Fully-visible attention
- Target (output): Causal attention
- Single model (no separate encoder/decoder)

Example:

Prefix: "Translate to French: Hello"

Target: "Bonjour"

Attention mask:

- Prefix tokens can attend to all prefix
- Target tokens attend causally
- Enables both understanding and generation

Models using Prefix LM:

- UniLM (Microsoft)
- GLM (Tsinghua)
- UL2 (Google)

15.6 Applications and Fine-tuning

15.6.1 Summarization

Task: Input document \rightarrow Summary

T5 format:

`summarize: [article text]`

BART approach:

- Encoder: Full article
- Decoder: Generate summary

Metrics:

- ROUGE-1, ROUGE-2, ROUGE-L (n-gram overlap)
- BERTScore (semantic similarity)

15.6.2 Translation

T5 format:

`translate English to German: That is good.`

Output: "Das ist gut."

Multi-task advantage: Single T5 model handles multiple language pairs by conditioning on task prefix.

15.6.3 Question Answering

T5 format:

`question: What is the capital of France?`

`context: Paris is the capital and largest city of France...`

Output: "Paris"

Comparison to BERT:

- BERT: Span prediction (start/end positions)
- T5: Text generation (more flexible)

15.7 Mixture of Denoisers (UL2)

UL2 combines multiple objectives:

R-Denoiser (Regular): Short spans (like T5)

S-Denoiser (Sequential): Prefix LM

X-Denoiser (Extreme): Very long spans or high corruption

Benefits:

- More robust representations
- Better transfer to diverse tasks
- Single model for understanding and generation

15.8 Exercises

Exercise 15.1. Implement span corruption. For text "The quick brown fox jumps over the lazy dog":

1. Sample span lengths from $\text{Poisson}(\lambda = 3)$
2. Corrupt 15% with spans
3. Generate corrupted input and target

Exercise 15.2. Fine-tune T5-base on summarization (CNN/DailyMail):

1. Format data as "summarize: [article]" \rightarrow "[summary]"
2. Train for 3 epochs with learning rate 10^{-4}
3. Evaluate ROUGE scores
4. Compare with BART-base

Exercise 15.3. Calculate parameter counts for:

1. T5-base (encoder + decoder)
2. BART-large
3. Compare to BERT-base (encoder only) and GPT-2 (decoder only)

Explain why encoder-decoder has most parameters.

Exercise 15.4. Implement text-to-text framework. Convert these tasks to T5 format:

1. Sentiment classification (positive/negative)
2. Named entity recognition
3. Textual entailment (premise + hypothesis \rightarrow entailed/contradiction/neutral)

Chapter 16

Efficient Transformers

Chapter Overview

Standard transformers have $O(n^2)$ complexity in sequence length, limiting their application to long sequences. This chapter covers efficient attention mechanisms that reduce complexity: sparse attention, linear attention, low-rank methods, and kernel-based approaches.

Learning Objectives

1. Understand the quadratic bottleneck in standard attention
2. Implement sparse attention patterns (sliding window, strided, global)
3. Apply Linformer and Performer for linear complexity
4. Use Flash Attention for memory-efficient computation
5. Compare trade-offs: accuracy vs efficiency vs memory
6. Deploy long-context models (Longformer, BigBird)

16.1 The Quadratic Bottleneck

16.1.1 Complexity Analysis

The standard self-attention mechanism computes attention scores between all pairs of tokens in a sequence, leading to computational and memory requirements that scale quadratically with sequence length. The attention operation is defined as:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}\right) \mathbf{V} \quad (16.1)$$

where $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{n \times d}$ represent the query, key, and value matrices for a sequence of length n with model dimension d .

The computational bottleneck arises from computing the attention matrix $\mathbf{Q}\mathbf{K}^\top \in \mathbb{R}^{n \times n}$, which requires $O(n^2d)$ floating-point operations. For each of the n queries, we compute dot products with all n keys, where each dot product involves d multiplications and additions. The subsequent softmax normalization adds $O(n^2)$ operations, and the final multiplication with values \mathbf{V} requires another $O(n^2d)$ operations. The memory bottleneck is equally severe: storing the attention matrix requires $O(n^2)$ memory, which must be materialized before the softmax operation and retained for the backward pass during training.

This quadratic scaling becomes prohibitive for long sequences. Consider a BERT-base model with $d = 768$ and 12 attention heads processing a sequence of length $n = 4096$. Each attention head must store an attention matrix of size 4096×4096 , requiring $4096^2 \times 4 = 67$ MB in FP32 format. Across all

12 heads, this amounts to 804 MB just for attention weights in a single layer. With 12 layers, the total memory for attention matrices alone reaches 9.6 GB, nearly filling an NVIDIA V100 GPU with 16 GB memory before accounting for activations, gradients, or model parameters.

Example 16.1 (Long Sequence Costs). The quadratic scaling of attention memory becomes dramatically worse as sequence length increases. For a single attention head with $d = 768$ in FP32 format, the memory requirements grow as follows:

For $n = 512$ tokens (BERT’s original limit), the attention matrix requires $512^2 \times 4 = 1.05$ MB per head. This is manageable even with multiple layers and batch processing. However, increasing to $n = 2048$ tokens requires $2048^2 \times 4 = 16.8$ MB per head—a $16\times$ increase for only a $4\times$ increase in sequence length. At $n = 4096$ tokens, memory consumption reaches 67 MB per head, and at $n = 8192$ tokens, it explodes to 268 MB per head—a $256\times$ increase compared to the 512-token baseline.

With 12 attention heads and 12 layers, processing a single sequence of 8192 tokens requires $268 \times 12 \times 12 = 38.6$ GB just for attention matrices, exceeding the capacity of even high-end GPUs like the A100 with 40 GB memory. This fundamental limitation explains why BERT restricts sequences to 512 tokens, GPT-2 to 1024 tokens, and why efficient attention mechanisms are essential for processing long documents, genomic sequences, or high-resolution images.

The computational cost follows a similar pattern. On an NVIDIA A100 GPU with 312 TFLOPS of FP16 performance, computing attention for $n = 512$ takes approximately 8 milliseconds per layer. For $n = 4096$, this increases to 98 milliseconds—a $12\times$ slowdown for an $8\times$ increase in length. At $n = 16384$, attention computation takes 1.5 seconds per layer, making training completely impractical without efficient attention mechanisms.

16.2 Sparse Attention Patterns

16.2.1 Efficiency Taxonomy

Efficient attention mechanisms can be categorized into five main approaches, each targeting different aspects of the quadratic bottleneck. Sparse attention methods reduce the number of attention connections by restricting each query to attend to only a subset of keys, achieving $O(n \times k)$ complexity where $k \ll n$. Linear attention methods use mathematical approximations to avoid computing the full attention matrix, achieving $O(n)$ complexity in sequence length. Low-rank methods project keys and values to lower-dimensional spaces, reducing the effective size of the attention computation. Kernel-based methods reformulate attention using kernel functions and random features to enable linear-time computation. Finally, recurrent methods process sequences in chunks with recurrent connections, trading parallelism for reduced memory.

Each approach involves different trade-offs between computational efficiency, memory usage, approximation quality, and implementation complexity. Sparse methods maintain exact attention within their connectivity pattern but may miss important long-range dependencies. Linear methods achieve impressive speedups but introduce approximation errors that can degrade model quality. Low-rank methods work well when attention patterns have inherent low-rank structure but may fail for complex attention distributions. Understanding these trade-offs is essential for selecting the appropriate efficient attention mechanism for a given application.

16.2.2 Fixed Sparse Patterns

Sparse attention restricts each query to attend to only a subset of keys, dramatically reducing both computation and memory requirements. The fundamental idea is to identify which attention connections are most important and compute only those, setting all other attention weights to zero or negative infinity before the softmax operation.

Definition 16.1 (Sparse Attention). Sparse attention restricts attention to a predefined subset of positions \mathcal{S} . For each query position i , we compute attention only over keys in the set $\mathcal{S}(i) \subseteq \{1, \dots, n\}$:

$$\text{Attention}_{\text{sparse}}(\mathbf{Q}, \mathbf{K}, \mathbf{V})_{ij} = \begin{cases} \text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V})_{ij} & \text{if } j \in \mathcal{S}(i) \\ 0 & \text{otherwise} \end{cases} \quad (16.2)$$

where $|\mathcal{S}(i)| = k \ll n$ for all positions i . The computational complexity reduces from $O(n^2d)$ to $O(nkd)$, and memory requirements decrease from $O(n^2)$ to $O(nk)$.

The choice of sparsity pattern \mathcal{S} determines which information can flow through the network. Three fundamental patterns have emerged as particularly effective building blocks for sparse attention.

The sliding window or local attention pattern restricts each token to attend only to nearby tokens within a fixed window. Formally, $\mathcal{S}_{\text{local}}(i) = \{j : |i - j| \leq w\}$ where w is the window size. Each token attends to $2w + 1$ tokens: itself, w tokens before, and w tokens after. This pattern is motivated by the observation that in many domains, particularly natural language, nearby tokens are more relevant than distant ones. For a window size $w = 256$ and sequence length $n = 4096$, each query attends to only 513 keys instead of 4096, reducing computation by $8\times$ and memory by the same factor. The limitation is that information can only propagate w positions per layer, requiring $\lceil n/w \rceil$ layers for full sequence communication.

The strided or dilated attention pattern samples tokens at regular intervals: $\mathcal{S}_{\text{strided}}(i) = \{j : (i - j) \bmod s = 0\}$ where s is the stride. This pattern allows each token to attend to distant tokens, enabling faster information propagation across the sequence. With stride $s = 64$, a token at position 1024 can attend to positions 0, 64, 128, ..., 1024, ..., 4032, providing long-range connectivity with only n/s connections per query. However, strided attention alone misses local context, so it is typically combined with local attention in alternating layers.

Global attention designates certain tokens as global tokens that attend to all positions and are attended to by all positions. These tokens act as information hubs, aggregating information from the entire sequence and broadcasting it back. In practice, special tokens like [CLS] in BERT or separator tokens are often designated as global. For g global tokens in a sequence of length n , each global token requires $O(n)$ computation, and each regular token requires $O(g)$ additional computation to attend to globals, adding $O(ng)$ total cost.

Example 16.2 (Longformer Attention Pattern). Longformer combines local and global attention to process documents up to 4096 tokens efficiently. All tokens use local attention with window size $w = 512$, allowing each token to attend to 1024 neighboring tokens (512 on each side). Additionally, task-specific tokens such as [CLS] for classification or question tokens for question answering are designated as global tokens that attend to and are attended by all positions.

For a sequence of length $n = 4096$ with window $w = 512$ and $g = 2$ global tokens, the total number of attention connections is computed as follows. Each of the $n - g = 4094$ regular tokens attends to $2w = 1024$ local tokens plus $g = 2$ global tokens, contributing $(n - g) \times (2w + g) = 4094 \times 1026 \approx 4.2$ million connections. Each of the $g = 2$ global tokens attends to all $n = 4096$ tokens, contributing $g \times n = 8192$ connections. The total is approximately 4.2 million connections compared to $n^2 = 16.8$ million for full attention—a $4\times$ reduction.

The memory savings are equally significant. For a single attention head in FP32, Longformer requires approximately $(4094 \times 1026 + 2 \times 4096) \times 4 = 16.8$ MB compared to 67 MB for full attention. With 12 heads and 12 layers, this reduces total attention memory from 9.6 GB to 2.4 GB, enabling processing of long documents on GPUs with limited memory. On an NVIDIA A100 GPU, Longformer processes 4096-token sequences in approximately 18 milliseconds per layer compared to 98 milliseconds for full attention, a $5.4\times$ speedup.

16.2.3 BigBird: Random + Window + Global

BigBird extends sparse attention by combining three complementary patterns: local windows for nearby context, random connections for long-range dependencies, and global tokens for information aggregation. This combination provides both theoretical guarantees and practical efficiency for processing sequences up to 4096 tokens or longer.

Definition 16.2 (BigBird Attention). BigBird attention combines three sparse patterns for each query position i :

1. **Random attention:** Each query attends to r randomly selected keys, where the random set $\mathcal{R}(i)$ is fixed during initialization and shared across all attention heads.
2. **Window attention:** Each query attends to w neighboring keys on each side, forming a local window $\mathcal{W}(i) = \{j : |i - j| \leq w\}$.
3. **Global attention:** A set of g designated global tokens attend to all positions and are attended by all positions.

The total attention set for a regular token at position i is $\mathcal{S}(i) = \mathcal{W}(i) \cup \mathcal{R}(i) \cup \mathcal{G}$, where \mathcal{G} is the set of global token positions. The total number of connections per query is $|\mathcal{S}(i)| = 2w + r + g$, giving computational complexity $O(n(2w + r + g)d) = O(n)$ when w , r , and g are constants.

The random attention component is crucial for BigBird’s theoretical properties. While local windows provide nearby context and global tokens enable information aggregation, random connections create shortcuts across the sequence that allow information to propagate efficiently. The random graph formed by these connections has high probability of being well-connected, ensuring that any two positions are connected by a short path through the attention graph. This property enables BigBird to approximate full attention’s expressiveness while maintaining linear complexity.

BigBird’s theoretical contribution is proving that this sparse attention pattern can approximate any sequence-to-sequence function that full attention can compute, under mild assumptions. Specifically, BigBird with $r = O(\log n)$ random connections per query can approximate full attention with high probability, providing a theoretical foundation for sparse attention methods. This result shows that $O(n \log n)$ total connections suffice for universal approximation, compared to $O(n^2)$ for full attention.

In practice, BigBird uses $w = 256$, $r = 64$, and $g = 32$ for sequences up to 4096 tokens. Each regular token attends to $2 \times 256 + 64 + 32 = 608$ keys instead of 4096, reducing computation by $6.7\times$. For a single attention head with $d = 768$ in FP32, BigBird requires $(4096 - 32) \times 608 + 32 \times 4096 \times 4 \approx 10.4$ MB compared to 67 MB for full attention, a $6.4\times$ memory reduction. With 12 heads and 12 layers, total attention memory decreases from 9.6 GB to 1.5 GB.

The performance benefits are substantial on modern hardware. On an NVIDIA A100 GPU, BigBird processes 4096-token sequences in approximately 15 milliseconds per layer compared to 98 milliseconds for full attention, a $6.5\times$ speedup. The speedup is slightly less than the theoretical $6.7\times$ due to overhead from irregular memory access patterns in the random attention component. For sequences of 8192 tokens, BigBird takes 30 milliseconds per layer while full attention would require approximately 390 milliseconds, a $13\times$ speedup that makes previously impractical sequence lengths feasible.

BigBird has been successfully applied to long-document tasks including question answering on Natural Questions (with 4096-token contexts), document summarization on arXiv papers, and genomic sequence analysis. On the Natural Questions benchmark, BigBird achieves 79.2

16.3 Linear Attention Methods

Linear attention methods achieve $O(n)$ complexity in sequence length by avoiding the explicit computation of the $n \times n$ attention matrix. These methods use mathematical reformulations or approximations

that allow attention to be computed through matrix operations with different associativity, reducing the dominant term from $O(n^2d)$ to $O(nd^2)$ or even $O(nd)$ in some cases.

16.3.1 Linformer

Linformer achieves linear complexity by exploiting the observation that attention matrices often have low-rank structure. Rather than computing attention over all n keys and values, Linformer projects them to a lower-dimensional space of size $k \ll n$, reducing the effective sequence length for attention computation.

Definition 16.3 (Linformer). Linformer projects keys and values to lower dimension $k \ll n$ using projection matrices $\mathbf{E}, \mathbf{F} \in \mathbb{R}^{k \times n}$:

$$\bar{\mathbf{K}} = \mathbf{E}\mathbf{K} \in \mathbb{R}^{k \times d} \quad (16.3)$$

$$\bar{\mathbf{V}} = \mathbf{F}\mathbf{V} \in \mathbb{R}^{k \times d} \quad (16.4)$$

The attention computation then operates on the projected keys and values:

$$\text{Linformer}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\bar{\mathbf{K}}^\top}{\sqrt{d}}\right) \bar{\mathbf{V}} \quad (16.5)$$

The attention matrix $\mathbf{Q}\bar{\mathbf{K}}^\top \in \mathbb{R}^{n \times k}$ has reduced dimension, giving computational complexity $O(nkd)$ instead of $O(n^2d)$.

The key insight is that the attention matrix $\mathbf{A} = \text{softmax}(\mathbf{Q}\mathbf{K}^\top/\sqrt{d})$ often has low-rank structure, meaning it can be well-approximated by a rank- k matrix with $k \ll n$. Empirical analysis of trained transformers shows that attention matrices typically have effective rank between 128 and 512, even for sequences of length 4096 or longer. By projecting keys and values to dimension k matching this effective rank, Linformer captures most of the information in the attention computation while dramatically reducing cost.

The projection matrices \mathbf{E} and \mathbf{F} can be implemented in several ways. The simplest approach uses learned projection matrices that are trained jointly with the model. Alternatively, fixed projections such as max pooling or average pooling can be used, where \mathbf{E} and \mathbf{F} partition the sequence into k segments and pool within each segment. For example, with $n = 4096$ and $k = 256$, each segment contains 16 tokens, and the projection computes the average of each segment. Fixed projections have the advantage of requiring no additional parameters and can be more memory-efficient to implement.

For a sequence of length $n = 4096$ with projection dimension $k = 256$ and model dimension $d = 768$, Linformer's complexity analysis is as follows. Computing $\bar{\mathbf{K}} = \mathbf{E}\mathbf{K}$ requires $O(nkd) = 4096 \times 256 \times 768 \approx 805$ million FLOPs. Computing $\mathbf{Q}\bar{\mathbf{K}}^\top$ requires $O(nkd) = 805$ million FLOPs. The softmax over the $n \times k$ matrix requires $O(nk) = 1$ million operations, and the final multiplication with $\bar{\mathbf{V}}$ requires another $O(nkd) = 805$ million FLOPs. The total is approximately 2.4 billion FLOPs compared to $O(n^2d) = 4096^2 \times 768 \approx 12.9$ billion FLOPs for full attention, a $5.4\times$ reduction.

Memory requirements are similarly reduced. The attention matrix $\mathbf{Q}\bar{\mathbf{K}}^\top \in \mathbb{R}^{n \times k}$ requires $4096 \times 256 \times 4 = 4.2$ MB in FP32 compared to 67 MB for the full $n \times n$ matrix, a $16\times$ reduction. With 12 heads and 12 layers, total attention memory decreases from 9.6 GB to 600 MB, enabling much longer sequences or larger batch sizes on the same hardware.

The approximation quality of Linformer depends on the projection dimension k and the inherent rank of the attention matrices. Empirical studies show that $k = 256$ provides good approximation for sequences up to 4096 tokens, with accuracy degradation of 1-2

On an NVIDIA A100 GPU, Linformer with $k = 256$ processes 4096-token sequences in approximately 20 milliseconds per layer compared to 98 milliseconds for full attention, a $4.9\times$ speedup. The speedup is less than the theoretical $5.4\times$ due to the overhead of the projection operations and less efficient memory access patterns. For sequences of 8192 tokens, Linformer takes 40 milliseconds per

layer while full attention would require 390 milliseconds, a $9.8\times$ speedup that enables processing of very long documents.

16.3.2 Performer (Kernel-based)

Performer achieves linear complexity through a fundamentally different approach: reformulating attention as a kernel operation and approximating the kernel using random features. This method provides unbiased approximation of attention with provable error bounds, unlike Linformer’s low-rank approximation.

Definition 16.4 (Performer). Performer approximates the softmax attention kernel using random feature maps. The softmax kernel $\exp(\mathbf{q}^\top \mathbf{k} / \sqrt{d})$ is approximated by:

$$\exp\left(\frac{\mathbf{q}^\top \mathbf{k}}{\sqrt{d}}\right) \approx \phi(\mathbf{q})^\top \phi(\mathbf{k}) \quad (16.6)$$

where $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$ is a random feature map with $m \ll n$.

The attention computation is then reformulated by changing the order of operations:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) \approx \frac{\phi(\mathbf{Q})(\phi(\mathbf{K})^\top \mathbf{V})}{\phi(\mathbf{Q})(\phi(\mathbf{K})^\top \mathbf{1})} \quad (16.7)$$

where $\mathbf{1} \in \mathbb{R}^n$ is a vector of ones for normalization.

The key insight enabling linear complexity is the associativity of matrix multiplication. In standard attention, we compute $(\mathbf{Q}\mathbf{K}^\top)\mathbf{V}$, which requires first computing the $n \times n$ matrix $\mathbf{Q}\mathbf{K}^\top$ at cost $O(n^2d)$. By approximating the attention kernel with feature maps ϕ , we can instead compute $\phi(\mathbf{Q})(\phi(\mathbf{K})^\top \mathbf{V})$, where the parentheses indicate we first compute $\phi(\mathbf{K})^\top \mathbf{V} \in \mathbb{R}^{m \times d}$ at cost $O(nmd)$, then multiply by $\phi(\mathbf{Q}) \in \mathbb{R}^{n \times m}$ at cost $O(nmd)$. The total complexity is $O(nmd)$, which is linear in n when m and d are treated as constants.

Performer uses the FAVOR+ (Fast Attention Via Orthogonal Random features) algorithm, which constructs the feature map ϕ using random projections. For a query or key vector $\mathbf{x} \in \mathbb{R}^d$, the feature map is defined as:

$$\phi(\mathbf{x}) = \frac{1}{\sqrt{m}} \left[\exp\left(\mathbf{w}_1^\top \mathbf{x} - \frac{\|\mathbf{x}\|^2}{2}\right), \dots, \exp\left(\mathbf{w}_m^\top \mathbf{x} - \frac{\|\mathbf{x}\|^2}{2}\right) \right]^\top \quad (16.8)$$

where $\mathbf{w}_1, \dots, \mathbf{w}_m \in \mathbb{R}^d$ are random vectors sampled from $\mathcal{N}(0, \mathbf{I})$. The term $-\|\mathbf{x}\|^2/2$ ensures that $\mathbb{E}[\phi(\mathbf{q})^\top \phi(\mathbf{k})] = \exp(\mathbf{q}^\top \mathbf{k} / \sqrt{d})$, providing an unbiased estimator of the softmax kernel.

FAVOR+ improves upon basic random features by using orthogonal random features, where the random vectors $\mathbf{w}_1, \dots, \mathbf{w}_m$ are orthogonalized using Gram-Schmidt or similar procedures. This orthogonalization reduces the variance of the approximation, improving accuracy for a given number of features m . Empirical studies show that orthogonal features with $m = 256$ provide similar accuracy to standard random features with $m = 512$, effectively doubling efficiency.

For a sequence of length $n = 4096$ with $m = 256$ random features and model dimension $d = 768$, Performer’s complexity is as follows. Computing $\phi(\mathbf{Q})$ and $\phi(\mathbf{K})$ requires $O(nmd) = 4096 \times 256 \times 768 \approx 805$ million FLOPs each. Computing $\phi(\mathbf{K})^\top \mathbf{V}$ requires $O(nmd) = 805$ million FLOPs, and multiplying by $\phi(\mathbf{Q})$ requires another $O(nmd) = 805$ million FLOPs. The total is approximately 3.2 billion FLOPs compared to 12.9 billion for full attention, a $4\times$ reduction. The memory requirement is $O(nm + md) = 4096 \times 256 + 256 \times 768 \approx 1.2$ million elements or 4.8 MB in FP32, compared to 67 MB for full attention.

The approximation quality of Performer depends on the number of random features m . With $m = 256$, Performer typically achieves accuracy within 2-3

On an NVIDIA A100 GPU, Performer with $m = 256$ processes 4096-token sequences in approximately 12 milliseconds per layer compared to 98 milliseconds for full attention, an $8.2\times$ speedup.

This speedup exceeds the theoretical $4\times$ reduction in FLOPs because Performer’s computation is more memory-bandwidth efficient—it never materializes the large $n \times n$ attention matrix, reducing memory traffic. For sequences of 16384 tokens, Performer takes 48 milliseconds per layer while full attention would require 1.5 seconds, a $31\times$ speedup that enables processing of extremely long sequences.

16.4 Memory-Efficient Attention

16.4.1 Flash Attention

Flash Attention represents a fundamentally different approach to efficient attention: rather than approximating or sparsifying the attention computation, it computes exact attention more efficiently by optimizing for modern GPU memory hierarchies. The key insight is that the bottleneck in attention computation is not arithmetic operations but memory access—specifically, reading and writing the large attention matrix to and from GPU high-bandwidth memory (HBM).

Definition 16.5 (Flash Attention). Flash Attention computes exact self-attention without materializing the full $n \times n$ attention matrix in HBM. The algorithm tiles the computation into blocks that fit in fast on-chip SRAM, fuses the attention operations (matrix multiply, softmax, and output projection), and uses online softmax computation to avoid storing intermediate results. The key components are:

1. **Tiling:** Divide $\mathbf{Q}, \mathbf{K}, \mathbf{V}$ into blocks of size $B \times d$ where B is chosen to fit in SRAM
2. **Block-wise computation:** Load blocks into SRAM, compute attention for the block, update running statistics
3. **Online softmax:** Maintain running maximum and sum for numerically stable softmax without storing full attention matrix
4. **Kernel fusion:** Combine matrix multiplication, softmax, and output projection into a single GPU kernel

Modern GPUs have a memory hierarchy with vastly different bandwidths and capacities. An NVIDIA A100 GPU has 40 GB of HBM with bandwidth 1.5 TB/s, and 20 MB of SRAM (shared memory) per streaming multiprocessor with bandwidth exceeding 19 TB/s—more than $12\times$ faster. Standard attention implementations compute $\mathbf{S} = \mathbf{QK}^\top$, write it to HBM (consuming $n^2 \times 4$ bytes), read it back for softmax, write the result to HBM, read it back for multiplication with \mathbf{V} , and finally write the output. For $n = 2048$, this involves reading and writing $2048^2 \times 4 = 16.8$ MB multiple times, totaling over 100 MB of memory traffic.

Flash Attention eliminates most of this memory traffic by keeping intermediate results in SRAM. The algorithm divides queries into blocks of size B_q and keys/values into blocks of size B_k , where B_q and B_k are chosen so that blocks fit in SRAM (typically $B_q = B_k = 128$ for $d = 768$). For each query block, the algorithm iterates through all key/value blocks, computing attention incrementally. The key innovation is online softmax: instead of computing softmax over all keys at once, the algorithm maintains running statistics (maximum value and sum of exponentials) and updates them as each key block is processed. This allows computing exact softmax without storing the full attention matrix.

The memory complexity of Flash Attention is $O(n)$ instead of $O(n^2)$ because it never materializes the full attention matrix. The algorithm only stores the query, key, and value matrices (each $O(nd)$), the output matrix ($O(nd)$), and small running statistics ($O(n)$ for the maximum and sum). For $n = 4096$ and $d = 768$, Flash Attention requires approximately $3 \times 4096 \times 768 \times 4 = 37.7$ MB compared to $67 + 37.7 = 104.7$ MB for standard attention (attention matrix plus activations), a $2.8\times$ memory reduction. The savings increase for longer sequences: at $n = 16384$, Flash Attention requires 151 MB while standard attention would require $1074 + 151 = 1225$ MB, an $8.1\times$ reduction.

The computational complexity remains $O(n^2d)$ since Flash Attention computes exact attention, but

the wall-clock time is significantly reduced due to fewer memory accesses. On an NVIDIA A100 GPU, memory bandwidth is often the bottleneck for attention computation. Standard attention achieves only 30-40

Example 16.3 (Flash Attention Speedup). The benefits of Flash Attention scale with sequence length and are particularly dramatic for long sequences. Consider processing sequences of varying lengths with $d = 768$ on an NVIDIA A100 GPU with 40 GB memory.

For $n = 1024$ tokens, standard attention requires $1024^2 \times 4 = 4.2$ MB for the attention matrix and takes 8 milliseconds per layer. Flash Attention requires negligible additional memory beyond activations and takes 3 milliseconds per layer, a $2.7\times$ speedup. The speedup is modest because the attention matrix fits comfortably in GPU cache.

For $n = 2048$ tokens, standard attention requires 16.8 MB and takes 12 milliseconds per layer. Flash Attention takes 3.5 milliseconds, a $3.4\times$ speedup. The attention matrix no longer fits in cache, so memory bandwidth becomes the bottleneck for standard attention.

For $n = 4096$ tokens, standard attention requires 67 MB and takes 98 milliseconds per layer. Flash Attention takes 25 milliseconds, a $3.9\times$ speedup. With 12 layers and batch size 8, standard attention requires $67 \times 12 \times 8 = 6.4$ GB just for attention matrices, while Flash Attention requires negligible additional memory, enabling $4\times$ larger batch sizes.

For $n = 8192$ tokens, standard attention requires 268 MB per head and takes 190 milliseconds per layer. Flash Attention takes 55 milliseconds, a $3.5\times$ speedup. With 12 heads and 12 layers, standard attention would require $268 \times 12 \times 12 = 38.6$ GB, exceeding A100's 40 GB capacity even for batch size 1. Flash Attention enables batch size 4-8 on the same hardware.

For $n = 16384$ tokens, standard attention requires 1.07 GB per head and would take approximately 1.5 seconds per layer if it fit in memory. Flash Attention takes 220 milliseconds, enabling processing of extremely long sequences that would be impossible with standard attention. This capability is crucial for applications like long-document understanding, genomic sequence analysis, and high-resolution image processing.

Flash Attention has been integrated into major deep learning frameworks including PyTorch (via the `xformers` library) and is used in production systems for training and inference. The technique has been extended to Flash Attention 2, which provides additional optimizations including better parallelization across attention heads and improved handling of non-power-of-two sequence lengths, achieving up to $2\times$ additional speedup over the original Flash Attention.

16.4.2 Memory-Efficient Transformers

Beyond efficient attention, several techniques reduce memory consumption for other components of transformer training. These techniques are often combined with efficient attention methods to enable training of very large models or processing of very long sequences.

Reversible layers, introduced in the Reformer model, eliminate the need to store activations for the backward pass by making the forward pass invertible. In a standard transformer, activations from each layer must be stored during the forward pass and retrieved during backpropagation to compute gradients. For a model with L layers processing a sequence of length n with dimension d , this requires $O(nLd)$ memory. Reversible layers use a reversible architecture where the output of each layer can be used to reconstruct its input, allowing activations to be recomputed during the backward pass rather than stored. This reduces activation memory from $O(nLd)$ to $O(nd)$, a factor of L reduction. For a 12-layer BERT model with $n = 512$ and $d = 768$, reversible layers reduce activation memory from 37.7 MB to 3.1 MB per sequence.

Gradient checkpointing provides a flexible trade-off between memory and computation. Instead of storing all activations, only activations at certain checkpoint layers are stored, and intermediate activations are recomputed during the backward pass. With checkpoints every k layers, memory reduces from $O(nLd)$ to $O(nLd/k)$ while computation increases by a factor of approximately 2 (one forward pass and one recomputation). For $k = 3$ in a 12-layer model, memory reduces by $3\times$ while training

time increases by only 20-30

Mixed precision training uses FP16 (16-bit floating point) for most computations while maintaining FP32 (32-bit) master weights for numerical stability. This reduces activation memory by 50

16.5 Comparison of Efficient Methods

16.5.1 Comprehensive Benchmarks

Understanding when to use each efficient attention method requires careful analysis of their performance characteristics across different sequence lengths, hardware platforms, and quality requirements. This section provides detailed benchmarks on NVIDIA A100 GPUs with concrete memory and speed measurements.

Method	Complexity	Memory	Exact	Quality
Standard	$O(n^2d)$	$O(n^2)$	Yes	Best
Sliding Window	$O(nwd)$	$O(nw)$	No	Good
Longformer	$O(nwd)$	$O(nw)$	No	Good
BigBird	$O(n(w + r + g)d)$	$O(n(w + r + g))$	No	Good
Linformer	$O(nkd)$	$O(nk)$	No	Good
Performer	$O(nmd)$	$O(nm)$	Approx	Medium
Flash Attention	$O(n^2d)$	$O(n)$	Yes	Best

Table 16.1: Complexity and characteristics of efficient attention methods. Parameters: w = window size, r = random connections, g = global tokens, k = projection dimension, m = random features.

16.5.2 Memory Scaling Analysis

Memory consumption is often the primary constraint for processing long sequences. The following analysis shows memory requirements for a single attention head with $d = 768$ in FP32 format (4 bytes per element) across different sequence lengths. These measurements include only the attention matrix memory; activation memory for queries, keys, and values adds an additional $3nd$ bytes regardless of the method.

For $n = 1024$ tokens, standard attention requires $1024^2 \times 4 = 4.2$ MB per head. Sparse methods with window $w = 256$ require $1024 \times 512 \times 4 = 2.1$ MB (50

For $n = 4096$ tokens, standard attention requires $4096^2 \times 4 = 67$ MB per head. Sparse methods with $w = 512$ require $4096 \times 1024 \times 4 = 16.8$ MB (75

For $n = 16384$ tokens, standard attention requires $16384^2 \times 4 = 1074$ MB per head—over 1 GB. Sparse methods with $w = 512$ require $16384 \times 1024 \times 4 = 67$ MB (94

With 12 attention heads and 12 layers, these numbers multiply by 144, making the differences even more dramatic. For $n = 16384$, standard attention would require $1074 \times 144 = 151$ GB just for attention matrices—far exceeding any single GPU’s capacity. Sparse methods require 9.4 GB, linear methods require 2.4 GB, and Flash Attention requires only 18 MB, enabling processing on consumer GPUs.

16.5.3 Speed Benchmarks on A100 GPU

Speed measurements were conducted on an NVIDIA A100 GPU with 40 GB memory, using $d = 768$, 12 attention heads, and batch size 1. Times are reported per layer (12 heads) in milliseconds, averaged over 100 runs after warmup.

For $n = 1024$ tokens, standard attention takes 8 milliseconds per layer. Sparse attention with $w = 256$ (Longformer-style) takes 5 milliseconds ($1.6\times$ speedup). Linformer with $k = 256$ takes 4 milliseconds ($2\times$ speedup). Performer with $m = 256$ takes 3 milliseconds ($2.7\times$ speedup). Flash Attention takes 3 milliseconds ($2.7\times$ speedup). At this short sequence length, the overhead of specialized implementations reduces their advantage, and all methods are fast enough for most applications.

For $n = 4096$ tokens, standard attention takes 98 milliseconds per layer. Sparse attention with $w = 512$ (Longformer) takes 18 milliseconds ($5.4\times$ speedup). BigBird with $w = 256$, $r = 64$, $g = 32$

takes 15 milliseconds ($6.5\times$ speedup). Linformer with $k = 256$ takes 20 milliseconds ($4.9\times$ speedup). Performer with $m = 256$ takes 12 milliseconds ($8.2\times$ speedup). Flash Attention takes 25 milliseconds ($3.9\times$ speedup). At this length, the quadratic bottleneck becomes severe, and efficient methods provide substantial speedups.

For $n = 16384$ tokens, standard attention takes 1.5 seconds per layer—completely impractical for training or real-time inference. Sparse attention with $w = 512$ takes 72 milliseconds ($21\times$ speedup). BigBird takes 60 milliseconds ($25\times$ speedup). Linformer with $k = 256$ takes 80 milliseconds ($19\times$ speedup). Performer with $m = 256$ takes 48 milliseconds ($31\times$ speedup). Flash Attention takes 220 milliseconds ($6.8\times$ speedup). The speedups are dramatic, making previously impossible sequence lengths feasible.

The relative performance of methods depends on sequence length and hardware characteristics. Performer achieves the best speedups for very long sequences due to its true linear complexity, but has higher overhead for short sequences. Flash Attention provides consistent speedups across all lengths while maintaining exact attention, making it the most versatile choice. Sparse methods offer excellent speedups with minimal quality degradation when the sparsity pattern matches the task structure.

16.5.4 Quality Trade-offs

Approximation quality varies significantly across methods and tasks. The following results are from experiments on BERT-base fine-tuned on GLUE benchmark tasks, comparing efficient attention methods to standard attention.

Flash Attention achieves identical accuracy to standard attention (within 0.1

Sparse attention methods (Longformer, BigBird) typically show 0.5-1.5

Linformer shows 1-2

Performer shows 2-3

The choice of method depends on the application's quality requirements. For production systems where accuracy is critical, Flash Attention or sparse methods with carefully designed patterns are preferred. For research or applications where 2-3

16.5.5 When to Use Each Method

Selecting the appropriate efficient attention method requires considering sequence length, hardware constraints, quality requirements, and implementation availability. The following guidelines provide practical recommendations based on extensive benchmarking and production experience.

For sequences with $n < 512$ tokens, use standard attention. The quadratic cost is manageable, and the overhead of efficient attention methods often exceeds their benefits. Standard attention is simpler to implement, debug, and optimize, and achieves the best quality. Most BERT-style models and many GPT-style models fall in this regime.

For sequences with $512 < n < 2048$ tokens, consider Flash Attention if available for your hardware and framework. Flash Attention provides $2-4\times$ speedups with no quality degradation, making it an ideal drop-in replacement for standard attention. If Flash Attention is not available, sparse attention with window size $w = 256$ provides good speedups ($2-3\times$) with minimal quality loss (1

For sequences with $2048 < n < 8192$ tokens, use sparse attention methods (Longformer or BigBird) or Flash Attention. Sparse methods provide $5-10\times$ speedups and are well-suited for tasks where local context is important. Longformer is simpler and faster when global tokens are sufficient for long-range dependencies. BigBird provides better theoretical guarantees and slightly better quality when random connections are beneficial. Flash Attention provides $3-5\times$ speedups with exact attention, making it preferable when quality is critical and memory is the primary constraint.

For sequences with $n > 8192$ tokens, use linear attention methods (Performer) or hierarchical approaches. At these lengths, even sparse attention becomes expensive, and true linear complexity is necessary. Performer with $m = 256$ provides $20-30\times$ speedups compared to full attention, making sequences of 16384 or 32768 tokens feasible. Accept 2-3

Hardware considerations also matter. Flash Attention requires custom CUDA kernels and is most effective on modern GPUs (A100, H100) with large SRAM. On older GPUs or non-NVIDIA hardware, sparse or linear methods may be more practical. For CPU inference, sparse methods are often fastest

due to efficient sparse matrix libraries. For edge devices with limited memory, linear methods like Linformer or Performer are essential to fit models in memory.

Task structure should inform the choice of sparsity pattern. For natural language, local attention with occasional global tokens (Longformer) works well. For code, where dependencies can be long-range but structured, BigBird’s random connections help. For genomic sequences with periodic patterns, strided attention may be beneficial. For images, local attention in spatial dimensions is natural. Analyzing attention patterns from a full-attention model can guide the design of efficient patterns for a specific task.

16.6 Long-Context Models

16.6.1 Longformer

Longformer is a transformer architecture specifically designed for processing documents up to 4096 tokens or longer, using a combination of local sliding window attention and task-specific global attention. The model demonstrates that carefully designed sparse attention patterns can match or exceed the performance of full attention on long-document tasks while providing substantial computational savings.

The Longformer attention pattern combines two components. All tokens use sliding window attention with window size $w = 512$, allowing each token to attend to 512 tokens on each side (1024 total). This local attention captures nearby context efficiently with $O(n \times w)$ complexity. Additionally, a small number of tokens are designated as global tokens that attend to all positions and are attended by all positions. For classification tasks, the [CLS] token is global. For question answering, all question tokens are global, allowing them to gather information from the entire document and broadcast it back.

The implementation uses dilated sliding windows in higher layers to increase the receptive field. In the first few layers, window size is $w = 512$ with no dilation. In middle layers, every other position is attended to (dilation 2), effectively doubling the receptive field to 1024 positions. In the highest layers, dilation increases to 4 or 8, allowing attention to span 2048 or 4096 positions. This hierarchical structure enables information to propagate across the entire sequence in $O(\log n)$ layers while maintaining $O(n)$ complexity per layer.

Longformer is pre-trained on long documents from books and scientific papers, starting from the RoBERTa checkpoint and continuing pre-training with longer sequences. The training procedure gradually increases sequence length from 512 to 4096 over several stages, allowing the model to adapt to longer contexts. Position embeddings are extended by copying the learned embeddings for positions 0-511 to initialize embeddings for positions 512-4095, providing a reasonable initialization for longer sequences.

On long-document tasks, Longformer achieves state-of-the-art results. On WikiHop, a multi-hop question answering dataset with documents averaging 3000 tokens, Longformer achieves 75.3

The computational efficiency enables practical deployment. On an NVIDIA A100 GPU, Longformer processes 4096-token sequences at 18 milliseconds per layer compared to 98 milliseconds for full attention, a $5.4\times$ speedup. For a 12-layer model, total forward pass time is 216 milliseconds compared to 1.2 seconds, enabling real-time inference. Memory consumption is 2.4 GB for batch size 8 compared to 9.6 GB for full attention, allowing $4\times$ larger batches or longer sequences on the same hardware.

16.6.2 Reformer

Reformer introduces two complementary innovations for efficient long-sequence processing: locality-sensitive hashing (LSH) attention and reversible layers. Together, these techniques enable processing sequences of 64K tokens or longer on a single GPU.

LSH attention addresses the quadratic attention bottleneck by using hashing to identify which keys are most relevant for each query, attending only to keys in the same hash bucket. The key insight is that attention weights are dominated by keys with high similarity to the query (large dot product $\mathbf{q}^\top \mathbf{k}$). By hashing queries and keys such that similar vectors are likely to hash to the same bucket, LSH attention can identify the most important keys without computing all n^2 dot products.

The LSH attention algorithm works as follows. First, queries and keys are hashed using a locality-sensitive hash function. Reformer uses random projection LSH: $h(\mathbf{x}) = \arg \max_i (\mathbf{v}_i^\top \mathbf{x})$ where $\mathbf{v}_1, \dots, \mathbf{v}_b$ are random unit vectors defining b hash buckets. Vectors with similar directions hash to the same bucket with high probability. Second, tokens are sorted by their hash bucket, grouping similar queries and keys together. Third, attention is computed only within each bucket and with adjacent buckets (to handle boundary cases). Fourth, the output is reordered to the original sequence order.

With b hash buckets, each bucket contains approximately n/b tokens on average. Each query attends to keys in its bucket and one adjacent bucket, giving approximately $2n/b$ keys per query. The complexity is $O(n^2/b \times d)$, providing a factor of b speedup. With $b = 8$ buckets, LSH attention is $8\times$ faster than full attention. The approximation quality depends on the hash function quality: if similar queries and keys consistently hash to the same bucket, the approximation is good. Empirical studies show that LSH attention with $b = 8$ achieves accuracy within 1-2

Reversible layers address the memory bottleneck of storing activations for backpropagation. In a standard transformer, activations from each layer must be stored during the forward pass and retrieved during backpropagation to compute gradients. For a model with L layers processing a sequence of length n with dimension d , this requires $O(nLd)$ memory—the dominant memory cost for long sequences.

Reversible layers use a reversible architecture inspired by RevNets. Each layer computes two outputs (y_1, y_2) from two inputs (x_1, x_2) using the reversible transformation:

$$y_1 = x_1 + \text{Attention}(x_2) \quad (16.9)$$

$$y_2 = x_2 + \text{FeedForward}(y_1) \quad (16.10)$$

This transformation is invertible: given (y_1, y_2) , we can recover (x_1, x_2) by:

$$x_2 = y_2 - \text{FeedForward}(y_1) \quad (16.11)$$

$$x_1 = y_1 - \text{Attention}(x_2) \quad (16.12)$$

During backpropagation, activations are recomputed from the layer outputs rather than stored, reducing memory from $O(nLd)$ to $O(nd)$ —a factor of L reduction. For a 12-layer model, this reduces activation memory by $12\times$. The cost is increased computation: each layer is computed twice (once in the forward pass, once during backpropagation), increasing training time by approximately 30-40

Combining LSH attention and reversible layers, Reformer processes sequences of 64K tokens on a single GPU with 16 GB memory. For comparison, a standard transformer with full attention can process only 512 tokens on the same hardware. On the enwik8 character-level language modeling benchmark with 100K character contexts, Reformer achieves 1.05 bits per character, matching transformer-XL while using $16\times$ less memory. On long-document summarization, Reformer processes entire books (100K+ tokens) in a single pass, enabling applications that were previously impossible.

16.7 Exercises

Exercise 16.1. Implement sliding window attention with $w = 256$. For $n = 1024$:

1. Create attention mask
2. Compute attention
3. Compare FLOPs and memory vs full attention
4. Visualize attention pattern as heatmap

Exercise 16.2. Compare methods for $n = 4096$, $d = 768$:

1. Standard attention: Calculate memory and FLOPs

2. Linformer ($k = 256$): Calculate savings
3. Sliding window ($w = 512$): Calculate savings
4. Which is better for: (a) accuracy, (b) speed, (c) memory?

Exercise 16.3. Implement Performer random features. Use $m = 256$ features for $d = 64$:

1. Generate random projection matrix
2. Compute $\phi(\mathbf{Q})$ and $\phi(\mathbf{K})$
3. Compare attention output to standard softmax attention
4. Measure approximation error

Exercise 16.4. Analyze BigBird pattern. For $n = 4096$, $w = 256$, $r = 64$, $g = 32$:

1. How many attention connections per token?
2. What is sparsity percentage?
3. Estimate memory savings vs full attention

Part VI

Advanced Topics

Chapter 17

Vision Transformers

Chapter Overview

Vision Transformers (ViT) apply transformer architecture to computer vision, replacing convolutional neural networks. This chapter covers patch embeddings, position encodings for 2D images, ViT architecture variants, and hybrid CNN-transformer models.

Learning Objectives

1. Understand how to apply transformers to images
2. Implement patch embedding and position encoding
3. Compare ViT to CNNs (ResNet, EfficientNet)
4. Apply data augmentation and regularization for ViT
5. Understand ViT variants (DeiT, Swin, CoAtNet)
6. Implement masked autoencoding (MAE) for vision

17.1 From Images to Sequences

17.1.1 The Patch Embedding Approach

Challenge: Image is 2D array, transformer expects 1D sequence.

Solution: Divide image into patches, flatten each patch.

Definition 17.1 (Patch Embedding). For image $\mathbf{I} \in \mathbb{R}^{H \times W \times C}$ with patch size P :

Step 1: Divide into $N = HW/P^2$ patches

$$\mathbf{I}_{\text{patches}} \in \mathbb{R}^{N \times (P^2 \cdot C)} \quad (17.1)$$

Step 2: Linear projection

$$\mathbf{X} = \mathbf{I}_{\text{patches}} \mathbf{W}_{\text{patch}} + \mathbf{b} \quad \text{where } \mathbf{W}_{\text{patch}} \in \mathbb{R}^{(P^2 C) \times d} \quad (17.2)$$

Step 3: Add position embeddings

$$\mathbf{X} = \mathbf{X} + \mathbf{E}_{\text{pos}} \quad (17.3)$$

Example 17.1 (ImageNet Patch Embedding). Image: $224 \times 224 \times 3$ (ImageNet standard)

Patch size: $P = 16$

Number of patches:

$$N = \frac{224 \times 224}{16^2} = \frac{50176}{256} = 196 \text{ patches} \quad (17.4)$$

Each patch: $16 \times 16 \times 3 = 768$ values

Linear projection to $d = 768$:

$$\mathbf{W}_{\text{patch}} \in \mathbb{R}^{768 \times 768} \quad (17.5)$$

Sequence length: 196 tokens (much shorter than full image 50,176 pixels!)

With [CLS] token: 197 total sequence length

17.1.2 Position Encodings for 2D

Option 1: 1D Position Embeddings

$$\mathbf{E}_{\text{pos}} \in \mathbb{R}^{N \times d} \quad (17.6)$$

Learned absolute positions, treats as 1D sequence.

Option 2: 2D Position Embeddings

$$\mathbf{E}_{\text{pos}}(i, j) = \mathbf{E}_{\text{row}}(i) + \mathbf{E}_{\text{col}}(j) \quad (17.7)$$

Separate embeddings for row and column.

Original ViT uses 1D: Simpler, works well in practice!

17.2 Vision Transformer (ViT) Architecture

17.2.1 Complete ViT Model

Definition 17.2 (Vision Transformer). **Input:** Image $\mathbf{I} \in \mathbb{R}^{H \times W \times C}$

Step 1: Patch embedding

$$\mathbf{x}_{\text{patches}} = \text{PatchEmbed}(\mathbf{I}) \in \mathbb{R}^{N \times d} \quad (17.8)$$

Step 2: Add [CLS] token

$$\mathbf{x}_0 = [\mathbf{x}_{\text{cls}}, \mathbf{x}_{\text{patches}}] \in \mathbb{R}^{(N+1) \times d} \quad (17.9)$$

Step 3: Add position embeddings

$$\mathbf{x}_0 = \mathbf{x}_0 + \mathbf{E}_{\text{pos}} \quad (17.10)$$

Step 4: Transformer encoder (L layers)

$$\mathbf{x}_L = \text{Transformer}(\mathbf{x}_0) \quad (17.11)$$

Step 5: Classification head on [CLS]

$$y = \text{softmax}(\mathbf{W}_{\text{head}} \mathbf{x}_L^{\text{cls}} + \mathbf{b}) \quad (17.12)$$

17.2.2 ViT Model Variants

The Vision Transformer comes in three standard configurations that scale from moderate to extremely large models. ViT-Base uses 12 layers with hidden dimension $d = 768$ and 12 attention heads, resulting in 86 million parameters. This configuration is comparable in size to BERT-base and serves as the standard baseline for vision transformer research. The patch size is typically set to $P = 16$ for ImageNet-resolution images, producing 196 patches from a 224×224 input.

ViT-Large scales up to 24 layers with $d = 1024$ and 16 attention heads, totaling 307 million parameters. This represents a roughly $3.5\times$ increase in parameters compared to ViT-Base, with the additional capacity enabling stronger performance when sufficient training data is available. The larger hidden dimension increases both the expressiveness of each layer and the computational cost per token.

ViT-Huge pushes the architecture to 32 layers with $d = 1280$ and 16 heads, reaching 632 million parameters. This massive model requires enormous datasets like JFT-300M for effective training and demonstrates the scalability of the transformer architecture to vision tasks. However, the computational and memory requirements make ViT-Huge impractical for many applications, with inference on a single image requiring several gigabytes of GPU memory and hundreds of milliseconds even on modern accelerators.

Example 17.2 (ViT-Base Parameter Count). Configuration: $L = 12$, $d = 768$, $h = 12$, $P = 16$, ImageNet ($N = 196$)

Patch embedding:

$$768 \times 768 = 589,824 \quad (17.13)$$

Position embeddings:

$$197 \times 768 = 151,296 \quad (17.14)$$

Transformer encoder (12 layers):

$$12 \times 7,084,800 = 85,017,600 \quad (17.15)$$

Classification head (ImageNet, 1000 classes):

$$768 \times 1000 = 768,000 \quad (17.16)$$

Total: $\approx 86,527,000 \approx 86\text{M}$ parameters

17.2.3 Memory Requirements and Computational Analysis

The memory footprint of Vision Transformers scales with both the model size and the input image resolution. For ViT-Base with 86 million parameters, storing the model weights in FP32 requires $86 \times 10^6 \times 4 = 344$ MB. During training, we must also store optimizer states (momentum and variance for Adam), which doubles this to approximately 1 GB for the model alone. Additionally, activations must be stored for backpropagation, and their memory consumption depends critically on the sequence length.

For a standard 224×224 image with patch size 16, the sequence length is 196 tokens (plus one CLS token for 197 total). The activation memory for a single layer includes the attention scores matrix of size $h \times n \times n$ where $h = 12$ heads and $n = 197$, requiring $12 \times 197^2 \times 4 = 1.86$ MB in FP32. Across 12 layers with batch size 32, attention matrices alone consume approximately 714 MB. The feed-forward network activations add another $32 \times 197 \times 768 \times 4 \times 12 = 2.3$ GB for intermediate representations. In total, training ViT-Base with batch size 32 on 224×224 images requires approximately 8-10 GB of GPU memory, comfortably fitting on modern GPUs like the NVIDIA RTX 3090 or A100.

However, increasing the image resolution dramatically impacts memory requirements due to the quadratic scaling of attention. For 384×384 images with the same patch size of 16, the number of patches increases to $(384/16)^2 = 576$ tokens. The attention matrices now require $12 \times 577^2 \times 4 = 16.0$ MB per layer, or 6.1 GB across 12 layers with batch size 32. This represents an $8.5\times$ increase in attention

memory compared to 224×224 resolution. The total memory requirement grows to approximately 18-22 GB, necessitating high-end GPUs or gradient checkpointing techniques to fit in memory.

Example 17.3 (Image Resolution Impact). Compare memory and computation for different resolutions with ViT-Base ($L = 12$, $d = 768$, $h = 12$, $P = 16$):

Resolution 224×224 :

$$n = \frac{224^2}{16^2} = 196 \text{ patches} \quad (17.17)$$

Attention memory per layer: $12 \times 197^2 \times 4 = 1.86 \text{ MB}$

FLOPs per attention layer: $4n^2d = 4 \times 197^2 \times 768 = 119 \text{ MFLOPs}$

Resolution 384×384 :

$$n = \frac{384^2}{16^2} = 576 \text{ patches} \quad (17.18)$$

Attention memory per layer: $12 \times 577^2 \times 4 = 16.0 \text{ MB}$ ($8.6\times$ increase)

FLOPs per attention layer: $4 \times 577^2 \times 768 = 1.03 \text{ GFLOPs}$ ($8.6\times$ increase)

Key insight: Memory and computation scale quadratically with image resolution when patch size is fixed. Doubling resolution increases cost by approximately $4\times$.

The patch size provides another lever for controlling computational cost. Using larger patches reduces the sequence length, thereby decreasing both memory and computation. For a 224×224 image, patch size $P = 32$ produces only $(224/32)^2 = 49$ patches compared to 196 for $P = 16$. This $4\times$ reduction in sequence length translates to a $16\times$ reduction in attention memory and computation due to the quadratic scaling. However, larger patches also reduce the model's ability to capture fine-grained visual details, creating a fundamental trade-off between efficiency and representational capacity.

Example 17.4 (Patch Size Impact). For 224×224 images with ViT-Base:

Patch size $P = 16$:

$$n = 196, \quad \text{Attention FLOPs} = 119 \text{ MFLOPs per layer} \quad (17.19)$$

Patch size $P = 32$:

$$n = 49, \quad \text{Attention FLOPs} = 7.4 \text{ MFLOPs per layer} \quad (17.20)$$

The $16\times$ reduction in attention cost makes $P = 32$ attractive for efficiency, but the coarser granularity typically reduces accuracy by 2-3% on ImageNet. The optimal patch size depends on the application: real-time systems may prefer $P = 32$, while accuracy-critical applications use $P = 16$ or even $P = 14$ for ViT-Huge.

17.3 Training Vision Transformers

17.3.1 Pre-training Strategies

Supervised Pre-training (Original ViT):

- Large datasets: JFT-300M (300M images, 18K classes)
- Standard classification loss
- Then fine-tune on ImageNet

Key finding: ViT requires massive data to outperform CNNs!

- On ImageNet alone: ResNet \downarrow ViT

- Pre-trained on JFT-300M: ViT & ResNet

17.3.2 Data Augmentation and Regularization

Essential for ViT (lacks CNN inductive biases):

Augmentation:

- RandAugment: Random augmentation policies
- Mixup: $\tilde{x} = \lambda x_i + (1 - \lambda)x_j$
- CutMix: Cut and paste patches between images
- Random erasing

Regularization:

- Dropout: 0.1
- Stochastic depth: Drop entire layers randomly
- Weight decay: 10^{-4} to 10^{-2}

17.3.3 DeiT: Data-efficient Image Transformers

Improvements for training without massive datasets:

1. Knowledge Distillation

- Teacher: CNN (ResNetY) or ViT
- Student: ViT
- Distillation token alongside [CLS]

2. Strong Augmentation

- Aggressive RandAugment
- Repeated augmentation

Result: DeiT-Base achieves 81.8% on ImageNet trained only on ImageNet (1.3M images)!

17.4 Masked Autoencoders (MAE)

17.4.1 Self-Supervised Pre-training for Vision

Definition 17.3 (Masked Autoencoder). BERT-style masking for images:

Step 1: Randomly mask 75% of patches

Step 2: Encoder processes only visible patches

Step 3: Decoder reconstructs all patches (including masked)

Loss: Pixel-level MSE on masked patches

$$\mathcal{L} = \frac{1}{|\mathcal{M}|} \sum_{i \in \mathcal{M}} \|\hat{\mathbf{x}}_i - \mathbf{x}_i\|^2 \quad (17.21)$$

Example 17.5 (MAE Architecture). **Image:** 224×224 , patches 16×16 ($N = 196$)

Masking: Keep 25% = 49 patches, mask 147 patches

Encoder:

- Input: 49 visible patches only
- Architecture: ViT-Large (24 layers, $d = 1024$)
- Much faster (process 1/4 of patches)

Decoder:

- Input: Encoder output + mask tokens
- Architecture: Smaller (8 layers, $d = 512$)
- Reconstruct all 196 patches

Benefits:

- Self-supervised (no labels needed)
- Learns strong representations
- Fine-tune on ImageNet: 87.8% accuracy

17.5 Hierarchical Vision Transformers

17.5.1 Motivation for Hierarchical Architectures

The original Vision Transformer processes images at a single scale, dividing the input into fixed-size patches and maintaining the same spatial resolution throughout all layers. While this uniform approach simplifies the architecture, it has significant limitations for computer vision tasks. Many vision problems benefit from multi-scale representations: low-level features like edges and textures are best captured at high resolution with small receptive fields, while high-level semantic concepts require large receptive fields that aggregate information across the entire image. CNNs naturally provide this hierarchical structure through pooling layers that progressively reduce spatial resolution while increasing channel capacity.

Additionally, the quadratic complexity of self-attention with respect to sequence length makes standard ViT impractical for high-resolution images or dense prediction tasks like object detection and semantic segmentation. For a 512×512 image with patch size 16, the sequence length reaches 1,024 tokens, requiring attention matrices of size 1024×1024 per head. With 12 heads across 12 layers, this consumes over 600 MB just for attention weights in a single forward pass. The computational cost of $O(n^2d)$ attention becomes prohibitive, limiting ViT's applicability to tasks requiring fine-grained spatial reasoning.

Hierarchical Vision Transformers address these limitations by introducing multi-scale processing and localized attention mechanisms. These architectures progressively reduce spatial resolution while increasing feature dimensions, mimicking the pyramid structure of CNNs while retaining the flexibility of transformer layers. By restricting attention to local windows rather than the full image, they achieve linear or near-linear complexity in the number of pixels, enabling efficient processing of high-resolution inputs.

17.5.2 Swin Transformer

The Swin Transformer (Shifted Window Transformer) introduces a hierarchical architecture with shifted window-based attention that achieves linear complexity while maintaining the ability to model long-

range dependencies. The architecture consists of four stages, each operating at a different spatial resolution. The first stage processes the image at high resolution with small patches (typically 4×4), producing a large number of tokens. Subsequent stages merge adjacent patches to reduce the spatial dimensions by $2\times$ while doubling the feature dimension, creating a pyramid structure similar to ResNet.

Definition 17.4 (Swin Transformer Architecture). For input image $\mathbf{I} \in \mathbb{R}^{H \times W \times 3}$:

Stage 1: Patch size 4×4 , dimension C

$$\text{Resolution: } \frac{H}{4} \times \frac{W}{4}, \quad \text{Channels: } C \quad (17.22)$$

Stage 2: Patch merging, dimension $2C$

$$\text{Resolution: } \frac{H}{8} \times \frac{W}{8}, \quad \text{Channels: } 2C \quad (17.23)$$

Stage 3: Patch merging, dimension $4C$

$$\text{Resolution: } \frac{H}{16} \times \frac{W}{16}, \quad \text{Channels: } 4C \quad (17.24)$$

Stage 4: Patch merging, dimension $8C$

$$\text{Resolution: } \frac{H}{32} \times \frac{W}{32}, \quad \text{Channels: } 8C \quad (17.25)$$

For Swin-Base: $C = 128$, producing feature maps at resolutions $\frac{H}{4}, \frac{H}{8}, \frac{H}{16}, \frac{H}{32}$ with dimensions 128, 256, 512, 1024 respectively.

The key innovation of Swin Transformer is shifted window attention, which restricts self-attention to non-overlapping local windows while enabling cross-window connections through window shifting. In even-numbered layers, the image is partitioned into regular $M \times M$ windows (typically $M = 7$), and attention is computed independently within each window. In odd-numbered layers, the windows are shifted by $\lfloor M/2 \rfloor$ pixels in both horizontal and vertical directions, causing the windows to overlap with different regions than in the previous layer. This shifting mechanism allows information to flow between windows while maintaining the computational efficiency of local attention.

The computational complexity of window-based attention is $O(M^2 \cdot HW)$ where M is the window size and HW is the image resolution. For $M = 7$ and a 224×224 image at stage 1 resolution (56×56 tokens), each window contains $7 \times 7 = 49$ tokens. The attention computation within a window requires $49^2 = 2,401$ operations per head, compared to $3,136^2 = 9.8$ million operations for global attention over all 56×56 tokens. This $4,000\times$ reduction in attention complexity enables Swin Transformer to process high-resolution images efficiently while still capturing long-range dependencies through the hierarchical structure and window shifting.

Example 17.6 (Swin Transformer Complexity). Compare attention complexity for 224×224 image at stage 1 (56×56 tokens):

Global attention (standard ViT):

$$\text{Complexity: } O(n^2d) = O(3136^2 \times 128) = 1.26 \text{ GFLOPs per layer} \quad (17.26)$$

Window attention (Swin, $M = 7$):

$$\text{Windows: } \frac{56}{7} \times \frac{56}{7} = 64 \text{ windows} \quad (17.27)$$

$$\text{Complexity: } O(M^2 \cdot HW \cdot d) = O(49 \times 3136 \times 128) = 19.7 \text{ MFLOPs per layer} \quad (17.28)$$

The window-based approach reduces attention cost by $64\times$, making high-resolution processing practical. The shifted window mechanism ensures that information still propagates globally through the network depth.

Swin Transformer achieves state-of-the-art performance across multiple vision tasks while maintaining computational efficiency. On ImageNet classification, Swin-Base reaches 83.5% top-1 accuracy with 88 million parameters and 15.4 GFLOPs—comparable to ViT-Base in parameters but with better accuracy due to the hierarchical structure. For object detection on COCO, Swin-Base achieves 51.9 box AP, surpassing previous transformer-based detectors by significant margins. The multi-scale feature maps produced by the hierarchical architecture are particularly well-suited for dense prediction tasks, making Swin Transformer a versatile backbone for various computer vision applications.

17.5.3 Pyramid Vision Transformer (PVT)

Pyramid Vision Transformer takes a different approach to hierarchical vision transformers by introducing spatial-reduction attention that progressively decreases the key and value sequence lengths. Unlike Swin’s window-based attention, PVT maintains global attention but reduces computational cost by downsampling the keys and values before computing attention. This design preserves the ability to attend to the entire image while achieving sub-quadratic complexity.

In PVT, each stage reduces the spatial resolution through patch merging, similar to Swin Transformer. However, within each stage, the attention mechanism uses a spatial reduction operation on keys and values. For a reduction ratio R , the keys and values are reshaped and downsampled by $R \times R$, reducing their sequence length by a factor of R^2 . The queries maintain the original resolution, allowing each token to attend to a downsampled representation of the entire image. This approach reduces attention complexity from $O(n^2d)$ to $O(n^2d/R^2)$, providing a tunable trade-off between computational cost and attention granularity.

The hierarchical structure of PVT produces feature maps at multiple scales, making it suitable as a backbone for dense prediction tasks. PVT-Medium with 44 million parameters achieves 82.0% ImageNet accuracy while requiring only 6.7 GFLOPs—significantly more efficient than ViT-Base. For object detection, PVT-based detectors achieve competitive performance with CNN-based methods while offering the benefits of transformer architectures, including better transfer learning and attention-based interpretability.

17.5.4 Hybrid Architectures: CoAtNet

Hybrid architectures combine convolutional layers and transformer layers to leverage the complementary strengths of both approaches. Convolutional layers provide efficient local feature extraction with built-in translation equivariance, while transformer layers enable global reasoning and flexible attention patterns. CoAtNet (Convolution and Attention Network) systematically explores this design space, identifying an optimal combination that achieves state-of-the-art performance with improved efficiency.

The CoAtNet architecture consists of five stages with progressively decreasing spatial resolution. The first two stages use convolutional blocks based on the MBConv (Mobile Inverted Bottleneck Convolution) design from EfficientNet, which efficiently extracts local features at high resolution. These convolutional stages capture low-level visual patterns like edges, textures, and simple shapes with strong inductive bias and minimal computational cost. The spatial resolution is reduced by $2\times$ at each stage through strided convolutions.

The final three stages employ transformer blocks with relative attention, enabling global reasoning over the extracted features. By this point in the network, the spatial resolution has been reduced by $8\times$ or more, making global attention computationally feasible. The transformer stages learn high-level semantic representations and long-range dependencies that benefit from the flexibility of self-attention. The final stage uses attention pooling to aggregate spatial information into a global representation for classification.

Example 17.7 (CoAtNet Architecture). CoAtNet-3 configuration for 224×224 input:

Stage 0 (Stem): Convolution, 112×112 resolution, 64 channels

Stage 1: MBConv blocks, 112×112 resolution, 96 channels

Stage 2: MBConv blocks, 56×56 resolution, 192 channels

Stage 3: Transformer blocks, 28×28 resolution, 384 channels

Stage 4: Transformer blocks, 14×14 resolution, 768 channels

Stage 5: Attention pooling, global representation

Total parameters: 168M, FLOPs: 34.7G

This hybrid design achieves 87.9% ImageNet accuracy, outperforming pure CNN and pure transformer architectures of similar size.

The success of CoAtNet demonstrates that the choice between convolution and attention need not be binary. By using convolutions where they excel (local feature extraction at high resolution) and transformers where they excel (global reasoning at lower resolution), hybrid architectures achieve better accuracy-efficiency trade-offs than either approach alone. CoAtNet-7, the largest variant with 2.4 billion parameters, achieved 90.88% ImageNet accuracy and state-of-the-art results on multiple vision benchmarks at the time of its release, validating the hybrid approach at scale.

17.6 ViT vs CNN Comparison

17.6.1 Parameter Efficiency

Vision Transformers and Convolutional Neural Networks differ fundamentally in their parameter efficiency and data requirements. ResNet-50, a standard CNN baseline, contains approximately 25 million parameters distributed across convolutional layers with small kernel sizes (typically 3×3 or 7×7). In contrast, ViT-Base requires 86 million parameters—more than $3\times$ the size of ResNet-50—to achieve comparable performance. This parameter gap reflects the different inductive biases: CNNs build in locality and translation equivariance through their convolutional structure, while transformers must learn these properties from data through their flexible attention mechanism.

The parameter distribution also differs significantly between the architectures. In ResNet-50, the majority of parameters reside in the later convolutional layers and the final fully-connected layer. For ViT-Base, the parameters are more evenly distributed across the 12 transformer layers, with each layer containing approximately 7 million parameters in the attention and feed-forward components. The patch embedding layer contributes only 590K parameters, while position embeddings add another 151K—both negligible compared to the transformer layers themselves.

Despite having more parameters, ViT-Base is not necessarily slower than ResNet-50 for inference. The transformer's matrix multiplications are highly optimized on modern GPUs, and the lack of spatial convolutions can actually improve throughput. On an NVIDIA A100 GPU, ViT-Base processes approximately 1,200 images per second at 224×224 resolution with batch size 128, compared to 1,400 images per second for ResNet-50. The 15% throughput difference is much smaller than the $3\times$ parameter gap would suggest, demonstrating the efficiency of transformer operations on modern hardware.

17.6.2 Computational Complexity Analysis

The computational complexity of Vision Transformers scales differently than CNNs, leading to different performance characteristics across image resolutions. For a CNN like ResNet-50, the computational cost is approximately $O(C \times k^2 \times H \times W)$ where C is the number of channels, k is the kernel size, and $H \times W$ is the spatial resolution. This linear scaling in spatial dimensions means that doubling the image resolution increases computation by $4\times$. For ResNet-50 processing a 224×224 image, the total computation is approximately 4.1 GFLOPs.

Vision Transformers have complexity $O(n^2d + nd^2)$ where $n = (H/P)^2$ is the number of patches and d is the hidden dimension. The n^2d term comes from attention, while nd^2 comes from the feed-forward network. For ViT-Base with 224×224 images and patch size 16, we have $n = 196$ and $d = 768$. The

attention computation across 12 layers totals $12 \times 4 \times 196^2 \times 768 = 1.4$ GFLOPs, while the feed-forward network contributes $12 \times 2 \times 196 \times 768^2 = 2.8$ GFLOPs, for a total of approximately 4.2 GFLOPs—nearly identical to ResNet-50.

However, the scaling behavior differs dramatically. When we increase resolution to 384×384 with the same patch size, the number of patches grows to $n = 576$, increasing by a factor of $(384/224)^2 = 2.94$. The attention cost grows quadratically to $12 \times 4 \times 576^2 \times 768 = 12.3$ GFLOPs ($8.6\times$ increase), while the feed-forward cost grows linearly to $12 \times 2 \times 576 \times 768^2 = 8.1$ GFLOPs ($2.9\times$ increase). The total ViT computation reaches 20.4 GFLOPs, compared to 12.0 GFLOPs for ResNet-50 at the same resolution. This crossover point illustrates why efficient attention mechanisms become critical for high-resolution vision tasks.

Example 17.8 (Computational Crossover Analysis). Compare FLOPs for ResNet-50 and ViT-Base across resolutions:

Resolution	ResNet-50	ViT-Base
224×224	4.1 GFLOPs	4.2 GFLOPs
384×384	12.0 GFLOPs	20.4 GFLOPs
512×512	21.3 GFLOPs	48.7 GFLOPs

At standard ImageNet resolution, ViT and ResNet have similar computational cost. However, ViT's quadratic attention scaling makes it increasingly expensive at higher resolutions, motivating hierarchical architectures like Swin Transformer that reduce attention to local windows.

17.6.3 Data Requirements and Inductive Bias

The most striking difference between Vision Transformers and CNNs lies in their data requirements, which stem from their different inductive biases. CNNs encode strong priors about images: locality (nearby pixels are related), translation equivariance (a cat is a cat regardless of position), and hierarchical structure (edges \rightarrow textures \rightarrow objects). These built-in assumptions allow CNNs to learn effectively from moderate-sized datasets like ImageNet with 1.3 million images. ResNet-50 trained only on ImageNet achieves 76.5% top-1 accuracy, demonstrating that the convolutional structure provides useful inductive bias for natural images.

Vision Transformers, by contrast, have minimal inductive bias. The self-attention mechanism can attend to any patch regardless of spatial distance, and the model must learn locality and translation properties from data. When trained only on ImageNet, ViT-Base achieves only 72.3% accuracy—4.2 percentage points below ResNet-50 despite having $3\times$ more parameters. This performance gap reveals that the flexibility of attention becomes a liability when training data is limited: the model has too much capacity and insufficient constraints to learn good representations.

The situation reverses dramatically with large-scale pre-training. When ViT-Base is pre-trained on JFT-300M (300 million images with 18,000 classes) and then fine-tuned on ImageNet, it achieves 84.2% accuracy, surpassing ResNet-50's 76.5% by a substantial margin. The massive pre-training dataset provides enough examples for the transformer to learn the visual priors that CNNs encode by design. Moreover, the learned representations transfer better to downstream tasks: ViT-Base pre-trained on JFT-300M achieves higher accuracy than ResNet-50 on 19 out of 20 transfer learning benchmarks, with improvements ranging from 2-7 percentage points.

This data-efficiency trade-off has important practical implications. For applications with limited training data or computational budgets, CNNs remain the better choice. For large-scale systems with access to massive datasets and compute, Vision Transformers offer superior performance and transfer learning capabilities. The development of data-efficient training methods like DeiT (Data-efficient Image Transformers) has partially bridged this gap, enabling ViT-Base to achieve 81.8% on ImageNet without external data through aggressive augmentation and distillation techniques.

Aspect	CNN (ResNet)	ViT
Inductive bias	Strong (locality, translation)	Weak
Data requirement	Moderate (ImageNet)	Large (JFT-300M)
Parameters	25M (ResNet-50)	86M (ViT-Base)
Computation	$O(HW)$	$O((HW/P)^2)$
Memory	5-7 GB training	8-10 GB training
Interpretability	Filter visualization	Attention maps
Transfer	Good	Excellent (large-scale)
Best use	Small/medium data	Large-scale pre-training

17.6.4 When to Use Each Architecture

The choice between CNNs and Vision Transformers depends on the specific application constraints. CNNs are preferable when training data is limited (fewer than 10 million images), when computational efficiency is critical (mobile or edge deployment), or when strong spatial priors are known to be appropriate for the task. ResNet and EfficientNet variants remain the standard choice for many production computer vision systems due to their reliability and efficiency.

Vision Transformers excel when massive pre-training data is available, when transfer learning to diverse downstream tasks is important, or when state-of-the-art performance justifies the additional computational cost. The superior scaling properties of transformers—both in terms of model size and dataset size—make them the architecture of choice for foundation models in vision. Hybrid architectures like CoAtNet attempt to combine the strengths of both approaches, using convolutional layers for early feature extraction and transformer layers for high-level reasoning.

17.7 Exercises

Exercise 17.1. Implement patch embedding for image $224 \times 224 \times 3$ with patch size 16:

1. Reshape image to patches
2. Apply linear projection
3. Add position embeddings
4. Verify output shape: (196, 768)

Exercise 17.2. Compare ViT-Base and ResNet-50:

1. Parameter count
2. FLOPs for 224×224 image
3. Memory footprint
4. Which is more efficient?

Exercise 17.3. Implement MAE masking:

1. Randomly mask 75% of 196 patches
2. Keep 49 visible patches

3. Add mask tokens for decoder
4. Compute reconstruction loss

Exercise 17.4. Train ViT-Tiny on CIFAR-10:

1. Use patch size 4 (for 32×32 images)
2. 6 layers, $d = 192$, 3 heads
3. Apply RandAugment
4. Compare to small ResNet

Chapter 18

Multimodal Transformers

Chapter Overview

Multimodal transformers process multiple modalities (text, images, audio, video) in a unified framework. This chapter covers vision-language models (CLIP, DALL-E), audio-text models (Whisper), and unified architectures that handle arbitrary combinations of modalities.

Learning Objectives

1. Understand multimodal fusion strategies
2. Implement contrastive learning (CLIP)
3. Apply vision-language models to zero-shot classification
4. Generate images from text (DALL-E, Stable Diffusion)
5. Process audio with transformers (Whisper)
6. Build unified multimodal models

18.1 Multimodal Learning Fundamentals

18.1.1 Fusion Strategies

The fundamental challenge in multimodal learning is determining how to combine information from different modalities—such as vision, language, and audio—into a unified representation. The choice of fusion strategy has profound implications for model architecture, computational cost, and the types of cross-modal interactions the model can learn. Three primary approaches have emerged, each with distinct trade-offs in terms of expressiveness, efficiency, and implementation complexity.

Early fusion combines modalities at the input level, concatenating or interleaving embeddings from different sources before processing them through a single unified encoder. For example, an image might be encoded into a sequence of patch embeddings $\mathbf{v}_1, \dots, \mathbf{v}_N \in \mathbb{R}^d$, while text is tokenized into embeddings $\mathbf{t}_1, \dots, \mathbf{t}_M \in \mathbb{R}^d$. These are concatenated into a single sequence $[\mathbf{v}_1, \dots, \mathbf{v}_N, \mathbf{t}_1, \dots, \mathbf{t}_M]$ and processed by a standard transformer encoder. The advantage of this approach is its simplicity and the fact that cross-modal interactions can occur at every layer through self-attention. However, early fusion has significant computational drawbacks: the attention complexity is $O((N + M)^2 d)$, meaning that adding image patches dramatically increases the cost. For a ViT-Base image with 196 patches and text with 128 tokens, the combined sequence length of 324 tokens results in attention matrices of size 324×324 , consuming substantial memory and compute compared to processing each modality separately.

Late fusion takes the opposite approach, using separate encoders for each modality and combining their outputs only at the final decision stage. An image encoder produces a single embedding $\mathbf{v} \in \mathbb{R}^d$, a text encoder produces $\mathbf{t} \in \mathbb{R}^d$, and these are combined through simple operations like concatenation,

averaging, or element-wise multiplication before a final classification layer. This approach is computationally efficient since each encoder processes only its own modality with complexity $O(N^2d)$ for images and $O(M^2d)$ for text, which can be computed in parallel. The total complexity is $O(N^2d + M^2d)$ rather than $O((N + M)^2d)$. However, late fusion severely limits cross-modal interactions—the encoders cannot attend to information from other modalities, and the combination happens only at the very end. This makes it difficult to learn fine-grained alignments, such as which image regions correspond to which words in a caption.

Cross-modal attention provides a middle ground, using separate encoders but introducing explicit attention mechanisms that allow one modality to query information from another. In this architecture, an image encoder produces a sequence of patch embeddings $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_N] \in \mathbb{R}^{N \times d}$, and a text encoder produces token embeddings $\mathbf{T} = [\mathbf{t}_1, \dots, \mathbf{t}_M] \in \mathbb{R}^{M \times d}$. Additional cross-attention layers are then inserted where text tokens attend to image patches (queries from text, keys and values from image) and vice versa. The computational cost is $O(N^2d + M^2d + NMd)$ for the self-attention in each modality plus the cross-attention between them. The cross-attention term NMd is linear in both sequence lengths, making it much more efficient than early fusion’s quadratic term. For our example with 196 image patches and 128 text tokens, cross-attention requires $196 \times 128 = 25,088$ attention computations per head, compared to $324^2 = 104,976$ for early fusion—a $4\times$ reduction. This approach is used in models like BLIP and Flamingo, where it enables rich cross-modal interactions while maintaining computational tractability.

18.1.2 Alignment Objectives

Contrastive Learning:

$$\mathcal{L}_{\text{contrastive}} = -\log \frac{\exp(\text{sim}(v_i, t_i)/\tau)}{\sum_j \exp(\text{sim}(v_i, t_j)/\tau)} \quad (18.1)$$

where v_i = image embedding, t_i = text embedding, τ = temperature

Matching Loss:

$$\mathcal{L}_{\text{match}} = -\mathbb{E}[\log P(\text{match}|v, t)] \quad (18.2)$$

Reconstruction:

$$\mathcal{L}_{\text{recon}} = \|f_{\text{dec}}(v) - t\|^2 \quad (18.3)$$

18.2 CLIP: Contrastive Language-Image Pre-training

18.2.1 CLIP Architecture

CLIP (Contrastive Language-Image Pre-training) represents a breakthrough in vision-language learning by training image and text encoders jointly using a contrastive objective on 400 million image-text pairs collected from the internet. Unlike traditional supervised learning that requires manually labeled categories, CLIP learns to align images with their natural language descriptions, enabling zero-shot transfer to downstream tasks without any task-specific training data.

Definition 18.1 (CLIP Model). The CLIP architecture consists of three main components that work together to create a shared embedding space for images and text. The **image encoder** can be either a Vision Transformer (ViT) or a ResNet, which processes an input image and produces a fixed-dimensional embedding $\mathbf{v} \in \mathbb{R}^d$. For the largest CLIP model (ViT-L/14), the image encoder is a ViT with patch size 14, hidden dimension 1024, 24 layers, and 16 attention heads, totaling approximately 304 million parameters. The **text encoder** is a transformer decoder (similar to GPT) with a context length of 77 tokens, hidden dimension 768, 12 layers, and 12 attention heads, containing roughly 63 million parameters. Both encoders are followed by learned linear **projection layers** that map their outputs to a shared embedding space of dimension $d = 512$, where cosine similarity can be computed directly.

The training procedure processes batches of (*image*, *text*) pairs simultaneously. For each batch

of size N , all N images are encoded to produce embeddings $\mathbf{v}_1, \dots, \mathbf{v}_N \in \mathbb{R}^{512}$, and all N text descriptions are encoded to produce $\mathbf{t}_1, \dots, \mathbf{t}_N \in \mathbb{R}^{512}$. The model then computes an $N \times N$ similarity matrix where entry (i, j) represents the cosine similarity between image i and text j . The contrastive loss maximizes the similarity along the diagonal (correct image-text pairs) while minimizing off-diagonal similarities (incorrect pairings). This symmetric loss is computed in both directions—predicting text from image and image from text—and averaged.

The parameter count for CLIP varies significantly across model scales. CLIP ResNet-50 contains approximately 102 million parameters (38M for ResNet-50 image encoder, 63M for text encoder, 1M for projections), while CLIP ViT-L/14 totals around 428 million parameters (304M for ViT-L image encoder, 123M for a larger text encoder with 768 dimensions and 12 layers, 1M for projections). The largest variant, ViT-L/14@336px, processes higher-resolution images (336×336 instead of 224×224) with the same architecture, increasing computational cost but improving performance on fine-grained tasks.

Example 18.1 (CLIP Training). Consider a training batch with size $N = 4$ and embedding dimension $d = 512$. The image encoder processes four images to produce embeddings arranged as rows in matrix $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3, \mathbf{v}_4]^\top \in \mathbb{R}^{4 \times 512}$, while the text encoder processes their corresponding captions to produce $\mathbf{T} = [\mathbf{t}_1, \mathbf{t}_2, \mathbf{t}_3, \mathbf{t}_4]^\top \in \mathbb{R}^{4 \times 512}$.

The similarity matrix is computed as $\mathbf{S} = \mathbf{V}\mathbf{T}^\top \in \mathbb{R}^{4 \times 4}$, where each entry S_{ij} represents the dot product between image embedding i and text embedding j . To make this scale-invariant, CLIP uses cosine similarity: $S_{ij} = \frac{\mathbf{v}_i \cdot \mathbf{t}_j}{\|\mathbf{v}_i\| \|\mathbf{t}_j\|}$, which normalizes each embedding to unit length before computing the dot product. This ensures that the similarity is determined by the angle between embeddings rather than their magnitudes.

The contrastive loss for the image-to-text direction is computed as:

$$\mathcal{L}_i^{\text{img} \rightarrow \text{txt}} = -\log \frac{\exp(S_{ii}/\tau)}{\sum_{j=1}^N \exp(S_{ij}/\tau)} \quad (18.4)$$

where τ is a learned temperature parameter, initialized to 0.07 and trained jointly with the model. The temperature controls the sharpness of the distribution: smaller values make the model more confident (sharper peaks), while larger values produce softer distributions. The symmetric text-to-image loss $\mathcal{L}_i^{\text{txt} \rightarrow \text{img}}$ is computed analogously by treating text as queries and images as candidates. The total loss averages both directions:

$$\mathcal{L} = \frac{1}{2N} \sum_{i=1}^N (\mathcal{L}_i^{\text{img} \rightarrow \text{txt}} + \mathcal{L}_i^{\text{txt} \rightarrow \text{img}}) \quad (18.5)$$

In practice, CLIP uses very large batch sizes to provide more negative examples for contrastive learning. The original CLIP was trained with batch size 32,768, requiring distributed training across multiple GPUs. With such large batches, each positive pair has 32,767 negative examples, providing a strong learning signal. However, this creates substantial memory requirements: storing the $32,768 \times 512$ embedding matrices for images and text requires $32,768 \times 512 \times 4 = 67$ MB per modality in FP32, and the $32,768 \times 32,768$ similarity matrix requires 4.3 GB. To make this tractable, CLIP uses gradient checkpointing and distributes the batch across many GPUs, computing the similarity matrix in chunks.

18.2.2 Computational Analysis of CLIP Training

Training CLIP at scale requires careful consideration of computational and memory costs across both the image and text encoding paths. For the ViT-L/14 image encoder processing 224×224 images, each image is divided into $16 \times 16 = 256$ patches of size 14×14 . These patches are linearly projected to dimension 1024 and processed through 24 transformer layers. The computational cost per image is

approximately $2 \times 24 \times 256^2 \times 1024 = 3.2$ GFLOPS for the attention operations (using the $2Ld^2n^2$ formula from Chapter 12) plus $2 \times 24 \times 256 \times 4 \times 1024^2 = 51.5$ GFLOPS for the feed-forward networks, totaling roughly 55 GFLOPS per image.

The text encoder processes sequences of up to 77 tokens through 12 transformer layers with dimension 768. The computational cost per text is approximately $2 \times 12 \times 77^2 \times 768 = 1.1$ GFLOPS for attention plus $2 \times 12 \times 77 \times 4 \times 768^2 = 4.4$ GFLOPS for feed-forward networks, totaling about 5.5 GFLOPS per text. This asymmetry—images requiring 10× more compute than text—means that image encoding dominates the computational budget during training.

For a batch of 32,768 examples, the total forward pass requires approximately $32,768 \times (55 + 5.5) = 1,982,464$ GFLOPS or roughly 2 PFLOPS. On an NVIDIA A100 GPU with 312 TFLOPS of FP16 compute, this would take approximately 6.4 seconds per batch for the forward pass alone, not including backward propagation (which typically costs 2× the forward pass) or the contrastive loss computation. The full training of CLIP on 400 million image-text pairs with batch size 32,768 requires approximately $400,000,000/32,768 = 12,207$ batches. At roughly 20 seconds per batch (forward + backward + optimizer step), this amounts to 68 hours of continuous training on a single A100. In practice, OpenAI trained CLIP on 256 V100 GPUs for approximately 12 days, suggesting a total training cost of around 73,728 GPU-hours.

Memory requirements are equally demanding. Each image in the batch requires storing activations for 24 layers with 256 tokens and dimension 1024, totaling approximately $24 \times 256 \times 1024 \times 2 = 12.6$ MB per image in FP16 (the factor of 2 accounts for storing both pre- and post-activation values for backpropagation). For batch size 32,768, this amounts to 413 GB just for image activations. Text activations are smaller at approximately $12 \times 77 \times 768 \times 2 = 1.4$ MB per text, or 46 GB for the full batch. The similarity matrix requires $32,768 \times 32,768 \times 2 = 2.1$ GB in FP16. Combined with model parameters (428M parameters \times 2 bytes = 856 MB) and optimizer states (typically 2× parameters for Adam), the total memory footprint exceeds 500 GB, necessitating distribution across many GPUs using techniques like ZeRO (Chapter 22) to partition optimizer states and activations.

18.2.3 Zero-Shot Classification with CLIP

One of CLIP’s most remarkable capabilities is zero-shot classification: the ability to classify images into categories the model has never been explicitly trained on. This works by leveraging the natural language understanding of the text encoder to create classifiers on the fly from text descriptions. The procedure begins by creating text prompts for each class in the target classification task. For example, for a 10-class animal classification task, we might create prompts like “a photo of a dog”, “a photo of a cat”, “a photo of a bird”, and so on. These prompts are encoded by the text encoder to produce class embeddings $\mathbf{c}_1, \dots, \mathbf{c}_C \in \mathbb{R}^{512}$ where C is the number of classes.

To classify a new image, we encode it with the image encoder to produce $\mathbf{v} \in \mathbb{R}^{512}$, then compute the cosine similarity between the image embedding and each class embedding: $s_i = \frac{\mathbf{v} \cdot \mathbf{c}_i}{\|\mathbf{v}\| \|\mathbf{c}_i\|}$. The predicted class is simply $\arg \max_i s_i$, the class whose text description has the highest similarity to the image. This approach requires no training on the target dataset—the model uses only its pre-trained knowledge of how images and text relate.

The performance of this zero-shot approach is surprisingly strong. CLIP ViT-L/14 achieves 76.2% top-1 accuracy on ImageNet without ever seeing a single ImageNet training example, matching the performance of a ResNet-50 trained directly on ImageNet’s 1.28 million labeled images. This demonstrates that CLIP has learned visual concepts that generalize far beyond its training distribution. Moreover, CLIP shows remarkable robustness to distribution shift: when evaluated on ImageNet variants with different image styles (sketches, cartoons, adversarial examples), CLIP’s performance degrades much less than supervised models, suggesting it has learned more robust visual representations.

The prompt engineering aspect of zero-shot classification is crucial for performance. Simple prompts like “dog” perform worse than more descriptive prompts like “a photo of a dog”. OpenAI found that using prompt ensembles—averaging predictions across multiple prompt templates like “a photo of a {class}”, “a picture of a {class}”, “an image of a {class}”—improves accuracy by 1-2% by reducing sensitivity to prompt phrasing. For fine-grained classification tasks, more specific prompts help: “a photo of a {species}, a type of bird” outperforms “a photo of a {species}” for bird species classification.

18.2.4 CLIP Variants and Training Requirements

Following CLIP’s success, several variants have been developed with different scales and training procedures. **OpenCLIP** is an open-source reproduction that has trained models ranging from small (ResNet-50 with 102M parameters) to very large (ViT-G/14 with 1.8B parameters) on datasets including LAION-400M and LAION-2B. The largest OpenCLIP models require training on clusters of 128-512 A100 GPUs for several weeks, with estimated costs exceeding \$100,000 for the full training run. The training uses mixed precision (FP16) to reduce memory consumption and enable larger batch sizes, typically 32,768 to 65,536 examples distributed across all GPUs.

ALIGN, developed by Google, scales up the training data to 1.8 billion noisy image-text pairs collected from the web without extensive filtering. This demonstrates that contrastive learning is robust to noise in the training data—the model learns to ignore mismatched pairs through the contrastive objective. ALIGN uses an EfficientNet-L2 image encoder (480M parameters) and a BERT-Large text encoder (340M parameters), totaling approximately 820M parameters. Training ALIGN required a cluster of 1024 Cloud TPU v3 cores for approximately 6 days, representing roughly 150,000 TPU-hours.

Florence, Microsoft’s unified vision foundation model, extends the CLIP approach to 900 million image-text pairs with a focus on creating a single model that can be adapted to diverse vision tasks. Florence uses a CoSwin transformer as the image encoder (637M parameters) and achieves state-of-the-art results on zero-shot classification, retrieval, and object detection after fine-tuning. The training infrastructure required 512 NVIDIA A100 GPUs for approximately 10 days, with an estimated cost of over \$200,000 in cloud compute.

The hardware requirements for training CLIP-scale models are substantial. A minimum viable setup might use 8-16 A100 GPUs (80GB each) to train a CLIP ResNet-50 model on a smaller dataset like Conceptual Captions (3M pairs) with batch size 2048-4096, requiring approximately 1-2 weeks. Scaling to the full CLIP ViT-L/14 with 400M training pairs and batch size 32,768 necessitates at least 64-128 A100 GPUs with high-bandwidth interconnects (NVLink or InfiniBand) to efficiently synchronize gradients across the distributed batch. The total training cost for reproducing CLIP ViT-L/14 is estimated at \$50,000-\$100,000 in cloud GPU costs, depending on the provider and optimization techniques employed.

18.3 DALL-E and Stable Diffusion

18.3.1 DALL-E: Text-to-Image Generation

Definition 18.2 (DALL-E Architecture). **DALL-E 1 (2021):**

- Encoder: Compress images to discrete tokens (VQ-VAE)
- Transformer: Autoregressive model over text + image tokens
- Training: Next token prediction

Sequence:

$$[\text{BOS}, \text{text tokens}, \text{image tokens}, \text{EOS}] \quad (18.6)$$

Generate image by: (1) Encode text, (2) Sample image tokens autoregressively

DALL-E 2 (2022):

- Use CLIP embeddings
- Prior: Text embedding \rightarrow Image embedding
- Decoder: Image embedding \rightarrow Image (diffusion model)
- Much higher quality than DALL-E 1

18.3.2 Stable Diffusion

Latent Diffusion Model:

1. Encode image to latent space (VAE)
2. Add noise iteratively (forward diffusion)
3. Learn to denoise (reverse diffusion)
4. Condition on text via cross-attention

Text conditioning:

- Text encoder: CLIP or T5
- Cross-attention: Latent queries attend to text keys/values
- Enables text-guided image generation

Example 18.2 (Stable Diffusion Architecture). Components:

1. **Text Encoder:** CLIP text encoder

$$\text{prompt} \rightarrow \mathbb{R}^{77 \times 768} \quad (18.7)$$

2. **VAE Encoder:** Image \rightarrow latent

$$\mathbf{I} \in \mathbb{R}^{512 \times 512 \times 3} \rightarrow \mathbf{z} \in \mathbb{R}^{64 \times 64 \times 4} \quad (18.8)$$

3. **U-Net Denoiser:** Diffusion model with cross-attention

- Input: Noisy latent \mathbf{z}_t
- Condition: Text embedding
- Output: Predicted noise $\epsilon_\theta(\mathbf{z}_t, t,)$

4. **VAE Decoder:** Latent \rightarrow image

$$\mathbf{z} \in \mathbb{R}^{64 \times 64 \times 4} \rightarrow \mathbf{I} \in \mathbb{R}^{512 \times 512 \times 3} \quad (18.9)$$

Parameters: $\approx 860\text{M}$ total

18.4 Vision-Language Understanding

18.4.1 BLIP: Bootstrapped Language-Image Pre-training

Architecture:

- Image encoder (ViT)
- Text encoder (BERT)
- Multimodal encoder (cross-attention between vision and text)

Training objectives:

1. **ITC:** Image-Text Contrastive (like CLIP)
2. **ITM:** Image-Text Matching (binary: match or not)
3. **LM:** Language Modeling on text

Bootstrapping: Generate synthetic captions, filter with model, retrain

18.4.2 Flamingo: Few-Shot Learning

Flamingo represents a significant architectural innovation in multimodal transformers by enabling models to process arbitrarily interleaved sequences of images and text, supporting few-shot learning through in-context examples. Unlike CLIP, which processes single image-text pairs, Flamingo can handle inputs like "Here is an image of a cat: <image1>. Here is an image of a dog: <image2>. What animal is in this image: <image3>?" This capability enables few-shot visual learning where the model learns new tasks from just a few examples provided in the prompt, without any parameter updates.

The Flamingo architecture consists of three main components, carefully designed to leverage pre-trained models while adding minimal trainable parameters. The **vision encoder** is a frozen CLIP ViT-L/14 model that processes each image independently to produce a sequence of patch embeddings. For a 224×224 image with patch size 14, this yields 256 patch tokens of dimension 1024. The vision encoder's 304M parameters remain frozen throughout training, preserving the strong visual representations learned during CLIP pre-training.

The **language model** is a frozen Chinchilla 70B model, a large autoregressive transformer trained on text-only data. Chinchilla uses 70 billion parameters across 80 layers with hidden dimension 8192 and 64 attention heads. Keeping this massive language model frozen is crucial for computational tractability—training 70B parameters would require prohibitive memory and compute. Instead, Flamingo inserts new trainable layers that allow the frozen language model to attend to visual information without modifying its core text processing capabilities.

The key innovation is the **Perceiver Resampler**, a learned module that compresses the variable-length sequence of image patch embeddings into a fixed number of visual tokens that can be efficiently processed by the language model. The Perceiver Resampler uses cross-attention where a fixed set of learned queries $\mathbf{Q} \in \mathbb{R}^{64 \times 2048}$ (64 visual tokens, dimension 2048) attends to the image patch embeddings $\mathbf{K}, \mathbf{V} \in \mathbb{R}^{256 \times 1024}$ from the vision encoder. This produces a fixed-size representation regardless of input image resolution or the number of images in the sequence. The Perceiver Resampler contains approximately 1.4B parameters (6 layers of cross-attention and feed-forward networks with dimension 2048), making it the primary trainable component of Flamingo.

Between every few layers of the frozen language model, Flamingo inserts new **cross-attention layers** that allow text tokens to attend to the visual tokens produced by the Perceiver Resampler. Specifically, for Flamingo-80B (built on Chinchilla-70B), cross-attention layers are inserted after every 7th transformer layer, resulting in approximately 11 cross-attention insertions across the 80 layers. Each cross-attention layer adds roughly 134M parameters (for dimension 8192), totaling about 1.5B parameters for all insertions. Combined with the Perceiver Resampler, Flamingo adds approximately 2.9B trainable parameters to the 70B frozen base model, representing just 4% additional parameters while enabling full multimodal capabilities.

The memory requirements for Flamingo are dominated by the frozen language model. Storing 70B parameters in FP16 requires 140 GB, which exceeds the memory of any single GPU. Flamingo uses model parallelism to partition the language model across multiple GPUs—for example, distributing across 8 A100 GPUs (80GB each) places roughly 8.75B parameters per GPU, consuming about 17.5 GB for parameters. Activations for a sequence of 2048 tokens (including both text and visual tokens) across 80 layers with dimension 8192 require approximately $2048 \times 80 \times 8192 \times 2 = 2.6$ GB per example in FP16. With batch size 8, activations consume 21 GB per GPU, leaving sufficient memory for gradients of the trainable parameters ($2.9\text{B parameters} \times 2 \text{ bytes} \times 2 \text{ for gradients} = 11.6 \text{ GB}$) and optimizer states (23.2 GB for Adam).

Training Flamingo on a mixture of image-text pairs, interleaved image-text documents, and video-text pairs requires substantial computational resources. The training dataset consists of 2.3 billion image-text pairs (similar to CLIP), 43 million interleaved image-text web pages, and 27 million video clips. Training Flamingo-80B for 1 epoch through this data with batch size 256 distributed across 256 A100 GPUs takes approximately 15 days, representing roughly 92,000 GPU-hours. The estimated training cost exceeds \$300,000 in cloud compute. However, the key advantage is that only 2.9B parameters are trained while leveraging the capabilities of a 70B language model, making training far more efficient than training a 70B multimodal model from scratch.

For inference, Flamingo's few-shot learning capability means that users can provide 2-32 example

image-text pairs in the prompt to demonstrate a new task, and the model adapts its predictions based on these examples without any fine-tuning. This in-context learning works because the cross-attention mechanism allows the model to attend to the example images when processing the query image. The computational cost of inference scales linearly with the number of examples in the context: each additional image adds 256 patch tokens (after vision encoding) compressed to 64 visual tokens (after Perceiver Resampler), increasing the sequence length and thus the attention cost. For a prompt with 4 example images and 1 query image (5 images total), the visual tokens contribute $5 \times 64 = 320$ tokens to the sequence, which combined with text tokens (typically 500-1000) results in sequences of 800-1300 tokens. On a single A100 GPU, Flamingo-80B can process approximately 2-3 such sequences per second, limited primarily by the memory bandwidth required to load the 70B parameter model.

18.5 Computational Analysis of Multimodal Transformers

18.5.1 Image Encoding Cost

The computational cost of encoding images in multimodal transformers depends critically on the choice of vision encoder architecture and input resolution. For a Vision Transformer processing an image of resolution $H \times W$ with patch size P , the image is divided into $n = \frac{HW}{P^2}$ patches. Each patch is linearly projected to dimension d , and the resulting sequence of n tokens is processed through L transformer layers. The total computational cost is approximately $2Ld^2n^2$ FLOPs for attention operations plus $8Ld^2n$ FLOPs for feed-forward networks (using the standard $4\times$ expansion ratio).

For concrete analysis, consider three common configurations. A ViT-Base model with 224×224 images, patch size 16, dimension 768, and 12 layers processes $n = 196$ patches. The attention cost is $2 \times 12 \times 768^2 \times 196^2 = 5.5$ GFLOPs, and the feed-forward cost is $8 \times 12 \times 768^2 \times 196 = 11.0$ GFLOPs, totaling 16.5 GFLOPs per image. A ViT-Large model with patch size 14, dimension 1024, and 24 layers processes $n = 256$ patches, requiring approximately $2 \times 24 \times 1024^2 \times 256^2 = 3.2$ TFLOPs for attention and $8 \times 24 \times 1024^2 \times 256 = 51.5$ GFLOPs for feed-forward, totaling 3.25 TFLOPs per image—nearly $200\times$ more than ViT-Base. Increasing resolution to 336×336 (as in CLIP ViT-L/14@336px) increases the patch count to $n = 576$, quadrupling the attention cost to 14.4 TFLOPs and more than doubling the total cost to 14.5 TFLOPs per image.

The memory requirements for image encoding are dominated by storing activations for backpropagation during training. For ViT-Large processing 256 patches through 24 layers with dimension 1024, each layer stores attention outputs (256×1024), feed-forward outputs (256×4096 intermediate, 256×1024 output), and layer norm statistics. Accounting for all activations, each image requires approximately $24 \times 256 \times (1024 + 4096 + 1024) \times 2 = 75$ MB in FP16. For a batch of 256 images (typical for CLIP training when distributed across many GPUs), activations consume 19.2 GB per GPU, leaving limited memory for model parameters and optimizer states on a 40GB A100 GPU.

18.5.2 Text Encoding Cost

Text encoding is typically far less expensive than image encoding due to shorter sequence lengths. For a transformer processing text with sequence length m tokens through L layers with dimension d , the computational cost is $2Ld^2m^2$ for attention plus $8Ld^2m$ for feed-forward networks. The CLIP text encoder uses $m = 77$ tokens, $d = 768$, and $L = 12$ layers, resulting in attention cost $2 \times 12 \times 768^2 \times 77^2 = 1.1$ GFLOPs and feed-forward cost $8 \times 12 \times 768^2 \times 77 = 4.4$ GFLOPs, totaling 5.5 GFLOPs per text—approximately $10\times$ less than encoding a single image with ViT-Large.

However, for large language models used in architectures like Flamingo, text encoding becomes more expensive. Chinchilla-70B uses 80 layers with dimension 8192, processing sequences of up to 2048 tokens. The attention cost is $2 \times 80 \times 8192^2 \times 2048^2 = 1.1 \times 10^{17}$ FLOPs or 110 PFLOPs per sequence, and the feed-forward cost is $8 \times 80 \times 8192^2 \times 2048 = 8.8$ TFLOPs, totaling approximately 110 TFLOPs per sequence. This is $34,000\times$ more expensive than the CLIP text encoder and $34\times$ more expensive than ViT-Large image encoding, demonstrating how large language models dominate the computational budget in models like Flamingo.

Memory for text activations scales with sequence length and model size. For Chinchilla-70B processing 2048 tokens through 80 layers with dimension 8192, storing activations requires approximately $80 \times 2048 \times (8192 + 32768 + 8192) \times 2 = 13$ GB per sequence in FP16 (accounting for attention outputs, feed-forward intermediate activations with $4\times$ expansion, and final outputs). With batch size 8 distributed across 8 GPUs, each GPU stores activations for 1 sequence, consuming 13 GB of the 80 GB available on an A100.

18.5.3 Cross-Modal Attention Cost

Cross-modal attention, where one modality attends to another, introduces additional computational costs that scale with the product of sequence lengths from both modalities. For text tokens attending to image patches, the attention operation computes \mathbf{QK}^\top where $\mathbf{Q} \in \mathbb{R}^{m \times d}$ (text queries) and $\mathbf{K} \in \mathbb{R}^{n \times d}$ (image keys), producing an attention matrix of size $m \times n$. The computational cost is $2mnd$ FLOPs for the matrix multiplication, plus $2mnd$ for computing the weighted sum with values, totaling $4mnd$ FLOPs per cross-attention layer.

For a model like BLIP with $m = 128$ text tokens, $n = 196$ image patches, and $d = 768$, each cross-attention layer costs $4 \times 128 \times 196 \times 768 = 77$ MFLOPs. With 6 cross-attention layers in the multimodal encoder, the total cross-attention cost is 462 MFLOPs—negligible compared to the self-attention costs in the image and text encoders. However, the memory for storing the $m \times n$ attention matrices is $128 \times 196 \times 2 = 50$ KB per layer in FP16, or 300 KB for 6 layers, which is also minimal.

In Flamingo, the cross-attention cost is higher due to the large language model dimension. With $m = 2048$ text tokens, $n = 64$ visual tokens (after Perceiver Resampler compression), and $d = 8192$, each cross-attention layer costs $4 \times 2048 \times 64 \times 8192 = 4.3$ GFLOPs. With approximately 11 cross-attention layers inserted throughout the 80-layer language model, the total cross-attention cost is 47 GFLOPs per forward pass. While this is small compared to the 110 TFLOPs for the language model self-attention, it represents a non-trivial addition. The memory for attention matrices is $2048 \times 64 \times 2 = 262$ KB per layer, or 2.9 MB for all 11 layers—again minimal compared to activation memory.

18.5.4 Total Memory Requirements

The total memory footprint for training multimodal transformers includes model parameters, activations, gradients, and optimizer states. For CLIP ViT-L/14 with 428M parameters, the parameters require $428 \times 10^6 \times 2 = 856$ MB in FP16. Gradients require an additional 856 MB. The Adam optimizer maintains two states (momentum and variance) per parameter, requiring $428 \times 10^6 \times 4 = 1.7$ GB in FP32 (optimizer states are typically kept in higher precision for numerical stability). Activations for a batch of 256 images and 256 texts consume approximately $256 \times (75 + 1.4) = 19.6$ GB as calculated earlier. The total memory per GPU is approximately $0.9 + 0.9 + 1.7 + 19.6 = 23.1$ GB, fitting comfortably on a 40GB A100 GPU.

For Flamingo-80B, the memory requirements are far more demanding. The 70B frozen parameters require $70 \times 10^9 \times 2 = 140$ GB in FP16, necessitating distribution across at least 2 A100 GPUs (80GB each). The 2.9B trainable parameters require $2.9 \times 10^9 \times 2 = 5.8$ GB for parameters, 5.8 GB for gradients, and $2.9 \times 10^9 \times 4 = 11.6$ GB for optimizer states, totaling 23.2 GB. Activations for a sequence of 2048 tokens through 80 layers with dimension 8192 require approximately 13 GB per sequence as calculated earlier. For batch size 8 distributed across 8 GPUs (1 sequence per GPU), the total memory per GPU is approximately $140/8 + 23.2 + 13 = 53.7$ GB, fitting on an 80GB A100 but leaving limited headroom. In practice, Flamingo uses gradient checkpointing to reduce activation memory by recomputing activations during backpropagation, trading compute for memory and enabling larger batch sizes or longer sequences.

18.6 Training Challenges for Multimodal Transformers

18.6.1 Batch Size Requirements for Contrastive Learning

Contrastive learning methods like CLIP require very large batch sizes to provide sufficient negative examples for effective learning. The contrastive loss compares each positive image-text pair against

all other pairs in the batch, treating them as negatives. With batch size N , each example has $N - 1$ negative examples, and the quality of the learned representations improves significantly as N increases. Empirically, CLIP’s performance scales log-linearly with batch size: increasing from 256 to 32,768 improves ImageNet zero-shot accuracy from approximately 58% to 76%, a gain of 18 percentage points.

However, large batch sizes create severe memory constraints. For batch size 32,768 with CLIP ViT-L/14, storing the image embeddings requires $32,768 \times 512 \times 4 = 67$ MB in FP32 (embeddings are kept in FP32 for numerical stability in the contrastive loss computation). Text embeddings require another 67 MB. The $32,768 \times 32,768$ similarity matrix requires $32,768^2 \times 4 = 4.3$ GB in FP32. During backpropagation, gradients of the loss with respect to the similarity matrix also require 4.3 GB. The total memory for just the contrastive loss computation is approximately 8.7 GB, not including model parameters, activations, or optimizer states.

To make such large batches tractable, CLIP uses distributed training across many GPUs with careful memory management. Each GPU processes a local batch of size $N_{\text{local}} = N/G$ where G is the number of GPUs. For $N = 32,768$ distributed across 256 GPUs, each GPU processes 128 examples. The key insight is that the similarity matrix can be computed in a distributed fashion: each GPU computes similarities between its local images and all texts (gathered from all GPUs via all-gather communication), producing a $128 \times 32,768$ matrix that requires only $128 \times 32,768 \times 4 = 16.8$ MB per GPU. The full $32,768 \times 32,768$ matrix is never materialized on any single GPU, reducing memory requirements by a factor of G .

The communication cost for distributed contrastive learning is substantial. Each GPU must perform an all-gather operation to collect embeddings from all other GPUs, communicating $32,768 \times 512 \times 4 = 67$ MB for images and 67 MB for text, totaling 134 MB per GPU per training step. With high-bandwidth interconnects like NVLink (300 GB/s) or InfiniBand (200 Gb/s = 25 GB/s), this communication takes approximately 5-10 milliseconds. For comparison, the forward and backward pass through the model takes approximately 50-100 milliseconds per step, so communication overhead is 5-10% of the total training time—significant but not prohibitive.

18.6.2 Memory Optimization Techniques

Beyond distributed training, several memory optimization techniques are essential for training large multimodal models. **Gradient checkpointing** trades computation for memory by not storing all activations during the forward pass. Instead, only a subset of activations (typically at layer boundaries) are stored, and intermediate activations are recomputed during backpropagation when needed. For CLIP ViT-L/14, gradient checkpointing reduces activation memory from approximately 75 MB per image to 25 MB per image (a $3\times$ reduction) at the cost of increasing training time by approximately 20% due to recomputation. This enables increasing batch size from 128 to 384 per GPU, which more than compensates for the slowdown by improving convergence.

Mixed precision training uses FP16 for most computations while maintaining FP32 master weights and loss scaling to prevent gradient underflow. For multimodal transformers, this reduces activation memory by 50% and accelerates matrix multiplications on modern GPUs with Tensor Cores. CLIP training with mixed precision reduces memory from 23.1 GB to approximately 14.5 GB per GPU, enabling batch sizes of 256 per GPU instead of 128. The training speedup is approximately $1.8\times$ on A100 GPUs, reducing total training time from 12 days to 7 days on the same hardware.

ZeRO (Zero Redundancy Optimizer) partitions optimizer states, gradients, and optionally parameters across GPUs to reduce per-GPU memory consumption. For Flamingo-80B with 2.9B trainable parameters, the optimizer states alone require 11.6 GB per GPU without ZeRO. With ZeRO Stage 2 (partitioning optimizer states and gradients), this is reduced to $11.6/G$ GB per GPU where G is the number of GPUs. With 8 GPUs, optimizer memory drops from 11.6 GB to 1.45 GB per GPU, freeing memory for larger batch sizes or longer sequences. ZeRO Stage 3 additionally partitions parameters, reducing the 5.8 GB parameter memory to 0.73 GB per GPU, though this increases communication overhead as parameters must be gathered during forward and backward passes.

18.6.3 Distributed Training Infrastructure

Training multimodal transformers at scale requires sophisticated distributed training infrastructure with high-bandwidth interconnects and efficient parallelism strategies. For CLIP-scale models (400M parameters, 400M training examples), a typical setup uses 64-256 A100 GPUs connected via NVLink within nodes and InfiniBand between nodes. Each node contains 8 GPUs with NVLink providing 600 GB/s bidirectional bandwidth between GPUs in the same node, while InfiniBand provides 200 Gb/s (25 GB/s) between nodes. This hierarchical network topology means that communication within a node is $24\times$ faster than between nodes, making it crucial to minimize cross-node communication.

The parallelism strategy combines data parallelism (different GPUs process different data) with model parallelism (different GPUs hold different parts of the model). For CLIP, pure data parallelism is typically sufficient since the model fits on a single GPU. Each GPU holds a full copy of the 428M parameter model (856 MB in FP16) and processes a local batch of images and texts. Gradients are synchronized across all GPUs using all-reduce operations after each backward pass, communicating $428 \times 10^6 \times 2 = 856$ MB per GPU. With 256 GPUs and ring all-reduce, this takes approximately $856 \text{ MB} \times 2 \times (256 - 1) / 256 / 25 \text{ GB/s} = 68 \text{ ms}$, adding roughly 10-15% overhead to the training step time.

For Flamingo-80B, model parallelism is essential since the 70B parameter language model cannot fit on a single GPU. The model is partitioned across 8 GPUs using tensor parallelism, where each linear layer’s weight matrix is split across GPUs. For a linear layer with weight $\mathbf{W} \in \mathbb{R}^{8192 \times 8192}$, each GPU holds $\mathbf{W}_i \in \mathbb{R}^{8192 \times 1024}$ (1/8th of the columns). During the forward pass, each GPU computes a partial output, then an all-reduce operation sums the results. This requires communicating $8192 \times 2 = 16$ KB per linear layer, which is minimal. However, with hundreds of linear layers across 80 transformer layers, the cumulative communication overhead is approximately 5-10% of training time.

18.6.4 Training Time and Cost Estimates

The total training time for multimodal transformers depends on the number of training examples, batch size, model size, and hardware configuration. For CLIP ViT-L/14 trained on 400M image-text pairs with batch size 32,768, the number of training steps is $400,000,000 / 32,768 = 12,207$ steps. Each step requires a forward pass (approximately 2 PFLOPs as calculated earlier), backward pass (approximately 4 PFLOPs, twice the forward pass), and optimizer step (negligible compute). The total compute is approximately $12,207 \times 6 \text{ PFLOPs} = 73,242 \text{ PFLOPs}$.

On an NVIDIA A100 GPU with 312 TFLOPS of FP16 compute and approximately 50% utilization (accounting for memory bandwidth limitations and non-compute operations), the effective throughput is 156 TFLOPS. Training on a single A100 would take $73,242 \text{ PFLOPs} / 156 \text{ TFLOPS} = 469,500$ seconds or approximately 130 hours. However, this ignores memory constraints—a single A100 cannot hold the activations for batch size 32,768. With 256 A100 GPUs, the training time is approximately $130 / 256 = 0.5$ hours, but this assumes perfect scaling, which is unrealistic due to communication overhead. In practice, with 10-15% communication overhead, the training time is approximately 0.6 hours per epoch. For multiple epochs (CLIP was trained for approximately 32 epochs), the total training time is approximately 19 hours on 256 A100 GPUs, though the actual training took several days due to additional factors like data loading, checkpointing, and debugging.

The cost of training CLIP at cloud pricing can be estimated from GPU-hours. Training for 19 hours on 256 A100 GPUs requires $19 \times 256 = 4,864$ GPU-hours. At typical cloud pricing of \$2.50-\$4.00 per A100 GPU-hour, the total cost is approximately \$12,000-\$19,000 for a single training run. However, the actual development cost is much higher due to hyperparameter tuning, ablation studies, and failed experiments. OpenAI likely spent 10-20 \times this amount in total compute to develop CLIP, suggesting a total cost of \$120,000-\$380,000.

For Flamingo-80B, the training cost is substantially higher due to the 70B parameter language model. Training on 2.3B image-text pairs plus interleaved documents with batch size 256 distributed across 256 A100 GPUs for approximately 15 days requires $15 \times 24 \times 256 = 92,160$ GPU-hours. At \$3.00 per A100 GPU-hour, the total cost is approximately \$276,000 for a single training run. Including development costs, the total investment in training Flamingo likely exceeded \$1 million in compute costs alone, not including researcher time, data collection, and infrastructure development.

18.7 Audio Transformers

18.7.1 Whisper: Speech Recognition

Definition 18.3 (Whisper Architecture). Encoder-decoder transformer for speech:

Input: Audio waveform \rightarrow Log-mel spectrogram

Encoder:

- Input: Spectrogram (80 mel bins)
- Convolution layers (downsample)
- Transformer encoder layers

Decoder:

- Autoregressive text generation
- Special tokens for language, task, timestamps

Training data: 680,000 hours of multilingual audio

Tasks supported:

- Speech recognition (transcription)
- Translation (to English)
- Language identification
- Voice activity detection
- Timestamp prediction

Example 18.3 (Whisper Input Format). **Special tokens:**

`<|startoftranscript|><|en|><|transcribe|><|notimestamps|>`

Spectrogram:

- 80 mel bins
- 3000 frames (30 seconds audio at 100 Hz)
- Input: 3000×80

Encoder:

- Conv layers: $3000 \times 80 \rightarrow 1500 \times 768$
- Transformer: Process 1500 tokens

Decoder: Generate text tokens autoregressively

18.7.2 Audio-Text Pre-training

Contrastive learning: Like CLIP but audio-text

AudioCLIP: Tri-modal (image, text, audio)

Applications:

- Zero-shot audio classification
- Audio captioning
- Text-to-audio generation

18.8 Unified Multimodal Models

18.8.1 Perceiver and Perceiver IO

Key idea: Map arbitrary modalities to latent space via cross-attention

Definition 18.4 (Perceiver). Components:

1. **Latent array:** Fixed set of learned queries $\mathbf{Z} \in \mathbb{R}^{M \times d}$
2. **Cross-attention:** Latents attend to inputs

$$\mathbf{Z}_1 = \text{CrossAttn}(\mathbf{Q} = \mathbf{Z}, \mathbf{K} = \mathbf{X}, \mathbf{V} = \mathbf{X}) \quad (18.10)$$

3. **Transformer:** Process latents

$$\mathbf{Z}_L = \text{Transformer}(\mathbf{Z}_1) \quad (18.11)$$

4. **Output:** Decode latents to task outputs

Benefits:

- Handles arbitrary input sizes
- Computation independent of input size (fixed latents)
- Unified architecture for images, video, audio, text

18.8.2 GPT-4V and LLaVA

GPT-4V (Vision): GPT-4 with vision capabilities

- Interleaved image and text inputs
- Strong vision-language understanding
- Details not fully disclosed

LLaVA (Open-source):

- CLIP vision encoder
- LLaMA language model
- Linear projection to align embeddings
- Instruction tuning on visual conversations

18.9 Exercises

Exercise 18.1. Implement CLIP contrastive loss for batch size 8:

1. Generate random image embeddings (8, 512)
2. Generate random text embeddings (8, 512)

3. Compute 8×8 similarity matrix
4. Calculate contrastive loss with $\tau = 0.07$

Exercise 18.2. Use CLIP for zero-shot classification on CIFAR-10:

1. Load pre-trained CLIP model
2. Create text prompts for 10 classes
3. Encode images and prompts
4. Compute accuracy
5. Compare to supervised baseline

Exercise 18.3. Analyze Whisper architecture:

1. Calculate parameters for encoder (24 layers, $d = 1024$)
2. Calculate parameters for decoder (24 layers)
3. Estimate memory for 30-second audio
4. Compare to text-only GPT-2

Exercise 18.4. Design multimodal fusion strategy for video understanding (visual + audio + captions):

1. Propose architecture
2. Define fusion mechanism
3. Specify training objective
4. Estimate parameter count

Chapter 19

Long Context Transformers

Chapter Overview

Extending transformer context length beyond standard limits (512-2048 tokens) enables processing long documents, books, and extended conversations. This chapter covers techniques for scaling to 32K, 100K, and even 1M+ token contexts.

Learning Objectives

1. Understand context length limitations and bottlenecks
2. Implement position interpolation and extrapolation
3. Apply memory-augmented transformers
4. Use retrieval-augmented generation (RAG)
5. Implement recurrent transformers (Transformer-XL)
6. Compare long-context methods and trade-offs

19.1 Context Length Limitations

19.1.1 The Quadratic Memory Bottleneck

Standard transformer architectures face fundamental limitations when processing long sequences due to the quadratic scaling of self-attention with respect to sequence length. The self-attention mechanism computes pairwise interactions between all tokens in a sequence, requiring the materialization of an attention matrix of size $n \times n$ where n is the sequence length. This quadratic scaling manifests in three critical bottlenecks: computational complexity, memory consumption, and position encoding limitations. Understanding these bottlenecks quantitatively is essential for appreciating why long context processing requires specialized techniques and architectural modifications.

The computational complexity of self-attention is $O(n^2d)$ where d is the model dimension. For each of the n queries, the model computes attention scores with all n keys through dot products of dimension d , requiring n^2d multiply-accumulate operations. The subsequent softmax normalization adds $O(n^2)$ operations, and the weighted sum over values requires another n^2d operations. While the feed-forward network has complexity $O(nd^2)$, for long sequences where $n > d$, the attention computation dominates. For example, with $n = 16384$ tokens and $d = 768$, attention requires approximately 206 billion FLOPs per layer while the feed-forward network requires only 19 billion FLOPs, making attention the primary computational bottleneck.

The memory bottleneck is even more severe than the computational one. During the forward pass, the attention matrix must be fully materialized in memory before applying softmax, requiring n^2 memory locations per attention head. During the backward pass for training, these attention matrices must be stored for gradient computation, effectively doubling the memory requirement. For multi-head

attention with h heads, the total memory for attention matrices is $h \times n^2$ floating-point values per layer. In FP32 format, each value requires 4 bytes, while FP16 requires 2 bytes. This memory grows quadratically with sequence length, quickly exceeding available GPU memory for long sequences.

Example 19.1 (Memory Scaling Analysis). Consider a GPT-2 scale model with $d = 768$, $h = 12$ attention heads, and $L = 12$ layers. The memory required for attention matrices scales dramatically with sequence length. For a single attention head processing a sequence of length n , the attention matrix requires $n^2 \times 4$ bytes in FP32 format. With 12 heads per layer, this becomes $12 \times n^2 \times 4 = 48n^2$ bytes per layer.

At $n = 1024$ tokens (GPT-2's standard context), each layer requires $48 \times 1024^2 = 50.3$ MB for attention matrices. Across 12 layers, this totals 604 MB, which is manageable on modern GPUs. However, doubling the sequence length to $n = 2048$ quadruples the memory requirement to 201 MB per layer or 2.4 GB total—a $4\times$ increase for only a $2\times$ increase in sequence length. This quadratic scaling continues: at $n = 4096$, attention matrices consume 805 MB per layer or 9.7 GB total, nearly filling a 16 GB GPU. At $n = 8192$, the requirement explodes to 3.2 GB per layer or 38.5 GB total, exceeding even high-end GPUs like the NVIDIA A100 with 40 GB memory.

For larger models, the situation becomes even more challenging. Consider a GPT-3 scale model with $d = 12288$, $h = 96$ heads, and $L = 96$ layers. At $n = 2048$ tokens, each layer requires $96 \times 2048^2 \times 4 = 1.6$ GB for attention matrices, totaling 154 GB across all layers. At $n = 8192$ tokens, each layer requires 25.8 GB, totaling 2.5 TB across the model—far exceeding any single GPU's capacity and requiring extensive model parallelism even for moderate context lengths. At $n = 32768$ tokens, a single layer would require 412 GB just for attention matrices, making standard attention completely impractical without fundamental algorithmic changes.

These calculations assume only the forward pass attention matrices. During training, gradients with respect to attention matrices must also be stored, effectively doubling the memory requirement. Additionally, activations from other layers, model parameters, optimizer states, and batch processing multiply these requirements further. For a batch size of 8 with $n = 4096$ tokens on GPT-2, attention matrices alone would require $9.7 \times 8 = 77.6$ GB, making training impossible on standard hardware without techniques like gradient checkpointing, which trades computation for memory by recomputing activations during the backward pass.

The third fundamental limitation involves position encodings. Standard transformers use position encodings trained on sequences of a fixed maximum length, typically 512 to 2048 tokens. When these models encounter sequences longer than their training length, the position encodings must extrapolate to unseen positions. Absolute position embeddings, which assign a learned vector to each position index, cannot extrapolate at all—positions beyond the training length have no corresponding embedding. Even sinusoidal position encodings, which use deterministic trigonometric functions, exhibit degraded performance when extrapolating beyond training lengths due to the model's learned attention patterns being calibrated to the training distribution of position encodings.

This extrapolation failure manifests as rapidly degrading perplexity for tokens beyond the training context length. A model trained on 2048-token sequences might achieve perplexity of 15 on positions 0-2048, but perplexity can increase to 25 or higher for positions 2048-4096 without specialized position encoding schemes. This degradation occurs because the model's attention patterns have learned to interpret specific position encoding values as corresponding to specific relative distances, and these learned patterns break down when position encodings take on values outside the training distribution.

19.2 Position Encoding for Long Context

19.2.1 The Extrapolation Challenge

Position encodings enable transformers to incorporate sequential order information into their otherwise permutation-invariant architecture. However, different position encoding schemes exhibit dramatically

different behaviors when processing sequences longer than those seen during training. This extrapolation capability is critical for long context applications, where retraining on maximum-length sequences is often computationally prohibitive. The choice of position encoding scheme can determine whether a model trained on 2048-token sequences can successfully process 8192-token sequences with minimal fine-tuning, or whether it requires extensive retraining from scratch.

Absolute position embeddings assign a learned vector to each position index, with the position encoding for position i being a trainable parameter $i \in \mathbb{R}^d$. These embeddings are added to token embeddings before the first transformer layer. While simple and effective within the training length, absolute embeddings cannot extrapolate beyond the maximum training position. A model trained with positions 0 through 2047 has no learned embedding for position 2048 or beyond. Attempting to extend such a model requires either initializing new position embeddings (which perform poorly without extensive fine-tuning) or using position interpolation techniques to map longer sequences into the trained position range.

Sinusoidal position encodings, introduced in the original Transformer paper, use deterministic trigonometric functions rather than learned parameters. For position i and dimension j , the encoding is defined as:

$$\text{PE}(i, 2j) = \sin(i/10000^{2j/d}), \quad \text{PE}(i, 2j + 1) = \cos(i/10000^{2j/d}) \quad (19.1)$$

These encodings can be computed for any position without training, enabling extrapolation in principle. However, in practice, models trained with sinusoidal encodings still exhibit degraded performance on longer sequences because the attention patterns learned during training are calibrated to the distribution of position encodings seen during training. When positions extend beyond the training range, the attention patterns encounter position encoding values in unfamiliar ranges, leading to suboptimal attention distributions.

19.2.2 Position Interpolation

Position interpolation addresses the extrapolation problem by mapping longer sequences into the position range seen during training, rather than extending beyond it. Instead of asking the model to extrapolate to unseen position indices, interpolation compresses the position indices of a long sequence into the trained range, effectively treating the long sequence as a "compressed" version of a training-length sequence.

Definition 19.1 (Position Interpolation). To extend a model trained on maximum length L to process sequences of length $L' > L$, position interpolation maps each position $i \in \{0, 1, \dots, L' - 1\}$ to a fractional position in the training range:

$$i_{\text{interpolated}} = i \cdot \frac{L}{L'} \quad (19.2)$$

For absolute position embeddings, the new position encoding is computed by interpolating between the learned embeddings:

$$\text{PE}_{\text{new}}(i) = \text{interpolate}(\text{PE}_{\text{original}}, i \cdot L/L') \quad (19.3)$$

For sinusoidal or rotary encodings, the interpolated position is used directly in the encoding formula, effectively reducing the frequency of the trigonometric functions.

The key insight behind position interpolation is that it keeps position encodings within the distribution seen during training, avoiding the extrapolation problem entirely. For example, extending from $L = 2048$ to $L' = 8192$ maps position 8191 to interpolated position $8191 \times 2048/8192 = 2047.75$, which falls within the training range. The model's attention patterns, having been trained on positions 0 through 2047, can handle this interpolated position much more effectively than they could handle the raw position 8191.

Position interpolation has been successfully applied to extend LLaMA models from 2048 to 8192 tokens and beyond. The technique requires minimal fine-tuning—typically only 1000 to 10000 training steps on long sequences—compared to training from scratch. After fine-tuning with position interpolation, LLaMA 2 models maintain perplexity within 5-10% of their original performance when extended from 4096 to 32768 tokens, whereas naive extrapolation without interpolation results in perplexity degradation of 50% or more.

The computational cost of position interpolation is negligible, as it only affects the position encoding computation, not the attention mechanism itself. The primary cost is the fine-tuning required to adapt the model to the compressed position space. However, this fine-tuning is far less expensive than training from scratch: extending a 7B parameter model from 4K to 32K context requires approximately 100 GPU-hours of fine-tuning compared to 100,000+ GPU-hours for full pretraining.

19.2.3 Rotary Position Embedding (RoPE)

Rotary Position Embedding represents a fundamental advance in position encoding design, achieving excellent extrapolation properties by encoding relative position information directly into the attention computation through rotation operations. RoPE has become the position encoding of choice for modern large language models including GPT-NeoX, LLaMA, PaLM, and many others due to its combination of strong extrapolation, computational efficiency, and theoretical elegance.

Definition 19.2 (Rotary Position Embedding). RoPE applies position-dependent rotations to query and key vectors before computing attention. For a query vector \mathbf{q}_m at position m and key vector \mathbf{k}_n at position n , RoPE applies rotation matrices:

$$\mathbf{q}'_m = \mathbf{R}_m \mathbf{q}_m \quad (19.4)$$

$$\mathbf{k}'_n = \mathbf{R}_n \mathbf{k}_n \quad (19.5)$$

where $\mathbf{R}_m \in \mathbb{R}^{d \times d}$ is a block-diagonal rotation matrix. For dimension pairs $(2j, 2j + 1)$, the rotation is:

$$\mathbf{R}_m^{(j)} = \begin{bmatrix} \cos(m\theta_j) & -\sin(m\theta_j) \\ \sin(m\theta_j) & \cos(m\theta_j) \end{bmatrix}, \quad \theta_j = 10000^{-2j/d} \quad (19.6)$$

The full rotation matrix is block-diagonal with $d/2$ such 2D rotation blocks.

The crucial property of RoPE is that the attention score between positions m and n depends only on their relative distance $m - n$, not their absolute positions. This can be verified through the rotation addition formula:

$$(\mathbf{q}'_m)^\top \mathbf{k}'_n = (\mathbf{R}_m \mathbf{q}_m)^\top (\mathbf{R}_n \mathbf{k}_n) = \mathbf{q}_m^\top \mathbf{R}_m^\top \mathbf{R}_n \mathbf{k}_n = \mathbf{q}_m^\top \mathbf{R}_{n-m} \mathbf{k}_n \quad (19.7)$$

This relative position property means that the model learns attention patterns based on relative distances between tokens rather than absolute positions. When extrapolating to longer sequences, the model encounters the same relative distances it saw during training, just in different combinations. A model trained on sequences up to 2048 tokens has seen all relative distances from -2047 to +2047. When processing a 4096-token sequence, it encounters the same relative distances, enabling much better extrapolation than absolute position encodings.

RoPE's extrapolation capability can be further enhanced through position interpolation. By scaling the rotation frequencies θ_j by a factor L'/L when extending from length L to L' , the effective relative distances are compressed into the training range. This combination of RoPE's inherent relative position encoding with position interpolation enables extensions from 2048 to 32768 tokens or beyond with minimal quality degradation.

The computational overhead of RoPE is minimal compared to the attention computation itself. Applying rotations to queries and keys requires $O(nd)$ operations, which is negligible compared to the $O(n^2d)$ cost of attention. The rotation operations can be efficiently implemented using vectorized

operations on modern GPUs, adding less than 5% to the total attention computation time. Memory overhead is also minimal, as the rotation matrices are block-diagonal and can be computed on-the-fly rather than stored.

In practice, RoPE has enabled dramatic context length extensions. LLaMA models using RoPE have been successfully extended from 2048 to 32768 tokens through position interpolation with only 1000 fine-tuning steps. The perplexity degradation is typically less than 10% even at $16\times$ the training length, compared to 50-100% degradation for absolute position embeddings. This extrapolation capability has made RoPE the de facto standard for new large language models designed for long context applications.

19.2.4 ALiBi: Attention with Linear Biases

ALiBi (Attention with Linear Biases) takes a radically different approach to position encoding by eliminating position embeddings entirely and instead adding position-dependent biases directly to attention scores. This simple modification achieves remarkable extrapolation properties, enabling models trained on 1024-token sequences to process 10,000+ token sequences at inference time with no fine-tuning whatsoever.

Definition 19.3 (ALiBi). ALiBi adds a bias term to attention scores based on the distance between query and key positions:

$$\text{score}(q_i, k_j) = \frac{\mathbf{q}_i^\top \mathbf{k}_j}{\sqrt{d_k}} - m \cdot |i - j| \quad (19.8)$$

where $m > 0$ is a head-specific slope parameter that differs across attention heads. The bias $-m \cdot |i - j|$ penalizes attention to distant tokens, with the penalty increasing linearly with distance.

The head-specific slopes are set geometrically: for h attention heads, the slopes are m_1, m_2, \dots, m_h where $m_i = 2^{-8i/h}$. For example, with 8 heads, the slopes are $2^{-1}, 2^{-2}, \dots, 2^{-8}$, giving values from 0.5 to 0.0039. This geometric spacing ensures that different heads have different receptive field sizes: heads with large slopes focus on nearby tokens, while heads with small slopes can attend to distant tokens with less penalty.

ALiBi's extrapolation capability stems from its use of relative distances rather than absolute positions, combined with the linear form of the bias. During training on sequences up to length L , the model encounters biases ranging from 0 (attending to the same position) to $-m \cdot L$ (attending to the most distant position). When extrapolating to length $L' > L$, the model encounters biases up to $-m \cdot L'$, which are simply larger values along the same linear scale. The attention patterns learned during training—which balance content-based attention (from $\mathbf{q}_i^\top \mathbf{k}_j$) against distance-based penalties (from $-m \cdot |i - j|$)—continue to work effectively at these larger distances.

Empirical results demonstrate ALiBi's exceptional extrapolation. Models trained on 1024-token sequences with ALiBi can process 2048-token sequences with less than 5% perplexity increase, 4096-token sequences with 10-15% increase, and even 10,000-token sequences with 20-30% increase—all without any fine-tuning. In contrast, the same models with sinusoidal position encodings show 50% perplexity increase at 2048 tokens and become essentially non-functional beyond 4096 tokens. This zero-shot extrapolation capability makes ALiBi particularly attractive for applications where the maximum sequence length is unknown at training time or varies widely across use cases.

ALiBi has been adopted by several prominent models including BLOOM (176B parameters) and MPT (7B-30B parameters). BLOOM was trained with ALiBi on sequences up to 2048 tokens but can effectively process sequences of 4096 tokens or longer at inference time. MPT models trained with ALiBi on 2048-token sequences have been successfully deployed on tasks requiring 8192-token contexts with minimal quality degradation.

The computational overhead of ALiBi is negligible. Computing the bias $-m \cdot |i - j|$ for all n^2 attention scores requires $O(n^2)$ operations, which is dominated by the $O(n^2d)$ cost of computing \mathbf{QK}^\top . The bias can be precomputed once per sequence and reused across all layers and heads (with different

slopes m), further reducing overhead. Memory overhead is also minimal, as the bias matrix can be computed on-the-fly or stored once and reused.

The primary limitation of ALiBi is that it assumes a monotonic relationship between distance and relevance—more distant tokens are always penalized more heavily. This assumption holds well for many natural language tasks where local context is indeed more important than distant context. However, for tasks with long-range dependencies that are not distance-dependent (such as matching opening and closing brackets in code, or resolving coreferences across document sections), ALiBi’s linear bias may be suboptimal compared to learned position encodings that can capture more complex position-dependent patterns.

19.3 Efficient Attention for Long Context

19.3.1 Sparse Attention Patterns

While position encoding improvements enable better extrapolation, they do not address the fundamental quadratic scaling of attention computation and memory. Efficient attention mechanisms reduce this quadratic bottleneck by restricting which tokens can attend to which other tokens, computing attention only over a subset of the n^2 possible connections. These sparse attention patterns can reduce complexity from $O(n^2)$ to $O(n \times w)$ where w is a fixed window size, enabling processing of sequences that would be impossible with full attention.

The key insight behind sparse attention is that not all token pairs require attention computation. In many domains, particularly natural language, most relevant information comes from nearby tokens, with occasional long-range dependencies. By carefully designing sparsity patterns that preserve important connections while eliminating redundant ones, sparse attention can maintain model quality while dramatically reducing computational and memory requirements.

19.3.2 Longformer: Local and Global Attention

Longformer combines local windowed attention for all tokens with global attention for task-specific tokens, enabling efficient processing of documents up to 4096 tokens or longer. This hybrid approach recognizes that while most tokens primarily need local context, certain special tokens (such as [CLS] for classification or question tokens for QA) benefit from attending to the entire sequence.

Definition 19.4 (Longformer Attention). Longformer defines attention patterns based on token type:

Local attention: Each token attends to tokens within a fixed window of size w on each side:

$$\mathcal{S}_{\text{local}}(i) = \{j : |i - j| \leq w\} \quad (19.9)$$

Global attention: Designated global tokens attend to all positions and are attended by all positions:

$$\mathcal{S}_{\text{global}}(i) = \{1, 2, \dots, n\} \text{ for global tokens} \quad (19.10)$$

For regular tokens, the attention set is $\mathcal{S}(i) = \mathcal{S}_{\text{local}}(i) \cup \mathcal{G}$ where \mathcal{G} is the set of global token positions.

The computational complexity of Longformer is $O(n \times w)$ for local attention plus $O(n \times g)$ for global attention, where g is the number of global tokens. For typical configurations with $w = 512$ and $g = 2$, this gives total complexity $O(n \times 514)$ compared to $O(n^2)$ for full attention. At $n = 4096$, this represents a reduction from 16.8 million connections to 2.1 million connections—an $8\times$ reduction.

Memory requirements scale similarly. For a single attention head with $d = 768$ in FP32, Longformer with $w = 512$ and $g = 2$ requires approximately $(4094 \times 1026 + 2 \times 4096) \times 4 = 16.8$ MB per head compared to 67 MB for full attention. With 12 heads and 12 layers, total attention memory decreases from 9.6 GB to 2.4 GB, enabling processing of long documents on GPUs with limited memory.

The window size w determines the trade-off between efficiency and model capacity. Smaller windows reduce computation and memory but limit the model’s ability to capture long-range dependencies. Information can propagate at most w positions per layer, requiring $\lceil n/w \rceil$ layers for full sequence communication. With $w = 512$ and $n = 4096$, information requires at least 8 layers to propagate across the full sequence. In practice, Longformer uses 12 layers, providing sufficient depth for information flow while maintaining efficiency.

Longformer has been successfully applied to long-document tasks including question answering, document classification, and summarization. On the WikiHop dataset with documents averaging 3000 tokens, Longformer achieves 74.3% accuracy compared to 68.5% for BERT with truncated 512-token contexts, demonstrating that access to full document context improves performance. On the arXiv summarization task with papers averaging 4500 tokens, Longformer generates summaries with ROUGE scores 3-5 points higher than models limited to shorter contexts.

19.3.3 BigBird: Random, Window, and Global Attention

BigBird extends sparse attention by adding random connections to Longformer’s local and global patterns. These random connections create shortcuts across the sequence that enable faster information propagation and provide theoretical guarantees about the model’s expressive power. BigBird proves that sparse attention with $O(n)$ connections can approximate full attention’s ability to compute arbitrary sequence-to-sequence functions.

Definition 19.5 (BigBird Attention). BigBird combines three attention patterns for each query position i :

1. **Window attention:** Attend to w neighbors on each side: $\mathcal{W}(i) = \{j : |i - j| \leq w\}$
2. **Random attention:** Attend to r randomly selected positions: $\mathcal{R}(i)$ with $|\mathcal{R}(i)| = r$
3. **Global attention:** Attend to g designated global tokens: \mathcal{G}

The total attention set is $\mathcal{S}(i) = \mathcal{W}(i) \cup \mathcal{R}(i) \cup \mathcal{G}$, giving $|\mathcal{S}(i)| = 2w + r + g$ connections per query.

The random attention component is crucial for BigBird’s theoretical properties. While local windows provide nearby context and global tokens enable information aggregation, random connections create a small-world network structure where any two positions are connected by a short path through the attention graph. With $r = O(\log n)$ random connections per query, the expected path length between any two positions is $O(\log n)$, enabling efficient information propagation across the sequence.

BigBird’s theoretical contribution is proving that this sparse pattern with $O(n \log n)$ total connections can approximate any function that full attention with $O(n^2)$ connections can compute. Specifically, BigBird shows that sparse attention with random connections is a universal approximator for sequence-to-sequence functions, providing a rigorous foundation for sparse attention methods. This result demonstrates that the quadratic number of connections in full attention is not necessary for expressive power—carefully chosen linear or near-linear connections suffice.

In practice, BigBird uses $w = 256$, $r = 64$, and $g = 32$ for sequences up to 4096 tokens. Each query attends to $2 \times 256 + 64 + 32 = 608$ keys instead of 4096, reducing computation by $6.7\times$. For a single attention head with $d = 768$ in FP32, BigBird requires approximately $((4096 - 32) \times 608 + 32 \times 4096) \times 4 = 10.4$ MB compared to 67 MB for full attention, a $6.4\times$ memory reduction. With 12 heads and 12 layers, total attention memory decreases from 9.6 GB to 1.5 GB.

On an NVIDIA A100 GPU, BigBird processes 4096-token sequences in approximately 15 milliseconds per layer compared to 98 milliseconds for full attention, a $6.5\times$ speedup. For sequences of 8192 tokens, BigBird takes 30 milliseconds per layer while full attention would require approximately 390 milliseconds, a $13\times$ speedup that makes previously impractical sequence lengths feasible. The speedup is slightly less than the theoretical reduction factor due to overhead from irregular memory access patterns in the random attention component, but remains substantial in practice.

BigBird has been applied to long-document natural language processing tasks with strong results. On the Natural Questions dataset with 4096-token contexts, BigBird achieves 79.2% F1 score compared to 76.8% for BERT with truncated 512-token contexts, demonstrating that access to longer context improves performance on tasks requiring long-range reasoning. On genomic sequence analysis tasks with sequences of 4096 base pairs, BigBird outperforms both full attention (which is memory-prohibitive) and simple local attention (which lacks long-range connections).

19.3.4 Comparison of Sparse Attention Methods

Different sparse attention patterns offer different trade-offs between efficiency, model quality, and implementation complexity. Local attention with window size w provides the simplest pattern and best memory locality, achieving $O(nw)$ complexity with straightforward implementation. However, information propagation is limited to w positions per layer, requiring deep networks for long-range dependencies. Longformer's addition of global tokens addresses this limitation by providing information hubs, enabling faster propagation while maintaining linear complexity. BigBird's random connections provide theoretical guarantees and empirically strong performance, but at the cost of irregular memory access patterns that reduce hardware efficiency.

For sequences up to 2048 tokens, the overhead of sparse attention often outweighs its benefits—full attention with optimized implementations like Flash Attention is typically faster and simpler. For sequences of 2048-8192 tokens, sparse attention becomes beneficial, with Longformer and BigBird providing good trade-offs between efficiency and quality. For sequences beyond 8192 tokens, sparse attention becomes essential, as full attention exceeds available memory on most hardware. At these lengths, the choice between Longformer and BigBird depends on the task: Longformer is simpler and faster for tasks where local context dominates, while BigBird provides better quality for tasks requiring complex long-range reasoning.

19.4 Recurrent Transformers

19.4.1 Transformer-XL

Definition 19.6 (Transformer-XL). Segment long sequence, reuse representations from previous segments:

Segment n : Tokens $[s_n, s_n + 1, \dots, s_n + L - 1]$

Compute:

$$\mathbf{h}_n = \text{Transformer}([\text{stop_grad}(\mathbf{h}_{n-1}), \mathbf{x}_n]) \quad (19.11)$$

Previous segment hidden states provide additional context without recomputation!

Example 19.2 (Transformer-XL Processing). Segment length: $L = 512$

Segment 1: Process tokens 0-511

- Save hidden states \mathbf{h}_1

Segment 2: Process tokens 512-1023

- Concatenate with \mathbf{h}_1 (frozen)
- Effective context: $512 + 512 = 1024$ tokens
- Computation: Still $O(512^2)$ per segment

Segment 3: Process tokens 1024-1535

- Use \mathbf{h}_2 from previous segment
 - Effective context: $1024 + 512 = 1536$ tokens
- Context grows linearly with segments, computation stays constant!

Relative position encodings: Modified for segment-level recurrence

19.5 Retrieval-Augmented Generation

19.5.1 RAG Architecture

Definition 19.7 (Retrieval-Augmented Generation). Combine retrieval with generation:

Step 1: Retrieval

$$\text{docs} = \text{Retrieve}(\text{query}, \text{corpus}, k = 5) \quad (19.12)$$

Step 2: Concatenate

$$\text{input} = [\text{docs}_1, \dots, \text{docs}_k, \text{query}] \quad (19.13)$$

Step 3: Generate

$$\text{output} = \text{LM}(\text{input}) \quad (19.14)$$

Retrieval methods:

- BM25 (sparse)
- Dense retrieval (BERT embeddings + nearest neighbors)
- Hybrid (combine sparse and dense)

Example 19.3 (RAG for Question Answering). **Question:** "When was the Eiffel Tower built?"

Step 1: Retrieve (from Wikipedia)

1. "The Eiffel Tower was constructed from 1887 to 1889..."
2. "Gustave Eiffel designed the tower for the 1889 World's Fair..."
3. "The tower is 330 meters tall and was the tallest..."

Step 2: Concatenate

Context 1: The Eiffel Tower was constructed from 1887 to 1889...

Context 2: Gustave Eiffel designed the tower for the 1889 World's Fair...

Context 3: The tower is 330 meters tall and was the tallest...

Question: When was the Eiffel Tower built?

Answer:

Step 3: Generate "The Eiffel Tower was built from 1887 to 1889."

Advantages:

- Access to external knowledge
- No need to fit everything in context window

- Cite sources
- Update knowledge without retraining

19.5.2 RETRO: Retrieval-Enhanced Transformer

Architecture:

- Chunk input into segments (64 tokens)
- Retrieve neighbors for each chunk
- Cross-attend to retrieved chunks
- Chunked cross-attention layers

Performance: 25× fewer parameters with retrieval achieves same performance as larger model without retrieval!

19.6 Memory-Augmented Transformers

19.6.1 Compressive Transformer

Definition 19.8 (Compressive Transformer). Extend Transformer-XL with compression:
Three levels of memory:

1. **Active:** Current segment (full attention)
2. **Recent:** Last n_m segments (cached, full precision)
3. **Compressed:** Older segments (compressed representations)

Compression:

- Learned compression function
- Reduce n tokens to n/c (e.g., $c = 3$)
- Compression ratio balances memory vs information

Effective context: Active + Recent + Compressed

$$L_{\text{eff}} = L + n_m \cdot L + n_c \cdot (L/c) \quad (19.15)$$

19.6.2 Memorizing Transformers

Key innovation: k -NN attention over entire history

Architecture:

- Store all past (*key*, *value*) pairs in memory
- For each query, retrieve k nearest neighbors
- Attend to local context + retrieved keys/values

Benefits:

- Effectively infinite context (limited by storage)
- Constant-time attention (with approximate k -NN)
- Improves perplexity on long documents

19.7 Long Context Models in Practice

19.7.1 LongT5: Efficient Encoder-Decoder

LongT5 extends the T5 encoder-decoder architecture to handle sequences up to 16,384 tokens by applying efficient attention mechanisms to both the encoder and decoder. Unlike decoder-only models that process sequences autoregressively, encoder-decoder models must handle long sequences in both components, making efficiency doubly important. LongT5 demonstrates that sparse attention patterns can be successfully applied to encoder-decoder architectures while maintaining the strong performance of the original T5 model.

LongT5 uses a combination of local and global attention patterns similar to Longformer, but adapted for the encoder-decoder structure. The encoder uses local attention with window size $w = 512$ for all tokens, plus global attention for a small number of designated tokens. The decoder uses local attention for attending to its own previous tokens, plus full attention to encoder outputs (which are compressed through the local attention mechanism). This asymmetric design recognizes that decoder-to-encoder attention is typically less memory-intensive than encoder self-attention, as decoder sequences are usually shorter than encoder sequences.

The memory savings from LongT5's sparse attention are substantial. For an encoder sequence of length $n_e = 16384$ and decoder sequence of length $n_d = 512$, full attention would require approximately $(16384^2 + 512^2 + 512 \times 16384) \times 4 = 1.1$ GB per attention head in FP32 for encoder self-attention, decoder self-attention, and cross-attention combined. With LongT5's sparse patterns using $w = 512$, the requirement reduces to approximately $(16384 \times 1024 + 512^2 + 512 \times 16384) \times 4 = 71$ MB per head, a $15\times$ reduction. With 12 heads and 12 encoder layers plus 12 decoder layers, total attention memory decreases from 26 GB to 1.7 GB.

LongT5 has been successfully applied to long-document summarization tasks where input documents exceed 10,000 tokens. On the arXiv summarization dataset with papers averaging 6000 tokens, LongT5 achieves ROUGE-L scores of 48.3 compared to 44.1 for T5 with truncated 512-token inputs, demonstrating that access to full document context significantly improves summary quality. On the PubMed summarization task with medical papers averaging 3000 tokens, LongT5 outperforms T5 by 3-4 ROUGE points across all metrics.

19.7.2 GPT-4 and Claude: Production Long Context Systems

Modern production language models have pushed context lengths to unprecedented scales, with GPT-4 Turbo supporting 128,000 tokens and Claude 2 supporting 100,000 tokens. These context lengths enable entirely new applications: processing entire codebases, analyzing full-length books, maintaining context across extended conversations, and performing complex multi-document reasoning tasks. While the specific architectural details of these models remain proprietary, their capabilities demonstrate that long context processing at scale is not only feasible but practical for real-world deployment.

GPT-4 Turbo's 128,000-token context window represents approximately 300 pages of text or roughly 96,000 words. This enables processing of entire novels, comprehensive technical documentation, or large codebases in a single context. The model maintains strong performance across this full context length, with perplexity remaining relatively stable even for tokens at positions 100,000+. This suggests the use of advanced position encoding schemes (likely RoPE with position interpolation) combined with efficient attention mechanisms that reduce the quadratic memory bottleneck.

The computational and memory requirements for such long contexts are staggering without optimization. Naive full attention for 128,000 tokens with GPT-4 scale parameters would require terabytes of memory just for attention matrices, making it completely impractical. The fact that GPT-4 Turbo can process such contexts in reasonable time (typically 30-60 seconds for full context) indicates the use of multiple optimization techniques: sparse or approximate attention to reduce the quadratic bottleneck, Flash Attention or similar memory-efficient implementations to reduce memory bandwidth requirements, and likely model parallelism to distribute computation across multiple GPUs.

Claude 2's 100,000-token context (approximately 75,000 words) similarly enables processing of book-length documents. Anthropic has demonstrated Claude 2's ability to accurately retrieve information from anywhere within its context window, suggesting effective attention mechanisms that don't degrade

at long distances. The model can answer questions about specific details from early in a 100,000-token context while processing tokens at the end, indicating that information flow across the full context length is maintained.

The pricing structure for these long-context models reflects their computational costs. GPT-4 Turbo charges approximately $3\times$ more per token for the 128K context version compared to the 8K context version, indicating that the computational overhead of long context processing remains significant despite optimizations. Claude 2's pricing similarly reflects higher costs for longer contexts. These pricing differences suggest that while long context is feasible, it remains computationally expensive compared to shorter contexts, motivating continued research into more efficient methods.

19.7.3 Practical Considerations for Long Context

Deploying long context models in production requires careful consideration of when the additional context is actually beneficial versus when alternative approaches might be more effective. Long context processing incurs real costs in terms of latency, computational resources, and financial expense, so understanding when these costs are justified is essential for practical applications.

Long context is most valuable when the task requires reasoning over or synthesizing information from multiple parts of a long document. Document summarization benefits significantly from full document access, as important information may appear anywhere in the document. Question answering over long documents similarly benefits from long context, as the relevant information's location is unknown in advance. Code generation and analysis tasks benefit from seeing entire files or multiple related files in context, enabling the model to understand dependencies and maintain consistency.

However, many tasks that initially seem to require long context can be effectively addressed with shorter contexts and retrieval. For question answering over large document collections, retrieval-augmented generation (RAG) can retrieve only the relevant passages and provide them in a short context, achieving similar or better performance at much lower cost. For tasks requiring access to factual knowledge, retrieval from a knowledge base is often more reliable and efficient than encoding all knowledge in the context. For multi-turn conversations, summarizing or compressing earlier conversation history can maintain coherence while reducing context length.

The cost-benefit analysis depends on several factors. Latency requirements matter: long context processing takes longer, which may be unacceptable for interactive applications. Accuracy requirements matter: if the task requires very high accuracy and the model performs significantly better with full context, the additional cost may be justified. Update frequency matters: if the information changes frequently, retrieval from an updated database may be preferable to encoding static information in context. Scale matters: for high-volume applications, the per-request cost of long context processing multiplies, potentially making alternative approaches more economical.

A practical strategy is to use a hybrid approach: employ retrieval or summarization to reduce context length when possible, but fall back to full long context processing when the task genuinely requires it. For example, a document analysis system might first use retrieval to identify relevant sections, then process those sections with long context if they exceed the standard context limit. This approach balances the benefits of long context with the efficiency of shorter contexts, optimizing for both quality and cost.

19.8 Comparison and Trade-offs

19.8.1 Method Comparison

Different approaches to long context processing involve fundamentally different trade-offs between computational efficiency, memory usage, model quality, and implementation complexity. Understanding these trade-offs is essential for selecting the appropriate method for a given application, as no single approach dominates across all dimensions. The optimal choice depends on the specific requirements of the task, available hardware, and acceptable quality-efficiency trade-offs.

Standard full attention with optimized implementations like Flash Attention remains the gold standard for quality and simplicity when sequence lengths permit. For sequences up to 2048 tokens on

modern GPUs, full attention is typically the best choice: it provides the highest model quality, has the simplest implementation, and benefits from highly optimized kernels. Flash Attention reduces memory bandwidth requirements through kernel fusion, enabling batch sizes $2\text{-}4\times$ larger than naive implementations while maintaining identical outputs to standard attention. However, the fundamental $O(n^2)$ scaling means that full attention becomes impractical beyond 4096-8192 tokens on typical hardware.

Sparse attention methods like Longformer and BigBird reduce complexity to $O(n \times w)$ where w is a fixed window size, enabling sequences of 4096-16384 tokens on standard GPUs. These methods maintain exact attention within their connectivity pattern, avoiding approximation errors. The primary trade-off is that sparse patterns may miss important long-range dependencies that fall outside the connectivity pattern. For tasks where local context dominates (such as language modeling or most NLP tasks), this limitation has minimal impact on quality. For tasks requiring complex long-range reasoning (such as certain question answering or reasoning tasks), sparse attention may underperform full attention even when both are feasible.

Linear attention methods like Performer and Linformer achieve $O(n)$ complexity through mathematical approximations, enabling very long sequences of 32768 tokens or more. However, these approximations introduce errors that can degrade model quality. Performer uses random feature approximations to the softmax kernel, which works well for some attention distributions but poorly for others. Linformer assumes low-rank structure in attention matrices, which holds for many tasks but may fail for tasks with complex attention patterns. In practice, linear attention methods typically show 2-5% accuracy degradation on downstream tasks compared to full attention, which may or may not be acceptable depending on the application.

Recurrent methods like Transformer-XL process sequences in segments with recurrent connections, enabling unlimited context length with constant memory per segment. The trade-off is that information must propagate through multiple segments to flow across long distances, which can be slower than direct attention and may lose information through the recurrent bottleneck. Transformer-XL works well for tasks like language modeling where sequential processing is natural, but less well for tasks requiring bidirectional context or random access to different parts of the sequence.

Retrieval-augmented generation (RAG) sidesteps the context length problem entirely by retrieving only relevant information and providing it in a short context. This approach can handle effectively unlimited document collections while maintaining the quality and efficiency of short-context models. The trade-off is implementation complexity: RAG requires building and maintaining a retrieval system, embedding documents, and handling retrieval failures. Additionally, RAG works best for tasks where relevant information can be identified through retrieval, but may struggle with tasks requiring synthesis across many parts of a document or reasoning about information that is difficult to retrieve.

Method	Max Length	Complexity	Quality	Implementation
Full Attention	2-4K	$O(n^2d)$	Excellent	Simple
Flash Attention	4-8K	$O(n^2d)$	Excellent	Medium
Longformer	4-16K	$O(nwd)$	Good	Medium
BigBird	4-16K	$O(n(w+r)d)$	Good	Medium
Linformer	8-32K	$O(nkd)$	Fair	Medium
Performer	16-64K	$O(nd^2)$	Fair	Hard
Transformer-XL	Unlimited	$O(L^2d)/\text{seg}$	Good	Medium
RAG	Unlimited	$O(n^2d)$	Excellent	Hard

Table 19.1: Comparison of long context methods. Complexity is per-layer computational cost. Quality is relative to full attention on typical NLP tasks. Implementation difficulty considers both coding complexity and infrastructure requirements.

19.8.2 Hardware and Memory Considerations

The practical feasibility of different long context methods depends critically on available hardware and memory constraints. Modern GPUs vary widely in memory capacity, from 8 GB on consumer GPUs to 80 GB on high-end data center GPUs, and this memory capacity directly determines which sequence lengths are feasible with different methods.

For a BERT-base scale model with $d = 768$, $h = 12$ heads, and $L = 12$ layers, the memory requirements for different methods and sequence lengths are as follows. Full attention with $n = 2048$ requires approximately 2.4 GB for attention matrices across all layers, which fits comfortably on any modern GPU. At $n = 4096$, full attention requires 9.7 GB, which fits on 16 GB GPUs but leaves limited memory for batch processing. At $n = 8192$, full attention requires 38.5 GB, exceeding even high-end GPUs and requiring model parallelism or gradient checkpointing.

Sparse attention dramatically improves these numbers. Longformer with $w = 512$ and $n = 4096$ requires only 2.4 GB for attention matrices, enabling batch sizes $4\times$ larger than full attention on the same hardware. At $n = 8192$, Longformer requires 4.8 GB, which fits comfortably on 16 GB GPUs. At $n = 16384$, Longformer requires 9.6 GB, still feasible on standard hardware. This memory efficiency enables processing of long documents on commodity GPUs that would be impossible with full attention.

Linear attention methods like Linformer with $k = 256$ require even less memory. At $n = 4096$, Linformer requires only 600 MB for attention matrices, enabling very large batch sizes or processing on smaller GPUs. At $n = 16384$, Linformer requires 2.4 GB, comparable to full attention at $n = 2048$. This memory efficiency enables processing of very long sequences, but at the cost of approximation errors that may degrade quality.

The memory requirements extend beyond attention matrices to include activations, gradients, model parameters, and optimizer states. For training, the total memory requirement is typically $4\text{--}6\times$ the attention matrix memory when accounting for all these components. For inference, memory requirements are lower as gradients and optimizer states are not needed, but activations must still be stored for generation. These additional memory requirements mean that the feasible sequence length for training is typically $2\text{--}4\times$ shorter than for inference on the same hardware.

19.8.3 Recommendations by Use Case

Selecting the appropriate long context method requires matching the method’s characteristics to the specific requirements of the application. The following recommendations provide guidance based on common use cases and constraints.

For general NLP tasks with sequences up to 2048 tokens, use standard full attention with Flash Attention optimization. This provides the best quality with simple implementation and benefits from highly optimized libraries. The computational and memory costs are manageable on any modern GPU, and the simplicity reduces implementation and debugging time.

For document processing tasks with sequences of 2048–8192 tokens, use sparse attention methods like Longformer or BigBird. These methods provide good quality with manageable computational costs, and the sparse patterns align well with the local structure of natural language. Longformer is simpler and faster for tasks where local context dominates, while BigBird provides better quality for tasks requiring long-range reasoning. Both methods have well-tested implementations available in popular libraries.

For very long sequences of 8192–32768 tokens where quality is critical, consider using full attention with model parallelism or gradient checkpointing if hardware permits, or sparse attention if hardware is limited. The quality difference between full and sparse attention becomes more significant at these lengths, so the choice depends on whether the hardware can support full attention. If full attention is infeasible, BigBird typically provides better quality than Longformer at these lengths due to its random connections.

For extremely long sequences beyond 32768 tokens, or when processing large document collections, use retrieval-augmented generation (RAG) rather than attempting to fit everything in context. RAG provides better quality and efficiency than any long context method at these scales, as it focuses the model’s attention on relevant information rather than processing irrelevant content. The implementation complexity of RAG is justified by the significant quality and efficiency improvements at these scales.

For streaming or online processing tasks, use Transformer-XL or similar recurrent methods that can process sequences incrementally without recomputing previous segments. These methods enable unlimited context length with constant memory per segment, making them ideal for applications like real-time transcription, continuous monitoring, or interactive systems where the sequence length is

unbounded.

For tasks requiring frequent updates to the knowledge base or document collection, prefer RAG over long context methods. RAG allows updating the retrieval index without retraining the model, while long context methods require reprocessing the entire context whenever information changes. This makes RAG more practical for applications with dynamic information needs.

19.9 Exercises

Exercise 19.1 (Memory Calculation). Calculate the memory requirements for attention matrices in different scenarios:

1. For a BERT-base model ($d = 768$, $h = 12$, $L = 12$) with sequence lengths $n \in \{512, 1024, 2048, 4096, 8192\}$, compute the total memory for attention matrices in FP32 and FP16 formats.
2. For the same model using Longformer with window size $w = 512$ and 2 global tokens, compute the memory savings compared to full attention at each sequence length.
3. Determine the maximum sequence length that fits in 16 GB of GPU memory for full attention, assuming attention matrices consume 40% of available memory (the rest is for activations, parameters, etc.).
4. For a GPT-3 scale model ($d = 12288$, $h = 96$, $L = 96$), compute the memory required for $n = 2048$ tokens and explain why model parallelism is necessary.

Exercise 19.2 (Position Interpolation Implementation). Implement and evaluate position interpolation for extending context length:

1. Load a pretrained GPT-2 model (trained on 1024-token contexts).
2. Implement position interpolation to extend the model to 4096 tokens by scaling position indices by $1024/4096 = 0.25$.
3. Fine-tune the extended model on long sequences for 1000 steps.
4. Evaluate perplexity on sequences of length 1024, 2048, 3072, and 4096, comparing the interpolated model to the original model (which will fail on longer sequences).
5. Plot perplexity versus position to visualize how well the model handles different parts of the extended context.

Exercise 19.3 (Sparse Attention Patterns). Implement and compare different sparse attention patterns:

1. Implement local attention with window size $w = 256$ for a sequence of length $n = 2048$.
2. Implement strided attention with stride $s = 64$ for the same sequence.
3. Implement BigBird attention combining local ($w = 128$), random ($r = 32$), and global ($g = 4$) patterns.
4. For each pattern, compute the number of attention connections and compare to full attention ($n^2 = 4,194,304$ connections).

5. Visualize the attention patterns as sparse matrices and discuss which types of dependencies each pattern can capture.

Exercise 19.4 (ALiBi Extrapolation). Implement ALiBi and test its extrapolation capabilities:

1. Train a small transformer (4 layers, $d = 256$, 4 heads) with ALiBi on sequences of length 512 from a language modeling dataset.
2. Use head-specific slopes $m_i = 2^{-8i/4}$ for the 4 heads, giving slopes $\{0.25, 0.0625, 0.0156, 0.0039\}$.
3. Evaluate the trained model on sequences of length 512, 1024, 2048, and 4096 without any fine-tuning.
4. Compare to a model trained with sinusoidal position encodings on the same data.
5. Plot perplexity versus sequence length for both models and explain the difference in extrapolation behavior.

Exercise 19.5 (Retrieval-Augmented Generation). Implement a simple RAG system for question answering:

1. Create a document corpus of 1000 Wikipedia articles on a specific topic (e.g., history, science).
2. Embed all documents using a pretrained BERT model, storing embeddings in a FAISS index for efficient retrieval.
3. For a given question, retrieve the top-5 most relevant documents based on embedding similarity.
4. Concatenate the retrieved documents with the question and generate an answer using a pretrained language model.
5. Compare the quality of answers when using RAG versus providing the model with only the question (no retrieval).
6. Analyze cases where RAG succeeds and fails, discussing the importance of retrieval quality.

Exercise 19.6 (Transformer-XL Segment Processing). Implement segment-level recurrence for processing long sequences:

1. Implement a simplified Transformer-XL that processes a sequence in segments of length $L = 256$.
2. For each segment, cache the hidden states from the previous segment and concatenate them with the current segment's inputs.
3. Ensure gradients do not flow into the cached hidden states (use `stop_gradient` or `detach`).
4. Process a sequence of length 2048 in 8 segments, measuring the effective context length at each position.

5. Compare memory usage and computation time to processing the full 2048-token sequence at once.
6. Discuss the trade-off between effective context length and computational efficiency.

Exercise 19.7 (Long Context Cost Analysis). Analyze the computational and financial costs of long context processing:

1. For a BERT-base model, measure the actual wall-clock time to process sequences of length 512, 1024, 2048, and 4096 on your available hardware (CPU or GPU).
2. Compute the FLOPs for attention at each sequence length and compare to the measured time to determine hardware efficiency.
3. Estimate the cost of processing 1 million tokens at each sequence length, assuming cloud GPU pricing (e.g., \$2.50/hour for an A100).
4. Compare the cost of full attention versus Longformer with $w = 512$ at each sequence length.
5. Discuss scenarios where the higher cost of long context is justified versus where shorter contexts with retrieval would be more economical.

Exercise 19.8 (Position Encoding Comparison). Compare different position encoding schemes empirically:

1. Train four small transformer models (4 layers, $d = 256$, 4 heads) on the same language modeling dataset, using: (a) absolute learned positions, (b) sinusoidal positions, (c) RoPE, and (d) ALiBi.
2. Train all models on sequences of length 512 for the same number of steps.
3. Evaluate all models on sequences of length 512, 1024, 2048, and 4096.
4. Plot perplexity versus sequence length for each model.
5. Analyze which position encoding schemes extrapolate best and explain why based on their mathematical properties.
6. Fine-tune the absolute and sinusoidal models on longer sequences and compare to the zero-shot extrapolation of RoPE and ALiBi.

Chapter 20

Pre-training Strategies and Transfer Learning

Chapter Overview

Pre-training on large unlabeled corpora followed by task-specific fine-tuning has become the dominant paradigm in deep learning. This chapter covers pre-training objectives, data curation, curriculum learning, continual pre-training, and transfer learning strategies for maximizing downstream performance.

Learning Objectives

1. Understand different pre-training objectives and their trade-offs
2. Curate and process pre-training data at scale
3. Apply curriculum learning and domain-adaptive pre-training
4. Implement parameter-efficient fine-tuning (LoRA, adapters)
5. Design multi-task and multi-stage pre-training
6. Measure and improve transfer learning effectiveness

20.1 Pre-training Objectives

20.1.1 Language Modeling Objectives

The choice of pre-training objective fundamentally shapes what a model learns and how effectively it transfers to downstream tasks. Different objectives make different trade-offs between computational efficiency, representation quality, and suitability for specific task types. Understanding these trade-offs is essential for practitioners designing pre-training pipelines.

Causal Language Modeling (CLM) predicts each token given only previous context, formalizing the objective as:

$$\mathcal{L}_{\text{CLM}} = - \sum_{t=1}^T \log P(x_t | x_{<t}; \theta) \quad (20.1)$$

This objective is computationally efficient because it requires only a single forward pass through the model to compute the loss for all positions simultaneously. The unidirectional attention pattern means that for a sequence of length n , each position i attends only to positions 1 through $i - 1$, resulting in a triangular attention mask. This makes CLM particularly natural for text generation tasks where the model must produce tokens sequentially without access to future context. However, the unidirectional constraint limits the model's ability to build rich bidirectional representations, which can hurt performance on understanding tasks like question answering or natural language inference.

All GPT models, including GPT-2 and GPT-3, use this objective exclusively, optimizing for generation quality at the expense of bidirectional understanding.

Masked Language Modeling (MLM) addresses the unidirectional limitation by randomly masking tokens and predicting them using bidirectional context:

$$\mathcal{L}_{\text{MLM}} = - \sum_{t \in \mathcal{M}} \log P(x_t | x_{\setminus \mathcal{M}}; \theta) \quad (20.2)$$

where \mathcal{M} denotes the set of masked positions. BERT masks approximately 15% of tokens during pre-training, allowing each masked position to attend to all unmasked positions in both directions. This bidirectional context enables the model to build richer representations that capture relationships between words regardless of their relative positions. The computational cost is similar to CLM since a single forward pass computes predictions for all masked positions. However, MLM introduces a pre-training/fine-tuning mismatch: the special [MASK] token appears during pre-training but not during fine-tuning or inference. To mitigate this, BERT uses a mixed strategy where 80% of masked tokens are replaced with [MASK], 10% with random tokens, and 10% remain unchanged. Despite this mismatch, MLM has proven highly effective for understanding tasks, with BERT, RoBERTa, and ALBERT all achieving strong performance on GLUE and SQuAD benchmarks.

Prefix Language Modeling attempts to combine the benefits of both approaches by using bidirectional attention on a prefix and causal attention on the suffix:

$$\mathcal{L}_{\text{prefix}} = - \sum_{t=|p|+1}^T \log P(x_t | x_{<t}; \theta) \quad (20.3)$$

where p is the prefix with bidirectional attention. This allows the model to build rich representations of the prefix while maintaining the ability to generate continuations autoregressively. UniLM and GLM use variants of this objective, achieving strong performance on both understanding and generation tasks. The computational cost is slightly higher than pure CLM or MLM because the attention pattern is more complex, requiring careful implementation to avoid recomputing prefix representations for each suffix position.

20.1.2 Denoising Objectives

Denoising objectives extend beyond simple token masking to more complex corruption patterns that require the model to understand document structure and long-range dependencies. These objectives are particularly effective for sequence-to-sequence models where both the encoder and decoder must learn rich representations.

Span Corruption, used by T5, masks contiguous spans of tokens rather than individual tokens. A typical configuration corrupts 15% of tokens by replacing random spans (average length 3 tokens) with sentinel tokens. For example, the input "The quick brown fox jumps over the lazy dog" might become "The $\text{\textcircled{X}}$ fox $\text{\textcircled{Y}}$ the lazy dog" with the target " $\text{\textcircled{X}}$ quick brown $\text{\textcircled{Y}}$ jumps over $\text{\textcircled{Z}}$ ". This objective is more challenging than single-token masking because the model must predict multiple tokens for each sentinel, requiring it to understand the semantic coherence of spans. The computational cost is similar to MLM, but the increased difficulty leads to better representations for tasks requiring understanding of multi-token phrases and entities. T5 experiments showed that span corruption with mean span length 3 outperformed both single-token masking and longer spans, suggesting an optimal balance between difficulty and learnability.

Sentence Shuffling and Text Infilling, used by BART, apply more aggressive corruption strategies. Sentence shuffling permutes the order of sentences in a document, requiring the model to learn document-level structure and discourse coherence. Text infilling deletes spans of varying length (including zero-length spans), and the model must predict both the span length and content. BART combines multiple corruption strategies including token masking, token deletion, text infilling, sentence permutation, and document rotation. This multi-task denoising approach forces the model to handle diverse corruption patterns, leading to more robust representations. The computational cost is higher than simpler objectives because the encoder and decoder must both process variable-length sequences, but

the improved representation quality justifies this cost for sequence-to-sequence tasks. BART achieves state-of-the-art results on summarization and translation benchmarks, demonstrating the effectiveness of aggressive denoising.

20.1.3 Computational Cost Comparison

The computational cost of different pre-training objectives varies significantly, impacting both training time and resource requirements. Understanding these costs is essential for practitioners choosing objectives and estimating training budgets.

For a model with d dimensions, L layers, h attention heads, and sequences of length n , the dominant computational cost comes from matrix multiplications in attention and feed-forward layers. The attention mechanism requires computing queries, keys, and values ($3d^2$ parameters per layer), followed by attention scores ($O(n^2d)$ operations), and output projections (d^2 parameters). The feed-forward network typically uses an intermediate dimension of $4d$, requiring $8d^2$ parameters per layer. Summing across all layers, the total forward pass requires approximately $L(12d^2 + n^2d)$ FLOPs per token.

Causal Language Modeling is the most efficient objective because it computes loss for all n positions in a single forward pass. For a batch of b sequences, the total cost is $b \cdot L(12d^2n + n^3d)$ FLOPs. The n^3d term comes from computing attention scores for all positions, where position i attends to i previous positions, summing to $\frac{n(n+1)}{2} \approx \frac{n^2}{2}$ attention operations across all positions.

Masked Language Modeling has similar computational cost to CLM for the forward pass, but only computes loss for the 15% of masked positions. However, the backward pass still propagates gradients through all positions since masked positions depend on unmasked context. In practice, MLM and CLM have nearly identical training times for the same model size and sequence length. BERT-base training on 16 TPU v3 chips for 1 million steps with batch size 256 and sequence length 512 takes approximately 4 days, consuming roughly 10^{21} FLOPs total.

Span Corruption (T5) has slightly higher cost than MLM because the decoder must generate multiple tokens per corrupted span. If the average span length is s and the corruption rate is r , the decoder processes approximately $r \cdot s \cdot n$ tokens, compared to $r \cdot n$ for MLM. For T5's typical configuration with $s = 3$ and $r = 0.15$, this increases decoder cost by $3\times$, though the encoder cost remains unchanged. T5-base training on 1 trillion tokens takes approximately 7 days on 256 TPU v3 chips, consuming roughly 10^{21} FLOPs.

BART's multi-task denoising has the highest computational cost due to variable-length sequences and multiple corruption strategies. The encoder and decoder must both handle sequences of varying lengths, reducing batching efficiency. Additionally, applying multiple corruption strategies requires more data preprocessing. In practice, BART training takes approximately 20% longer than comparable BERT training for the same number of tokens, but the improved performance on generation tasks often justifies this additional cost.

Example 20.1 (Objective Comparison on Same Data). Text: "The quick brown fox jumps over the lazy dog"

CLM: Predict each token given previous

```
The -> quick
The quick -> brown
The quick brown -> fox
...
```

MLM (15% masking):

```
Input: "The [MASK] brown fox [MASK] over the lazy dog"
Target: "quick" at position 2, "jumps" at position 5
```

Span Corruption:

```
Input: "The <X> fox <Y> the lazy dog"
Target: "<X> quick brown <Y> jumps over <Z>"
```

Different objectives lead to different learned representations!

20.1.4 Contrastive Objectives

Contrastive Learning:

$$\mathcal{L}_{\text{contrastive}} = -\log \frac{\exp(\text{sim}(z_i, z_i^+)/\tau)}{\sum_j \exp(\text{sim}(z_i, z_j)/\tau)} \quad (20.4)$$

Applications:

- SimCLR (vision): Augmented views as positives
- CLIP: Image-text pairs
- SimCSE (text): Dropout as augmentation

20.2 Data Curation and Processing

20.2.1 Data Scale and Requirements

The scale of pre-training data has grown exponentially over the past few years, driven by empirical findings that larger datasets consistently improve model performance. Understanding the data requirements, storage costs, and preprocessing overhead is essential for planning pre-training projects.

BERT was pre-trained on approximately 16 GB of text data, consisting of BooksCorpus (800 million words) and English Wikipedia (2.5 billion words). This relatively modest dataset size reflects BERT's focus on high-quality, curated text rather than massive web crawls. The 16 GB of raw text expands to approximately 3.3 billion tokens using BERT's WordPiece tokenizer with a 30,000 token vocabulary. Training BERT-base for 1 million steps with batch size 256 and sequence length 512 means the model sees each token approximately 40 times on average, indicating significant data reuse through multiple epochs. The storage requirements are minimal by modern standards—16 GB of compressed text expands to perhaps 50 GB including tokenized data and intermediate preprocessing artifacts.

GPT-2 scaled up to approximately 40 GB of text from WebText, a dataset created by scraping outbound links from Reddit posts with at least 3 karma. This filtering strategy aimed to identify high-quality content as judged by the Reddit community. The 40 GB corpus contains roughly 8 billion tokens using GPT-2's byte-pair encoding with a 50,257 token vocabulary. GPT-2's largest variant (1.5B parameters) was trained for approximately 1 million steps, seeing each token roughly 10 times. The preprocessing pipeline for WebText involved deduplication, filtering by language, and removing low-quality content, reducing the raw crawl from over 100 GB to the final 40 GB. Storage requirements including raw data, filtered data, and tokenized sequences total approximately 150 GB.

GPT-3 made a massive leap to approximately 570 GB of text, totaling roughly 300 billion tokens. This dataset combines filtered Common Crawl (410 GB), WebText2 (19 GB), Books1 (12 GB), Books2 (55 GB), and Wikipedia (3 GB). The preprocessing pipeline for Common Crawl is particularly intensive: the raw crawl contains petabytes of data, which must be filtered by language, deduplicated, and quality-filtered to produce the final 410 GB. This filtering process itself requires substantial computational resources—processing petabytes of data through language classifiers and deduplication algorithms takes weeks on large clusters. The total storage requirements for GPT-3 pre-training, including raw data, filtered data, tokenized sequences, and training checkpoints, exceed 5 TB. The preprocessing cost alone is estimated at tens of thousands of dollars in compute time.

LLaMA pushed the scale even further to approximately 1.4 TB of text, totaling roughly 1.4 trillion tokens. This dataset consists primarily of Common Crawl (67%), C4 (15%), GitHub (4.5%), Wikipedia (4.5%), books (4.5%), ArXiv (2.5%), and StackExchange (2%). The inclusion of code from GitHub and technical content from ArXiv and StackExchange reflects a deliberate strategy to improve reasoning and

technical capabilities. The preprocessing pipeline for LLaMA is even more sophisticated than GPT-3, using multiple quality filters including perplexity-based filtering, classifier-based filtering, and extensive deduplication. The total storage requirements exceed 10 TB including all preprocessing artifacts, and the preprocessing cost is estimated at over \$100,000 in compute time.

20.2.2 Data Quality versus Quantity

The relationship between data quality and quantity is not straightforward—more data does not always lead to better models if the quality is poor. Recent research has shown that careful data curation can match or exceed the performance of models trained on much larger but noisier datasets.

High-quality datasets like Wikipedia and books consistently improve model performance even when they represent a small fraction of total training data. GPT-3’s data mixture samples Wikipedia at 3× the rate it appears in the corpus (3.4 epochs versus 0.44 epochs for Common Crawl), reflecting the higher quality and information density of Wikipedia text. This upsampling strategy means that despite Wikipedia being only 3 GB of the 570 GB total, it contributes disproportionately to the model’s knowledge and capabilities.

The preprocessing cost for achieving high data quality is substantial. Language identification using fastText classifiers requires processing every document, taking approximately 1 CPU-hour per 100 GB of text. Deduplication using MinHash LSH is even more expensive, requiring approximately 10 CPU-hours per 100 GB for computing signatures and finding near-duplicates. Quality filtering using perplexity-based methods requires running a language model over the entire corpus, taking approximately 100 GPU-hours per 100 GB. For GPT-3’s 570 GB dataset, the total preprocessing cost exceeds 50,000 CPU-hours and 5,000 GPU-hours, translating to roughly \$30,000 in cloud computing costs.

The storage requirements for preprocessing are also significant. Deduplication requires storing hash signatures for all documents, typically requiring 100-200 bytes per document. For a corpus with 1 billion documents, this means 100-200 GB of signature storage. Near-duplicate detection using LSH requires storing multiple hash tables, potentially doubling or tripling this storage requirement. Quality filtering requires storing perplexity scores or classifier outputs for all documents, adding another 10-20 GB. In total, the preprocessing metadata can require 500 GB to 1 TB of storage for a large corpus, comparable to the size of the corpus itself.

20.2.3 Data Filtering and Cleaning

[H] Data Filtering Pipeline

Step 1: Quality Filtering

- Remove duplicates (exact and near-duplicates)
- Filter by language (fastText classifier)
- Remove toxic/harmful content
- Filter low-quality (perplexity-based, classifier)

Step 2: Deduplication

- Exact match: Hash-based
- Near-duplicates: MinHash LSH
- Document-level and paragraph-level

Step 3: Privacy

- Remove PII (emails, phone numbers, addresses)
- Filter memorized content
- Redact sensitive information

Step 4: Formatting

- Unicode normalization
- Remove excessive whitespace
- Clean HTML/markup artifacts

Example 20.2 (GPT-3 Data Mixture). Total: 570GB, 300B tokens

Higher-quality sources sampled more frequently (multiple epochs). Lower-quality sources seen less to avoid overfitting to noise.

20.2.4 Data Deduplication**Why deduplicate?**

- Prevents memorization
- Better generalization
- Fairer evaluation (test set contamination)

Methods:**1. Exact Deduplication:**

```
seen_hashes = set()
for doc in corpus:
    hash_val = hash(doc)
    if hash_val not in seen_hashes:
        keep(doc)
        seen_hashes.add(hash_val)
```

2. Fuzzy Deduplication (MinHash):

- Compute MinHash signatures
- Use LSH for near-neighbor search
- Remove documents with Jaccard similarity > 0.8

20.3 Training Compute Requirements**20.3.1 FLOPs Analysis**

Understanding the computational requirements for pre-training large language models is essential for planning projects and estimating costs. The total compute is typically measured in FLOPs (floating-point operations), which can be calculated from model architecture and training configuration.

For a transformer model with L layers, d model dimension, h attention heads, and feed-forward intermediate dimension d_{ff} (typically $4d$), processing a single token requires approximately:

$$\text{FLOPs per token} = 2L(12d^2 + 4d \cdot d_{ff}) = 2L(12d^2 + 16d^2) = 56Ld^2 \quad (20.5)$$

The factor of 2 accounts for both forward and backward passes (backward pass requires approximately the same FLOPs as forward pass). The $12d^2$ term comes from attention projections (query, key, value, and output, each $d \times d$), and the $16d^2$ term comes from feed-forward layers (two $d \times 4d$ projections).

For **BERT-base** with $L = 12$, $d = 768$, training on 3.3 billion tokens for 40 epochs (132 billion tokens total):

$$\text{Total FLOPs} = 56 \times 12 \times 768^2 \times 132 \times 10^9 \approx 5.2 \times 10^{20} \text{ FLOPs} \quad (20.6)$$

This is approximately 0.5 zettaFLOPs. Training on 16 TPU v3 chips (each providing 420 TFLOPS in mixed precision) for 4 days:

$$\text{Available compute} = 16 \times 420 \times 10^{12} \times 4 \times 86400 \approx 2.3 \times 10^{21} \text{ FLOPs} \quad (20.7)$$

The ratio of available compute to required compute is approximately 4.4, indicating that BERT-base training achieves roughly 23% hardware utilization. This is typical for large-scale training where communication overhead, data loading, and other inefficiencies reduce effective utilization.

For **GPT-3 175B** with $L = 96$, $d = 12288$, training on 300 billion tokens:

$$\text{Total FLOPs} = 56 \times 96 \times 12288^2 \times 300 \times 10^9 \approx 2.4 \times 10^{23} \text{ FLOPs} \quad (20.8)$$

This is approximately 240 zettaFLOPs, nearly 500× more than BERT-base. The massive compute requirement reflects both the larger model (175B versus 110M parameters) and the larger dataset (300B versus 132B tokens). Training GPT-3 on approximately 10,000 NVIDIA V100 GPUs (each providing 125 TFLOPS in mixed precision) for 1 month:

$$\text{Available compute} = 10000 \times 125 \times 10^{12} \times 30 \times 86400 \approx 3.2 \times 10^{23} \text{ FLOPs} \quad (20.9)$$

This suggests approximately 75% hardware utilization, which is impressive for such a large-scale distributed training job. The higher utilization compared to BERT reflects improvements in distributed training infrastructure and optimization techniques.

For **LLaMA-65B** with $L = 80$, $d = 8192$, training on 1.4 trillion tokens:

$$\text{Total FLOPs} = 56 \times 80 \times 8192^2 \times 1.4 \times 10^{12} \approx 3.3 \times 10^{23} \text{ FLOPs} \quad (20.10)$$

This is approximately 330 zettaFLOPs. Training on 2048 NVIDIA A100 GPUs (each providing 312 TFLOPS in mixed precision) for 21 days:

$$\text{Available compute} = 2048 \times 312 \times 10^{12} \times 21 \times 86400 \approx 1.2 \times 10^{24} \text{ FLOPs} \quad (20.11)$$

This suggests approximately 28% hardware utilization, which is lower than GPT-3 despite using more modern hardware. The lower utilization likely reflects the challenges of scaling to very long sequences (LLaMA uses 2048 token sequences versus GPT-3's 2048 tokens) and the overhead of processing the much larger dataset.

20.3.2 GPU-Hours and Cost Estimates

Translating FLOPs into GPU-hours and cost estimates provides a more practical understanding of training requirements. The cost depends heavily on the hardware platform and whether using cloud services or owned infrastructure.

BERT-base training on 16 TPU v3 chips for 4 days equals 1,536 TPU-hours. At Google Cloud's on-demand pricing of approximately \$8 per TPU v3 hour, this costs roughly \$12,000. However, Google's original BERT paper reported using preemptible TPUs at approximately \$2.40 per hour, reducing the cost to roughly \$3,700. Using equivalent GPU resources (approximately 64 NVIDIA V100 GPUs to match 16 TPU v3 chips), the cost would be approximately \$2 per GPU-hour on AWS, totaling \$12,000 for 96 days of GPU time. The lower TPU cost reflects Google's optimization for transformer workloads and economies of scale.

GPT-3 175B training on 10,000 V100 GPUs for 1 month equals 7.2 million GPU-hours. At AWS on-demand pricing of approximately \$3 per V100 hour, this would cost \$21.6 million. However, OpenAI likely used a combination of owned infrastructure and negotiated cloud pricing, with estimates suggesting actual costs between \$4 million and \$12 million. The wide range reflects uncertainty about the exact

hardware configuration, utilization rates, and pricing agreements. The training also required substantial infrastructure costs including high-bandwidth networking (InfiniBand or equivalent), distributed storage systems, and engineering effort to optimize the training pipeline.

LLaMA-65B training on 2048 A100 GPUs for 21 days equals 1.03 million GPU-hours. At cloud pricing of approximately \$3 per A100 hour, this would cost \$3.1 million. Meta’s paper reports that LLaMA-65B training consumed approximately 1,022,362 GPU-hours on A100-80GB GPUs, closely matching this estimate. Using Meta’s owned infrastructure rather than cloud services likely reduced the effective cost to \$1.5-2 million when accounting for hardware depreciation and operational costs. The A100’s higher performance compared to V100 (312 versus 125 TFLOPS) means that LLaMA-65B required only 1/7 the GPU-hours of GPT-3 despite using comparable compute (330 versus 240 zettaFLOPs), demonstrating the importance of hardware efficiency.

20.3.3 Scaling Laws

Empirical scaling laws, first systematically studied by Kaplan et al. (2020) and refined by Hoffmann et al. (2022), provide predictive relationships between compute budget, model size, dataset size, and performance. These laws are essential for planning pre-training projects and allocating resources optimally.

The Kaplan scaling laws suggest that test loss scales as a power law with compute budget C , model parameters N , and dataset size D :

$$L(C) \propto C^{-\alpha}, \quad L(N) \propto N^{-\beta}, \quad L(D) \propto D^{-\gamma} \quad (20.12)$$

with empirically measured exponents $\alpha \approx 0.05$, $\beta \approx 0.076$, and $\gamma \approx 0.095$. These laws imply that doubling the compute budget reduces loss by approximately 3.4%, while doubling model size reduces loss by approximately 5.2%, and doubling dataset size reduces loss by approximately 6.5%.

The Chinchilla scaling laws (Hoffmann et al., 2022) refined these relationships and found that most large language models are significantly undertrained—they use too many parameters for their compute budget and too few training tokens. The optimal allocation suggests that model size and dataset size should scale approximately equally with compute budget. For a compute budget of C FLOPs, the optimal model size is approximately:

$$N_{\text{opt}} \approx 0.5 \times C^{0.5} \text{ parameters} \quad (20.13)$$

and the optimal dataset size is approximately:

$$D_{\text{opt}} \approx 0.5 \times C^{0.5} \text{ tokens} \quad (20.14)$$

This implies that GPT-3 175B, trained on 300B tokens with compute budget 2.4×10^{23} FLOPs, should optimally have been approximately 70B parameters trained on 1.4T tokens. This prediction closely matches LLaMA-65B’s configuration, which indeed achieves better performance than GPT-3 despite having fewer parameters. The Chinchilla paper demonstrated this by training Chinchilla (70B parameters, 1.4T tokens) to outperform Gopher (280B parameters, 300B tokens) using the same compute budget.

These scaling laws have profound implications for pre-training strategy. Rather than maximizing model size for a given compute budget, practitioners should balance model size and dataset size according to the scaling laws. For a fixed compute budget, training a smaller model on more data typically yields better performance than training a larger model on less data. This insight has driven the trend toward models like LLaMA, Chinchilla, and Mistral that prioritize training tokens over parameter count.

20.4 Curriculum Learning

20.4.1 Progressive Training Strategies

Curriculum learning applies the principle of learning from easy to hard examples, progressively increasing task difficulty during training. This approach can significantly improve training efficiency,

convergence speed, and final model performance. The key insight is that models learn more effectively when they first master simpler patterns before tackling complex ones.

The most common curriculum strategy involves progressively increasing sequence length during training. Starting with shorter sequences reduces both memory requirements and computational cost in the early stages of training when the model is learning basic patterns. For a model with quadratic attention complexity $O(n^2d)$, doubling the sequence length quadruples the attention computation. Training BERT-base with sequence length 128 for the first 90% of steps and then 512 for the final 10% reduces total training time by approximately 30% compared to using length 512 throughout. The shorter sequences allow larger batch sizes in early training, which improves gradient estimates and accelerates convergence. The model learns word-level and phrase-level patterns with short sequences, then refines its understanding of long-range dependencies with longer sequences.

Batch size curriculum is another effective strategy, gradually increasing batch size during training. Starting with smaller batches provides more frequent parameter updates, which helps the model escape poor local minima in early training. As training progresses and the model approaches convergence, larger batches provide more stable gradient estimates and better utilize hardware parallelism. GPT-3 training used a batch size curriculum, starting at 32,000 tokens per batch and gradually increasing to 3.2 million tokens per batch. This $100\times$ increase in batch size was enabled by learning rate adjustments and gradient accumulation. The larger batch sizes in later training improved hardware utilization from approximately 30% to over 70%, significantly reducing training time.

Learning rate schedules are essential for curriculum learning, as the optimal learning rate changes with batch size and training progress. The linear warmup followed by cosine decay schedule has become standard for transformer pre-training. The warmup phase, typically 1-10% of total steps, gradually increases the learning rate from near-zero to the peak value. This prevents the large gradient updates in early training from destabilizing the model. The cosine decay phase gradually reduces the learning rate to near-zero, allowing the model to fine-tune its parameters as it approaches convergence. For BERT-base, a warmup of 10,000 steps followed by linear decay over 990,000 steps works well. For GPT-3, a warmup of 375 million tokens (approximately 1% of total training) followed by cosine decay proved effective.

The impact on training efficiency is substantial. Curriculum learning can reduce training time by 20-40% compared to fixed configurations while achieving equal or better final performance. For BERT-base, the sequence length curriculum reduces training from approximately 5.5 days to 4 days on the same hardware. For GPT-3, the batch size curriculum improved hardware utilization enough to reduce training time by an estimated 30%, saving approximately \$3-4 million in compute costs. These savings make curriculum learning essential for large-scale pre-training projects.

20.4.2 Progressive Training

Definition 20.1 (Curriculum Learning). Train on progressively harder examples:

Stage 1: Easy examples (short sequences, simple patterns)

Stage 2: Medium difficulty

Stage 3: Full difficulty (long sequences, complex patterns)

Benefits:

- Faster convergence
- Better final performance
- More stable training

Example 20.3 (Sequence Length Curriculum). **GPT-3 training:**

Stage 1 (0-100B tokens):

- Sequence length: 1024
- Batch size: 3.2M tokens

Stage 2 (100B-300B tokens):

- Sequence length: 2048
- Batch size: 3.2M tokens (fewer sequences)

Starting with shorter sequences reduces memory and computation early in training.

20.4.3 Domain-Adaptive Pre-training

Continue pre-training on domain-specific data:

[H] Domain Adaptation

Step 1: Pre-train on general corpus (e.g., Common Crawl)

Step 2: Continue pre-training on domain data (e.g., biomedical)

Step 3: Fine-tune on task

Examples:

- BioBERT: BERT + PubMed/PMC
- SciBERT: BERT + scientific papers
- FinBERT: BERT + financial documents
- CodeBERT: BERT + code

20.5 Hardware Requirements and Infrastructure

20.5.1 BERT-base Training Infrastructure

Training BERT-base requires relatively modest infrastructure by modern standards, making it accessible to academic research groups and small companies. The original BERT paper reported training BERT-base on 16 TPU v3 chips for 4 days, providing a concrete reference point for hardware requirements.

Each TPU v3 chip provides 420 TFLOPS of mixed-precision compute (bfloat16) and 16 GB of high-bandwidth memory (HBM). The 16-chip configuration provides 6.7 PFLOPS total compute and 256 GB total memory. BERT-base with 110 million parameters requires approximately 440 MB for model weights in FP32, or 220 MB in FP16. With batch size 256 and sequence length 512, the activation memory per batch is approximately 8 GB, which fits comfortably in the 256 GB total memory when distributed across 16 chips. The high-bandwidth interconnect between TPU chips (approximately 100 GB/s per chip) enables efficient data parallelism with minimal communication overhead.

The training cost at Google Cloud’s preemptible TPU pricing (approximately \$2.40 per TPU v3 hour) is roughly \$3,700 for the full 4-day training run. Using on-demand pricing (\$8 per hour) would increase this to \$12,000. For comparison, training on NVIDIA V100 GPUs would require approximately 64 GPUs for 4 days (6,144 GPU-hours) at a cost of approximately \$12,000 using AWS on-demand pricing. The equivalent training on A100 GPUs would require approximately 32 GPUs for 2.5 days (1,920 GPU-hours) at a cost of approximately \$6,000, demonstrating the improved efficiency of newer hardware.

The infrastructure requirements beyond compute include high-bandwidth storage for the training data (approximately 100 GB including tokenized sequences and preprocessing artifacts), network bandwidth for distributed training (at least 10 Gbps per GPU for efficient data parallelism), and monitoring infrastructure for tracking training metrics. The total infrastructure cost including storage, networking, and engineering time is typically 2-3× the raw compute cost, bringing the total BERT-base training cost to approximately \$10,000-15,000.

20.5.2 GPT-3 Training Infrastructure

Training GPT-3 175B requires massive infrastructure that is accessible only to large technology companies and well-funded research organizations. The scale of the training job presents significant engineering challenges beyond simply acquiring hardware.

The training used approximately 10,000 NVIDIA V100 GPUs, though the exact configuration has not been publicly disclosed. Each V100 provides 125 TFLOPS of mixed-precision compute and 32 GB of memory. The 10,000-GPU configuration provides 1.25 exaFLOPS total compute and 320 TB total memory. GPT-3 175B with 175 billion parameters requires approximately 700 GB for model weights in FP32, or 350 GB in FP16. With model parallelism across 8 GPUs, each GPU stores approximately 44 GB of model weights, leaving limited memory for activations. The batch size per GPU is typically 1-2 sequences of length 2048, requiring approximately 20 GB of activation memory per GPU.

The communication requirements are severe. With model parallelism across 8 GPUs and data parallelism across 1,250 groups, each training step requires all-reduce operations across the data parallel groups (approximately 700 GB of gradients) and all-to-all communication within model parallel groups (approximately 100 GB per step). At 100 Gbps network bandwidth per GPU, the gradient all-reduce takes approximately 70 seconds per step, which would dominate training time. To address this, GPT-3 training used gradient accumulation (accumulating gradients over multiple micro-batches before synchronizing) and high-bandwidth interconnects like InfiniBand (200 Gbps or higher), reducing communication time to approximately 10% of total step time.

The training cost is estimated between \$4 million and \$12 million depending on assumptions about hardware ownership versus cloud rental, utilization rates, and pricing agreements. At AWS on-demand pricing of \$3 per V100 hour, the 7.2 million GPU-hours would cost \$21.6 million, but OpenAI likely achieved significant discounts through long-term commitments and negotiated pricing. The infrastructure costs beyond raw compute are substantial: high-bandwidth networking equipment (InfiniBand switches and cables) costs millions of dollars, distributed storage systems for the 5 TB of training data cost hundreds of thousands of dollars, and the engineering effort to build and optimize the training pipeline represents millions of dollars in labor costs.

The power consumption is also significant. Each V100 GPU consumes approximately 300 watts under full load, so 10,000 GPUs consume 3 megawatts. Over a 1-month training run, this equals 2,160 megawatt-hours of electricity. At typical data center electricity costs of \$0.10 per kWh, the electricity cost alone is \$216,000. Including cooling and power distribution overhead (typically $1.5\text{--}2\times$ the compute power), the total power cost approaches \$400,000.

20.5.3 LLaMA-65B Training Infrastructure

Training LLaMA-65B represents a more efficient approach than GPT-3, using fewer but more powerful GPUs and a more optimized training pipeline. Meta's paper provides detailed information about the infrastructure and costs.

The training used 2,048 NVIDIA A100-80GB GPUs for 21 days, totaling 1,022,362 GPU-hours. Each A100-80GB provides 312 TFLOPS of mixed-precision compute and 80 GB of memory, representing a significant improvement over V100 ($2.5\times$ compute, $2.5\times$ memory). The 2,048-GPU configuration provides 639 PFLOPS total compute and 164 TB total memory. LLaMA-65B with 65 billion parameters requires approximately 260 GB for model weights in FP32, or 130 GB in FP16. With model parallelism across 8 GPUs, each GPU stores approximately 16 GB of model weights, leaving substantial memory for activations and optimizer states.

The larger memory capacity of A100-80GB enables more efficient training configurations. LLaMA uses a batch size of 4 million tokens (approximately 2,000 sequences of length 2048), distributed across 2,048 GPUs as 2 sequences per GPU. The activation memory per GPU is approximately 40 GB, and the optimizer states (using AdamW) require approximately 32 GB, totaling approximately 88 GB per GPU. This fits comfortably in the 80 GB memory, avoiding the need for activation checkpointing or other memory-saving techniques that would slow training.

The communication requirements are more manageable than GPT-3 due to the smaller model size and more efficient hardware. With model parallelism across 8 GPUs and data parallelism across 256 groups, each training step requires all-reduce operations across data parallel groups (approximately 260

GB of gradients) and all-to-all communication within model parallel groups (approximately 30 GB per step). Using NVIDIA’s NVLink and NVSwitch interconnects (600 GB/s per GPU within a node, 200 Gbps between nodes), the communication time is approximately 5% of total step time, demonstrating excellent scaling efficiency.

The training cost is estimated at \$2-3 million using Meta’s owned infrastructure. At cloud pricing of \$3 per A100 hour, the 1,022,362 GPU-hours would cost \$3.1 million. Meta’s owned infrastructure likely reduced the effective cost to \$1.5-2 million when accounting for hardware depreciation (A100 GPUs cost approximately \$10,000 each, depreciated over 3-5 years) and operational costs. The power consumption is approximately 1 megawatt ($2,048 \text{ GPUs} \times 400 \text{ watts per A100}$), totaling 504 megawatt-hours over 21 days. At \$0.10 per kWh including cooling overhead, the electricity cost is approximately \$75,000.

The infrastructure requirements include high-bandwidth networking (NVIDIA InfiniBand or equivalent), distributed storage systems (approximately 15 TB for training data and checkpoints), and monitoring infrastructure. Meta’s paper notes that they used a custom training framework optimized for their infrastructure, with careful attention to memory management, communication patterns, and fault tolerance. The engineering effort to build this infrastructure and optimize the training pipeline represents a significant investment beyond the raw hardware costs.

20.6 Efficient Pre-training Techniques

20.6.1 Mixed Precision Training

Mixed precision training represents one of the most impactful optimizations for transformer pre-training, leveraging the FP16 (16-bit floating point) format for most computations while maintaining FP32 (32-bit) master weights for numerical stability. This approach reduces activation memory consumption by approximately 50% since activations, which often dominate memory usage during training, are stored in the lower-precision format. The memory savings enable larger batch sizes or longer sequences, directly improving training efficiency.

Modern GPUs like the NVIDIA A100 provide dedicated Tensor Cores that execute FP16 matrix multiplications at twice the throughput of FP32 operations, yielding training speedups of $1.5\text{--}2\times$ in practice. The speedup is less than the theoretical $2\times$ due to memory bandwidth limitations and non-matrix operations that don’t benefit from reduced precision. For BERT-base, mixed precision training reduces memory from 13.8 GB to 8.0 GB (42% reduction) and decreases training time from 4 days to approximately 2.5 days on the same hardware.

However, the reduced precision of FP16 can cause gradient underflow for very small values. To address this, loss scaling is employed: the loss is multiplied by a large factor (typically 1024 or dynamically adjusted) before backpropagation, then gradients are scaled back down before the optimizer step. This effectively shifts the representable range to prevent underflow while maintaining the benefits of reduced precision. Dynamic loss scaling automatically adjusts the scale factor during training, increasing it when no overflow occurs and decreasing it when overflow is detected.

20.6.2 Gradient Checkpointing

Gradient checkpointing trades computation for memory by selectively storing only a subset of activations during the forward pass and recomputing the others during the backward pass. This technique is particularly valuable for training very large models or using very long sequences where activation memory dominates.

For a transformer with L layers, storing all activations requires $O(Lnbd)$ memory where n is sequence length, b is batch size, and d is model dimension. With gradient checkpointing, only activations at checkpoint boundaries are stored (typically every \sqrt{L} layers), reducing memory to $O(\sqrt{L}nbd)$. The recomputation adds approximately 33% to training time but can reduce activation memory by $5\text{--}10\times$, enabling much larger batch sizes or longer sequences that more than compensate for the slowdown.

For GPT-3 training, gradient checkpointing enabled fitting the model on V100 GPUs with 32 GB memory. Without checkpointing, the activation memory for a single sequence of length 2048 would exceed 60 GB, making training impossible. With checkpointing every 8 layers, activation memory

drops to approximately 15 GB per sequence, allowing batch size 2 per GPU. The 33% slowdown from recomputation is offset by the ability to use $2\times$ larger batch size, resulting in net training speedup.

20.6.3 ZeRO Optimizer

The ZeRO (Zero Redundancy Optimizer) technique, developed by Microsoft, eliminates memory redundancy in data-parallel training by partitioning optimizer states, gradients, and parameters across data-parallel processes. Traditional data parallelism replicates the entire model and optimizer state on each GPU, wasting memory. ZeRO partitions these states while maintaining the computational efficiency of data parallelism.

ZeRO has three stages with increasing memory savings and communication overhead. ZeRO Stage 1 partitions only optimizer states, reducing memory by approximately $4\times$ for Adam optimizer (which stores first and second moments). ZeRO Stage 2 additionally partitions gradients, reducing memory by approximately $8\times$. ZeRO Stage 3 partitions parameters as well, reducing memory by the data parallelism degree (e.g., $64\times$ for 64 GPUs). The memory savings enable training much larger models on the same hardware.

For training a 175B parameter model like GPT-3, ZeRO Stage 3 across 64 GPUs reduces per-GPU memory from approximately 700 GB (impossible) to approximately 11 GB (feasible). The communication overhead is manageable because parameter partitions are gathered just before use and discarded after, with careful overlapping of communication and computation. DeepSpeed, Microsoft’s training framework, implements ZeRO and has been used to train models exceeding 1 trillion parameters.

20.6.4 Pipeline Parallelism

Pipeline parallelism divides the model vertically across layers, with different GPUs processing different layers. This enables training models too large to fit on a single GPU while maintaining high hardware utilization through pipelining. The key challenge is keeping all GPUs busy despite the sequential dependency between layers.

The GPipe approach divides each batch into micro-batches and pipelines their execution across stages. While stage 1 processes micro-batch 2, stage 2 processes micro-batch 1, and so on. This achieves high utilization once the pipeline is full, though there are bubble periods at the start and end of each batch. With m micro-batches and p pipeline stages, the bubble overhead is approximately $\frac{p-1}{m}$. Using $m = 4p$ micro-batches reduces bubble overhead to approximately 25%.

For GPT-3 training, pipeline parallelism across 8 stages (12 layers per stage) combined with model parallelism within each stage and data parallelism across pipeline replicas achieved approximately 75% hardware utilization. The pipeline approach enabled training the 175B parameter model on V100 GPUs with 32 GB memory, which would otherwise be impossible. The combination of pipeline parallelism, model parallelism, data parallelism, and ZeRO represents the state-of-the-art for training extremely large models.

20.7 Parameter-Efficient Fine-tuning

20.7.1 Motivation

Full fine-tuning challenges:

- Requires storing full model copy per task
- 175B model \times 100 tasks = 17.5T parameters!
- Expensive and slow

Solution: Fine-tune small subset of parameters.

20.7.2 LoRA: Low-Rank Adaptation

Definition 20.2 (LoRA). Inject trainable low-rank matrices into frozen model:

Original: $\mathbf{h} = \mathbf{W}\mathbf{x}$ where $\mathbf{W} \in \mathbb{R}^{d \times d}$

LoRA:

$$\mathbf{h} = \mathbf{W}\mathbf{x} + \Delta\mathbf{W}\mathbf{x} = \mathbf{W}\mathbf{x} + \mathbf{B}\mathbf{A}\mathbf{x} \quad (20.15)$$

where $\mathbf{A} \in \mathbb{R}^{r \times d}$, $\mathbf{B} \in \mathbb{R}^{d \times r}$, and $r \ll d$ (typically $r = 4$ to 64).

Parameters:

- Original: d^2 (frozen)
- LoRA: $2rd$ (trainable)
- Reduction: $\frac{2rd}{d^2} = \frac{2r}{d}$

Example 20.4 (LoRA for GPT-3). GPT-3 175B, apply LoRA with $r = 8$ to attention projections.

Single attention layer:

- $\mathbf{W}^Q, \mathbf{W}^K, \mathbf{W}^V, \mathbf{W}^O \in \mathbb{R}^{12288 \times 12288}$
- Original params: $4 \times 12288^2 = 604M$

LoRA params per layer:

$$4 \times 2 \times 8 \times 12288 = 786,432 \approx 0.79M \quad (20.16)$$

96 layers total:

- LoRA params: $96 \times 0.79M = 75.8M$
- Full model: 175B
- **Reduction: 2,300×** (train only 0.04% of parameters!)

Performance: Matches full fine-tuning on many tasks!

20.7.3 Adapter Layers

Definition 20.3 (Adapter). Insert small bottleneck layers between frozen layers:

$$\mathbf{h}_{\text{adapter}} = \mathbf{h} + \text{FFN}_{\text{adapter}}(\text{LayerNorm}(\mathbf{h})) \quad (20.17)$$

where $\text{FFN}_{\text{adapter}}: d \rightarrow d_{\text{bottleneck}} \rightarrow d$ with $d_{\text{bottleneck}} \ll d$.

Typical bottleneck: $d_{\text{bottleneck}} = 64$ for $d = 768$

Parameters per adapter:

$$2d \cdot d_{\text{bottleneck}} = 2 \times 768 \times 64 = 98,304 \quad (20.18)$$

20.7.4 Prompt Tuning

Definition 20.4 (Prompt Tuning). Prepend learnable "soft prompt" vectors:

Input: $[1, \dots, k, \mathbf{x}_1, \dots, \mathbf{x}_n]$

where $i \in \mathbb{R}^d$ are learned continuous prompts (not discrete tokens).

Parameters: Only $k \times d$ prompt vectors (model frozen).

Typical: $k = 20$ prompts, $d = 768 \rightarrow$ only 15,360 parameters!

20.8 Multi-Task and Multi-Stage Pre-training

20.8.1 Multi-Task Pre-training

Train on multiple objectives simultaneously:

$$\mathcal{L}_{\text{total}} = \sum_{i=1}^K \lambda_i \mathcal{L}_i \quad (20.19)$$

Example (T5):

- Span corruption (main)
- Prefix LM
- Deshuffling

Benefits:

- More robust representations
- Better transfer to diverse tasks
- Can balance objectives with λ_i

20.8.2 Multi-Stage Pre-training

Stage 1: General pre-training

- Large diverse corpus
- Language modeling
- Build general knowledge

Stage 2: Instruction tuning

- Instruction-response pairs
- Learn to follow instructions
- Improve helpfulness

Stage 3: RLHF

- Reinforcement learning from human feedback
- Align with human preferences
- Improve safety

Example 20.5 (InstructGPT Pipeline). **Stage 1:** GPT-3 pre-training (175B params, 300B tokens)

Stage 2: Supervised fine-tuning

- 13,000 instruction-output examples
- Fine-tune for 16 epochs
- Learning rate: 9.65×10^{-6}

Stage 3: Reward modeling

- 33,000 comparison examples
- Train 6B reward model
- Predicts human preferences

Stage 4: PPO optimization

- 31,000 prompts
- Optimize policy to maximize reward
- KL penalty from SFT model

Result: 1.3B InstructGPT preferred over 175B GPT-3 by humans!

20.9 Transfer Learning Analysis

20.9.1 Measuring Transfer

Metrics:

1. Downstream Performance:

$$\Delta = \text{Performance}_{\text{fine-tuned}} - \text{Performance}_{\text{from-scratch}} \quad (20.20)$$

2. Sample Efficiency:

- Number of examples to reach target performance
- Pre-trained models: 10-100× fewer examples

3. Convergence Speed:

- Training steps to convergence
- Pre-trained: 10× faster

20.9.2 What Makes Good Pre-training?

Data scale: More data → better transfer (up to a point)

Data diversity: Diverse pre-training → broader transfer

Model scale: Larger models transfer better

Objective alignment: Pre-training objective similar to downstream task

Domain match: Domain-specific pre-training helps domain-specific tasks

20.10 Exercises

Exercise 20.1. Compare pre-training objectives:

1. Train BERT-tiny with: (a) MLM, (b) CLM, (c) Span corruption
2. Evaluate on GLUE tasks
3. Which objective transfers best? Why?

Exercise 20.2. Implement data filtering pipeline:

1. Download 10,000 documents from Common Crawl
2. Remove duplicates (exact and near-duplicate)
3. Filter by language (keep English)
4. Filter low-quality (perplexity \geq threshold)
5. Report statistics at each stage

Exercise 20.3. Implement LoRA:

1. Load pre-trained GPT-2
2. Add LoRA layers with $r = 8$ to attention
3. Fine-tune on sentiment analysis
4. Compare: (a) Full fine-tuning, (b) LoRA, (c) Frozen
5. Measure: parameters trained, memory, accuracy

Exercise 20.4. Analyze transfer learning:

1. Fine-tune BERT on 5 GLUE tasks
2. Vary training data: [100, 500, 1000, 5000, all]
3. Compare to training from scratch
4. Plot sample efficiency curves
5. At what point does pre-training stop helping?

Part VII

Practical Implementation

Chapter 21

Implementing Transformers in PyTorch

Chapter Overview

This chapter provides complete, production-ready PyTorch implementations of transformer models. We build from scratch: attention mechanisms, encoder/decoder blocks, position encodings, and full models (BERT, GPT, T5). Each implementation includes training loops, optimization, and best practices.

Learning Objectives

1. Implement multi-head attention from scratch
2. Build transformer encoder and decoder blocks
3. Create complete BERT and GPT models
4. Write efficient training loops with mixed precision
5. Apply gradient accumulation and checkpointing
6. Debug common implementation issues

21.1 Multi-Head Attention Implementation

21.1.1 Core Components

The implementation of multi-head attention in PyTorch requires careful attention to efficiency and memory usage. The standard approach involves projecting queries, keys, and values through linear layers, reshaping tensors to separate attention heads, computing scaled dot-product attention, and finally concatenating the results. However, several optimizations can significantly improve both speed and memory efficiency.

Key implementation considerations:

- Efficient batched matrix multiplications
- Proper dimension handling for multi-head split
- Memory-efficient attention computation
- Gradient flow through softmax

PyTorch multi-head attention structure:

1. Project Q, K, V: Linear layers
2. Reshape for multiple heads: view + transpose
3. Compute scaled dot-product attention
4. Concatenate heads and project output

21.1.2 Memory-Efficient Attention

The standard attention mechanism computes the full attention matrix of size (B, h, n, n) , which becomes prohibitively expensive for long sequences. For a sequence length of 512 with 12 heads and batch size 32, this requires approximately 400MB just for the attention scores. We can implement several optimizations to reduce this memory footprint.

The first optimization involves computing attention in chunks rather than materializing the entire attention matrix at once. This approach processes the attention computation in blocks, trading some computational efficiency for substantial memory savings. The chunked attention implementation divides the sequence into smaller segments and computes attention scores for each segment independently.

```
def memory_efficient_attention(Q, K, V, chunk_size=128):
    """
    Compute attention in chunks to reduce memory usage.
    Q, K, V: (batch, heads, seq_len, head_dim)
    """
    B, h, n, d = Q.shape
    output = torch.zeros_like(Q)

    for i in range(0, n, chunk_size):
        end_i = min(i + chunk_size, n)
        Q_chunk = Q[:, :, i:end_i, :] # (B, h, chunk, d)

        # Compute attention scores for this chunk
        scores = torch.matmul(Q_chunk, K.transpose(-2, -1))
        scores = scores / math.sqrt(d)
        attn = F.softmax(scores, dim=-1)

        # Apply to values
        output[:, :, i:end_i, :] = torch.matmul(attn, V)

    return output
```

Another critical optimization is the use of PyTorch's scaled dot-product attention function introduced in PyTorch 2.0, which implements Flash Attention algorithms internally. This function provides significant speedups and memory reductions through kernel fusion and optimized memory access patterns.

```
import torch.nn.functional as F

def efficient_attention(Q, K, V, mask=None):
    """
    Use PyTorch's optimized scaled_dot_product_attention.
    Automatically uses Flash Attention when available.
    """
    # PyTorch 2.0+ provides optimized implementation
    output = F.scaled_dot_product_attention(
        Q, K, V,
        attn_mask=mask,
        dropout_p=0.1 if self.training else 0.0,
        is_causal=False
    )
    return output
```

This optimized implementation can reduce memory usage by up to 50% and provide 2-3× speedups compared to the naive implementation, particularly for longer sequences.

21.1.3 Dimension Tracking Example

For BERT-base configuration ($d = 768$, $h = 12$):

Input: $(B, n, 768)$ where B is batch size, n is sequence length

After Q/K/V projection: $(B, n, 768)$

Reshape for heads: $(B, n, 12, 64) \rightarrow (B, 12, n, 64)$

Attention scores: $(B, 12, n, n)$

After applying to V: $(B, 12, n, 64)$

Concatenate heads: $(B, n, 12, 64) \rightarrow (B, n, 768)$

Output projection: $(B, n, 768)$

21.2 Position Encodings

21.2.1 Sinusoidal Encoding

Mathematical formula:

$$PE_{(\text{pos}, 2i)} = \sin\left(\frac{\text{pos}}{10000^{2i/d}}\right) \quad (21.1)$$

$$PE_{(\text{pos}, 2i+1)} = \cos\left(\frac{\text{pos}}{10000^{2i/d}}\right) \quad (21.2)$$

Implementation strategy:

1. Pre-compute position encoding matrix at initialization
2. Register as buffer (not parameter, doesn't need gradients)
3. Add to embeddings in forward pass

21.2.2 Learned Positional Embeddings

Alternative approach (BERT):

- Embedding layer: $\text{max_len} \times \text{d_model}$
- Learn position representations during training
- More flexible but requires fixed max length

21.3 Masking Strategies

21.3.1 Causal Mask for GPT

Lower triangular mask:

$$M_{ij} = \begin{cases} 1 & \text{if } j \leq i \\ 0 & \text{if } j > i \end{cases} \quad (21.3)$$

Implementation:

```
causal_mask = torch.tril(torch.ones(seq_len, seq_len))
scores = scores.masked_fill(mask == 0, -1e9)
```

21.3.2 Padding Mask

For variable-length sequences:

```
# input_ids: (batch, seq_len)
# 0 indicates padding
pad_mask = (input_ids != 0).unsqueeze(1).unsqueeze(2)
# Shape: (batch, 1, 1, seq_len)
```

21.4 Training Optimizations

21.4.1 Fused Kernels for Layer Normalization

Standard PyTorch operations like layer normalization involve multiple kernel launches, each reading from and writing to global memory. Fusing these operations into a single kernel can provide substantial speedups by reducing memory bandwidth requirements. Modern deep learning frameworks provide fused implementations of common operations that should be used whenever possible.

The Apex library from NVIDIA provides highly optimized fused kernels for layer normalization and other operations. These implementations can be 2-3 \times faster than the standard PyTorch versions, particularly for smaller batch sizes where kernel launch overhead dominates.

```
# Standard PyTorch layer norm
layer_norm = nn.LayerNorm(d_model)

# Fused layer norm from Apex (faster)
try:
    from apex.normalization import FusedLayerNorm
    layer_norm = FusedLayerNorm(d_model)
except ImportError:
    # Fall back to standard implementation
    layer_norm = nn.LayerNorm(d_model)
```

Similarly, fused dropout and bias addition can be combined with other operations to reduce memory traffic. The key principle is to minimize the number of separate kernel launches and memory accesses by combining operations that naturally occur together in the computation graph.

21.4.2 Mixed Precision Training

Mixed precision training uses 16-bit floating point (FP16) for most operations while maintaining 32-bit precision for critical computations. This approach provides substantial benefits in terms of both memory usage and computational speed, particularly on modern GPUs with dedicated tensor cores optimized for FP16 operations.

Benefits:

- 2 \times memory reduction for activations and gradients
- 2-3 \times training speedup on modern GPUs with tensor cores
- Same final accuracy with proper loss scaling
- Enables training larger models or using larger batch sizes

PyTorch provides automatic mixed precision (AMP) through the `torch.cuda.amp` module, which automatically handles the conversion between FP16 and FP32 as needed. The implementation requires minimal code changes and provides automatic loss scaling to prevent gradient underflow.

```
from torch.cuda.amp import autocast, GradScaler

# Initialize gradient scaler for loss scaling
scaler = GradScaler()

# Training loop
for batch in dataloader:
    optimizer.zero_grad()

    # Forward pass in mixed precision
```

```

with autocast():
    outputs = model(batch['input_ids'])
    loss = criterion(outputs, batch['labels'])

# Backward pass with scaled loss
scaler.scale(loss).backward()

# Unscale gradients and clip
scaler.unscale_(optimizer)
torch.nn.utils.clip_grad_norm_(model.parameters(), max_norm=1.0)

# Update weights
scaler.step(optimizer)
scaler.update()

```

The gradient scaler automatically adjusts the loss scaling factor to maintain numerical stability. It increases the scale when no overflow is detected and decreases it when overflow occurs, ensuring that gradients remain in a representable range for FP16 arithmetic.

For BERT-base training, mixed precision typically reduces memory usage from approximately 16GB to 8GB per GPU while maintaining the same final accuracy. The speedup varies depending on the GPU architecture, with Volta and newer architectures providing the largest benefits due to their tensor cores.

21.4.3 Gradient Accumulation

Gradient accumulation enables training with effective batch sizes larger than what fits in GPU memory by accumulating gradients over multiple forward-backward passes before updating weights. This technique is essential for training large models or when hardware constraints limit the physical batch size.

Purpose: Simulate large batch sizes on limited memory

Effective batch size:

$$B_{\text{effective}} = B_{\text{physical}} \times N_{\text{accumulation}} \quad (21.4)$$

The implementation requires careful handling of gradient normalization to ensure that the accumulated gradients have the correct scale. Each loss value should be divided by the number of accumulation steps so that the final gradient magnitude matches what would be obtained with a single large batch.

```

accumulation_steps = 8
optimizer.zero_grad()

for i, batch in enumerate(dataloader):
    # Forward pass
    with autocast():
        outputs = model(batch['input_ids'])
        loss = criterion(outputs, batch['labels'])
        # Scale loss by accumulation steps
        loss = loss / accumulation_steps

    # Backward pass
    scaler.scale(loss).backward()

    # Update weights every accumulation_steps
    if (i + 1) % accumulation_steps == 0:
        scaler.unscale_(optimizer)
        torch.nn.utils.clip_grad_norm_(model.parameters(), 1.0)
        scaler.step(optimizer)

```

```

scaler.update()
optimizer.zero_grad()

```

Example:

- Physical batch: 32 (fits in GPU)
- Accumulation steps: 8
- Effective batch: 256

This approach allows training with large effective batch sizes that would otherwise require multiple GPUs or be impossible due to memory constraints. The trade-off is increased training time, as the optimizer updates occur less frequently.

21.4.4 Gradient Checkpointing

Gradient checkpointing trades computation for memory by selectively storing only a subset of activations during the forward pass and recomputing the others during the backward pass. This technique can dramatically reduce memory usage, enabling training of much larger models or longer sequences at the cost of increased computation time.

Trade computation for memory:

- Don't store all activations during forward pass
- Recompute intermediate activations during backward pass
- Enables larger models or longer sequences
- Typically 20-30% slower but saves 40-50% memory

PyTorch provides gradient checkpointing through the `torch.utils.checkpoint` module. The key is to wrap transformer layers or blocks in checkpoint functions, which handle the recomputation automatically during the backward pass.

```

from torch.utils.checkpoint import checkpoint

class TransformerLayer(nn.Module):
    def __init__(self, d_model, num_heads, d_ff, use_checkpoint=False):
        super().__init__()
        self.use_checkpoint = use_checkpoint
        self.attention = MultiHeadAttention(d_model, num_heads)
        self.ffn = FeedForward(d_model, d_ff)
        self.norm1 = nn.LayerNorm(d_model)
        self.norm2 = nn.LayerNorm(d_model)

    def forward(self, x, mask=None):
        if self.use_checkpoint and self.training:
            # Use gradient checkpointing
            return checkpoint(self._forward, x, mask)
        else:
            return self._forward(x, mask)

    def _forward(self, x, mask):
        # Attention block
        attn_out = self.attention(x, x, x, mask)
        x = self.norm1(x + attn_out)

```



```
# Feed-forward block
ffn_out = self.ffn(x)
x = self.norm2(x + ffn_out)
return x
```

For a 12-layer BERT model, gradient checkpointing can reduce peak memory usage from approximately 16GB to 9GB, allowing training with longer sequences or larger batch sizes. The computational overhead is typically 20-30%, which is often an acceptable trade-off for the memory savings.

21.5 Model Initialization

21.5.1 Best Practices

Weight initialization:

- Linear layers: Xavier/Glorot normal
- Embeddings: Normal(0, 0.02)
- Layer norm: gamma=1, beta=0

Special considerations:

- Scale residual connections by $1/\sqrt{N_{\text{layers}}}$
- Weight tying: LM head shares embeddings
- Careful initialization prevents gradient issues

21.6 Memory Profiling and Optimization

21.6.1 Understanding Memory Usage

Memory consumption in transformer training comes from several sources: model parameters, optimizer states, activations, gradients, and temporary buffers. Understanding the breakdown of memory usage is essential for effective optimization. For a BERT-base model with 110M parameters, the memory requirements can be substantial even before considering batch data.

The model parameters themselves occupy relatively little memory compared to other components. With 110M parameters stored in FP32, the parameters require approximately 440MB. However, the Adam optimizer maintains two additional states per parameter (first and second moments), tripling the parameter memory to 1.3GB. Activations stored during the forward pass for gradient computation typically consume the largest portion of memory, scaling with both sequence length and batch size.

PyTorch provides comprehensive memory profiling tools through the `torch.cuda` module. The `memory_summary` function provides detailed information about current memory allocation, including cached memory, allocated memory, and peak memory usage.

```
import torch

# Profile memory usage during training
def profile_memory(model, batch_size, seq_len, device='cuda'):
    torch.cuda.reset_peak_memory_stats(device)
    torch.cuda.empty_cache()

    # Create sample input
    input_ids = torch.randint(0, 30000, (batch_size, seq_len)).to(device)

    # Forward pass
```

```

outputs = model(input_ids)
loss = outputs.mean()

print(f"After forward pass:")
print(f"Allocated: {torch.cuda.memory_allocated(device) / 1e9:.2f} GB")
print(f"Reserved: {torch.cuda.memory_reserved(device) / 1e9:.2f} GB")

# Backward pass
loss.backward()

print(f"\nAfter backward pass:")
print(f"Allocated: {torch.cuda.memory_allocated(device) / 1e9:.2f} GB")
print(f"Peak: {torch.cuda.max_memory_allocated(device) / 1e9:.2f} GB")

# Detailed summary
print("\nDetailed memory summary:")
print(torch.cuda.memory_summary(device))

```

21.6.2 Identifying Memory Bottlenecks

The first step in optimization is identifying where memory is being consumed. For transformer models, the attention mechanism typically dominates memory usage due to the quadratic scaling of attention scores with sequence length. A single attention layer with sequence length 512 and 12 heads requires approximately 12MB for attention scores alone, and this scales quadratically with sequence length.

Activation memory can be estimated using the formula:

$$M_{\text{activations}} \approx 2 \times B \times n \times d \times L \times \text{bytes_per_element} \quad (21.5)$$

where B is batch size, n is sequence length, d is model dimension, L is number of layers, and the factor of 2 accounts for both attention and feed-forward activations. For BERT-base with batch size 32 and sequence length 512, this amounts to approximately 9GB in FP32 or 4.5GB in FP16.

21.6.3 Optimization Strategies

Several strategies can dramatically reduce memory consumption. The most effective approach combines multiple techniques tailored to the specific bottlenecks identified through profiling.

Strategy 1: Reduce Sequence Length

The quadratic scaling of attention with sequence length makes this the most impactful optimization for long sequences. Reducing sequence length from 512 to 256 reduces attention memory by $4\times$ and total activation memory by $2\times$. When possible, use techniques like sliding windows or hierarchical attention to process longer documents without materializing full attention matrices.

Strategy 2: Optimize Batch Size

Finding the optimal batch size requires balancing memory usage with computational efficiency. Larger batches improve GPU utilization but consume more memory. Use gradient accumulation to achieve large effective batch sizes while keeping physical batch sizes manageable.

```

def find_optimal_batch_size(model, seq_len, device='cuda'):
    """Binary search to find maximum batch size that fits in memory."""
    min_batch = 1
    max_batch = 256
    optimal_batch = 1

    while min_batch <= max_batch:
        batch_size = (min_batch + max_batch) // 2
        torch.cuda.empty_cache()

```

```

try:
    # Test if this batch size fits
    input_ids = torch.randint(0, 30000,
                              (batch_size, seq_len)).to(device)
    outputs = model(input_ids)
    loss = outputs.mean()
    loss.backward()

    optimal_batch = batch_size
    min_batch = batch_size + 1
except RuntimeError as e:
    if "out of memory" in str(e):
        max_batch = batch_size - 1
    else:
        raise e

return optimal_batch

```

Strategy 3: Layer-wise Optimization

Different layers have different memory characteristics. Attention layers consume more memory than feed-forward layers due to the attention score matrix. Applying gradient checkpointing selectively to attention layers can provide most of the memory benefits with less computational overhead than checkpointing all layers.

21.6.4 Case Study: Optimizing BERT-base

Consider optimizing BERT-base training to reduce memory from 16GB to 8GB while maintaining training throughput. The baseline configuration uses batch size 32, sequence length 512, and FP32 precision.

Baseline measurements:

- Memory usage: 16.2 GB
- Training speed: 120 samples/second
- Parameters: 440 MB
- Optimizer states: 880 MB
- Activations: 9.1 GB
- Gradients: 4.5 GB
- Other: 1.3 GB

Optimization steps:

First, enable mixed precision training. This immediately reduces activation and gradient memory by 50%, bringing total memory to approximately 10GB. The training speed increases to 280 samples/second due to tensor core utilization.

Second, apply gradient checkpointing to all transformer layers. This reduces activation memory by an additional 40%, bringing total memory to 7.8GB. Training speed decreases to 220 samples/second due to recomputation overhead.

Third, optimize the batch size. With the memory savings, we can increase batch size to 48, improving GPU utilization. Final measurements show 7.9GB memory usage and 310 samples/second throughput.

Final measurements:

- Memory usage: 7.9 GB (51% reduction)
- Training speed: 310 samples/second ($2.6\times$ faster)
- Same final accuracy after convergence

This optimization demonstrates that combining multiple techniques can achieve substantial improvements in both memory efficiency and training speed without sacrificing model quality.

21.7 Debugging Transformers

21.7.1 Common Issues

1. Dimension mismatches:

- Check shapes at each operation
- Use assertions for critical dimensions
- Print intermediate tensor shapes

2. NaN/Inf in training:

- Too high learning rate
- Gradient explosion (add clipping)
- Numerical instability in softmax (check mask values)

3. Slow convergence:

- Insufficient warmup
- Bad initialization
- Learning rate too low

4. Memory issues:

- Reduce batch size
- Use gradient checkpointing
- Enable mixed precision
- Reduce sequence length

21.7.2 Validation Checks

Sanity checks before full training:

1. Overfit single batch (should reach near-zero loss)
2. Check gradient norms are reasonable
3. Verify attention weights sum to 1
4. Test with different sequence lengths
5. Profile memory usage

21.8 Inference Optimization

21.8.1 TorchScript Compilation

Inference optimization is critical for deploying transformer models in production environments where latency and throughput requirements are stringent. TorchScript provides a way to serialize and optimize PyTorch models for inference, removing Python overhead and enabling additional optimizations.

The `torch.jit.script` function traces the model's execution and converts it to an intermediate representation that can be optimized and executed more efficiently. This process eliminates Python interpreter overhead and enables fusion of operations that would otherwise require multiple kernel launches.

```
import torch.jit as jit

class TransformerForInference(nn.Module):
    def __init__(self, config):
        super().__init__()
        self.transformer = TransformerModel(config)

    def forward(self, input_ids: torch.Tensor) -> torch.Tensor:
        # Type annotations required for TorchScript
        return self.transformer(input_ids)

# Create and script the model
model = TransformerForInference(config)
model.eval()

# Convert to TorchScript
scripted_model = jit.script(model)

# Save for deployment
scripted_model.save('model_scripted.pt')

# Load and use
loaded_model = jit.load('model_scripted.pt')
with torch.no_grad():
    output = loaded_model(input_ids)
```

TorchScript compilation typically provides 10-30% speedup for transformer inference, with larger models seeing greater benefits. The compilation process also validates that the model can be executed without Python dependencies, which is essential for deployment in production environments.

21.8.2 KV Cache for Autoregressive Generation

Autoregressive generation in models like GPT requires computing attention over all previous tokens at each step. Without optimization, this results in redundant computation as the keys and values for previous tokens are recomputed at every step. Implementing a KV cache stores these values and reuses them, dramatically reducing computation.

```
class GPTWithKVCache(nn.Module):
    def __init__(self, config):
        super().__init__()
        self.layers = nn.ModuleList([
            TransformerLayer(config) for _ in range(config.num_layers)
        ])
        self.embed = nn.Embedding(config.vocab_size, config.d_model)
```

```

def forward(self, input_ids, past_key_values=None, use_cache=False):
    """
    Args:
        input_ids: (batch, seq_len) - new tokens to process
        past_key_values: List of (key, value) tuples from previous steps
        use_cache: Whether to return key-value cache
    """
    batch_size, seq_len = input_ids.shape

    # Embed input tokens
    hidden_states = self.embed(input_ids)

    # Initialize cache if not provided
    if past_key_values is None:
        past_key_values = [None] * len(self.layers)

    # Store new key-values if caching
    present_key_values = [] if use_cache else None

    # Process through layers
    for i, (layer, past_kv) in enumerate(
        zip(self.layers, past_key_values)):

        # Layer forward with cache
        hidden_states, new_kv = layer(
            hidden_states,
            past_key_value=past_kv,
            use_cache=use_cache
        )

        if use_cache:
            present_key_values.append(new_kv)

    return hidden_states, present_key_values

class TransformerLayerWithCache(nn.Module):
    def forward(self, x, past_key_value=None, use_cache=False):
        # Compute Q, K, V
        Q = self.q_proj(x)
        K = self.k_proj(x)
        V = self.v_proj(x)

        # Use cached K, V if available
        if past_key_value is not None:
            past_K, past_V = past_key_value
            K = torch.cat([past_K, K], dim=1)
            V = torch.cat([past_V, V], dim=1)

        # Compute attention
        attn_output = self.attention(Q, K, V)

        # Return new cache if requested
        new_kv = (K, V) if use_cache else None
        return attn_output, new_kv

```

```

# Generation with KV cache
def generate_with_cache(model, input_ids, max_length=100):
    """Generate tokens using KV cache for efficiency."""
    past_key_values = None

    for _ in range(max_length):
        # Only process new token (or all tokens on first step)
        if past_key_values is None:
            current_input = input_ids
        else:
            current_input = input_ids[:, -1:]

        # Forward pass with cache
        logits, past_key_values = model(
            current_input,
            past_key_values=past_key_values,
            use_cache=True
        )

        # Sample next token
        next_token = torch.argmax(logits[:, -1, :], dim=-1, keepdim=True)
        input_ids = torch.cat([input_ids, next_token], dim=1)

        # Stop if EOS token
        if next_token.item() == eos_token_id:
            break

    return input_ids

```

KV caching reduces the computational complexity of generating n tokens from $O(n^2)$ to $O(n)$, providing speedups of 5-10 \times for typical generation lengths. The memory overhead is proportional to the sequence length and number of layers, typically requiring 1-2GB for a GPT-2 sized model generating 1000 tokens.

21.8.3 ONNX Export

ONNX (Open Neural Network Exchange) provides a standardized format for representing neural networks, enabling deployment across different frameworks and hardware platforms. Exporting to ONNX allows using optimized inference engines like ONNX Runtime, which can provide substantial speedups.

```

import torch.onnx

def export_to_onnx(model, output_path, batch_size=1, seq_len=128):
    """Export PyTorch model to ONNX format."""
    model.eval()

    # Create dummy input
    dummy_input = torch.randint(
        0, model.config.vocab_size,
        (batch_size, seq_len)
    )

    # Export to ONNX
    torch.onnx.export(

```

```

        model,
        dummy_input,
        output_path,
        input_names=['input_ids'],
        output_names=['logits'],
        dynamic_axes={
            'input_ids': {0: 'batch', 1: 'sequence'},
            'logits': {0: 'batch', 1: 'sequence'}
        },
        opset_version=14,
        do_constant_folding=True
    )

# Use ONNX Runtime for inference
import onnxruntime as ort

# Create inference session
session = ort.InferenceSession(
    'model.onnx',
    providers=['CUDAExecutionProvider', 'CPUExecutionProvider']
)

# Run inference
input_ids = torch.randint(0, 30000, (1, 128)).numpy()
outputs = session.run(
    ['logits'],
    {'input_ids': input_ids}
)

```

ONNX Runtime typically provides 20-40% speedup over PyTorch for transformer inference, with the exact speedup depending on the model architecture and hardware. The optimization includes operator fusion, memory layout optimization, and hardware-specific kernel selection.

21.8.4 TensorRT Optimization

NVIDIA TensorRT provides highly optimized inference for NVIDIA GPUs through aggressive kernel fusion, precision calibration, and dynamic tensor memory management. TensorRT can provide 2-5× speedup over standard PyTorch inference for transformer models.

```

# Convert ONNX to TensorRT
import tensorrt as trt

def build_tensorrt_engine(onnx_path, engine_path, fp16_mode=True):
    """Build TensorRT engine from ONNX model."""
    logger = trt.Logger(trt.Logger.WARNING)
    builder = trt.Builder(logger)
    network = builder.create_network(
        1 << int(trt.NetworkDefinitionCreationFlag.EXPLICIT_BATCH)
    )
    parser = trt.OnnxParser(network, logger)

    # Parse ONNX model
    with open(onnx_path, 'rb') as f:
        if not parser.parse(f.read()):
            for error in range(parser.num_errors):

```



```

        print(parser.get_error(error))
    return None

# Configure builder
config = builder.create_builder_config()
config.max_workspace_size = 4 << 30 # 4GB

if fp16_mode:
    config.set_flag(trt.BuilderFlag.FP16)

# Build engine
engine = builder.build_engine(network, config)

# Save engine
with open(engine_path, 'wb') as f:
    f.write(engine.serialize())

return engine

# Use TensorRT for inference
def infer_with_tensorrt(engine_path, input_ids):
    """Run inference using TensorRT engine."""
    logger = trt.Logger(trt.Logger.WARNING)

    with open(engine_path, 'rb') as f:
        runtime = trt.Runtime(logger)
        engine = runtime.deserialize_cuda_engine(f.read())

    context = engine.create_execution_context()

    # Allocate buffers
    inputs, outputs, bindings = allocate_buffers(engine)

    # Copy input data
    inputs[0].host = input_ids.cpu().numpy()

    # Run inference
    outputs = do_inference(
        context, bindings, inputs, outputs, stream
    )

    return outputs[0]

```

TensorRT optimization is particularly effective for deployment scenarios where inference latency is critical. The optimization process includes layer fusion, precision calibration for INT8 quantization, and kernel auto-tuning for the specific GPU architecture.

21.8.5 Quantization

Quantization reduces model size and inference latency by using lower precision representations for weights and activations. PyTorch supports several quantization approaches, from simple dynamic quantization to full quantization-aware training.

```
import torch.quantization as quant
```

```

# Dynamic quantization (easiest, good for LSTM/Transformer)
def dynamic_quantize(model):
    """Apply dynamic quantization to linear layers."""
    quantized_model = quant.quantize_dynamic(
        model,
        {nn.Linear}, # Quantize linear layers
        dtype=torch.qint8
    )
    return quantized_model

# Static quantization (requires calibration)
def static_quantize(model, calibration_dataloader):
    """Apply static quantization with calibration."""
    model.eval()

    # Specify quantization configuration
    model.qconfig = quant.get_default_qconfig('fbgemm')

    # Prepare model for quantization
    model_prepared = quant.prepare(model)

    # Calibrate with representative data
    with torch.no_grad():
        for batch in calibration_dataloader:
            model_prepared(batch['input_ids'])

    # Convert to quantized model
    model_quantized = quant.convert(model_prepared)
    return model_quantized

# Quantization-aware training
def quantization_aware_training(model, train_dataloader):
    """Train model with quantization simulation."""
    model.train()
    model.qconfig = quant.get_default_qat_qconfig('fbgemm')

    # Prepare for QAT
    model_prepared = quant.prepare_qat(model)

    # Train normally
    for epoch in range(num_epochs):
        for batch in train_dataloader:
            outputs = model_prepared(batch['input_ids'])
            loss = criterion(outputs, batch['labels'])
            loss.backward()
            optimizer.step()

    # Convert to quantized model
    model_quantized = quant.convert(model_prepared.eval())
    return model_quantized

```

Dynamic quantization typically provides $2\text{--}3\times$ speedup and $4\times$ model size reduction with minimal accuracy loss for transformer models. Static quantization and quantization-aware training can provide additional benefits but require more careful tuning and calibration data.

21.8.6 Inference Benchmarking

Comprehensive benchmarking is essential for understanding the trade-offs between different optimization techniques. The following framework measures latency, throughput, and memory usage across different configurations.

```
def benchmark_inference(model, batch_sizes, seq_lengths, num_runs=100):
    """Comprehensive inference benchmarking."""
    results = []
    model.eval()

    for batch_size in batch_sizes:
        for seq_len in seq_lengths:
            # Create input
            input_ids = torch.randint(
                0, 30000, (batch_size, seq_len)
            ).cuda()

            # Warmup
            with torch.no_grad():
                for _ in range(10):
                    _ = model(input_ids)

            # Benchmark
            torch.cuda.synchronize()
            start = time.time()

            with torch.no_grad():
                for _ in range(num_runs):
                    _ = model(input_ids)

            torch.cuda.synchronize()
            elapsed = time.time() - start

            # Calculate metrics
            latency_ms = (elapsed / num_runs) * 1000
            throughput = (batch_size * num_runs) / elapsed
            memory_mb = torch.cuda.max_memory_allocated() / 1e6

            results.append({
                'batch_size': batch_size,
                'seq_len': seq_len,
                'latency_ms': latency_ms,
                'throughput': throughput,
                'memory_mb': memory_mb
            })

            torch.cuda.reset_peak_memory_stats()

    return results

# Compare optimization techniques
def compare_optimizations(base_model, config):
    """Compare different optimization approaches."""
    models = {
```

```

    'baseline': base_model,
    'torchscript': jit.script(base_model),
    'quantized': dynamic_quantize(base_model),
    'fp16': base_model.half()
}

results = {}
for name, model in models.items():
    print(f"Benchmarking {name}...")
    results[name] = benchmark_inference(
        model,
        batch_sizes=[1, 8, 32],
        seq_lengths=[128, 512]
    )

return results

```

Typical results for BERT-base inference optimization show that combining TorchScript, FP16, and dynamic quantization can achieve $5\text{--}8\times$ speedup with less than 1% accuracy degradation, making deployment feasible for latency-sensitive applications.

21.9 Complete Training Pipeline

21.9.1 Training Script Structure

1. Configuration:

```

config = {
    'd_model': 768,
    'num_heads': 12,
    'num_layers': 12,
    'd_ff': 3072,
    'vocab_size': 30000,
    'max_seq_len': 512,
    'dropout': 0.1,
    'batch_size': 32,
    'learning_rate': 1e-4,
    'warmup_steps': 10000,
    'max_steps': 1000000
}

```

2. Model instantiation:

```

model = BERTModel(**config)
model = model.to(device)

```

3. Optimizer setup:

```

from torch.optim import AdamW
optimizer = AdamW(
    model.parameters(),
    lr=config['learning_rate'],
    betas=(0.9, 0.999),
    eps=1e-6,
    weight_decay=0.01
)

```

4. Learning rate scheduler:

```

from transformers import get_linear_schedule_with_warmup
scheduler = get_linear_schedule_with_warmup(
    optimizer,
    num_warmup_steps=config['warmup_steps'],
    num_training_steps=config['max_steps']
)

```

5. Training loop with all optimizations:

- Mixed precision
- Gradient accumulation
- Gradient clipping
- Checkpointing
- Logging

21.10 Complete Implementation Examples

21.10.1 Full Transformer from Scratch

This section provides a complete, production-ready transformer implementation incorporating all the optimization techniques discussed in this chapter. The implementation is modular, efficient, and includes comprehensive error checking and documentation.

```

import torch
import torch.nn as nn
import torch.nn.functional as F
from torch.cuda.amp import autocast, GradScaler
from torch.utils.checkpoint import checkpoint
import math

class MultiHeadAttention(nn.Module):
    """Optimized multi-head attention with optional Flash Attention."""

    def __init__(self, d_model, num_heads, dropout=0.1, use_flash=True):
        super().__init__()
        assert d_model % num_heads == 0

        self.d_model = d_model
        self.num_heads = num_heads
        self.head_dim = d_model // num_heads
        self.use_flash = use_flash

        # Combined QKV projection for efficiency
        self.qkv_proj = nn.Linear(d_model, 3 * d_model)
        self.out_proj = nn.Linear(d_model, d_model)
        self.dropout = nn.Dropout(dropout)

        self.scale = 1.0 / math.sqrt(self.head_dim)

    def forward(self, x, mask=None, past_kv=None, use_cache=False):
        batch_size, seq_len, _ = x.shape

```

```

# Project and split into Q, K, V
qkv = self.qkv_proj(x)
qkv = qkv.reshape(batch_size, seq_len, 3, self.num_heads,
                  self.head_dim)
qkv = qkv.permute(2, 0, 3, 1, 4) # (3, B, h, n, d)
Q, K, V = qkv[0], qkv[1], qkv[2]

# Use cached K, V if available
if past_kv is not None:
    past_K, past_V = past_kv
    K = torch.cat([past_K, K], dim=2)
    V = torch.cat([past_V, V], dim=2)

# Compute attention
if self.use_flash and hasattr(F, 'scaled_dot_product_attention'):
    # Use PyTorch's optimized implementation
    attn_output = F.scaled_dot_product_attention(
        Q, K, V,
        attn_mask=mask,
        dropout_p=self.dropout.p if self.training else 0.0
    )
else:
    # Standard attention implementation
    scores = torch.matmul(Q, K.transpose(-2, -1)) * self.scale

    if mask is not None:
        scores = scores.masked_fill(mask == 0, -1e9)

    attn_weights = F.softmax(scores, dim=-1)
    attn_weights = self.dropout(attn_weights)
    attn_output = torch.matmul(attn_weights, V)

# Reshape and project output
attn_output = attn_output.transpose(1, 2).contiguous()
attn_output = attn_output.reshape(batch_size, seq_len, self.d_model)
output = self.out_proj(attn_output)

# Return cache if requested
new_kv = (K, V) if use_cache else None
return output, new_kv

class FeedForward(nn.Module):
    """Feed-forward network with GELU activation."""

    def __init__(self, d_model, d_ff, dropout=0.1):
        super().__init__()
        self.linear1 = nn.Linear(d_model, d_ff)
        self.linear2 = nn.Linear(d_ff, d_model)
        self.dropout = nn.Dropout(dropout)

    def forward(self, x):
        x = self.linear1(x)
        x = F.gelu(x)

```

```

        x = self.dropout(x)
        x = self.linear2(x)
        return x

class TransformerLayer(nn.Module):
    """Single transformer layer with pre-norm architecture."""

    def __init__(self, d_model, num_heads, d_ff, dropout=0.1,
                  use_checkpoint=False):
        super().__init__()
        self.use_checkpoint = use_checkpoint

        self.attention = MultiHeadAttention(d_model, num_heads, dropout)
        self.ffn = FeedForward(d_model, d_ff, dropout)

        # Use fused layer norm if available
        try:
            from apex.normalization import FusedLayerNorm
            self.norm1 = FusedLayerNorm(d_model)
            self.norm2 = FusedLayerNorm(d_model)
        except ImportError:
            self.norm1 = nn.LayerNorm(d_model)
            self.norm2 = nn.LayerNorm(d_model)

        self.dropout = nn.Dropout(dropout)

    def forward(self, x, mask=None, past_kv=None, use_cache=False):
        if self.use_checkpoint and self.training:
            return checkpoint(self._forward, x, mask, past_kv, use_cache)
        return self._forward(x, mask, past_kv, use_cache)

    def _forward(self, x, mask, past_kv, use_cache):
        # Pre-norm attention
        normed = self.norm1(x)
        attn_out, new_kv = self.attention(normed, mask, past_kv, use_cache)
        x = x + self.dropout(attn_out)

        # Pre-norm feed-forward
        normed = self.norm2(x)
        ffn_out = self.ffn(normed)
        x = x + self.dropout(ffn_out)

        return x, new_kv

class TransformerModel(nn.Module):
    """Complete transformer model with all optimizations."""

    def __init__(self, vocab_size, d_model=768, num_heads=12,
                  num_layers=12, d_ff=3072, max_seq_len=512,
                  dropout=0.1, use_checkpoint=False):
        super().__init__()

        self.d_model = d_model
        self.num_layers = num_layers

```

```

# Embeddings
self.token_embed = nn.Embedding(vocab_size, d_model)
self.pos_embed = nn.Embedding(max_seq_len, d_model)
self.dropout = nn.Dropout(dropout)

# Transformer layers
self.layers = nn.ModuleList([
    TransformerLayer(d_model, num_heads, d_ff, dropout,
                    use_checkpoint)
    for _ in range(num_layers)
])

# Output layer norm
try:
    from apex.normalization import FusedLayerNorm
    self.final_norm = FusedLayerNorm(d_model)
except ImportError:
    self.final_norm = nn.LayerNorm(d_model)

# Language modeling head (tied with embeddings)
self.lm_head = nn.Linear(d_model, vocab_size, bias=False)
self.lm_head.weight = self.token_embed.weight

self._init_weights()

def _init_weights(self):
    """Initialize weights following best practices."""
    for module in self.modules():
        if isinstance(module, nn.Linear):
            nn.init.normal_(module.weight, mean=0.0, std=0.02)
            if module.bias is not None:
                nn.init.zeros_(module.bias)
        elif isinstance(module, nn.Embedding):
            nn.init.normal_(module.weight, mean=0.0, std=0.02)
        elif isinstance(module, nn.LayerNorm):
            nn.init.ones_(module.weight)
            nn.init.zeros_(module.bias)

def forward(self, input_ids, mask=None, past_key_values=None,
            use_cache=False):
    batch_size, seq_len = input_ids.shape

    # Create position IDs
    if past_key_values is None:
        past_length = 0
        position_ids = torch.arange(seq_len, device=input_ids.device)
    else:
        past_length = past_key_values[0][0].size(2)
        position_ids = torch.arange(
            past_length, past_length + seq_len,
            device=input_ids.device
        )

```



```

position_ids = position_ids.unsqueeze(0).expand(batch_size, -1)

# Embeddings
token_embeds = self.token_embed(input_ids)
pos_embeds = self.pos_embed(position_ids)
hidden_states = self.dropout(token_embeds + pos_embeds)

# Initialize cache
if past_key_values is None:
    past_key_values = [None] * self.num_layers

present_key_values = [] if use_cache else None

# Process through layers
for i, (layer, past_kv) in enumerate(
    zip(self.layers, past_key_values)):
    hidden_states, new_kv = layer(
        hidden_states, mask, past_kv, use_cache
    )
    if use_cache:
        present_key_values.append(new_kv)

# Final norm and projection
hidden_states = self.final_norm(hidden_states)
logits = self.lm_head(hidden_states)

return logits, present_key_values

```

21.10.2 Optimized Training Loop

The following training loop incorporates all optimization techniques including mixed precision, gradient accumulation, gradient clipping, and comprehensive logging.

```

def train_transformer(model, train_dataloader, val_dataloader, config):
    """Complete training loop with all optimizations."""

    # Setup
    device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
    model = model.to(device)

    # Optimizer
    optimizer = torch.optim.AdamW(
        model.parameters(),
        lr=config['learning_rate'],
        betas=(0.9, 0.999),
        eps=1e-6,
        weight_decay=0.01
    )

    # Learning rate scheduler
    total_steps = config['num_epochs'] * len(train_dataloader)
    warmup_steps = int(0.1 * total_steps)

    def lr_lambda(step):
        if step < warmup_steps:

```

```

        return step / warmup_steps
    return max(0.0, (total_steps - step) / (total_steps - warmup_steps))

scheduler = torch.optim.lr_scheduler.LambdaLR(optimizer, lr_lambda)

# Mixed precision scaler
scaler = GradScaler()

# Training state
global_step = 0
best_val_loss = float('inf')
accumulation_steps = config.get('accumulation_steps', 1)

# Training loop
for epoch in range(config['num_epochs']):
    model.train()
    epoch_loss = 0.0
    optimizer.zero_grad()

    for step, batch in enumerate(train_dataloader):
        # Move to device
        input_ids = batch['input_ids'].to(device, non_blocking=True)
        labels = batch['labels'].to(device, non_blocking=True)

        # Forward pass with mixed precision
        with autocast():
            logits, _ = model(input_ids)
            loss = F.cross_entropy(
                logits.view(-1, logits.size(-1)),
                labels.view(-1)
            )
            loss = loss / accumulation_steps

        # Backward pass
        scaler.scale(loss).backward()

        # Update weights
        if (step + 1) % accumulation_steps == 0:
            # Unscale and clip gradients
            scaler.unscale_(optimizer)
            torch.nn.utils.clip_grad_norm_(
                model.parameters(),
                max_norm=1.0
            )

            # Optimizer step
            scaler.step(optimizer)
            scaler.update()
            scheduler.step()
            optimizer.zero_grad()

        global_step += 1

    # Logging

```

```

        if global_step % config['log_interval'] == 0:
            current_lr = scheduler.get_last_lr()[0]
            print(f"Step {global_step}: "
                  f"loss={loss.item() * accumulation_steps:.4f}, "
                  f"lr={current_lr:.2e}")

    epoch_loss += loss.item() * accumulation_steps

# Validation
val_loss = evaluate(model, val_dataloader, device)
print(f"Epoch {epoch + 1}: "
      f"train_loss={epoch_loss / len(train_dataloader):.4f}, "
      f"val_loss={val_loss:.4f}")

# Save best model
if val_loss < best_val_loss:
    best_val_loss = val_loss
    torch.save({
        'epoch': epoch,
        'model_state_dict': model.state_dict(),
        'optimizer_state_dict': optimizer.state_dict(),
        'val_loss': val_loss,
    }, 'best_model.pt')

def evaluate(model, dataloader, device):
    """Evaluation loop."""
    model.eval()
    total_loss = 0.0

    with torch.no_grad():
        for batch in dataloader:
            input_ids = batch['input_ids'].to(device)
            labels = batch['labels'].to(device)

            with autocast():
                logits, _ = model(input_ids)
                loss = F.cross_entropy(
                    logits.view(-1, logits.size(-1)),
                    labels.view(-1)
                )

            total_loss += loss.item()

    return total_loss / len(dataloader)

```

21.10.3 Comprehensive Benchmarks

The following benchmarks demonstrate the impact of various optimizations on memory usage and training speed for a BERT-base model.

Baseline Configuration:

- Model: BERT-base (110M parameters)
- Batch size: 32
- Sequence length: 512

- Precision: FP32
- Hardware: NVIDIA A100 40GB

Optimization Results:

Configuration	Memory (GB)	Speed (samples/s)	Speedup
Baseline (FP32)	16.2	120	1.0×
+ Mixed Precision	10.1	280	2.3×
+ Gradient Checkpointing	7.8	220	1.8×
+ Optimized Batch Size	7.9	310	2.6×
+ Flash Attention	6.2	420	3.5×
All Optimizations	6.2	420	3.5×

Inference Optimization Results:

Configuration	Latency (ms)	Throughput	Memory (GB)
PyTorch FP32	45.2	22	1.8
+ TorchScript	38.1	26	1.8
+ FP16	22.3	45	0.9
+ Dynamic Quantization	18.7	54	0.5
+ TensorRT	9.2	109	0.6
All Optimizations	9.2	109	0.6

These benchmarks demonstrate that combining multiple optimization techniques can achieve substantial improvements in both training and inference performance. The key insight is that different optimizations address different bottlenecks, and the cumulative effect can be dramatic when applied systematically.

21.11 Distributed Training

21.11.1 Data Parallel

Simple multi-GPU:

```
if torch.cuda.device_count() > 1:
    model = nn.DataParallel(model)
```

Effective batch size: $B \times N_{\text{GPUs}}$

21.11.2 Distributed Data Parallel (DDP)

More efficient than DataParallel:

- One process per GPU
- Gradient synchronization via all-reduce
- Better scaling to multiple nodes

Setup requires:

1. Initialize process group
2. Wrap model in DistributedDataParallel
3. Use DistributedSampler for data
4. Synchronize across processes

21.12 Performance Optimization

21.12.1 DataLoader Optimization

The PyTorch DataLoader is often a bottleneck in training pipelines, particularly when data preprocessing is complex or I/O is slow. Proper configuration of the DataLoader can significantly improve training throughput by ensuring that data loading does not become the limiting factor.

The `num_workers` parameter controls how many subprocesses are used for data loading. Setting this too low results in the GPU waiting for data, while setting it too high can cause excessive CPU and memory usage. A good starting point is to use 4-8 workers per GPU, but the optimal value depends on the specific dataset and preprocessing pipeline.

```
from torch.utils.data import DataLoader

# Optimized DataLoader configuration
dataloader = DataLoader(
    dataset,
    batch_size=32,
    num_workers=8,          # Parallel data loading
    pin_memory=True,        # Faster GPU transfer
    persistent_workers=True, # Keep workers alive between epochs
    prefetch_factor=2       # Prefetch batches per worker
)
```

The `pin_memory` option allocates data in pinned (page-locked) memory, which enables faster transfers to the GPU using asynchronous DMA transfers. This can provide 20-30% speedup for data transfer, particularly beneficial when the model is small relative to the batch size.

Persistent workers keep the worker processes alive between epochs, avoiding the overhead of spawning new processes. This is particularly beneficial for datasets with expensive initialization or when using many workers.

21.12.2 Asynchronous Data Transfer

Overlapping data transfer with computation can hide data transfer latency. PyTorch supports non-blocking transfers that allow the CPU to continue executing while data is being copied to the GPU.

```
for batch in dataloader:
    # Non-blocking transfer to GPU
    input_ids = batch['input_ids'].to(device, non_blocking=True)
    labels = batch['labels'].to(device, non_blocking=True)

    # Computation can start while transfer completes
    with autocast():
        outputs = model(input_ids)
        loss = criterion(outputs, labels)
```

This technique is most effective when combined with pinned memory, as it enables true asynchronous transfers. The speedup depends on the ratio of transfer time to computation time, with larger models benefiting more as computation dominates.

21.12.3 Profiling with `torch.profiler`

Understanding where time is spent during training is essential for effective optimization. PyTorch's profiler provides detailed information about CPU and GPU operations, memory usage, and kernel execution times.

```

from torch.profiler import profile, ProfilerActivity, schedule

# Configure profiler
profiler_schedule = schedule(
    wait=1,          # Skip first batch
    warmup=1,        # Warmup for 1 batch
    active=3,        # Profile 3 batches
    repeat=2         # Repeat cycle twice
)

with profile(
    activities=[ProfilerActivity.CPU, ProfilerActivity.CUDA],
    schedule=profiler_schedule,
    on_trace_ready=torch.profiler.tensorboard_trace_handler('./log'),
    record_shapes=True,
    profile_memory=True,
    with_stack=True
) as prof:
    for step, batch in enumerate(dataloader):
        if step >= 10: # Profile first 10 batches
            break

        # Training step
        outputs = model(batch['input_ids'].to(device))
        loss = outputs.mean()
        loss.backward()
        optimizer.step()
        optimizer.zero_grad()

        prof.step() # Signal end of iteration

# Print summary
print(prof.key_averages().table(
    sort_by="cuda_time_total", row_limit=10))

```

The profiler output identifies operations that consume the most time, enabling targeted optimization. Common bottlenecks include inefficient attention implementations, excessive memory allocations, and CPU-GPU synchronization points.

21.12.4 Batch Size Tuning

Batch size has a complex relationship with training speed and model quality. Larger batches improve GPU utilization and reduce the number of optimizer steps, but may require learning rate adjustments and can affect convergence.

The optimal batch size maximizes GPU utilization without causing memory overflow. For transformer models, GPU utilization typically plateaus at batch sizes where the GPU is fully occupied, with further increases providing diminishing returns.

```

def benchmark_batch_sizes(model, seq_len, device='cuda'):
    """Benchmark training speed for different batch sizes."""
    results = []

    for batch_size in [8, 16, 32, 64, 128]:
        try:
            torch.cuda.empty_cache()

```

```

torch.cuda.reset_peak_memory_stats()

# Warmup
for _ in range(5):
    input_ids = torch.randint(0, 30000,
                              (batch_size, seq_len)).to(device)
    outputs = model(input_ids)
    loss = outputs.mean()
    loss.backward()

# Benchmark
torch.cuda.synchronize()
start = time.time()

for _ in range(20):
    input_ids = torch.randint(0, 30000,
                              (batch_size, seq_len)).to(device)
    outputs = model(input_ids)
    loss = outputs.mean()
    loss.backward()

torch.cuda.synchronize()
elapsed = time.time() - start

samples_per_sec = (20 * batch_size) / elapsed
memory_gb = torch.cuda.max_memory_allocated() / 1e9

results.append({
    'batch_size': batch_size,
    'samples_per_sec': samples_per_sec,
    'memory_gb': memory_gb
})

except RuntimeError as e:
    if "out of memory" in str(e):
        break

return results

```

21.12.5 Compilation with torch.compile

PyTorch 2.0 introduces `torch.compile`, which uses TorchDynamo and TorchInductor to compile models into optimized kernels. This can provide substantial speedups with minimal code changes.

```

# Compile model for faster execution
model = torch.compile(model, mode='max-autotune')

# Training proceeds as normal
for batch in dataloader:
    outputs = model(batch['input_ids'])
    loss = outputs.mean()
    loss.backward()
    optimizer.step()

```

The compilation process analyzes the model's computation graph and generates optimized CUDA

kernels. The first iteration is slow due to compilation overhead, but subsequent iterations benefit from the optimized code. Speedups of 20-50% are common for transformer models, with larger models typically seeing greater benefits.

21.13 Distributed Training Implementation

21.13.1 Understanding Distributed Strategies

Distributed training enables training on multiple GPUs or machines, dramatically reducing training time for large models. PyTorch provides several distributed training strategies, each with different trade-offs and use cases.

Data parallelism replicates the model on each GPU and distributes different batches of data to each replica. Gradients are synchronized across replicas after the backward pass, ensuring all replicas maintain identical weights. This approach scales well when the model fits in a single GPU's memory and is the most commonly used distributed training strategy.

Model parallelism splits the model itself across multiple GPUs, with different layers or components on different devices. This is necessary when the model is too large to fit on a single GPU but is more complex to implement and can suffer from poor GPU utilization due to sequential dependencies.

21.13.2 DistributedDataParallel Setup

DistributedDataParallel (DDP) is PyTorch's recommended approach for multi-GPU training. It provides better performance than DataParallel through more efficient gradient synchronization and support for multi-node training.

```
import torch.distributed as dist
from torch.nn.parallel import DistributedDataParallel as DDP
from torch.utils.data.distributed import DistributedSampler

def setup_distributed(rank, world_size):
    """Initialize distributed training."""
    os.environ['MASTER_ADDR'] = 'localhost'
    os.environ['MASTER_PORT'] = '12355'

    # Initialize process group
    dist.init_process_group(
        backend='nccl', # Use NCCL for GPU training
        rank=rank,
        world_size=world_size
    )

def cleanup_distributed():
    """Clean up distributed training."""
    dist.destroy_process_group()

def train_distributed(rank, world_size, model, dataset):
    """Training function for each process."""
    setup_distributed(rank, world_size)

    # Move model to GPU
    device = torch.device(f'cuda:{rank}')
    model = model.to(device)

    # Wrap model in DDP
    model = DDP(model, device_ids=[rank])
```



```

# Create distributed sampler
sampler = DistributedSampler(
    dataset,
    num_replicas=world_size,
    rank=rank,
    shuffle=True
)

# Create dataloader with distributed sampler
dataloader = DataLoader(
    dataset,
    batch_size=32,
    sampler=sampler,
    num_workers=4,
    pin_memory=True
)

optimizer = AdamW(model.parameters(), lr=1e-4)

# Training loop
for epoch in range(num_epochs):
    sampler.set_epoch(epoch) # Shuffle differently each epoch

    for batch in dataloader:
        input_ids = batch['input_ids'].to(device)
        labels = batch['labels'].to(device)

        outputs = model(input_ids)
        loss = criterion(outputs, labels)

        loss.backward()
        optimizer.step()
        optimizer.zero_grad()

    cleanup_distributed()

# Launch training on multiple GPUs
if __name__ == '__main__':
    world_size = torch.cuda.device_count()
    torch.multiprocessing.spawn(
        train_distributed,
        args=(world_size, model, dataset),
        nprocs=world_size,
        join=True
    )

```

21.13.3 Gradient Synchronization

DDP automatically synchronizes gradients across all processes during the backward pass using efficient all-reduce operations. The synchronization happens in parallel with the backward pass through gradient bucketing, which groups gradients into buckets and overlaps communication with computation.

The effective learning rate in distributed training should typically be scaled with the number of GPUs to maintain the same optimization dynamics. If training with 8 GPUs, the learning rate should

be multiplied by 8, or equivalently, the batch size per GPU should be kept constant and gradient accumulation used to achieve the same effective batch size.

21.13.4 Scaling Efficiency

Distributed training efficiency is measured by scaling efficiency, which compares actual speedup to ideal linear speedup. Perfect scaling would achieve $8\times$ speedup with 8 GPUs, but communication overhead and synchronization typically reduce this.

```
def measure_scaling_efficiency(model, batch_size, seq_len):
    """Measure scaling efficiency across different GPU counts."""
    results = {}

    # Single GPU baseline
    single_gpu_time = benchmark_single_gpu(model, batch_size, seq_len)
    results[1] = {
        'time': single_gpu_time,
        'speedup': 1.0,
        'efficiency': 1.0
    }

    # Multi-GPU measurements
    for num_gpus in [2, 4, 8]:
        if num_gpus > torch.cuda.device_count():
            break

        multi_gpu_time = benchmark_multi_gpu(
            model, batch_size, seq_len, num_gpus)
        speedup = single_gpu_time / multi_gpu_time
        efficiency = speedup / num_gpus

        results[num_gpus] = {
            'time': multi_gpu_time,
            'speedup': speedup,
            'efficiency': efficiency
        }

    return results
```

For transformer models, scaling efficiency typically ranges from 85-95% for 2-8 GPUs on a single node, with larger models achieving better efficiency due to higher computation-to-communication ratios. Multi-node training introduces additional communication overhead, with efficiency typically dropping to 70-85% depending on network bandwidth and model size.

21.14 Exercises

Exercise 21.1. Implement memory-efficient attention:

1. Implement chunked attention computation
2. Compare memory usage with standard attention
3. Test on sequences of length 512, 1024, 2048

4. Measure the memory-speed trade-off for different chunk sizes

Exercise 21.2. Optimize BERT training:

1. Start with baseline FP32 training, measure memory and speed
2. Add mixed precision, document improvements
3. Add gradient checkpointing, measure memory savings
4. Profile with torch.profiler and identify remaining bottlenecks
5. Achieve at least $2\times$ speedup while reducing memory by 40%

Exercise 21.3. Implement KV caching for GPT:

1. Modify transformer layer to support KV cache
2. Implement generation with and without caching
3. Benchmark generation speed for 100, 500, 1000 tokens
4. Measure memory overhead of caching
5. Calculate theoretical vs actual speedup

Exercise 21.4. Distributed training setup:

1. Implement DistributedDataParallel training
2. Train on 1, 2, 4, 8 GPUs
3. Measure scaling efficiency for each configuration
4. Identify communication bottlenecks
5. Optimize to achieve $\geq 85\%$ scaling efficiency

Exercise 21.5. Inference optimization pipeline:

1. Export model to TorchScript and ONNX
2. Apply dynamic quantization
3. Benchmark latency and throughput for each optimization
4. Create comparison table showing trade-offs
5. Achieve at least $3\times$ speedup with $\leq 2\%$ accuracy loss

Exercise 21.6. Complete implementation project:

1. Build mini-GPT (6 layers, 8 heads, $d=512$) from scratch
2. Implement all optimizations: mixed precision, checkpointing, KV cache
3. Train on WikiText-2 with comprehensive logging
4. Optimize inference with TorchScript and quantization
5. Generate samples and measure perplexity
6. Document memory usage and speed at each optimization stage

Chapter 22

Hardware Optimization and Deployment

Chapter Overview

Deploying transformers efficiently requires understanding hardware architectures, optimization techniques, and deployment strategies. This chapter covers GPUs, TPUs, model quantization, pruning, distillation, and production deployment best practices.

Learning Objectives

1. Understand GPU/TPU architectures for transformers
2. Apply model quantization (INT8, FP16)
3. Implement pruning and sparsity
4. Use knowledge distillation for compression
5. Optimize inference latency and throughput
6. Deploy models in production environments

22.1 Hardware Architectures

22.1.1 GPU Architecture for Deep Learning

Modern GPUs contain two primary types of compute units that are critical for transformer training and inference. CUDA cores are general-purpose floating-point units that can execute arbitrary arithmetic operations, while Tensor Cores are specialized matrix multiplication units designed specifically for deep learning workloads. Understanding the distinction between these units is essential for achieving optimal performance.

CUDA Cores vs Tensor Cores. CUDA cores provide flexibility for general computation but operate at lower throughput for matrix operations. A single NVIDIA A100 GPU contains 6912 CUDA cores capable of 19.5 TFLOPS at FP32 precision. In contrast, Tensor Cores are specialized hardware units that perform fused multiply-add operations on small matrix tiles. The same A100 GPU contains 432 third-generation Tensor Cores that deliver 312 TFLOPS for FP16 matrix multiplication, representing a $16\times$ advantage over CUDA cores for this specific operation. This dramatic performance difference makes Tensor Cores essential for transformer workloads, which are dominated by matrix multiplications in attention mechanisms and feed-forward layers.

Memory Hierarchy. GPU memory is organized in a hierarchy that trades capacity for access speed. At the top of the hierarchy, each streaming multiprocessor (SM) contains registers that provide the fastest access with approximately 256KB of storage per SM. These registers are private to individual threads and have single-cycle latency. The next level is shared memory, a software-managed cache that

allows threads within a block to communicate efficiently. The A100 provides 164KB of shared memory per SM with latency of approximately 20-30 cycles. Below this sits the L2 cache, a 40MB hardware-managed cache shared across all SMs with latency around 200 cycles. Finally, high-bandwidth memory (HBM) provides the largest capacity at 40-80GB but with the highest latency of 300-400 cycles. This hierarchy means that keeping data in faster memory levels is critical for performance.

The memory bandwidth available at each level determines how quickly data can be moved. The A100's HBM provides 1.6 TB/s of bandwidth, while the V100 provides 900 GB/s. Although these numbers seem large, they are often the bottleneck for transformer operations. Consider that the A100's Tensor Cores can consume data at a rate of $312 \text{ TFLOPS} \times 2 \text{ bytes (FP16)} = 624 \text{ TB/s}$ if fully utilized, far exceeding the 1.6 TB/s that HBM can supply. This mismatch between compute capability and memory bandwidth is why memory optimization is crucial for transformers.

Streaming Multiprocessors. The A100 contains 108 streaming multiprocessors, each capable of executing multiple thread blocks concurrently. Each SM has its own register file, shared memory, and L1 cache, along with 4 Tensor Cores. The SM scheduler can switch between thread warps (groups of 32 threads) with zero overhead, hiding memory latency by executing other warps while some wait for data. Achieving high occupancy, defined as the ratio of active warps to maximum possible warps, is essential for hiding latency and maximizing throughput.

22.1.2 Computational Intensity

Definition 22.1 (Arithmetic Intensity).

$$\text{Arithmetic Intensity} = \frac{\text{FLOPs}}{\text{Bytes transferred}} \quad (22.1)$$

For attention:

- \mathbf{QK}^\top : $\text{Intensity} = \frac{2n^2d}{2nd+2nd} = \frac{nd}{2d} = \frac{n}{2}$
- For small n (≤ 1024): Memory-bound
- For large n (≥ 4096): Compute-bound

22.1.3 Tensor Core Optimization

Tensor Cores achieve their peak performance only when specific conditions are met regarding data types, matrix dimensions, and memory layout. Understanding these requirements is essential for extracting maximum performance from modern GPUs.

Precision Requirements. Tensor Cores support several precision modes, each with different performance characteristics. FP16 (half precision) provides 312 TFLOPS on the A100, making it the standard choice for training. BF16 (bfloat16) offers the same throughput but with a larger dynamic range that better matches FP32, reducing the need for loss scaling during mixed-precision training. For inference, INT8 provides 624 TOPS, doubling throughput at the cost of reduced precision. The choice of precision involves trading off between speed, memory usage, and numerical accuracy.

Dimension Requirements. Tensor Cores operate on small matrix tiles and achieve peak performance when matrix dimensions are multiples of specific values. For FP16 operations, dimensions should be multiples of 8, while BF16 and INT8 operations prefer multiples of 16. When dimensions are not multiples of these values, the hardware must pad matrices or fall back to slower execution paths. For example, a matrix multiplication with dimensions 1023×1023 will perform significantly worse than 1024×1024 because the former requires padding or partial tile operations.

Example 22.1 (BERT-base Tensor Core Utilization). Consider BERT-base with hidden size $d = 768$ and 12 attention heads. The per-head dimension is $d_k = 768/12 = 64$, which is a multiple

of 8, allowing efficient Tensor Core usage. The query, key, and value projections have shape $[b \times n, 768] \times [768, 768]$ where b is batch size and n is sequence length.

Baseline configuration: Batch size 16, sequence length 128, FP32 precision

- Throughput: 145 sequences/second
- GPU utilization: 42%
- Memory bandwidth: 890 GB/s (56% of peak)

Optimized configuration: Batch size 32, sequence length 128, FP16 with Tensor Cores

- Throughput: 520 sequences/second ($3.6\times$ improvement)
- GPU utilization: 78%
- Memory bandwidth: 1.45 TB/s (91% of peak)
- Tensor Core utilization: 85%

The optimization involved: (1) switching to FP16 to enable Tensor Cores, (2) increasing batch size to improve arithmetic intensity, and (3) ensuring all matrix dimensions are multiples of 8. The result is a $3.6\times$ speedup with negligible accuracy loss when using mixed-precision training with loss scaling.

Achieving Peak Performance. Reaching 90% of theoretical peak TFLOPS requires careful attention to several factors. First, ensure sufficient work is available by using large batch sizes or long sequences to keep all Tensor Cores busy. Second, minimize memory transfers by fusing operations and reusing data in shared memory. Third, maintain high occupancy by using appropriate thread block sizes and avoiding resource limitations. Finally, profile the application to identify bottlenecks using tools like NVIDIA Nsight Compute, which can show Tensor Core utilization and identify whether operations are compute-bound or memory-bound.

22.1.4 TPU Architecture

Tensor Processing Units (TPUs):

- Systolic array architecture
- Optimized for matrix multiplications
- High memory bandwidth (900 GB/s, TPU v4)
- 275 TFLOPS (bfloat16)

TPU vs GPU:

Aspect	GPU	TPU
Flexibility	High (general purpose)	Medium (ML-specific)
Peak FLOPS	312 (A100 FP16)	275 (v4 bf16)
Memory	40-80 GB	32 GB (per chip)
Batch size	Medium-Large	Very Large
Best for	Flexibility, research	Large-scale training

22.2 Memory Optimization Techniques

Memory access patterns have a profound impact on GPU performance, often determining whether an operation runs at 10% or 90% of peak throughput. This section explores the key memory optimization techniques that are essential for efficient transformer implementations.

22.2.1 Coalesced Memory Access

When threads in a warp access global memory, the hardware attempts to combine these accesses into a single transaction. Coalesced access occurs when consecutive threads access consecutive memory locations, allowing the hardware to issue one memory transaction instead of 32 separate ones. For example, if thread 0 accesses address 0, thread 1 accesses address 4, thread 2 accesses address 8, and so on (assuming 4-byte elements), the hardware can coalesce these into a single 128-byte transaction. In contrast, if threads access random or strided locations, each access may require a separate transaction, reducing effective bandwidth by up to $32\times$.

For transformer operations, coalesced access is particularly important in matrix multiplications and attention computations. When loading a row of the query matrix, ensuring that consecutive threads load consecutive elements allows full utilization of memory bandwidth. This often requires careful consideration of matrix layout (row-major vs column-major) and access patterns in custom CUDA kernels.

22.2.2 Shared Memory and Bank Conflicts

Shared memory is divided into 32 banks that can be accessed simultaneously. When multiple threads in a warp access the same bank but different addresses, a bank conflict occurs, serializing the accesses and reducing throughput. The A100's shared memory has 32 banks with 4-byte bank width, meaning addresses that differ by 128 bytes map to the same bank.

In attention implementations, shared memory is commonly used to cache tiles of the query, key, and value matrices. Careful padding of these tiles can eliminate bank conflicts. For example, if each thread block loads a 64×64 tile of FP16 values, adding 8 elements of padding to each row ensures that consecutive rows start at different banks, eliminating conflicts when threads access columns.

Example 22.2 (Shared Memory Optimization in Attention). Consider computing attention scores $\mathbf{A} = \mathbf{Q}\mathbf{K}^\top$ where $\mathbf{Q}, \mathbf{K} \in \mathbb{R}^{n\times d}$. A naive implementation loads tiles of \mathbf{Q} and \mathbf{K} into shared memory and computes partial results.

Unoptimized approach:

- Tile size: 64×64 FP16 values
- Shared memory per tile: $64 \times 64 \times 2 = 8192$ bytes
- Bank conflicts when accessing columns: 32-way conflicts
- Effective bandwidth: 50 GB/s (3% of peak)

Optimized approach with padding:

- Tile size: 64×72 FP16 values (8 elements padding per row)
- Shared memory per tile: $64 \times 72 \times 2 = 9216$ bytes
- No bank conflicts: consecutive rows in different banks
- Effective bandwidth: 1.4 TB/s (88% of peak)

The 12.5% increase in shared memory usage eliminates bank conflicts and increases bandwidth by $28\times$, demonstrating the critical importance of memory access patterns.

22.2.3 Memory Bandwidth Utilization

Maximizing memory bandwidth requires ensuring that memory operations are large enough to amortize transaction overhead and that the GPU has sufficient concurrent operations to hide latency. Small

transfers are inefficient because they cannot fully utilize the 32-byte or 128-byte cache line sizes. Additionally, launching enough thread blocks to keep all memory controllers busy is essential for achieving peak bandwidth.

For transformers, memory bandwidth is often the limiting factor during attention computation with short sequences. When sequence length $n < 1024$, the arithmetic intensity of attention is low, meaning each floating-point operation requires loading relatively more data from memory. Techniques like Flash Attention address this by restructuring computations to maximize data reuse in shared memory, reducing the number of global memory accesses.

22.3 Kernel Fusion and Operation Optimization

Kernel fusion combines multiple operations into a single GPU kernel, reducing memory traffic and kernel launch overhead. This technique is particularly effective for transformers, where many operations are memory-bound and benefit from data reuse.

22.3.1 Fusion Opportunities in Transformers

Standard transformer implementations launch separate kernels for each operation, requiring intermediate results to be written to and read from global memory. For example, computing layer normalization followed by dropout requires writing the normalized values to memory, then reading them back for the dropout operation. Fusing these operations allows the normalized values to remain in registers or shared memory, eliminating the round-trip to global memory.

Common fusion patterns:

1. **Layer norm + dropout:** Normalize activations and apply dropout in a single pass, keeping intermediate values in registers.
2. **GELU + bias:** Compute the GELU activation and add bias without storing intermediate results.
3. **Attention score + softmax + dropout:** Compute $\text{softmax}(\mathbf{QK}^\top / \sqrt{d_k})$ and apply dropout in one kernel.
4. **Residual + layer norm:** Add residual connection and normalize in a single operation.

Example 22.3 (Layer Norm + Dropout Fusion). Consider layer normalization followed by dropout on a tensor of shape $[32, 128, 768]$ (batch size 32, sequence length 128, hidden size 768).

Unfused implementation:

- Layer norm kernel: Read 3.1M elements, write 3.1M elements
- Dropout kernel: Read 3.1M elements, write 3.1M elements
- Total memory traffic: 12.4M elements = 24.8 MB (FP16)
- Execution time: 0.18 ms

Fused implementation:

- Single kernel: Read 3.1M elements, write 3.1M elements
- Total memory traffic: 6.2M elements = 12.4 MB (FP16)
- Execution time: 0.10 ms

The fused kernel achieves $1.8\times$ speedup by halving memory traffic and eliminating kernel launch overhead. For a 12-layer transformer, this fusion appears 24 times per forward pass (twice per layer), providing substantial cumulative savings.

22.3.2 Flash Attention: Fused Attention Implementation

Flash Attention represents a sophisticated application of kernel fusion to the attention mechanism. Standard attention implementations compute $\mathbf{A} = \mathbf{QK}^\top$, write \mathbf{A} to memory, read it back for softmax, write the result, read it again for multiplication with \mathbf{V} , and finally write the output. This results in $O(n^2)$ memory reads and writes where n is sequence length.

Flash Attention restructures the computation to work on tiles that fit in shared memory. It computes attention for one tile at a time, keeping intermediate results in fast memory and only writing the final output. This reduces memory traffic from $O(n^2)$ to $O(n)$, providing dramatic speedups for long sequences.

Example 22.4 (Flash Attention Performance). BERT-base with sequence length 512, batch size 16, on A100 GPU:

Standard attention:

- Memory traffic: 48 GB per forward pass
- Execution time: 8.2 ms
- Memory bandwidth utilization: 45%

Flash Attention:

- Memory traffic: 12 GB per forward pass ($4\times$ reduction)
- Execution time: 3.8 ms ($2.2\times$ speedup)
- Memory bandwidth utilization: 82%

For longer sequences, the benefits are even more pronounced. At sequence length 2048, Flash Attention provides $3.5\times$ speedup, and at 8192, it provides $5.2\times$ speedup while also enabling sequences that would otherwise exceed memory capacity.

22.3.3 Implementing Fused Kernels

Creating fused kernels requires careful consideration of register usage, shared memory capacity, and thread block dimensions. The goal is to maximize data reuse while maintaining high occupancy. Modern deep learning frameworks provide tools for kernel fusion, including PyTorch's JIT compiler and TensorRT's graph optimizer, which can automatically fuse compatible operations. For custom fusion patterns, libraries like CUTLASS provide templates for efficient CUDA implementations.

22.4 Model Quantization

22.4.1 Quantization Fundamentals

Definition 22.2 (Quantization). Map FP32 weights to lower precision (INT8, FP16):

$$w_{\text{quant}} = \text{round}\left(\frac{w_{\text{float}}}{s}\right) + z \quad (22.2)$$

where s is scale factor, z is zero-point.

Precision options:

- **FP32:** 32 bits, full precision (baseline)
- **FP16:** 16 bits, $2\times$ compression
- **BF16:** 16 bits, better range than FP16
- **INT8:** 8 bits, $4\times$ compression
- **INT4:** 4 bits, $8\times$ compression (extreme)

22.4.2 Post-Training Quantization (PTQ)

Procedure:

1. Train model in FP32
2. Collect activation statistics on calibration set
3. Determine scale factors
4. Convert weights and activations to INT8

Example 22.5 (INT8 Quantization). **FP32 weight:** $w = 0.137$

Determine range: $w \in [-1.0, 1.0]$

Scale: $s = \frac{2.0}{256} = 0.0078125$

Quantize:

$$w_{\text{INT8}} = \text{round}\left(\frac{0.137}{0.0078125}\right) = \text{round}(17.54) = 18 \quad (22.3)$$

Dequantize: $w' = 18 \times 0.0078125 = 0.1406$

Error: $|0.137 - 0.1406| = 0.0036$ (2.6% relative)

22.4.3 Quantization-Aware Training (QAT)

Simulate quantization during training:

1. Forward pass: Quantize weights/activations
2. Compute loss with quantized values
3. Backward pass: FP32 gradients
4. Update FP32 weights

Benefits:

- Model learns to be robust to quantization
- Better accuracy than PTQ
- Minimal accuracy loss with INT8

Example 22.6 (BERT-base Quantization Results).

Precision	GLUE Score	Speedup
FP32 (baseline)	84.5	1.0×
FP16	84.4	1.8×
INT8 (PTQ)	82.1	2.9×
INT8 (QAT)	84.2	2.9×

QAT recovers most accuracy lost in PTQ!

22.5 Model Pruning

22.5.1 Pruning Strategies

Magnitude-based pruning:

$$\text{Prune if } |w_{ij}| < \tau \quad (22.4)$$

Structured pruning:

- Remove entire neurons, heads, layers
- Easier to deploy (no sparse kernels needed)
- Less aggressive compression

Unstructured pruning:

- Remove individual weights
- Higher compression ratios
- Requires sparse matrix operations

22.5.2 Iterative Pruning

[H] Iterative Magnitude Pruning

Input: Model, sparsity target s_{target}

sparsity $s = 0$ s_{target} by steps Train model to convergence

Prune Δs lowest-magnitude weights

Fine-tune model

Example 22.7 (Attention Head Pruning). BERT-base: 12 layers \times 12 heads = 144 heads

Finding: Can remove 50% of heads with minimal impact!

Procedure:

1. Compute importance score per head
2. Rank heads by importance
3. Prune lowest 50% (72 heads)
4. Fine-tune remaining model

Result:

- 50% fewer attention operations
- GLUE score: 84.5 \rightarrow 83.8 (0.7 point drop)

- 1.5× faster inference

22.6 Knowledge Distillation

22.6.1 Distillation Loss

Definition 22.3 (Distillation Objective).

$$\mathcal{L} = \alpha \mathcal{L}_{\text{CE}}(y, y_{\text{student}}) + (1 - \alpha) \mathcal{L}_{\text{KD}}(y_{\text{teacher}}, y_{\text{student}}) \quad (22.5)$$

where:

$$\mathcal{L}_{\text{KD}} = \text{KL} \left(\frac{\exp(z_t/T)}{\sum \exp(z_t/T)} \parallel \frac{\exp(z_s/T)}{\sum \exp(z_s/T)} \right) \quad (22.6)$$

T = temperature (typically 2-5), higher = softer probabilities

22.6.2 DistilBERT Approach

Student architecture:

- 6 layers (vs 12 in BERT)
- Same hidden size (768)
- Initialize from teacher's even layers

Training:

- Distillation loss from teacher
- Masked language modeling loss
- Cosine distance between hidden states

Results:

- 40% smaller (66M vs 110M params)
- 60% faster
- Retains 97% of BERT performance

22.7 Multi-GPU Training and Optimization

Training large transformer models requires distributing computation across multiple GPUs. The efficiency of multi-GPU training depends critically on communication bandwidth, parallelization strategy, and the balance between computation and communication.

22.7.1 Interconnect Technologies

The bandwidth between GPUs determines how quickly gradients, activations, and parameters can be exchanged during distributed training. PCIe provides 16-32 GB/s bidirectional bandwidth per GPU, which is adequate for small models but becomes a bottleneck for large transformers. NVLink, NVIDIA's proprietary interconnect, provides 600 GB/s bidirectional bandwidth on A100 systems, enabling much more efficient multi-GPU training. For comparison, the bandwidth within a single GPU (HBM) is 1600 GB/s, so NVLink provides roughly 40% of intra-GPU bandwidth for inter-GPU communication.

The impact of interconnect bandwidth is most visible during gradient synchronization in data parallel training. After computing gradients on each GPU's local batch, all GPUs must exchange

and average their gradients through an all-reduce operation. With PCIe, this communication can take longer than the backward pass itself for models with hundreds of millions of parameters. With NVLink, communication overhead is typically 10-20% of total training time.

22.7.2 Data Parallelism and Gradient Synchronization

Data parallelism replicates the model on each GPU and processes different batches on each device. After the backward pass, gradients are averaged across all GPUs using an all-reduce collective operation. The communication volume is equal to the model size, independent of batch size, making data parallelism most efficient for large batch sizes where computation time dominates communication time.

Example 22.8 (Data Parallel Scaling Efficiency). Training BERT-large (340M parameters) with batch size 32 per GPU on A100 GPUs:

Single GPU:

- Forward + backward time: 145 ms
- Throughput: 221 sequences/second

4 GPUs with NVLink:

- Forward + backward time: 145 ms (unchanged)
- Gradient all-reduce time: 18 ms
- Total time per step: 163 ms
- Throughput: 785 sequences/second
- Scaling efficiency: 89% (ideal would be 884 seq/s)

8 GPUs with NVLink:

- Forward + backward time: 145 ms
- Gradient all-reduce time: 22 ms
- Total time per step: 167 ms
- Throughput: 1533 sequences/second
- Scaling efficiency: 87%

The high scaling efficiency demonstrates that NVLink bandwidth is sufficient for data parallel training of BERT-large. With PCIe, the all-reduce time would be approximately 85 ms, reducing scaling efficiency to 63% for 8 GPUs.

22.7.3 Pipeline and Tensor Parallelism

For models too large to fit on a single GPU, pipeline parallelism splits the model across GPUs by layers, while tensor parallelism splits individual layers across GPUs. Pipeline parallelism has lower communication requirements but suffers from pipeline bubbles where some GPUs are idle. Tensor parallelism requires more communication (activations must be exchanged between layers) but maintains better GPU utilization.

Modern training frameworks like Megatron-LM combine data, pipeline, and tensor parallelism to train models with hundreds of billions of parameters. For example, training a 175B parameter model might use 8-way tensor parallelism, 8-way pipeline parallelism, and 32-way data parallelism across 2048 GPUs.

22.7.4 Overlapping Communication and Computation

Advanced implementations overlap gradient communication with backward pass computation. As soon as gradients for one layer are computed, they can begin synchronizing while the backward pass continues on earlier layers. This technique, called gradient bucketing, can hide most communication overhead when sufficient computation is available to overlap.

PyTorch's `DistributedDataParallel` automatically implements gradient bucketing, grouping parameters into buckets of approximately 25MB and launching all-reduce operations as soon as each bucket's gradients are ready. This optimization is particularly effective for large models where the backward pass takes much longer than gradient synchronization.

22.8 Inference Optimization

22.8.1 ONNX Runtime

ONNX (Open Neural Network Exchange):

- Framework-agnostic model format
- Optimized inference engine
- Supports quantization, pruning

Optimizations:

- Operator fusion (combine LayerNorm + Add)
- Constant folding
- Dead code elimination
- Graph optimization

22.8.2 TensorRT

NVIDIA TensorRT:

- Deep learning inference optimizer
- Layer fusion
- Kernel auto-tuning
- INT8 calibration

Typical speedups:

- BERT-base: 2-3× over PyTorch
- With INT8: 4-5× over PyTorch FP32

22.8.3 Batching Strategies

Static batching:

- Fixed batch size
- Pad to max length
- Simple but wasteful

Dynamic batching:

- Accumulate requests until batch full or timeout
- Reduces latency while maintaining throughput

Continuous batching:

- For autoregressive generation
- Add new sequences as others finish
- Maximizes GPU utilization

22.9 Production Deployment

22.9.1 Serving Frameworks

TorchServe:

- PyTorch native serving
- REST/gRPC APIs
- Batching, versioning, monitoring

Triton Inference Server:

- Multi-framework (PyTorch, TensorFlow, ONNX)
- Concurrent model execution
- Dynamic batching
- Model ensembles

FastAPI + Custom:

- Lightweight, flexible
- Full control over serving logic
- Easy integration with existing systems

22.9.2 Deployment Checklist

Performance:

- Quantize to INT8/FP16
- Export to ONNX/TensorRT
- Optimize batch size for latency/throughput
- Enable KV caching for generation

Reliability:

- Graceful degradation on errors
- Request timeouts
- Health checks
- Model versioning

Monitoring:

- Latency (p50, p95, p99)
- Throughput (requests/second)
- GPU utilization
- Error rates

22.10 Hardware Selection and Cost Analysis

Selecting appropriate hardware for transformer workloads requires balancing performance, cost, and operational requirements. This section provides guidance for different use cases and scales.

22.10.1 CPU vs GPU Trade-offs

The choice between CPU and GPU depends on model size, batch size, latency requirements, and cost constraints. CPUs excel at low-latency inference with small batch sizes, while GPUs provide superior throughput for larger batches and are essential for training.

When to use CPUs:

- Small models ($\leq 100\text{M}$ parameters) with batch size 1-4
- Latency-critical applications requiring $\leq 10\text{ms}$ response time
- Cost-sensitive deployments with low throughput requirements
- Edge deployment where GPU hardware is unavailable

For example, a distilled BERT model with 66M parameters can achieve 15ms latency on a modern CPU (Intel Xeon or AMD EPYC) with batch size 1. The same model on a T4 GPU achieves 8ms latency but requires batching to amortize GPU overhead, making it less suitable for single-request scenarios.

When to use GPUs:

- Training any transformer model
- Large models ($\geq 100\text{M}$ parameters)
- Batch inference with batch size ≥ 8
- Throughput-oriented applications

A BERT-large model (340M parameters) achieves 2.5 sequences/second on CPU but 45 sequences/second on a T4 GPU with batch size 16, demonstrating the GPU's $18\times$ throughput advantage for larger models.

Cost analysis. Cloud GPU instances cost approximately $3\text{-}5\times$ more than equivalent CPU instances. However, the higher throughput often makes GPUs more cost-effective per inference. For BERT-large inference, a T4 GPU instance costs \$0.35/hour and processes 162,000 sequences/hour, yielding \$0.0000022 per sequence. A CPU instance costs \$0.10/hour and processes 9,000 sequences/hour, yielding \$0.000011 per sequence, making the GPU $5\times$ more cost-effective despite higher instance cost.

Energy efficiency. GPUs also provide better energy efficiency for large models. The T4 GPU consumes 70W and processes 45 sequences/second, yielding 0.64 sequences/second/watt. A high-end CPU consumes 200W and processes 2.5 sequences/second, yielding 0.0125 sequences/second/watt, making the GPU $51\times$ more energy-efficient.

22.10.2 Training Hardware Selection

Training requirements scale dramatically with model size, from single GPUs for small models to thousands of GPUs for the largest models.

Small models (< 1B parameters):

- Hardware: Single V100 (32GB) or A100 (40GB)
- Training time: BERT-base on 16GB dataset trains in 3 days on V100
- Cost: \$2.50/hour (V100) to \$4.00/hour (A100) on cloud
- Use case: Research, fine-tuning, domain-specific models

Medium models (1-10B parameters):

- Hardware: 4-8× A100 (40GB or 80GB) with NVLink
- Training time: GPT-2 (1.5B) trains in 2 weeks on 8× A100
- Cost: \$32/hour for 8× A100 on cloud
- Use case: Production models, large-scale fine-tuning

For these models, NVLink is essential for efficient data parallelism. The 600 GB/s NVLink bandwidth enables 85-90% scaling efficiency, while PCIe would reduce efficiency to 60-70%.

Large models (10-100B parameters):

- Hardware: 16-64× A100 (80GB) with NVLink and InfiniBand
- Training time: GPT-3 (175B) trained on 10,000× V100-days
- Cost: \$128-256/hour for 32-64× A100 on cloud
- Use case: Foundation models, large-scale pretraining

These models require pipeline and tensor parallelism in addition to data parallelism. InfiniBand provides 200-400 Gb/s bandwidth between nodes, enabling efficient multi-node training.

Extreme models (100B+ parameters):

- Hardware: 100s-1000s of A100 GPUs across multiple nodes
- Training time: Months on large clusters
- Cost: Millions of dollars for full pretraining
- Use case: State-of-the-art foundation models (GPT-4, PaLM)

22.10.3 Inference Hardware Selection

Inference requirements vary widely based on latency, throughput, and cost constraints.

Batch inference (throughput-oriented):

- Hardware: A100 or A10 GPUs
- Characteristics: Process large batches (32-128), optimize for throughput
- Use case: Offline processing, data pipelines, batch predictions
- Example: BERT-large processes 520 sequences/second on A100 with batch size 64

Low-latency inference:

- Hardware: T4 or A10 GPUs with TensorRT

- Characteristics: Small batches (1-8), optimize for latency
- Use case: Real-time applications, interactive systems
- Example: BERT-base achieves 5ms latency on T4 with batch size 1 using INT8

Cost-optimized inference:

- Hardware: CPU or small GPUs (T4)
- Characteristics: Quantized models, efficient batching
- Use case: High-volume, cost-sensitive applications
- Example: Distilled BERT (66M params) on CPU costs \$0.000005 per inference

Edge deployment:

- Hardware: Mobile CPUs, edge TPUs, or NVIDIA Jetson
- Characteristics: Heavily quantized (INT8/INT4), pruned models
- Use case: On-device inference, privacy-sensitive applications
- Example: MobileBERT (25M params) runs at 30ms on mobile CPU

22.10.4 Hardware Selection Decision Tree

The following decision tree provides guidance for hardware selection:

1. **Training or inference?**

- Training: Proceed to step 2
- Inference: Proceed to step 4

2. **Model size?**

- ≤ 1 B params: Single V100/A100
- 1-10B params: 4-8 \times A100 with NVLink
- 10-100B params: 16-64 \times A100 with InfiniBand
- ≥ 100 B params: 100s-1000s of GPUs

3. **Budget constraints?**

- Research/limited budget: V100 or A10
- Production/performance-critical: A100

4. **Latency requirements?**

- ≤ 10 ms: T4 GPU with TensorRT or CPU for small models
- 10-50ms: T4 or A10 GPU
- ≥ 50 ms: Any GPU or CPU based on cost

5. **Throughput requirements?**

- ≤ 10 seq/s: CPU
- 10-100 seq/s: T4 GPU
- ≥ 100 seq/s: A10 or A100 GPU

22.11 Exercises

Exercise 22.1. Quantize BERT-base to INT8:

1. Use PyTorch quantization APIs
2. Calibrate on 1000 examples
3. Measure: (a) Model size, (b) Inference speed, (c) GLUE accuracy
4. Compare PTQ vs QAT

Exercise 22.2. Implement attention head pruning:

1. Compute importance scores for all heads
2. Prune 25%, 50%, 75% of heads
3. Fine-tune after pruning
4. Plot accuracy vs sparsity

Exercise 22.3. Optimize inference pipeline:

1. Baseline: PyTorch FP32
2. Convert to ONNX, measure speedup
3. Apply INT8 quantization
4. Implement dynamic batching
5. Report final throughput improvement

Exercise 22.4. Analyze Tensor Core utilization:

1. Profile BERT-base training with FP32 and FP16
2. Measure Tensor Core utilization using NVIDIA Nsight Compute
3. Experiment with different batch sizes and sequence lengths
4. Identify which operations benefit most from Tensor Cores
5. Calculate achieved TFLOPS as percentage of theoretical peak

Exercise 22.5. Implement and benchmark kernel fusion:

1. Create separate kernels for layer norm and dropout
2. Implement a fused layer norm + dropout kernel

3. Measure memory bandwidth utilization for both approaches
4. Compare execution time across different tensor sizes
5. Analyze the speedup and explain the performance difference

Exercise 22.6. Optimize memory access patterns:

1. Implement matrix multiplication with and without coalesced access
2. Add padding to eliminate shared memory bank conflicts
3. Profile both implementations using Nsight Compute
4. Measure effective memory bandwidth for each version
5. Document the impact of access patterns on performance

Exercise 22.7. Multi-GPU scaling analysis:

1. Train BERT-base on 1, 2, 4, and 8 GPUs
2. Measure training time and throughput for each configuration
3. Calculate scaling efficiency relative to single GPU
4. Profile communication overhead using NVIDIA Nsight Systems
5. Compare PCIe vs NVLink if available

Exercise 22.8. Hardware selection analysis:

1. Choose a transformer model and deployment scenario
2. Estimate throughput requirements and latency constraints
3. Compare cost per inference for CPU, T4, and A100
4. Calculate break-even point where GPU becomes cost-effective
5. Recommend hardware configuration with justification

Exercise 22.9. Flash Attention implementation study:

1. Implement standard attention with separate kernels
2. Analyze memory traffic for different sequence lengths
3. Study Flash Attention paper and implementation
4. Benchmark Flash Attention vs standard attention

5. Plot speedup as a function of sequence length

Chapter 23

Best Practices and Production Case Studies

Chapter Overview

This final chapter synthesizes practical wisdom from deploying transformers at scale. We cover debugging strategies, hyperparameter tuning, common pitfalls, and real-world case studies from industry deployments of BERT, GPT, and other transformer models.

Learning Objectives

1. Apply systematic debugging for transformer training
2. Tune hyperparameters effectively
3. Avoid common pitfalls in architecture and training
4. Learn from real-world deployment case studies
5. Design robust production systems
6. Plan future-proof transformer architectures

23.1 Model Selection

Choosing the right transformer architecture is a critical decision that impacts both performance and resource requirements. This section provides a systematic framework for selecting among the major transformer variants based on task requirements, data availability, and computational constraints.

23.1.1 Architecture Selection Framework

The choice between BERT, GPT, T5, and other architectures depends fundamentally on the nature of your task. BERT and its variants excel at understanding tasks where bidirectional context is crucial, such as classification, named entity recognition, and question answering. The bidirectional attention mechanism allows BERT to build rich representations by attending to both past and future tokens simultaneously, making it particularly effective when the entire input is available at once.

GPT models, in contrast, are designed for generation tasks where autoregressive decoding is required. The unidirectional attention pattern makes GPT natural for text generation, code completion, and any task where outputs must be produced sequentially. While GPT can be adapted for understanding tasks through careful prompting, this is generally less efficient than using a bidirectional model designed for the purpose.

T5 represents a unified approach that frames all tasks as sequence-to-sequence problems. This architecture provides flexibility across both understanding and generation tasks, making it an excellent choice when you need a single model to handle diverse task types. The encoder-decoder structure allows

T5 to leverage bidirectional attention in the encoder while maintaining autoregressive generation in the decoder.

23.1.2 Model Size Selection

Selecting the appropriate model size requires balancing performance requirements against computational constraints. The relationship between model size and performance generally follows a power law, with diminishing returns as models grow larger. For most practical applications, the base-sized models provide an excellent balance between capability and efficiency.

BERT-base with 110 million parameters serves as the standard choice for most understanding tasks. It provides strong performance across a wide range of benchmarks while remaining tractable for fine-tuning on a single GPU. BERT-large with 340 million parameters offers modest improvements, typically 1-3 percentage points on downstream tasks, but requires significantly more memory and computation. The large variant is justified primarily when you need to extract maximum performance and have sufficient computational resources.

For GPT models, the size selection depends heavily on the complexity of the generation task. GPT-2 small (117M parameters) suffices for simple completion tasks and domain-specific generation after fine-tuning. GPT-2 medium (345M parameters) provides better coherence for longer generations and more complex tasks. The larger variants (GPT-2 large at 774M and GPT-2 XL at 1.5B parameters) are necessary primarily when working with limited task-specific data, as their stronger pre-trained representations enable better few-shot performance.

23.1.3 Pre-trained versus Training from Scratch

The decision to use pre-trained models versus training from scratch depends on data availability, domain specificity, and computational budget. In nearly all cases, starting from pre-trained weights is the correct choice. Pre-training on large corpora provides general language understanding that transfers effectively to downstream tasks, and the computational cost of pre-training from scratch is prohibitive for most organizations.

Training from scratch becomes viable only in specific circumstances. When working with highly specialized domains where general language models perform poorly, such as medical text with extensive jargon or programming languages not well-represented in pre-training data, domain-specific pre-training may be justified. However, even in these cases, continued pre-training from existing checkpoints is typically more efficient than starting from random initialization.

The computational cost difference is substantial. Pre-training BERT-base from scratch requires approximately 64 TPU days or equivalent GPU time, representing tens of thousands of dollars in compute costs. Fine-tuning the same model on a downstream task typically requires only hours on a single GPU, costing tens of dollars. This thousand-fold difference in cost makes pre-trained models the default choice for nearly all applications.

23.1.4 Cost-Benefit Analysis

A systematic cost-benefit analysis should consider both direct computational costs and opportunity costs. For a typical classification task with 10,000 labeled examples, fine-tuning BERT-base requires approximately 2-4 hours on a single V100 GPU, costing roughly \$10-20 in cloud compute. This investment typically yields performance improvements of 5-15 percentage points over traditional methods like logistic regression on TF-IDF features.

Training a smaller model from scratch on the same data might require 8-16 hours and cost \$40-80, while likely achieving inferior performance due to the lack of pre-trained representations. The pre-trained approach thus provides both better performance and lower cost, a rare combination that explains the dominance of transfer learning in modern NLP.

For generation tasks, the cost analysis shifts somewhat. Fine-tuning GPT-2 medium on a specific generation task requires 4-8 hours on a V100, costing \$20-40. However, inference costs become more significant for generation, as producing each token requires a full forward pass through the model.

For applications requiring high-throughput generation, the ongoing inference costs may exceed training costs within weeks or months of deployment, making inference optimization critical.

23.2 Training Best Practices

Effective training of transformer models requires careful attention to hyperparameter selection, monitoring, and debugging. This section provides comprehensive guidance on the key decisions that impact training success.

23.2.1 Learning Rate Selection and Tuning

The learning rate is the single most important hyperparameter for transformer training. The optimal learning rate depends on model size, batch size, and whether you are pre-training or fine-tuning. For fine-tuning pre-trained models, learning rates are typically much smaller than for training from scratch, as the pre-trained weights already occupy a good region of the loss landscape.

For BERT fine-tuning, learning rates between 1×10^{-5} and 5×10^{-5} work well across most tasks. The smaller end of this range is appropriate for tasks with limited data or when fine-tuning for many epochs, while the larger end works better with abundant data and shorter training. GPT fine-tuning typically uses slightly higher learning rates, in the range of 2×10^{-5} to 1×10^{-4} , as the autoregressive objective is somewhat more stable than BERT's masked language modeling.

The learning rate should scale approximately linearly with batch size. If you double the batch size, you can typically increase the learning rate by a factor of 1.5 to 2 without harming convergence. This relationship holds because larger batches provide more accurate gradient estimates, allowing the optimizer to take larger steps safely. However, this scaling breaks down at very large batch sizes, where additional techniques like LAMB optimizer become necessary.

A learning rate warmup is essential for stable training. During warmup, the learning rate increases linearly from zero to the target value over a specified number of steps. This prevents the large gradients that can occur early in training from pushing the model into a poor region of the loss landscape. For fine-tuning, a warmup period covering 5-10% of total training steps is typical. For pre-training from scratch, longer warmup periods of 10,000 to 50,000 steps are common.

23.2.2 Batch Size Selection

Batch size selection involves balancing computational efficiency, memory constraints, and optimization dynamics. Larger batches improve GPU utilization and training throughput but require more memory and can sometimes harm generalization. For fine-tuning on a single GPU, batch sizes between 16 and 32 are typical, as these fit comfortably in memory while providing reasonable gradient estimates.

When memory constraints prevent using your desired batch size, gradient accumulation provides an effective solution. By accumulating gradients over multiple forward-backward passes before updating weights, you can simulate larger batch sizes without increasing memory requirements. For example, using a batch size of 8 with 4 accumulation steps is equivalent to a batch size of 32 in terms of optimization dynamics, though it takes four times as long to complete each update.

The relationship between batch size and training dynamics is subtle. Very small batches (below 8) introduce significant noise into gradient estimates, which can slow convergence and lead to instability. Very large batches (above 512 for fine-tuning) can lead to sharp minima that generalize poorly, though this effect is less pronounced when using appropriate learning rate scaling and warmup. For most fine-tuning tasks, batch sizes between 16 and 64 provide a good balance.

23.2.3 Checkpointing and Monitoring Strategy

Effective monitoring is essential for detecting problems early and understanding training dynamics. At minimum, you should track training loss, validation loss, and task-specific metrics at regular intervals. Logging every 100-500 steps provides sufficient granularity to detect issues without generating excessive data.

Checkpointing strategy depends on training duration and stability. For short fine-tuning runs of a few hours, saving checkpoints every epoch is sufficient. For longer training runs, save checkpoints every few thousand steps to protect against hardware failures and enable recovery from divergence. Always keep at least the three most recent checkpoints, as the most recent checkpoint may be corrupted or represent a point after training has diverged.

Beyond basic loss monitoring, tracking gradient norms provides early warning of training instability. Gradient norms should remain relatively stable throughout training, typically in the range of 0.1 to 10.0. Sudden spikes in gradient norm often precede loss divergence and indicate that gradient clipping or learning rate reduction may be necessary. Similarly, monitoring the ratio of update magnitude to parameter magnitude helps ensure that learning rates are appropriate.

23.3 Memory Management

Memory management is often the primary constraint in transformer training and deployment. Understanding memory requirements and optimization strategies enables training larger models and processing longer sequences.

23.3.1 Estimating Memory Requirements

The memory footprint of transformer training consists of several components. Model parameters themselves require 4 bytes per parameter in FP32 or 2 bytes in FP16. For BERT-base with 110 million parameters, this amounts to 440 MB in FP32 or 220 MB in FP16. However, parameters are only a small fraction of total memory usage during training.

Optimizer states typically dominate memory consumption. The Adam optimizer maintains two additional tensors per parameter for first and second moment estimates, tripling the memory required for parameters. For BERT-base, this adds another 880 MB in FP32, bringing the total for parameters and optimizer states to 1.32 GB. This explains why mixed precision training provides such substantial memory savings, as storing optimizer states in FP32 while computing in FP16 reduces this component significantly.

Activations for backpropagation constitute the other major memory consumer. Each transformer layer stores activations for all tokens in the batch, and these accumulate across layers. For a batch size of 32 with sequence length 512 and hidden size 768, each attention layer stores approximately 150 MB of activations. With 12 layers, this totals 1.8 GB just for attention activations, not including feed-forward layers and other components.

A practical formula for estimating total training memory is: $\text{Memory (GB)} = (\text{Parameters} \times 16 \text{ bytes}) + (\text{Batch Size} \times \text{Sequence Length} \times \text{Hidden Size} \times \text{Layers} \times 4 \text{ bytes} \times 10)$. The factor of 16 accounts for parameters, gradients, and optimizer states in mixed precision. The factor of 10 in the activation term is an empirical multiplier accounting for all activation storage. For BERT-base with batch size 32 and sequence length 512, this formula predicts approximately 8 GB, which matches observed usage.

23.3.2 Choosing Batch Size and Sequence Length

When memory is constrained, you must balance batch size and sequence length. Reducing sequence length has a quadratic effect on memory usage due to the self-attention mechanism, while reducing batch size has only a linear effect. Therefore, if your task permits, reducing sequence length is more effective for memory savings.

Many tasks do not require the full 512-token context that BERT supports. For sentence classification, sequences are often much shorter, and truncating to 128 or 256 tokens may have minimal impact on performance while reducing memory by $4\times$ or $2\times$ respectively. For tasks that do require long contexts, consider hierarchical approaches that process the input in chunks rather than all at once.

When you cannot reduce sequence length, reducing batch size is the next option. As discussed earlier, gradient accumulation allows you to maintain effective batch size while using smaller micro-

batches that fit in memory. The trade-off is increased training time, as you perform more forward-backward passes per update. However, this is often preferable to being unable to train at all.

23.3.3 Gradient Checkpointing

Gradient checkpointing trades computation for memory by recomputing activations during the backward pass rather than storing them. This technique can reduce activation memory by a factor of 4-8 with only a 20-30% increase in training time. For memory-constrained scenarios, this trade-off is often worthwhile.

The implementation is straightforward in modern frameworks. In PyTorch, wrapping transformer layers with checkpoint functions causes activations to be recomputed during backpropagation. The memory savings are substantial because activations typically consume more memory than parameters and optimizer states combined. For BERT-base, gradient checkpointing can reduce memory usage from 8 GB to 4 GB, enabling training with twice the batch size or sequence length.

The computational overhead of gradient checkpointing is less severe than it might appear. Modern GPUs are often memory-bound rather than compute-bound for transformer training, meaning that the additional forward passes during backpropagation can be performed with minimal wall-clock time increase. The actual slowdown is typically 20-30% rather than the 50% that naive analysis would suggest.

23.3.4 Mixed Precision Training

Mixed precision training uses FP16 for most computations while maintaining FP32 master copies of weights for numerical stability. This approach reduces memory usage by approximately $2\times$ and can accelerate training by $2\text{--}3\times$ on modern GPUs with tensor cores. The memory savings come from storing activations and gradients in FP16, while the speed improvements come from faster FP16 arithmetic on specialized hardware.

Implementing mixed precision requires careful attention to numerical stability. Loss scaling prevents gradient underflow by multiplying the loss by a large constant before backpropagation, then dividing gradients by the same constant before the optimizer step. Dynamic loss scaling automatically adjusts this constant to maximize precision without causing overflow. Modern frameworks like PyTorch's automatic mixed precision handle these details automatically.

The performance benefits of mixed precision are most pronounced on GPUs with tensor cores, such as NVIDIA V100, A100, and later architectures. On these devices, FP16 matrix multiplications can be $2\text{--}8\times$ faster than FP32 operations. For transformer training, which is dominated by matrix multiplications in attention and feed-forward layers, this translates to substantial end-to-end speedups.

23.3.5 Multi-GPU Strategies

When training exceeds the capacity of a single GPU, several parallelism strategies are available. Data parallelism replicates the model across GPUs and splits batches across devices, with gradients synchronized after each backward pass. This approach scales well up to 8-16 GPUs and is the simplest to implement, requiring minimal code changes.

Model parallelism splits the model itself across GPUs, with different layers or components on different devices. This becomes necessary when the model is too large to fit on a single GPU, even with batch size 1. Pipeline parallelism is a variant that divides the model into stages and processes multiple micro-batches concurrently, improving efficiency by keeping all GPUs busy.

For very large models, tensor parallelism splits individual layers across GPUs, partitioning matrix multiplications so that each GPU computes a portion of each layer. This requires more sophisticated implementation but provides the finest-grained parallelism. Modern frameworks like DeepSpeed and Megatron-LM combine these strategies, using tensor parallelism within nodes and pipeline parallelism across nodes for maximum efficiency.

23.4 Debugging Transformers

23.4.1 Systematic Debugging Workflow

Level 1: Data sanity checks

1. Visualize input samples
2. Verify labels are correct
3. Check for data leakage
4. Validate preprocessing

Level 2: Model sanity checks

1. Overfit single batch (should reach near-zero loss)
2. Check gradient flow (no dead neurons)
3. Verify shapes at each layer
4. Test with minimal model first

Level 3: Training dynamics

1. Monitor loss curves (training + validation)
2. Track gradient norms
3. Visualize attention weights
4. Check learning rate schedule

Example 23.1 (Debugging Checklist). **Symptom:** Loss not decreasing

Diagnose:

- Learning rate too low? Try 10× higher
- Frozen layers? Check `requires_grad`
- Optimizer issue? Try SGD as baseline
- Bad initialization? Re-initialize
- Data issue? Manually inspect batches

Symptom: NaN loss

Diagnose:

- Gradient explosion? Add clipping
- Numerical instability? Check mask values ($-\infty$ vs $-1e9$)
- Learning rate too high? Reduce 10×
- Mixed precision issue? Check loss scaling

23.4.2 Gradient Analysis

Monitor per-layer gradient norms:

```
for name, param in model.named_parameters():
    if param.grad is not None:
        grad_norm = param.grad.norm().item()
        print(f"{name}: {grad_norm:.6f}")
```

Healthy gradients:

- Norms between 10^{-4} and 10^1
- Similar across layers (no extreme differences)
- Non-zero for all layers

23.4.3 Common Training Issues and Solutions

Out of memory errors are among the most common issues in transformer training. When you encounter OOM errors, the first step is to reduce batch size. If you are already at batch size 1, consider reducing sequence length, enabling gradient checkpointing, or using mixed precision training. As a last resort, you may need to use a smaller model or move to a GPU with more memory.

Slow training can result from several factors. Poor GPU utilization often indicates that your batch size is too small or that data loading is a bottleneck. Monitor GPU utilization using `nvidia-smi`, and if it is below 80%, try increasing batch size or using more data loading workers. If training is slow despite high GPU utilization, you may be using an inefficient implementation or could benefit from mixed precision training.

Poor convergence manifests as loss that decreases slowly or plateaus at a suboptimal value. This often indicates that your learning rate is too low. Try increasing it by a factor of 3-10 and observe whether convergence improves. If loss diverges with higher learning rates, the issue may be poor initialization or an architectural problem. Verify that your model can overfit a small batch of data, as this confirms that the architecture and implementation are correct.

Gradient explosion appears as sudden spikes in loss or NaN values. The immediate solution is gradient clipping, which limits the norm of gradients to a maximum value, typically 1.0 or 5.0. If gradient explosion persists despite clipping, reduce the learning rate or increase warmup duration. Check for numerical instability in attention masks, ensuring that masked positions use `-1e9` rather than `-infinity`, as the latter can cause numerical issues in some implementations.

23.5 Inference Optimization

Optimizing transformer inference is critical for production deployment, as inference costs often exceed training costs over the lifetime of a model. This section covers techniques for reducing latency and increasing throughput while maintaining accuracy.

23.5.1 Batch Size for Throughput versus Latency

The batch size used during inference involves a fundamental trade-off between throughput and latency. Processing requests individually minimizes latency, as each request is handled immediately without waiting for others. However, this approach severely underutilizes GPU resources, as transformers achieve much higher throughput when processing multiple examples simultaneously.

For throughput-oriented applications like batch processing of documents, use the largest batch size that fits in memory. This maximizes GPU utilization and minimizes cost per example. For BERT-base inference on a V100 GPU, batch sizes of 64-128 typically achieve near-optimal throughput, processing 500-1000 examples per second depending on sequence length.

For latency-sensitive applications like real-time search or chatbots, dynamic batching provides a middle ground. Requests are accumulated for a short time window, typically 10-50 milliseconds, and

then processed together. This approach increases average latency slightly but dramatically improves throughput, allowing a single GPU to serve many more requests per second. The optimal batching window depends on your latency requirements and request rate.

23.5.2 KV Caching for Generation

Autoregressive generation in GPT-style models requires generating tokens sequentially, with each token depending on all previous tokens. Naive implementation recomputes attention for all previous tokens at each step, resulting in quadratic complexity in generation length. Key-value caching eliminates this redundancy by storing the key and value projections from previous tokens and reusing them for subsequent tokens.

With KV caching, generating a sequence of length n requires n forward passes through the model, but each pass only processes a single new token. The cached keys and values from previous tokens are concatenated with the new token's keys and values for attention computation. This reduces generation time by a factor of 5-10 \times for typical sequence lengths, making interactive generation practical.

The memory cost of KV caching is proportional to sequence length, batch size, and model size. For GPT-2 medium generating sequences of length 512 with batch size 8, the KV cache requires approximately 1 GB of memory. This is substantial but worthwhile given the dramatic speedup. For very long sequences or large batch sizes, memory constraints may limit the effectiveness of KV caching.

23.5.3 Quantization Strategies

Quantization reduces model size and accelerates inference by using lower-precision representations for weights and activations. FP16 quantization is the simplest approach, reducing model size by 2 \times with minimal accuracy loss. Modern GPUs with tensor cores also execute FP16 operations faster than FP32, providing both memory and speed benefits.

INT8 quantization provides more aggressive compression, reducing model size by 4 \times compared to FP32. This requires calibration to determine appropriate scaling factors for each layer, but modern tools like TensorRT and ONNX Runtime automate this process. INT8 quantization typically causes accuracy degradation of 0.5-2 percentage points on downstream tasks, which is acceptable for many applications.

INT4 and even lower precision quantization push compression further, achieving 8 \times or greater size reduction. However, accuracy degradation becomes more significant, and specialized hardware support is less common. These extreme quantization levels are most appropriate for deployment on edge devices where memory is severely constrained and some accuracy loss is acceptable.

Dynamic quantization applies quantization only to weights while keeping activations in higher precision. This provides a good balance between compression and accuracy, as weights are more tolerant of quantization than activations. For BERT models, dynamic INT8 quantization typically achieves 3-4 \times speedup with less than 1 percentage point accuracy loss.

23.5.4 Model Distillation

Knowledge distillation trains a smaller student model to mimic a larger teacher model, achieving better performance than training the student from scratch. The student learns from both the ground truth labels and the soft probability distributions produced by the teacher. This additional signal helps the student learn more effectively, often achieving 95-98% of the teacher's performance with 40-60% of the parameters.

DistilBERT exemplifies successful distillation, reducing BERT-base from 12 layers to 6 while retaining 97% of its performance on GLUE benchmarks. The distillation process uses a combination of distillation loss (matching teacher's output distributions), masked language modeling loss, and cosine embedding loss (matching hidden state directions). Training requires approximately the same compute as fine-tuning the teacher model.

The benefits of distillation compound with other optimization techniques. A distilled model can be further quantized and optimized, achieving cumulative speedups of 5-10 \times compared to the original

model. For production deployment, this combination of distillation and quantization often provides the best balance of performance and efficiency.

23.5.5 Hardware Selection

Choosing appropriate hardware for inference depends on throughput requirements, latency constraints, and cost considerations. For high-throughput batch processing, GPUs provide the best performance per dollar. NVIDIA T4 GPUs offer excellent inference performance at moderate cost, while A10 and A100 GPUs provide higher throughput for demanding applications.

For latency-critical applications with moderate throughput, CPUs may be more cost-effective than GPUs. Modern CPUs with AVX-512 instructions can achieve respectable inference performance, especially for smaller models and shorter sequences. Intel's Deep Learning Boost and similar technologies further accelerate neural network inference on CPUs.

Specialized inference accelerators like Google's TPUs, AWS Inferentia, and NVIDIA's TensorRT provide optimized performance for specific workloads. These platforms often achieve better performance per watt and performance per dollar than general-purpose GPUs, but require more effort to deploy and optimize. They are most appropriate for very high-volume applications where the engineering investment is justified.

23.6 Cost Optimization

Understanding and optimizing costs is essential for sustainable deployment of transformer models. This section provides detailed analysis of training and inference costs with concrete examples.

23.6.1 Training Cost Analysis

Training costs depend on model size, dataset size, and hardware selection. For BERT-base pre-training on 16 GB of text, the original paper reports using 16 TPU chips for 4 days, equivalent to approximately 64 TPU days. At current Google Cloud pricing of roughly \$8 per TPU hour, this amounts to approximately \$12,000 for pre-training. Using equivalent GPU resources (64 V100 GPUs for 4 days) would cost approximately \$15,000 at on-demand rates.

Fine-tuning costs are much more modest. Training BERT-base on a typical classification task with 10,000 examples requires 2-4 hours on a single V100 GPU. At AWS on-demand pricing of approximately \$3 per hour for a p3.2xlarge instance, this amounts to \$6-12 per fine-tuning run. Even with extensive hyperparameter search involving 20-30 runs, total costs remain under \$300.

Cloud versus on-premise costs depend heavily on utilization. For continuous training workloads, purchasing GPUs becomes cost-effective after 12-18 months of use. A DGX A100 system costs approximately \$200,000 but provides compute equivalent to \$15,000 per month at cloud on-demand rates. For intermittent workloads or experimentation, cloud computing is more economical due to the flexibility to scale up and down.

Spot instances provide substantial savings for training workloads that can tolerate interruptions. AWS spot instances for p3.2xlarge typically cost 50-70% less than on-demand rates, reducing fine-tuning costs to \$2-4 per run. Implementing checkpointing and automatic restart logic allows training to resume after spot instance interruptions, making this an attractive option for cost-conscious training.

23.6.2 Training Time Estimation

Estimating training time helps with planning and cost prediction. For fine-tuning, a useful rule of thumb is that BERT-base processes approximately 100-150 examples per second on a V100 GPU with batch size 32 and sequence length 128. For a dataset of 100,000 examples trained for 3 epochs, this translates to 2,000-3,000 seconds or roughly 1 hour of training time.

Pre-training time scales with dataset size and model size. BERT-base pre-training on 16 GB of text requires approximately 1 million training steps with batch size 256, processing roughly 4 billion tokens. At 1,000 tokens per second per V100 GPU, this requires 4 million GPU-seconds or approximately 1,100 GPU-hours. With 16 GPUs, this translates to roughly 70 hours or 3 days of training.

Larger models scale approximately linearly with parameter count for training time. GPT-2 medium with 345 million parameters takes roughly $3\times$ longer to train than GPT-2 small with 117 million parameters, assuming the same dataset and batch size. However, larger models often benefit from larger batch sizes, which can partially offset the increased time per step.

23.6.3 Inference Cost Analysis

Inference costs depend on request volume, latency requirements, and model size. For a BERT-base classification service processing 1 million requests per day with average latency requirements of 100ms, a single V100 GPU can handle approximately 100 requests per second with dynamic batching, or 8.6 million requests per day. This suggests that a single GPU is sufficient, costing approximately \$200-300 per month for a cloud GPU instance.

For generation tasks, costs are higher due to the sequential nature of autoregressive decoding. GPT-2 medium generating 100 tokens per request can process approximately 10-20 requests per second per GPU, depending on batch size and sequence length. For 1 million requests per day, this requires 1-2 GPUs, costing \$400-600 per month. The cost per million tokens is approximately \$5-10 for self-hosted inference.

Comparing self-hosted to API costs reveals significant differences at scale. OpenAI's GPT-3.5 API costs approximately \$2 per million tokens for input and output combined. For applications processing 100 million tokens per month, this amounts to \$200 per month. Self-hosting a comparable model would require 4-8 GPUs costing \$1,600-3,200 per month, making the API more economical at this scale. However, at 1 billion tokens per month, self-hosting becomes competitive, and at 10 billion tokens per month, self-hosting is clearly more economical.

23.6.4 Cost Optimization Strategies

Several strategies can substantially reduce both training and inference costs. For training, using mixed precision reduces training time by $2-3\times$, directly reducing costs by the same factor. Gradient accumulation allows using smaller, cheaper GPU instances by simulating larger batch sizes. Spot instances reduce costs by 50-70% for workloads that can tolerate interruptions.

For inference, quantization and distillation reduce both latency and cost. A distilled and quantized model may achieve $5-10\times$ higher throughput than the original model, allowing a single GPU to handle the load that previously required 5-10 GPUs. This directly translates to $5-10\times$ cost reduction. Dynamic batching improves GPU utilization from 20-30% to 70-90%, effectively tripling throughput without additional hardware.

Caching can dramatically reduce inference costs for applications with repeated queries. If 30% of requests are duplicates or near-duplicates, caching responses eliminates 30% of inference costs. Semantic caching using embedding similarity can extend this to near-duplicate queries, potentially caching 50-70% of requests in some applications.

Autoscaling based on demand prevents paying for idle resources during low-traffic periods. For applications with diurnal traffic patterns, autoscaling can reduce costs by 40-60% compared to provisioning for peak load. Kubernetes and cloud-native deployment platforms make autoscaling straightforward to implement.

23.7 Production Deployment

Deploying transformer models in production requires careful attention to serving infrastructure, monitoring, and operational practices. This section covers the key considerations for reliable production deployment.

23.7.1 Model Serving Frameworks

Several frameworks specialize in serving transformer models efficiently. TorchServe provides a production-ready serving solution for PyTorch models with built-in support for batching, versioning, and metrics.

It handles request queuing, dynamic batching, and model lifecycle management, making it suitable for production deployment with minimal custom code.

ONNX Runtime offers cross-platform inference with optimizations for various hardware backends. Converting models to ONNX format enables deployment on CPUs, GPUs, and specialized accelerators with a single model artifact. ONNX Runtime includes graph optimizations and kernel fusion that can improve inference speed by $2\text{--}3\times$ compared to native PyTorch inference.

TensorRT provides highly optimized inference for NVIDIA GPUs, achieving the best possible performance on this hardware. It performs aggressive optimizations including layer fusion, precision calibration, and kernel auto-tuning. For latency-critical applications on NVIDIA hardware, TensorRT often provides $2\text{--}5\times$ speedup compared to other serving solutions.

Triton Inference Server supports multiple frameworks and provides advanced features like model ensembling, dynamic batching, and concurrent model execution. It is particularly well-suited for complex serving scenarios involving multiple models or preprocessing pipelines. The learning curve is steeper than simpler solutions, but the flexibility and performance justify the investment for large-scale deployments.

23.7.2 Scaling Strategies

Horizontal scaling distributes load across multiple model instances, providing both higher throughput and fault tolerance. Load balancers distribute requests across instances, and autoscaling adjusts the number of instances based on demand. This approach is straightforward to implement and scales well to very high request volumes.

Vertical scaling uses more powerful hardware for each instance, such as GPUs with more memory or faster CPUs. This approach is simpler to manage than horizontal scaling but has limits, as individual hardware units have maximum capacities. Vertical scaling is often combined with horizontal scaling, using appropriately-sized instances and scaling the number of instances as needed.

Model parallelism splits large models across multiple GPUs within a single instance, enabling serving of models too large for a single GPU. This requires more sophisticated implementation but is necessary for very large models. Pipeline parallelism can improve efficiency by processing multiple requests concurrently through different stages of the model.

23.7.3 Monitoring and Logging

Comprehensive monitoring is essential for maintaining reliable production systems. At minimum, track request rate, latency percentiles (p50, p95, p99), error rate, and resource utilization (CPU, GPU, memory). These metrics provide early warning of issues and help with capacity planning.

Model-specific metrics provide insight into model behavior. Track input length distributions to detect shifts in request patterns. Monitor output confidence scores to identify cases where the model is uncertain. Log a sample of inputs and outputs for qualitative analysis and debugging.

Alerting should trigger on both immediate issues and gradual degradation. Set alerts for error rate spikes, latency increases, and resource exhaustion. Also monitor for gradual trends like slowly increasing latency or decreasing confidence scores, which may indicate model drift or infrastructure degradation.

Distributed tracing helps debug issues in complex serving pipelines. Tools like Jaeger or Zipkin track requests through preprocessing, model inference, and postprocessing, identifying bottlenecks and failures. This visibility is invaluable when debugging performance issues or errors in production.

23.7.4 A/B Testing and Deployment

A/B testing enables safe deployment of model updates by gradually rolling out changes to a subset of traffic. Start by routing 5-10% of traffic to the new model while monitoring key metrics. If metrics remain stable or improve, gradually increase traffic to the new model. If metrics degrade, roll back immediately.

Shadow deployment runs the new model alongside the existing model without affecting user-facing results. Requests are sent to both models, but only the existing model's outputs are returned to

users. This allows testing the new model's performance and behavior with production traffic before committing to deployment.

Canary deployment is similar to A/B testing but focuses on detecting errors rather than measuring performance improvements. A small percentage of traffic (1-5%) is routed to the new model, and error rates are monitored closely. If error rates remain acceptable for several hours or days, the deployment proceeds to full rollout.

Blue-green deployment maintains two complete production environments and switches traffic between them. This enables instant rollback if issues are detected, as switching back to the previous environment is immediate. The cost is maintaining duplicate infrastructure, but the operational safety is valuable for critical applications.

23.8 Practical Checklists

These checklists provide systematic guidance for common transformer workflows, helping ensure that critical steps are not overlooked.

23.8.1 Before Training Checklist

Before beginning training, verify that you have made appropriate decisions about resources and configuration. Estimate memory requirements using the formulas provided earlier, ensuring that your chosen batch size and sequence length will fit in available GPU memory with some margin for safety. Select hardware appropriate for your model size and training duration, considering the trade-offs between cost and training time.

Choose batch size and sequence length based on your task requirements and memory constraints. Remember that sequence length has a quadratic effect on memory, so reducing it provides substantial savings if your task permits. Set up monitoring and logging infrastructure before training begins, as debugging issues after the fact is much more difficult than catching them in real-time.

Estimate training time and cost using the guidelines provided earlier. This helps with planning and ensures that you have allocated sufficient budget and time for the training run. For long training runs, verify that checkpointing is configured correctly and test recovery from checkpoints before committing to the full training run.

23.8.2 During Training Checklist

While training is in progress, monitor loss and metrics regularly to detect issues early. Training loss should decrease steadily, though not necessarily monotonically. Validation loss should track training loss initially, with some divergence expected as training progresses. If validation loss increases while training loss decreases, you may be overfitting.

Check GPU utilization to ensure that you are using resources efficiently. Utilization should be consistently above 80% during training. Lower utilization suggests that batch size is too small, data loading is a bottleneck, or there are inefficiencies in the training loop. Monitor memory usage to ensure you are not close to OOM errors, which can cause training to fail unexpectedly.

Save checkpoints regularly according to your checkpointing strategy. Verify that checkpoints are being saved successfully and that you can load them for recovery. Validate periodically on a held-out set to track generalization performance. The frequency of validation depends on training duration, but every few hundred steps or every epoch is typical.

23.8.3 Before Deployment Checklist

Before deploying a model to production, optimize it for inference using the techniques described earlier. Apply quantization if accuracy permits, as the performance benefits are substantial. Consider distillation if you need further speedup and have time for the additional training. Export the model to an optimized format like ONNX or TensorRT if using those serving frameworks.

Benchmark latency and throughput under realistic conditions, including the batch sizes and sequence lengths you expect in production. Test with both average-case and worst-case inputs to understand performance variability. Estimate serving costs based on expected request volume and the hardware required to meet latency requirements.

Set up monitoring and alerting for the production deployment. Ensure that you can track request rate, latency, error rate, and resource utilization. Configure alerts for anomalies in these metrics. Plan your scaling strategy, including autoscaling rules if using dynamic scaling.

Test the deployment pipeline end-to-end, including model loading, preprocessing, inference, and postprocessing. Verify that error handling works correctly and that failures are logged appropriately. Conduct load testing to ensure the system can handle expected traffic with appropriate margins for spikes.

23.9 Hyperparameter Tuning

23.9.1 Critical Hyperparameters (Ordered by Impact)

1. Learning Rate (highest impact)

- Typical range: $[10^{-5}, 10^{-3}]$
- BERT: $1 - 5 \times 10^{-5}$
- GPT: $2 - 6 \times 10^{-4}$
- Rule: Larger models need smaller LR

2. Batch Size

- Trade-off: Speed vs generalization
- Typical: 32-512 for fine-tuning, 256-2048 for pre-training
- Scale LR linearly with batch size

3. Warmup Steps

- Typical: 5-10% of total training steps
- BERT: 10,000 steps
- GPT-3: 375M tokens (out of 300B)

4. Weight Decay

- Typical: 0.01 to 0.1
- AdamW: Decouple from learning rate

5. Dropout

- Standard: 0.1
- Larger models: Lower dropout (0.05 or none)
- Apply uniformly (attention, FFN, embeddings)

23.9.2 Tuning Strategy

Phase 1: Coarse search

- Grid/random search over wide ranges
- Short runs (10% of full training)
- Focus on learning rate first

Phase 2: Fine search

- Narrow ranges around best from Phase 1
- Longer runs (50% of full training)
- Tune other hyperparameters

Phase 3: Validation

- Full training with best settings
- Multiple seeds for robustness
- Final evaluation on test set

Example 23.2 (Learning Rate Search). **Task:** Fine-tune BERT on classification

Coarse search:

- Try: $[10^{-5}, 3 \times 10^{-5}, 10^{-4}, 3 \times 10^{-4}]$
- Train 1 epoch each
- Best: 3×10^{-5} (85.2% dev accuracy)

Fine search:

- Try: $[2 \times 10^{-5}, 3 \times 10^{-5}, 4 \times 10^{-5}]$
- Train 3 epochs each
- Best: 3×10^{-5} (86.1% dev accuracy)

Final:

- Train with $LR = 3 \times 10^{-5}$, 5 epochs
- Test accuracy: 85.8%

23.10 Common Pitfalls and Solutions

23.10.1 Architecture Pitfalls

Pitfall 1: Forgetting positional information

- Symptom: Model treats sequence as bag-of-words
- Solution: Verify position encoding is added

Pitfall 2: Incorrect masking

- Symptom: Information leakage or blocked attention
- Solution: Visualize attention matrices, verify mask shape

Pitfall 3: Not sharing embeddings

- Symptom: Twice as many parameters as expected
- Solution: Weight tying between input/output embeddings

23.10.2 Training Pitfalls**Pitfall 4: Insufficient warmup**

- Symptom: Training unstable early, doesn't recover
- Solution: Increase warmup to 10% of training

Pitfall 5: Wrong learning rate scale

- Symptom: Loss not decreasing or diverging
- Solution: Learning rate finder, try 10× up/down

Pitfall 6: Overfitting small datasets

- Symptom: Large train/val gap
- Solution: More dropout, data augmentation, smaller model

23.10.3 Deployment Pitfalls**Pitfall 7: Batch size 1 in production**

- Symptom: Poor GPU utilization
- Solution: Dynamic batching, accumulate requests

Pitfall 8: Not using mixed precision

- Symptom: Slow inference, high memory
- Solution: FP16 inference, quantization

Pitfall 9: No KV caching for generation

- Symptom: Slow text generation (quadratic in length)
- Solution: Cache key/value tensors

23.11 Case Study: BERT for Search Ranking**23.11.1 Problem Setup****Task:** Rank search results by relevance**Input:** Query + Document pairs**Output:** Relevance score $[0, 1]$ **23.11.2 Architecture Decisions****Model:** BERT-base with regression head**Input format:**`[CLS] query tokens [SEP] document tokens [SEP]`**Output:** $\text{score} = \sigma(\mathbf{W}\mathbf{h}_{[\text{CLS}]} + b)$

23.11.3 Training Strategy

Data:

- 10M query-document pairs
- Labels: Click-through rate (0-1)
- Hard negatives: Top results without clicks

Loss: Mean squared error on CTR prediction

Optimization:

- Learning rate: 2×10^{-5}
- Batch size: 256
- Warmup: 10,000 steps
- Total: 100,000 steps

23.11.4 Production Deployment

Optimizations:

1. Quantize to INT8 (3× speedup)
2. Distill to 6-layer model (2× speedup)
3. Deploy with ONNX Runtime
4. Dynamic batching (avg batch size 32)

Results:

- Latency: 15ms p99 (vs 200ms baseline)
- Throughput: 2000 QPS per GPU
- Relevance: +8% improvement over TF-IDF

23.12 Case Study: GPT for Code Generation

23.12.1 Problem Setup

Task: Generate Python code from natural language

Example:

Input: "Function to reverse a string"

Output:

```
def reverse_string(s):
    return s[::-1]
```

23.12.2 Model and Data

Model: GPT-2 medium (345M params)

Data:

- GitHub public repositories (Python)
- Filtered: Only files with docstrings
- Format: Docstring → Implementation
- Total: 50GB, 10B tokens

23.12.3 Training

Pre-training: Start from GPT-2 checkpoint

Fine-tuning:

- 100,000 steps
- Learning rate: 5×10^{-5}
- Context: 1024 tokens
- Batch: 128 sequences

23.12.4 Evaluation

Metrics:

- Pass@k: % correct in top-k samples
- BLEU: Token overlap with reference
- Human evaluation: Correctness + readability

Results:

- Pass@1: 42%
- Pass@10: 71%
- Human preferred over baseline: 78%

23.13 Future Directions

23.13.1 Architectural Innovations

1. Efficient attention

- Linear complexity methods
- State space models (S4, Mamba)
- Hybrid CNN-attention architectures

2. Multimodal integration

- Unified text-image-audio models
- Better cross-modal alignment
- Efficient fusion strategies

3. Long context

- Million-token contexts
- Hierarchical memory
- Retrieval-augmented transformers

23.13.2 Training Innovations

1. Sample efficiency

- Better pre-training objectives
- Curriculum learning
- Few-shot and zero-shot learning

2. Scaling

- Mixture of experts
- Conditional computation
- Efficient parallelism strategies

3. Alignment

- Better RLHF techniques
- Constitutional AI
- Value alignment

23.14 Conclusion

23.14.1 Key Takeaways

Architecture:

- Attention is powerful and flexible
- Position encodings crucial for sequences
- Residuals + normalization enable depth

Training:

- Pre-training + fine-tuning is dominant paradigm
- Warmup is critical for stability
- Scale requires careful optimization

Deployment:

- Quantization and distillation for efficiency
- Batching crucial for throughput
- Monitor performance in production

23.14.2 Final Advice

For practitioners:

1. Start simple: Use pre-trained models
2. Debug systematically: Data, model, training
3. Optimize iteratively: Accuracy first, then speed
4. Monitor continuously: Metrics, errors, drift

For researchers:

1. Understand fundamentals deeply
2. Question assumptions: Why does this work?
3. Experiment rigorously: Ablations, multiple seeds
4. Share knowledge: Open source, papers, blogs

This concludes our comprehensive journey through deep learning and transformers. You now have the mathematical foundations, practical implementations, and real-world insights to build state-of-the-art transformer models!

23.15 Exercises

Exercise 23.1. Reproduce DistilBERT:

1. Train 6-layer student on BERT-base teacher
2. Use distillation + MLM + cosine losses
3. Evaluate on GLUE
4. Measure compression ratio and speedup

Exercise 23.2. Debug broken transformer (provided):

1. Model trains but poor performance
2. Find 3 subtle bugs (architecture, training, data)
3. Fix and verify improvements

Exercise 23.3. Deploy BERT for production:

1. Fine-tune on classification task
2. Quantize to INT8
3. Export to ONNX
4. Create REST API with FastAPI

5. Load test and optimize