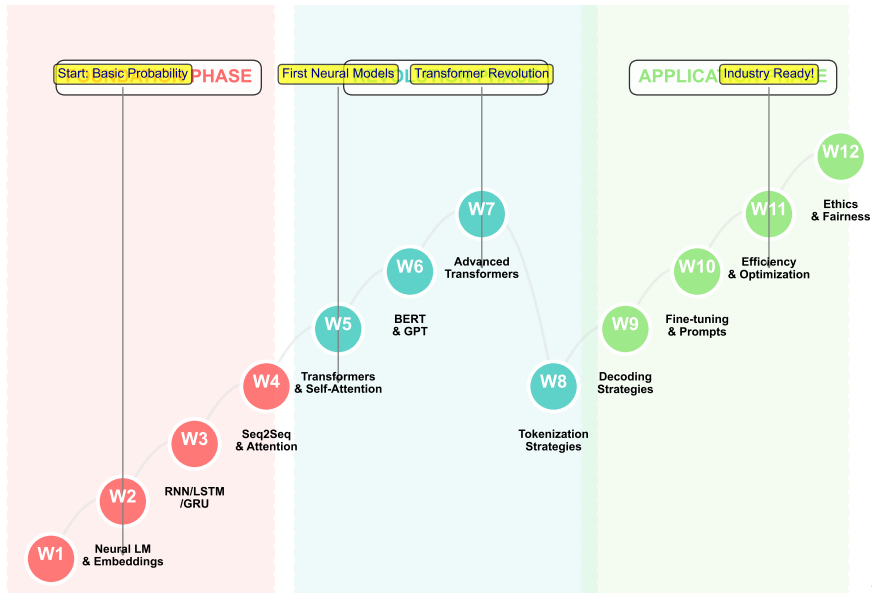# Natural Language Processing: Complete Course Overview
## From Foundations to State-of-the-Art
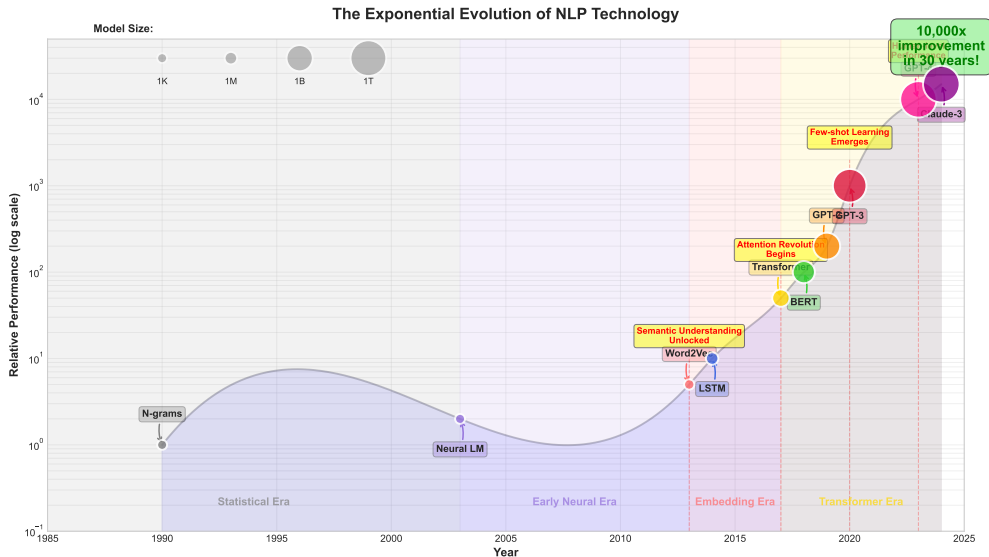
Joerg R. Osterrieder
www.joergosterrieder.com

BSc Computer Science - 12 Week Journey
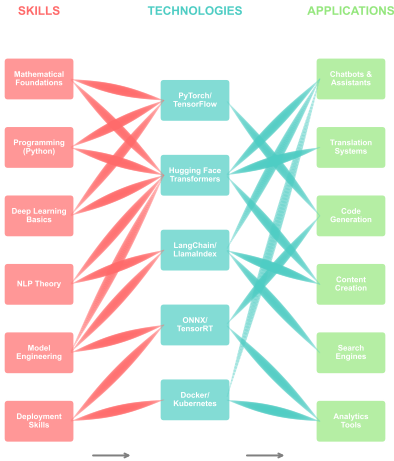
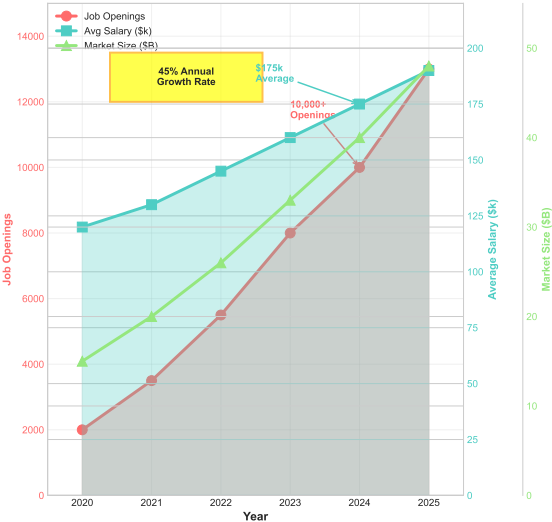# The NLP Journey: From Counting to Understanding



Start: Basic Probability PHASE

First Neural Models | Transformer Revolution

APPLICA | Industry Ready!

- W12 — Ethics & Fairness
- W11 — Efficiency & Optimization
- W10
- W9 — Fine-tuning & Prompts
- W8 — Tokenization Strategies
- W7 — Advanced Transformers
- W6 — BERT & GPT
- W5 — Transformers & Self-Attention
- W4
- W3 — Seq2Seq & Attention
- W2 — RNN/LSTM /GRU
- W1 — Neural LM & Embeddings
- W8 — Decoding Strategies

**Complexity:**

⭐ Low

The Exponential Evolution of NLP Technology

# Learning Outcomes & Real-World Impact



From Learning to Real-World Impact

SKILLS — TECHNOLOGIES — APPLICATIONS



NLP Industry Growth Metrics

Core Skills          Technologies          Applications

Week 1 Learning Journey

**N-gram Extraction**

**Probability Calculation**

**P(jumps | brown fox)**

Unigrams:
| The | quick | brown | fox | jumps |

Bigrams:
| The quick | quick brown | brown fox | fox jumps |

Trigrams:
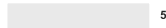| The quick brown | quick brown fox | brown fox jumps |

brown fox — 10

... jumps — 3

... runs — 2

... walks — 5

P(jumps | brown fox) = 3/10 = 0.3

## Markov Assumption
- Future depends on recent past
- Makes computation tractable

## Probability
- Count n-grams
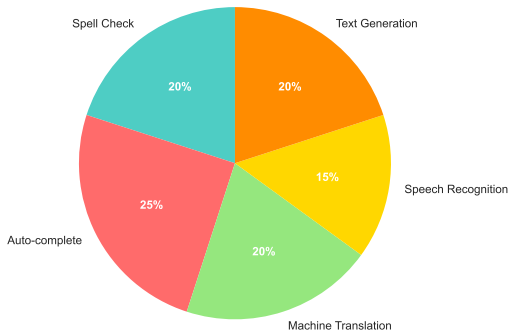- Maximum likelihood

## Challenges
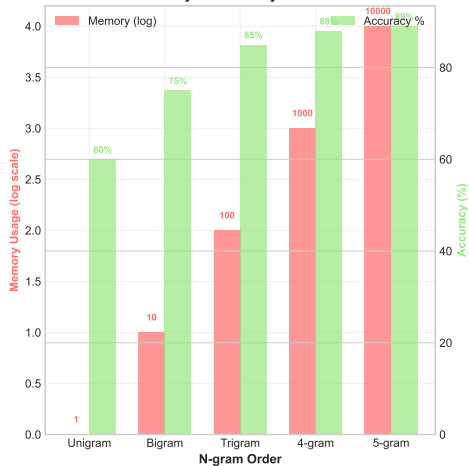- Data sparsity
- Context limitations

N-gram Applications



Memory vs Accuracy Trade-off

**Where N-grams Excel:**
- Spell checking
- Auto-complete

**Historical Impact:**
- Dominated 1980s-2000s
- Still used in hybrid systems

## From Counts to Continuous



## Neural LM Architecture



**Paradigm Shift:**
- From counting to learning
- Continuous space

**Key Innovations:**
- Distributed representations
- Backpropagation training

# Week 2: Word Embeddings



**Word Embedding Space (2D projection)**

**High-Dimensional Semantic Space**

**Properties**
- Similar words cluster
- Analogies work

**Methods**
- Word2Vec
- GloVe, FastText

**Dimensions**
- 50-300 dims typical
- Semantic features

**Word2Vec Algorithms**

**CBOW**           **Skip-gram**

The           sat

?           ?

on           ?

mat           ?

                                           ?

Predict center from context           Predict context from center

**Representation Quality**

- Similarity Task
- Analogy Task

One-hot: 20% / 10%
TF-IDF: 35% / 15%
Word2Vec: 85% / 75%
GloVe: 82% / 73%
FastText: 88% / 80%

**Neural Breakthrough!**

Accuracy (%) — Method

**Algorithms:**
- CBOW: Context → center
- Skip-gram: Center → context

**Applications:**
- Sentiment analysis
- Named entity recognition

**RNN: The Sequential Processing Bottleneck**
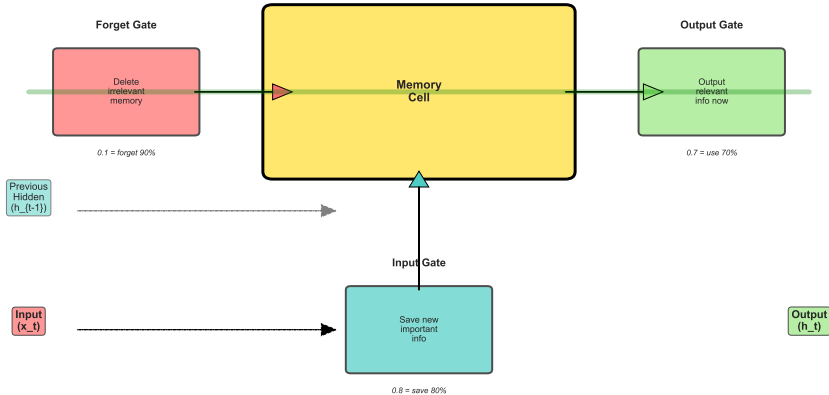


**What We'll Learn:**
- RNN architecture
- LSTM gates and memory

**Core Concepts:**
- **Vanishing Gradients**
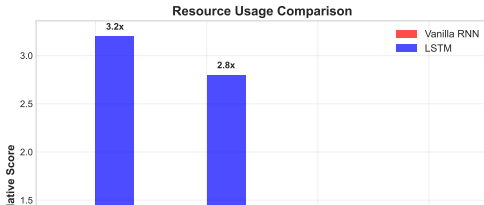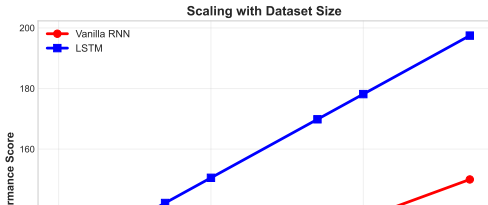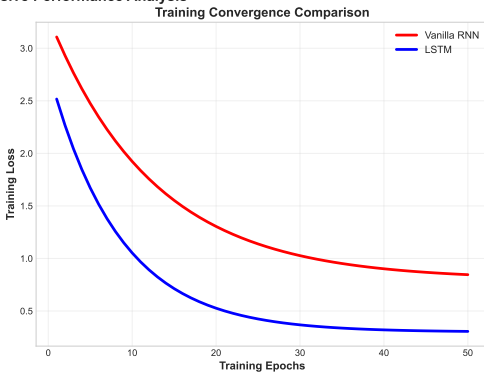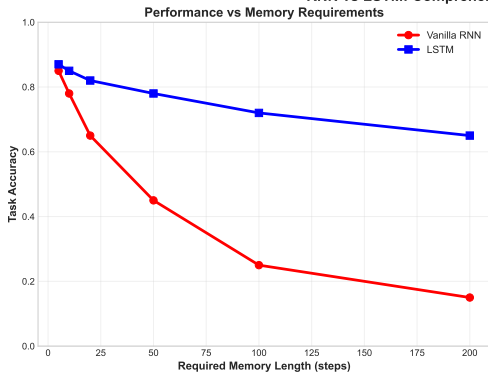  - Problem with long sequences
  - LSTM/GRU solutions

**LSTM Gates: Controlling Information Flow**
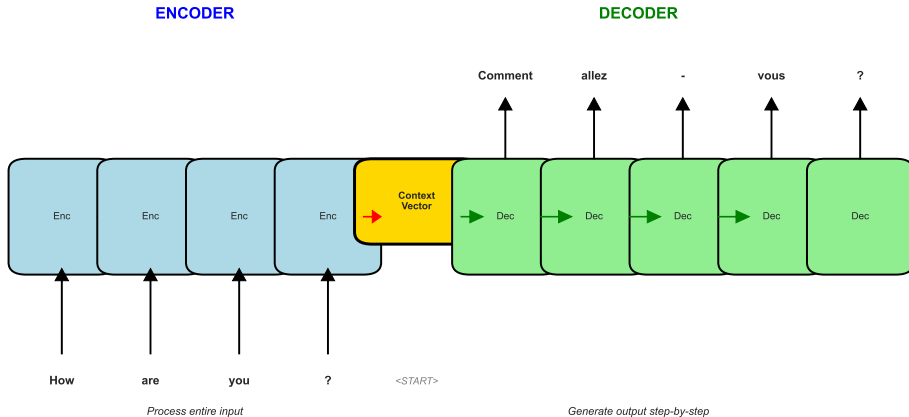
**Cell State Highway (Long-term Memory)**

**Forget Gate**

Delete
irrelevant
memory

*0.1 = forget 90%*

**Memory
Cell**

**Output Gate**

Output
relevant
info now

*0.7 = use 70%*

Previous
Hidden
(h_{t-1})

**Input Gate**

Input
(x_t)

Save new
important
info

*0.8 = save 80%*

Output
(h_t)

RNN vs LSTM: Comprehensive Performance Analysis

## Sequence-to-Sequence Architecture

*Variable input length → Fixed context → Variable output length*

**ENCODER**

**DECODER**

Comment     allez     -     vous     ?



| Enc | Enc | Enc | Enc | Context Vector | Dec | Dec | Dec | Dec | Dec |

How     are     you     ?     <START>

*Process entire input*          *Generate output step-by-step*

What We'll Learn:          Key Concepts:

Attention Weights: English to French Translation

Machine Translation Quality Over Time
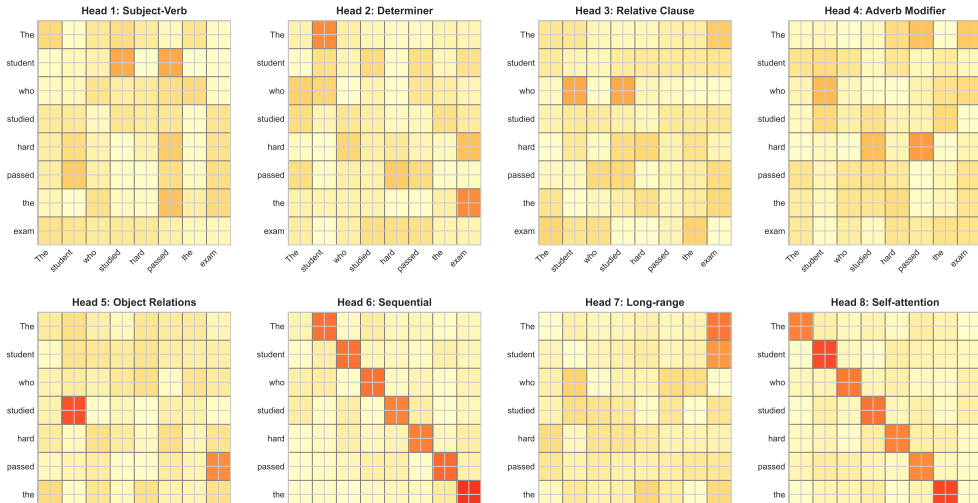The Seq2Seq Revolution
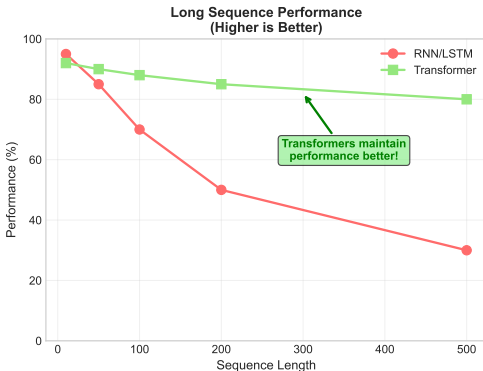
# Transformer Architecture

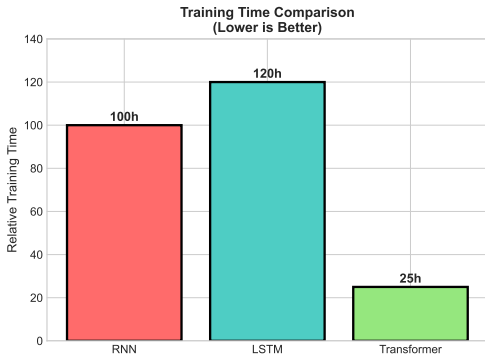*(Simplified for BSc Understanding)*

## Multi-Head Attention: 8 Different Perspectives on the Same Sentence
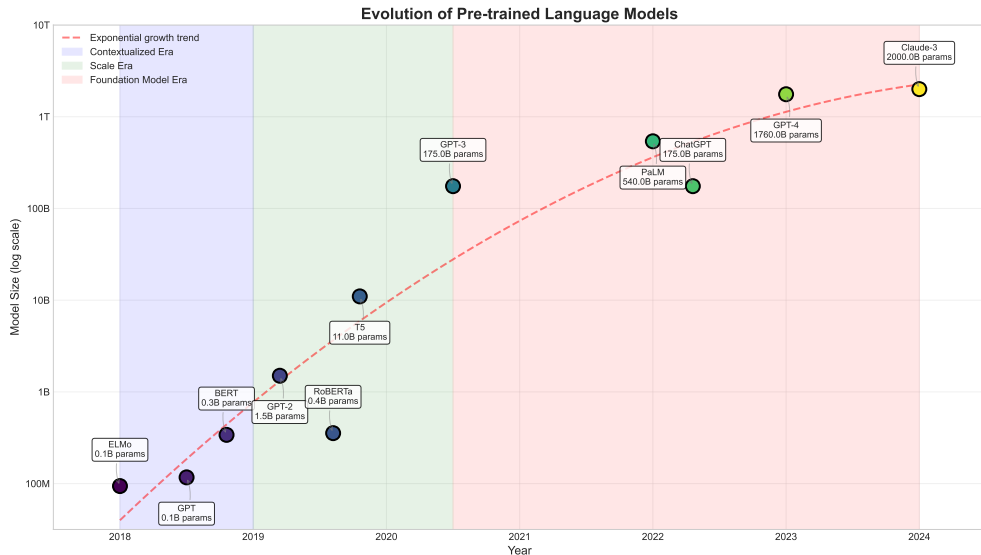
**Why Transformers Beat RNNs**

**Performance:**
- SOTA on all NLP tasks
- 100x faster training
- Scales to billions of parameters
- Transfer learning enabled

**Applications:**
- Foundation for BERT/GPT
- Computer Vision (ViT)
- Protein Folding
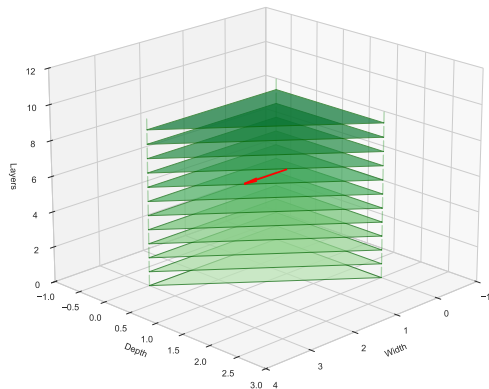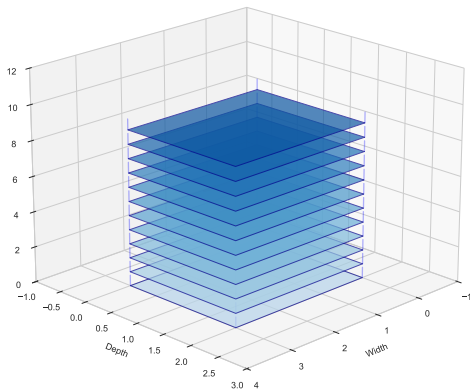- Multimodal models

Evolution of Pre-trained Language Models

BERT: Bidirectional Encoder    BERT vs GPT: Architectural Differences    GPT: Unidirectional Decoder

**BERT (Bidirectional):**
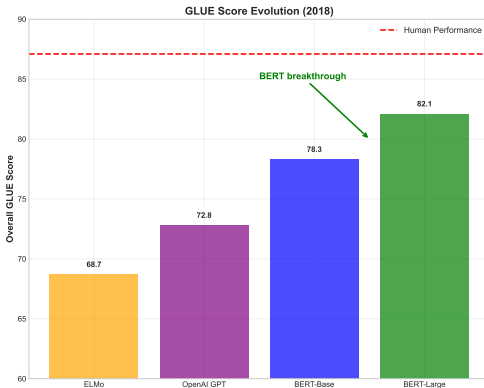
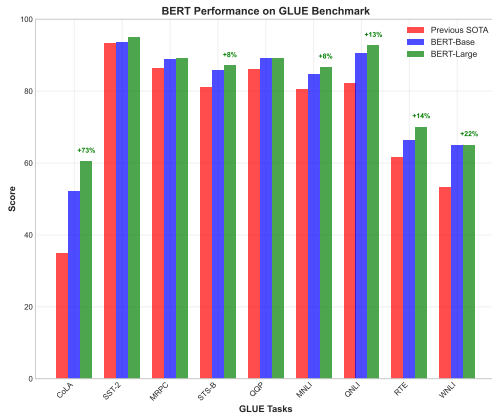- Masked Language Model
- Sees full context

**GPT (Autoregressive):**

- Next token prediction
- Left to right only

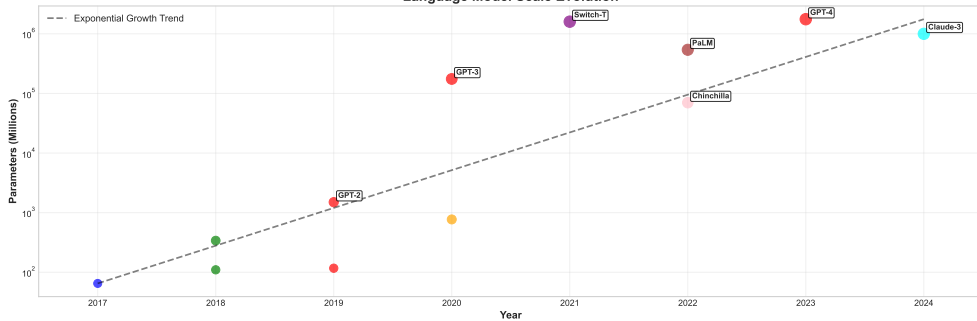BERT: Dominating Natural Language Understanding

**Performance:**
- GLUE benchmark SOTA
- Human-level on many tasks
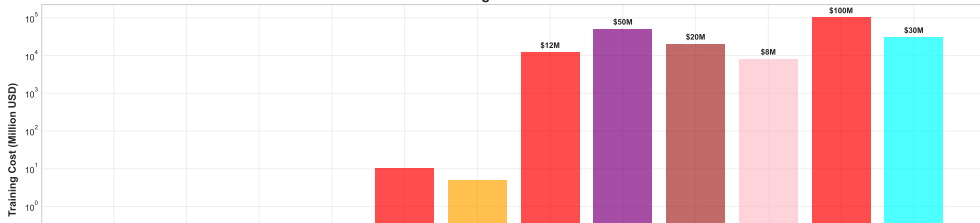- Few-shot learning
- Zero-shot capabilities

**Ecosystem:**
- Hugging Face hub
- 100,000+ models
- Easy fine-tuning
- Production ready

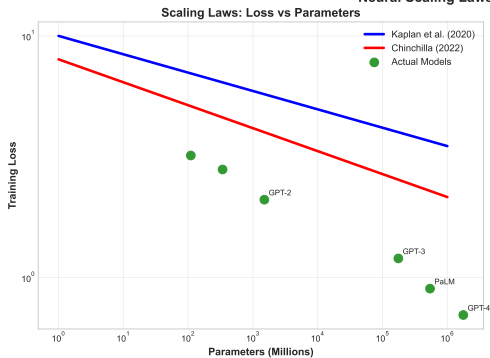The Scale Revolution: From Millions to Trillions

Language Model Scale Evolution

Training Cost Evolution

Neural Scaling Laws: The Science of Scale
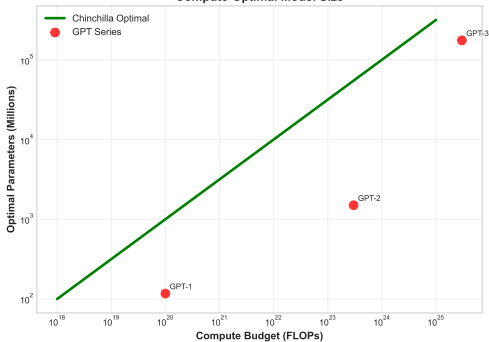
GPT-3: The Generalist AI That Changed Everything

The Hidden Impact of Tokenization

**Byte Pair Encoding: Learning Subwords**

*Corpus: "low lower lowest"*

**Tokenization Across Languages**

*Same meaning, different token counts*

English baseline

English: Hello world — 2 tokens

Chinese: 你好世界 — 4 tokens

Arabic: مرحبا بالعالم — 4 tokens

Japanese: こんにちは世界 — 4 tokens

Korean: 안녕 세계 — 3 tokens

Emoji: Hello — 5 tokens

Cost Impact:
Chinese text costs
2x more than English
on GPT-4!

Use Cases:

Best Practices:

## Decoding Strategy Selection Guide

Creativity

| Factual Q&A | Greedy | $T=0.1$ | Maximum accuracy |
| Code Generation | Beam-5 | $T=0.2, p=0.95$ | Syntactic correctness |
| Translation | Beam-4 | $T=0.3$ | Preserve meaning |
| Summarization | Top-p=0.9 | $T=0.5$ | Balance accuracy/fluency |
| Dialogue | Top-p=0.9 | $T=0.7$ | Natural conversation |
| Creative Writing | Top-k=50 | $T=0.9, p=0.95$ | Maximum creativity |

Temperature Controls Probability Distribution Sharpness

**Greedy**
- Fastest
- Deterministic
- Can be repetitive

**Beam Search**
- Better quality
- Multiple hypotheses
- More compute

**Sampling**
- Creative
- Diverse outputs
- Temperature control

Decoding Strategy Performance Analysis

**Adaptation Strategy Decision Tree**

## LoRA: Low-Rank Adaptation

*Instead of updating 16M parameters, update only 32K!*



**W**
**(d×d)**
**Frozen**

Original Weights
(e.g., 4096×4096)

*Freeze*

**A**
**(d×r)**

× **B (r×d)**

+

**W + ΔW**
**(d×d)**
**Adapted**

Example: d=4096, r=8
Original: 4096×4096 = 16,777,216 parameters
LoRA: (4096×8) + (8×4096) = 65,536 parameters (0.39%!)

Fine-tuning vs Prompting: Multi-dimensional Comparison

*Prompting excels at flexibility and speed*

*LoRA balances performance and efficiency*

*Full fine-tuning for maximum accuracy*

**Domains:**   **Best Practices:**

Model Compression Landscape: Size vs Performance Trade-off

What We'll Learn:

Goals:

Quantization: Trading Precision for Efficiency

**Model Deployment Pipeline: From Research to Production**

| Pre-trained Model | Distillation | Quantization | Pruning | Deployment |
|---|---|---|---|---|
| 440MB BERT-base | 66MB DistilBERT | 17MB INT8 | 5MB Sparse | Mobile Ready |

**Size:** 440MB    66MB    17MB    5MB    5MB

**Latency:** 100ms    80ms    40ms    20ms    15ms

**Deployment Targets:**

**Accuracy:** 100%    97%    96%    94%    94%

| Cloud GPU | Edge Server | Mobile Phone | IoT Device |
|---|---|---|---|

*Each stage trades model size for deployment flexibility*

**Platforms:**

**Metrics:**

**The AI Ethics Landscape: Interconnected Challenges**

*Each ethical dimension affects and is affected by others*



Racial bias

Fairness

Gender bias

Re-identification

Adversarial attacks

Privacy

Economic bias

Data leaks

Safety

Misuse

Surveillance

Explainability

Alignment

Governance

Transparency

Black box

Accountability

Liability

Interpretability

Resource use

Regulation

## How Bias Enters AI Systems: From Society to Model to Impact

| Historical Bias | Data Collection | Annotation Bias | Model Bias | Deployment Bias |
|---|---|---|---|---|
| *Society reflects past injustices* | *Who/what is included?* | *Labeler subjectivity* | *Algorithm choices* | *Usage context* |

Example: Crime data reflects policing bias

Example: Internet text excludes many voices

Example: "Professional" means different things

Example: Objective functions encode values

Example: Hiring tool amplifies inequality

**Feedback Loop: Biased outputs reinforce societal biases**

**Fairness Interventions Across the ML Pipeline**

**Fairness Metrics**

| Demographic Parity | Equal Opportunity | Equalized Odds | Individual Fairness |

**Pre-processing** · **In-processing** · **Post-processing**

**Pre-processing**
- Synthetic data
- Re-sampling
- Data augmentation

**In-processing**
- Constraints
- Adversarial debiasing
- Fair objectives

**Post-processing**
- Fair ranking
- Threshold optimization
- Output calibration

*Note: Different fairness metrics often conflict - perfect fairness is mathematically impossible*

**Technical:**　　　　**Societal:**